

# Lineare und Ganzzahlige Programmierung

(Algorithmische Diskrete Mathematik II)

Skriptum zur Vorlesung im WS 2009/2010

Prof. Dr. Martin Grötschel  
Institut für Mathematik  
Technische Universität Berlin

Version vom 9. Februar 2010



## Vorwort

Die Vorlesung „Lineare und Ganzzahlige Programmierung“ ist die zweite Vorlesung (ADM II) im Zyklus des Studienschwerpunktes „Algorithmische Diskrete Mathematik“ an der TU Berlin in den Studiengängen Mathematik, Techno- und Wirtschaftsmathematik. Detaillierte Kenntnis der vorausgegangenen Vorlesung „Graphen- und Netzwerkalgorithmen“ wird nicht vorausgesetzt. Graphen- und Netzwerke und damit zusammenhängende Probleme werden lediglich als Beispielmaterial verwendet.

In dieser Vorlesung werden die Grundlagen der linearen Optimierung (auch lineare Programmierung genannt) behandelt, also einige Aspekte der linearen Algebra und der Polyedertheorie. Natürlich steht die Entwicklung von Algorithmen zur Lösung derartiger Optimierungsprobleme im Vordergrund. Sehr ausführlich werden der Simplexalgorithmus und seine Varianten behandelt, aber ebenso werden wir Innere-Punkte-Verfahren und die Ellipsoidmethode diskutieren.

Im zweiten Teil der Vorlesung geht es um ganzzahlige und gemischt-ganzzahlige Optimierung. Wir beginnen mit der Behandlung von Relaxierungen und dualen Heuristiken, besprechen dann Branch&Bound-Verfahren und führen Schnittebenenverfahren ein.

Oktober 2009

M. Grötschel

## Vorbemerkungen

### Literatur

Die Literatur zum Thema „Lineare Programmierung“ ist überwältigend umfangreich. Hierzu rufe man einfach einmal die Datenbank MATH des Zentralblatts für Mathematik auf und starte eine Suche nach Büchern, die die Wörter *linear* und *programming* im Titel enthalten. Man erhält derzeit über 200 Referenzen. Fast 40 Bücher enthalten die beiden Wörter *lineare Optimierung* und fast 30 die Wörter *lineare Programmierung* in ihrem Titel.

Eine Google-Suche am 11. Oktober 2009 nach „linear programming“ liefert derzeit ungefähr 1.6 Millionen Ergebnisse. Der (bei meiner Suche) am höchsten bewertete Eintrag war die Wikipedia-Erläuterung. Unter den ersten zehn Erwähnungen war auch „Linear Programming FAQs“, siehe

<http://www-unix.mcs.anl.gov/otc/Guide/faq/>,

welche auf den NEOS Wiki-Server führt:

[http://wiki.mcs.anl.gov/NEOS/index.php/NEOS\\_Wiki](http://wiki.mcs.anl.gov/NEOS/index.php/NEOS_Wiki), über den man u. a. zu der folgenden von Bob Fourer gepflegten Seite kommt:

<http://www.faqs.org/faqs/linear-programming-faq/>,

die eine wirklich gute Quelle mit vielfältigen und ausgezeichneten Informationen zur linearen Optimierung ist. Unter der Rubrik „Textbooks“ ist hier eine Reihe guter Bücher zum Thema zusammengestellt worden. Ich empfehle aus dieser Auswahl:

- Dantzig, George B., *Linear Programming and Extensions*, Princeton University Press, 1963. Dies ist der „Klassiker“ des Gebiets, auch heute noch eine interessante Quelle, und 1998 im Paperback-Format erneut erschienen.
- Chvátal, Vasek, *Linear Programming*, Freeman, 1983. Dies ist ein ausgezeichnetes Einführungsbuch, das sich insbesondere an Leser mit geringer mathematischer Vorbildung richtet.
- Padberg, Manfred, *Linear Optimization and Extensions*, Springer, 2001. Wenn sie wissen möchten, was die Berliner Luftbrücke mit linearer Optimierung zu tun hat, dann sollten Sie einmal in dieses Buch schauen.
- Schrijver, Alexander, *Theory of Linear and Integer Programming*, Wiley, 1986. Dieses Buch fasst den Kenntnisstand bis zum Jahre 1985 so gut wie vollständig zusammen, ein herausragendes Buch, das sich an den fortgeschrittenen Leser richtet.

- Vanderbei, Robert J., *Linear Programming: Foundations and Extensions*. Kluwer Academic Publishers, 1996. Eine schöne, Praxis bezogene Darstellung von Simplex- und Innere-Punkte-Methoden, Source-Codes zu den im Buch präsentierten Verfahren sind online verfügbar, siehe <http://www.princeton.edu/~rvdb/LPbook/>
- Wright, Stephen J., *Primal-Dual Interior-Point Methods*. SIAM Publications, 1997. Innere-Punkte-Methoden sind erst Mitte der 80er Jahre entdeckt und populär geworden. Wrights Buch behandelt die wichtigsten Verfahrensklassen dieses Ansatzes.

Leider wird die letztgenannte Webseite seit 2005 nicht mehr fortgeschrieben, so dass einige der Links inzwischen veraltet und einige neuere Referenzen nicht vorhanden sind.

Polyedertheorie ist ein wichtiger Aspekt der Vorlesung. Sehr gute Bücher hierzu sind:

- Grünbaum, Branko, *Convex Polytopes*, Springer-Verlag, Second Edition, 2003
- Ziegler, Günter M., *Lectures on Polytopes*, Springer-Verlag, Revised Edition, 1998

Ein Artikel, der die enormen Fortschritte bei der praktischen Lösung linearer Programme seit Anfang der 90er Jahre beschreibt, ist:

- Robert E. Bixby, *Solving Real-World Linear Programs: A Decade and More of Progress*, Operations Research 50 (2002) 3-15.

Fortschritte in den letzten Jahren im Bereich der ganzzahligen und gemischt-ganzzahligen Optimierung sind dokumentiert in

- Robert E. Bixby, Mary Fenelon, Zonghao Gu, Ed Rothberg, und Roland Wunderling, *Mixed-Integer Programming: A Progress Report*, erscheint 2004 in: Grötschel, M. (ed.), *The Sharpest Cut*, MPS/SIAM Series on Optimization.

## Lineare Optimierung im Internet: Vermischtes

Umfangreiche Literaturverweise, Preprint-Sammlungen, Hinweise auf Personen, die sich mit Optimierung und Spezialgebieten davon beschäftigen, verfügbare Computercodes, Beispiele zur Modellierung von praktischen Problemen, Testbeispiele mit Daten linearer Programme, etc, findet man u. a. auf den folgenden Internet-Seiten:

- <http://www.optimization-online.org/>  
(Optimization Online is a repository of e-prints about optimization and related topics.)
- <http://www-neos.mcs.anl.gov/>  
(Dies ist die Hauptseite des NEOS-Servers mit vielen Links zu Optimierungsseiten)
- <http://miplib.zib.de>  
(Testbeispiele für gemischt-ganzzahlige Optimierungsprobleme)
- <http://elib.zib.de/pub/Packages/mp-testdata/> (Verweise auf Sammlungen von Test-Beispielen wie die Netlib-LP-Testprobleme.)

Im Internet gibt es auch verschiedene Einführungen in die lineare Optimierung, insbesondere solche, die sich an Anfänger richten. Leider sind die URLs selten stabil, sodass sie hier nicht aufgelistet werden.

Wie oben erwähnt, wird in dieser Vorlesung besonders viel Wert auf die geometrische Fundierung der linearen Optimierung gelegt. Die Lösungsmengen von linearen Programmen sind Polyeder. Im Internet findet man viele Seiten, die Polyeder bildlich darstellen. Hier sind zwei Webseiten mit animierten Visualisierungen von Polyedern:

- <http://www.eg-models.de/> (gehe z. B. zu Classical Models, Regular Solids oder Discrete Mathematics, Polytopes)
- <http://mathworld.wolfram.com/topics/Polyhedra.html>

## Software

Häufig möchte man Polyeder von einer Darstellungsform (z.B. Lösungsmenge von Ungleichungssystemen) in eine andere Darstellungsform (z.B. konvexe Hülle

von endlich vielen Punkten) verwandeln, um gewisse Eigenschaften besser zu verstehen. Auch hierzu gibt es eine Reihe von Codes (die meisten stehen kostenlos zum Download bereit). Eine Sammlung solcher Verfahren ist zu finden unter:

- <http://elib.zib.de/pub/Packages/mathprog/polyth/index.html>

Die Programmsammlung PORTA stammt aus einer Diplomarbeit, die Thomas Christof bei mir geschrieben und dann mit anderen Kollegen (u. a. A. Löbel (ZIB)) weiterentwickelt hat. Das deutlich umfangreichere Programmpaket Polymake von E. Gawrilow und M. Joswig wurde im mathematischen Institut der TU Berlin entwickelt und wird hier weiter gepflegt und ausgebaut.

Es gibt eine Vielzahl von Codes zur Lösung linearer Programme und Modellierungssprachen, mit Hilfe derer man praktische Probleme (einigermaßen) bequem in Datensätze verwandeln kann, die LP-Löser akzeptieren. Eine kommentierte Übersicht hierzu ist zu finden unter:

- [http://wiki.mcs.anl.gov/NEOS/index.php/Optimization\\_Software\\_Guide](http://wiki.mcs.anl.gov/NEOS/index.php/Optimization_Software_Guide)

Eine weitere Liste findet man unter:

- <http://elib.zib.de/pub/Packages/mathprog/linprog/index.html>

Die Zeitschrift OR/MS Today veröffentlicht regelmäßig Übersichten über kommerzielle und frei verfügbare Software zur Lösung linearer Programme. Die letzte Änderung der Übersicht datiert vom 25.09.2009:

- <http://lionhrtpub.com/orms/surveys/LP/LP-survey.html>

In dieser Vorlesung haben wir Zugang zu dem kommerziellen Code CPLEX der Firma ILOG (2009 von IBM übernommen). Nach Ansicht vieler ist dies derzeit der weltweit beste Code zur Lösung linearer und ganzzahliger Programme (Ihnen steht allerdings nicht die allerneueste Version zur Verfügung).

Die seit kurzem auf dem Markt agierende Firma *Gurobi* bietet für Universitätsangehörige freien Zugang zu ihren Codes, siehe:

- [http://www.gurobi.com/html/download\\_academic.html](http://www.gurobi.com/html/download_academic.html).

Neueste Tests von Hans Mittelmann zeigen, dass der Gurobi Code zur weltweiten Spitzenklasse gehört, siehe:

- <http://plato.asu.edu/ftp/lpcom.html>.

Vom Konrad-Zuse-Zentrum (ZIB) wird ein für akademische Zwecke kostenlos verfügbarer Code namens SoPlex zur Lösung von LPs angeboten, siehe:

- <http://www.zib.de/Optimization/Software/Soplex>

Dieses Programm stammt aus der Dissertation von Roland Wunderling aus dem Jahre 1997 an der TU Berlin und wird am ZIB weiterentwickelt. Es ist sogar als Source Code verfügbar.

Als Modellierungssprache (z. B. zur Erzeugung von Daten im lp- oder mps-Format, die von den meisten LP-Lösern gelesen werden können) wird in der Vorlesung Zimpl angeboten. Dieses Programm wird von Thorsten Koch (ZIB) entwickelt und gepflegt, siehe:

- <http://www.zib.de/koch/zimpl/>

Zum Schluss noch ein Hinweis auf allgemeine Optimierungsprobleme. Hans Mittelmann unterhält eine Webseite, siehe

- <http://plato.la.asu.edu/guide.html>

die dabei hilft, für ein vorliegendes Optimierungsproblem geeignete Software ausfindig zu machen.

# Inhaltsverzeichnis

<b>1</b>	<b>Einführung</b>	<b>1</b>
1.1	Ein Beispiel . . . . .	1
1.2	Optimierungsprobleme . . . . .	11
1.3	Notation . . . . .	15
1.3.1	Grundmengen . . . . .	15
1.3.2	Vektoren und Matrizen . . . . .	16
1.3.3	Kombinationen von Vektoren, Hüllen, Unabhängigkeit . .	21
<b>2</b>	<b>Polyeder</b>	<b>23</b>
<b>3</b>	<b>Fourier-Motzkin-Elimination und Projektion</b>	<b>31</b>
3.1	Fourier-Motzkin-Elimination . . . . .	31
3.2	Lineare Optimierung und Fourier-Motzkin-Elimination . . . . .	39
<b>4</b>	<b>Das Farkas-Lemma und seine Konsequenzen</b>	<b>43</b>
4.1	Verschiedene Versionen des Farkas-Lemmas . . . . .	43
4.2	Alternativ- und Transpositionssätze . . . . .	45
4.3	Trennsätze . . . . .	47
<b>5</b>	<b>Der Dualitätssatz der linearen Programmierung</b>	<b>51</b>

<b>6</b>	<b>Grundlagen der Polyedertheorie</b>	<b>63</b>
6.1	Transformationen von Polyedern . . . . .	63
6.2	Kegelpolarität . . . . .	66
6.3	Darstellungssätze . . . . .	69
<b>7</b>	<b>Seitenflächen von Polyedern</b>	<b>73</b>
7.1	Die $\gamma$ -Polare und gültige Ungleichungen . . . . .	74
7.2	Seitenflächen . . . . .	76
7.3	Dimension . . . . .	81
7.4	Facetten und Redundanz . . . . .	83
<b>8</b>	<b>Ecken und Extremalen</b>	<b>91</b>
8.1	Rezessionskegel, Linienraum, Homogenisierung . . . . .	92
8.2	Charakterisierungen von Ecken . . . . .	96
8.3	Spitze Polyeder . . . . .	98
8.4	Extremalen von spitzen Polyedern . . . . .	101
8.5	Einige Darstellungssätze . . . . .	103
<b>9</b>	<b>Die Grundversion der Simplex-Methode</b>	<b>109</b>
9.1	Basen, Basislösungen, Entartung . . . . .	111
9.2	Basisaustausch (Pivoting), Simplexkriterium . . . . .	116
9.3	Das Simplexverfahren . . . . .	123
9.4	Die Phase I . . . . .	135
<b>10</b>	<b>Varianten der Simplex-Methode</b>	<b>141</b>
10.1	Der revidierte Simplexalgorithmus . . . . .	141
10.2	Spalten- und Zeilenauswahlregeln . . . . .	143
10.3	Die Behandlung oberer Schranken . . . . .	147

10.4	Das duale Simplexverfahren . . . . .	156
10.5	Zur Numerik des Simplexverfahrens . . . . .	163
10.6	Spezialliteratur zum Simplexalgorithmus . . . . .	171
<b>11</b>	<b>Postoptimierung und parametrische Programme</b>	<b>175</b>
11.1	Änderung der rechten Seite $b$ . . . . .	176
11.2	Änderungen der Zielfunktion $c$ . . . . .	178
11.3	Änderungen des Spaltenvektors $A_{\cdot,j}, j = N(s)$ . . . . .	179
11.4	Hinzufügung einer neuen Variablen . . . . .	181
11.5	Hinzufügung einer neuen Restriktion . . . . .	182
<b>12</b>	<b>Die Ellipsoidmethode</b>	<b>185</b>
12.1	Polynomiale Algorithmen . . . . .	186
12.2	Reduktionen . . . . .	188
12.3	Beschreibung der Ellipsoidmethode . . . . .	194
12.4	Laufzeit der Ellipsoidmethode . . . . .	203
12.5	Ein Beispiel . . . . .	205
<b>13</b>	<b>Der Karmarkar-Algorithmus und Innere-Punkte-Methoden</b>	<b>211</b>
13.1	13.1 Reduktionen . . . . .	212
13.2	Die Grundversion des Karmarkar-Algorithmus . . . . .	214
<b>14</b>	<b>Duale Heuristiken, Relaxierungen</b>	<b>227</b>
14.1	Zwei Relaxierungstechniken . . . . .	229
14.2	Lagrange-Relaxierungen . . . . .	235
<b>15</b>	<b>Branch &amp; Bound-Verfahren</b>	<b>249</b>
<b>16</b>	<b>Theorie der ganzzahligen und</b>	

<b>gemischt-ganzzahligen Optimierung</b>	<b>265</b>
16.1 Einführung . . . . .	265
16.2 Ganzzahlige Punkte in rationalen Polyedern . . . . .	271
16.3 Schnittebentheorie . . . . .	276
<b>17 Allgemeine Schnittebenenverfahren der ganzzahligen Programmierung</b>	<b>283</b>
17.1 Ein Schnittebenenverfahren für ganzzahlige Programme . . . . .	284
17.2 Ein Schnittebenenverfahren für gemischt-ganzzahlige Programme	307

# Kapitel 1

## Einführung

Dieses erste Kapitel dient der Einführung in das Gebiet “Optimierung”. Optimierungsaufgaben werden zunächst ganz allgemein definiert und dann immer mehr spezialisiert, um zu Aufgaben zu kommen, über die mathematisch interessante Aussagen gemacht werden können. Anwendungen der mathematischen Modelle werden nur andeutungsweise behandelt. Wir beschäftigen uns zunächst hauptsächlich mit der Theorie, denn ohne theoretische Vorkenntnisse kann man keine Anwendungen betreiben.

Eine für die Praxis wichtige Tätigkeit ist die Modellbildung oder -formulierung, also die Aufgabe ein reales Problem der Praxis in eine mathematische Form zu “kleiden”. Diese Aufgabe wird von Mathematikern häufig ignoriert oder als trivial betrachtet. Aber es ist gar nicht so einfach, für ein wichtiges Praxisproblem ein “gutes” mathematisches Modell zu finden. Dies muss anhand vieler Beispiele geübt werden, um „Modellbildungserfahrung“ zu gewinnen. Wir werden dies vornehmlich in den Übungen zur Vorlesung abhandeln.

### 1.1 Ein Beispiel

Was ist ein Optimierungsproblem? Bevor wir diese Frage durch eine allgemeine Definition beantworten, wollen wir ein Beispiel ausführlich besprechen, um zu zeigen, wie Optimierungsprobleme entstehen, und wie man sie mathematisch behandeln kann. Wir beginnen mit einem (natürlich für Vorlesungszwecke aufbereiteten) Beispiel aus der Ölindustrie.

**(1.1) Beispiel (Ölraffinierung).** In Ölraffinerien wird angeliefertes Rohöl durch

Anwendung von chemischen und (oder) physikalischen Verfahren in gewisse gewünschte Komponenten zerlegt. Die Ausbeute an verschiedenen Komponenten hängt von dem eingesetzten Verfahren (Crackprozess) ab. Wir nehmen an, dass eine Raffinerie aus Rohöl drei Komponenten (schweres Öl  $S$ , mittelschweres Öl  $M$ , leichtes Öl  $L$ ) herstellen will. Sie hat zwei Crackverfahren zur Verfügung, die die folgende Ausbeute liefern und die die angegebenen Kosten (Energie, Maschinenabschreibung, Arbeit) verursachen. (Zur Verkürzung der Schreibweise kürzen wir das Wort "Mengeinheit" durch ME und das Wort "Geldeinheit" durch GE ab. 10 ME Rohöl ergeben in:

Crackprozess 1 :	2 ME	$S$	
	2 ME	$M$	Kosten: 3 GE
	1 ME	$L$	

Crackprozeß 2 :	1 ME	$S$	
	2 ME	$M$	Kosten: 5 GE
	4 ME	$L$	

Aufgrund von Lieferverpflichtungen muss die Raffinerie folgende Mindestproduktion herstellen:

3 ME	$S$
5 ME	$M$
4 ME	$L$

Diese Mengen sollen so kostengünstig wie möglich produziert werden.

□

Wir wollen nun die obige Aufgabe (1.1) mathematisch formulieren. Unser Ergebnis wird eine mathematische Aufgabe sein, die wir später als lineares Programm oder lineares Optimierungsproblem bezeichnen werden.

Zunächst stellen wir fest, dass die Crackprozesse unabhängig voneinander arbeiten und (im Prinzip) mit jeder beliebigen Rohölmenge beschickt werden können. Wir führen daher zwei Variable  $x_1$  und  $x_2$  ein, die das Produktionsniveau der beiden Prozesse beschreiben. Wird z. B.  $x_1 = 1.5$  gewählt, so bedeutet dies, dass der Crackprozess 1 mit 15 ME Rohöl beschickt wird, während bei  $x_2 = 0.5$  der Crackprozess 2 nur mit 5 ME Rohöl gefahren wird. Natürlich macht es keinen Sinn, einen Crackprozess mit negativen Rohölmengen zu beschicken, aber jeder Vektor  $x = (x_1, x_2) \in \mathbb{R}^2$ ,  $x \geq 0$ , beschreibt ein mögliches Produktionsniveau der beiden Prozesse.

Angenommen durch  $(x_1, x_2) \geq 0$  sei ein Produktionsniveau beschrieben, dann folgt aus den technologischen Bedingungen der beiden Crackprozesse, dass sich

der Ausstoß an schwerem Öl  $S$  auf

$$2x_1 + x_2 \quad \text{ME} \quad S$$

beläuft, und analog erhalten wir

$$\begin{array}{rcl} 2x_1 & +2x_2 & \text{ME} \quad M \\ x_1 & +4x_2 & \text{ME} \quad L. \end{array}$$

Die Kosten betragen offenbar

$$z = 3x_1 + 5x_2 \quad \text{GE}.$$

Die eingegangenen Lieferverpflichtungen besagen, dass gewisse Mindestmengen an  $S$ ,  $M$  und  $L$  produziert werden müssen, das heißt, ein Vektor  $(x_1, x_2)$  beschreibt nur dann ein Produktionsniveau, bei dem auch alle Lieferverpflichtungen erfüllt werden können, wenn gilt:

$$(1.2) \quad \begin{array}{rcl} 2x_1 + x_2 & \geq & 3 \\ 2x_1 + 2x_2 & \geq & 5 \\ x_1 + 4x_2 & \geq & 4 \\ x_1 & \geq & 0 \\ x_2 & \geq & 0 \end{array}$$

Unsere Aufgabe besteht nun darin, unter allen Vektoren  $(x_1, x_2) \in \mathbb{R}^2$ , die die 5 Ungleichungen in (1.2) erfüllen, einen Vektor zu finden, der minimale Kosten verursacht, d. h. bei dem

$$z = 3x_1 + 5x_2$$

einen möglichst kleinen Wert annimmt.

Es hat sich eingebürgert, eine derartige Aufgabe ein **lineares Programm** oder ein **lineares Optimierungsproblem** zu nennen und es wie folgt zu schreiben:

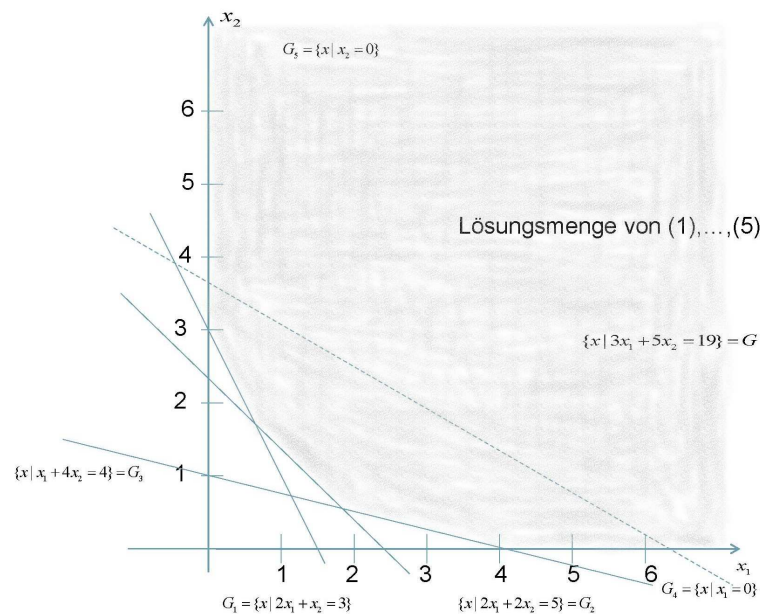
$$\min 3x_1 + 5x_2$$

unter den Nebenbedingungen:

$$(1.3) \quad \begin{array}{rcl} (1) & 2x_1 + x_2 \geq 3 & (S\text{-Ungleichung}) \\ (2) & 2x_1 + 2x_2 \geq 5 & (M\text{-Ungleichung}) \\ (3) & x_1 + 4x_2 \geq 4 & (L\text{-Ungleichung}) \\ (4) & x_1 \geq 0 & \\ (5) & x_2 \geq 0 & \end{array}$$

Dabei nennt man die zu minimierende Funktion  $3x_1 + 5x_2$  die **Zielfunktion** des Problems. Jeder Vektor, der (1), ..., (5) erfüllt, heißt **zulässige Lösung** des Problems.

Da wir nur 2 Variablen haben, können wir unsere Aufgabe graphisch analysieren. Ersetzen wir das “ $\geq$ ”-Zeichen in (1), ..., (5) durch ein Gleichheitszeichen, so erhalten wir 5 Gleichungen. Die Lösungsmenge einer Gleichung im  $\mathbb{R}^2$  ist bekanntlich eine Gerade. Diese 5 Geraden “beranden” die Lösungsmenge des Systems der Ungleichungen (1), ..., (5). Diese Lösungsmenge ist in Abbildung 1.1 graphisch dargestellt. (Der Leser möge sich unbedingt mit der geometrischen Interpretation von Ungleichungen vertraut machen.)



**Abb. 1.1**

Die Zielfunktion ist natürlich keine Gerade, sie repräsentiert eine Schar paralleler Geraden. Nehmen wir z. B. den Vektor  $x = (3, 2)$ , der alle Ungleichungen (1), ..., (5) erfüllt. Sein Zielfunktionswert ist 19, d. h. bei diesem Produktionsniveau treten Gesamtkosten in Höhe von 19 GE auf. In Abbildung 1.1 ist die Gerade  $G = \{x \mid 3x_1 + 5x_2 = 19\}$  gestrichelt gezeichnet. Die Menge aller derjenigen Punkte, die auf dieser Geraden liegen und (1), ..., (5) erfüllen, stellen Produktionsniveaus mit Gesamtkosten 19 GE dar.

Geometrisch ist nun klar, wie wir einen Punkt finden können, der (1), ..., (5) erfüllt und die Zielfunktion minimiert. Wir verschieben die Gerade  $G$  so lange

parallel in Richtung auf den Punkt  $(0, 0)$ , bis die verschobene Gerade die Lösungsmenge von (1), ..., (5) nur noch tangiert. Führen wir dies graphisch durch, so sehen wir, dass wir die Tangentialstellung im Punkt  $x^* = (2, \frac{1}{2})$  erreichen. Die zu  $G$  parallele Gerade

$$G' = \{x \mid 3x_1 + 5x_2 = 8.5\}$$

berührt die Lösungsmenge von (1), ..., (5) in nur einem Punkt, nämlich  $x^*$ , jede weitere Parallelverschiebung würde zu einem leeren Durchschnitt mit dieser Lösungsmenge führen. Wir schließen daraus, dass

$$x^* = (2, \frac{1}{2})$$

die Optimallösung unseres Problems ist, d. h. alle Lieferverpflichtungen können bei diesem Produktionsniveau erfüllt werden, und alle anderen Produktionsniveaus führen zu Gesamtkosten, die höher sind als die Kosten von 8.5 GE, die beim Produktionsniveau  $x^*$  anfallen.

Die vorgeführte Lösung des Beispiels (1.1) ist anschaulich und sieht recht einfach aus. Aber es ist klar, dass sich diese graphische Argumentation nicht für mehr als zwei Variable durchführen lässt. Die graphische Methode ist zwar nützlich, um die geometrische Intuition anzuregen, aber zur Lösung von linearen Programmen taugt sie nicht. Und dann stellt sich natürlich die Frage, ob unsere Schlussfolgerungen als ein korrekter Beweis für die Optimalität von  $x^*$  angesehen werden können.

Wir wollen daher einen zweiten Versuch zur Lösung von (1.1) unternehmen und betrachten das Problem etwas allgemeiner. Wählen wir statt der speziellen Zielfunktion  $3x_1 + 5x_2$  eine beliebige, sagen wir

$$ax_1 + bx_2,$$

so sehen wir, dass unser Problem gar keine Optimallösung haben muss. Ist nämlich einer der Parameter negativ, sagen wir  $a < 0$ , so können wir den Zielfunktionswert so klein machen, wie wir wollen. Z. B. gilt für die Folge der Punkte  $x^{(n)} := (n, 0)$ ,  $n \in \mathbb{N}$ ,

$$x^{(n)} \text{ erfüllt } (1), \dots, (5) \text{ für alle } n \geq 4, \\ ax_1^{(n)} + bx_2^{(n)} = na.$$

Also ist der Wert der Zielfunktion über der Lösungsmenge  $(1), \dots, (5)$  nicht nach unten beschränkt. Von unserem Beispiel ausgehend ist eine derartige Zielfunktion natürlich unsinnig, aber mathematisch müssen derartige Fälle behandelt und klassifiziert werden. Kommt so etwas in der Praxis vor? Natürlich, und zwar öfter als man denkt! Gründe hierfür liegen in der Regel bei fehlerhafter Dateneingabe, falschem Zugriff auf Datenbanken oder irgendwelchen Übertragungsfehlern.

Auch Lösungsmengen, die aufgrund der eigentlichen Daten nicht leer sein sollten, können auf diese Weise plötzlich leer sein.

Sind  $a$  und  $b$  nichtnegativ, so ist klar, dass es eine Optimallösung gibt. Führen wir unser Parallelverschiebungsexperiment mehrmals durch, so werden wir feststellen, dass wir immer eine Optimallösung finden können, die den Durchschnitt von zwei der fünf Geraden  $G_1, \dots, G_5$  bildet. Man könnte also vermuten, dass dies immer so ist, d. h. dass immer Optimallösungen gefunden werden können, die (bei  $n$  Variablen) den Durchschnitt von  $n$  Hyperebenen bilden.

Wir werden später sehen, dass dies (bis auf Pathologien) in der Tat der Fall ist, und dass dies eine der fruchtbarsten Beobachtungen in Bezug auf die Lösungstheorie von Problemen des Typs (1.1) ist.

Die Lösungsmenge von (1.3) (1), ..., (5) enthält unendlich viele Punkte. Unsere obige Überlegung hat die Suche auf endlich viele Punkte reduziert, und diese Punkte können wir sogar effektiv angeben. Wir schreiben die Ungleichungen (1), ..., (5) als Gleichungen und betrachten alle möglichen 10 Systeme von 2 Gleichungen (in 2 Variablen) dieses Systems. Jedes dieser 10 Gleichungssysteme können wir mit Hilfe der Gauß-Elimination lösen. Z. B. ist der Durchschnitt von  $G_1$  und  $G_3$  der Punkt  $y_1 = (\frac{8}{7}, \frac{5}{7})$ , und der Durchschnitt von  $G_1$  und  $G_2$  der Punkt  $y_2 = (\frac{1}{2}, 2)$ . Der Punkt  $y_1$  erfüllt jedoch die Ungleichung (2) nicht. Wir müssen also zunächst alle Durchschnitte ausrechnen, dann überprüfen, ob der jeweils errechnete Punkt tatsächlich alle Ungleichungen erfüllt, und unter den so gefundenen Punkten den besten auswählen.

Dieses Verfahren ist zwar endlich, aber bei  $m$  Ungleichungen und  $n$  Variablen erfordert es  $\binom{m}{n}$  Aufrufe der Gauß-Elimination und weitere Überprüfungen. Es ist also praktisch nicht in vernünftiger Zeit ausführbar. Dennoch steckt hierin der Kern des heutzutage wichtigsten Verfahrens zur Lösung von linearen Programmen, des Simplexverfahrens. Man berechnet zunächst eine Lösung, die im Durchschnitt von  $n$  Hyperebenen liegt, und dann versucht man "systematisch" weitere derartige Punkte zu produzieren und zwar dadurch, dass man jeweils nur eine Hyperebene durch eine andere austauscht und so von einem Punkt auf "zielstrebige" Weise über "Nachbarpunkte" zum Optimalpunkt gelangt. Überlegen Sie sich einmal selbst, wie Sie diese Idee verwirklichen würden!

Eine wichtige Komponente des Simplexverfahrens beruht auf der Tatsache, dass man effizient erkennen kann, ob ein gegebener Punkt tatsächlich der gesuchte Optimalpunkt ist oder nicht. Dahinter steckt die sogenannte *Dualitätstheorie*, die wir anhand unseres Beispiels (1.1) erläutern wollen. Wir wollen also zeigen, dass der Punkt  $x^* = (2, \frac{1}{2})$  minimal ist.

Ein Weg, dies zu beweisen, besteht darin, untere Schranken für das Minimum zu finden. Wenn man sogar in der Lage ist, eine untere Schranke zu bestimmen, die mit dem Zielfunktionswert einer Lösung des Ungleichungssystems übereinstimmt, ist man fertig (denn der Zielfunktionswert irgendeiner Lösung ist immer eine obere Schranke, und wenn untere und obere Schranke gleich sind, ist der optimale Wert gefunden).

Betrachten wir also unser lineares Program (1.3). Wir bemerken zunächst, dass man alle Ungleichungen (1),  $\dots$ , (5) skalieren kann, d. h. wir können z. B. die Ungleichung

$$(1) \quad 2x_1 + x_2 \geq 3$$

mit 2 multiplizieren und erhalten:

$$(1') \quad 4x_1 + 2x_2 \geq 6.$$

Die Menge der Punkte des  $\mathbb{R}^2$ , die (1) erfüllen, stimmt offenbar überein mit der Menge der Punkte, die (1') erfüllen. Wir können auch mit negativen Faktoren skalieren, aber Achtung, dann dreht sich das Ungleichungszeichen um!

Eine weitere simple Beobachtung ist die folgende. Erfüllt ein Vektor zwei Ungleichungen, so erfüllt er auch die Summe der beiden Ungleichungen. Das gilt natürlich nur, wenn die Ungleichungszeichen gleichgerichtet sind. Addieren wir z. B. zur Ungleichung (1) die Ungleichung (3), so erhalten wir

$$3x_1 + 5x_2 \geq 7.$$

Die linke Seite dieser Ungleichung ist (natürlich nicht rein zufällig) gerade unsere Zielfunktion. Folglich können wir aus der einfachen Überlegung, die wir gerade gemacht haben, schließen, dass jeder Vektor, der (1.3) (1),  $\dots$ , (5) erfüllt, einen Zielfunktionswert hat, der mindestens 7 ist. Der Wert 7 ist somit eine untere Schranke für das Minimum in (1.3). Eine tolle Idee, oder?

Das Prinzip ist damit dargestellt. Wir skalieren die Ungleichungen unseres Systems und addieren sie, um untere Schranken zu erhalten. Da unsere Ungleichungen alle in der Form " $\geq$ " geschrieben sind, dürfen wir nur positiv skalieren, denn zwei Ungleichungen  $a_1x_1 + a_2x_2 \leq \alpha$  und  $b_1x_1 + b_2x_2 \geq \beta$  kann man nicht addieren. Gibt jede beliebige Skalierung und Addition eine untere Schranke? Machen wir noch zwei Versuche. Wir addieren (1) und (2) und erhalten

$$4x_1 + 3x_2 \geq 8.$$

Das hilft uns nicht weiter, denn die linke Seite der Ungleichung kann nicht zum Abschätzen der Zielfunktion benutzt werden. Betrachten wir nun, das 1.5-fache

der Ungleichung (2), so ergibt sich

$$3x_1 + 3x_2 \geq 7.5.$$

Hieraus schließen wir: Ist  $3x_1 + 3x_2 \geq 7.5$ , dann ist erst recht  $3x_1 + 5x_2 \geq 7.5$ , denn nach (5) gilt ja  $x_2 \geq 0$ .

Daraus können wir eine Regel für die Bestimmung unterer Schranken ableiten. Haben wir eine neue Ungleichung als Summe von positiv skalierten Ausgangsungleichungen erhalten und hat die neue Ungleichung die Eigenschaft, dass jeder ihrer Koeffizienten höchstens so groß ist wie der entsprechende Koeffizient der Zielfunktion, so liefert die rechte Seite der neuen Ungleichung eine untere Schranke für den Minimalwert von (1.3).

Diese Erkenntnis liefert uns ein neues mathematisches Problem. Wir suchen nicht-negative Multiplikatoren (Skalierungsfaktoren) der Ungleichungen (1), (2), (3) mit gewissen Eigenschaften. Multiplizieren wir (1) mit  $y_1$ , (2) mit  $y_2$  und (3) mit  $y_3$ , so darf die Summe  $2y_1 + 2y_2 + y_3$  nicht den Wert des ersten Koeffizienten der Zielfunktion (also 3) überschreiten. Analog darf die Summe  $y_1 + 2y_2 + 4y_3$  nicht den Wert 5 überschreiten. Und die  $y_i$  sollen nicht negativ sein. Ferner soll die rechte Seite der Summenungleichung so groß wie möglich werden. Die rechte Seite kann man durch  $3y_1 + 5y_2 + 4y_3$  berechnen. Daraus folgt, dass wir die folgende Aufgabe lösen müssen. Bestimme

$$\max 3y_1 + 5y_2 + 4y_3$$

unter den Nebenbedingungen:

$$(1.4) \quad \begin{aligned} 2y_1 + 2y_2 + y_3 &\leq 3 \\ y_1 + 2y_2 + 4y_3 &\leq 5 \\ y_1, y_2, y_3 &\geq 0 \end{aligned}$$

Auch (1.4) ist ein lineares Programm. Aus unseren Überlegungen folgt, dass der Maximalwert von (1.4) höchstens so groß ist wie der Minimalwert von (1.3), da jede Lösung von (1.4) eine untere Schranke für (1.3) liefert. Betrachten wir z. B. den Punkt

$$y^* = (y_1^*, y_2^*, y_3^*) = \left(0, \frac{7}{6}, \frac{2}{3}\right).$$

Durch Einsetzen in die Ungleichungen von (1.4) sieht man dass  $y^*$  alle Ungleichungen erfüllt. Addieren wir also zum  $\frac{7}{6}$ -fachen der Ungleichung (2) das  $\frac{2}{3}$ -fache der Ungleichung (3), so erhalten wir

$$3x_1 + 5x_2 \geq 8.5.$$

Damit wissen wir, dass der Minimalwert von (1.3) mindestens 8.5 ist. Der Punkt  $x^*$  liefert gerade diesen Wert, er ist also optimal. Und außerdem ist  $y^*$  eine Optimallösung von (1.4).

Die hier dargestellten Ideen bilden die Grundgedanken der Dualitätstheorie der linearen Programmierung, und sie sind außerordentlich nützlich bei der Entwicklung von Algorithmen und in Beweisen. In der Tat haben wir bereits ein kleines Resultat erzielt, das wir wie folgt formal zusammenfassen wollen.

**(1.5) Satz.** *Es seien  $c \in \mathbb{R}^n$ ,  $b \in \mathbb{R}^m$ , und  $A$  sei eine reelle  $(m, n)$ -Matrix. Betrachten wir die Aufgaben*

$$(P) \quad \min \quad c^T x \quad \text{unter} \quad Ax \geq b, \quad x \geq 0,$$

(lies: bestimme einen Vektor  $x^*$ , der  $Ax^* \geq b$  und  $x^* \geq 0$  erfüllt und dessen Wert  $c^T x^*$  so klein wie möglich ist),

$$(D) \quad \max \quad y^T b \quad \text{unter} \quad y^T A \leq c^T, \quad y \geq 0$$

(lies: bestimme einen Vektor  $y^*$ , der  $y^{*T} A \leq c^T$  und  $y^* \geq 0$  erfüllt und dessen Wert  $y^{*T} b$  so groß wie möglich ist),

dann gilt Folgendes:

Seien  $x_0 \in \mathbb{R}^n$  und  $y_0 \in \mathbb{R}^m$  Punkte mit  $Ax_0 \geq b$ ,  $x_0 \geq 0$  und  $y_0^T A \leq c^T$ ,  $y_0 \geq 0$ , dann gilt

$$y_0^T b \leq c^T x_0.$$

**Beweis :** Durch einfaches Einsetzen:

$$y_0^T b \leq y_0^T (Ax_0) = (y_0^T A)x_0 \leq c^T x_0.$$

□

Satz (1.5) wird **schwacher Dualitätssatz** genannt. Für Optimallösungen  $x^*$  und  $y^*$  von (P) bzw. (D), gilt nach (1.5)  $y^{*T} b \leq c^T x^*$ . Wir werden später zeigen, dass in diesem Falle sogar immer  $y^{*T} b = c^T x^*$  gilt. Dies ist allerdings nicht ganz so einfach zu beweisen wie der obige Satz (1.5).

Unser Beispiel (1.1) stammt aus einer ökonomischen Anwendung. Das lineare Programm (1.3) haben wir daraus durch mathematische Formalisierung abgeleitet und daraus durch mathematische Überlegungen ein neues lineares Programm (1.4) gewonnen. Aber auch dieses Programm hat einen ökonomischen Sinn.

**(1.6) Ökonomische Interpretation von (1.4).** Als Manager eines Ölkonzerns sollte man darüber nachdenken, ob es überhaupt sinnvoll ist, die drei Ölsorten  $S$ ,  $M$ ,  $L$  selbst herzustellen. Vielleicht ist es gar besser, die benötigten Ölmengen auf dem Markt zu kaufen. Die Manager müssen sich also überlegen, welche Preise sie auf dem Markt für  $S$ ,  $M$ ,  $L$  gerade noch bezahlen würden, ohne die Produktion selbst aufzunehmen. Oder anders ausgedrückt, bei welchen Marktpreisen für  $S$ ,  $M$ ,  $L$  lohnt sich die Eigenproduktion gerade noch.

Bezeichnen wir die Preise für eine Mengeneinheit von  $S$ ,  $M$ ,  $L$  mit  $y_1$ ,  $y_2$ ,  $y_3$ , so sind beim Ankauf der benötigten Mengen für  $S$ ,  $M$ ,  $L$  genau  $3y_1 + 5y_2 + 4y_3$  GE zu bezahlen. Der erste Crackprozess liefert bei 10 ME Rohöleinsatz 2 ME  $S$ , 2 ME  $M$  und 1 ME  $L$  und verursacht 3 GE Kosten. Würden wir den Ausstoss des ersten Crackprozesses auf dem Markt kaufen, so müssten wir dafür  $2y_1 + 2y_2 + y_3$  GE bezahlen. Ist also  $2y_1 + 2y_2 + y_3 > 3$ , so ist es auf jeden Fall besser, den Crackprozess 1 durchzuführen als auf dem Markt zu kaufen. Analog gilt für den zweiten Prozess: Erfüllen die Marktpreise die Bedingung  $y_1 + 2y_2 + 4y_3 > 5$ , so ist die Durchführung des Crackprozesses 2 profitabel.

Nur solche Preise  $y_1$ ,  $y_2$ ,  $y_3$ , die die Bedingungen

$$\begin{aligned} 2y_1 + 2y_2 + y_3 &\leq 3 \\ y_1 + 2y_2 + 4y_3 &\leq 5 \\ y_1, y_2, y_3 &\geq 0 \end{aligned}$$

erfüllen, können also das Management veranlassen, auf dem Markt zu kaufen, statt selbst zu produzieren. Der beim Kauf der durch Lieferverträge benötigten Mengen zu bezahlende Preis beläuft sich auf

$$3y_1 + 5y_2 + 4y_3.$$

Wenn wir diese Funktion unter den obigen Nebenbedingungen maximieren, erhalten wir Preise (in unserem Falle  $y_1^* = 0$ ,  $y_2^* = \frac{7}{6}$ ,  $y_3^* = \frac{2}{3}$ ), die die Eigenproduktion gerade noch profitabel erscheinen lassen. Diese Preise werden **Schattenpreise** genannt. Sie sagen nichts über den tatsächlichen Marktpreis aus, sondern geben in gewisser Weise den Wert des Produktes für den Hersteller an.

Das Resultat  $y_1^* = 0$  besagt z. B., dass das Schweröl  $S$  für den Hersteller wertlos ist. Warum? Das liegt daran, dass bei der optimalen Kombination der Crackprozesse mehr Schweröl anfällt (nämlich 4.5 ME) als benötigt wird (3 ME). Das heißt, bei jeder Lieferverpflichtung für Schweröl, die zwischen 0 und 4.5 ME liegt, würde sich die optimale Prozesskombination nicht ändern, und die gesamten Kosten wären in jedem Falle die gleichen. Ein gutes Management wird in einer solchen Situation die Verkaufsmitarbeiter dazu anregen, Schweröl zu verkaufen.

Bei jedem beliebigen Verkaufspreise winkt (innerhalb eines begrenzten Lieferumfangs) Gewinn! Und es gibt weniger Entsorgungsprobleme. Müsste der Hersteller dagegen eine ME mehr an leichtem Öl  $L$  liefern, so müsste die Kombination der Crackprozesse geändert werden, und zwar würden dadurch insgesamt Mehrkosten in Höhe von  $\frac{2}{3}$  GE anfallen. Analog würde bei einer Erhöhung der Produktion von 5 ME mittleren Öls auf 6 ME eine Kostensteigerung um  $\frac{7}{6}$  GE erfolgen.

Diesen Sachverhalt können wir dann auch so interpretieren. Angenommen, die Ölfirma hat die Prozessniveaus  $x_1^* = 2$ ,  $x_2^* = 0.5$  gewählt. Ein Mineralölhändler möchte 1 ME leichtem Öl  $L$  zusätzlich kaufen. Ist der Marktpreis mindestens  $\frac{2}{3}$  GE pro ME  $L$  wird die Firma das Geschäft machen, andernfalls wäre es besser, die zusätzliche Einheit nicht selbst zu produzieren sondern von einem anderen Hersteller zu kaufen. Würde der Mineralölhändler 1 ME schweres Öl nachfragen, würde die Firma bei jedem Preis verkaufen, denn sie hat davon zu viel.

□

## 1.2 Optimierungsprobleme

Wir wollen zunächst ganz informell, ohne auf technische Spitzfindigkeiten einzugehen, Optimierungsprobleme einführen. Sehr viele Probleme lassen sich wie folgt formulieren.

Gegeben seien eine Menge  $S$  und eine geordnete Menge  $(T, \leq)$ , d. h. zwischen je zwei Elementen  $s, t \in T$  gilt genau eine der folgenden Beziehungen  $s < t$ ,  $s > t$  oder  $s = t$ . Ferner sei eine Abbildung  $f : S \rightarrow T$  gegeben. Gesucht ist ein Element  $x^* \in S$  mit der Eigenschaft  $f(x^*) \geq f(x)$  für alle  $x \in S$  (Maximierungsproblem) oder  $f(x^*) \leq f(x)$  für alle  $x \in S$  (Minimierungsproblem). Es ist üblich hierfür eine der folgenden Schreibweisen zu benutzen:

$$(1.7) \quad \begin{array}{ll} \max_{x \in S} f(x) & \text{oder} \quad \max\{f(x) \mid x \in S\}, \\ \min_{x \in S} f(x) & \text{oder} \quad \min\{f(x) \mid x \in S\}. \end{array}$$

In der Praxis treten als geordnete Mengen  $(T, \leq)$  meistens die reellen Zahlen  $\mathbb{R}$ , die rationalen Zahlen  $\mathbb{Q}$  oder die ganzen Zahlen  $\mathbb{Z}$  auf, alle mit der natürlichen Ordnung versehen (die wir deswegen auch gar nicht erst notieren). Die Aufgabe (1.7) ist viel zu allgemein, um darüber etwas Interessantes sagen zu können. Wenn  $S$  durch die Auflistung aller Elemente gegeben ist, ist das Problem entweder sinnlos oder trivial (man rechnet ganz einfach  $f(x)$  für alle  $x \in S$  aus). Das heißt,  $S$  muss irgendwie (explizit oder implizit) strukturiert sein, so dass vernünftige Aussagen über  $S$  möglich sind, ohne dass man alle Elemente in  $S$  einzeln kennt. Das

gleiche gilt für die Funktion  $f : S \rightarrow T$ . Ist sie nur punktweise durch  $x \mapsto f(x)$  gegeben, lohnt sich das Studium von (1.7) nicht. Erst wenn  $f$  durch hinreichend strukturierte “Formeln” bzw. “Eigenschaften” bestimmt ist, werden tieferliegende mathematische Einsichten möglich.

Die Optimierungsprobleme, die in der Praxis auftreten, haben fast alle irgendeine “vernünftige” Struktur. Das muss nicht unbedingt heißen, dass die Probleme dadurch auf einfache Weise lösbar sind, aber immerhin ist es meistens möglich, sie in das zur Zeit bekannte und untersuchte Universum der verschiedenen Typen von Optimierungsproblemen einzureihen und zu klassifizieren.

Im Laufe des Studiums werden Ihnen noch sehr unterschiedliche Optimierungsaufgaben begegnen. Viele werden von einem der folgenden Typen sein.

**(1.8) Ein Kontrollproblem.** Gegeben sei ein Steuerungsprozess (z. B. die Bewegungsgleichung eines Autos), etwa der Form

$$\dot{x}(t) = f(t, x(t), u(t)),$$

wobei  $u$  eine Steuerung ist (Benzinzufuhr). Ferner seien eine Anfangsbedingung

$$x(0) = x_0,$$

(z. B.: das Auto steht) sowie eine Endbedingung

$$x(T) = x_1$$

(z.B. das Auto hat eine Geschwindigkeit von 50 km/h) gegeben. Gesucht ist eine Steuerung  $u$  für den Zeitraum  $[0, T]$ , so dass z. B.

$$\int_0^T |u|^2 dt$$

minimal ist (etwa minimaler Benzinverbrauch).

□

**(1.9) Ein Approximationsproblem.** Gegeben sei eine (numerisch schwierig auszuwertende) Funktion  $f$ , finde eine Polynom  $p$  vom Grad  $n$ , so dass

$$\|f - p\| \quad \text{oder} \quad \|f - p\|_\infty$$

minimal ist.

□

**(1.10) Nichtlineares Optimierungsproblem.** Es seien  $f, g_i$  ( $i = 1, \dots, m$ ),  $h_j$  ( $j = 1, \dots, p$ ) differenzierbare Funktionen von  $\mathbb{R}^n \rightarrow \mathbb{R}$ , dann heißt

$$\begin{aligned} \min f(x) \\ g_i(x) \leq 0 \quad i = 1, \dots, m \\ h_j(x) = 0 \quad j = 1, \dots, p \\ x \in \mathbb{R}^n \end{aligned}$$

ein nichtlineares Optimierungsproblem. Ist eine der Funktionen nicht differenzierbar, so spricht man von einem nichtdifferenzierbaren Optimierungsproblem. Im Allgemeinen wird davon ausgegangen, daß alle betrachteten Funktionen zumindest stetig sind.

□

**(1.11) Konvexes Optimierungsproblem.** Eine Menge  $S \subseteq \mathbb{R}^n$  heißt **konvex**, falls gilt: Sind  $x, y \in S$  und ist  $\lambda \in \mathbb{R}, 0 \leq \lambda \leq 1$ , dann gilt  $\lambda x + (1 - \lambda)y \in S$ . Eine Funktion  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  heißt **konvex**, falls für alle  $\lambda \in \mathbb{R}, 0 \leq \lambda \leq 1$  und alle  $x, y \in \mathbb{R}^n$  gilt

$$\lambda f(x) + (1 - \lambda)f(y) \geq f(\lambda x + (1 - \lambda)y).$$

Ist  $S \subseteq \mathbb{R}^n$  konvex (z. B. kann  $S$  wie folgt gegeben sein  $S = \{x \in \mathbb{R}^n \mid g_i(x) \leq 0, i = 1, \dots, m\}$  wobei die  $g_i$  konvexe Funktionen sind), und ist  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  eine konvexe Funktion, dann ist

$$\min_{x \in S} f(x)$$

ein konvexes Minimierungsproblem.

□

**(1.12) Lineares Optimierungsproblem (Lineares Programm).**

Gegeben seien  $c \in \mathbb{R}^n, A \in \mathbb{R}^{(m,n)}, b \in \mathbb{R}^m$ , dann heißt

$$\begin{aligned} \max c^T x \\ Ax \leq b \\ x \in \mathbb{R}^n \end{aligned}$$

lineares Optimierungsproblem.

□

**(1.13) Lineares ganzzahliges Optimierungsproblem.**

Gegeben seien  $c \in \mathbb{R}^n$ ,  $A \in \mathbb{R}^{(m,n)}$ ,  $b \in \mathbb{R}^m$ , dann heißt

$$\begin{aligned} \max c^T x \\ Ax \leq b \\ x \in \mathbb{Z}^n \end{aligned}$$

lineares ganzzahliges (oder kurz: ganzzahliges) Optimierungsproblem.

□

Selbst bei Optimierungsproblemen wie (1.8), . . . , (1.13), die nicht sonderlich allgemein erscheinen mögen, kann es sein, dass (bei spezieller Wahl der Zielfunktion  $f$  und der Nebenbedingungen) die Aufgabenstellung (finde ein  $x^* \in S$ , so dass  $f(x^*)$  so groß (oder klein) wie möglich ist) keine vernünftige Antwort besitzt. Es mag sein, dass  $f$  über  $S$  unbeschränkt ist;  $f$  kann beschränkt sein über  $S$ , aber ein Maximum kann innerhalb  $S$  nicht erreicht werden, d.h., das "max" müßte eigentlich durch "sup" ersetzt werden.  $S$  kann leer sein, ohne dass dies a priori klar ist, etc.. Der Leser möge sich Beispiele mit derartigen Eigenschaften überlegen! Bei der Formulierung von Problemen dieser Art muss man sich also Gedanken darüber machen, ob die betrachtete Fragestellung überhaupt eine sinnvolle Antwort erlaubt.

In unserer Vorlesung werden wir uns lediglich mit den Problemen (1.12) und (1.13) beschäftigen. Das lineare Optimierungsproblem (1.12) ist sicherlich das derzeit für die Praxis bedeutendste Problem, da sich außerordentlich viele und sehr unterschiedliche reale Probleme als lineare Programme formulieren lassen. Außerdem liegt eine sehr ausgefeilte Theorie vor, und mit den modernen Verfahren der linearen Optimierung können derartige Probleme mit Hunderttausenden (und manchmal sogar mehr) von Variablen und Ungleichungen gelöst werden. Dagegen ist Problem (1.13), viel schwieriger. Die Einschränkung der Lösungsmenge auf die zulässigen ganzzahligen Lösungen führt direkt zu einem Sprung im Schwierigkeitsgrad des Problems. Verschiedene spezielle lineare ganzzahlige Programme können in beliebiger Größenordnung gelöst werden. Bei wenig strukturierten allgemeinen Problemen des Typs (1.13) versagen dagegen die Lösungsverfahren manchmal bereits bei weniger als 100 Variablen und Nebenbedingungen.

Über Kontrolltheorie (Probleme des Typs (1.8)), Approximationstheorie (Probleme des Typs (1.9)), Nichtlineare Optimierung (Probleme des Typs (1.10) und (1.11)) werden an der TU Berlin Spezialvorlesungen angeboten. Es ist anzumerken, dass sich sowohl die Theorie als auch die Algorithmen zur Lösung von Pro-

blemen des Typs (1.8) bis (1.11) ganz erheblich von denen zur Lösung von Problemen des Typs (1.12) und (1.13) unterscheiden.

## 1.3 Notation

Wir gehen in der Vorlesung davon aus, dass die Hörer die Grundvorlesungen Lineare Algebra und Analysis kennen. Wir werden hauptsächlich Methoden der linearen Algebra benutzen und setzen eine gewisse Vertrautheit mit der Vektor- und Matrizenrechnung voraus.

### 1.3.1 Grundmengen

Wir benutzen folgende Bezeichnungen:

$$\begin{aligned}\mathbb{N} &= \{1, 2, 3, \dots\} = \text{Menge der } \mathbf{\text{natürlichen Zahlen}}, \\ \mathbb{Z} &= \text{Menge der } \mathbf{\text{ganzen Zahlen}}, \\ \mathbb{Q} &= \text{Menge der } \mathbf{\text{rationalen Zahlen}}, \\ \mathbb{R} &= \text{Menge der } \mathbf{\text{reellen Zahlen}},\end{aligned}$$

und mit  $M_+$  bezeichnen wir die Menge der nichtnegativen Zahlen in  $M$  für  $M \in \{\mathbb{Z}, \mathbb{Q}, \mathbb{R}\}$ . Wir betrachten  $\mathbb{Q}$  und  $\mathbb{R}$  als Körper mit der üblichen Addition und Multiplikation und der kanonischen Ordnung " $\leq$ ". Desgleichen betrachten wir  $\mathbb{N}$  und  $\mathbb{Z}$  als mit den üblichen Rechenarten versehen. Wir werden uns fast immer in diesen Zahlenuniversen bewegen, da diese die für die Praxis relevanten sind. Manche Sätze gelten jedoch nur, wenn wir uns auf  $\mathbb{Q}$  oder  $\mathbb{R}$  beschränken. Um hier eine saubere Trennung zu haben, treffen wir die folgende Konvention. Wenn wir das Symbol

$$\mathbb{K}$$

benutzen, so heißt dies immer das  $\mathbb{K}$  einer der angeordneten Körper  $\mathbb{R}$  oder  $\mathbb{Q}$  ist. Sollte ein Satz nur für  $\mathbb{R}$  oder nur für  $\mathbb{Q}$  gelten, so treffen wir die jeweils notwendige Einschränkung.

Für diejenigen, die sich für möglichst allgemeine Sätze interessieren, sei an dieser Stelle folgendes vermerkt. Jeder der nachfolgend angegebenen Sätze bleibt ein wahrer Satz, wenn wir als Grundkörper  $\mathbb{K}$  einen archimedisch angeordneten

Körper wählen. Ein bekannter Satz besagt, dass jeder archimedisch angeordnete Körper isomorph zu einem Unterkörper von  $\mathbb{R}$  ist, der  $\mathbb{Q}$  enthält. Unsere Sätze bleiben also richtig, wenn wir statt  $\mathbb{K} \in \{\mathbb{Q}, \mathbb{R}\}$  irgendeinen archimedisch angeordneten Körper  $\mathbb{K}$  mit  $\mathbb{Q} \subseteq \mathbb{K} \subseteq \mathbb{R}$  wählen.

Wir können in fast allen Sätzen (insbesondere bei denen, die keine Ganzzahligkeitsbedingungen haben) auch die Voraussetzung “archimedisch” fallen lassen, d. h. fast alle Sätze gelten auch für angeordnete Körper. Vieles, was wir im  $\mathbb{K}^n$  beweisen, ist auch in beliebigen metrischen Räumen oder Räumen mit anderen als euklidischen Skalarprodukten richtig. Diejenigen, die Spaß an derartigen Verallgemeinerungen haben, sind eingeladen, die entsprechenden Beweise in die allgemeinere Sprache zu übertragen.

In dieser Vorlesung interessieren wir uns für so allgemeine Strukturen nicht. Wir verbleiben in den (für die Praxis besonders wichtigen) Räumen, die über den reellen oder rationalen Zahlen errichtet werden. Also, nochmals, wenn immer wir das Symbol  $\mathbb{K}$  im weiteren gebrauchen, gilt

$$\mathbb{K} \in \{\mathbb{R}, \mathbb{Q}\},$$

und  $\mathbb{K}$  ist ein Körper mit den üblichen Rechenoperationen und Strukturen. Natürlich ist  $\mathbb{K}_+ = \{x \in \mathbb{K} \mid x \geq 0\}$

Die Teilmengenbeziehung zwischen zwei Mengen  $M$  und  $N$  bezeichnen wir wie üblich mit  $M \subseteq N$ . Gilt  $M \subseteq N$  und  $M \neq N$ , so schreiben wir  $M \subset N$ .  $M \setminus N$  bezeichnet die mengentheoretische Differenz  $\{x \in M \mid x \notin N\}$ .

### 1.3.2 Vektoren und Matrizen

Ist  $R$  eine beliebige Menge,  $n \in \mathbb{N}$ , so bezeichnen wir mit

$$R^n$$

die Menge aller **n-Tupel** oder Vektoren der Länge  $n$  mit Komponenten aus  $R$ . (Aus technischen Gründen ist es gelegentlich nützlich, Vektoren  $x \in R^0$ , also Vektoren ohne Komponenten, zu benutzen. Wenn wir dies tun, werden wir es explizit erwähnen, andernfalls setzen wir immer  $n \geq 1$  voraus.) Wir betrachten Vektoren  $x = (x_i)_{i=1, \dots, n} \in R^n$  immer als **Spaltenvektoren** d. h.

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix},$$

wollen wir mit Zeilenvektoren rechnen, so schreiben wir  $x^T$  (lies:  $x$  **transponiert**). Die Menge  $\mathbb{K}^n$  ist bekanntlich ein  $n$ -dimensionaler Vektorraum über  $\mathbb{K}$ . Mit

$$y^T x := \sum_{i=1}^n x_i y_i$$

bezeichnen wir das **innere Produkt** zweier Vektoren  $x, y \in \mathbb{K}^n$ . Wir nennen  $x$  und  $y$  **senkrecht (orthogonal)**, falls  $x^T y = 0$  gilt. Der  $\mathbb{K}^n$  ist für uns immer (wenn nichts anderes gesagt wird) mit der **euklidischen Norm**

$$\|x\| := \sqrt{x^T x}$$

ausgestattet.

Für Mengen  $S, T \subseteq \mathbb{K}^n$  und  $\alpha \in \mathbb{K}$  benutzen wir die folgenden Standardbezeichnungen für Mengenoperationen

$$\begin{aligned} S + T &:= \{x + y \in \mathbb{K}^n \mid x \in S, y \in T\}, \\ S - T &:= \{x - y \in \mathbb{K}^n \mid x \in S, y \in T\}, \\ \alpha S &:= \{\alpha x \in \mathbb{K}^n \mid x \in S\}. \end{aligned}$$

Einige Vektoren aus  $\mathbb{K}^n$  werden häufig auftreten, weswegen wir sie mit besonderen Symbolen bezeichnen. Mit  $e_j$  bezeichnen wir den Vektor aus  $\mathbb{K}^n$ , dessen  $j$ -te Komponente 1 und dessen übrige Komponenten 0

sind. Mit  $0$  bezeichnen wir den Nullvektor, mit  $\mathbb{1}$  den Vektor, dessen Komponenten alle 1 sind. Also

$$e_j = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad 0 = \begin{pmatrix} 0 \\ \vdots \\ \vdots \\ 0 \end{pmatrix}, \quad \mathbb{1} = \begin{pmatrix} 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix}.$$

Welche Dimension die Vektoren  $e_j$ ,  $0$ ,  $\mathbb{1}$  haben, ergibt sich jeweils aus dem Zusammenhang.

Für eine Menge  $R$  und  $m, n \in \mathbb{N}$  bezeichnet

$$R^{(m,n)} \text{ oder } R^{m \times n}$$

die Menge der  $(m, n)$ -**Matrizen** ( $m$  Zeilen,  $n$  Spalten) mit Einträgen aus  $R$ . (Aus technischen Gründen werden wir gelegentlich auch  $n = 0$  oder  $m = 0$  zulassen, d. h. wir werden auch Matrizen mit  $m$  Zeilen und ohne Spalten bzw.  $n$  Spalten und ohne Zeilen betrachten. Dieser Fall wird jedoch immer explizit erwähnt, somit ist in der Regel  $n \geq 1$  und  $m \geq 1$  vorausgesetzt.) Ist  $A \in R^{(m,n)}$ , so schreiben wir

$$A = (a_{ij})_{\substack{i=1,\dots,m \\ j=1,\dots,n}}$$

und meinen damit, dass  $A$  die folgende Form hat

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix},$$

Wenn nichts anderes gesagt wird, hat  $A$  die Zeilenindexmenge  $M = \{1, \dots, m\}$  und die Spaltenindexmenge  $N = \{1, \dots, n\}$ . Die  $j$ -te **Spalte** von  $A$  ist ein  $m$ -Vektor, den wir mit  $A_{.j}$  bezeichnen,

$$A_{.j} = \begin{pmatrix} a_{1j} \\ \vdots \\ a_{mj} \end{pmatrix}.$$

Die  $i$ -te **Zeile** von  $A$  ist ein Zeilenvektor der Länge  $n$ , den wir mit  $A_{i.}$  bezeichnen, d. h.

$$A_{i.} = (a_{i1}, a_{i2}, \dots, a_{in}).$$

Wir werden in dieser Vorlesung sehr häufig Untermatrizen von Matrizen konstruieren, sie umsortieren und mit ihnen rechnen müssen. Um dies ohne Zweideutigkeiten durchführen zu können, führen wir zusätzlich eine etwas exaktere als die oben eingeführte Standardbezeichnungsweise ein.

Wir werden — wann immer es aus technischen Gründen notwendig erscheint — die Zeilen- und Spaltenindexmengen einer  $(m, n)$ -Matrix  $A$  nicht als Mengen sondern als Vektoren auffassen. Wir sprechen dann vom

$$\begin{aligned} \text{vollen Zeilenindexvektor } M &= (1, 2, \dots, m) \\ \text{vollen Spaltenindexvektor } N &= (1, 2, \dots, n) \end{aligned}$$

von  $A$ . Ein **Zeilenindexvektor** von  $A$  ist ein Vektor mit höchstens  $m$  Komponenten, der aus  $M$  durch Weglassen einiger Komponenten von  $M$  und Permutation der übrigen Komponenten entsteht. Analog entsteht ein **Spaltenindexvektor**

durch Weglassen von Komponenten von  $N$  und Permutation der übrigen Komponenten. Ist also

$$\begin{aligned} I &= (i_1, i_2, \dots, i_p) && \text{ein Zeilenindexvektor von } A \text{ und} \\ J &= (j_1, j_2, \dots, j_q) && \text{ein Spaltenindexvektor von } A, \end{aligned}$$

so gilt immer  $i_s, i_t \in \{1, \dots, m\}$  und  $i_s \neq i_t$  für  $1 \leq s < t \leq p$ , und analog gilt  $j_s, j_t \in \{1, \dots, n\}$  und  $j_s \neq j_t$  für  $1 \leq s < t \leq q$ . Wir setzen

$$A_{IJ} := \begin{pmatrix} a_{i_1 j_1} & a_{i_1 j_2} & \cdots & a_{i_1 j_q} \\ a_{i_2 j_1} & a_{i_2 j_2} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ a_{i_p j_1} & a_{i_p j_2} & \cdots & a_{i_p j_q} \end{pmatrix}$$

und nennen  $A_{IJ}$  **Untermatrix** von  $A$ .

$A_{IJ}$  ist also eine  $(p, q)$ -Matrix, die aus  $A$  dadurch entsteht, dass man die Zeilen, die zu Indizes gehören, die nicht in  $I$  enthalten sind, und die Spalten, die zu Indizes gehören, die nicht in  $J$  enthalten sind, streicht und dann die so entstehende Matrix umsortiert.

Ist  $I = (i)$  und  $J = (j)$ , so erhalten wir zwei Darstellungsweisen für Zeilen bzw. Spalten von  $A$ :

$$\begin{aligned} A_{IN} &= A_{i.}, \\ A_{MJ} &= A_{.j}. \end{aligned}$$

Aus Gründen der Notationsvereinfachung werden wir auch die folgende (etwas unsaubere) Schreibweise benutzen. Ist  $M$  der volle Zeilenindexvektor von  $A$  und  $I$  ein Zeilenindexvektor, so schreiben wir auch

$$I \subseteq M,$$

obwohl  $I$  und  $M$  keine Mengen sind, und für  $i \in \{1, \dots, m\}$  benutzen wir

$$i \in I \quad \text{oder} \quad i \notin I,$$

um festzustellen, dass  $i$  als Komponente von  $I$  auftritt oder nicht. Analog verfahren wir bezüglich der Spaltenindizes.

Gelegentlich spielt die tatsächliche Anordnung der Zeilen und Spalten keine Rolle. Wenn also z. B.  $I \subseteq \{1, \dots, m\}$  und  $J \subseteq \{1, \dots, n\}$  gilt, dann werden wir auch einfach schreiben

$$A_{IJ},$$

obwohl diese Matrix dann nur bis auf Zeilen- und Spaltenpermutationen definiert ist. Wir werden versuchen, diese Bezeichnungen immer so zu verwenden, dass keine Zweideutigkeiten auftreten. Deshalb treffen wir ab jetzt die Verabredung, dass wir — falls wir Mengen (und nicht Indexvektoren)  $I$  und  $J$  benutzen — die Elemente von  $I = \{i_1, \dots, i_p\}$  und von  $J = \{j_1, \dots, j_q\}$  kanonisch anordnen, d. h. die Indizierung der Elemente von  $I$  und  $J$  sei so gewählt, dass  $i_1 < i_2 < \dots < i_p$  und  $j_1 < j_2 < \dots < j_q$  gilt. Bei dieser Verabredung ist dann  $A_{IJ}$  die Matrix die aus  $A$  durch Streichen der Zeilen  $i, i \notin I$ , und der Spalten  $j, j \notin J$ , entsteht.

Ist  $I \subseteq \{1, \dots, m\}$  und  $J \subseteq \{1, \dots, n\}$  oder sind  $I$  und  $J$  Zeilen- bzw. Spaltenindexvektoren, dann schreiben wir auch

$$A_I \quad \text{statt} \quad A_{IN}$$

$$A_{.J} \quad \text{statt} \quad A_{MJ}$$

$A_I$  entsteht also aus  $A$  durch Streichen der Zeilen  $i, i \notin I$ ,  $A_{.J}$  durch Streichen der Spalten  $j, j \notin J$ .

Sind  $A, B \in \mathbb{K}^{(m,n)}$ ,  $C \in \mathbb{K}^{(n,s)}$ ,  $\alpha \in \mathbb{K}$ , so sind

$$\begin{array}{ll} \text{die Summe} & A + B, \\ \text{das Produkt} & \alpha A, \\ \text{das Matrixprodukt} & AC \end{array}$$

wie in der linearen Algebra üblich definiert.

Für einige häufig auftretende Matrizen haben wir spezielle Symbole reserviert. Mit  $0$  bezeichnen wir die Nullmatrix (alle Matrixelemente sind Null), wobei sich die Dimension der Nullmatrix jeweils aus dem Zusammenhang ergibt. (Das Symbol  $0$  kann also sowohl eine Zahl, einen Vektor als auch eine Matrix bezeichnen). Mit  $I$  bezeichnen wir die Einheitsmatrix. Diese Matrix ist quadratisch, die Hauptdiagonalelemente von  $I$  sind Eins, alle übrigen Null. Wollen wir die Dimension von  $I$  betonen, so schreiben wir auch  $I_n$  und meinen damit die  $(n, n)$ -Einheitsmatrix. Diejenige  $(m, n)$ -Matrix, bei der alle Elemente Eins sind, bezeichnen wir mit  $E$ . Wir schreiben auch  $E_{m,n}$  bzw.  $E_n$ , um die Dimension zu spezifizieren ( $E_n$  ist eine  $(n, n)$ -Matrix). Ist  $x$  ein  $n$ -Vektor, so bezeichnet  $\text{diag}(x)$  diejenige  $(n, n)$ -Matrix  $A = (a_{ij})$  mit  $a_{ii} = x_i$  ( $i = 1, \dots, n$ ) und  $a_{ij} = 0$  ( $i \neq j$ ).

Wir halten an dieser Stelle noch einmal Folgendes fest: Wenn wir von einer Matrix  $A$  sprechen, ohne anzugeben, welche Dimension sie hat und aus welchem Bereich sie ist, dann nehmen wir implizit an, dass  $A \in \mathbb{K}^{(m,n)}$  gilt. Analog gilt immer  $x \in \mathbb{K}^n$ , wenn sich nicht aus dem Zusammenhang anderes ergibt.

### 1.3.3 Kombinationen von Vektoren, Hüllen, Unabhängigkeit

Ein Vektor  $x \in \mathbb{K}^n$  heißt **Linearkombination** der Vektoren  $x_1, \dots, x_k \in \mathbb{K}^n$ , falls es einen Vektor  $\lambda = (\lambda_1, \dots, \lambda_k)^T \in \mathbb{K}^k$  gibt mit

$$x = \sum_{i=1}^k \lambda_i x_i .$$

Gilt zusätzlich:

$$\left. \begin{array}{l} \lambda \geq 0 \\ \lambda^T \mathbf{1} = 1 \\ \lambda \geq 0 \quad \text{und} \quad \lambda^T \mathbf{1} = 1 \end{array} \right\} \text{ so heißt } x \left\{ \begin{array}{l} \text{konische} \\ \text{affine} \\ \text{konvexe} \end{array} \right\} \text{ Kombination}$$

der Vektoren  $x_1, \dots, x_k$ . Diese Kombinationen heißen **echt**, falls weder  $\lambda = 0$  noch  $\lambda = e_j$  für ein  $j \in \{1, \dots, k\}$  gilt.

Für eine nichtleere Teilmenge  $S \subseteq \mathbb{K}^n$  heißt:

$$\left. \begin{array}{l} \text{lin}(S) \\ \text{cone}(S) \\ \text{aff}(S) \\ \text{conv}(S) \end{array} \right\} \text{ die } \left\{ \begin{array}{l} \text{lineare} \\ \text{konische} \\ \text{affine} \\ \text{konvexe} \end{array} \right\} \text{ Hülle von } S, \text{ d. h.}$$

die Menge aller Vektoren, die als lineare (konische, affine oder konvexe) Kombination von endlich vielen Vektoren aus  $S$  dargestellt werden können. Wir setzen außerdem

$$\begin{aligned} \text{lin}(\emptyset) &:= \text{cone}(\emptyset) := \{0\}, \\ \text{aff}(\emptyset) &:= \text{conv}(\emptyset) := \emptyset. \end{aligned}$$

Ist  $A$  eine  $(m, n)$ -Matrix, so schreiben wir auch

$$\text{lin}(A), \text{cone}(A), \text{aff}(A), \text{conv}(A)$$

und meinen damit die lineare, konische, affine bzw. konvexe Hülle der Spaltenvektoren  $A_{.1}, A_{.2}, \dots, A_{.n}$  von  $A$ . Eine Teilmenge  $S \subseteq \mathbb{K}^n$  heißt

$$\left\{ \begin{array}{l} \text{linearer Raum} \\ \text{Kegel} \\ \text{affiner Raum} \\ \text{konvexe Menge} \end{array} \right\} \text{ falls } \left\{ \begin{array}{l} S = \text{lin}(S) \\ S = \text{cone}(S) \\ S = \text{aff}(S) \\ S = \text{conv}(S) \end{array} \right\} .$$

Die hier benutzten Begriffe sind üblicher Standard, wobei für „linearen Raum“ in der linearen Algebra in der Regel *Vektorraum* oder *Untervektorraum* benutzt wird. Der Begriff *Kegel* wird jedoch – u. a. in verschiedenen Zweigen der Geometrie – allgemeiner verwendet. „Unser Kegel“ ist in der Geometrie ein „abgeschlossener konvexer Kegel“.

Eine nichtleere endliche Teilmenge  $S \subseteq \mathbb{K}^n$  heißt **linear** (bzw. **affin**) **unabhängig**, falls kein Element von  $S$  als echte Linearkombination (bzw. Affinkombination) von Elementen von  $S$  dargestellt werden kann. Die leere Menge ist affin, jedoch nicht linear unabhängig. Jede Menge  $S \subseteq \mathbb{K}^n$ , die nicht linear bzw. affin unabhängig ist, heißt **linear** bzw. **affin abhängig**. Aus der linearen Algebra wissen wir, dass eine linear (bzw. affin) unabhängige Teilmenge des  $\mathbb{K}^n$  höchstens  $n$  (bzw.  $n + 1$ ) Elemente enthält. Für eine Teilmenge  $S \subseteq \mathbb{K}^n$  heißt die Kardinalität einer größten linear (bzw. affin) unabhängigen Teilmenge von  $S$  der **Rang** (bzw. **affine Rang**) von  $S$ . Wir schreiben dafür  $\text{rang}(S)$  bzw.  $\text{arang}(S)$ . Die **Dimension** einer Teilmenge  $S \subseteq \mathbb{K}^n$ , Bezeichnung:  $\dim(S)$ , ist die Kardinalität einer größten affin unabhängigen Teilmenge von  $S$  minus 1, d. h.  $\dim(S) = \text{arang}(S) - 1$ .

Der **Rang einer Matrix**  $A$ , bezeichnet mit  $\text{rang}(A)$ , ist der Rang ihrer Spaltenvektoren. Aus der linearen Algebra wissen wir, dass  $\text{rang}(A)$  mit dem Rang der Zeilenvektoren von  $A$  übereinstimmt. Gilt für eine  $(m, n)$ -Matrix  $A$ ,  $\text{rang}(A) = \min\{m, n\}$ , so sagen wir, dass  $A$  **vollen Rang** hat. Eine  $(n, n)$ -Matrix mit vollem Rang ist **regulär**, d. h. sie besitzt eine (eindeutig bestimmte) inverse Matrix (geschrieben  $A^{-1}$ ) mit der Eigenschaft  $AA^{-1} = I$ .

# Kapitel 2

## Polyeder

Das Hauptanliegen dieser Vorlesung ist die Behandlung von Problemen der linearen Programmierung. Man kann die Verfahren zur Lösung derartiger Probleme rein algebraisch darstellen als Manipulation von linearen Gleichungs- und Ungleichungssystemen. Die meisten Schritte dieser Algorithmen haben jedoch auch eine geometrische Bedeutung; und wenn man erst einmal die geometrischen Prinzipien und Regeln verstanden hat, die hinter diesen Schritten liegen, ist es viel einfacher, die Algorithmen und ihre Korrektheitsbeweise zu verstehen. Wir wollen daher in dieser Vorlesung einen geometrischen Zugang wählen und zunächst die Lösungsmengen linearer Programme studieren. Speziell wollen wir unser Augenmerk auf geometrische Sachverhalte und ihre algebraische Charakterisierung legen. Dies wird uns helfen, aus geometrischen Einsichten algebraische Verfahren abzuleiten.

### (2.1) Definition.

- (a) Eine Teilmenge  $G \subseteq \mathbb{K}^n$  heißt **Hyperebene**, falls es einen Vektor  $a \in \mathbb{K}^n \setminus \{0\}$  und  $\alpha \in \mathbb{K}$  gibt mit

$$G = \{x \in \mathbb{K}^n \mid a^T x = \alpha\}.$$

Der Vektor  $a$  heißt **Normalenvektor** zu  $G$ .

- (b) Eine Teilmenge  $H \subseteq \mathbb{K}^n$  heißt **Halbraum**, falls es einen Vektor  $a \in \mathbb{K}^n \setminus \{0\}$  und  $\alpha \in \mathbb{K}$  gibt mit

$$H = \{x \in \mathbb{K}^n \mid a^T x \leq \alpha\}.$$

Wir nennen  $a$  den **Normalenvektor** zu  $H$ . Die Hyperebene  $G = \{x \in \mathbb{K}^n \mid a^T x = \alpha\}$  heißt die zum Halbraum  $H$  gehörende Hyperebene oder die **berandende Hyperebene**, und  $H$  heißt der zu  $G$  gehörende Halbraum.

- (c) Eine Teilmenge  $P \subseteq \mathbb{K}^n$  heißt **Polyeder**, falls es ein  $m \in \mathbb{Z}_+$ , eine Matrix  $A \in \mathbb{K}^{(m,n)}$  und einen Vektor  $b \in \mathbb{K}^m$  gibt mit

$$P = \{x \in \mathbb{K}^n \mid Ax \leq b\}.$$

Um zu betonen, dass  $P$  durch  $A$  und  $b$  definiert ist, schreiben wir auch

$$P = \mathbf{P}(A, b) := \{x \in \mathbb{K}^n \mid Ax \leq b\}.$$

- (d) Ein Polyeder  $P$  heißt **Polytop**, wenn es beschränkt ist, d. h. wenn es ein  $B \in \mathbb{K}$ ,  $B > 0$  gibt mit  $P \subseteq \{x \in \mathbb{K}^n \mid \|x\| \leq B\}$ .

□

Polyeder können wir natürlich auch in folgender Form schreiben

$$P(A, b) = \bigcap_{i=1}^m \{x \in \mathbb{K}^n \mid A_i \cdot x \leq b_i\}.$$

Halbräume sind offensichtlich Polyeder. Aber auch die leere Menge ist ein Polyeder, denn  $\emptyset = \{x \mid 0^T x \leq -1\}$ , und der gesamte Raum ist ein Polyeder, denn  $\mathbb{K}^n = \{x \mid 0^T x \leq 0\}$ . Sind alle Zeilenvektoren  $A_i$  von  $A$  vom Nullvektor verschieden, so sind die bei der obigen Durchschnittsbildung beteiligten Mengen Halbräume. Ist ein Zeilenvektor von  $A$  der Nullvektor, sagen wir  $A_1 = 0^T$ , so ist  $\{x \in \mathbb{K}^n \mid A_1 \cdot x \leq b_1\}$  entweder leer (falls  $b_1 < 0$ ) oder der gesamte Raum  $\mathbb{K}^n$  (falls  $b_1 \geq 0$ ). Das heißt, entweder ist  $P(A, b)$  leer oder die Mengen  $\{x \mid A_i \cdot x \leq b_i\}$  mit  $A_i = 0$  können bei der obigen Durchschnittsbildung weggelassen werden. Daraus folgt

Jedes Polyeder  $P \neq \mathbb{K}^n$  ist Durchschnitt von endlich vielen Halbräumen.

Gilt  $P = P(A, b)$ , so nennen wir das Ungleichungssystem  $Ax \leq b$  ein **P definierendes System** (von linearen Ungleichungen). Sind  $\alpha > 0$  und  $1 \leq i < j \leq m$ , so gilt offensichtlich

$$P(A, b) = P(A, b) \cap \{x \mid \alpha A_i \cdot x \leq \alpha b_i\} \cap \{x \mid (A_i + A_j) \cdot x \leq b_i + b_j\}.$$

Daraus folgt, dass  $A$  und  $b$  zwar  $P = P(A, b)$  eindeutig bestimmen, dass aber  $P$  unendlich viele Darstellungen der Form  $P(D, d)$  hat.

**(2.2) Beispiel.** Wir betrachten das Ungleichungssystem

$$(1) \quad 2x_1 \leq 5$$

$$(2) \quad -2x_2 \leq 1$$

$$(3) \quad -x_1 - x_2 \leq -1$$

$$(4) \quad 2x_1 + 9x_2 \leq 23$$

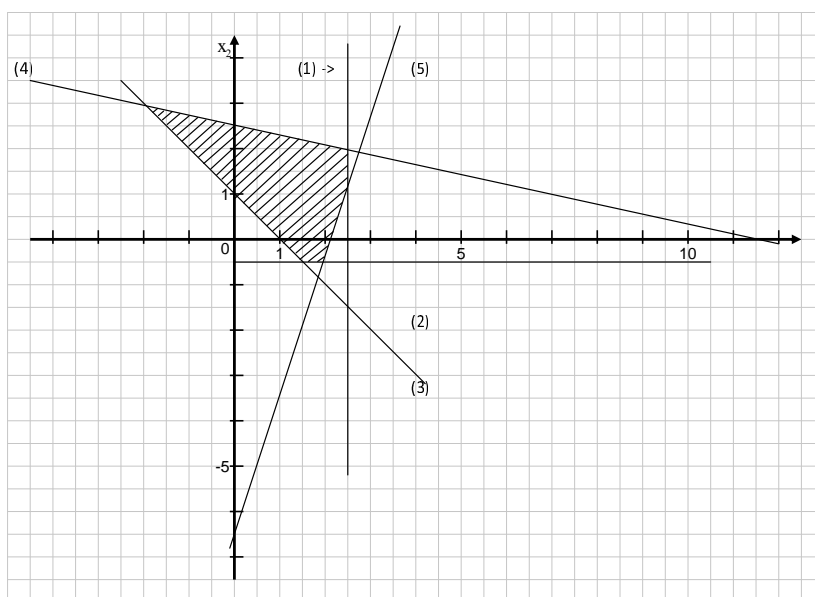
$$(5) \quad 6x_1 - 2x_2 \leq 13$$

Hieraus erhalten wir die folgende Matrix  $A$  und den Vektor  $b$ :

$$A = \begin{pmatrix} 2 & 0 \\ 0 & -2 \\ -1 & -1 \\ 2 & 9 \\ 6 & -2 \end{pmatrix} \quad b = \begin{pmatrix} 5 \\ 1 \\ -1 \\ 23 \\ 13 \end{pmatrix}$$

Das Polyeder  $P = P(A, b)$  ist die Lösungsmenge des obigen Ungleichungssystems (1), ..., (5) und ist in Abbildung 2.1 graphisch dargestellt.

□



**Abb. 2.1**

Die Mengen zulässiger Lösungen linearer Programme treten nicht immer in der Form  $Ax \leq b$  auf. Häufig gibt es auch Gleichungen und vorzeichenbeschränkte Variable. Vorzeichenbeschränkungen sind natürlich auch lineare Ungleichungssysteme, und ein Gleichungssystem  $Dx = d$  kann in der Form von zwei Ungleichungssystemen  $Dx \leq d$  und  $-Dx \leq -d$  geschrieben werden. Allgemeiner gilt

**(2.3) Bemerkung.** Die Lösungsmenge des Systems

$$\begin{aligned} Bx + Cy &= c \\ Dx + Ey &\leq d \\ x &\geq 0 \\ x \in \mathbb{K}^p, \quad y &\in \mathbb{K}^q \end{aligned}$$

ist ein Polyeder.

**Beweis :** Setze  $n := p + q$  und

$$A := \begin{pmatrix} B & C \\ -B & -C \\ D & E \\ -I & 0 \end{pmatrix}, \quad b := \begin{pmatrix} c \\ -c \\ d \\ 0 \end{pmatrix}$$

dann ist  $P(A, b)$  die Lösungsmenge des vorgegebenen Gleichungs- und Ungleichungssystems.  $\square$

Ein spezieller Polyedertyp wird uns häufig begegnen, weswegen wir für ihn eine besondere Bezeichnung wählen wollen. Für  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$  setzen wir

$$P^=(A, b) := \{x \in \mathbb{K}^n \mid Ax = b, x \geq 0\}.$$

Nicht alle Polyeder können in der Form  $P^=(A, b)$  dargestellt werden, z. B. nicht  $P = \{x \in \mathbb{K} \mid x \leq 1\}$ .

Wir werden später viele Sätze über Polyeder  $P$  beweisen, deren Aussagen **darstellungsabhängig** sind, d. h. die Art und Weise, wie  $P$  gegeben ist, geht explizit in die Satzaussage ein. So werden sich z. B. die Charakterisierungen gewisser Polyedereigenschaften von  $P(A, b)$  (zumindest formal) von den entsprechenden Charakterisierungen von  $P^=(A, b)$  unterscheiden. Darstellungsabhängige Sätze wollen wir jedoch nur einmal beweisen (normalerweise für Darstellungen, bei denen die Resultate besonders einprägsam oder einfach sind), deshalb werden wir uns nun Transformationsregeln überlegen, die angeben, wie man von einer Darstellungsweise zu einer anderen und wieder zurück kommt.

**(2.4) Transformationen.****Regel I: Einführung von Schlupfvariablen**

Gegeben seien  $a \in \mathbb{K}^n$ ,  $\alpha \in \mathbb{K}$ . Wir schreiben die Ungleichung

$$(i) \quad a^T x \leq \alpha$$

in der Form einer Gleichung und einer Vorzeichenbeschränkung

$$(ii) \quad a^T x + y = \alpha, \quad y \geq 0.$$

$y$  ist eine neue Variable, genannt **Schlupfvariable**.

Es gilt:

$$\begin{aligned} x \text{ erfüllt (i)} &\implies \begin{pmatrix} x \\ y \end{pmatrix} \text{ erfüllt (ii) für } y = \alpha - a^T x. \\ \begin{pmatrix} x \\ y \end{pmatrix} \text{ erfüllt (ii)} &\implies x \text{ erfüllt (i).} \end{aligned}$$

Allgemein: Ein Ungleichungssystem  $Ax \leq b$  kann durch Einführung eines Schlupfvariablenvektors  $y$  transformiert werden in ein Gleichungssystem mit Vorzeichenbedingung  $Ax + y = b$ ,  $y \geq 0$ . Zwischen den Lösungsmengen dieser beiden Systeme bestehen die oben angegebenen Beziehungen.  $P(A, b)$  und  $\left\{ \begin{pmatrix} x \\ y \end{pmatrix} \in \mathbb{K}^{n+m} \mid Ax + Iy = b, y \geq 0 \right\}$  sind jedoch zwei durchaus verschiedene Polyeder in verschiedenen Vektorräumen.

**Regel II: Einführung von vorzeichenbeschränkten Variablen**

Ist  $x$  eine (eindimensionale) nicht vorzeichenbeschränkte Variable, so können wir zwei vorzeichenbeschränkte Variablen  $x^+$  und  $x^-$  einführen, um  $x$  darzustellen. Wir setzen

$$x := x^+ - x^- \quad \text{mit} \quad x^+ \geq 0, \quad x^- \geq 0.$$

□

Mit den Regeln I und II aus (2.4) können wir z. B. jedem Polyeder  $P(A, b) \subseteq \mathbb{K}^n$  ein Polyeder  $P=(D, d) \subseteq \mathbb{K}^{2n+m}$  wie folgt zuordnen. Wir setzen

$$D := (A, -A, I_m), \quad d := b,$$

d. h. es gilt

$$P^=(D, d) = \left\{ \begin{pmatrix} x \\ y \\ z \end{pmatrix} \in \mathbb{K}^{2n+m} \mid Ax - Ay + z = b, x, y, z \geq 0 \right\}.$$

Es ist üblich, die oben definierten Polyeder  $P(A, b)$  und  $P^=(D, d)$  **äquivalent** zu nennen. Hierbei hat "Äquivalenz" folgende Bedeutung.

Für  $x \in \mathbb{K}^n$  sei

$$\begin{aligned} x^+ &:= (x_1^+, \dots, x_n^+)^T \quad \text{mit} \quad x_i^+ = \max\{0, x_i\}, \\ x^- &:= (x_1^-, \dots, x_n^-)^T \quad \text{mit} \quad x_i^- = \max\{0, -x_i\}, \end{aligned}$$

dann gilt:

$$\begin{aligned} x \in P(A, b) &\implies \begin{pmatrix} x^+ \\ x^- \\ z \end{pmatrix} \in P^=(D, d) \quad \text{für } z = b - Ax \\ \begin{pmatrix} u \\ v \\ w \end{pmatrix} \in P^=(D, d) &\implies x := u - v \in P(A, b). \end{aligned}$$

Besonders einfache Polyeder, auf die sich jedoch fast alle Aussagen der Polyedertheorie zurückführen lassen, sind polyedrische Kegel.

**(2.5) Definition.** Ein Kegel  $C \subseteq \mathbb{K}^n$  heißt **polyedrisch** genau dann, wenn  $C$  ein Polyeder ist.

□

**(2.6) Bemerkung.** Ein Kegel  $C \subseteq \mathbb{K}^n$  ist genau dann polyedrisch, wenn es eine Matrix  $A \in \mathbb{K}^{(m,n)}$  gibt mit

$$C = P(A, 0).$$

**Beweis :** Gilt  $C = P(A, 0)$ , so ist  $C$  ein Polyeder und offensichtlich ein Kegel.

Sei  $C$  ein polyedrischer Kegel, dann existieren nach Definition (2.1) (c) eine Matrix  $A$  und ein Vektor  $b$  mit  $C = P(A, b)$ . Da jeder Kegel den Nullvektor enthält,

gilt  $0 = A0 \leq b$ . Angenommen, es existiert ein  $\bar{x} \in C$  mit  $A\bar{x} \not\leq 0$ , d. h. es existiert eine Zeile von  $A$ , sagen wir  $A_{i\cdot}$ , mit  $t := A_{i\cdot}\bar{x} > 0$ . Da  $C$  ein Kegel ist, gilt  $\lambda\bar{x} \in C$  für alle  $\lambda \in \mathbb{K}_+$ . Jedoch für  $\bar{\lambda} := \frac{b_i}{t} + 1$  gilt einerseits  $\bar{\lambda}\bar{x} \in C$  und andererseits  $A_{i\cdot}(\bar{\lambda}\bar{x}) = \bar{\lambda}t > b_i$ , ein Widerspruch. Daraus folgt, für alle  $x \in C$  gilt  $Ax \leq 0$ . Hieraus ergibt sich  $C = P(A, 0)$ .  $\square$

In der folgenden Tabelle 2.1 haben wir alle möglichen Transformationen aufgelistet. Sie soll als “Nachschlagewerk” dienen.

<p>Transformation</p> <p>nach</p> <p>von</p>	$By \leq d$	$By = d$ $y \geq 0$	$By \leq d$ $y \geq 0$
$Ax = b$ hierbei ist $x_i^+ = \max\{0, x_i\}$ $x_i^- = \max\{0, -x_i\}$	$\begin{pmatrix} A \\ -A \end{pmatrix} y \leq \begin{pmatrix} b \\ -b \end{pmatrix}$ $y = x$	$(A, -A) y = b$ $y \geq 0$ $y = \begin{pmatrix} x^+ \\ x^- \end{pmatrix}$	$\begin{pmatrix} A & -A \\ -A & A \end{pmatrix} y \leq \begin{pmatrix} b \\ -b \end{pmatrix}$ $y \geq 0$ $y = \begin{pmatrix} x^+ \\ x^- \end{pmatrix}$
$Ax \leq b$	$Ay \leq b$ $y = x$	$(A, -A, I) y = b$ $y \geq 0$ $y = \begin{pmatrix} x^+ \\ x^- \\ b - Ax \end{pmatrix}$	$(A, -A) y \leq b$ $y \geq 0$ $y = \begin{pmatrix} x^+ \\ x^- \end{pmatrix}$
$Ax = b$ $x \geq 0$	$\begin{pmatrix} A \\ -A \\ -I \end{pmatrix} y \leq \begin{pmatrix} b \\ -b \\ 0 \end{pmatrix}$ $y = x$	$Ay = b$ $y \geq 0$ $y = x$	$\begin{pmatrix} A \\ -A \end{pmatrix} y \leq \begin{pmatrix} b \\ -b \end{pmatrix}$ $y \geq 0$ $y = x$
$Ax \leq b$ $x \geq 0$	$\begin{pmatrix} A \\ -I \end{pmatrix} y \leq \begin{pmatrix} b \\ 0 \end{pmatrix}$ $y = x$	$(A, I) y = b$ $y \geq 0$ $y = \begin{pmatrix} x \\ b - Ax \end{pmatrix}$	$Ay \leq b$ $y \geq 0$ $y = x$

# Kapitel 3

## Fourier-Motzkin-Elimination und Projektion

In diesem Kapitel untersuchen wir die folgende Frage: Wie kann man herausfinden, ob ein lineares Ungleichungssystem  $Ax \leq b$  eine Lösung hat oder nicht? Nehmen wir an, dass  $A$  und  $b$  rational sind (wenn man Computer benutzen will, muss man diese Annahmen sowieso machen), dann können wir, da  $\mathbb{Q}^n$  abzählbar ist, die Elemente von  $\mathbb{Q}^n$  durchnummerieren und iterativ jedes Element dieser Folge in das System  $Ax \leq b$  einsetzen. Falls  $Ax \leq b$  lösbar ist, werden wir nach endlich vielen Schritten einen zulässigen Vektor finden. Hat  $Ax \leq b$  keine zulässige Lösung, so läuft diese Enumeration unendlich lange. Kann man die Unlösbarkeit vielleicht auch mit einem endlichen Verfahren prüfen? Wir werden hierauf in diesem Kapitel eine positive Antwort geben.

Wir beginnen mit der Beschreibung eines Algorithmus, der aus einer  $(m, n)$ -Matrix  $A$  und einem Vektor  $b \in \mathbb{K}^m$  eine  $(r, n)$ -Matrix  $D$  und einen Vektor  $d \in \mathbb{K}^r$  macht, so dass eine Spalte von  $D$  aus lauter Nullen besteht und dass gilt:  $Ax \leq b$  hat eine Lösung genau dann, wenn  $Dx \leq d$  eine Lösung hat.

### 3.1 Fourier-Motzkin-Elimination

**(3.1) Algorithmus. Fourier-Motzkin-Elimination (der  $j$ -ten Variablen).**

**Input:** Eine  $(m, n)$ -Matrix  $A = (a_{ij})$ , ein Vektor  $b \in \mathbb{K}^m$  und ein Spaltenindex  $j \in \{1, \dots, n\}$  der Matrix  $A$ .

**Output:** Eine  $(r, n)$ -Matrix  $D = (d_{ij})$  ( $r$  wird im Algorithmus berechnet) und ein

Vektor  $d \in \mathbb{K}^r$ , so dass  $D_{.j}$  (die  $j$ -te Spalte von  $D$ ) der Nullvektor ist.

**Schritt 1.** Partitioniere die Menge der Zeilenindizes  $M = \{1, \dots, m\}$  von  $A$  wie folgt:

$$\begin{aligned} N &:= \{i \in M \mid a_{ij} < 0\}, \\ Z &:= \{i \in M \mid a_{ij} = 0\}, \\ P &:= \{i \in M \mid a_{ij} > 0\}. \end{aligned}$$

(Die Menge  $Z \cup (N \times P)$  wird die Zeilenindexmenge von  $D$ .)

**Schritt 2.** Setze  $r := |Z \cup (N \times P)|$ ,  $R := \{1, \dots, r\}$ , sei  $p : R \rightarrow Z \cup (N \times P)$  eine Bijektion (d.h. eine kanonische Indizierung der Elemente von  $Z \cup (N \times P)$ ).

**Schritt 3.** Führe für  $i = 1, 2, \dots, r$  aus:

(a) Falls  $p(i) \in Z$ , dann setze  $D_{i.} := A_{p(i).}$ ,  $d_i := b_{p(i)}$

(d.h., die  $i$ -te Zeile von  $D$  ist gleich der  $p(i)$ -ten Zeile von  $A$ ).

(b) Falls  $p(i) = (s, t) \in N \times P$ , dann setze

$$\begin{aligned} D_{i.} &:= a_{tj}A_{s.} - a_{sj}A_{t.}, \\ d_i &:= a_{tj}b_s - a_{sj}b_t \end{aligned}$$

(d.h., das  $a_{sj}$ -fache der  $t$ -ten Zeile von  $A$  wird vom  $a_{tj}$ -fachen der  $s$ -ten Zeile von  $A$  abgezogen und als  $i$ -te Zeile von  $D$  betrachtet).

**Ende.**

**Hinweis:** Die in (3.1) konstruierte Menge  $Z \cup (N \times P)$  kann leer und somit  $r = 0$  sein. Die Matrix  $D$  ist dann eine Matrix mit 0 Zeilen und  $n$  Spalten und somit  $Dx \leq d$  ein „leeres Ungleichungssystem“. Die Menge der zulässigen Lösungen  $P(D, d) = \{x \in \mathbb{K}^n \mid Dx \leq d\}$  unterliegt folglich keiner Einschränkung und ist daher der ganze Raum  $\mathbb{K}^n$ .

Bevor wir den Algorithmus analysieren, betrachten wir ein Beispiel:

$$\begin{aligned} (1) \quad & -7x_1 + 6x_2 \leq 25 \\ (2) \quad & +x_1 - 5x_2 \leq 1 \\ (3) \quad & +x_1 \leq 7 \\ (4) \quad & -x_1 + 2x_2 \leq 12 \\ (5) \quad & -x_1 - 3x_2 \leq 1 \\ (6) \quad & +2x_1 - x_2 \leq 10 \end{aligned}$$

Wir wollen die zweite Variable  $x_2$  eliminieren und erhalten in Schritt 1 von (3.1):

$$N = \{2, 5, 6\}, \quad Z = \{3\}, \quad P = \{1, 4\}.$$

Daraus ergibt sich  $|Z \cup (N \times P)| = 7$ . Die Nummerierung der Zeilen von  $D$  ist dann  $R = \{1, \dots, 7\}$ , und wir setzen:

$$p(1) = 3, \quad p(2) = (2, 1), \quad p(3) = (2, 4), \quad p(4) = (5, 1), \quad p(5) = (5, 4), \quad p(6) = (6, 1), \quad p(7) = (6, 7).$$

Die zweite Zeile von  $D$  ergibt sich dann als die Summe vom 6-fachen der zweiten Zeile von  $A$  und dem 5-fachen der ersten Zeile von  $A$ . Die zweite Ungleichung des Systems  $Dx \leq d$  hat somit die Form  $-29x_1 + 0x_2 \leq 131$ .

Insgesamt ergibt sich (wenn wir die aus lauter Nullen bestehende zweite Spalte weglassen), das folgende System  $Dx \leq d$ :

$$\begin{array}{rcl} (1) & + & x_1 \leq 7 \\ (2) & - & 29x_1 \leq 131 \\ (3) & - & 3x_1 \leq 62 \\ (4) & - & 27x_1 \leq 81 \\ (5) & - & 5x_1 \leq 38 \\ (6) & + & 5x_1 \leq 85 \\ (7) & + & 3x_1 \leq 32 \end{array}$$

Ganz offensichtlich können einige der Koeffizienten gekürzt werden. Wir sehen, dass die Zeilen (1), (6), (7) obere und die Ungleichungen (2), (3), (4), (5) untere Schranken auf den für  $x_1$  zulässigen Wert liefern. Wir sehen aber auch, dass die Eliminationsmethode viele überflüssige (wir werden das später *redundante* nennen) Ungleichungen liefert.

Rechentechnisch ist es offensichtlich geschickter, die Zeilen der Matrix  $A$  vor Beginn des Algorithmus zu skalieren. In einem Vorbereitungsschritt multiplizieren wir jede Zeile  $A_i$  der Matrix  $A$  und die Zahl  $b_i$  mit dem Wert  $|\frac{1}{a_{ij}}|$ . Nennen wir die so skalierte Matrix ebenfalls  $A$ , so haben die Elemente der  $j$ -ten Spalte den Wert 0, +1 oder -1. Die Berechnung der Zeile  $D_i$  der Matrix  $D$  und der rechten Seite  $d_i$  des neuen Ungleichungssystems  $Dx \leq d$  in Schritt 3.(b) geschieht dann durch

$$\begin{aligned} D_i &:= A_s - A_t. \\ d_i &:= b_s - b_t. \end{aligned}$$

Wir wollen nun das Ergebnis der Elimination der  $j$ -ten Variablen analysieren. Wir schauen zunächst einige Trivialfälle an.

Wenn alle Elemente der Spalte  $A_{.j}$  Null sind, so ist  $A = D$ . Gilt  $M = P$  oder  $M = N$ , sind also alle Elemente der Spalte  $A_{.j}$  von Null verschieden und haben dasselbe Vorzeichen, dann ist  $r = 0$  und somit  $D$  die leere Matrix mit 0 Zeilen und  $n$  Spalten.  $P(D, d)$  ist dann der gesamte Raum  $\mathbb{K}^n$ .

**(3.2) Satz.** Sei  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$ ,  $j$  ein Spaltenindex von  $A$ , und seien  $D \in \mathbb{K}^{(r,n)}$ ,  $d \in \mathbb{K}^r$  die in (3.1) konstruierte Matrix bzw. Vektor. Es gilt:

- (a) Die  $j$ -te Spalte von  $D$  ist der Nullvektor.
- (b) Jede Ungleichung  $D_i \cdot x \leq d_i$ ,  $i = 1, \dots, r$ , ist eine konische Kombination der Ungleichungen  $A_k \cdot x \leq b_k$ ,  $k = 1, \dots, m$ , d. h. es existiert ein Vektor  $u \in \mathbb{R}^m$ ,  $u \geq 0$  mit  $u^T A = D_i$  und  $u^T b = d_i$ .
- (c) Daraus folgt, es existiert eine  $(r, m)$ -Matrix  $U$ ,  $U \geq 0$ , mit  $UA = D$  und  $Ub = d$ .
- (d) Sei  $\bar{x} \in \mathbb{K}^n$  mit  $\bar{x}_j = 0$ , sei  $e_j$  der  $j$ -te Einheitsvektor und seien  $N, Z, P$  die in Schritt 1 von (3.1) definierten Mengen.

Setzen wir:

$$\lambda_i := \frac{1}{a_{ij}}(b_i - A_i \cdot \bar{x}) \quad \text{für alle } i \in P \cup N$$

$$L := \begin{cases} -\infty & \text{falls } N = \emptyset \\ \max\{\lambda_i \mid i \in N\} & \text{falls } N \neq \emptyset \end{cases}$$

$$U := \begin{cases} +\infty & \text{falls } P = \emptyset. \\ \min\{\lambda_i \mid i \in P\} & \text{falls } P \neq \emptyset. \end{cases}$$

Dann gilt

$$(d_1) \quad \bar{x} \in P(D, d) \quad \Rightarrow \quad L \leq U \quad \text{und} \quad \bar{x} + \lambda e_j \in P(A, b) \quad \forall \lambda \in [L, U]$$

$$(d_2) \quad \bar{x} + \lambda e_j \in P(A, b) \quad \Rightarrow \quad \lambda \in [L, U] \quad \text{und} \quad \bar{x} \in P(D, d).$$

**Beweis :**

- (a) Falls  $D$  die leere Matrix ist, so ist die Behauptung natürlich richtig. Wurde eine Zeile  $i$  von  $D$  in Schritt 3 (a) definiert, so ist  $d_{ij} = 0$ , da  $a_{p(i)j} = 0$ . Andernfalls gilt  $d_{ij} = a_{tj}a_{sj} - a_{sj}a_{tj} = 0$  nach Konstruktion. Das heißt, die  $j$ -te Spalte  $D_{.j}$  von  $D$  ist der Nullvektor. Man beachte auch, dass im Falle  $N = \emptyset$  oder  $P = \emptyset$ ,  $D$  nur aus den Zeilen  $A_i$  von  $A$  besteht mit  $i \in Z$ .

(b) ist offensichtlich nach Definition und impliziert (c).

(d) Sei  $\bar{x} \in \mathbb{K}^n$  mit  $\bar{x}_j = 0$  ein beliebiger Vektor.

(d<sub>1</sub>) Angenommen  $\bar{x} \in P(D, d)$ . Wir zeigen zunächst, dass  $L \leq U$  gilt.

Ist  $P = \emptyset$  oder  $N = \emptyset$ , so gilt offensichtlich  $L \leq U$  nach Definition.

Wir können also annehmen, dass  $P \neq \emptyset$ ,  $N \neq \emptyset$  gilt. Es seien  $s \in N$ ,  $t \in P$  und  $q \in R$  so gewählt, dass  $p(q) = (s, t)$  und  $\lambda_s = L$ ,  $\lambda_t = U$  gelten.

Dann gilt

$$a_{tj}A_s\bar{x} - a_{sj}A_t\bar{x} = D_q\bar{x} \leq d_q = a_{tj}b_s - a_{sj}b_t,$$

woraus folgt

$$a_{sj}(b_t - A_t\bar{x}) \leq a_{tj}(b_s - A_s\bar{x})$$

und somit (wegen  $a_{tj} > 0$  und  $a_{sj} < 0$ )

$$U = \lambda_t = \frac{1}{a_{tj}}(b_t - A_t\bar{x}) \geq \frac{1}{a_{sj}}(b_s - A_s\bar{x}) = \lambda_s = L.$$

Wir zeigen nun  $A_i(\bar{x} + \lambda e_j) \leq b_i$  für alle  $\lambda \in [L, U]$  und alle  $i \in M = P \cup N \cup Z$ , wobei  $e_j$  den  $j$ -ten Einheitsvektor bezeichnet.

Ist  $i \in Z$ , so gibt es ein  $j \in R$  mit  $p(j) = i$  und  $D_j = A_i$ ,  $d_j = b_i$ . Aus  $D_j\bar{x} \leq d_j$  folgt daher die Behauptung.

Ist  $i \in P$ , dann gilt  $U < +\infty$  und somit

$$A_i(\bar{x} + \lambda e_j) = A_i\bar{x} + a_{ij}\lambda \leq A_i\bar{x} + a_{ij}U \leq A_i\bar{x} + a_{ij}\lambda_i = b_i.$$

Analog folgt die Behauptung für  $i \in N$ .

Der Beweis (d<sub>2</sub>) sei dem Leser überlassen. □

**(3.3) Folgerung.**  $P(A, b) \neq \emptyset \iff P(D, d) \neq \emptyset$ . □

Sprechweise für (3.3):

$Ax \leq b$  ist konsistent (lösbar)  $\iff Dx \leq d$  ist konsistent (lösbar).

Die Fourier-Motzkin-Elimination der  $j$ -ten Variablen kann man natürlich auch auf ein System von normalen und strikten Ungleichungen anwenden. Eine Nachvollziehung des obigen Beweises liefert:

**(3.4) Satz.** Seien  $A$  eine  $(m, n)$ -Matrix,  $b \in \mathbb{K}^m$ ,  $j \in \{1, \dots, n\}$  und  $I, J$  eine Partition der Zeilenindexmenge  $M = \{1, \dots, m\}$ . Seien  $D \in \mathbb{K}^{r \times n}$ ,  $d \in \mathbb{K}^r$  die durch Fourier-Motzkin-Elimination der Variablen  $j$  gewonnene Matrix bzw. Vektor.

Setze  $E := p^{-1}((Z \cap I) \cup ((N \times P) \cap (I \times I)))$ ,  $F := R \setminus E$ , (wobei  $p$  die in (3.1) Schritt 2 definierte Abbildung ist),

dann gilt:

Das System

$A_I x \leq b_I, A_J x < b_J$  hat eine Lösung genau dann, wenn

das System

$D_E x \leq d_E, D_F x < d_F$  eine Lösung hat.

**Beweis :** Hausaufgabe. □

Wir können nun die Fourier-Motzkin-Elimination sukzessive auf alle Spaltenindizes anwenden. Ist in einem Schritt die neu konstruierte Matrix leer, so ist das System konsistent. Wir wollen herausfinden, was passiert, wenn  $Ax \leq b$  nicht konsistent ist. Wir nehmen daher an, dass keine der sukzessiv konstruierten Matrizen leer ist.

Eliminieren wir die 1. Spalte von  $A$ , so erhalten wir nach Satz (3.2) eine  $(r_1, n)$ -Matrix  $D_1$  und einen Vektor  $d_1 \in \mathbb{K}^{r_1}$  mit

$$Ax \leq b \text{ konsistent} \iff D_1 x \leq d_1 \text{ konsistent.}$$

Die erste Spalte von  $D_1$  ist eine Nullspalte. Nun eliminieren wir die zweite Spalte von  $D_1$  und erhalten mit Satz (3.2) eine  $(r_2, n)$ -Matrix  $D_2$  und einen Vektor  $d_2 \in \mathbb{K}^{r_2}$  mit

$$D_1 x \leq d_1 \text{ konsistent} \iff D_2 x \leq d_2 \text{ konsistent}$$

Die erste und die zweite Spalte von  $D_2$  bestehen aus lauter Nullen. Wir eliminieren nun die dritte Spalte von  $D_2$  und fahren auf diese Weise fort, bis wir die  $n$ -te Spalte eliminiert haben. Insgesamt erhalten wir:

$$Ax \leq b \text{ konsistent} \iff D_1 x \leq d_1 \text{ konsistent.} \iff \dots \iff D_n x \leq d_n \text{ konsistent.}$$

Was haben wir durch diese äquivalenten Umformungen gewonnen? Nach Konstruktion hat die  $n$ -te Matrix  $D_n$  dieser Folge von Matrizen nur noch Nullspalten, ist also eine Nullmatrix. Das System  $D_n x \leq d_n$  ist also nichts anderes als  $0x \leq d_n$ , und die Konsistenz dieses letzten Ungleichungssystems ist äquivalent zur Konsistenz von  $Ax \leq b$ . Die Konsistenz von  $0x \leq d_n$  ist trivial überprüfbar, denn  $0x \leq d_n$  hat genau dann keine Lösung, wenn der Vektor  $d_n \in \mathbb{K}^{r_n}$  eine negative Komponente hat; das heißt,  $d_n \geq 0$  ist äquivalent zur Konsistenz von  $Ax \leq b$ .

Aus Satz (3.3) folgt außerdem Folgendes: Es existiert eine Matrix  $U_1 \in \mathbb{K}^{(r_1, m)}$  mit  $U_1 \geq 0$  und

$$D_1 = U_1 A \quad , \quad d_1 = U_1 b$$

Erneute Anwendung von (3.3) liefert die Existenz einer Matrix  $U_2 \in \mathbb{K}^{(r_2, r_1)}$  mit  $U_2 \geq 0$  und

$$D_2 = U_2 D_1 \quad , \quad d_2 = U_2 d_1$$

Setzen wir die Schlussweise fort, so erhalten wir eine Matrix  $U_n \in \mathbb{K}^{(r_n, r_{n-1})}$  mit  $U_n \geq 0$  und

$$0 = D_n = U_n D_{n-1} \quad , \quad d_n = U_n d_{n-1}$$

Für die Matrix

$$U := U_n \cdot \dots \cdot U_1 \in \mathbb{K}^{(r_n, m)}$$

gilt dann

$$U \geq 0 \quad , \quad 0 = UA \quad , \quad d_n = Ub$$

Daraus folgt, dass jede Zeile der Matrix  $D_n$ , welche die Nullmatrix ist, eine konische Kombination von Zeilen von  $A$  ist.

Wir haben uns oben überlegt, dass  $Ax \leq b$  genau dann nicht konsistent ist, wenn  $D_n x \leq d_n$  nicht konsistent ist. Dies ist genau dann der Fall, wenn es eine Komponente von  $d_n$  gibt, die negativ ist, anders ausgedrückt, wenn es einen Vektor  $u$  (eine Zeile der Matrix  $U$ ) gibt mit  $0^T = u^T A$  und  $u^T b < 0$ . Fassen wir diese Beobachtung zusammen.

**(3.5) Folgerung.** *Es seien  $A \in \mathbb{K}^{(m, n)}$  und  $b \in \mathbb{K}^m$ , dann gilt: Das Ungleichungssystem  $Ax \leq b$  hat genau dann keine Lösung, wenn es einen Vektor  $u \in \mathbb{K}^m$ ,  $u \geq 0$  gibt mit  $u^T A = 0^T$  und  $u^T b < 0$ .*

□

Diese (einfache) Konsequenz aus unserem Eliminationsverfahren (3.1) ist eines der nützlichsten Resultate der Polyedertheorie.

**Ergänzende Abschweifung**

Das Fourier-Motzkin-Eliminationsverfahren ist ein Spezialfall eines allgemeinen Projektionsverfahrens auf lineare Teilräume des  $\mathbb{K}^n$ .

**(3.6) Definition.** Sei  $L$  ein linearer Unterraum von  $\mathbb{K}^n$ . Ein Vektor  $\bar{x} \in \mathbb{K}^n$  heißt **orthogonale Projektion** eines Vektors  $x \in \mathbb{K}^n$  auf  $L$ , falls  $\bar{x} \in L$  und  $(x - \bar{x})^T \bar{x} = 0$  gilt.

**(3.7) Bemerkung.** Sei  $B \in K^{(m,n)}$  eine Matrix mit vollem Zeilenrang,  $L = \{x \in \mathbb{K}^n \mid Bx = 0\}$ , dann ist für jeden Vektor  $x \in \mathbb{K}^n$ , der Vektor

$$\bar{x} := (I - B^T(BB^T)^{-1}B)x$$

die orthogonale Projektion auf  $L$ .

Die Fourier-Motzkin-Elimination der  $j$ -ten Variablen kann man, wie wir jetzt zeigen, als orthogonale Projektion von  $P(A, b)$  auf den Vektorraum  $L = \{x \mid x_j = 0\}$  deuten.

Wir zeigen nun, wie man ein Polyeder  $P(A, b)$  auf  $L = \{y \mid c^T y = 0\}, c \neq 0$ , projiziert.

**(3.8) Algorithmus.** Orthogonale Projektion eines Polyeders auf  $L = \{y \in \mathbb{K}_n \mid c^T y = 0\}, c \neq 0$  (kurz: Projektion entlang  $c$ ).

**Input:** Ein Vektor  $c \in \mathbb{K}^n, c \neq 0$  (die Projektionsrichtung),  
eine Matrix  $A \in \mathbb{K}^{(m,n)}$ ,  
ein Vektor  $b \in \mathbb{K}^m$ .

**Output:** Eine Matrix  $D \in \mathbb{K}^{(r,n)}$  ( $r$  wird im Algorithmus berechnet),  
ein Vektor  $d \in \mathbb{K}^r$

(1) Partitioniere die Menge  $M = \{1, \dots, m\}$  der Zeilenindizes von  $A$  wie folgt

$$\begin{aligned} N &:= \{i \in M \mid A_i \cdot c < 0\}, \\ Z &:= \{i \in M \mid A_i \cdot c = 0\}, \\ P &:= \{i \in M \mid A_i \cdot c > 0\}. \end{aligned}$$

(2) Setze

$$\begin{aligned} r &:= |Z \cup N \times P|, \\ R &:= \{1, 2, \dots, r\}, \text{ und} \\ p &: R \rightarrow Z \cup N \times P \text{ sei eine Bijektion.} \end{aligned}$$

(3) Für  $i = 1, 2, \dots, r$  führe aus:

(a) Ist  $p(i) \in Z$ , dann setze

$$D_{i\cdot} := A_{p(i)\cdot}, \quad d_i = b_{p(i)}.$$

(D. h. die  $i$ -te Zeile von  $D$  ist gleich der  $p(i)$ -ten Zeile von  $A$ .)

(b) Ist  $p(i) = (s, t) \in N \times P$ , dann setze

$$\begin{aligned} D_{i\cdot} &:= (A_t \cdot c)A_{s\cdot} - (A_s \cdot c)A_{t\cdot}, \\ d_i &:= (A_t \cdot c)b_s - (A_s \cdot c)b_t. \end{aligned}$$

(4) Gib  $D$  und  $d$  aus.

**Ende.**

□

**(3.9) Satz.** Seien  $D \in \mathbb{K}^{(r,n)}$  und  $d \in \mathbb{K}^r$  die in (3.1) konstruierte Matrix bzw. Vektor. Dann gilt

(a) Die Zeilen von  $D$  sind orthogonal zu  $c$ , d.h., die Zeilenvektoren von  $D$  sind in  $L = \{y \in \mathbb{K}_n \mid c^T y = 0\}$  enthalten und konische Kombinationen der Zeilen von  $A$ .

(b)  $P(D, d) \cap L$  ist die orthogonale Projektion von  $P(A, b)$  auf  $L$ .

**(3.10) Satz.** Die Projektion eines Polyeders ist ein Polyeder!

Durch sukzessive Projektion kann man ein Polyeder auf einen beliebigen Unterraum  $\{y \in \mathbb{K}_n \mid By = 0\}$  des  $\mathbb{K}_n$  projizieren.

## 3.2 Lineare Optimierung und Fourier-Motzkin-Elimination

Die Fourier-Motzkin-Elimination ist nicht nur eine Projektionsmethode, man kann sie auch zur Lösung linearer Programme verwenden. Dies geht wie folgt.

Wir beginnen mit einem linearen Programm der Form

$$\begin{aligned} \max \quad & c^T x \\ \text{s.t.} \quad & Ax \leq a \end{aligned}$$

mit  $A \in \mathbb{K}^{m \times n}$ ,  $a \in \mathbb{K}^m$  und setzen  $P = P(A, b)$ . Wir führen nun eine zusätzliche Variable  $x_{n+1}$  ein und fügen zur Matrix  $A$  eine zusätzliche Spalte und Zeile hinzu. Diese Matrix nennen wir  $B$ . Die  $(m+1)$ -te Zeile von  $B$  enthält in den ersten  $n$  Spalten den Vektor  $c$ , d. h.  $b_{m+1,j} = c_j$ ,  $j = 1, \dots, n$ ; wir setzen  $b_{i,n+1} = 0$  für  $i = 1, \dots, m$  und  $b_{m+1,n+1} = -1$ . Wir verlängern den Vektor  $a \in \mathbb{K}^m$  um eine Komponente, die den Wert 0 enthält, und nennen diesen  $(m+1)$ -Vektor  $b$ :

$$B = \begin{pmatrix} A & 0 \\ c^T & -1 \end{pmatrix}, \quad b = \begin{pmatrix} a \\ 0 \end{pmatrix}$$

Wenden wir die Fourier-Motzkin-Elimination auf  $Ax \leq a$  an, so können wir feststellen, ob die Lösungsmenge  $P$  des linearen Programms leer ist oder nicht. Führen wir die Fourier-Motzkin-Elimination mit  $Bx \leq b$  durch und eliminieren wir nur die ersten  $n$  Variablen, so erhalten wir am Ende ein System  $D_n x \leq d_n$ , welches nichts anderes ist als ein Ungleichungssystem der Form  $\alpha_i \leq x_{n+1} \leq \beta_j$ , wobei die  $\beta_j$  die positiven und die  $\alpha_i$  die negativen Einträge des Vektors  $d_n$  sind. Das Minimum über die Werte  $\beta_j$  liefert den maximalen Wert der Zielfunktion  $c^T x$  des gegebenen linearen Programms.

Setzen wir  $\bar{x}_i := 0$ ,  $i = 1, \dots, n$  und  $\bar{x}_{n+1} := \min\{\beta_j\}$ , so können wir mit diesem Vektor  $\bar{x} = (\bar{x}_i)$  beginnend durch Rücksubstitution – wie in Satz (3.2) (d) beschrieben – eine Optimallösung von  $\max c^T x$ ,  $Ax \leq b$  berechnen.

Wenn wir statt des Maximums von  $c^T x$  über  $P$  das Minimum finden wollen, so fügen wir bei der Konstruktion der Matrix  $B$  statt der Zeile  $(c^T, -1)$  die Zeile  $(-c^T, 1)$  ein. Wollen wir das Maximum und Minimum gleichzeitig bestimmen, so führen wir die Fourier-Motzkin-Elimination der ersten  $n$  Variablen auf der Matrix durch, die um die beiden Zeilen  $(c^T, -1)$  und  $(-c^T, 1)$  erweitert wurde.

Dieses Verfahren zur Lösung linearer Programme ist konzeptionell extrem einfach und klar. Wenn es in der Praxis funktionieren würde, könnten wir das Thema „lineare Optimierung“ beenden. Als Hausaufgabe haben Sie die Fourier-Motzkin-Elimination implementiert. Versuchen Sie einmal, damit lineare Programme „anständiger Größenordnung“ zu lösen. Sie werden sich wundern und wissen dann, warum die Vorlesung weitergeht.

Wir beschließen die Diskussion der Fourier-Motzkin-Elimination als Methode zur Lösung linearer Programme durch numerische Berechnung eines Beispiels.

**(3.11) Beispiel.**

Wir beginnen mit folgendem System von 5 Ungleichungen mit 2 Variablen

$$\begin{array}{rclcl}
 (1) & & -x_2 & \leq & 0 \\
 (2) & -x_1 & -x_2 & \leq & -1 \\
 (3) & -x_1 & +x_2 & \leq & 3 \\
 (4) & +x_1 & & \leq & 3 \\
 (5) & +x_1 & +2x_2 & \leq & 9
 \end{array}$$

Wir wollen die Zielfunktion

$$x_1 + 3x_2$$

über der Menge der zulässigen Lösungen sowohl maximieren als auch minimieren. Dazu schreiben wir (wie oben erläutert) das folgende Ungleichungssystem

$$\begin{array}{rclcl}
 (1) & & -x_2 & \leq & 0 \\
 (2) & -x_1 & -x_2 & \leq & -1 \\
 (3) & -x_1 & +x_2 & \leq & 3 \\
 (4) & +x_1 & & \leq & 3 \\
 (5) & +x_1 & +2x_2 & \leq & 9 \\
 (6) & +x_1 & +3x_2 & -x_3 & \leq & 0 \\
 (7) & -x_1 & -3x_2 & +x_3 & \leq & 0
 \end{array}$$

auf und eliminieren die Variable  $x_1$ . Das Ergebnis ist das folgende Ungleichungssystem:

$$\begin{array}{rclcl}
 (1) & (1) & -x_2 & \leq & 0 \\
 (2, 4) & (2) & -x_2 & \leq & 2 \\
 (2, 5) & (3) & +x_2 & \leq & 8 \\
 (2, 6) & (4) & +2x_2 & -x_3 & \leq & -1 \\
 (3, 4) & (5) & +x_2 & \leq & 6 \\
 (3, 5) & (6) & +x_2 & \leq & 4 \\
 (3, 6) & (7) & +4x_2 & -x_3 & \leq & 3 \\
 (7, 4) & (8) & -3x_2 & +x_3 & \leq & 3 \\
 (7, 5) & (9) & -x_2 & +x_3 & \leq & 9
 \end{array}$$

Wir haben oben die Variable  $x_1$  weggelassen. Die erste Spalte zeigt an, woher die neue Ungleichung kommt. So ist z. B. die 5. Ungleichung aus der Kombination der 3. und 4. Ungleichung des vorhergehenden Systems entstanden. Ungleichungen der Form  $0x_2 + 0x_3 \leq \alpha$  haben wir weggelassen. Eine solche ist z. B. die Ungleichung , die aus der Kombination der 7. und 6. Ungleichung entsteht.

Die Elimination der Variablen  $x_2$  liefert nun analog

$$\begin{array}{llll} (1,4) & (1) & -x_3 & \leq -1 \\ (1,7) & (2) & -x_3 & \leq 3 \\ (2,4) & (3) & -x_3 & \leq 3 \\ (2,7) & (4) & -x_3 & \leq 11 \\ (8,3) & (5) & +x_3 & \leq 27 \\ (8,4) & (6) & -x_3 & \leq 3 \\ (8,5) & (7) & +x_3 & \leq 21 \\ (8,6) & (8) & +x_3 & \leq 15 \\ (8,7) & (9) & +x_3 & \leq 21 \\ (9,3) & (10) & +x_3 & \leq 17 \\ (9,4) & (11) & +x_3 & \leq 17 \\ (9,5) & (12) & +x_3 & \leq 15 \\ (9,6) & (13) & +x_3 & \leq 13 \\ (9,7) & (14) & +3x_3 & \leq 39 \end{array}$$

Die größte untere Schranke für  $x_3$  ist 1. Sie wird von der ersten Ungleichung geliefert. Die kleinste obere Schranke hat den Wert 13, dieser folgt aus den Ungleichungen 13 und 14. Damit ist 13 der Maximalwert der Zielfunktion, 1 ist der Minimalwert. Setzen wir  $x_3 = 13$  in das vorhergehende Ungleichungssystem ein, so ergibt sich der Wert  $x_2 = 4$ . Durch Substitution dieser beiden Werte in das Ausgangssystem erhalten wir  $x_1 = 1$ . Der maximale Zielfunktionswert wird also im zulässigen Punkt  $x_1 = 1, x_2 = 4$  angenommen.  $\square$

# Kapitel 4

## Das Farkas-Lemma und seine Konsequenzen

Wir werden zunächst Folgerung (3.5) auf verschiedene Weisen formulieren, um sie für den späteren Gebrauch “aufzubereiten”. Dann werden wir einige interessante Folgerungen daraus ableiten.

### 4.1 Verschiedene Versionen des Farkas-Lemmas

**(4.1) (allgemeines) Farkas-Lemma.** Für dimensionsverträgliche Matrizen  $A, B, C, D$  und Vektoren  $a, b$  gilt:

$$\left. \begin{array}{l} \text{Es existieren } x, y \text{ mit} \\ Ax + By \leq a \\ Cx + Dy = b \\ x \geq 0 \end{array} \right\} \dot{\vee} \left\{ \begin{array}{l} \text{Es existieren } u, v \text{ mit} \\ u^T A + v^T C \geq 0^T \\ u^T B + v^T D = 0^T \\ u \geq 0 \\ u^T a + v^T b < 0. \end{array} \right.$$

(Hierbei bezeichnet “ $\dot{\vee}$ ” das “entweder oder”, d. h. eine der beiden Aussagen gilt, aber niemals beide gleichzeitig, also eine Alternative.)

**Beweis :** Durch die Transformationsregeln (2.4) führen wir die obige Aussage auf Folgerung (3.5) zurück. Die linke Aussage der Alternative lässt sich schreiben als

“ $\exists z$  mit  $\bar{A}z \leq \bar{a}$ ” wobei

$$\bar{A} := \begin{pmatrix} A & B \\ C & D \\ -C & -D \\ -I & 0 \end{pmatrix} \quad \text{und} \quad \bar{a} = \begin{pmatrix} a \\ b \\ -b \\ 0 \end{pmatrix}.$$

Nach (3.5) hat dieses System genau dann keine Lösung, wenn gilt:

$$\exists y^T = (u^T, \bar{v}^T, \overline{\bar{v}}^T, w^T) \geq 0$$

mit  $y^T \bar{A} = 0^T$  und  $y^T \bar{a} < 0$ . Ausführlich geschrieben heißt dies:

$$\exists \begin{pmatrix} u \\ \bar{v} \\ \overline{\bar{v}} \\ w \end{pmatrix} \geq 0 \quad \text{mit} \quad \begin{array}{l} u^T A + \bar{v}^T C - \overline{\bar{v}}^T C - w^T = 0^T \\ u^T B + \bar{v}^T D - \overline{\bar{v}}^T D = 0^T \\ u^T a + \bar{v}^T b - \overline{\bar{v}}^T b < 0. \end{array}$$

Setzen wir  $v := \bar{v} - \overline{\bar{v}}$  und betrachten wir  $w$  als einen Vektor von Schlupfvariablen, so lässt sich dieses System äquivalent schreiben als:

$$\exists u, v \quad \text{mit} \quad \begin{array}{l} u^T A + v^T C \geq 0^T \\ u^T B + v^T D = 0^T \\ u \geq 0 \\ u^T a + v^T b < 0, \end{array}$$

und dies ist das gewünschte Resultat.  $\square$

Durch Spezialisierung von (4.1) erhalten wir:

**(4.2) (spezielles) Farkas-Lemma.** Es seien  $A \in \mathbb{K}^{(m,n)}$  und  $b \in \mathbb{K}^m$ , dann gilt:

- (a)  $\exists x$  mit  $Ax \leq b$   $\dot{\vee}$   $\exists u \geq 0$  mit  $u^T A = 0^T$  und  $u^T b < 0$ .
- (b)  $\exists x \geq 0$  mit  $Ax \leq b$   $\dot{\vee}$   $\exists u \geq 0$  mit  $u^T A \geq 0^T$  und  $u^T b < 0$ .
- (c)  $\exists x \geq 0$  mit  $Ax = b$   $\dot{\vee}$   $\exists u$  mit  $u^T A \geq 0^T$  und  $u^T b < 0$ .
- (d)  $\exists x$  mit  $Ax = b$   $\dot{\vee}$   $\exists u$  mit  $u^T A = 0^T$  und  $u^T b < 0$ .

$\square$

Das Farkas-Lemma (4.2) in seiner Version (a) charakterisiert also

- (a) Die Lösbarkeit eines linearen Ungleichungssystems, in seiner Version (4.2) (b)

- (b) *die nichtnegative Lösbarkeit eines linearen Ungleichungssystems*,  
in seiner Version (4.2) (c)
- (c) *die nichtnegative Lösbarkeit eines linearen Gleichungssystems*, in seiner  
Version (4.2) (d)
- (d) *die Lösbarkeit eines linearen Gleichungssystems*,

wobei (d) natürlich bereits aus der linearen Algebra bekannt ist. Die zweite Alternative in (d) ist nichts anderes als eine Umformulierung von  $\text{rang}(A) \neq \text{rang}(A, b)$ .

Die linken Alternativen in (4.1) und (4.2) (a), (b), (c), (d) sagen aus, dass gewisse Polyeder nicht leer sind. Die Lösungsmengen der rechten Alternativen sind dagegen keine Polyeder, weil jeweils eine strikte Ungleichung in den Gleichungs- und Ungleichungssystemen vorkommt. Da in allen Fällen die rechten Seiten der rechten Alternativen Nullvektoren sind, sind die Lösungsmengen (ohne die strikte Ungleichung) nach (2.6) Kegel. Folglich können die Lösungsvektoren skaliert werden. Aus dieser Beobachtung folgt:

**(4.3) Farkas-Lemma (polyedrische Version).** *Es seien  $A \in \mathbb{K}^{(m,n)}$  und  $b \in \mathbb{K}^m$ , dann gilt:*

$$(a) P(A, b) \neq \emptyset \quad \dot{\vee} \quad P^= \left( \begin{pmatrix} A^T \\ b^T \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right) \neq \emptyset.$$

$$(b) P^+(A, b) \neq \emptyset \quad \dot{\vee} \quad P^+ \left( \begin{pmatrix} -A^T \\ b^T \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right) \neq \emptyset.$$

$$(c) P^=(A, b) \neq \emptyset \quad \dot{\vee} \quad P \left( \begin{pmatrix} -A^T \\ b^T \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right) \neq \emptyset.$$

□

Polyedrische Formulierungen von (4.2) (b), (d) und (4.1) seien dem Leser überlassen.

## 4.2 Alternativ- und Transpositionssätze

Wir haben das Farkas-Lemma nun in verschiedenen Versionen formuliert. Folgerung (3.5) beschreibt ein zur Unlösbarkeit eines Ungleichungssystems äquivalentes Kriterium. Die Sätze (4.1) und (4.2) haben die Form von Alternativen,

während in (4.3) die Lösbarkeit eines Systems mit der Unlösbarkeit eines durch Transposition gewonnenes System festgestellt wird. Man nennt deshalb Sätze des Typs (4.1), (4.2) bzw. (4.3) Alternativ- bzw. Transpositionssätze. Die Literatur ist außerordentlich reich an derartigen Sätzen, da sie — wie wir später noch sehen werden — auf vielerlei Weise nützlich sind. Wir wollen hier noch einige weitere Sätze dieser Art angeben.

**(4.4) Satz.** *Es seien  $A \in \mathbb{K}^{(p,n)}$ ,  $B \in \mathbb{K}^{(q,n)}$ ,  $a \in \mathbb{K}^p$ ,  $b \in \mathbb{K}^q$ , dann gilt genau eine der beiden folgenden Alternativen:*

- (a)  $\exists x \in \mathbb{K}^n$  mit  $Ax \leq a$ ,  $Bx < b$
- (b)  $(b_1) \quad \exists u \in \mathbb{K}_+^p, v \in \mathbb{K}_+^q \setminus \{0\}$  mit  $u^T A + v^T B = 0^T$ ,  $u^T a + v^T b \leq 0$ , oder  
 $(b_2) \quad \exists u \in \mathbb{K}_+^p$ , mit  $u^T A = 0^T$ ,  $u^T a < 0$ .

**Beweis :** Angenommen (a) und  $(b_2)$  gelten gleichzeitig, dann gibt es also einen Vektor  $x$  mit  $Ax \leq a$  und einen Vektor  $u \in \mathbb{K}_+^p$  mit  $u^T A = 0^T$ ,  $u^T a < 0$ , was (4.2) (a) widerspricht. Angenommen (a) und  $(b_1)$  gelten gleichzeitig, dann gilt:  $0^T x = (u^T A + v^T B)x = u^T (Ax) + v^T (Bx) < u^T a + v^T b \leq 0$ , ein Widerspruch.

Wir nehmen nun an, dass (a) nicht gilt und wollen zeigen, dass dann (b) gilt. Wir wenden die Fourier-Motzkin-Elimination  $n$ -mal iterativ auf  $\begin{pmatrix} A \\ B \end{pmatrix}$  und  $\begin{pmatrix} a \\ b \end{pmatrix}$  an. Nach (3.4) und erhalten wir nach  $n$  Schritten Matrizen  $C$ ,  $D$  und Vektoren  $c$ ,  $d$ , so dass  $Ax \leq a$ ,  $Bx < b$  genau dann lösbar ist, wenn  $Cx \leq c$ ,  $Dx < d$  lösbar ist, wobei  $C = 0$ ,  $D = 0$  gilt. Nach Annahme gilt (a) nicht, also muss es einen Zeilenindex  $i$  von  $C$  geben mit  $c_i < 0$  oder einen Zeilenindex  $j$  von  $D$  mit  $d_j \leq 0$ .

Wendet man die Fourier-Motzkin-Elimination nur auf  $A$ ,  $a$  und  $n$ -mal sukzessiv an, so erhält man nach (3.4) das System  $C$ ,  $c$ . Das heißt, die Lösbarkeit von  $Ax \leq a$  ist äquivalent zur Lösbarkeit von  $Cx \leq c$ . Ist also  $Ax \leq a$  nicht lösbar, so gibt es nach (3.5) einen Vektor  $u \geq 0$  mit  $u^T A = 0^T$ ,  $u^T a < 0$ , das heißt  $(b_2)$  gilt. Ist  $Ax \leq a$  lösbar, so gibt es keinen Vektor  $u \geq 0$  mit  $u^T A = 0^T$ ,  $u^T a < 0$ , d. h.  $c_i \geq 0$  für alle  $i$ . Folglich muss  $d_j \leq 0$  für ein  $j$  gelten.  $D_j \cdot = 0^T$  ist eine konische Kombination von Zeilen von  $A$  und Zeilen von  $B$ . Nach Definition, siehe (3.4), muss dabei mindestens eine Zeile von  $B$  mit einem positiven Multiplikator beteiligt sein, also gibt es  $u \in \mathbb{K}_+^p$ ,  $v \in \mathbb{K}_+^q \setminus \{0\}$  mit  $u^T A + v^T B = D_j \cdot = 0^T$  und  $u^T a + v^T b = d_j \leq 0$ , das heißt  $(b_1)$  gilt.  $\square$

**(4.5) Folgerung.** Ist  $P(A, a) \neq \emptyset$ , dann gilt genau eine der beiden folgenden Alternativen:

(a)  $\exists x \in \mathbb{K}^n$  mit  $Ax \leq a, Bx < b$ ,

(b)  $\exists \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{K}_+^{p+q}, v \neq 0$  mit  $u^T A + v^T B = 0^T$  und  $u^T a + v^T b \leq 0$ .

□

**(4.6) Satz** (Gordan (1873)). *Es gilt genau eine der beiden folgenden Alternativen:*

(a)  $\exists x$  mit  $Ax < 0$ ,

(b)  $\exists u \geq 0, u \neq 0$  mit  $u^T A = 0^T$ .

**Beweis :** Setze in (4.5)  $B := A$  (wobei  $A$  die obige Matrix ist),  $b := 0, A := 0, a := 0$ . □

**(4.7) Satz** (Stiemke (1915)). *Es gilt genau eine der folgenden Alternativen:*

(a)  $\exists x > 0$  mit  $Ax = 0$ ,

(b)  $\exists u$  mit  $u^T A \geq 0^T, u^T A \neq 0$ .

**Beweis :** Setze in (4.5)  $B := -I, b = 0, A := \begin{pmatrix} A \\ -A \end{pmatrix}, a := 0$ . □

Der Satz von Stiemke charakterisiert also die strikt positive Lösbarkeit eines homogenen Gleichungssystems, während der Satz von Gordan die semipositive Lösbarkeit des homogenen Gleichungssystems  $u^T A = 0$  kennzeichnet. Eine Sammlung weiterer Alternativsätze kann man in

O. L. Mangasarian, “Nonlinear Programming”, McGraw-Hill, New York, 1969 (80 QH 400 M 277)

finden.

### 4.3 Trennsätze

So genannte Trennsätze spielen in der Konvexitätstheorie und — in noch allgemeinerer Form — in der Funktionalanalysis eine wichtige Rolle. Diese Sätze sind

i. a. vom folgenden Typ: Seien  $S$  und  $T$  zwei Mengen, dann gibt es unter gewissen Voraussetzungen eine Hyperebene  $H$ , die  $S$  und  $T$  trennt. Geometrisch heißt dies, dass  $S$  auf der einen,  $T$  auf der anderen Seite von  $H$  liegt. Genauer: Sei  $H = \{x | c^T x = \gamma\}$ ,  $c \neq 0$ , eine Hyperebene, dann sagen wir, dass  $H$  zwei Mengen  $P$  und  $Q$  **trennt**, falls  $P \cup Q \not\subseteq H$  und  $P \subseteq \{x | c^T x \leq \gamma\}$ ,  $Q \subseteq \{x | c^T x \geq \gamma\}$  gilt. Wir sagen, dass  $H$  die Mengen  $P$  und  $Q$  **strikt trennt**, falls  $P \subseteq \{x | c^T x < \gamma\}$  und  $Q \subseteq \{x | c^T x > \gamma\}$ .

Das Farkas-Lemma impliziert einige interessante Trennsätze für Polyeder  $P \subseteq \mathbb{K}^n$ . Wir wollen hier einen Satz dieser Art angeben, der als Prototyp für die allgemeineren Trennsätze angesehen werden kann.

**(4.8) Trennsatz.** *Es seien  $P = P(A, a)$  und  $Q = P(B, b)$  zwei Polyeder im  $\mathbb{K}^n$ . Es gibt eine  $P$  und  $Q$  strikt trennende Hyperebene genau dann, wenn  $P \cap Q = \emptyset$  gilt, es sei denn, eines der Polyeder ist leer und das andere der  $\mathbb{K}^n$ .*

**Beweis :** Gibt es eine  $P$  und  $Q$  strikt trennende Hyperebene, dann ist  $P \cap Q$  offensichtlich leer.

Sei  $P \cap Q = \emptyset$ . Zu zeigen ist die Existenz eines Vektors  $c \in \mathbb{K}^n$ ,  $c \neq 0$ , und einer Zahl  $\gamma \in \mathbb{K}$ , so dass  $P \subseteq \{x | c^T x < \gamma\}$  und  $Q \subseteq \{x | c^T x > \gamma\}$  gilt. Sind  $P$  und  $Q$  leer, so leistet jede beliebige Hyperebene das Gewünschte. Ist eines der Polyeder leer, sagen wir  $Q = \emptyset$ ,  $P \neq \emptyset$ , dann ist wegen  $P \neq \mathbb{K}^n$  eine der Zeilen von  $A$  von Null verschieden, sagen wir  $A_i. \neq 0$ . Dann haben  $c^T := A_i.$ ,  $\gamma := a_i + 1$  die geforderten Eigenschaften. Sind  $P$  und  $Q$  nicht leer, dann gilt nach Voraussetzung

$$P \left( \begin{pmatrix} A \\ B \end{pmatrix}, \begin{pmatrix} a \\ b \end{pmatrix} \right) = \emptyset.$$

Folglich existiert nach (4.2) (a) ein Vektor  $\begin{pmatrix} u \\ v \end{pmatrix} \geq 0$  mit  $u^T A + v^T B = 0^T$  und  $u^T a + v^T b < 0$ . Wir setzen

$$c^T := u^T A (= -v^T B) \quad \text{und} \quad \gamma := \frac{1}{2}(u^T a - v^T b).$$

Daraus folgt:

$$x \in P \Rightarrow c^T x = u^T A x \leq u^T a < u^T a - \frac{1}{2}(u^T a + v^T b) = \frac{1}{2}(u^T a - v^T b) = \gamma.$$

$$x \in Q \Rightarrow c^T x = -v^T B x \geq -v^T b > -v^T b + \frac{1}{2}(u^T a + v^T b) = \frac{1}{2}(u^T a - v^T b) = \gamma.$$

Der Vektor  $c$  ist offenbar von Null verschieden.  $\square$

Der folgende Satz, den wir nicht mit den bisher entwickelten Methoden beweisen können, soll als Beispiel für die oben angesprochene allgemeine Art von Trennsätzen dienen.

**(4.9) Trennsatz für konvexe Mengen.** *Es seien  $P$  und  $Q$  zwei nichtleere konvexe Teilmengen des  $\mathbb{R}^n$ . Es gibt eine  $P$  und  $Q$  trennende Hyperebene  $H = \{x \in \mathbb{R}^n \mid c^T x = \gamma\}$  (d. h.  $P \cup Q \not\subseteq H$ ,  $P \subseteq \{x \in \mathbb{R}^n \mid c^T x \leq \gamma\}$ ,  $Q \subseteq \{x \in \mathbb{R}^n \mid c^T x \geq \gamma\}$ ) genau dann, wenn der Durchschnitt des relativ Inneren von  $P$  mit dem relativ Inneren von  $Q$  leer ist.*

**Beweis :** Siehe Stoer & Witzgall (1970), Seite 98, Satz (3.3.9) oder Leichtweiss, "Konvexe Mengen", VEB Deutscher Verlag der Wissenschaften, Berlin, 1980, Seite 28, Satz 3.1.  $\square$

Der aufmerksame Leser wird bemerkt haben, dass in (4.9) der Vektorraum  $\mathbb{R}^n$  (und nicht wie üblich  $\mathbb{K}^n$ ) benutzt wurde. In der Tat ist Satz (4.9) für konvexe Teilmengen von  $\mathbb{Q}^n$  im folgenden Sinne falsch. Es ist nicht immer möglich, zwei konvexe Mengen in  $\mathbb{Q}^n$ , deren relativ innere Mengen kein gemeinsames Element besitzen, durch eine Hyperebene  $c^T x = \gamma$  zu trennen, so dass der Vektor  $(c^T, \gamma)^T$  rational ist. Man betrachte nur  $P = \{x \in \mathbb{Q} \mid x < \sqrt{2}\}$ ,  $Q = \{x \in \mathbb{Q} \mid x > \sqrt{2}\}$ . Dieses Beispiel zeigt auch, dass die bisher entwickelten konstruktiven Beweistechniken, die ja für  $\mathbb{Q}$  und  $\mathbb{R}$  gleichermaßen arbeiten, nicht mächtig genug sind, um (4.9) abzuleiten. Offenbar muss das Vollständigkeitsaxiom der reellen Zahlen auf irgendeine Weise benutzt werden.

Satz (4.9) können wir allerdings für Polyeder beweisen. Zu seiner Formulierung müssen wir jedoch schärfere Anforderungen an die Darstellung der involvierten Polyeder stellen.

#### **(4.10) Schwacher Trennsatz.**

*Seien  $P = \{x \in \mathbb{K}^n \mid Ax = a, Bx \leq b\}$ ,  $Q = \{x \in \mathbb{K}^n \mid Cx = c, Dx \leq d\}$  zwei nichtleere Polyeder, so dass keine der Ungleichungen  $Bx \leq b$  von allen Punkten in  $P$  und keine der Ungleichungen  $Dx \leq d$  von allen Punkten in  $Q$  mit Gleichheit erfüllt wird. Es gibt eine  $P$  und  $Q$  trennende Hyperebene genau dann, wenn  $R := \{x \in \mathbb{K}^n \mid Ax = a, Cx = c, Bx < b, Dx < d\}$  leer ist.*

**Beweis :** Ist  $P \cap Q = \emptyset$ , so ist  $R = \emptyset$ , und die Behauptung folgt (in schärferer Form) aus Satz (4.8). Sei also  $P \cap Q \neq \emptyset$ .

Die Menge  $R$  ist offenbar die Lösungsmenge des Systems

$$Ax \leq a, -Ax \leq -a, Cx \leq c, -Cx \leq -c, Bx < b, Dx < d.$$

Dieses hat nach Folgerung (4.5) genau dann keine Lösung, wenn es einen Vektor  $(\bar{u}^T, \bar{u}^T, \bar{v}^T, \bar{v}^T, w^T, y^T)$  mit  $(w^T, y^T) \neq 0^T$  gibt mit der Eigenschaft

$$\bar{u}^T A - \underline{\bar{u}}^T A + \bar{v}^T C - \underline{\bar{v}}^T C + w^T B + y^T D = 0^T,$$

$$\bar{u}^T a - \underline{\bar{u}}^T a + \bar{v}^T c - \underline{\bar{v}}^T c + w^T b + y^T d \leq 0.$$

Wir setzen  $u := \bar{u} - \underline{\bar{u}}, v := \bar{v} - \underline{\bar{v}}$  und

$$\begin{aligned} r^T &:= u^T A + w^T B (= -v^T C - y^T D), \\ \rho &:= \frac{1}{2}(u^T a + w^T b - v^T c - y^T d). \end{aligned}$$

Sei  $H =: \{x \mid r^T x = \rho\}$ , dann ist  $H$  eine  $P$  und  $Q$  trennende Hyperebene, denn

$$\begin{aligned} (1) \quad x \in P &\implies r^T x = u^T Ax + w^T Bx \leq u^T a + w^T b \\ &\leq u^T a + w^T b - \frac{1}{2}(u^T a + w^T b + v^T c + y^T d) \\ &= \rho \\ (2) \quad x \in Q &\implies r^T x = -v^T Cx - y^T Dx \geq -v^T c - y^T d \\ &\geq -v^T c - y^T d + \frac{1}{2}(u^T a + w^T b + v^T c + y^T d) \\ &= \rho. \end{aligned}$$

$R = \emptyset$  ist also äquivalent zur Existenz der Hyperebene  $H = \{x \mid r^T x = \rho\}$ . Aus der Voraussetzung folgt, dass Punkte  $p \in P$  und  $q \in Q$  existieren mit  $Bp < b$  und  $Dq < d$ . Die Abschätzungen (1), (2) zeigen, dass aus  $(w^T, y^T) \neq 0^T$  folgt, dass  $r^T p < \rho$  oder  $r^T q > \rho$  gilt. Somit gilt  $r \neq 0$  und  $P \cup Q \not\subseteq H$ ; das aber heißt,  $H$  ist eine trennende Hyperebene.

Gibt es umgekehrt eine  $P$  und  $Q$  trennende Hyperebene  $H = \{x \mid r^T x = p\}$ , dann gibt es einen Punkt in  $P \cup Q$ , der nicht in  $H$  ist; sagen wir  $P \not\subseteq H$ . Daraus folgt  $P = \{x \mid Ax = a, Bx \leq b, r^T x \leq \rho\}$  und  $\{x \mid Ax = a, Bx < b\} = \{x \mid Ax = q, Bx < b, r^T x < \rho\}$  und somit ergibt sich, dass  $R = \{x \mid Ax = a, Cx = c, Bx < b, r^T x < \rho, Dx < d\}$  leer ist, denn  $\{x \mid Cx = c, Dx < d, r^T x < \rho\} = \emptyset$ , weil  $Q = \{x \mid Cx = c, Dx \leq d\} \subseteq \{x \mid r^T x \geq \rho\}$ .  $\square$

Satz (4.10) folgt aus (4.9), wenn man sich überlegt, dass für ein Polytop  $P = \{x \mid Ax = a, Bx \leq b\}$  mit der Eigenschaft, dass keine der Ungleichungen in  $Bx \leq b$  durch alle Punkte aus  $P$  mit Gleichheit erfüllt ist, das relativ Innere von  $P$  genau die Menge  $\{x \mid Ax = a, Bx < b\}$  ist.

# Kapitel 5

## Der Dualitätssatz der linearen Programmierung

In Satz (1.5) haben wir bereits den sogenannten schwachen Dualitätssatz angegeben und bewiesen. Die Gleichheit der optimalen Zielfunktionswerte kann man (unter anderem) mit Hilfe des Farkas-Lemmas beweisen. Man kann den (nachfolgend formulierten) Dualitätssatz als eine Optimierungsversion des Farkas-Lemmas ansehen. Beide Sätze sind in dem Sinne äquivalent, dass der eine ohne Mühe aus dem anderen gefolgert werden kann und umgekehrt. So wie das Farkas Lemma für die Polyedertheorie von zentraler Bedeutung ist, ist der Dualitätssatz von außerordentlicher Tragweite in der linearen Optimierung und zwar sowohl in algorithmisch-praktischer als auch in theoretischer Hinsicht. Wenn man heutzutage in mathematischen (Forschungs-)Aufsätzen einen Satz wie “. . . it follows by an LP-argument . . .” liest, dann heißt das so gut wie immer, dass der Dualitätssatz in einer seiner Versionen geschickt anzuwenden ist. Wegen seiner Bedeutung widmen wir dem Dualitätssatz ein eigenes Kapitel. (Er hätte natürlich auch in Kapitel 4 eingearbeitet werden können.)

Wir wollen zunächst das folgende lineare Programm

$$(5.1) \quad \begin{aligned} \max \quad & c^T x \\ & Ax \leq b \end{aligned}$$

wobei  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$  und  $c \in \mathbb{K}^n$  gegeben sind, untersuchen. Unsere Aufgabe besteht darin, unter den Vektoren  $x \in \mathbb{K}^n$ , die **zulässig** für (5.1) sind, d. h. die  $Ax \leq b$  erfüllen, einen Vektor  $x^* \in \mathbb{K}^n$  zu finden, der **optimal** ist, d. h.  $c^T x^* \geq c^T x$  für alle zulässigen  $x$ . In Kapitel 1 haben wir bereits eine Motiva-

tion für die Dualitätstheorie gegeben, hier wollen wir noch einen weiteren Zugang beschreiben. Zunächst wollen wir (5.1) polyedrisch darstellen. Wir führen dazu eine neue Variable  $z$  ein und schreiben (5.1) in der Form

$$(5.2) \quad \max z$$

$$\begin{pmatrix} 1 & -c^T \\ 0 & A \end{pmatrix} \begin{pmatrix} z \\ x \end{pmatrix} \leq \begin{pmatrix} 0 \\ b \end{pmatrix}$$

Daraus folgt, dass jedes lineare Programm in ein LP umgeformt werden kann, bei dem lediglich Lösungsvektoren gesucht werden, deren erste Komponente so groß wie möglich ist. Schreiben wir (5.2) noch einmal um und bringen wir  $z$  auf die rechte Seite, so kann (5.2) folgendermaßen geschrieben werden.

$$(5.3) \quad \max z$$

$$\begin{pmatrix} -c^T \\ A \end{pmatrix} x \leq \begin{pmatrix} -z \\ b \end{pmatrix}.$$

Für festes  $z \in \mathbb{K}$  ist die Lösungsmenge von (5.3) ein Polyeder im  $\mathbb{K}^n$ . Wir setzen für alle  $z \in \mathbb{K}$

$$(5.4) \quad P_z := \left\{ x \in \mathbb{K}^n \mid \begin{pmatrix} -c^T \\ A \end{pmatrix} x \leq \begin{pmatrix} -z \\ b \end{pmatrix} \right\}.$$

Nunmehr können wir (5.3) wie folgt darstellen:

$$(5.5) \quad \max \{z \mid P_z \neq \emptyset\}.$$

Wir suchen also ein  $z \in \mathbb{K}$ , das so groß wie möglich ist unter der Nebenbedingung, dass das Polyeder  $P_z \neq \emptyset$  ist.

Üblicherweise nennt man die Aufgabe zu zeigen, dass eine Menge nicht leer ist, ein **primales Problem**. Zu zeigen, dass eine Menge leer ist, ist ein **duales Problem**. Diese beiden Aufgaben sind i. a. nicht von gleicher Schwierigkeit. Betrachten wir z. B. ein Polyeder  $P \subseteq \mathbb{Q}^n$ . Legen wir einfach eine Abzählung  $x_1, x_2, x_3, \dots$  der Elemente von  $\mathbb{Q}^n$  fest und überprüfen wir, ob  $x_i \in P$  ist für  $i = 1, 2, \dots$ , so sind wir sicher, dass dieses Verfahren nach endlich vielen Schritten abbricht, falls  $P \neq \emptyset$  ist. Jedoch kann dieses Verfahren niemals nach endlich vielen Schritten folgern, dass  $P$  leer ist. Ein Ziel von Dualitätstheorien ist es, "gute Kriterien" für die Lösung der dualen Probleme zu finden. Das Farkas-Lemma liefert ein solches. Gilt  $P = \{x \mid Ax \leq b\}$ , so ist nach (4.3) (a)  $P$  genau dann leer,

wenn  $Q = P = \left( \begin{pmatrix} A^T \\ b^T \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right)$  nicht leer ist. Durch Anwendung des obigen Abzählverfahrens auf  $P$  und  $Q$  gleichzeitig können wir also nach endlich vielen Schritten entscheiden, ob  $P$  leer ist oder nicht.

Das angegebene Abzählverfahren sollte nicht als ein ernsthafter Vorschlag zur algorithmischen Lösung des primalen und dualen Problems angesehen werden. Wir werden später sehen, dass dies alles viel besser und schneller gemacht werden kann. Diese Bemerkungen sollten lediglich das Problembewusstsein des Lesers schärfen!

Zurück zu (5.5). Es kann sein, dass  $P_z \neq \emptyset$  gilt für alle  $z \in \mathbb{K}$ . Dann ist (5.1) unbeschränkt und hat keine Optimallösung. Andernfalls können wir versuchen den Maximalwert in (5.5) nach oben zu beschränken. Die bestmögliche Schranke ist natürlich gegeben durch

$$(5.6) \quad \inf \{ z \mid P_z = \emptyset \}.$$

Nach (4.2) (a) ist  $P_z = \emptyset$  äquivalent dazu, dass das System  $y^T A = \lambda c^T$ ,  $y^T b < \lambda z$ ,  $y \geq 0$ ,  $\lambda \geq 0$  eine Lösung hat. Man überlegt sich leicht, dass dieses System — im Falle der Lösbarkeit von (5.1) — genau dann lösbar ist, wenn  $y^T A = c^T$ ,  $y^T b < z$ ,  $y \geq 0$  lösbar ist. Wenn wir das kleinste  $z$  in (5.6) finden wollen, müssen wir also einen möglichst kleinen Wert  $y^T b$  bestimmen. (5.6) ist somit äquivalent zu

$$(5.7) \quad \begin{aligned} \min \quad & y^T b \\ & y^T A = c^T \\ & y \geq 0 \end{aligned}$$

Problem (5.7) ist offenbar wiederum ein lineares Programm. Wir nennen es das zum **primalen Problem** (5.1) **duale lineare Programm**. Wir wollen nun zeigen, dass (5.1) und (5.7) den gleichen Optimalwert haben.

**(5.8) Dualitätssatz.** *Es seien  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$  und  $c \in \mathbb{K}^n$ , dann haben die beiden linearen Programme*

$$(P) \quad \max_{Ax \leq b} c^T x \quad \text{und} \quad (D) \quad \begin{aligned} \min \quad & y^T b \\ & y^T A = c^T \\ & y \geq 0 \end{aligned}$$

*optimale Lösungen, deren Zielfunktionswerte gleich sind, genau dann, wenn sie zulässige Lösungen besitzen.*

**Beweis:** Wenn (P) und (D) optimale Lösungen haben, dann haben sie auch zulässige. Haben (P) und (D) zulässige Lösungen, so gilt nach (1.5), dass der Wert von (P) nicht größer als der von (D) ist. (P) und (D) haben zulässige Lösungen mit gleichem Zielfunktionswert genau dann, wenn das System

$$(1) \quad \begin{aligned} \begin{pmatrix} 0 \\ b^T \end{pmatrix} y + \begin{pmatrix} A \\ -c^T \end{pmatrix} x &\leq \begin{pmatrix} b \\ 0 \end{pmatrix} \\ A^T y &= c \\ y &\geq 0 \end{aligned}$$

eine Lösung hat. Dies ist nach (4.1) äquivalent dazu, dass das System

$$(2) \quad \begin{aligned} zb^T + v^T A^T &\geq 0^T \\ u^T A - zc^T &= 0^T \\ u &\geq 0 \\ z &\geq 0 \\ u^T b + v^T c &< 0 \end{aligned}$$

keine Lösung hat. Nehmen wir an, dass (2) eine Lösung  $(u^T, v^T, z)$  hat. Gibt es eine Lösung von (2) mit  $z = 0$ , so hat nach (4.1) das System

$$\begin{aligned} Ax &\leq b \\ A^T y &= c \\ y &\geq 0 \end{aligned}$$

keine Lösung. Das aber heißt, dass (P) oder (D) keine zulässige Lösung hat. Widerspruch! Gibt es eine Lösung von (2) mit  $z > 0$ , so können wir durch Skalieren eine Lösung finden mit  $z = 1$ . Daraus folgt

$$0 > u^T b + v^T c \geq -u^T A v + v^T A^T u = 0.$$

Widerspruch! Dies impliziert, dass (2) inkonsistent ist. Also hat nach (4.1) das System (1) eine Lösung, und Satz (5.8) ist bewiesen.  $\square$

**(5.9) Folgerung.**  $P$  bzw.  $D$  seien die Lösungsmengen von (P) bzw. (D). Wir setzen

$$z^* := \begin{cases} +\infty, & \text{falls (P) unbeschränkt,} \\ -\infty, & \text{falls } P = \emptyset, \\ \max\{c^T x \mid x \in P\}, & \text{andernfalls,} \end{cases}$$

$$u^* := \begin{cases} -\infty, & \text{falls (D)unbeschränkt,} \\ +\infty, & \text{falls } D = \emptyset, \\ \min\{y^T b, y \in D\}, & \text{andernfalls.} \end{cases}$$

Dann gilt

- (a)  $-\infty < z^* = u^* < +\infty \iff z^*$  endlich  $\iff u^*$  endlich.
- (b)  $z^* = +\infty \Rightarrow D = \emptyset$ .
- (c)  $u^* = -\infty \Rightarrow P = \emptyset$ .
- (d)  $P = \emptyset \Rightarrow D = \emptyset$  oder  $u^* = -\infty$ .
- (e)  $D = \emptyset \Rightarrow P = \emptyset$  oder  $z^* = +\infty$ .

**Beweis :** (a) Aus  $z^* = u^*$  endlich folgt natürlich  $u^*$  und  $z^*$  endlich.

Sei  $z^*$  endlich, dann ist  $P \neq \emptyset$  und

$$(1) \quad \begin{array}{r} -c^T x \leq -z \\ Ax \leq b \end{array}$$

ist für  $z = z^*$  lösbar, jedoch unlösbar für jedes  $z > z^*$ . Sei also  $z > z^*$ , dann ist nach (4.2) (a) das duale Problem

$$(2) \quad \begin{array}{r} -uc^T + y^T A = 0^T \\ -uz + y^T b < 0 \\ u, y \geq 0 \end{array}$$

lösbar. Hat dieses System eine Lösung mit  $u = 0$ , so hat  $y^T A = 0$ ,  $y^T b < 0$ ,  $y \geq 0$  eine Lösung. Also folgt aus (4.2) (a), dass  $P = \emptyset$ , ein Widerspruch. Folglich muss es eine Lösung von (2) geben mit  $u > 0$ . Durch Skalieren können wir eine derartige Lösung mit  $u = 1$  finden. Daraus folgt, dass das System

$$\begin{array}{r} y^T b < z \\ y^T A = c^T \\ y \geq 0 \end{array}$$

eine Lösung besitzt. Folglich ist  $D \neq \emptyset$ , und aus Satz (5.8) folgt die Behauptung.

Ist  $u^*$  endlich, so verläuft der Beweis analog.

(b) Ist  $D \neq \emptyset$ , so impliziert der schwache Dualitätssatz (1.5)  $z^* \leq \min\{y^T b \mid y \in D\} = u^* < +\infty$ .

(c) Analog zu (b).

(d) Angenommen  $P = \emptyset$  und  $u^* > -\infty$ . Dann gilt entweder  $u^* = +\infty$  und somit  $D = \emptyset$  oder nach (a) folgt  $z^* = u^*$  endlich, also  $P \neq \emptyset$ . Widerspruch!

(e) Analog zu (d). □

Man könnte vermuten, dass in (5.9) (b) und (c) auch die umgekehrte Richtung gilt. Die Aussagen (d) und (e) zeigen aber, dass beim Versuch eines Beweises der Rückrichtung Komplikationen auftreten können. Das nachfolgende Beispiel zeigt, dass es in der Tat zueinander duale lineare Programme gibt, die beide leere Lösungsmengen haben.

**(5.10) Beispiel.** Es seien

$$P = \left\{ x \in \mathbb{K}^2 \mid \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} x \leq \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right\},$$

$$D = \left\{ y \in \mathbb{K}^2 \mid \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} y = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, y \geq 0 \right\}.$$

Dann sind die linearen Programme

$$\max_{x \in P} x_1 + x_2 \quad \text{und} \quad \min_{y \in D} -y_2$$

zueinander dual, und es gilt:

(a)  $P = \emptyset$ , da nach (4.2) (a) das System  $(u_1, u_2) \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} = (0, 0)$ ,

$u_1, u_2 \geq 0, -u_2 < 0$  eine Lösung hat (z. B.  $u_1 = u_2 = 1$ ).

(b)  $D = \emptyset$ , da nach (4.2) (c) das System  $(v_1, v_2) \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \geq (0, 0)$ ,

$v_1 + v_2 < 0$ , eine Lösung hat (z. B.  $v_1 = v_2 = -1$ ). □

Wir haben bereits die Sprechweise “zueinander duale lineare Programme” benutzt, ohne dafür eine Begründung zu geben. Sie erfolgt hiermit:

**(5.11) Bemerkung.** Gegeben seien dimensionsverträgliche Matrizen  $A, B, C, D$  und Vektoren  $a, b, c, d$ , dann gilt: Das zum **primalen linearen Programm**

$$(5.12) \quad \begin{aligned} \max \quad & c^T x + d^T y \\ Ax + By & \leq a \\ Cx + Dy & = b \\ x & \geq 0 \end{aligned}$$

duale lineare Programm ist

$$(5.13) \quad \begin{aligned} \min \quad & u^T a + v^T b \\ u^T A + v^T C & \geq c^T \\ u^T B + v^T D & = d^T \\ u & \geq 0. \end{aligned}$$

Das zum linearen Programm (5.13) duale Programm ist (5.12).

**Beweis :** Wir benutzen die Transformationsregeln (2.4) und schreiben (5.12) in der Form

$$\begin{aligned} \max \quad & c^T x + d^T y \\ Ax + By & \leq a \\ Cx + Dy & \leq b \\ -Cx - Dy & \leq -b \\ -Ix + 0y & \leq 0 \end{aligned}$$

Das hierzu duale Programm ist nach Definition

$$\begin{aligned}
& \min u^T a + v_1^T b - v_2^T b \\
& u^T A + v_1^T C - v_2^T C - w^T = c^T \\
& u^T B + v_1^T D - v_2^T D = d^T \\
& u, v_1, v_2, w \geq 0.
\end{aligned}$$

Setzen wir  $v := v_1 - v_2$  und lassen wir  $w$  weg, so erhalten wir (5.13). Analog folgt, dass (5.12) dual zu (5.13) ist.  $\square$

Von nun an können wir also sagen, dass das duale Programm zum dualen Programm das primale ist, bzw. von einem Paar dualer Programme sprechen. Wir wollen nun zur Übersicht und als Nachschlagewerk in Tabelle 5.1 eine Reihe von Paaren dualer Programme auflisten. Die Korrektheit der Aussagen ergibt sich direkt aus (5.11).

primales LP	duales LP
(P <sub>1</sub> ) $\max c^T x, Ax \leq b, x \geq 0$	(D <sub>1</sub> ) $\min y^T b, y^T A \geq c^T, y \geq 0$
(P <sub>2</sub> ) $\min c^T x, Ax \geq b, x \geq 0$	(D <sub>1</sub> ) $\max y^T b, y^T A \leq c^T, y \geq 0$
(P <sub>3</sub> ) $\max c^T x, Ax = b, x \geq 0$	(D <sub>3</sub> ) $\min y^T b, y^T A \geq c^T$
(P <sub>4</sub> ) $\min c^T x, Ax = b, x \geq 0$	(D <sub>4</sub> ) $\max y^T b, y^T A \leq c^T$
(P <sub>5</sub> ) $\max c^T x, Ax \leq b$	(D <sub>5</sub> ) $\min y^T b, y^T A = c^T, y \geq 0$
(P <sub>6</sub> ) $\min c^T x, Ax \geq b$	(D <sub>6</sub> ) $\max y^T b, y^T A = c^T, y \geq 0$

Tabelle 5.1

Aus den obigen primal-dual Relationen kann man die folgenden Transformationsregeln in Tabelle 5.2 ableiten.

primal	dual
Gleichung oder Ungleichung	Variable
Ungleichung	nichtnegative Variable
Gleichung	nicht vorzeichenbeschränkte Variable
nichtnegative Variable	Ungleichung
nicht vorzeichenbeschränkte Variable	Gleichung

Tabelle 5.2

Speziell bedeutet dies z. B. bezüglich des LP  $\max c^T x, Ax \leq b, x \geq 0$

- Jeder primalen Ungleichung  $A_i x \leq b_i$  ist eine nichtnegative duale Variable  $y_i$  zugeordnet.
- Jeder primalen nichtnegativen Variablen  $x_j$  ist die duale Ungleichung  $y^T A_{.j} \geq c_j$  zugeordnet.

Aus dem Vorhergehenden dürfte klar sein, dass der Dualitätssatz nicht nur für das spezielle Paar linearer Programme (P), (D) aus Satz (5.8) gilt, sondern für alle Paare dualer linearer Programme. Um später darauf Bezug nehmen zu können, wollen wir den Dualitätssatz in voller Allgemeinheit formulieren.

**(5.14) Allgemeiner Dualitätssatz.** *Es seien (P) ein lineares Programm und (D) das zu (P) duale Programm.*

- (a) *Haben (P) und (D) zulässige Lösungen, so haben (P) und (D) Optimallösungen, und die Zielfunktionswerte aller Optimallösungen stimmen überein.*
- (b) *Hat eines der beiden Programme (P) oder (D) eine Optimallösung, so auch das andere.*
- (c) *Hat eines der Programme (P) oder (D) keine zulässige Lösung, so ist das andere unbeschränkt oder hat keine zulässige Lösung.*
- (d) *Ist eines der Programme (P) oder (D) unbeschränkt, so hat das andere keine zulässige Lösung.*

Wir wollen nun zwei weitere wichtige Sätze zur Charakterisierung von Optimallösungen linearer Programme formulieren und beweisen.

**(5.15) Satz vom schwachen komplementären Schlupf.** *Es seien  $A, B, C, D$  und  $a, b, c, d$  dimensionsverträglich und (5.12), (5.13) das zugehörige Paar dualer linearer Programme. Die Vektoren  $\begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix}$  bzw.  $\begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix}$  seien zulässig für (5.12) bzw. (5.13), dann sind die folgenden Aussagen äquivalent:*

- (a)  $\begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix}$  ist optimal für (5.12) und  $\begin{pmatrix} \bar{u} \\ \bar{v} \end{pmatrix}$  ist optimal für (5.13).
- (b)  $(c^T - (\bar{u}^T A + \bar{v}^T C))\bar{x} - \bar{u}^T(a - (A\bar{x} + B\bar{y})) = 0$ .
- (c) *Für alle Komponenten  $\bar{u}_i$  der Duallösung gilt:  $\bar{u}_i > 0 \Rightarrow A_i \bar{x} + B_i \bar{y} = a_i$ .  
Für alle Komponenten  $\bar{x}_j$  der Primallösung gilt:  $\bar{x}_j > 0 \Rightarrow \bar{u}^T A_{.j} + \bar{v}^T C_{.j} = c_j$ .*

- (d) Für alle Zeilenindizes  $i$  von  $A$  und  $B$  gilt:  $A_{i.}\bar{x} + B_{i.}\bar{y} < a_i \Rightarrow \bar{u}_i = 0$ .  
Für alle Spaltenindizes  $j$  von  $A$  und  $C$  gilt:  $\bar{u}^T A_{.j} + \bar{v}^T C_{.j} > c_j \Rightarrow \bar{x}_j = 0$ .

**Beweis :**

$$\begin{aligned} \text{(a)} &\iff c^T \bar{x} + d^T \bar{y} = \bar{u}^T a + \bar{v}^T b \quad (\text{nach (5.14)}) \\ &\iff c^T \bar{x} + (\bar{u}^T B + \bar{v}^T D) \bar{y} - \bar{u}^T a - \bar{v}^T (C \bar{x} + D \bar{y}) = 0 \\ &\iff c^T \bar{x} + \bar{u}^T B \bar{y} + \bar{v}^T D \bar{y} - \bar{u}^T a - \bar{v}^T C \bar{x} - \bar{v}^T D \bar{y} - \bar{u}^T A \bar{x} + \bar{u}^T A \bar{x} = 0 \\ &\iff c^T \bar{x} - (\bar{u}^T A + \bar{v}^T C) \bar{x} - \bar{u}^T a + \bar{u}^T (A \bar{x} + B \bar{y}) = 0 \\ &\iff \text{(b)} \end{aligned}$$

$$\text{(b)} \implies \text{(c)}$$

Es seien  $t^T := c^T - (\bar{u}^T A + \bar{v}^T C)$  und  $s := a - (A \bar{x} + B \bar{y})$ . Nach Voraussetzung gilt also  $t \leq 0$  und  $s \geq 0$ . Aus  $\bar{x} \geq 0$  und  $\bar{u} \geq 0$  folgt daher  $t^T \bar{x} \leq 0$  und  $\bar{u}^T s \geq 0$ , mithin  $t^T \bar{x} - \bar{u}^T s \leq 0$ . Es kann also  $t^T \bar{x} - \bar{u}^T s = 0$  nur dann gelten, wenn  $t^T \bar{x} = 0$  und  $\bar{u}^T s = 0$ . Daraus ergibt sich (c).

$$\text{(c)} \implies \text{(b)} \quad \text{trivial.}$$

$$\text{(c)} \iff \text{(d)} \quad \text{durch Negation.}$$

□

Aus (5.15) folgt für jedes Paar dualer linearer Programme durch Spezialisierung die Gültigkeit eines Satzes vom schwachen komplementären Schlupf.

**(5.16) Satz vom starken komplementären Schlupf.** *Besitzen beide dualen linearen Programme*

$$\text{(P)} \quad \max_{Ax \leq b} c^T x \quad \text{und} \quad \text{(D)} \quad \begin{array}{l} \min u^T b \\ u^T A = c^T \\ u \geq 0 \end{array}$$

zulässige Lösungen, so existieren optimale Lösungen  $\bar{x}, \bar{u}$ , so dass für alle Zeilenindizes  $i$  von  $A$  gilt:

$$\begin{aligned} \bar{u}_i > 0 &\iff A_{i.}\bar{x} = b_i \\ \text{bzw. } (A_{i.}\bar{x} < b_i) &\iff \bar{u}_i = 0 \end{aligned}$$

**Beweis :** Aufgrund von Satz (5.15) genügt es zu zeigen, dass optimale Lösungen  $\bar{x}, \bar{u}$  existieren mit

$$\begin{aligned} \text{(i)} \quad A_{i.}\bar{x} = b_i &\implies \bar{u}_i > 0 \\ \text{(bzw. (ii))} \quad \bar{u}_i = 0 &\implies A_{i.}\bar{x} < b_i \end{aligned}$$

Setzen wir

$$\bar{A} := \begin{pmatrix} -c^T & b^T \\ A & 0 \\ 0 & A^T \\ 0 & -A^T \\ 0 & -I \end{pmatrix}, \quad \bar{a} := \begin{pmatrix} 0 \\ b \\ c \\ -c \\ 0 \end{pmatrix}, \quad \bar{B} := (A, -I), \quad \bar{b} := b$$

dann gilt:

$$\begin{aligned} \bar{x}, \bar{u} \text{ optimal} &\iff \bar{A} \begin{pmatrix} \bar{x} \\ \bar{u} \end{pmatrix} \leq \bar{a}, \\ \bar{x}, \bar{u} \text{ erfüllen (i) und (ii)} &\iff \bar{B} \begin{pmatrix} \bar{x} \\ \bar{u} \end{pmatrix} < \bar{b}. \end{aligned}$$

Zum Beweis der Behauptung genügt es also zu zeigen, dass

$$\bar{A} \begin{pmatrix} \bar{x} \\ \bar{u} \end{pmatrix} \leq \bar{a}, \quad \bar{B} \begin{pmatrix} \bar{x} \\ \bar{u} \end{pmatrix} < \bar{b}$$

konsistent ist. Mit Satz (5.14) hat nach Voraussetzung  $\bar{A} \begin{pmatrix} \bar{x} \\ \bar{u} \end{pmatrix} \leq \bar{a}$  eine Lösung, also ist nach Folgerung (4.5) die Konsistenz dieses Systems äquivalent zur Inkonsistenz von

$$p \geq 0, \quad 0 \neq q \geq 0, \quad p^T \bar{A} + q^T \bar{B} = 0^T, \quad p^T \bar{a} + q^T \bar{b} \leq 0.$$

Angenommen, dieses System besitzt eine Lösung, sagen wir  $p^T = (\lambda, p_1^T, p_2^T, p_3^T, p_4^T) \geq 0, 0 \neq q \geq 0$ , dann gilt:

$$\begin{aligned} (p_1 + q)^T A &= \lambda c^T \\ (p_2 - p_3)^T A^T - (p_4 + q)^T &= -\lambda b^T \\ (p_1 + q)^T b + (p_2 - p_3)^T c &\leq 0. \end{aligned}$$

Daraus folgt

$$\begin{aligned} 0 &\geq \lambda((p_1 + q)^T b + (p_2 - p_3)^T c) = (p_1 + q)^T \lambda b + (p_2 - p_3)^T \lambda c \\ &= (p_1 + q)^T A(p_3 - p_2) + (p_1 + q)^T (p_4 + q) + (p_2 - p_3)^T A^T (p_1 + q) \\ &= p_1^T p_4 + p_1^T q + q^T p_4 + q^T q \\ &\geq q^T q \\ &> 0. \end{aligned}$$

Daraus folgt  $(p_1 + q)^T b + (p_2 - p_3)^T c > 0$ , ein Widerspruch.  $\square$

**(5.17) Hausaufgabe.** Formulieren und beweisen Sie den Satz vom starken komplementären Schlupf für das Paar dualer linearer Programme (5.12), (5.13).

□

Es gibt eine natürliche Merkregel für die verschiedenen Sätze vom komplementären Schlupf. Wir wissen bereits, dass zu jeder primalen Variablen eine duale Restriktion (auch komplementäre Restriktion genannt) und zu jeder primalen Restriktion eine duale (komplementäre) Variable gehören. Die Sätze (5.15) und (5.16) zeigen nun, dass zur Charakterisierung von Optimallösungen vorzeichenbeschränkte Variable und Ungleichungen besonders wichtig sind. Zu jeder vorzeichenbeschränkten primalen (dualen) Variablen gehört eine duale (primale) Ungleichung und umgekehrt. Ist  $a^T x \leq \alpha$  eine Ungleichung und  $\bar{x}$  ein Vektor, so nennen wir die Ungleichung **straff** (bezüglich  $\bar{x}$ ), falls  $a^T \bar{x} = \alpha$  gilt, andernfalls nennen wir sie **locker**. Der Satz vom schwachen komplementären Schlupf (5.15) sagt dann aus: Gewisse Vektoren sind optimal für (5.12) bzw. (5.13) genau dann, wenn für jede lockere Nichtnegativitätsbedingung die komplementäre Ungleichung straff ist, d. h. wenn in der komplementären Ungleichung kein Schlupf auftritt. Satz (5.16) kann wie folgt formuliert werden. Es gibt Optimallösungen, bei denen Schlüpfen komplementär auftreten, bzw. bei denen die duale Nichtnegativitätsbedingung genau dann locker ist, wenn die komplementäre primale Ungleichung straff ist.

# Kapitel 6

## Grundlagen der Polyedertheorie

Wir wollen nun die Polyedertheorie weiter ausbauen, weitere Darstellungen von Polyedern angeben und neue Operationen mit Polyedern einführen. Wir werden dabei meistens von der geometrischen Anschauung ausgehen und die Begriffe von dieser Sicht aus motivieren. Wir werden jedoch unser besonderes Augenmerk auf die analytische Beschreibung der geometrischen Konzepte legen.

### 6.1 Transformationen von Polyedern

Wir haben bereits in Kapitel 3 Projektionen von Polyedern entlang eines Richtungsvektors  $c$  untersucht und in Folgerung (3.9) festgestellt, dass eine derartige Projektion auf ein Polyeder wieder ein Polyeder ist. Wenn man mehrfach hintereinander projiziert, bleibt diese Eigenschaft natürlich erhalten. Sind  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$  und  $k, r \geq 0$  mit  $k + r = n$ , so nennt man die Menge

$$Q := \left\{ x \in \mathbb{K}^k \mid \exists y \in \mathbb{K}^r \text{ mit } \begin{pmatrix} x \\ y \end{pmatrix} \in P(\tilde{A}, b) \right\}$$

eine **Projektion von**  $P(A, b)$  **auf**  $\mathbb{K}^k$ . Hierbei ist  $\tilde{A} \in \mathbb{K}^{(m,n)}$  eine Matrix, die durch Spaltenvertauschung aus  $A$  hervorgeht. Offensichtlich folgt aus unseren Resultaten aus Kapitel 3:

**(6.1) Bemerkung.** *Jede Projektion eines Polyeders  $P(A, b) \subseteq \mathbb{K}^n$  auf  $\mathbb{K}^k$ ,  $k \leq n$ , ist ein Polyeder.*

□

Dieser Sachverhalt ist in größerer Allgemeinheit gültig. Erinnern wir uns daran, dass jede **affine Abbildung**  $f : \mathbb{K}^n \rightarrow \mathbb{K}^k$  gegeben ist durch eine Matrix  $D \in \mathbb{K}^{(k,n)}$  und einen Vektor  $d \in \mathbb{K}^k$ , so dass

$$f(x) = Dx + d \quad \forall x \in \mathbb{K}^n.$$

Für derartige Abbildungen gilt:

**(6.2) Satz.** *Affine Bilder von Polyedern sind Polyeder.*

**Beweis :** Seien  $P = P(A, b) \subseteq \mathbb{K}^n$  ein Polyeder und  $f(x) = Dx + d$  eine affine Abbildung von  $\mathbb{K}^n$  in den  $\mathbb{K}^k$ , dann gilt

$$\begin{aligned} f(P) &= \{y \in \mathbb{K}^k \mid \exists x \in \mathbb{K}^n \text{ mit } Ax \leq b \text{ und } y = Dx + d\} \\ &= \{y \in \mathbb{K}^k \mid \exists x \in \mathbb{K}^n \text{ mit } B \begin{pmatrix} x \\ y \end{pmatrix} \leq \bar{b}\}, \end{aligned}$$

wobei

$$B := \begin{pmatrix} A & 0 \\ D & -I \\ -D & I \end{pmatrix}, \quad \bar{b} := \begin{pmatrix} b \\ -d \\ d \end{pmatrix}.$$

Wenden wir nun unser Verfahren (3.6) iterativ auf  $B, \bar{b}$  und die Richtungsvektoren  $e_1, e_2, \dots, e_n$  an, so erhalten wir nach Satz (3.8) ein System  $C \begin{pmatrix} x \\ y \end{pmatrix} \leq c$  mit  $C = (0, \bar{C})$ , und es gilt:

$$\begin{aligned} \forall y \in \mathbb{K}^k : (\exists x \in \mathbb{K}^n \text{ mit } B \begin{pmatrix} x \\ y \end{pmatrix} \leq \bar{b}) &\iff \exists x \in \mathbb{K}^n \text{ mit } C \begin{pmatrix} x \\ y \end{pmatrix} \leq c \\ &\iff \bar{C}y \leq c. \end{aligned}$$

Daraus folgt  $f(P) = \{y \in \mathbb{K}^k \mid \bar{C}y \leq c\}$  ist ein Polyeder.  $\square$

Man beachte, dass der Beweis von (6.2) sogar ein Verfahren zur expliziten Konstruktion des affinen Bildes von  $P(A, b)$  beinhaltet. Aus (6.2) ergibt sich direkt die folgende (auch aus anderen Gründen unmittelbar einsichtige) Beobachtung.

**(6.3) Folgerung** (Satz von Weyl). *Für jede Matrix  $A \in \mathbb{K}^{(m,n)}$  gilt:*

$$\left. \begin{array}{l} \text{lin}(A) \\ \text{aff}(A) \\ \text{conv}(A) \\ \text{cone}(A) \end{array} \right\} \text{ ist ein Polyeder.}$$

**Beweis :** Wir zeigen die Behauptung für die konische Hülle. Alle anderen Fälle beweist man analog.

$$\begin{aligned} \text{cone}(A) &= \{x \in \mathbb{K}^m \mid \exists y \geq 0 \text{ mit } x = Ay\} \\ &= \{x \in \mathbb{K}^m \mid \exists y \in \mathbb{K}^n \text{ mit } \begin{pmatrix} I & -A \\ -I & A \\ 0 & -I \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \leq 0\}. \end{aligned}$$

Die letzte Menge ist die Projektion eines Polyeders im  $\mathbb{K}^{m+n}$  auf den  $\mathbb{K}^m$ , also nach (6.1) bzw. (6.2) ein Polyeder.  $\square$

Offenbar besagt die obige Folgerung nichts anderes als: Die lineare, affine, konvexe oder konische Hülle einer endlichen Teilmenge des  $\mathbb{K}^n$  ist ein Polyeder. Für die konische Hülle hat dies WEYL (1935) gezeigt (daher der Name für Folgerung (6.3)).

**(6.4) Folgerung.** Die Summe  $P = P_1 + P_2$  zweier Polyeder  $P_1, P_2$  ist ein Polyeder.

**Beweis :** Es seien  $P_1 = P(A, a), P_2 = P(B, b)$ , dann gilt:

$$\begin{aligned} P = P_1 + P_2 &= \{x + y \in \mathbb{K}^n \mid Ax \leq a, By \leq b\} \\ &= \{z \in \mathbb{K}^n \mid \exists x, y \in \mathbb{K}^n \text{ mit } Ax \leq a, By \leq b, z = x + y\} \\ &= \{z \in \mathbb{K}^n \mid \exists x, y \in \mathbb{K}^n \text{ mit } A(z - y) \leq a, B(z - x) \leq b\} \\ &= \{z \in \mathbb{K}^n \mid \exists x, y \in \mathbb{K}^n \text{ mit } D \begin{pmatrix} x \\ y \\ z \end{pmatrix} \leq \begin{pmatrix} a \\ b \end{pmatrix}\} \end{aligned}$$

mit

$$D = \begin{pmatrix} 0 & -A & A \\ -B & 0 & B \end{pmatrix}.$$

Also ist  $P$  die Projektion eines Polyeders des  $\mathbb{K}^{3n}$  auf den  $\mathbb{K}^n$ , und somit nach (6.1) ein Polyeder.  $\square$

Verbinden wir nun die Erkenntnis aus (6.3), dass  $\text{conv}(A)$  und  $\text{cone}(B)$  Polyeder sind, mit (6.4), so erhalten wir:

**(6.5) Folgerung.** Es seien  $A \in \mathbb{K}^{(m,n)}, B \in \mathbb{K}^{(m,n')}$ , dann gilt

$$P = \text{conv}(A) + \text{cone}(B)$$

ist ein Polyeder.  $\square$

Die obige Folgerung erscheint (durch geschickte Vorbereitung) völlig trivial, sie ist jedoch eine durchaus beachtenswerte Erkenntnis, denn wir werden bald zeigen, dass in der Tat alle Polyeder von der Form  $\text{conv}(A) + \text{cone}(B)$  sind.

## 6.2 Kegelpolarität

Es gibt mehrere Möglichkeiten die Umkehrung von (6.5) zu beweisen. Eine besondere elegante, die eine geometrische Version des Farkas-Lemmas benutzt, führt über die Kegelpolarität. Diese Operation mag zunächst nur als technisches Hilfsmittel erscheinen. Sie und ihre Verallgemeinerungen (allgemeine Polaritäten, Blocker, Antiblocker) sind jedoch bedeutende Methoden in der Polyedertheorie und der linearen sowie ganzzahligen Optimierung.

Wir beginnen mit einer Neuinterpretation des Farkas-Lemmas (4.2) (c). Dieses besagt

$$\exists x \geq 0, Ax = b \iff \forall u (A^T u \geq 0 \Rightarrow u^T b \geq 0).$$

Durch diese Aussage sind offenbar auch alle rechten Seiten  $b$  charakterisiert, für die  $x \geq 0, Ax = b$  eine Lösung hat. Nach Definition gilt  $\text{cone}(A) = \{b \in \mathbb{K}^m \mid \exists x \geq 0 \text{ mit } Ax = b\}$ , also können wir aus (4.2) (c) folgern:

**(6.6) Bemerkung.** Für alle Matrizen  $A \in \mathbb{K}^{(m,n)}$  gilt:

$$\text{cone}(A) = \{b \in \mathbb{K}^m \mid u^T b \leq 0 \forall u \in P(A^T, 0)\}.$$

□

Bemerkung (6.6) kann man geometrisch wie folgt beschreiben. Die Menge der zulässigen rechten Seiten  $b$  von  $Ax = b, x \geq 0$  ist genau die Menge aller Vektoren  $b \in \mathbb{K}^m$ , welche einen stumpfen Winkel mit allen Vektoren des Kegels  $P(A^T, 0)$  bilden. Allgemeiner definieren wir nun für jede beliebige Menge  $S \subseteq \mathbb{K}^n$

$$S^\circ := \{y \in \mathbb{K}^n \mid y^T x \leq 0 \forall x \in S\}.$$

$S^\circ$  ist die Menge aller Vektoren, die einen stumpfen Winkel mit allen Vektoren aus  $S$  bilden.  $S^\circ$  heißt **polarer Kegel** von  $S$ . (Überzeugen Sie sich, dass  $S^\circ$  ein Kegel ist!) Wir erinnern hier an das in der linearen Algebra definierte **orthogonale Komplement**

$$S^\perp := \{y \in \mathbb{K}^n \mid y^T x = 0 \forall x \in S\}.$$

Offensichtlich gilt  $S^\perp \subseteq S^\circ$ . Unter Benutzung der obigen Definition können wir Bemerkung (6.6) nun auch wie folgt aufschreiben.

**(6.7) Folgerung.** Für alle Matrizen  $A \in \mathbb{K}^{(m,n)}$  gilt

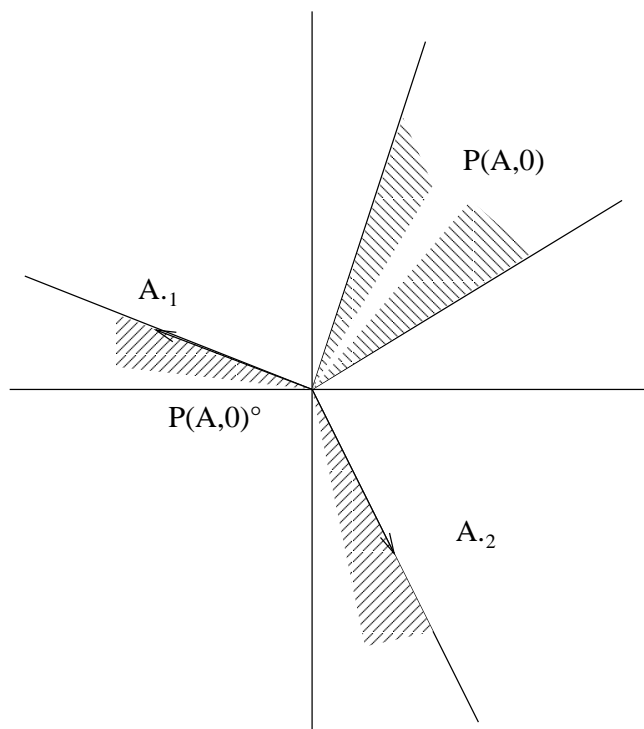
$$P(A, 0)^\circ = \text{cone}(A^T).$$

□

Folgerung (6.7) und die vorher gemachten Beobachtungen wollen wir an einem Beispiel erläutern. Es sei

$$A = \begin{pmatrix} -3 & 1 \\ 1 & -2 \end{pmatrix},$$

dann sind die Kegel  $P(A, 0)$  und  $P(A, 0)^\circ$  in Abbildung 6.1 gezeichnet.



**Abb. 6.1**

$P(A, 0)^\circ$  besteht also aus allen Vektoren, die mit den Elementen des Kegels  $P(A, 0)$  einen stumpfen Winkel bilden, und das sind gerade diejenigen Vektoren, die als konische Kombination der Normalenvektoren  $A_i$  dargestellt werden können, also  $P(A, 0)^\circ = \text{cone}\left(\begin{pmatrix} -3 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -2 \end{pmatrix}\right)$ . Ferner gilt:  $Ax = b, x \geq 0$  ist genau dann lösbar, wenn  $b \in P(A, 0)^\circ$  gilt. Daraus folgt z. B., dass  $Ax = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, x \geq 0$  nicht lösbar ist, während  $Ax = \begin{pmatrix} -1 \\ 0 \end{pmatrix}, x \geq 0$  eine Lösung hat.

Aus der Definition des polaren Kegels und des orthogonalen Komplements ergeben sich unmittelbar einige triviale Beziehungen, deren Beweis wir dem Leser zur Übung überlassen. Wir schreiben im Weiteren

$$S^{\circ\circ} := (S^\circ)^\circ.$$

**(6.8) Hausaufgabe.** Für  $S, S_i \subseteq \mathbb{K}^n, i = 1, \dots, k$  gilt:

- (a)  $S_i \subseteq S_j \implies S_j^\circ \subseteq S_i^\circ$   
 (b)  $S \subseteq S^{\circ\circ}$   
 (c)  $\left(\bigcup_{i=1}^k S_i\right)^\circ = \bigcap_{i=1}^k S_i^\circ$   
 (d)  $S^\circ = \text{cone}(S^\circ) = (\text{cone}(S))^\circ$   
 (e)  $S = \text{lin}(S) \implies S^\circ = S^\perp$ . Gilt die Umkehrung?  
 (f) Ersetzen wir in (a), ..., (d) "o" durch " $\perp$ ", sind dann auch noch alle Behauptungen wahr?  $\square$

**(6.9) Hausaufgabe.** Für welche Mengen  $S \subseteq \mathbb{K}^n$  gilt

- (a)  $S^\circ = S^{\circ\circ\circ}$ ,  
 (b)  $S = S^\circ$  ?

$\square$

Die Aussage (6.8) (d) impliziert insbesondere:

**(6.10) Folgerung.**  $\text{cone}(A^T)^\circ = P(A, 0)$ .

**Beweis:**  $(\text{cone}(A^T))^\circ = \text{cone}((A^T)^\circ) = (A^T)^\circ = \{x \mid Ax \leq 0\} = P(A, 0)$ .  $\square$

Das folgende Korollar aus (6.7) und (6.10) wird in der Literatur häufig mit einem Namen belegt.

**(6.11) Polarensatz.** Für jede Matrix  $A \in \mathbb{K}^{(m,n)}$  gilt:

$$\begin{aligned} P(A, 0)^{\circ\circ} &= P(A, 0), \\ \text{cone}(A)^{\circ\circ} &= \text{cone}(A). \end{aligned}$$

**Beweis:**

$$\begin{aligned} P(A, 0) &\stackrel{(6.10)}{=} \text{cone}(A^T)^\circ \stackrel{(6.7)}{=} P(A, 0)^{\circ\circ}, \\ \text{cone}(A) &\stackrel{(6.7)}{=} P(A^T, 0)^\circ \stackrel{(6.10)}{=} \text{cone}(A)^{\circ\circ}. \end{aligned}$$

$\square$

Unser kurzer Exkurs über Kegelpolarität ist damit beendet.

### 6.3 Darstellungssätze

Wir wollen nun zeigen, dass Polyeder nicht nur in der Form  $P(A, b)$  dargestellt werden können und benutzen dazu die bisher entwickelte Maschinerie.

**(6.12) Satz (Minkowski (1896)).** Eine Teilmenge  $K \subseteq \mathbb{K}^n$  ist genau dann ein polyedrischer Kegel, wenn  $K$  die konische Hülle von endlich vielen Vektoren ist. Mit anderen Worten: Zu jeder Matrix  $A \in \mathbb{K}^{(m,n)}$  gibt es eine Matrix  $B \in \mathbb{K}^{(n,k)}$ , so dass

$$P(A, 0) = \text{cone}(B)$$

gilt und umgekehrt.

**Beweis :**  $P(A, 0) \stackrel{(6.10)}{=} \text{cone}(A^T)^\circ \stackrel{(6.3)}{=} P(B^T, 0)^\circ \stackrel{(6.7)}{=} \text{cone}(B)$ . □

**(6.13) Satz.** Es seien  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$ , dann existieren endliche Mengen  $V, E \subseteq \mathbb{K}^n$  mit

$$P(A, b) = \text{conv}(V) + \text{cone}(E).$$

**Beweis :** Setze

$$H := P\left(\begin{pmatrix} A & -b \\ 0^T & -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix}\right),$$

dann gilt:  $x \in P(A, b) \iff \begin{pmatrix} x \\ 1 \end{pmatrix} \in H$ .  $H$  ist nach Definition ein polyedrischer Kegel. Also gibt es nach Satz (6.12) eine Matrix  $B \in \mathbb{K}^{(n+1,k)}$  mit  $H = \text{cone}(B)$ . Aufgrund der Definition von  $H$  hat die letzte Zeile von  $B$  nur nichtnegative Elemente. Durch Skalieren der Spalten von  $B$  und Vertauschen von Spalten können wir  $B$  in eine Matrix  $\bar{B}$  so umformen, dass gilt

$$\bar{B} = \begin{pmatrix} V & E \\ \mathbf{1}^T & 0^T \end{pmatrix}, \quad \text{cone}(\bar{B}) = H.$$

Daraus folgt:

$$\begin{aligned} x \in P(A, b) &\iff \begin{pmatrix} x \\ 1 \end{pmatrix} \in H \\ &\iff x = V\lambda + E\mu \text{ mit } \lambda^T \mathbf{1} = 1, \lambda, \mu \geq 0 \\ &\iff x \in \text{conv}(V) + \text{cone}(E). \end{aligned}$$

□

**(6.14) Folgerung.** Eine Teilmenge  $P \subseteq \mathbb{K}^n$  ist genau dann ein Polytop, wenn  $P$  die konvexe Hülle endlich vieler Vektoren ist.

**Beweis :** Sei  $V \subseteq \mathbb{K}^n$  endlich und  $P = \text{conv}(V)$ , dann ist  $P$  nach (6.3) ein Polyeder. Ist  $x \in P$ , so gilt  $x = \sum_{i=1}^k \lambda_i v_i$ ,  $v_i \in V$ ,  $\lambda_i \geq 0$ ,  $\sum_{i=1}^k \lambda_i = 1$ , und somit  $\|x\| \leq \sum_{i=1}^k \|v_i\|$ , d. h.  $P \subseteq \{x \mid \|x\| \leq \sum_{v \in V} \|v\|\}$ . Also ist  $P$  beschränkt, d. h.  $P$  ist ein Polytop.

Ist umgekehrt  $P$  ein Polytop, so gibt es nach Satz (6.13) endliche Mengen  $V, E$  mit  $P = \text{conv}(V) + \text{cone}(E)$ . Gibt es einen Vektor  $e \in E$  mit  $e \neq 0$ , so gilt für alle  $n \in \mathbb{N}$ ,  $x + ne \in P$  für alle  $x \in \text{conv}(V)$ . Also ist  $P$  unbeschränkt, falls  $E \setminus \{0\} \neq \emptyset$ . Daraus folgt  $E \in \{\emptyset, \{0\}\}$ , und dann gilt trivialerweise  $\text{conv}(V) = \text{conv}(V) + \text{cone}(E) = P$ .  $\square$

**(6.15) Darstellungssatz.** Eine Teilmenge  $P \subseteq \mathbb{K}^n$  ist genau dann ein Polyeder, wenn  $P$  die Summe eines Polytops und eines polyedrischen Kegels ist, d. h. wenn es endliche Mengen  $V, E \subseteq \mathbb{K}^n$  gibt mit

$$P = \text{conv}(V) + \text{cone}(E).$$

**Beweis :** Kombiniere (6.12), (6.13), (6.14) und (6.5).  $\square$

Ist  $P \subseteq \mathbb{K}^n$  ein Polyeder, so wissen wir nunmehr, dass es für  $P$  zwei mögliche Darstellungen gibt. Es gilt nämlich

$$P = P(A, b) = \text{conv}(V) + \text{cone}(E),$$

wobei  $A$  eine  $(m, n)$ -Matrix,  $b \in \mathbb{K}^m$  und  $V, E$  endliche Mengen sind. Diese beiden Darstellungen sind grundsätzlich verschieden, was natürlich in vielerlei Hinsicht nützlich sein kann. Manche Aussagen über Polyeder sind völlig trivial, wenn man von der einen Beschreibung ausgeht, während sie aus der anderen Beschreibung nicht unmittelbar folgen.

Die Darstellung  $P(A, b)$  nennt man auch **äußere Beschreibung** von  $P$ . Der Grund für diese Bezeichnung liegt darin, dass man wegen

$$P = \bigcap_{i=1}^m \{x \mid A_i \cdot x \leq b_i\} \subseteq \{x \mid A_i \cdot x \leq b_i\},$$

das Polyeder  $P$  als Durchschnitt von größeren Mengen betrachten kann.  $P$  wird sozusagen „von außen“ durch sukzessives Hinzufügen von Ungleichungen (bzw. Halbräumen) konstruiert.

Hingegen nennt man  $\text{conv}(V) + \text{cone}(E)$  eine **innere Beschreibung** von  $P$ . Ist  $E = \emptyset$ , so ist die Bezeichnung offensichtlich, denn  $V \subseteq P$  und somit wird  $P$

durch konvexe Hüllenbildung von Elementen von sich selbst erzeugt. Analoges gilt, wenn  $P$  ein polyedrischer Kegel ist. Sind jedoch  $V$  und  $E$  nicht leer, dann ist  $E$  nicht notwendigerweise eine Teilmenge von  $P$ , jedoch gelingt es eben aus den Vektoren  $v \in V$  zusammen mit den Vektoren  $e \in E$  das Polyeder  $P$  „von innen her“ zu konstruieren.

Die Sätze (6.12), (6.14) und (6.15) beinhalten weitere wichtige Charakterisierungen von polyedrischen Kegeln, Polytopen und Polyedern. Wir erinnern uns aus der linearen Algebra daran, dass jeder lineare Teilraum  $L$  des  $\mathbb{K}^n$  eine endliche Basis hat, d. h. eine endliche Teilmenge  $B$  besitzt, so dass  $B$  linear unabhängig ist und  $L = \text{lin}(B)$  gilt. Die linearen Teilräume des  $\mathbb{K}^n$  sind also diejenigen Teilmengen des  $\mathbb{K}^n$ , deren Elemente durch Linearkombinationen einer endlichen Menge erzeugt werden können. Nach (6.14) sind Polytope genau diejenigen Teilmengen des  $\mathbb{K}^n$ , die durch Konvexkombinationen einer endlichen Menge erzeugt werden können.

Wir werden in Kapitel 8 sehen, dass es sogar eine eindeutig bestimmte minimale (im Sinne der Mengeneinklusion) endliche Menge  $V \subseteq \mathbb{K}^n$  gibt mit  $P = \text{conv}(V)$ , d. h. Polytope haben sogar eine eindeutig bestimmte „konvexe Basis“. Nach (6.12) sind polyedrische Kegel genau diejenigen Teilmengen des  $\mathbb{K}^n$ , die ein endliches „Kegelerzeugendensystem“ haben. Auch hier gibt es natürlich minimale endliche Mengen, die die Kegel konisch erzeugen. Aber nur unter zusätzlichen Voraussetzungen haben zwei minimale konische Erzeugendensysteme auch gleiche Kardinalität, und Eindeutigkeit gilt lediglich bis auf Multiplikation mit positiven Skalaren. Häufig nennt man eine Teilmenge  $T$  des  $\mathbb{K}^n$  **endlich erzeugt**, falls  $T = \text{conv}(V) + \text{cone}(E)$  für endliche Mengen  $V, E$  gilt. Nach (6.15) sind also die Polyeder gerade die endlich erzeugten Teilmengen des  $\mathbb{K}^n$ . Fassen wir zusammen, so gilt:

**(6.16) Bemerkung.** Ist  $T \subseteq \mathbb{K}^n$ , so gilt

- (a)  $T$  ist ein linearer Teilraum  $\iff T$  ist die lineare Hülle einer endlichen Menge.
- (b)  $T$  ist ein affiner Teilraum  $\iff T$  ist die affine Hülle einer endlichen Menge.
- (c)  $T$  ist ein polyedrischer Kegel  $\iff T$  ist die konische Hülle einer endlichen Menge.
- (d)  $T$  ist ein Polytop  $\iff T$  ist die konvexe Hülle einer endlichen Menge.
- (e)  $T$  ist ein Polyeder  $\iff T$  ist endlich erzeugt.

□



# Kapitel 7

## Seitenflächen von Polyedern

Wir wollen uns nun mit den “Begrenzungsflächen” von Polyedern  $P$  beschäftigen und diese charakterisieren. Wir gehen dabei so vor, dass wir die Objekte, die wir untersuchen wollen, zunächst “abstrakt”, d. h. durch eine allgemeine Definition einführen, und dann werden wir diese Objekte **darstellungsabhängig** kennzeichnen. Das heißt, wir werden jeweils zeigen, wie die betrachteten Objekte “aussehen”, wenn  $P$  durch  $P(A, b)$  oder durch  $\text{conv}(V) + \text{cone}(E)$  gegeben ist. Wir werden uns meistens auf die Darstellung  $P(A, b)$  beziehen, da diese für die Optimierung die wichtigere ist. In Abschnitt 7.1 werden wir jedoch zeigen, wie man Charakterisierungen bezüglich einer Darstellung auf Charakterisierungen bezüglich der anderen Darstellung transformieren kann, so dass aus den hier vorgestellten Resultaten die Kennzeichnungen bezüglich der Darstellung  $\text{conv}(V) + \text{cone}(E)$  (mit etwas technischem Aufwand) folgen.

Wir wollen nochmals auf eine die Schreibtechnik vereinfachende Konvention hinweisen. Wir betrachten Matrizen und endliche Mengen als (im Wesentlichen) ein und dieselben Objekte. Ist z. B.  $A$  eine  $(m, n)$ -Matrix, so schreiben wir

$$\text{conv}(A)$$

und meinen damit die konvexe Hülle der Spaltenvektoren von  $A$ . Ist umgekehrt z. B.  $V \subseteq \mathbb{K}^n$  eine endliche Menge, so fassen wir  $V$  auch als eine  $(n, |V|)$ -Matrix auf und schreiben

$$V^T,$$

um die Matrix zu bezeichnen, deren Zeilen die Vektoren aus  $V$  sind.  $V^T$  ist natürlich nicht eindeutig definiert, da wir keine Reihenfolge der Vektoren aus  $V$  angegeben haben. Wir werden diese Schreibweise jedoch nur dann benutzen, wenn es auf die Reihenfolge der Zeilen nicht ankommt.

## 7.1 Die $\gamma$ -Polare und gültige Ungleichungen

In Abschnitt 6 haben wir bereits einen Polarentyp, die Kegelpolare, zur Beweisvereinfachung eingeführt. Hier wollen wir eine weitere Polare betrachten, die es uns ermöglichen wird, eine Charakterisierung bezüglich  $P(A, b)$  in eine Charakterisierung bezüglich  $\text{conv}(V) + \text{cone}(E)$  zu übertragen.

**(7.1) Definition.** Es seien  $S \subseteq \mathbb{K}^n$ ,  $a \in \mathbb{K}^n$ ,  $\alpha \in \mathbb{K}$ .

- (a) Die Ungleichung  $a^T x \leq \alpha$  heißt **gültig** bezüglich  $S$ , falls  $S \subseteq \{x \in \mathbb{K}^n \mid a^T x \leq \alpha\}$ .
- (b) Die Menge  $S^\gamma := \left\{ \begin{pmatrix} a \\ \alpha \end{pmatrix} \in \mathbb{K}^{n+1} \mid a^T x \leq \alpha \forall x \in S \right\}$  heißt  **$\gamma$ -Polare** von  $S$ .
- (c) Eine Hyperebene  $H = \{x \mid a^T x = \alpha\}$  heißt **Stützhyperebene** von  $S$ , falls  $\begin{pmatrix} a \\ \alpha \end{pmatrix} \in S^\gamma$  und  $S \cap H \neq \emptyset$ .  $\square$

Zu sagen, dass  $a^T x \leq \alpha$ ,  $a \neq 0$ , gültig bezüglich  $S$  ist, heißt nichts anderes als:  $S$  ist im Halbraum  $a^T x \leq \alpha$  enthalten. Die  $\gamma$ -Polare  $S^\gamma$  kann als die „Menge aller gültigen Ungleichungen bezüglich  $S$ “ betrachtet werden. Wir wollen nun die  $\gamma$ -Polare eines Polyeders charakterisieren.

**(7.2) Satz.** Es sei  $P \subseteq \mathbb{K}^n$ ,  $P \neq \emptyset$ , ein Polyeder mit den Darstellungen  $P = P(A, b) = \text{conv}(V) + \text{cone}(E)$ , dann gilt:

$$\begin{aligned} (a) \quad P^\gamma &= \left\{ \begin{pmatrix} a \\ \alpha \end{pmatrix} \in \mathbb{K}^{n+1} \mid \exists u \geq 0, u^T A = a^T, u^T b \leq \alpha \right\} \\ &= \text{cone} \left( \begin{pmatrix} A^T & 0 \\ b^T & 1 \end{pmatrix} \right). \end{aligned}$$

$$\begin{aligned} (b) \quad P^\gamma &= \left\{ \begin{pmatrix} a \\ \alpha \end{pmatrix} \in \mathbb{K}^{n+1} \mid \begin{pmatrix} V^T & -\mathbf{1} \\ E^T & 0 \end{pmatrix} \cdot \begin{pmatrix} a \\ \alpha \end{pmatrix} \leq 0 \right\} \\ &= P \left( \left( \begin{pmatrix} V^T & -\mathbf{1} \\ E^T & 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right) \right). \end{aligned}$$

**Beweis :**

$$\begin{aligned}
(a) \quad \begin{pmatrix} a \\ \alpha \end{pmatrix} \in P^\gamma &\iff Ax \leq b, a^T x > \alpha \quad \text{inkonsistent} \\
&\stackrel{(4.5)}{\iff} \exists u \geq 0, v > 0 \text{ mit } u^T A - v a^T = 0, u^T b - v \alpha \leq 0 \\
&\iff \exists u \geq 0 \text{ mit } u^T A = a^T, u^T b \leq \alpha \\
&\iff \exists u \geq 0, \lambda \geq 0 \text{ mit } \begin{pmatrix} a \\ \alpha \end{pmatrix} = \begin{pmatrix} A^T & 0 \\ b^T & 1 \end{pmatrix} \begin{pmatrix} u \\ \lambda \end{pmatrix} \\
&\iff \begin{pmatrix} a \\ \alpha \end{pmatrix} \in \text{cone} \begin{pmatrix} A^T & 0 \\ b^T & 1 \end{pmatrix}.
\end{aligned}$$

$$\begin{aligned}
(b) \quad \begin{pmatrix} a \\ \alpha \end{pmatrix} \in P^\gamma &\implies a^T v \leq \alpha \quad \forall v \in V \quad \text{und} \\
&\quad a^T(v + \lambda e) \leq \alpha \quad \forall v \in V, e \in E, \lambda \geq 0 \\
&\implies a^T e \leq 0 \\
&\quad (\text{andernfalls w\"are } a^T(v + \lambda e) > \alpha \text{ f\"ur gen\"ugend gro\sses } \lambda) \\
&\implies \begin{pmatrix} V^T & -\mathbb{1} \\ E^T & 0 \end{pmatrix} \begin{pmatrix} a \\ \alpha \end{pmatrix} \leq 0.
\end{aligned}$$

Gilt umgekehrt das letztere Ungleichungssystem, und ist  $x \in P$ , so existieren  $v_1, \dots, v_p \in V$  und  $e_1, \dots, e_q \in E$ ,  $\lambda_1, \dots, \lambda_p \geq 0$ ,  $\sum_{i=1}^p \lambda_i = 1$ ,  $\mu_1, \dots, \mu_q \geq 0$ , so dass

$$x = \sum_{i=1}^p \lambda_i v_i + \sum_{j=1}^q \mu_j e_j.$$

Und daraus folgt

$$a^T x = \sum_{i=1}^p \lambda_i a^T v_i + \sum_{j=1}^q \mu_j a^T e_j \leq \sum_{i=1}^p \lambda_i \alpha + \sum_{j=1}^q \mu_j 0 = \alpha,$$

also gilt  $\begin{pmatrix} a \\ \alpha \end{pmatrix} \in P^\gamma$ . □

**(7.3) Folgerung.** Die  $\gamma$ -Polare eines Polyeders  $\emptyset \neq P \subseteq \mathbb{K}^n$  ist ein polyedrischer Kegel im  $\mathbb{K}^{n+1}$ . □

**(7.4) Folgerung.** Ist  $\emptyset \neq P = P(A, b) = \text{conv}(V) + \text{cone}(E)$  ein Polyeder und  $a^T x \leq \alpha$  eine Ungleichung, dann sind äquivalent:

- (i)  $a^T x \leq \alpha$  ist gültig bezüglich  $P$ .
- (ii)  $\exists u \geq 0$  mit  $u^T A = a^T, u^T b \leq \alpha$ .
- (iii)  $a^T v \leq \alpha \forall v \in V$  und  $a^T e \leq 0 \forall e \in E$ .
- (iv)  $\begin{pmatrix} a \\ \alpha \end{pmatrix} \in P^\gamma$

□

## 7.2 Seitenflächen

Wir werden nun diejenigen Teilmengen von Polyedern untersuchen, die als Durchschnitte des Polyeders mit Stützhyperbenen entstehen.

**(7.5) Definition.**  $P \subseteq \mathbb{K}^n$  sei ein Polyeder. Eine Menge  $F \subseteq P$  heißt **Seitenfläche** von  $P$ , wenn es eine gültige Ungleichung  $c^T x \leq \gamma$  bezüglich  $P$  gibt mit

$$F = P \cap \{x \mid c^T x = \gamma\}.$$

Eine Seitenfläche  $F$  von  $P$  heißt **echt**, wenn  $F \neq P$  gilt.  $F$  heißt **nichttrivial**, wenn  $\emptyset \neq F \neq P$  gilt. Ist  $c^T x \leq \gamma$  gültig bezüglich  $P$ , dann heißt

$$F := P \cap \{x \mid c^T x = \gamma\}$$

die von  $c^T x \leq \gamma$  **induzierte oder definierte Seitenfläche**.

□

Offenbar sind Seitenflächen von Polyedern wiederum Polyeder. Eine Seitenfläche  $F$  kann durchaus von verschiedenen Ungleichungen definiert werden. Trivialerweise gilt  $F = P \cap \{x \mid c^T x = \gamma\} = P \cap \{x \mid \lambda c^T x = \lambda \gamma\}$  für alle  $\lambda > 0$ , aber auch zwei Ungleichungen, die nicht durch Skalierung auseinander hervorgehen, können  $F$  definieren.

**(7.6) Beispiel.** Sei  $P(A, b) \subseteq \mathbb{K}^2$  das durch

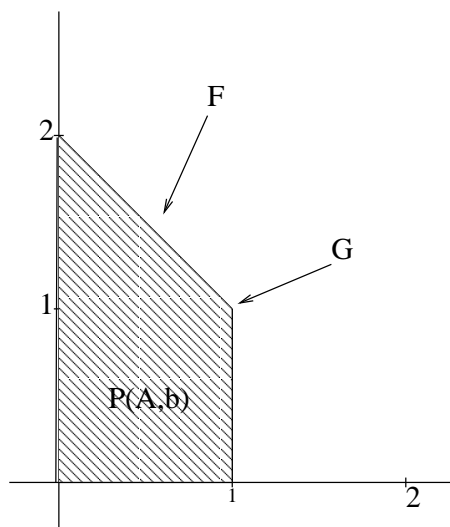
$$A = \begin{pmatrix} 1 & 1 \\ 1 & 0 \\ -1 & 0 \\ 0 & -1 \end{pmatrix} \quad b = \begin{pmatrix} 2 \\ 1 \\ 0 \\ 0 \end{pmatrix},$$

definierte Polyeder (siehe Abbildung 7.1), dann ist das Geradenstück zwischen  $(1, 1)^T$  und  $(0, 2)^T$  eine Seitenfläche  $F$  definiert durch  $x_1 + x_2 \leq 2$ . Die Menge  $G = \{(1, 1)^T\}$  ist ebenfalls eine Seitenfläche von  $P$ , denn

$$\begin{aligned} G &= P(A, b) \cap \{x \mid 2x_1 + x_2 = 3\} \\ &= P(A, b) \cap \{x \mid 3x_1 + x_2 = 4\}, \end{aligned}$$

und, wie man sieht, kann  $\{(1, 1)^T\}$  durch zwei völlig unterschiedliche Ungleichungen definiert werden.

□



**Abb. 7.1**

**(7.7) Folgerung.** Sei  $P \subseteq \mathbb{K}^n$  ein Polyeder, dann gilt:

- (a)  $P$  ist eine Seitenfläche von sich selbst.
- (b)  $\emptyset$  ist eine Seitenfläche von  $P$ .
- (c) Ist  $F = \{x \in P \mid c^T x = \gamma\}$  eine Seitenfläche von  $P$ , dann gilt  $c \neq 0$ .

- Beweis :** (a)  $P = \{x \in P \mid 0^T x = 0\}$   
 (b)  $\emptyset = \{x \in P \mid 0^T x = 1\}$   
 (c) klar.

□

**(7.8) Satz.** Seien  $P = P(A, b) \neq \emptyset$  ein Polyeder,  $c \in \mathbb{K}^n$ ,  $\gamma \in \mathbb{K}$  und

$$z^* = \begin{cases} +\infty, & \text{falls } c^T x \text{ auf } P \text{ unbeschränkt,} \\ \max\{c^T x \mid x \in P\}, & \text{andernfalls.} \end{cases}$$

- (a)  $c^T x \leq \gamma$  ist gültig bezüglich  $P \iff \gamma \geq z^*$ .  
 (b)  $z^* < +\infty \Rightarrow F = \{x \in P \mid c^T x = z^*\}$  ist eine nichtleere Seitenfläche von  $P$ , und  $\{x \mid c^T x = z^*\}$  ist eine Stützhyperebene von  $P$ , falls  $c \neq 0$ .  
 (c) Die Menge der Optimallösungen eines linearen Programms ist eine Seitenfläche des durch die Nebenbedingungen definierten Polyeders.

**Beweis :** (a) Ist  $\gamma < z^*$ , dann existiert ein  $x^* \in P$  mit  $c^T x^* > \gamma$ , also ist  $c^T x \leq \gamma$  nicht gültig. Ist  $\gamma \geq z^*$ , so gilt  $c^T x \leq z^* \leq \gamma \forall x \in P$ , also ist  $c^T x \leq \gamma$  gültig.

(b) Nach (a) ist  $c^T x \leq z^*$  gültig bezüglich  $P$  und nach Definition existiert ein  $x^* \in P$  mit  $c^T x^* = z^*$ , also ist  $F$  eine nichtleere Seitenfläche von  $P$ . Offenbar ist  $\{x \mid c^T x = z^*\}$  eine Stützhyperebene, falls  $c \neq 0$ .

(c) folgt aus (b). □

**(7.9) Definition.** Sei  $P = P(A, b) \subseteq \mathbb{K}^n$  ein Polyeder und  $M$  die Zeilenindexmenge von  $A$ . Für  $F \subseteq P$  sei

$$\text{eq}(F) := \{i \in M \mid A_i \cdot x = b_i \forall x \in F\},$$

d. h.  $\text{eq}(F)$  (genannt **Gleichheitsmenge** oder **equality set von  $F$** ) ist die Menge der für alle  $x \in F$  bindenden Restriktionen. Für  $I \subseteq M$  sei

$$\text{fa}(I) = \{x \in P \mid A_I \cdot x = b_I\}.$$

Wir nennen  $\text{fa}(I)$  die von  $I$  induzierte Seitenfläche.

□

Offenbar ist  $\text{fa}(I)$  tatsächlich eine Seitenfläche von  $P = P(A, b)$ , denn setzen wir  $c^T := \sum_{i \in I} A_i$ ,  $\gamma := \sum_{i \in I} b_i$ , so ist  $c^T x \leq \gamma$  gültig bezüglich  $P(A, b)$ , und, wie man leicht sieht, gilt

$$\text{fa}(I) = \{x \in P \mid c^T x = \gamma\}.$$

Zur Veranschaulichung der oben definierten Abbildungen  $\text{fa}$  und  $\text{eq}$  betrachten wir Beispiel (7.6). Für das in (7.6) definierte Polyeder  $P(A, b)$  gilt  $M = \{1, 2, 3, 4\}$  und

$$\text{fa}(\{1, 2\}) = G,$$

$$\text{eq}(\left\{\binom{0}{2}\right\}) = \{1, 3\}.$$

Wir zeigen nun, dass man die Gleichheitsmenge einer Seitenfläche explizit berechnen kann.

**(7.10) Satz.** Seien  $P = P(A, b)$  ein Polyeder und  $\emptyset \neq F = \{x \in P \mid c^T x = \gamma\}$  eine Seitenfläche von  $P$ . Dann gilt

$$\text{eq}(F) = \{i \in M \mid \exists u \geq 0 \text{ mit } u_i > 0 \text{ und } u^T A = c^T, u^T b = \gamma\}.$$

**Beweis:** Nach Voraussetzung und Folgerung (5.9) haben die beiden linearen Programme

$$(P) \quad \max_{Ax \leq b} c^T x \quad \text{und} \quad (D) \quad \begin{array}{l} \min u^T b \\ u^T A = c^T \\ u \geq 0 \end{array}$$

optimale Lösungen mit gleichem Zielfunktionswert  $\gamma$ , und  $F$  ist die Menge der Optimallösungen von (P). Sei nun  $i \in \text{eq}(F)$ . Aufgrund des Satzes (5.16) vom starken komplementären Schlupf existieren Optimallösungen  $\bar{x}$ ,  $\bar{u}$  von (P), (D) mit  $\bar{u}_j > 0 \Leftrightarrow A_j \bar{x} = b_j$ . Wegen  $\bar{x} \in F$  gilt  $A_i \bar{x} = b_i$ , also gibt es einen Vektor  $u$  mit den geforderten Eigenschaften.

Gibt es umgekehrt einen Vektor  $u \geq 0$  mit  $u_i > 0$  und  $u^T A = c^T$ ,  $u^T b = \gamma$ , so ist  $u$  optimal für (D), und aus dem Satz vom schwachen komplementären Schlupf (5.15) folgt  $A_i x = b_i$  für alle  $x \in F$ , d. h.  $i \in \text{eq}(F)$ .  $\square$

**(7.11) Satz.** Seien  $P = P(A, b)$  ein Polyeder und  $F \neq \emptyset$  eine Teilmenge von  $P$ , dann sind äquivalent:

- (i)  $F$  ist eine Seitenfläche von  $P$ .
- (ii)  $\exists I \subseteq M$  mit  $F = \text{fa}(I) = \{x \in P \mid A_I x = b_I\}$ .
- (iii)  $F = \text{fa}(\text{eq}(F))$ .

**Beweis :** Gelten (ii) oder (iii), dann ist  $F$  offenbar eine Seitenfläche von  $P$ .

(i)  $\Rightarrow$  (ii) Sei  $F = \{x \in P \mid c^T x = \gamma\}$  eine Seitenfläche. Setzen wir  $I := \text{eq}(F)$ , dann gilt nach Definition  $F \subseteq \{x \in P \mid A_I x = b_I\} =: F'$ . Ist  $x \in F'$ , so gilt  $x \in P$ ; es bleibt zu zeigen, dass  $c^T x = \gamma$  gilt. Zu jedem  $i \in I$  gibt es nach (7.10) einen Vektor  $u^{(i)}$  mit  $u^{(i)} \geq 0$ ,  $u_i^{(i)} > 0$ ,  $(u^{(i)})^T A = c^T$  und  $(u^{(i)})^T b = \gamma$ . Setze

$$u := \sum_{i \in I} \frac{1}{|I|} u^{(i)},$$

dann gilt nach Konstruktion  $u_i > 0 \forall i \in I$  und ferner  $u_i = 0 \forall i \in M \setminus I$  (andernfalls wäre  $i \in I$  nach (7.10)). Die Vektoren  $x$  und  $u$  sind zulässig für die linearen Programme (P), (D) des Beweises von (7.10). Aus dem Satz vom schwachen komplementären Schlupf (5.15) folgt, dass sie auch optimal sind. Daraus folgt  $c^T x = \gamma$  und somit  $x \in F$ .

(ii)  $\Rightarrow$  (iii) folgt direkt aus dem obigen Beweis. □

Aus Satz (7.11) folgt, dass zur Darstellung einer Seitenfläche von  $P(A, b)$  keine zusätzliche Ungleichung benötigt wird. Man braucht lediglich in einigen der Ungleichungen des Systems  $Ax \leq b$  Gleichheit zu fordern. Da jede nichtleere Seitenfläche auf diese Weise erzeugt werden kann, folgt:

**(7.12) Folgerung.** Sind  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$ , dann hat das Polyeder  $P(A, b)$  höchstens  $2^m + 1$  Seitenflächen.

**Beweis :**  $M = \{1, \dots, m\}$  hat  $2^m$  Teilmengen. Für jede Teilmenge  $I \subseteq M$  ist  $P \cap \{x \mid A_I x = b_I\}$  eine Seitenfläche von  $P$ . Dazu kommt u. U. noch die leere Seitenfläche. □

Man kann Seitenflächen auf ähnliche Weise durch Einführung von Abbildungen analog zu  $\text{eq}$  bzw.  $\text{fa}$  bezüglich der Darstellung  $P = \text{conv}(V) + \text{cone}(E)$  charakterisieren. Diese Kennzeichnungen von Seitenflächen sind jedoch technisch aufwendiger. Der interessierte Leser sei dazu auf Bachem & Grötschel, “New Aspects of Polyhedral Theory”, in: B. Korte (ed.), “Modern Applied Mathematics — Optimization and Operations Research”, North-Holland, Amsterdam 1982 verwiesen.

## 7.3 Dimension

Wir wollen nun zeigen, dass man auch die Dimension einer Seitenfläche eines Polyeders explizit berechnen kann. Zunächst führen wir einen Hilfsbegriff ein.

**(7.13) Definition.** Ein Element  $x$  eines Polyeders  $P$  heißt **innerer Punkt** von  $P$ , wenn  $x$  in keiner echten Seitenfläche von  $P$  enthalten ist.  $\square$

**Achtung!** Innere Punkte eines Polyeders  $P$  sind nicht notwendig auch topologisch innere Punkte im Sinne der natürlichen Topologie des  $\mathbb{K}^n$ . Unsere inneren Punkte sind topologisch innere Punkte im Sinne der Relativtopologie auf  $P$ .

**(7.14) Satz.** Jedes nichtleere Polyeder besitzt innere Punkte.

**Beweis :** Sei  $P = P(A, b)$  und  $I = \text{eq}(P(A, b))$ ,  $J = M \setminus I$ . Gilt  $I = M$ , so hat  $P$  keine echten Seitenflächen, also ist jedes Element von  $P$  ein innerer Punkt. Andernfalls ist das System  $Ax \leq b$  äquivalent zu

$$A_I x = b_I,$$

$$A_J x \leq b_J.$$

$P$  hat innere Punkte heißt dann, dass es ein  $x$  gibt mit  $A_I x = b_I$  und  $A_J x < b_J$ . Zu jedem  $i \in J$  existiert nach Definition ein Vektor  $y^{(i)} \in P$  mit  $A_i y^{(i)} < b_i$ . Setze  $y := \frac{1}{|J|} \sum_{i \in J} y^{(i)}$ , dann ist  $y$  Konvexkombination von Elementen von  $P$ , also  $y \in P$ , und es gilt  $A_J y < b_J$ . Mithin ist  $y$  ein innerer Punkt von  $P$ .  $\square$

**(7.15) Satz.** Sei  $F$  Seitenfläche eines Polyeders  $P(A, b)$  und  $\bar{x} \in F$ . Der Vektor  $\bar{x}$  ist ein innerer Punkt von  $F$  genau dann, wenn  $\text{eq}(\{\bar{x}\}) = \text{eq}(F)$ .

**Beweis :**  $\bar{x}$  ist genau dann ein innerer Punkt von  $F$ , wenn die kleinste (im Sinne der Mengeninklusion) Seitenfläche von  $F$ , die  $\bar{x}$  enthält,  $F$  selbst ist. Offenbar ist  $\text{fa}(\text{eq}(\{\bar{x}\}))$  die minimale Seitenfläche von  $P$ , die  $\bar{x}$  enthält. Daraus folgt die Behauptung.  $\square$

**(7.16) Satz.** Ist  $F \neq \emptyset$  eine Seitenfläche des Polyeders  $P(A, b) \subseteq \mathbb{K}^n$ , dann gilt

$$\dim(F) = n - \text{rang}(A_{\text{eq}(F)}).$$

**Beweis:** Sei  $I := \text{eq}(F)$ . Aus der linearen Algebra wissen wir, dass  $n = \text{rang}(A_I) + \dim(\text{Kern}(A_I))$  gilt. Zu zeigen ist also:  $\dim(F) = \dim(\text{Kern}(A_I))$ . Seien  $r := \dim(\text{Kern}(A_I))$  und  $s := \dim(F)$ .

“ $r \geq s$ ”: Da  $\dim(F) = s$ , gibt es  $s+1$  affin unabhängige Vektoren  $x_0, x_1, \dots, x_s \in F$ . Dann sind die Vektoren  $x_1 - x_0, \dots, x_s - x_0$  linear unabhängig und erfüllen  $A_I(x_i - x_0) = 0$ . Kern  $(A_I)$  enthält also mindestens  $s$  linear unabhängige Vektoren, also gilt  $r \geq s$ .

“ $s \geq r$ ”: Nach (7.14) besitzt  $F$  einen inneren Punkt  $\bar{x} \in F$ . Nach (7.15) gilt  $\text{eq}(\{\bar{x}\}) = \text{eq}(F) = I$ , und daraus folgt für  $J := M \setminus I$ :

$$A_I \cdot \bar{x} = b_I,$$

$$A_J \cdot \bar{x} < b_J.$$

Ist  $r = 0$ , so gilt  $s \geq 0$  wegen  $\bar{x} \in F$ . Sei also  $r \geq 1$ , und  $\{x_1, \dots, x_r\}$  sei eine Basis von Kern  $(A_I)$ . Für  $p = 1, \dots, r$  und  $j \in J$  setze:

$$\delta_{jp} := \begin{cases} \infty, & \text{falls } A_j \cdot x_p = 0 \\ \frac{b_j - A_j \cdot \bar{x}}{A_j \cdot x_p} & \text{andernfalls} \end{cases}$$

$$\varepsilon := \min\{\delta_{jp} \mid j \in J, p \in \{1, \dots, r\}\}.$$

(Setze  $\varepsilon \neq 0$  beliebig, falls  $\delta_{jp} = \infty$  für alle  $j, p$ .) Für  $i \in I$  und alle  $p \in \{1, \dots, r\}$  gilt nun

$$A_i \cdot (\bar{x} + \varepsilon x_p) = A_i \cdot \bar{x} + \varepsilon A_i \cdot x_p = A_i \cdot \bar{x} = b_i,$$

da  $A_i \cdot x_p = 0$ . Für  $j \in J$  gilt

$$\begin{aligned} A_j \cdot (\bar{x} + \varepsilon x_p) &= A_j \cdot \bar{x} + \varepsilon A_j \cdot x_p \\ &\leq A_j \cdot \bar{x} + \delta_{jp} A_j \cdot x_p \\ &= A_j \cdot \bar{x} + b_j - A_j \cdot \bar{x} \\ &= b_j. \end{aligned}$$

Daraus folgt  $\bar{x} + \varepsilon x_p \in F$  für alle  $p \in \{1, \dots, r\}$ . Da die Vektoren  $\varepsilon x_1, \dots, \varepsilon x_r$  linear unabhängig sind, sind die Vektoren  $\bar{x}, \bar{x} + \varepsilon x_1, \dots, \bar{x} + \varepsilon x_r$  affin unabhängig. Das heißt,  $F$  enthält mindestens  $r + 1$  affin unabhängige Vektoren, und somit gilt  $\dim(F) = s \geq r$ .  $\square$

**(7.17) Folgerung.**  $P = P(A, b) \subseteq \mathbb{K}^n$  sei ein nichtleeres Polyeder, dann gilt:

- (a)  $\dim(P) = n - \text{rang}(A_{\text{eq}(P)})$ .
- (b) Ist  $\text{eq}(P) = \emptyset$ , dann ist  $P$  volldimensional (d. h.  $\dim(P) = n$ ).
- (c) Ist  $F$  eine echte Seitenfläche von  $P$ , dann gilt  $\dim(F) \leq \dim(P) - 1$ .

□

Mit Hilfe von Satz (7.16) kann man auch die affine Hülle einer Seitenfläche auf einfache Weise bestimmen.

**(7.18) Satz.** Sei  $F \neq \emptyset$  eine Seitenfläche des Polyeders  $P(A, b)$ , dann gilt

$$\text{aff}(F) = \{x \mid A_{\text{eq}(F)} \cdot x = b_{\text{eq}(F)}\}.$$

**Beweis :** Es seien  $I := \text{eq}(F)$  und  $T := \{x \mid A_I \cdot x = b_I\}$ . Offenbar ist  $T$  ein affiner Raum und wegen  $F \subseteq T$  gilt  $\text{aff}(F) \subseteq \text{aff}(T) = T$ . Sei  $s = \dim(F)$ , dann folgt aus Satz (7.16), dass  $\dim(\text{Kern}(A_I)) = s$  und somit  $\dim(T) = s$  gilt. Aus  $\dim(\text{aff}(F)) = \dim T$  und  $\text{aff}(F) \subseteq T$  folgt  $\text{aff}(F) = T$ . □

## 7.4 Facetten und Redundanz

Wie wir bereits bemerkt haben, kann man zu einem Ungleichungssystem  $Ax \leq b$  beliebig viele Ungleichungen hinzufügen, ohne die Lösungsmenge des Systems zu ändern. Wir wollen nun untersuchen, wie man ein gegebenes Polyeder mit möglichst wenigen Ungleichungen darstellen kann. Dies ist speziell für die lineare Optimierung wichtig, da der Rechenaufwand zur Auffindung einer Optimallösung in der Regel von der Anzahl der vorgelegten Ungleichungen abhängt. Gesucht wird also eine Minimaldarstellung eines Polyeders, um rechentechnische Vorteile zu haben. Es wird sich zeigen, dass hierbei diejenigen Ungleichungen, die maximale echte Seitenflächen eines Polyeders definieren, eine wesentliche Rolle spielen. Deshalb wollen wir derartige Seitenflächen untersuchen.

**(7.19) Definition.**  $Ax \leq b$  sei ein Ungleichungssystem, und  $M$  sei die Zeilenindexmenge von  $A$ .

- (a) Sei  $I \subseteq M$ , dann heißt das System  $A_I \cdot x \leq b_I$  **unwesentlich** oder **redundant bezüglich**  $Ax \leq b$ , wenn  $P(A, b) = P(A_{M \setminus I}, b_{M \setminus I})$  gilt.

- (b) Enthält  $Ax \leq b$  ein unwesentliches Teilsystem  $A_I x \leq b_I$ , dann heißt  $Ax \leq b$  **redundant**, andernfalls **irredundant**.
- (c) Eine Ungleichung  $A_i x \leq b_i$  heißt **wesentlich** oder **nicht redundant** bezüglich  $Ax \leq b$ , wenn  $P(A, b) \neq P(A_{M \setminus \{i\}}, b_{M \setminus \{i\}})$  gilt.
- (d) Eine Ungleichung  $A_i x \leq b_i$  heißt **implizite Gleichung** bezüglich  $Ax \leq b$ , wenn  $i \in \text{eq}(P(A, b))$  gilt.
- (e) Ein System  $Ax \leq a, Bx = b$  heißt **irredundant**, wenn  $Ax \leq a$  keine unwesentliche Ungleichung bezüglich des Systems  $Ax \leq a, Bx \leq b, -Bx \leq -b$  enthält und  $B$  vollen Zeilenrang hat.
- (f) Eine nichttriviale Seitenfläche  $F$  von  $P(A, b)$  heißt **Facette** von  $P(A, b)$ , falls  $F$  in keiner anderen echten Seitenfläche von  $P(A, b)$  enthalten ist.

□

Wir weisen ausdrücklich darauf hin, dass Redundanz bzw. Irredundanz keine Eigenschaft des Polyeders  $P(A, b)$  ist, sondern eine Eigenschaft des Ungleichungssystems  $Ax \leq b$ . Wir werden sehen, dass ein Polyeder viele irredundante Beschreibungen haben kann. Ferner ist auch die Annahme falsch, dass man immer durch das gleichzeitige Weglassen aller unwesentlichen Ungleichungen eines Systems  $Ax \leq b$  eine irredundante Beschreibung von  $P(A, b)$  enthält. Wir wollen nun zunächst unwesentliche Ungleichungen charakterisieren.

**(7.20) Satz.** Ein Ungleichungssystem  $A_I x \leq b_I$  ist unwesentlich bezüglich  $Ax \leq b$  genau dann, wenn es eine Matrix  $U \in \mathbb{K}_+^{(I, m)}$  gibt mit  $UA = A_I$ ,  $Ub \leq b_I$  und  $U_I = 0$ .

**Beweis:** Für jedes  $i \in I$  ist nach (7.4) die Ungleichung  $A_i x \leq b_i$  gültig bezüglich  $P(A_{M \setminus I}, b_{M \setminus I})$  genau dann, wenn es einen Vektor  $\bar{u}_i \geq 0$  gibt mit  $\bar{u}_i^T A_{M \setminus I} = A_i$  und  $\bar{u}_i^T b \leq b_i$ . Dies ist genau dann der Fall, wenn es eine Matrix  $U \in \mathbb{K}_+^{(I, m)}$  gibt mit  $UA = A_I$ ,  $Ub \leq b_I$  und  $U_I = 0$ . □

**(7.21) Folgerung.**  $A_i x \leq b_i$  ist genau dann redundant bezüglich  $Ax \leq b$ , wenn es einen Vektor  $u \in \mathbb{K}_+^m$  gibt mit  $u^T A = A_i$ ,  $u^T b \leq b_i$ ,  $u_i = 0$ .

□

Der nächste Satz zeigt, wann man Ungleichungen nicht mehr weglassen kann, ohne das Polyeder zu ändern.

**(7.22) Satz.**  $P = P(A, b) \neq \emptyset$  sei ein Polyeder. Sei  $\emptyset \neq I \subseteq M \setminus \text{eq}(P)$  und  $P' := P(A_{M \setminus I}, b_{M \setminus I})$ . Dann gilt

$$P \neq P' \iff \exists \text{ nichttriviale Seitenfläche } F \subseteq P \text{ mit } \text{eq}(F) \subseteq I \cup \text{eq}(P).$$

**Beweis:** Im Weiteren bezeichnen wir mit  $\text{eq}_{P'}$  die “equality set”-Abbildung bezüglich  $P'$ . Es gilt offenbar  $\text{eq}_{P'}(F) \subseteq \text{eq}(F)$ .

“ $\Leftarrow$ ” Angenommen, es gilt  $P = P'$ , und  $F$  sei eine beliebige nichttriviale Seitenfläche von  $P$  (und somit auch von  $P'$ ). Da  $F$  eine nichttriviale Seitenfläche von  $P'$  ist, gibt es ein  $i \in (M \setminus I) \setminus \text{eq}_{P'}(P')$  mit  $i \in \text{eq}_{P'}(F)$ . Daraus folgt  $\text{eq}(F) \not\subseteq I \cup \text{eq}(P)$ .

“ $\Rightarrow$ ” Angenommen, es gilt  $P \neq P'$ . Wegen  $P \subseteq P'$  heißt dies, es existiert ein Vektor  $v \in P' \setminus P$ , und somit gibt es eine Indexmenge  $\emptyset \neq K \subseteq I$  mit der Eigenschaft  $A_i \cdot v \leq b_i \ \forall i \in M \setminus K$  und  $A_i \cdot v > b_i \ \forall i \in K$ . Nach Satz (7.14) hat  $P$  einen inneren Punkt, sagen wir  $w$ , d. h. es gilt  $A_i \cdot w = b_i \ \forall i \in \text{eq}(P)$  und  $A_i \cdot w < b_i \ \forall i \in M \setminus \text{eq}(P)$ .

Wir betrachten nun einen Punkt  $y$  auf der Strecke zwischen  $v$  und  $w$ , d. h.  $y = \lambda w + (1 - \lambda)v$  mit  $0 \leq \lambda \leq 1$ . Aufgrund der Voraussetzungen gilt:

$$\begin{aligned} A_i \cdot y &= b_i \quad \forall i \in \text{eq}(P) \\ A_i \cdot y &< b_i \quad \forall i \in M \setminus (\text{eq}(P) \cup K), \text{ falls } \lambda > 0. \end{aligned}$$

Ist  $i \in K$ , so gilt

$$\begin{aligned} A_i \cdot y \leq b_i &\iff \lambda A_i \cdot w + (1 - \lambda) A_i \cdot v \leq b_i \\ &\iff \lambda A_i \cdot (w - v) \leq b_i - A_i \cdot v \\ &\iff \lambda \geq \frac{b_i - A_i \cdot v}{A_i \cdot (w - v)} \quad (\text{da } A_i \cdot (w - v) < 0). \end{aligned}$$

Setzen wir

$$\begin{aligned} \mu &:= \max \left\{ \frac{b_i - A_i \cdot v}{A_i \cdot (w - v)} \mid i \in K \right\}, \\ L &:= \left\{ i \in K \mid \mu = \frac{b_i - A_i \cdot v}{A_i \cdot (w - v)} \right\}, \end{aligned}$$

dann gilt  $z := \mu w + (1 - \mu)v \in P$ ,  $\emptyset \neq L \subseteq K \subseteq I$  und

$$\begin{aligned} A_i \cdot z &= b_i \quad \forall i \in L \\ A_i \cdot z &< b_i \quad \forall i \in K \setminus L. \end{aligned}$$

Daraus folgt,  $z$  ist ein innerer Punkt von  $F := \text{fa}(L \cup \text{eq}(P))$ . Nach (7.15) gilt dann  $\text{eq}(F) = \text{eq}(\{z\}) = L \cup \text{eq}(P)$ , und das bedeutet, dass  $F$  eine nichttriviale Seitenfläche von  $P$  mit  $\text{eq}(F) \subseteq I \cup \text{eq}(P)$  ist.  $\square$

Wir werden nun wichtige Eigenschaften von Facetten bestimmen, Nichtredundanz kennzeichnen und Facetten charakterisieren.

**(7.23) Satz.** Sei  $F$  eine Facette von  $P = P(A, b)$ , dann gilt:

(a)  $\text{eq}(P) \subset \text{eq}(F)$ .

(b) Für alle  $i \in \text{eq}(F) \setminus \text{eq}(P)$  gilt

$$F = \text{fa}(\{i\}) = \{x \in P \mid A_i \cdot x = b_i\}.$$

**Beweis :** (a) gilt offensichtlich für alle nichttrivialen Seitenflächen von  $P$ .

(b) Die Abbildung  $\text{fa}$  ist Inklusionsumkehrend, d. h.

$$I \subseteq J \Rightarrow \text{fa}(I) \supseteq \text{fa}(J).$$

Daraus folgt  $F = \text{fa}(\text{eq}(F)) \subseteq \text{fa}(\{i\})$ . Da  $i \notin \text{eq}(P)$ , muss  $\text{fa}(\{i\})$  eine echte Seitenfläche von  $P$  sein. Aus der Maximalität von  $F$  folgt die Behauptung.  $\square$

**(7.24) Folgerung.** Sei  $P = P(A, b)$  ein Polyeder und  $\mathcal{F}$  die Menge der Facetten von  $P$ . Dann gilt:

(a)  $F_1, F_2 \in \mathcal{F}, F_1 \neq F_2 \implies \text{eq}(F_1) \cap \text{eq}(F_2) = \text{eq}(P)$ .

(b)  $|\mathcal{F}| \leq m - |\text{eq}(P)|$ .

(c) Es gibt eine Menge  $I \subseteq M$  mit folgenden Eigenschaften

(c<sub>1</sub>)  $I \subseteq M \setminus \text{eq}(P)$

(c<sub>2</sub>)  $|I| = |\mathcal{F}|$

(c<sub>3</sub>)  $F \in \mathcal{F} \iff \exists$  genau ein  $i \in I$  mit  $F = \text{fa}(\{i\})$ .

$\square$

Jede Menge  $I \subseteq M$  mit den Eigenschaften (c<sub>1</sub>), (c<sub>2</sub>), (c<sub>3</sub>) wollen wir **Facetten-Indexmenge** nennen. Satz (7.23) (b) zeigt, dass man Facetten von  $P$  dadurch erhält, dass man in nur einer Ungleichung  $A_i \cdot x \leq b_i$  des Systems  $Ax \leq b$  Gleichheit fordert. Jedoch ist es keineswegs so, dass für alle  $i \in M$  die Menge  $\text{fa}(\{i\})$  eine Facette von  $P$  ist! Dies gilt nur für solche  $i \in M$ , die in einer Facettenindexmenge enthalten sind.

**(7.25) Satz.** Seien  $P = P(A, b) \neq \emptyset$  ein Polyeder und  $\mathcal{F}$  die Menge der Facetten von  $P$ . Seien  $M$  die Zeilenindexmenge von  $A$ ,  $I \subseteq M \setminus \text{eq}(P)$  und  $J \subseteq \text{eq}(P)$ . Sei  $P' := \{x \mid A_J x = b_J, A_I x \leq b_I\}$ , dann gilt:

- (a)  $P = P' \iff$  (a<sub>1</sub>)  $\forall F \in \mathcal{F}$  gilt  $\mathcal{I} \cap \text{eq}(F) \neq \emptyset$  und  
 (a<sub>2</sub>)  $\text{rang}(A_J) = \text{rang}(A_{\text{eq}(P)})$ .
- (b)  $P = P(A_{I \cup \text{eq}(P)}, b_{I \cup \text{eq}(P)}) \iff \forall F \in \mathcal{F}$  gilt  $\mathcal{I} \cap \text{eq}(F) \neq \emptyset$ .

**Beweis :** Mit  $J = \text{eq}(P)$  folgt (b) direkt aus (a). Wir beweisen (a).

“ $\implies$ ” Nach Definition gilt offenbar  $J = \text{eq}_{P'}(P')$ . Angenommen (a<sub>2</sub>) ist nicht erfüllt, d. h.  $\text{rang}(A_J) < \text{rang}(A_{\text{eq}(P)})$ . Dann folgt aus der Dimensionsformel (7.17) (a)  $\dim(P') > \dim(P)$  und somit muss  $P \neq P'$  gelten. Widerspruch !

Angenommen (a<sub>1</sub>) ist nicht erfüllt. Dann gibt es eine Facette  $F$  von  $P$  mit  $\text{eq}(F) \subseteq M \setminus I = (M \setminus I) \cup \text{eq}(P)$ . Folglich gilt  $P \neq P'$  nach Satz (7.22). Widerspruch !

“ $\impliedby$ ” Wir zeigen zunächst, dass unter der Voraussetzung (a<sub>2</sub>) gilt:

$$A_J x = b_J \implies A_{\text{eq}(P)} x = b_{\text{eq}(P)}.$$

Da  $P' \neq \emptyset$ , gilt  $\text{rang}(A_J, b_J) = \text{rang}(A_J) = \text{rang}(A_{\text{eq}(P)}) = \text{rang}(A_{\text{eq}(P)}, b_{\text{eq}(P)})$ . Das heißt, für alle  $i \in \text{eq}(P)$  existieren  $K \subseteq J$  und  $\lambda_k, k \in K$ , mit  $A_i = \sum_{k \in K} \lambda_k A_k$ ,  $b_i = \sum_{k \in K} \lambda_k b_k$ . Erfüllt also der Vektor  $x$  das System  $A_J x = b_J$ , so gilt für alle  $i \in \text{eq}(P)$

$$A_i x = \sum_{k \in K} \lambda_k A_k x = \sum_{k \in K} \lambda_k b_k = b_i.$$

Nach (a<sub>1</sub>) gilt für jede Facette  $F$  von  $P$ :  $\text{eq}(F) \not\subseteq M \setminus I$ , und da Facetten maximale echte Seitenflächen sind und  $\text{eq}$  inklusionsumkehrend ist, folgt daraus  $\text{eq}(G) \not\subseteq M \setminus I$  für alle echten Seitenflächen  $G$  von  $P$ . Aus Satz (7.22) folgt daher  $P = P'$ .  $\square$

**(7.26) Folgerung.** Seien  $P = P(A, b) \neq \emptyset$  ein Polyeder,  $I \subseteq M \setminus \text{eq}(P)$ ,  $J \subseteq \text{eq}(P)$  und  $P = \{x \mid A_J x = b_J, A_I x \leq b_I\}$ . Diese Darstellung von  $P$  ist genau dann irredundant, wenn gilt:

- (a)  $I$  ist eine Facetten-Indexmenge von  $P$ .  
 (b)  $A_J$  ist eine  $(\text{rang}(A_{\text{eq}(P)}), n)$ -Matrix mit vollem Zeilenrang.

$\square$

**(7.27) Folgerung.** Sei  $P = P(A, b) \subseteq \mathbb{K}^n$  ein volldimensionales Polyeder (also  $\text{eq}(P) = \emptyset$ , bzw.  $\dim(P) = n$ ), dann gilt für alle  $I \subseteq M$

$$P(A_I, b_I) \text{ ist eine irredundante Beschreibung von } P \\ \iff I \text{ ist Facetten-Indexmenge von } P.$$

□

**(7.28) Satz.** Sei  $P = P(A, b)$  ein Polyeder, und  $F$  sei eine nichttriviale Seitenfläche von  $P$ . Dann sind äquivalent:

- (i)  $F$  ist eine Facette von  $P$ .
- (ii)  $F$  ist eine maximale echte Seitenfläche von  $P$ .
- (iii)  $\dim(F) = \dim(P) - 1$ .
- (iv)  $F$  enthält  $\dim(P)$  affin unabhängige Vektoren.
- (v) Sei  $c^T x \leq \gamma$  eine bezüglich  $P$  gültige Ungleichung mit  $F = \{x \in P \mid c^T x = \gamma\}$ , dann gilt für alle gültigen Ungleichungen  $d^T x \leq \delta$  mit  $F \subseteq \{x \in P \mid d^T x = \delta\}$ : Es gibt einen Vektor  $u \in \mathbb{K}^{\text{eq}(P)}$  und  $\alpha \in \mathbb{K}, \alpha \geq 0$  mit

$$\begin{aligned} d^T &= \alpha c^T + u^T A_{\text{eq}(P)}, \\ \delta &= \alpha \gamma + u^T b_{\text{eq}(P)}. \end{aligned}$$

**Beweis :** (i)  $\iff$  (ii): nach Definition.

(iv)  $\iff$  (iii): trivial.

(iii)  $\implies$  (ii): Angenommen  $F$  ist keine Facette, dann existiert eine echte Seitenfläche  $G$  von  $P$  mit  $F \subset G \subset P$ . Aus (7.17) (c) folgt dann  $\dim(F) \leq \dim(G) - 1 \leq \dim(P) - 2$ , Widerspruch!

(i)  $\implies$  (v): Sei  $F$  eine beliebige Facette von  $P$ . Wir nehmen zunächst an, dass  $A_I x \leq b_I, A_J x = b_J$  eine irredundante Darstellung von  $P(A, b)$  ist mit  $1 \in I$  und dass  $F = \{x \in P \mid A_1 x = b_1\}$  gilt. Sei nun  $d^T x \leq \delta$  eine gültige Ungleichung mit  $F \subseteq \{x \in P \mid d^T x = \delta\}$ . Aufgrund von Folgerung (7.4) gibt es Vektoren  $v \geq 0$  und  $w$  mit  $v^T A_I + w^T A_J = d^T$  und  $v^T b_I + w^T b_J \leq \delta$  (in der Tat gilt hier Gleichheit, da  $\{x \mid d^T x = \delta\}$  eine Stützhyperebene ist). Angenommen, es gibt einen Index  $i \in I \setminus \{1\}$  mit  $v_i > 0$ , dann gilt nach (7.10)  $i \in \text{eq}(F)$ . Dies ist aber ein Widerspruch dazu, dass  $I$  eine Facettenindexmenge ist. Hieraus folgt (v).

(v)  $\implies$  (iii): Da  $F$  eine echte Seitenfläche von  $P$  ist, gilt  $\dim(F) \leq \dim(P) - 1$ . Angenommen  $\dim(F) \leq \dim(P) - 2$ . O. b. d. A. können wir annehmen, dass  $F = \{x \in P \mid A_1 x = b_1\}$  gilt. Aus (7.16) folgt

$$\text{rang}(A_{\text{eq}(F)}) \geq \text{rang}(A_{\text{eq}(P)}) + 2.$$

Mithin gibt es einen Index  $i \in \text{eq}(F) \setminus (\text{eq}(P) \cup \{1\})$ , so dass der Zeilenvektor  $A_i$  linear unabhängig von den Zeilenvektoren  $A_j, j \in \text{eq}(P) \cup \{1\}$ , ist. Das aber heißt, dass das System

$$A_i = \alpha A_1 + u^T A_{\text{eq}(P)}$$

keine Lösung  $\alpha, u$  hat. Wegen  $F \subseteq \{x \in P \mid A_i x = b_i\}$  ist dies ein Widerspruch zu (v).  $\square$

**(7.29) Folgerung.** Seien  $P = P(A, b) \subseteq \mathbb{R}^n$  ein volldimensionales Polyeder und  $F = \{x \in P \mid c^T x = \gamma\}$  eine Seitenfläche von  $P$ . Dann sind äquivalent:

(i)  $F$  ist Facette von  $P$ .

(ii)  $\dim(F) = n - 1$ .

(iii) Für alle gültigen Ungleichungen  $d^T x \leq \delta, d \neq 0$ , mit  $F \subseteq \{x \in P \mid d^T x = \delta\}$  gilt: Es existiert ein  $\alpha > 0$  mit

$$\begin{aligned} d^T &= \alpha c^T, \\ \delta &= \alpha \gamma. \end{aligned}$$

$\square$

**(7.30) Beispiel.** Wir betrachten das Polyeder  $P = P(A, b) \subseteq \mathbb{R}^2$ , das wie folgt gegeben ist.

$$A = \begin{pmatrix} A_1. \\ A_2. \\ \cdot \\ \cdot \\ \cdot \\ A_6. \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ -2 & 2 \\ 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ -1 & -1 \end{pmatrix} \quad b = \begin{pmatrix} 0 \\ 0 \\ 2 \\ 2 \\ -1 \\ -2 \end{pmatrix}.$$

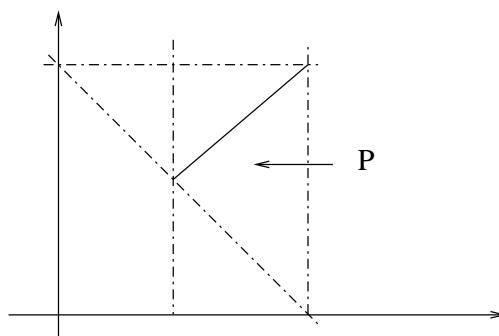


Abb. 7.2

$P$  hat 4 Seitenflächen, nämlich  $\emptyset$ ,  $P$  und  $F_1 = \left\{ \binom{2}{2} \right\}$ ,  $F_2 = \left\{ \binom{1}{1} \right\}$ .  $F_1$  und  $F_2$  sind Facetten von  $P$ . Es gilt  $\text{eq}(P) = \{1, 2\}$ ,  $\text{eq}(F_1) = \{1, 2, 3, 4\}$ ,  $\text{eq}(F_2) = \{1, 2, 5, 6\}$ ,  $\text{eq}(\emptyset) = \{1, \dots, 6\}$ . Die Mengen  $\{3, 5\}$ ,  $\{3, 6\}$ ,  $\{4, 5\}$ ,  $\{4, 6\}$  sind die Facettenindexmengen von  $P$ . Eine irredundante Beschreibung von  $P$  ist z. B. gegeben durch

$$P = \{x \mid A_1 \cdot x = 0, A_3 \cdot x \leq b_3, A_5 \cdot x \leq b_5\}.$$

Übrigens sind die Ungleichungen  $A_i \cdot x \leq b_i$ ,  $i = 3, 4, 5, 6$  redundant bezüglich  $P(A, b)$ . Die Ungleichungssysteme  $A_I \cdot x \leq b_I$  mit  $I = \{3, 5\}$  oder  $I = \{4, 6\}$  sind z. B. ebenfalls redundant. Aber  $A_I \cdot x \leq b_I$  ist nicht redundant bezüglich  $P(A, b)$ , falls  $I = \{3, 4, 5\}$ .  $\square$

# Kapitel 8

## Ecken und Extremalen

Im vorhergehenden Kapitel haben wir allgemeine Seitenflächen und speziell maximale Seitenflächen (Facetten) von Polyedern behandelt. Die Facetten sind bei der äußeren Beschreibung von Polyedern wichtig. Wir werden nun minimale Seitenflächen untersuchen und sehen, dass sie bei der inneren Darstellung von Bedeutung sind.

**(8.1) Definition.** Es seien  $P \subseteq \mathbb{K}^n$  ein Polyeder und  $F$  eine Seitenfläche von  $P$ . Gibt es  $x \in \mathbb{K}^n$  und  $0 \neq z \in \mathbb{K}^n$  mit

$$\left\{ \begin{array}{l} F = \{x\} \\ F = x + \text{lin}(\{z\}) \\ F = x + \text{cone}(\{z\}) \end{array} \right\}, \text{ so heißt } F \left\{ \begin{array}{l} \text{Ecke} \\ \text{Extremallinie} \\ \text{Extremalstrahl} \end{array} \right\}.$$

□

Ist  $F = \{x\}$  eine Ecke von  $P$ , so sagen wir einfach “ $x$  ist eine Ecke von  $P$ ”. Wir werden uns im weiteren insbesondere für Ecken interessieren und sie charakterisieren. Offenbar sind Ecken Seitenflächen der Dimension 0, während Extremallinien und Extremalstrahlen Seitenflächen der Dimension 1 sind. Allgemein heißen Seitenflächen der Dimension 1 **Kanten**. Kanten sind entweder Extremallinien, Extremalstrahlen oder Verbindungsstrecken zwischen zwei Ecken. Sind zwei Ecken  $x, y$  eines Polyeders  $P$  durch eine Kante verbunden, d. h.  $\text{conv}(\{x, y\})$  ist eine Seitenfläche von  $P$ , so nennt man  $x$  und  $y$  **adjazent** auf  $P$ .

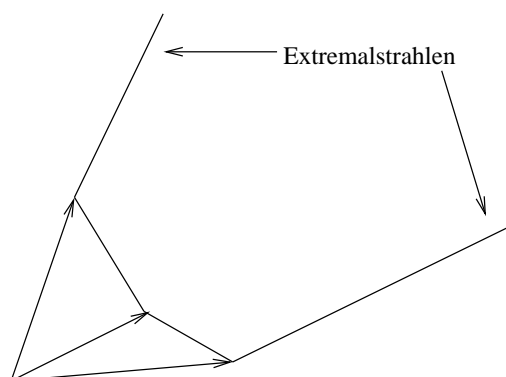


Abb. 8.1

## 8.1 Rezessionskegel, Linienraum, Homogenisierung

An dieser Stelle ist es nützlich einige weitere Objekte einzuführen, die man Polyedern (bzw. allgemeinen Mengen) zuordnen kann. Das Studium dieser Objekte ist für sich selbst betrachtet sehr interessant. Wir wollen diese Mengen jedoch nur als Hilfsmittel zur Vereinfachung von Beweisen verwenden, weswegen wir nicht weiter auf theoretische Untersuchungen dieser Mengen eingehen werden.

**(8.2) Definition.** Sei  $S \subseteq \mathbb{K}^n$  eine beliebige Menge. Wir definieren

- (a)  $\text{rec}(S) := \{y \in \mathbb{K}^n \mid \exists x \in S, \text{ so dass } \forall \lambda \geq 0 \text{ gilt } x + \lambda y \in S\}$ ,
- (b)  $\text{lineal}(S) := \{y \in \mathbb{K}^n \mid \exists x \in S, \text{ so dass } \forall \lambda \in \mathbb{K} \text{ gilt } x + \lambda y \in S\}$ ,
- (c)  $\text{hog}(S) := \left\{ \begin{pmatrix} x \\ 1 \end{pmatrix} \in \mathbb{K}^{n+1} \mid x \in S \right\}^{\circ\circ}$ .

Die Menge  $\text{rec}(S)$  heißt **Rezessionskegel** von  $S$ ,  $\text{lineal}(S)$  heißt **Linealitätsraum** oder **Linienraum** von  $S$ , und  $\text{hog}(S)$  heißt **Homogenisierung** von  $S$ .

□

Wir wollen nun die oben eingeführten Mengen bezüglich Polyedern charakterisieren. Nennen wir einen Vektor  $y$  mit  $x + \lambda y \in S$  für alle  $\lambda \geq 0$  eine ‘Richtung nach Unendlich’, so besteht der Rezessionskegel einer Menge  $S$  aus allen Richtungen nach Unendlich. Für Polyeder gilt Folgendes:

**(8.3) Satz.** Sei  $P = P(A, b) = \text{conv}(V) + \text{cone}(E)$  ein nichtleeres Polyeder, dann gilt

$$\text{rec}(P) = P(A, 0) = \text{cone}(E).$$

**Beweis :** (a)  $\text{rec}(P) = P(A, 0)$ .

Ist  $y \in \text{rec}(P)$ , so existiert ein  $x \in P$  mit  $x + \lambda y \in P \forall \lambda \geq 0$ . Daraus folgt  $b \geq A(x + \lambda y) = Ax + \lambda Ay$ . Gäbe es eine Komponente von  $Ay$ , die größer als Null ist, sagen wir  $(Ay)_i > 0$ , so wäre der Vektor  $x + \lambda_0 y$  mit

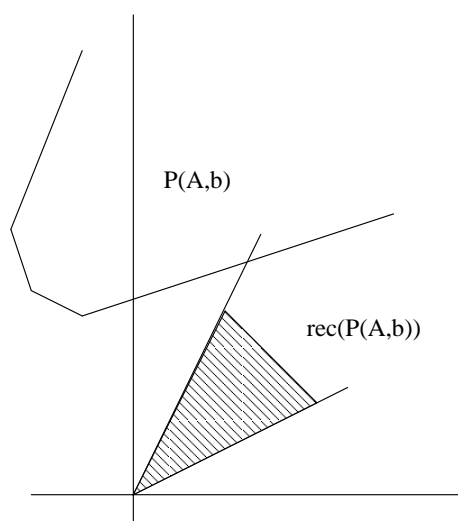
$$\lambda_0 = \frac{b_i - (Ax)_i}{(Ay)_i} + 1$$

nicht in  $P(A, b)$ , Widerspruch!

Ist  $y \in P(A, 0)$ , so gilt für alle  $x \in P(A, b)$  und  $\lambda \geq 0$ ,  $A(x + \lambda y) = Ax + \lambda Ay \leq b + 0 = b$ , also ist  $y \in \text{rec}(P)$ .

(b)  $P(A, 0) = \text{cone}(E)$ . Trivial!  $\square$

Insbesondere folgt aus (8.3), dass  $\text{rec}(P) = \{y \in \mathbb{K}^n \mid \forall x \in P \text{ und } \forall \lambda \geq 0 \text{ gilt } x + \lambda y \in P\}$ . Ist  $P$  ein Kegel, so gilt natürlich  $P = \text{rec}(P)$ , und offenbar ist ein Polyeder  $P$  genau dann ein Polytop, wenn  $\text{rec}(P) = \{0\}$ . Abbildung 8.2 zeigt ein Polyeder und seinen Rezessionskegel.



**Abb. 8.2**

Aus Definition (8.2) folgt  $\text{lineal}(P) = \text{rec}(P) \cap (-\text{rec}(P))$ . Offenbar ist  $\text{lineal}(P)$  ein linearer Teilraum des  $\mathbb{K}^n$ , und zwar ist es der größte lineare Teilraum  $L \subseteq \mathbb{K}^n$ , so dass  $x + L \subseteq P$  für alle  $x \in P$  gilt. Analytisch können wir  $\text{lineal}(P)$  wie folgt darstellen.

**(8.4) Satz.** Sei  $P = P(A, b) = \text{conv}(V) + \text{cone}(E)$  ein nichtleeres Polyeder,

dann gilt

$$\text{lineal}(P) = \{x \mid Ax = 0\} = \text{cone}(\{e \in E \mid -e \in \text{cone}(E)\}).$$

**Beweis :** Wegen  $\text{lineal}(P) = \text{rec}(P) \cap (-\text{rec}(P))$  folgt die Behauptung direkt aus (8.3).  $\square$

Wir kommen nun zur Homogenisierung. Die Definition der Homogenisierung erscheint etwas kompliziert: Man wende zweimal die Kegelpolarität auf die Menge  $S$  an! Geometrisch betrachtet ist  $\text{hog}(S)$  der Durchschnitt aller Ungleichungen mit rechter Seite 0, die gültig bezüglich  $\{(x) \mid x \in S\}$  sind.

**(8.5) Satz.** Sei  $P = P(A, b) = \text{conv}(V) + \text{cone}(E)$  ein nichtleeres Polyeder. Sei

$$B = \begin{pmatrix} A, & -b \\ 0, & -1, \end{pmatrix}$$

dann gilt

$$\text{hog}(P) = P(B, 0) = \text{cone}(\{(v) \mid v \in V\}) + \text{cone}(\{(e) \mid e \in E\}).$$

**Beweis :** Setzen wir  $P_1 := \{(x) \in \mathbb{K}^{n+1} \mid x \in P\}$ , so gilt offensichtlich

$$P_1 = \text{conv}\{(v) \mid v \in V\} + \text{cone}\{(e) \mid e \in E\}.$$

Aus Folgerung (7.4) (iii) ergibt sich dann:

$$\begin{aligned} P_1^\circ &= \{z \in \mathbb{K}^{n+1} \mid z^T u \leq 0 \forall u \in P_1\} \\ &= \{z \in \mathbb{K}^{n+1} \mid z^T (v) \leq 0 \forall v \in V, z^T (e) \leq 0 \forall e \in E\} \\ &= \left\{ z \mid \begin{pmatrix} V^T, & \mathbb{1} \\ E^T, & 0 \end{pmatrix} z \leq 0 \right\} = P \left( \begin{pmatrix} V^T, & \mathbb{1} \\ E^T, & 0 \end{pmatrix}, 0 \right). \end{aligned}$$

Mit Folgerung (6.7)  $P(A, 0)^\circ = \text{cone}(A^T)$  erhalten wir nun

$$\text{hog}(P) = P_1^{\circ\circ} = P \left( \begin{pmatrix} V^T, & \mathbb{1} \\ E^T, & 0 \end{pmatrix}, 0 \right)^\circ = \text{cone} \begin{pmatrix} V & E \\ \mathbb{1} & 0^T \end{pmatrix}.$$

Die zweite Charakterisierung von  $\text{hog}(P)$  folgt aus einer anderen Darstellung von  $P_1^\circ$ . Es gilt nämlich mit Satz (7.2):

$$\begin{aligned}
P_1^\circ &= \left\{ \begin{pmatrix} y \\ \lambda \end{pmatrix} \in \mathbb{K}^{n+1} \mid y^T x + \lambda 1 \leq 0 \forall x \in P \right\} = \left\{ \begin{pmatrix} y \\ \lambda \end{pmatrix} \in \mathbb{K}^{n+1} \mid y^T x \leq -\lambda \forall x \in P \right\} \\
&= \left\{ \begin{pmatrix} y \\ \lambda \end{pmatrix} \in \mathbb{K}^{n+1} \mid \begin{pmatrix} y \\ -\lambda \end{pmatrix} \in P^\gamma \right\} = \left\{ \begin{pmatrix} y \\ \lambda \end{pmatrix} \in \mathbb{K}^{n+1} \mid \begin{pmatrix} y \\ -\lambda \end{pmatrix} \in \text{cone} \begin{pmatrix} A^T & 0 \\ b^T & 1 \end{pmatrix} \right\} \\
&= \text{cone} \begin{pmatrix} A^T & 0 \\ -b^T & -1 \end{pmatrix}.
\end{aligned}$$

Folgerung (6.10) impliziert nun

$$\text{hog}(P) = P_1^{\circ\circ} = \left( \text{cone} \begin{pmatrix} A^T & 0 \\ -b^T & -1 \end{pmatrix} \right)^\circ = P \left( \begin{pmatrix} A & -b \\ 0 & -1 \end{pmatrix}, 0 \right).$$

□

In Abbildung 8.3 sind ein Polyeder  $P \subseteq \mathbb{R}^1$ , die im obigen Beweis definierte Menge  $P_1$  und  $\text{hog}(P)$  dargestellt.

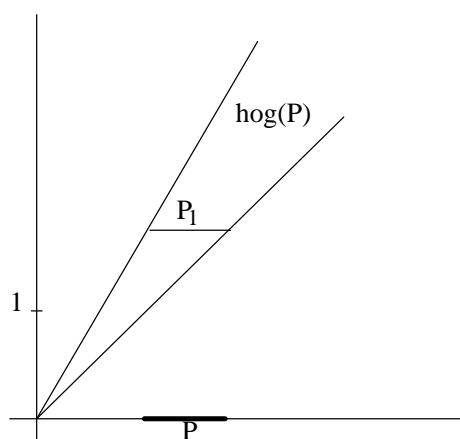


Abb. 8.3

**(8.6) Bemerkung.** Sei  $P \subseteq \mathbb{K}^n$  ein Polyeder, dann gilt:

$$(a) \quad x \in P \iff \begin{pmatrix} x \\ 1 \end{pmatrix} \in \text{hog}(P).$$

$$(b) \quad x \in \text{rec}(P) \iff \begin{pmatrix} x \\ 0 \end{pmatrix} \in \text{hog}(P).$$

**Beweis:** (a) ist trivial.

(b) Sei  $P = P(A, b)$  eine Darstellung von  $P$ , dann gilt  $\text{hog}(P) = P(B, 0)$  mit  $B = \begin{pmatrix} A & -b \\ 0 & -1 \end{pmatrix}$ . Folglich gilt nach (8.5) und (8.3)

$$\begin{pmatrix} x \\ 0 \end{pmatrix} \in \text{hog}(P) \iff \begin{pmatrix} x \\ 0 \end{pmatrix} \in P(B, 0) \iff Ax \leq 0 \iff x \in \text{rec}(P).$$

□

## 8.2 Charakterisierungen von Ecken

Wir beginnen mit der Kennzeichnung von Ecken und Extremalstrahlen von Kegeln.

**(8.7) Satz.** Sei  $K \subseteq \mathbb{K}^n$  ein polyedrischer Kegel, dann gilt:

- (a) Ist  $x$  Ecke von  $K$ , dann gilt  $x = 0$ .
- (b)  $F$  ist ein Extremalstrahl von  $K \iff \exists z \in \mathbb{K}^n \setminus \{0\}$  so dass  $F = \text{cone}(\{z\})$  eine Seitenfläche von  $K$  ist.

**Beweis :** Nach (2.6) gibt es eine Matrix  $A$  mit  $K = P(A, 0)$ . Aufgrund von (7.11) ist daher jede nichtleere Seitenfläche von  $K$  ein Kegel, der den Nullvektor enthält.

(a) Ist  $x$  Ecke, dann gilt also  $0 \in \{x\}$  und somit  $x = 0$ .

(b) " $\Leftarrow$ " nach Definition.

" $\Rightarrow$ " Nach (8.1) existieren  $x \in \mathbb{K}^n$  und  $z \in \mathbb{K}^n \setminus \{0\}$  mit  $F = x + \text{cone}(\{z\})$ . Nach obiger Bemerkung gilt  $0 \in F$ , und somit existiert ein  $\bar{\lambda} \geq 0$  mit  $0 = x + \bar{\lambda}z$ , d. h.  $x = -\bar{\lambda}z$ . Da  $F$  ein Kegel ist, gilt  $\{\lambda x \mid \lambda \in \mathbb{K}_+\} = \{\lambda(-\bar{\lambda}z) \mid \lambda \in \mathbb{K}_+\} \subseteq F$ . Falls  $\bar{\lambda} \neq 0$ , gilt aber  $\{-\lambda z \mid \lambda \geq 0\} \not\subseteq x + \text{cone}(\{z\})$ , ein Widerspruch. Daraus folgt  $x = 0$ .  $\square$

Der folgende Hilfssatz über innere Punkte wird im Weiteren benötigt.

**(8.8) Lemma.** Ist  $F$  eine nichtleere Seitenfläche von  $P = P(A, b)$ , gilt  $I = \text{eq}(F)$ , und ist  $B = \{y^1, \dots, y^k\}$  eine Basis des Kerns von  $A_I$ , dann gibt es zu jedem inneren Punkt  $x \in F$  von  $F$  ein  $\varepsilon > 0$ , so dass  $x \pm \varepsilon y^j \in P$  für alle  $j = 1, \dots, k$  gilt.

**Beweis :** Sei  $x$  innerer Punkt von  $F$  und  $J = \{1, \dots, m\} \setminus I$ . Nach (7.15) gilt  $\text{eq}(\{x\}) = I$ , also  $A_J x < b_J$  und ferner  $A_I(x \pm \varepsilon y^j) = b_I$  für alle  $\varepsilon \in \mathbb{K}$ . Für  $\varepsilon$  genügend klein gilt dann offenbar auch  $A_J(x \pm \varepsilon y^j) \leq b_J$ .  $\square$

Wir wollen nun die Ecken eines Polyeders charakterisieren.

**(8.9) Satz.** Seien  $P = P(A, b) \subseteq \mathbb{K}^n$  ein Polyeder und  $x \in P$ . Dann sind die folgenden Aussagen äquivalent:

- (1)  $x$  ist eine Ecke von  $P$ .

- (2)  $\{x\}$  ist eine nulldimensionale Seitenfläche von  $P$ .
- (3)  $x$  ist keine echte Konvexkombination von Elementen von  $P$ , d. h.  $y, z \in P$ ,  $y \neq z$ ,  $0 < \lambda < 1 \implies x \neq \lambda y + (1 - \lambda)z$ .
- (4)  $P \setminus \{x\}$  ist konvex.
- (5)  $\text{rang}(A_{\text{eq}(\{x\})}) = n$ .
- (6)  $\exists c \in \mathbb{K}^n \setminus \{0\}$ , so dass  $x$  die eindeutig bestimmte Optimallösung des linearen Programms  $\max c^T y$ ,  $y \in P$  ist.

**Beweis :** (1)  $\iff$  (2). Definition!

(1)  $\implies$  (6). Nach Definition ist  $\{x\}$  eine Seitenfläche, also existiert eine bezüglich  $P$  gültige Ungleichung  $c^T y \leq \gamma$ , so dass  $\{y \in P \mid c^T y = \gamma\} = \{x\}$  gilt. Folglich ist  $x$  die eindeutig bestimmte Optimallösung von  $\max c^T y$ ,  $y \in P$ . Ist  $P \neq \{x\}$ , so ist  $c \neq 0$ , andernfalls kann  $c \neq 0$  gewählt werden.

(6)  $\implies$  (3). Sei  $x$  die eindeutige Optimallösung von  $\max c^T u$ ,  $u \in P$  mit Wert  $\gamma$ . Gilt  $x = \lambda y + (1 - \lambda)z$  für  $y, z \in P$ ,  $y \neq z$ ,  $0 < \lambda < 1$ , dann folgt  $\gamma = c^T x = c^T(\lambda y + (1 - \lambda)z) = \lambda c^T y + (1 - \lambda)c^T z \leq \lambda \gamma + (1 - \lambda)\gamma = \gamma$ . Dies impliziert  $c^T y = \gamma = c^T z$ . Widerspruch!

(3)  $\iff$  (4) trivial.

(3)  $\implies$  (5). Sei  $I = \text{eq}(\{x\})$ , dann ist  $x$  innerer Punkt von  $F = \{y \in P \mid A_I y = b_I\}$ . Angenommen  $\text{rang}(A_I) < n$ , dann existiert ein Vektor  $y \neq 0$  im Kern von  $A_I$ , und nach Lemma (8.8) ein  $\varepsilon > 0$ , so dass  $x \pm \varepsilon y \in P$ . Daraus folgt  $x = \frac{1}{2}(x + \varepsilon y) + \frac{1}{2}(x - \varepsilon y)$ , also ist  $x$  echte Konvexkombination von Elementen von  $P$ , Widerspruch!

(5)  $\implies$  (2). Sei  $I = \text{eq}(\{x\})$ , dann hat  $A_I y = b_I$  nach Voraussetzung eine eindeutig bestimmte Lösung, nämlich  $x$ . Daraus folgt  $F = \{y \in P \mid A_I y = b_I\} = \{x\}$ , also ist  $\{x\}$  eine nulldimensionale Seitenfläche.  $\square$

Für Polyeder der Form  $P^=(A, b) = \{x \in \mathbb{K}^n \mid Ax = b, x \geq 0\}$  gibt es eine weitere nützliche Kennzeichnung der Ecken. Ist  $x \in \mathbb{K}^n$ , so setzen wir

$$\text{supp}(x) := \{i \in \{1, \dots, n\} \mid x_i \neq 0\}.$$

Die Indexmenge  $\text{supp}(x)$  heißt **Träger** von  $x$ .

**(8.10) Satz.** Für  $x \in P^=(A, b) \subseteq \mathbb{K}^n$  sind folgende Aussagen äquivalent:

- (1)  $x$  ist Ecke von  $P^=(A, b)$ .  
 (2)  $\text{rang}(A_{\text{supp}(x)}) = |\text{supp}(x)|$ .  
 (3) Die Spaltenvektoren  $A_{\cdot, j}$ ,  $j \in \text{supp}(x)$ , sind linear unabhängig.

**Beweis :** (2)  $\iff$  (3) trivial.

(1)  $\iff$  (2). Sei  $D = \begin{pmatrix} A \\ -A \\ -I \end{pmatrix}$ ,  $d = \begin{pmatrix} b \\ -b \\ 0 \end{pmatrix}$ , dann gilt  $P(D, d) = P^=(A, b)$ . Mit

Satz (8.9) gilt nun:

$$\begin{aligned}
 x \text{ Ecke von } P^=(A, b) &\iff x \text{ Ecke von } P(D, d) \\
 &\iff \text{rang}(D_{\text{eq}(\{x\})}) = n \\
 &\iff \text{rang} \begin{pmatrix} A \\ -A \\ I_J \end{pmatrix} = n \text{ mit } J = \{j \mid x_j = 0\} \\
 &\quad = \{1, \dots, n\} \setminus \text{supp}(x) \\
 (\text{Streichen der Spalten } j \in J) &\iff \text{rang} \begin{pmatrix} A_{\cdot, \text{supp}(x)} \\ -A_{\cdot, \text{supp}(x)} \end{pmatrix} = n - |J| = |\text{supp}(x)| \\
 &\iff \text{rang}(A_{\cdot, \text{supp}(x)}) = |\text{supp}(x)|.
 \end{aligned}$$

□

### 8.3 Spitze Polyeder

Nicht alle Polyeder haben Ecken, so z. B. der Durchschnitt von weniger als  $n$  Halbräumen des  $\mathbb{K}^n$ . Polyeder mit Ecken sind besonders für die lineare Optimierung wichtig.

**(8.11) Definition.** Ein Polyeder heißt **spitz**, wenn es eine Ecke besitzt.

□

Wir wollen im nachfolgenden spitze Polyeder charakterisieren und einige Aussagen über spitze Polyeder beweisen.

**(8.12) Satz.** Sei  $P = P(A, b) \subseteq \mathbb{K}^n$  ein nichtleeres Polyeder, dann sind äquivalent:

- (1)  $P$  ist spitz.
- (2)  $\text{rang}(A) = n$ .
- (3)  $\text{rec}(P)$  ist spitz, d. h.  $0$  ist eine Ecke von  $\text{rec}(P)$ .
- (4) Jede nichtleere Seitenfläche von  $P$  ist spitz.
- (5)  $\text{hog}(P)$  ist spitz.
- (6)  $P$  enthält keine Gerade.
- (7)  $\text{rec}(P)$  enthält keine Gerade.
- (8)  $\text{lineal}(P) = \{0\}$ .

**Beweis:** (1)  $\implies$  (2). Ist  $x$  eine Ecke von  $P$ , so gilt nach (8.9)  $n = \text{rang}(A_{\text{eq}(\{x\})}) \leq \text{rang}(A) \leq n$ , also  $\text{rang}(A) = n$ .

(2)  $\implies$  (1). Sei  $x \in P$  so gewählt, dass  $I := \text{eq}(\{x\})$  maximal bezüglich Mengeninklusion ist. Sei  $F = \{y \in P \mid A_I \cdot y = b_I\}$ , dann ist  $x$  innerer Punkt von  $F$ . Ist  $\text{rang}(A_I) < n$ , dann enthält der Kern von  $A_I$  einen Vektor  $y \neq 0$ , und mit Lemma (8.8) gibt es ein  $\varepsilon > 0$ , so dass  $x \pm \varepsilon y \in P$ . Die Gerade  $G = \{x + \lambda y \mid \lambda \in \mathbb{K}\}$  trifft mindestens eine der Hyperebenen  $H_j = \{y \mid A_j \cdot y = b_j\}$ ,  $j \notin I$ , da  $\text{rang}(A) > \text{rang}(A_I)$ . Also muss es ein  $\delta \in \mathbb{K}$  geben, so dass  $x + \delta y \in P$  und  $\text{eq}(\{x + \delta y\}) \supset I$ . Dies widerspricht der Maximalität von  $I$ .

Aus der Äquivalenz von (1) und (2) folgt direkt die Äquivalenz der Aussagen (3), (4), (5), da

$$\begin{aligned} \text{rec}(P) &= P(A, 0), \text{ nach (8.3),} \\ \text{hog}(P) &= P\left(\begin{pmatrix} A & -b \\ 0 & -1 \end{pmatrix}, 0\right), \text{ nach (8.5),} \\ F &= \{x \in \mathbb{K}^n \mid Ax \leq b, A_{\text{eq}(F)} \cdot x \leq b_{\text{eq}(F)}, -A_{\text{eq}(F)} \cdot x \leq -b_{\text{eq}(F)}\} \end{aligned}$$

für alle Seitenflächen  $F$  von  $P$  gilt.

(3)  $\implies$  (6). Angenommen  $P$  enthält eine Gerade  $G = \{u + \lambda v \mid \lambda \in \mathbb{K}\}$ ,  $v \neq 0$ , dann gilt  $b \geq A(u + \lambda v) = Au + \lambda Av$  für alle  $\lambda \in \mathbb{K}$ . Daraus folgt  $A(\lambda v) \leq 0$  für alle  $\lambda \in \mathbb{K}$  und somit  $v, -v \in \text{rec}(P)$ , d. h.  $0 = \frac{1}{2}v + \frac{1}{2}(-v)$  ist eine Konvexkombination. Also ist  $0$  keine Ecke von  $\text{rec}(P)$ .

(6)  $\implies$  (3). Ist  $\text{rec}(P)$  nicht spitz, so ist  $0$  echte Konvexkombination von Vektoren aus  $\text{rec}(P)$ , sagen wir  $0 = \lambda u + (1 - \lambda)v$ ,  $u \neq 0 \neq v$ ,  $0 < \lambda < 1$ . Dann aber

ist neben  $u$  auch  $-\lambda u = (1 - \lambda)v \in \text{rec}(P)$  und folglich ist  $G = \{\mu u \mid \mu \in \mathbb{K}\}$  eine Gerade in  $\text{rec}(P)$ , und für alle  $x \in P$  ist  $x + G$  eine Gerade in  $P$ .

Die Äquivalenz von (7) und (8) zu den übrigen Aussagen ist nun offensichtlich.  $\square$

**(8.13) Folgerung.** Sei  $P = P^=(A, b) \subseteq \mathbb{K}^n$ , dann gilt

$$P \neq \emptyset \iff P \text{ spitz.}$$

**Beweis:** Es ist  $P = P(D, d)$  mit  $D = \begin{pmatrix} A \\ -A \\ -I \end{pmatrix}$ ,  $d = \begin{pmatrix} b \\ -b \\ 0 \end{pmatrix}$ , und offenbar hat  $D$  den Rang  $n$ . Aus (8.12) folgt dann die Behauptung.  $\square$

**(8.14) Folgerung.** Sei  $P \subseteq \mathbb{K}^n$  ein Polytop, dann gilt

$$P \neq \emptyset \iff P \text{ spitz.}$$

**Beweis:** Da  $P$  beschränkt ist, gibt es einen Vektor  $u$  mit  $P \subseteq \{x \mid x \leq u\}$ . Gilt  $P = P(A, b)$ , so folgt daraus  $P = P(D, d)$  mit  $D = \begin{pmatrix} A \\ I \end{pmatrix}$ ,  $d = \begin{pmatrix} b \\ u \end{pmatrix}$ .  $D$  hat den Rang  $n$ , daraus folgt mit (8.12) die Behauptung.  $\square$

**(8.15) Folgerung.** Das Polyeder  $P$  sei spitz, und das lineare Programm  $\max c^T x$ ,  $x \in P$  habe eine Optimallösung, dann hat dieses lineare Programm auch eine optimale Lösung, die eine Ecke von  $P$  ist (eine sogenannte optimale **Ecklösung**).

**Beweis:**  $F = \{x \in P \mid c^T x = \max\{c^T y \mid y \in P\}\}$  ist eine nichtleere Seitenfläche von  $P$ , die nach Satz (8.12) eine Ecke hat.  $\square$

**(8.16) Folgerung.** Ist  $P$  ein nichtleeres Polytop, dann hat jedes lineare Programm  $\max c^T x$ ,  $x \in P$  eine optimale Ecklösung.  $\square$

Das folgende Korollar wird sich als sehr wichtig erweisen.

**(8.17) Folgerung.** Lineare Programme der Form

$$\begin{array}{l} \max c^T x \\ Ax = b \\ x \geq 0 \end{array} \quad \text{oder} \quad \begin{array}{l} \max c^T x \\ Ax \leq b \\ x \geq 0 \end{array}$$

besitzen genau dann eine endliche Optimallösung, wenn sie eine optimale Ecklösung besitzen.

**Beweis :** Nach (8.13) ist  $P^=(A, b)$  spitz (falls nicht leer) und somit hat das erste der obigen LP's nach (8.15) eine optimale Ecklösung, falls es überhaupt eine optimale Lösung hat. Analog folgt die Behauptung für das zweite LP.  $\square$

## 8.4 Extremalen von spitzen Polyedern

Wir wollen nachfolgend einige Aussagen über spitze Polyeder beweisen, die sich — entsprechend modifiziert — auch für allgemeine Polyeder zeigen lassen. Dabei treten jedoch einige unschöne technische Komplikationen auf, so dass wir hier auf die Behandlung dieser Verallgemeinerung verzichten.

**(8.18) Definition.** Sei  $P$  ein Polyeder. Ein Vektor  $z \in \text{rec}(P) \setminus \{0\}$  heißt **Extremale** (oder **Extremalvektor**) von  $P$ , wenn  $\text{cone}(\{z\})$  ein Extremalstrahl von  $\text{rec}(P)$  ist.

$\square$

Nur spitze Polyeder haben Extremalen. Denn ist  $P$  nicht spitz, so ist nach (8.12)  $\text{rec}(P)$  nicht spitz, also ist der Nullvektor eine echte Konvexkombination zweier von Null verschiedener Vektoren, sagen wir  $0 = \lambda u + (1 - \lambda)v$ ,  $0 < \lambda < 1$ . Ist  $F = \text{cone}(\{z\})$  ein Extremalstrahl von  $\text{rec}(P)$ , so gibt es eine bezüglich  $\text{rec}(P)$  gültige Ungleichung  $c^T x \leq 0$  mit  $F = \{x \in \text{rec}(P) \mid c^T x = 0\}$ . Nun gilt  $0 = c^T 0 = c^T(\lambda u + (1 - \lambda)v) = \lambda c^T u + (1 - \lambda)c^T v \leq 0$ . Aus  $c^T u \leq 0$ ,  $c^T v \leq 0$  folgt  $c^T u = c^T v = 0$  und somit  $u, v \in F$ , ein Widerspruch. Aussagen über Extremalen machen also nur für spitze Polyeder Sinn.

Ist  $K$  speziell ein spitzer polyedrischer Kegel, so ist (wegen  $\text{rec}(K) = K$ ) eine Extremale von  $K$  ein Vektor  $z \in K$ , so dass  $\text{cone}(\{z\})$  ein Extremalstrahl von  $K$  ist. Das heißt, jeder auf einem Extremalstrahl von  $K$  gelegener und von Null verschiedener Vektor ist eine Extremale von  $K$ .

**(8.19) Satz.** Seien  $P = P(A, b) \subseteq \mathbb{K}^n$  ein spitzes Polyeder und  $z \in \text{rec}(P) \setminus \{0\}$ . Dann sind äquivalent:

- (1)  $z$  ist eine Extremale von  $P$ .
- (2)  $\text{cone}(\{z\})$  ist ein Extremalstrahl von  $\text{rec}(P)$ .

- (3)  $z$  lässt sich nicht als echte konische Kombination zweier linear unabhängiger Elemente von  $\text{rec}(P)$  darstellen.
- (4)  $(\text{rec}(P) \setminus (\text{cone}\{z\})) \cup \{0\}$  ist ein Kegel.
- (5)  $\text{rang}(A_{\text{eq}(\{z\})}) = n - 1$  (eq bezüglich  $Ax \leq 0$ ).

**Beweis :** (1)  $\iff$  (2) Definition.

(2)  $\implies$  (3). Ist  $F := \text{cone}(\{z\})$  ein Extremalstrahl von  $\text{rec}(P)$ , so ist  $F$  eine eindimensionale Seitenfläche von  $\text{rec}(P)$ , d. h.  $F$  kann keine zwei linear unabhängigen Vektoren enthalten. Insbesondere gibt es eine bezüglich  $\text{rec}(P)$  gültige Ungleichung  $c^T x \leq 0$  mit  $F = \{x \in \text{rec}(P) \mid c^T x = 0\}$ . Gibt es zwei linear unabhängige Vektoren  $u, v \in \text{rec}(P)$  und  $\lambda, \mu > 0$  mit  $z = \lambda u + \mu v$ , so gilt  $0 = c^T z = c^T(\lambda u + \mu v) = \lambda c^T u + \mu c^T v \leq 0$ . Daraus folgt  $c^T u = c^T v = 0$ , d. h.  $u, v \in F$ , ein Widerspruch.

(3)  $\iff$  (4) trivial.

(3)  $\implies$  (5). Sei  $I = \text{eq}(\{z\})$ , dann ist  $z$  innerer Punkt von  $F := \{x \in \text{rec}(P) \mid A_I x = 0\}$ . Ist  $\text{rang}(A_I) < n - 1$ , dann enthält der Kern von  $A_I$  einen von  $z$  linear unabhängigen Vektor  $u$ . Nach Lemma (8.8) gibt es ein  $\varepsilon > 0$ , so dass  $z \pm \varepsilon u \in \text{rec}(P)$  gilt. Dann aber gilt  $z = \frac{1}{2}(z + \varepsilon u) + \frac{1}{2}(z - \varepsilon u)$ , d. h.  $z$  ist echte konische Kombination von zwei linear unabhängigen Elementen von  $\text{rec}(P)$ , Widerspruch. Offenbar ist  $\text{rang}(A_I) \neq n$ .

(5)  $\implies$  (2). Sei  $I = \text{eq}(\{z\})$ , dann folgt für  $F = \{x \in \text{rec}(P) \mid A_I x = 0\}$  aus der Voraussetzung, dass  $\dim(F) = 1$  gilt. Da  $\text{rec}(P)$  spitz ist, enthält nach (8.12)  $\text{rec}(P)$  keine Gerade, also muss die eindimensionale Seitenfläche  $F$  der Strahl  $\text{cone}(\{z\})$  sein.  $\square$

Wir wollen nun noch eine Beziehung zwischen den Ecken und Extremalen eines spitzen Polyeders und den Extremalen seiner Homogenisierung aufzeigen.

**(8.20) Satz.** Sei  $P \subseteq \mathbb{K}^n$  ein spitzes Polyeder, dann gilt:

- (a)  $x$  ist Ecke von  $P \iff \begin{pmatrix} x \\ 1 \end{pmatrix}$  ist Extremale von  $\text{hog}(P)$ .
- (b)  $z$  ist Extremale von  $P \iff \begin{pmatrix} z \\ 0 \end{pmatrix}$  ist Extremale von  $\text{hog}(P)$ .

**Beweis :** Sei  $P = P(A, b)$ , dann gilt nach Satz (8.5)  $\text{hog}(P) = P(B, 0)$  mit

$$B = \begin{pmatrix} A, & -b \\ 0, & -1 \end{pmatrix}.$$

(a) Sei  $I = \text{eq}(\{x\})$  bezüglich  $P(A, b)$ .  $x$  ist Ecke von  $P \stackrel{(8.9)}{\iff} \text{rang}(A_I) = n \iff \text{rang}(B_{\text{eq}(\{\begin{pmatrix} x \\ 1 \end{pmatrix}\})}) = n$  (denn die neu hinzugekommene Ungleichung ist nicht mit Gleichheit erfüllt)  $\stackrel{(8.19)}{\iff} \begin{pmatrix} x \\ 1 \end{pmatrix}$  ist Extremale von  $\text{hog}(P)$ .

(b)  $z$  ist Extremale von  $P \iff z$  ist Extremale von  $\text{rec}(P) \stackrel{(8.19)}{\iff} \text{rang}(A_{\text{eq}(\{z\})}) = n - 1 \iff \text{rang}(B_{\text{eq}(\{\begin{pmatrix} z \\ 0 \end{pmatrix}\})}) = n \stackrel{(8.19)}{\iff} \begin{pmatrix} z \\ 0 \end{pmatrix}$  ist Extremale von  $\text{hog}(P)$ .  $\square$

## 8.5 Einige Darstellungssätze

Wir knüpfen hier an Abschnitt 6.3 an, wo wir bereits verschiedene Darstellungssätze bewiesen haben. Einige dieser Sätze können wir nun mit Hilfe der in diesem Kapitel gewonnenen Erkenntnisse verschärfen.

Ist  $K$  ein polyedrischer Kegel und gilt  $K = \text{cone}(E)$ , dann nennen wir  $E$  eine **Kegelbasis** von  $K$ , wenn es keine echte Teilmenge  $E'$  von  $E$  gibt mit  $K = \text{cone}(E')$  und wenn jede andere minimale Menge  $F$  mit  $K = \text{cone}(F)$  dieselbe Kardinalität wie  $E$  hat. Ist  $P$  ein Polytop, dann heißt eine Menge  $V$  mit  $P = \text{conv}(V)$  **konvexe Basis** von  $P$ , wenn  $V$  keine echte Teilmenge besitzt, deren konvexe Hülle  $P$  ist, und wenn jede andere minimale Menge  $W$  mit  $P = \text{conv}(W)$  dieselbe Kardinalität wie  $V$  hat.

Trivialerweise sind die Elemente einer Kegelbasis  $E$  **konisch unabhängig**, d. h. kein  $e \in E$  ist konische Kombination der übrigen Elemente von  $E$ ; und die Elemente einer konvexen Basis sind **konvex unabhängig**, d. h. kein Element von  $V$  ist Konvexkombination der übrigen Elemente von  $V$ . Es gilt aber keineswegs, dass jeder Vektor  $x \in \text{cone}(E)$  bzw.  $x \in \text{conv}(V)$  eine eindeutige konische bzw. konvexe Darstellung durch Vektoren aus  $E$  bzw.  $V$  hat. In dieser Hinsicht unterscheiden sich also Kegelbasen und konvexe Basen von Vektorraumbasen.

**(8.21) Satz.** Sei  $\{0\} \neq K \subseteq \mathbb{K}^n$  ein spitzer polyedrischer Kegel, dann sind äquivalent:

- (1)  $E$  ist eine Kegelbasis von  $K$ .

- (2)  $E$  ist eine Menge, die man dadurch erhält, dass man aus jedem Extremalstrahl von  $K$  genau einen von Null verschiedenen Vektor (also eine Extremale von  $K$ ) auswählt.

**Beweis:** Ist  $z$  Extremale von  $K$ , so ist  $K' := (K \setminus \text{cone}\{z\}) \cup \{0\}$  nach (8.19) ein Kegel. Folglich gilt  $\text{cone}(E) \subseteq K'$  für alle Teilmengen  $E$  von  $K'$ . Also muss jede Kegelbasis von  $K$  mindestens ein (aus der Basiseigenschaft folgt sofort "genau ein") Element von  $\text{cone}\{z\} \setminus \{0\}$  enthalten.

Zum Beweis, dass jede wie in (2) spezifizierte Menge eine Kegelbasis ist, benutzen wir Induktion über  $d = \dim K$ . Für Kegel der Dimension 1 ist die Behauptung trivial. Sei die Behauptung für Kegel der Dimension  $d$  richtig und  $K$  ein Kegel mit  $\dim K = d + 1$ . Sei  $y \in K \setminus \{0\}$  beliebig und  $c \in \mathbb{K}^n \setminus \{0\}$  ein Vektor, so dass die Ecke 0 von  $K$  die eindeutig bestimmte Lösung von  $\max\{c^T x \mid x \in K\}$  ist ( $c$  existiert nach (8.9)). Sei  $z \in \{x \mid c^T x = 0\} \setminus \{0\}$ . Dann ist für die Gerade

$$G = \{y + \lambda z \mid \lambda \in \mathbb{K}\}$$

die Menge  $K \cap G$  ein endliches Streckenstück (andernfalls wäre  $z \in \text{rec}(K) = K$ , und wegen  $c^T z = 0$  wäre 0 nicht der eindeutige Maximalpunkt). Folglich gibt es zwei Punkte  $z_1$  und  $z_2$ , die auf echten Seitenflächen, sagen wir  $F_1$  und  $F_2$ , von  $K$  liegen, so dass  $K \cap G = \text{conv}(\{z_0, z_1\})$ . Die Dimensionen der Seitenflächen  $F_1, F_2$  sind höchstens  $d$ ,  $F_1$  und  $F_2$  sind Kegel, und die Extremalstrahlen von  $F_1$  und  $F_2$  sind Extremalstrahlen von  $K$ . Nach Induktionsvoraussetzung werden  $z_1$  und  $z_2$  durch die in (2) festgelegten Mengen bezüglich  $F_1$  und  $F_2$  konisch erzeugt. Daraus folgt, dass  $y$  durch jede Menge des Typs (2) konisch erzeugt werden kann. Dies impliziert die Behauptung.  $\square$

**(8.22) Folgerung.** Jeder spitze polyedrische Kegel besitzt eine — bis auf positive Skalierung der einzelnen Elemente — eindeutige Kegelbasis.  $\square$

**(8.23) Folgerung.** Jeder spitze polyedrische Kegel  $K \neq \{0\}$  ist die Summe seiner Extremalstrahlen, d. h. sind  $\text{cone}(\{e_i\})$ ,  $i = 1, \dots, k$  die Extremalstrahlen von  $K$ , so gilt

$$K = \text{cone}(\{e_1, \dots, e_k\}) = \sum_{i=1}^k \text{cone}(\{e_i\}).$$

**Beweis :** Klar. □

Der folgende Satz verschärft (6.13) für spitze Polyeder.

**(8.24) Satz.** *Jedes spitze Polyeder  $P$  lässt sich darstellen als die Summe der konvexen Hülle seiner Ecken und der konischen Hülle seiner Extremalen, d. h. sind  $V$  die Eckenmenge von  $P$  und  $\text{cone}(\{e\})$ ,  $e \in E$ , die Extremalstrahlen von  $(\text{rec}(P))$  (bzw. ist  $E$  eine Kegelbasis von  $\text{rec}(P)$ ), so gilt*

$$P = \text{conv}(V) + \text{cone}(E).$$

**Beweis :** Sei  $\text{hog}(P)$  die Homogenisierung von  $P$ . Da  $P$  spitz ist, ist  $\text{hog}(P)$  nach (8.12) ein spitzer Kegel. Nach (8.23) ist  $\text{hog}(P)$  die Summe seiner Extremalstrahlen  $\text{cone}(\{e'_i\})$   $i = 1, \dots, k$ . O. b. d. A. können wir annehmen, dass  $e'_i = \begin{pmatrix} v_i \\ 1 \end{pmatrix}$ ,  $i = 1, \dots, p$ , und  $e'_i = \begin{pmatrix} e_i \\ 0 \end{pmatrix}$ ,  $i = p+1, \dots, k$  gilt. Aus (8.20) folgt:  $V = \{v_1, \dots, v_p\}$  ist die Eckenmenge von  $P$  und  $E = \{e_{p+1}, \dots, e_k\}$  ist die Extremalenmenge von  $P$ . Nach (8.6) gilt  $x \in P \iff \begin{pmatrix} x \\ 1 \end{pmatrix} \in \text{hog}(P)$  und somit folgt  $x \in P \iff x = \sum_{i=1}^p \lambda_i v_i + \sum_{i=p+1}^k \mu_i e_i$ ,  $\lambda_i, \mu_i \geq 0$ ,  $\sum_{i=1}^p \lambda_i = 1$ , d. h.  $x \in \text{conv}(V) + \text{cone}(E)$ . □

**(8.25) Folgerung** (Satz von Krein-Milman). *Sei  $P$  ein Polyeder und  $V$  die Menge seiner Ecken, dann gilt*

$$P \text{ ist ein Polytop} \iff P = \text{conv}(V).$$

□

**(8.26) Folgerung.** *Polytope haben eine eindeutige konvexe Basis.*

□

Für die lineare Optimierung ist die folgende Beobachtung wichtig.

**(8.27) Satz.** *Seien  $P \subseteq \mathbb{K}^n$  ein spitzes Polyeder und  $c \in \mathbb{K}^n$ . Das lineare Programm  $\max c^T x$ ,  $x \in P$  ist genau dann unbeschränkt, wenn es eine Extremale  $e$  von  $P$  gibt mit  $c^T e > 0$ .*

**Beweis :** Seien  $V$  die Eckenmenge von  $P$  und  $E$  eine Kegelbasis von  $\text{rec}(P)$ , dann gilt nach (8.24)  $P = \text{conv}(V) + \text{cone}(E)$ . Es ist  $\gamma := \max\{c^T v \mid v \in V\} < \infty$ , da  $V$  endlich und  $\gamma = \max\{c^T x \mid x \in \text{conv}(V)\}$ . Gilt  $c^T e \leq 0, \forall e \in E$ , so

ist  $\gamma = \max\{c^T x \mid x \in P\} < \infty$ . Falls also  $\max\{c^T x \mid x \in P\}$  unbeschränkt ist, muss für mindestens ein  $e \in E$  gelten  $c^T e > 0$ . Die umgekehrte Richtung ist trivial.  $\square$

Wie bereits bemerkt, haben Elemente von spitzen Polyedern  $P$  i. a. keine eindeutige Darstellung als konvexe und konische Kombination von Ecken und Extremalen. Man kann jedoch zeigen, dass zu einer derartigen Darstellung von Elementen von  $P$  nicht allzu viele Ecken und Extremalen benötigt werden.

**(8.28) Satz.** *Es seien  $K \subseteq \mathbb{K}^n$  ein spitzer Kegel und  $0 \neq x \in K$ . Dann gibt es Extremalen  $y_1, \dots, y_d$  von  $K$ , wobei  $d \leq \dim(K) \leq n$  gilt, mit*

$$x = \sum_{i=1}^d y_i.$$

**Beweis :** Es seien  $\text{cone}(\{e_i\})$  die Extremalstrahlen von  $K$ ,  $i = 1, \dots, k$ , dann gibt es nach (8.23) Skalare  $\lambda_i \geq 0$ ,  $i = 1, \dots, k$ , mit

$$x = \sum_{i=1}^k \lambda_i e_i.$$

Unter allen möglichen Darstellungen von  $x$  dieser Art sei die obige eine, so dass  $I = \{i \in \{1, \dots, k\} \mid \lambda_i > 0\}$  minimal ist. Sagen wir, es gilt  $I = \{1, \dots, d\}$ . Angenommen die Vektoren  $e_i$ ,  $i \in I$  sind linear abhängig, dann gibt es  $\mu_1, \dots, \mu_d \in \mathbb{K}$ , nicht alle  $\mu_i$  Null, so dass gilt

$$\sum_{i=1}^d \mu_i e_i = 0.$$

Angenommen  $\mu_i \geq 0$  für  $i = 1, \dots, d$ , und o. B. d. A. sei  $\mu_1 > 0$ . Dann ist  $-e_1 = \sum_{i=2}^d \frac{\mu_i}{\mu_1} e_i$  eine konische Kombination, und nach Satz (8.4) gilt  $e_1 \in \text{lineal}(K)$ . Dies widerspricht nach (8.12) der Voraussetzung  $K$  ist spitz.

O. B. d. A. können wir daher annehmen, dass gilt  $\mu_1 < 0$  und

$$\frac{\lambda_1}{\mu_1} = \max\left\{\frac{\lambda_i}{\mu_i} \mid \mu_i < 0\right\}.$$

Daraus folgt

$$\begin{aligned} e_1 &= -\sum_{i=2}^d \frac{\mu_i}{\mu_1} e_i, \\ x &= \sum_{i=2}^d \left(\lambda_i - \frac{\lambda_1}{\mu_1} \mu_i\right) e_i. \end{aligned}$$

Diese Darstellung von  $x$  ist eine konische Kombination, denn

$$\begin{aligned}\mu_i \geq 0 &\implies (\lambda_i - \frac{\lambda_1}{\mu_1} \mu_i) \geq 0, \\ \mu_i < 0 &\implies \frac{\lambda_i}{\mu_i} \leq \frac{\lambda_1}{\mu_1} \implies \lambda_i \geq \frac{\lambda_1}{\mu_1} \mu_i,\end{aligned}$$

also kann  $x$  mit weniger als  $d$  Extremalen konisch dargestellt werden, Widerspruch zur Minimalität! Da ein Kegel höchstens  $\dim(K)$  linear unabhängige Vektoren enthält, folgt  $d \leq \dim(K)$ . Setzen wir  $y_i = \lambda_i e_i$ ,  $i = 1, \dots, d$ , so sind die Vektoren  $y_i$  Extremalen von  $K$  mit der gewünschten Eigenschaft.  $\square$

**(8.29) Folgerung.** *Ist  $P \subseteq \mathbb{K}^n$  ein spitzes Polyeder und ist  $x \in P$ , dann gibt es Ecken  $v_0, v_1, \dots, v_k$  und Extremalen  $e_{k+1}, e_{k+2}, \dots, e_d$  von  $P$  mit  $d \leq \dim(P)$  und nichtnegative Skalare  $\lambda_0, \dots, \lambda_k$  mit  $\sum_{i=0}^k \lambda_i = 1$ , so dass gilt*

$$x = \sum_{i=0}^k \lambda_i v_i + \sum_{i=k+1}^d e_i.$$

**Beweis:** Nach (8.12) ist  $\text{hog}(P)$  ein spitzer Kegel, und die Dimension von  $\text{hog}(P)$  ist  $\dim(P) + 1$ . Nach (8.6) gilt  $x \in P \iff \begin{pmatrix} x \\ 1 \end{pmatrix} \in \text{hog}(P)$ . Nach Satz (8.28) ist  $\begin{pmatrix} x \\ 1 \end{pmatrix}$  konische Kombination von  $d + 1 \leq \dim(P) + 1$  Extremalen von  $\text{hog}(P)$ . O. B. d. A. können wir annehmen, dass gilt

$$\begin{pmatrix} x \\ 1 \end{pmatrix} = \sum_{i=0}^k \begin{pmatrix} y_i \\ \lambda_i \end{pmatrix} + \sum_{i=k+1}^d \begin{pmatrix} e_i \\ 0 \end{pmatrix},$$

wobei  $\lambda_i > 0$ ,  $i = 0, \dots, k$ . Für  $v_i := \frac{1}{\lambda_i} y_i$  gilt dann

$$x = \sum_{i=0}^k \lambda_i v_i + \sum_{i=k+1}^d e_i, \quad \sum_{i=0}^k \lambda_i = 1.$$

Ferner sind nach (8.20) die Vektoren  $v_i$  Ecken von  $P$  und die Vektoren  $e_i$  Extremalen von  $P$ . Also haben wir die gewünschte Kombination von  $x$  gefunden.  $\square$

Das folgende direkte Korollar von (8.29) wird in der Literatur häufig **Satz von Caratheodory** genannt.

**(8.30) Folgerung.** *Ist  $P \subseteq \mathbb{K}^n$  ein Polytop, dann ist jedes Element von  $P$  Konvexkombination von höchstens  $\dim(P) + 1$  Ecken von  $P$ .*

$\square$



# Kapitel 9

## Die Grundversion der Simplex-Methode

In Satz (7.8) haben wir gezeigt, dass die Menge der Optimallösungen eines linearen Programms eine Seitenfläche des durch die Restriktionen definierten Polyeders ist. Dieses Ergebnis haben wir in Folgerung (8.15) dahingehend verschärft, dass wir nachweisen konnten, dass bei spitzen Polyedern das Optimum, falls es existiert, stets in einer Ecke angenommen wird. Wollen wir also ein Verfahren entwerfen, das lineare Programme über spitzen Polyedern löst, so genügt es, ein Verfahren anzugeben, das eine optimale Ecke derartiger Polyeder findet.

Aus Kapitel 2 wissen wir, dass wir uns auf lineare Programme über Polyedern der Form  $P^=(A, b) = \{x \in \mathbb{R}^n \mid Ax = b, x \geq 0\}$  beschränken können. In Folgerung (8.13) haben wir gezeigt, dass ein solches Polyeder stets spitz ist, falls es nicht leer ist. Wenn wir also einen Algorithmus angeben können, der die optimale Ecke eines linearen Programms der Form  $\max c^T x, x \in P^=(A, b)$  findet, so können wir mit diesem alle linearen Programmierungsprobleme lösen.

**(9.1) Definition.** *Ein lineares Programm der Form*

$$\begin{aligned} \max c^T x \\ Ax = b \\ x \geq 0 \end{aligned}$$

*heißt ein lineares Programm in **Standardform**. Wenn nichts anderes gesagt wird, nehmen wir immer an, dass  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$ , und  $c \in \mathbb{K}^n$  gilt.*

□

Ist der Rang von  $A$  gleich  $n$ , so wissen wir, dass das Gleichungssystem  $Ax = b$  entweder gar keine oder eine eindeutig bestimmte Lösung hat, die mit dem Gauß'schen Eliminationsverfahren berechnet werden kann. Da also entweder keine Lösung existiert oder der einzige zulässige Punkt leicht bestimmt werden kann, liegt kein wirkliches Optimierungsproblem vor. Wir wollen daher voraussetzen, dass  $\text{rang}(A) < n$  gilt. Ferner nehmen wir an, um die Darstellung einfacher zu machen, dass die Zeilen von  $A$  linear unabhängig sind und dass  $P^=(A, b) \neq \emptyset$  ist. Wir treffen also folgende Verabredung:

**(9.2) Generalvoraussetzungen für Kapitel 9.**

- (a)  $A$  ist eine  $(m, n)$ -Matrix mit  $m < n$ ,
- (b)  $\text{rang}(A) = m$ ,
- (c)  $P^=(A, b) \neq \emptyset$ .

□

Wir werden später zeigen, dass lineare Programme, bei denen die Voraussetzungen (9.2) nicht gegeben sind, auf solche zurückgeführt werden können, die diese erfüllen, dass also alle Voraussetzungen (9.2) o. B. d. A. gemacht werden können.

Um die weitere Darstellung schreibtechnisch zu vereinfachen, legen wir für dieses Kapitel folgendes fest (vergleiche Abschnitt 1.3.):

**(9.3) Konventionen in Kapitel 9.** *Es sei  $A \in \mathbb{K}^{(m,n)}$ , wobei nach (9.2)  $m < n$  gelten soll.*

- (a) *Mit  $\{1, \dots, m\}$  bezeichnen wir die Zeilenindexmenge von  $A$ , mit  $\{1, \dots, n\}$  die Spaltenindexmenge.*
- (b)  *$B$  und  $N$  bezeichnen stets Spaltenindexvektoren von  $A$ , wobei  $B = (p_1, \dots, p_m) \in \{1, \dots, n\}^m$  und  $N = (q_1, \dots, q_{n-m}) \in \{1, \dots, n\}^{n-m}$  gilt. Ferner kommt kein Spaltenindex, der in  $B$  vorkommt, auch in  $N$  vor, und umgekehrt. (Um Transpositionszeichen zu sparen, schreiben wir  $B$  und  $N$  immer als Zeilenvektoren.)*
- (c) *Wir werden aber  $B$  und  $N$  auch einfach als Mengen auffassen (wenn die Anordnung der Indizes keine Rolle spielt), dann gilt also nach (b)  $B \cap N = \emptyset$ ,  $B \cup N = \{1, \dots, n\}$ . Insbesondere können wir dann  $i \in B$  und  $j \in N$  schreiben.*
- (d) *Zur schreibtechnischen Vereinfachung schreiben wir im folgenden  $A_B$  und  $A_N$  statt  $A_{.B}$  und  $A_{.N}$ .*

- (e) Ist ein Spaltenindexvektor  $B$  wie oben angegeben definiert, so bezeichnet, falls nicht anders spezifiziert,  $N = (q_1, \dots, q_{n-m})$  immer den Spaltenindexvektor, der alle übrigen Spaltenindizes enthält, wobei  $q_i < q_j$  für  $i < j$  gilt.  $\square$

## 9.1 Basen, Basislösungen, Entartung

**(9.4) Definition.** Gegeben sei ein Gleichungssystem  $Ax = b$  mit  $A \in \mathbb{K}^{(m,n)}$ ,  $b \in \mathbb{K}^m$  und  $\text{rang}(A) = m$ . Spaltenindexvektoren  $B$  und  $N$  seien entsprechend Konvention (9.3) gegeben.

- (a) Ist  $A_B$  regulär, so heißt  $A_B$  **Basismatrix** oder kurz **Basis** von  $A$ , und  $A_N$  heißt **Nichtbasismatrix** oder **Nichtbasis** von  $A$ . Der Vektor  $x \in \mathbb{K}^n$  mit  $x_N = 0$  und  $x_B = A_B^{-1}b$  heißt **Basislösung von  $Ax = b$  zur Basis  $A_B$**  oder die zu  $A_B$  **gehörige Basislösung** oder kurz **Basislösung**.
- (b) Ist  $A_B$  eine Basis, dann heißen die Variablen  $x_j, j \in B$ , **Basisvariable** und die Variablen  $x_j, j \in N$ , **Nichtbasisvariable**.
- (c) Ist  $A_B$  eine Basis, dann heißen  $A_B$  und die zugehörige Basislösung  $x$  **zulässig** (bezüglich  $P^=(A, b)$ ), wenn  $A_B^{-1}b \geq 0$  gilt.
- (d) Eine Basislösung  $x$  zur Basis  $A_B$  heißt **nichtentartet** (oder **nichtdegeneriert**), falls  $x_B = A_B^{-1}b > 0$ , andernfalls heißt sie **entartet** (oder **degeneriert**).

$\square$

Die oben eingeführten Begriffe gehören zur Standardterminologie der linearen Programmierung. Eine erste begriffliche Brücke zur Polyedertheorie schlägt der nächste Satz.

**(9.5) Satz.** Seien  $P = P^=(A, b) \subseteq \mathbb{K}^n$  ein Polyeder mit  $\text{rang}(A) = m < n$  und  $x \in P$ . Dann sind äquivalent

- (1)  $x$  ist eine Ecke von  $P^=(A, b)$ .
- (2)  $x$  ist eine zulässige Basislösung (d. h. es gibt eine Basis  $A_B$  von  $A$  mit der Eigenschaft  $x_B = A_B^{-1}b \geq 0$  und  $x_N = 0$ ).

**Beweis :** (1)  $\implies$  (2). Sei  $I := \text{supp}(x)$ . Ist  $x$  eine Ecke von  $P^=(A, b)$ , so sind die Spaltenvektoren  $A_j, j \in I$ , nach (8.10) linear unabhängig. Wegen  $\text{rang}(A) = m$  gibt es eine Menge  $J \in \{1, \dots, n\} \setminus I, |J| = m - |I|$ , so dass die Spalten  $A_j, j \in I \cup J$ , linear unabhängig sind. Folglich ist  $A_B$  mit  $B = I \cup J$  eine Basis von  $A$ . Nehmen wir o. B. d. A. an, dass  $B$  aus den ersten  $m$  Spalten von  $A$  besteht, also  $A = (A_B, A_N)$  mit  $N = \{m+1, \dots, n\}$  gilt, dann erhalten wir aus  $x^T = (x_B^T, x_N^T) = (x_B^T, 0)$  folgendes:  $A_B^{-1}b = A_B^{-1}(Ax) = A_B^{-1}(A_B, A_N) \begin{pmatrix} x_B \\ 0 \end{pmatrix} = (I, A_B^{-1}A_N) \begin{pmatrix} x_B \\ 0 \end{pmatrix} = x_B$ . Da  $x \geq 0$  gilt, ist  $x$  eine Basislösung.

(2)  $\implies$  (1) folgt direkt aus (8.10).  $\square$

Damit sehen wir, dass Ecken von  $P^=(A, b)$  zulässigen Basislösungen von  $Ax = b, x \geq 0$  entsprechen. Also gibt es zu jeder Ecke eine zulässige Basis von  $A$ . Sind zwei Ecken verschieden, so sind natürlich die zugehörigen Basen verschieden. Dies gilt aber nicht umgekehrt!

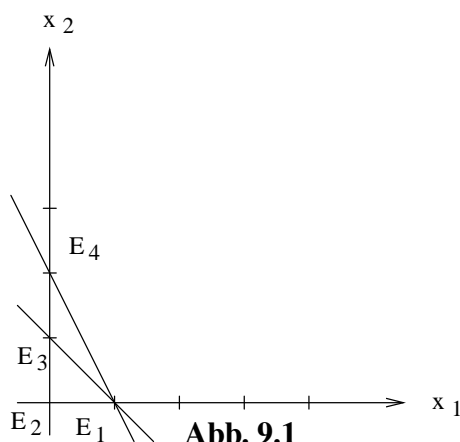
Ecke  $\overset{\text{eindeutig}}{\longleftrightarrow}$  zulässige Basislösung  $\overset{\text{nichteindeutig}}{\longleftrightarrow}$  zulässige Basis.

Wir wollen nun Entartung und ihre Ursachen untersuchen und betrachten zunächst drei Beispiele.

### (9.6) Beispiel (Degeneration).

- (a) Im  $\mathbb{K}^2$  gibt es keine "richtigen" degenerierten Probleme, da man solche Probleme durch Entfernung von redundanten Ungleichungen nichtdegeneriert machen kann. Im folgenden Beispiel ist die Ungleichung  $2x_1 + x_2 \leq 2$  redundant.

$$\begin{aligned} x_1 + x_2 &\leq 1 \\ 2x_1 + x_2 &\leq 2 \\ x_1, x_2 &\geq 0 \end{aligned}$$



Das durch das Ungleichungssystem definierte Polyeder (siehe Abbildung 9.1) hat die drei Ecken  $E_1, E_2, E_3$ . Wir formen durch Einführung von Schlupfvariablen um:

$$\begin{aligned} x_1 + x_2 + s_1 &= 1 \\ 2x_1 + x_2 + s_2 &= 2 \\ x_1, x_2, s_1, s_2 &\geq 0. \end{aligned}$$

Also ist  $P=(A, b)$  gegeben durch:

$$A = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 2 & 1 & 0 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 2 \end{pmatrix}.$$

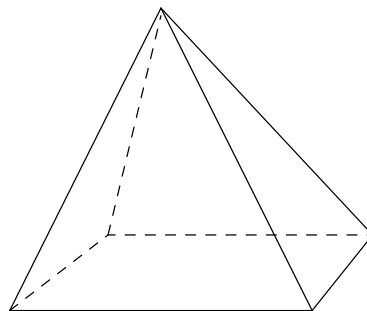
Die folgende Tabelle 9.1 enthält die Basen von  $A$  und die zugehörigen Basislösungen.

Spalten in $B$	$A_B$	$x_B$	$x_N$	$x$	Ecke
1,2	$\begin{pmatrix} 1 & 1 \\ 2 & 1 \end{pmatrix}$	$(x_1, x_2) = (1, 0)$	$(s_1, s_2)$	$(1, 0, 0, 0)$	$E_1$
1,3	$\begin{pmatrix} 1 & 1 \\ 2 & 0 \end{pmatrix}$	$(x_1, s_1) = (1, 0)$	$(x_2, s_2)$	$(1, 0, 0, 0)$	$E_1$
1,4	$\begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix}$	$(x_1, s_2) = (1, 0)$	$(x_2, s_1)$	$(1, 0, 0, 0)$	$E_1$
2,3	$\begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$	$(x_2, s_1) = (2, -1)$	$(x_1, s_2)$	$(0, 2, -1, 0)$	$E_4$
2,4	$\begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$	$(x_2, s_2) = (1, 1)$	$(x_1, s_1)$	$(0, 1, 0, 1)$	$E_3$
3,4	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$	$(s_1, s_2) = (1, 2)$	$(x_1, x_2)$	$(0, 0, 1, 2)$	$E_2$

**Tabelle 9.1**

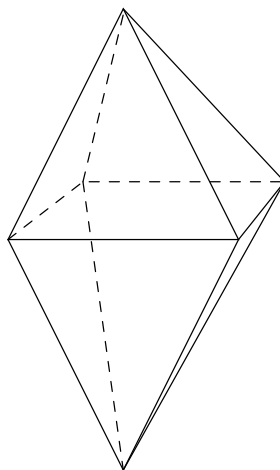
Im obigen Falle ist jede  $(2, 2)$ -Untermatrix von  $A$  eine Basis. Dies ist natürlich nicht immer so! Zur Ecke  $E_1$  gehören drei verschiedene Basen. Hier liegt also Entartung vor. Alle anderen Basislösungen sind nichtdegeneriert. Die zu  $B = (2, 3)$  gehörige Basis ist unzulässig. Zu unzulässigen Basislösungen sagt man manchmal auch unzulässige Ecke von  $P=(A, b)$ ,  $E_4$  ist also eine unzulässige Ecke.

- (b) Bei dreidimensionalen Polyedern im  $\mathbb{K}^3$  tritt echte Entartung auf. So ist z. B. die Spitze der folgenden Pyramide entartet, während alle übrigen vier Ecken nichtentartet sind.

**Abb. 9.2**

- (c) Dagegen sind alle Ecken der Doppelpyramide in Abbildung 9.3 entartet.

□

**Abb. 9.3**

Die Sätze (9.5) und (8.10) (und ihre Beweise) sagen uns, wie Ecken von  $P^=(A, b)$  mit den Basen von  $A_B$  zusammenhängen. Aus der linearen Algebra wissen wir zunächst, dass jeder Punkt  $x$  im  $\mathbb{K}^n$  als Durchschnitt von  $n$  Hyperebenen dargestellt werden kann, deren Normalenvektoren linear unabhängig sind. Algebraisch ausgedrückt, jeder Punkt  $x \in \mathbb{K}^n$  ist eindeutig bestimmte Lösung eines regulären  $(n, n)$ -Gleichungssystems  $Dy = d$  (mit  $\text{rang}(D) = n$ ).

Ecken von  $P^=(A, b)$  erhalten wir dadurch, dass wir nur Lösungen von regulären Gleichungssystemen

$$\begin{aligned} Ay &= b \\ e_j^T y &= 0, \quad j \in J \end{aligned}$$

mit  $J \subseteq \{1, \dots, n\}$  betrachten (und prüfen, ob die Lösungen in  $P^=(A, b)$  enthalten sind). Dabei kann es zu folgenden Situationen kommen:

Ist  $A_B$  eine Basis von  $A$  mit  $(A_B^{-1}b)_i \neq 0$  für alle  $i \in B$ , so gibt es nur eine einzige Möglichkeit das System  $Ay = b$  durch ein System  $e_j^T y = 0, j \in J$ , so zu ergänzen, dass das Gesamtsystem regulär wird. Man muss  $J = \{1, \dots, n\} \setminus B$  wählen, andernfalls sind nicht genügend Gleichungen vorhanden. Daraus folgt speziell:

**(9.7) Bemerkung.** Ist  $x$  eine nichtdegenerierte Basislösung von  $Ay = b, y \geq 0$ , dann gibt es eine eindeutig bestimmte zu  $x$  gehörige Basis  $A_B$  von  $A$  ( $x_B = A_B^{-1}b, x_N = 0$ ).

□

Die Umkehrung gilt i. a. nicht, wie das folgende Beispiel zeigt.

**(9.8) Beispiel.**  $P^=(A, b) \subseteq \mathbb{R}^3$  sei gegeben durch

$$A = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Die Matrix  $A$  besitzt zwei Basen  $A_{B_1}, A_{B_2}$  mit  $B_1 = (1, 3), B_2 = (2, 3)$ . Die zu  $A_{B_1}$  gehörige Basislösung ist  $x_1^T = (1, 0, 0)$ , die zu  $A_{B_2}$  gehörige Basislösung ist  $x_2^T = (0, 1, 0)$ . Beide Basislösungen sind entartet, aber die zugehörigen Basen sind eindeutig bestimmt. □

Ist  $A_B$  eine Basis von  $A$  und  $x$  die zugehörige Basislösung, so dass einige Komponenten von  $A_B^{-1}b$  Null sind (d. h.  $x$  ist entartet), so enthält das Gleichungssystem  $Ay = b, e_j^T y = 0$  für alle  $j \in J = \{j \mid x_j = 0\}$  mehr als  $n$  Gleichungen (also mehr als notwendig). I. a. ist dann  $x$  Lösung von mehreren regulären  $(n, n)$ -Untersystemen dieses Gleichungssystems. Man hat also keine eindeutige Zuordnung mehr von der Ecke  $x$  zu einem Gleichungssystem der Form  $Ay = b, e_j^T y = 0, j \in J$ . Wir werden sehen, dass Entartung zu rechentechnischen Problemen führt. Entartung kann drei Ursachen haben

- überflüssige Variable (im Beispiel (9.8) kann  $x_3$  weggelassen werden, dann sind alle Basislösungen nichtentartet),

- redundante Ungleichungen (Beispiel (9.6) (a)),
- geometrische Gründe (Beispiele (9.6) (b) und (c)). Die durch überflüssige Variablen oder redundante Ungleichungen hervorgerufene Entartung kann i. a. beseitigt werden durch Weglassen der überflüssigen Variablen bzw. der redundanten Ungleichungen. Es gibt jedoch Polyeder — wie z. B. die Doppelpyramide in (9.6) (c) — die genuin entartet sind. Der Grund liegt hier darin, dass die Ecken überbestimmt sind, d. h. es gibt Ecken  $x$  von  $P$ , in denen mehr als  $\dim(P)$  Facetten zusammentreffen. Dann bestimmen je  $\dim(P)$  der diese Facetten definierenden Ungleichungen eine Basis von  $A$ , die  $x$  als zugehörige Basislösung liefert.

## 9.2 Basisaustausch (Pivoting), Simplexkriterium

Da jede zulässige Basis der Matrix  $A$  eine Ecke von  $P=(A, b)$  definiert, kann man ein lineares Programm der Form (9.1) dadurch lösen, dass man alle Basen von  $A$  bestimmt — dies sind endlich viele —, den Zielfunktionswert der zugehörigen Basislösung errechnet und die beste Lösung auswählt. Da eine  $(m, n)$ -Matrix bis zu  $\binom{n}{m}$  zulässige Basen haben kann, ist dieses Verfahren aufgrund des gewaltigen Rechenaufwandes (in jedem Schritt Bestimmung einer inversen Matrix) so gut wie undurchführbar.

Man sollte also versuchen, z. B. ein Verfahren so zu entwerfen, dass man von einer zulässigen Basislösung ausgehend eine neue Basislösung bestimmt, die zulässig ist und einen besseren Zielfunktionswert hat. Ferner muss das Verfahren so angelegt sein, dass der Übergang von einer Basis zur nächsten nicht allzuviel Rechenaufwand erfordert.

Im Prinzip kann man damit nicht gewährleisten, dass nicht alle Basen enumeriert werden, aber aufgrund heuristischer Überlegungen scheint ein solcher Ansatz, wenn er realisiert werden kann, besser als das obige Enumerationsverfahren funktionieren zu können. Im weiteren werden wir zeigen, wie man einen Basisaustausch, der von einer zulässigen Basislösung zu einer besseren zulässigen Basislösung mit relativ wenig Rechenaufwand führt, durchführen kann.

**(9.9) Satz.** *Gegeben sei ein Gleichungssystem  $Ay = b$  mit  $\text{rang}(A) = m$ , und  $A_B$  sei eine Basis von  $A$ . Dann gilt*

$$x \text{ erfüllt } Ay = b \iff x \text{ erfüllt } x_B = A_B^{-1}b - A_B^{-1}A_Nx_N.$$

**Beweis :** O. B. d. A. sei  $B = (1, \dots, m)$ , d. h.  $A = (A_B, A_N)$ ,  $x = \begin{pmatrix} x_B \\ x_N \end{pmatrix}$ , dann gilt

$$\begin{aligned} Ax = b &\iff (A_B, A_N) \begin{pmatrix} x_B \\ x_N \end{pmatrix} = b &\iff A_B^{-1}(A_B, A_N) \begin{pmatrix} x_B \\ x_N \end{pmatrix} = A_B^{-1}b \\ & &\iff (I, A_B^{-1}A_N) \begin{pmatrix} x_B \\ x_N \end{pmatrix} = A_B^{-1}b \\ & &\iff x_B + A_B^{-1}A_N x_N = A_B^{-1}b \\ & &\iff x_B = A_B^{-1}b - A_B^{-1}A_N x_N \end{aligned}$$

□

Die obige Beziehung ist simpel, von rechentechnischer Bedeutung ist jedoch die Tatsache, dass wir die Basisvariablen  $x_j$ ,  $j \in B$ , in Abhängigkeit von den Nichtbasisvariablen darstellen können.

Wir wollen uns nun überlegen, wie wir den Übergang von einer Ecke zu einer anderen Ecke vollziehen können, d. h. wie wir aus einer Basis eine andere Basis ohne viel Aufwand konstruieren können. Um dies rechentechnisch exakt vorführen zu können, müssen wir die Indexmengen  $B$ ,  $N$  wie in (9.3) (b) angegeben als Vektoren auffassen.

**(9.10) Satz (Basisaustausch, Pivotoperation).** Gegeben sei ein Gleichungssystem  $Ay = b$  mit  $\text{rang}(A) = m$ .  $B = (p_1, \dots, p_m)$  und  $N = (q_1, \dots, q_{n-m})$  seien Spaltenindexvektoren von  $A$  wie in (9.3) (b) festgelegt, so dass  $A_B$  eine Basis von  $A$  ist. Wir setzen

$$\begin{aligned} \bar{A} &:= A_B^{-1}A_N = (\bar{a}_{rs})_{\substack{1 \leq r \leq m \\ 1 \leq s \leq n-m}}, \\ \bar{b} &:= A_B^{-1}b \in \mathbb{K}^m. \end{aligned}$$

Ist  $\bar{a}_{rs} \neq 0$ , so ist  $A_{B'}$  mit  $B' := (p_1, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m)$  eine Basis von  $A$ . Ferner gilt:

$$A_{B'}^{-1} = EA_B^{-1},$$

wobei  $E$  eine sogenannte  $(m, m)$ -**Elementarmatrix** ist, die wie folgt definiert ist:

$$E = \begin{pmatrix} 1 & & & & & & & & & \\ & \ddots & & & & & & & & \\ & & \eta_1 & & & & & & & \\ & & \vdots & & & & & & & \\ & & 1\eta_{r-1} & & & & & & & \\ & & \eta_r & & & & & & & \\ & & \eta_{r+1} & & & & & & & \\ & & \vdots & & \ddots & & & & & \\ & & \eta_m & & & & & & & 1 \end{pmatrix}.$$

Dabei ist die  $r$ -te Spalte von  $E$  der Vektor  $\eta = (\eta_1, \dots, \eta_m)^T$  (genannt **Eta-Spalte**) gegeben durch

$$\begin{aligned}\eta_r &:= \frac{1}{\bar{a}_{rs}} \\ \eta_i &:= \frac{\bar{a}_{is}}{\bar{a}_{rs}} \quad i \in \{1, \dots, m\} \setminus \{r\}.\end{aligned}$$

Das Element  $\bar{a}_{rs}$  heißt **Pivot-Element**.

Sind  $x, x'$  die zu  $A_B, A'_B$  gehörigen Basislösungen, so gilt:

$$\begin{aligned}x'_{p_i} &= \bar{b}_i - \frac{\bar{a}_{is}}{\bar{a}_{rs}} \bar{b}_r = E_i \cdot \bar{b} = E_i \cdot x_B, \quad 1 \leq i \leq m, \quad i \neq r \quad (\text{neue Basisvariable}) \\ x'_{q_s} &= \frac{1}{\bar{a}_{rs}} \bar{b}_r (= E_r \cdot \bar{b}_r = E_r \cdot x_B) \quad (\text{neue Basisvariable}) \\ x'_j &= 0 \quad \text{andernfalls} \quad (\text{neue Nichtbasisvariable}).\end{aligned}$$

**Beweis :** Es seien

$$d := A_{q_s}, \quad \bar{d} := A_B^{-1} d = \bar{A}_{.s}, \quad F := A_B^{-1} A_{B'}.$$

$A_{B'}$  erhält man aus  $A_B$  dadurch, dass die  $r$ -te Spalte  $A_{p_r}$  von  $A_B$  durch  $d$  ersetzt wird.  $F$  ist daher eine  $(m, m)$ -Elementarmatrix, deren  $r$ -te Spalte der Vektor  $\bar{d}$  ist. Offenbar gilt daher

$$F \cdot E = \begin{pmatrix} 1 & & & \bar{d}_1 & & & & & \\ & \ddots & & \vdots & & & & & \\ & & & 1 & \bar{d}_{r-1} & & & & \\ & & & & \bar{d}_r & & & & \\ & & & & \bar{d}_{r+1} & 1 & & & \\ & & & \vdots & & & \ddots & & \\ & & & \bar{d}_m & & & & & 1 \end{pmatrix} \begin{pmatrix} 1 & & & \eta_1 & & & & & \\ & \ddots & & \vdots & & & & & \\ & & & 1 \eta_{r-1} & & & & & \\ & & & \eta_r & & & & & \\ & & & \eta_{r+1} & 1 & & & & \\ & & & \vdots & & \ddots & & & \\ & & & \eta_m & & & & & 1 \end{pmatrix} = I,$$

daraus folgt  $E = F^{-1}$  ( $\bar{d}_r \neq 0$  war vorausgesetzt). Da  $F$  und  $A_B$  invertierbar sind, ergibt sich aus  $A_{B'} = A_B F$  sofort  $A_{B'}^{-1} = F^{-1} A_B^{-1} = E A_B^{-1}$ . Die übrigen Formeln ergeben sich durch einfache Rechnung.  $\square$

Bei dem in Satz (9.10) beschriebenen Basisaustausch wird also von einer Basis zu einer anderen Basis übergegangen, indem eine Nichtbasisvariable  $x_{q_s}$  und eine Basisvariable  $x_{p_r}$  gesucht werden, so dass  $\bar{a}_{rs} \neq 0$  gilt. Wenn diese Bedingung erfüllt ist, kann die Nichtbasisvariable  $x_{q_s}$  zu einer Basisvariablen gemacht werden, wobei  $x_{p_r}$  nichtbasisch wird. Alle übrigen Basisvariablen bleiben erhalten.

**(9.11) Beispiel.** Wir betrachten das Gleichungssystem  $Ay = b$  gegeben durch

$$A = \begin{pmatrix} 1 & 0 & 2 & -\frac{3}{2} & -4 & \frac{11}{2} & \frac{11}{2} \\ 0 & 1 & 0 & 0 & 3 & 0 & 2 \\ 0 & 0 & -1 & \frac{1}{2} & 2 & -\frac{3}{2} & -\frac{1}{2} \\ 0 & 0 & 0 & \frac{1}{2} & 1 & -\frac{1}{2} & -\frac{1}{2} \end{pmatrix} \quad b = \begin{pmatrix} 1 \\ 2 \\ -1 \\ 2 \end{pmatrix}.$$

Es seien  $B = (1, 2, 3, 4)$ ,  $N = (5, 6, 7)$ , dann ist  $A_B$  eine Basis von  $A$ , und es gilt

$$A_B^{-1} = \begin{pmatrix} 1 & 0 & 2 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 1 \\ 0 & 0 & 0 & 2 \end{pmatrix}, \quad \bar{A} = A_B^{-1}A_N = \begin{pmatrix} 1 & 2 & 4 \\ 3 & 0 & 2 \\ -1 & 1 & 0 \\ 2 & -1 & -1 \end{pmatrix}.$$

Die zugehörige Basislösung ist  $x_B^T = (1, 2, 3, 4)$ ,  $x_N^T = (0, 0, 0)$ , d. h.  $x^T = (1, 2, 3, 4, 0, 0, 0)$ .

Das Element  $\bar{a}_{23} = 2$  ist von Null verschieden. Wir können es also als Pivotelement  $\bar{a}_{rs}$  ( $r = 2$ ,  $s = 3$ ) wählen. Daraus ergibt sich:

$$B' = (1, 7, 3, 4), \quad N' = (5, 6, 2)$$

$$E = \begin{pmatrix} 1 & -2 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & \frac{1}{2} & 0 & 1 \end{pmatrix}, \quad EA_B^{-1} = A_{B'}^{-1} = \begin{pmatrix} 1 & -2 & 2 & 1 \\ 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & -1 & 1 \\ 0 & \frac{1}{2} & 0 & 2 \end{pmatrix}$$

und

$$A_{B'}^{-1}A = EA_B^{-1}A = \begin{pmatrix} 1 & -2 & 0 & 0 & -5 & 2 & 0 \\ 0 & \frac{1}{2} & 0 & 0 & \frac{3}{2} & 0 & 1 \\ 0 & 0 & 1 & 0 & -1 & 1 & 0 \\ 0 & \frac{1}{2} & 0 & 1 & \frac{7}{2} & -1 & 0 \end{pmatrix}, \quad \bar{A}' = \begin{pmatrix} -5 & 2 & -2 \\ \frac{3}{2} & 0 & \frac{1}{2} \\ -1 & 1 & 0 \\ \frac{7}{2} & -1 & \frac{1}{2} \end{pmatrix}.$$

Die zugehörige Basislösung ist  $x_{B'}^T = (-3, 1, 3, 5)$ ,  $x_{N'}^T = (0, 0, 0)$ , daraus ergibt sich  $x^T = (-3, 0, 3, 5, 0, 0, 1)$ .

□

**(9.12) Bemerkung.** Die neue Basisinverse  $A_{B'}^{-1}$  kann also wegen  $A_{B'}^{-1} = EA_B^{-1}$  leicht aus  $A_B^{-1}$  berechnet werden. Tatsächlich ist jedoch dieses "Pivoting" der numerisch aufwendigste und schwierigste Schritt im Simplexverfahren. Es gibt hier

zahllose “Update”-Varianten, von denen wir einige (Produktform der Inversen, LU-Dekomposition, Cholesky-Zerlegung) später bzw. in den Übungen noch kurz besprechen wollen, durch die eine schnelle und numerisch stabile Berechnung von  $A_{B'}^{-1}$  erreicht werden soll. Insbesondere für große und dünn besetzte Matrizen ( $n, m \geq 5000$ ) (large-scale-linear-programming) gibt es spezielle Techniken.

□

Die äquivalente Darstellung des Gleichungssystems  $Ax = b$  in Bemerkung (9.9) gibt uns ebenso die Möglichkeit, die Kosten auf einfache Weise über die Nichtbasisvariablen zu berechnen, woraus sich ein sehr einfaches Optimalitätskriterium gewinnen lässt.

**(9.13) Satz.** Gegeben sei ein lineares Programm in Standardform (9.1), und  $A_B$  sei eine Basis von  $A$ .

(a) Für den Zielfunktionswert  $c^T x$  von  $x \in P^=(A, b)$  gilt

$$c^T x = c_B^T A_B^{-1} b + (c_N^T - c_B^T A_B^{-1} A_N) x_N.$$

Der Term  $\bar{c}^T := c_N^T - c_B^T A_B^{-1} A_N$  heißt **reduzierte Kosten** (von  $x$ ), die Komponenten  $\bar{c}_i$  von  $\bar{c}$  heißen **reduzierte Kostenkoeffizienten**.

(b) (Simplexkriterium)

Ist  $A_B$  eine zulässige Basis und sind die reduzierten Kosten nicht-positiv, d. h.

$$\bar{c}^T = c_N^T - c_B^T A_B^{-1} A_N \leq 0,$$

dann ist die zugehörige Basislösung  $x$  mit  $x_B = A_B^{-1} b$ ,  $x_N = 0$  optimal für (9.1).

(c) Ist  $A_B$  eine zulässige Basis, und ist die zugehörige Basislösung nichtdegeneriert und optimal, dann gilt  $\bar{c} \leq 0$ .

**Beweis :**

(a) Nach (9.9) gilt  $Ax = b \iff x_B = A_B^{-1} b - A_B^{-1} A_N x_N$  und damit gilt

$$\begin{aligned} c^T x &= c_B^T x_B + c_N^T x_N = c_B^T (A_B^{-1} b - A_B^{-1} A_N x_N) + c_N^T x_N \\ &= c_B^T A_B^{-1} b + (c_N^T - c_B^T A_B^{-1} A_N) x_N. \end{aligned}$$

- (b) Sei  $y \in P^=(A, b)$  beliebig. Wenn wir zeigen können, dass  $c^T x \geq c^T y$  gilt, sind wir fertig.

$$c^T y = c_B^T A_B^{-1} b + \underbrace{\bar{c}^T}_{\leq 0} \underbrace{y_N}_{\geq 0} \leq c_B^T A_B^{-1} b = c_B^T x_B = c^T x.$$

- (c) Sei die Basislösung  $x$  mit  $x_B = A_B^{-1} b$ ,  $x_N = 0$  optimal. Dann gilt für alle  $y \in P^=(A, b)$ :

$$\begin{aligned} c^T x \geq c^T y &\iff c_B^T A_B^{-1} b + \bar{c}^T x_N \geq c_B^T A_B^{-1} b + \bar{c}^T y_N \\ &\iff 0 \geq \bar{c}^T y_N. \end{aligned}$$

Angenommen die  $i$ -te Komponente  $\bar{c}_i$  von  $\bar{c}^T$  wäre größer als Null. Nach Voraussetzung ist  $x$  nichtdegeneriert, also  $A_B^{-1} b > 0$ . Sei  $e_i$  der  $i$ -te Einheitsvektor in  $\mathbb{K}^{n-m}$ , dann gibt es somit ein  $\lambda > 0$  mit

$$A_B^{-1} b \geq A_B^{-1} A_N \lambda e_i.$$

Der Vektor  $x^\lambda$  mit  $x_B^\lambda = A_B^{-1} b - A_B^{-1} A_N \lambda e_i$ ,  $x_N = \lambda e_i$  ist somit ein für (9.1) zulässiger Vektor, und es gilt

$$c^T x^\lambda = c_B^T A_B^{-1} b + \bar{c}^T \lambda e_i = c^T x + \lambda \bar{c}_i > c^T x.$$

Widerspruch! □

Aus Satz (9.10) wissen wir, wie ein Basis- bzw. Eckenaustausch vorgenommen werden kann, Satz (9.13) besagt, unter welchen Umständen wir eine zulässige Basis verbessern können, bzw. wann sie optimal ist. Setzen wir diese Informationen zusammen, so gilt:

**(9.14) Satz (Basisverbesserung).** Gegeben sei ein lineares Programm in Standardform (9.1). Sei  $A_B$  eine zulässige Basis mit Basislösung  $x$ . Sei  $\bar{A} = A_B^{-1} A_N$ ,  $\bar{b} = A_B^{-1} b$ , und  $\bar{c}^T = c_N^T - c_B^T A_B^{-1} A_N$  seien die reduzierten Kosten. Sei  $q_s \in N$  ein Index mit  $\bar{c}_s > 0$ , dann gilt

- (a) Ist  $\bar{A}_{\cdot, s} \leq 0$ , dann ist  $c^T x$  auf  $P^=(A, b)$  unbeschränkt.

- (b) Ist  $\bar{A}_{\cdot, s} \not\leq 0$ , dann setzen wir

$$\lambda_0 := \min \left\{ \frac{\bar{b}_i}{\bar{a}_{is}} \mid i = 1, \dots, m, \bar{a}_{is} > 0 \right\},$$

und wählen einen Index

$$r \in \left\{ i \in \{1, \dots, m\} \mid \frac{\bar{b}_i}{\bar{a}_{is}} = \lambda_0 \right\}.$$

Dann ist  $A_{B'}$  mit  $B' = (p_1, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m)$  eine zulässige Basis mit Basislösung  $x'$ , so dass  $c^T x' \geq c^T x$  gilt.

(c) Gelten die Voraussetzungen von (b), und ist  $A_B$  nichtdegeneriert, dann gilt  $c^T x' > c^T x$ .

**Beweis:** (a) Nach (9.9) gilt:  $y \in P^=(A, b) \iff y_B = A_B^{-1}b - A_B^{-1}A_N y_N \geq 0$  und  $y_N \geq 0$ . Für beliebiges  $\lambda \geq 0$  ist daher der Vektor  $y^\lambda \in \mathbb{K}^n$  mit  $y_N^\lambda := \lambda e_s$  und  $y_B^\lambda := A_B^{-1}b - \bar{A}y_N^\lambda = A_B^{-1}b - \lambda A_{.s}$  wegen  $e_s \geq 0$  und  $A_B^{-1}b - \lambda A_{.s} \geq A_B^{-1}b \geq 0$  in  $P^=(A, b)$ . Da  $c^T y^\lambda = c_B^T A_B^{-1}b + \bar{c}y_N^\lambda = c_B^T A_B^{-1}b + \lambda \bar{c}_s$  mit  $\lambda$  über alle Schranken wächst, ist  $c^T x$  auf  $P^=(A, b)$  unbeschränkt.

(b) Wählen wir  $r$  und  $s$  wie im Satz angegeben, so ist  $\bar{a}_{rs} > 0$ . Nach Satz (9.10) ist dann  $B' = (p_1, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m)$  eine Basis mit Basislösung  $x'$ , wobei

$$x'_{p_i} = \bar{b}_i - \frac{\bar{b}_r}{\bar{a}_{rs}} \bar{a}_{is} \begin{cases} \geq \bar{b}_i - \frac{\bar{b}_i}{\bar{a}_{is}} \bar{a}_{is}, & \text{falls } \bar{a}_{is} > 0, i \neq r \quad (\text{Wahl von } \lambda_0!), \\ \geq \bar{b}_i, & \text{falls } \bar{a}_{is} \leq 0, \end{cases}$$

$$x'_{q_s} = \frac{\bar{b}_r}{\bar{a}_{rs}} \geq 0,$$

$$x'_i = 0 \quad \text{andernfalls.}$$

Also ist  $x'$  eine zulässige Basislösung.

Wir zeigen nun  $c^T x' \geq c^T x$ . Zunächst gilt für den in Satz (9.10) definierten Vektor  $\eta$ :

$$c_{B'}^T \eta - c_{p_r} = \frac{\bar{c}_s}{\bar{a}_{rs}} > 0,$$

denn

$$\begin{aligned} 0 < \bar{c}_s &= (c_N^T - c_B^T A_B^{-1} A_N)_s = c_{q_s} - c_B^T \bar{A}_{.s} = c_{q_s} - \sum_{i=1}^m c_{p_i} \bar{a}_{is} \\ &= (c_{q_s} - \sum_{i \neq r} c_{p_i} \bar{a}_{is}) - c_{p_r} \bar{a}_{rs} = \bar{a}_{rs} (c_{B'}^T \eta - c_{p_r}), \end{aligned}$$

und aus  $\bar{a}_{rs} > 0$  folgt nun die Behauptung. Somit gilt:

$$\begin{aligned} c^T x' &= c_{B'}^T x_{B'} = c_{B'}^T A_{B'}^{-1} b \stackrel{(9.10)}{=} c_{B'}^T E A_B^{-1} b = c_{B'}^T E x_B = \\ &= \sum_{i \neq r} c_{p_i} x_{p_i} + c_{B'}^T \eta x_{p_r} \\ &= c_B^T x_B + (c_{B'}^T \eta - c_{p_r}) x_{p_r} \\ &\geq c_B^T x_B = c^T x. \end{aligned}$$

Hieraus folgt (b).

(c) Ist  $x$  nichtdegeneriert, so gilt  $x_{p_r} > 0$ , und somit ist die letzte Ungleichung in der obigen Abschätzung strikt. Daraus ergibt sich (c).  $\square$

Der obige Beweis macht den geometrischen Inhalt des Satzes (9.14) nicht deutlich. Dem Leser wird dringend empfohlen, sich zu überlegen, welche geometrische Bedeutung der Parameter  $\lambda_0$  in (b) hat und welche Beziehung  $\bar{A}_{.s}$  im Falle  $\bar{A}_{.s} \leq 0$  zum Rezessionskegel von  $P^=(A, b)$  hat.

### 9.3 Das Simplexverfahren

Die Grundidee des Simplexverfahrens haben wir schon erläutert, sie kann kurz wie folgt beschrieben werden:

- Finde eine zulässige Basis  $A_B$ .
- Konstruiere aus  $A_B$  eine zulässige Basis  $A_{B'}$ , so dass die zulässige Basislösung von  $A_{B'}$  besser als die von  $A_B$  ist.
- Gibt es keine bessere Basislösung, so ist die letzte optimal.

Die vorhergehenden Sätze geben uns nun die Möglichkeit, die oben informell beschriebene Idee, analytisch darzustellen und algorithmisch zu realisieren. Die Grundversion des Simplexverfahrens lautet dann wie folgt:

#### (9.15) Grundversion des Simplexverfahrens.

**Input:**  $A' \in \mathbb{K}^{(m,n)}$ ,  $b' \in \mathbb{K}^m$ ,  $c \in \mathbb{K}^n$

**Output:**

Eine optimale Lösung  $x$  des linearen Programms

$$(9.1) \quad \begin{aligned} & \max c^T x \\ & A'x = b' \\ & x \geq 0, \end{aligned}$$

falls eine solche existiert. Andernfalls stellt das Verfahren fest, dass (9.1) unbeschränkt ist oder keine zulässige Lösung besitzt.

**Phase I:** (Test auf Zulässigkeit, Bestimmung einer zulässigen Basislösung)

Bestimme ein Subsystem  $Ax = b, x \geq 0$  von  $A'x = b', x \geq 0$ , so dass  $P^=(A, b) = P^=(A', b')$  gilt und die Zusatzvoraussetzungen  $\text{rang}(A) = \text{Zeilenzahl von } A < \text{Spaltenzahl von } A$  erfüllt sind (falls möglich). Sei  $\text{rang}(A) = m < n$ . Bestimme weiterhin eine zulässige Basis, d. h. finde  $B = (p_1, \dots, p_m) \in \{1, \dots, n\}^m$ , so dass  $A_B$  eine zulässige Basis von  $A$  ist, und sei  $N = (q_1, \dots, q_{n-m})$  wie üblich definiert. Berechne

$$\begin{aligned} A_B^{-1} \\ \bar{A} &= A_B^{-1} A_N, \\ \bar{b} &= A_B^{-1} b, \\ \bar{c}^T &= c_N^T - c_B^T A_B^{-1} A_N, \\ c_0 &= c_B^T \bar{b}. \end{aligned}$$

Wie die gesamte Phase I durchgeführt wird und wie zu verfahren ist, wenn die gewünschten Resultate nicht erzielt werden können, wird in (9.24) angegeben. Im weiteren gehen wir davon aus, dass Phase I erfolgreich war.

**Phase II:** (Optimierung)**(II.1)** (Optimalitätsprüfung)

Gilt  $\bar{c}_i \leq 0$  für  $i = 1, \dots, n - m$ , so ist die gegenwärtige Basislösung optimal (siehe (9.13) (b)). Gib den Vektor  $x$  mit  $x_B = \bar{b}$ ,  $x_N = 0$  aus. Der Optimalwert von (9.1) ist  $c^T x = c_B^T \bar{b} = c_0$ . **STOP!**

Falls  $\bar{c}_i > 0$  für ein  $i \in \{1, \dots, n - m\}$  gehe zu (II.2).

**(II.2)** (Bestimmung der Austausch- oder Pivotspalte)

Wähle einen Index  $s \in \{1, \dots, n - m\}$ , so dass  $\bar{c}_s > 0$  gilt. (Zur konkreten Realisierung dieser Auswahl, falls mehrere reduzierte Kostenkoeffizienten positiv sind, gibt es sehr viele verschiedene Varianten, die wir in Abschnitt 10.2 besprechen werden.)

**(II.3)** (Prüfung auf Beschränktheit des Optimums)

Gilt  $\bar{A}_{\cdot s} \leq 0$ , so ist das lineare Programm unbeschränkt (siehe (9.14) (a)), STOP!

Falls  $\bar{a}_{is} > 0$  für ein  $i \in \{1, \dots, m\}$ , gehe zu (II.4).

**(II.4)** (Berechne)

$$\lambda_0 := \min \left\{ \frac{\bar{b}_i}{\bar{a}_{is}} \mid \bar{a}_{is} > 0, i = 1, \dots, m \right\}$$

**(II.5)** (Bestimmung der Austausch- oder Pivotzeile)

Wähle einen Index  $r \in \{1, \dots, m\}$ , so dass gilt

$$\frac{\bar{b}_r}{\bar{a}_{rs}} = \lambda_0.$$

(Hierzu gibt es verschiedene Möglichkeiten, die wir in 10.2 erläutern werden.)

**(II.6)** (Pivotoperation, Basisaustausch)

Setze

$$\begin{aligned} B' &:= (p_1, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m), \\ N' &:= (q_1, \dots, q_{s-1}, p_r, q_{s+1}, \dots, q_{n-m}), \\ A_{B'}^{-1} &:= EA_B^{-1} \text{ (siehe (9.10)).} \end{aligned}$$

**(II.7)** (Updating)

Berechne  $\bar{A}$ ,  $\bar{b}$ ,  $\bar{c}$ , und  $c_0$  neu!

(Hierzu gibt es viele numerische Varianten. Wir geben nachfolgend ein didaktisch klares Verfahren an, das aber numerisch umständlich und höchstens für Handrechnung und kleine lineare Programme geeignet ist. Bessere Methoden werden in den Übungen besprochen.)

(a) Neuberechnung von  $\bar{A}$ :

$$\text{Setze } \bar{a}_{rs} := \frac{1}{a_{rs}}.$$

Für  $i = 1, \dots, m, i \neq r$  führe aus

$$\bar{a}_{is} := -\bar{a}_{rs}\bar{a}_{is}.$$

Für  $j = 1, \dots, n - m, j \neq s$  führe aus

$$\bar{a}_{rj} := \bar{a}_{rs}\bar{a}_{rj}.$$

Für  $j = 1, \dots, n - m, j \neq s$ , und  $i = 1, \dots, m, i \neq r$  führe aus

$$\bar{a}_{ij} := \bar{a}_{ij} - \frac{\bar{a}_{is}}{\bar{a}_{rs}}\bar{a}_{rj} = \bar{a}_{ij} + \bar{a}_{is}\bar{a}_{rj}.$$

(b) Neuberechnung von  $\bar{b}$

$$\bar{b}_r := \frac{\bar{b}_r}{a_{rs}}.$$

Für  $i = 1, \dots, m, i \neq r$  führe aus

$$\bar{b}_i := \bar{b}_i - \frac{\bar{a}_{is}}{\bar{a}_{rs}}\bar{b}_r.$$

(c) Neuberechnung von  $\bar{c}$  und  $c_0$ .

Setze für  $j = 1, \dots, n - m, j \neq s$

$$\bar{c}_j := \bar{c}_j - \frac{\bar{a}_{rj}}{\bar{a}_{rs}}\bar{c}_s, \quad \bar{c}_s := -\frac{\bar{c}_s}{\bar{a}_{rs}}, \quad c_0 := c_0 + \bar{c}_s\frac{\bar{b}_r}{\bar{a}_{rs}}.$$

(d) Setze  $\bar{A} := \bar{A}$ ,  $\bar{b} := \bar{b}$ ,  $\bar{c} := \bar{c}$  und gehe zu (II.1).

□

**(9.16) Die Tableauform der Simplexmethode.** Zur Veranschaulichung der Simplexmethode (nicht zur Verwendung in kommerziellen Codes) benutzt man manchmal ein Tableau  $T$ . Das Tableau entspricht dem Gleichungssystem

$$\begin{pmatrix} c^T \\ A \end{pmatrix} x = \begin{pmatrix} c^T x \\ b \end{pmatrix}$$

bzw. bei gegebener Basis  $A_B$  von  $A$  und geeigneter Permutation der Spalten dem System

$$\begin{array}{|c|c|} \hline c_N^T - c_B^T A_B^{-1} A_N & 0 \\ \hline A_B^{-1} A_N & I \\ \hline \end{array} \begin{pmatrix} x_N \\ x_B \end{pmatrix} = \begin{array}{|c|} \hline c^T x - c_B^T A_B^{-1} b \\ \hline A_B^{-1} b \\ \hline \end{array}$$

□

**(9.17) Definition.** Ist  $\max\{c^T x \mid Ax = b, x \geq 0\}$  ein LP in Standardform und ist  $A_B$  eine zulässige Basis von  $A$ , so heißt die  $(m+1, n+1)$ -Matrix

$$T_B := \begin{pmatrix} c^T - c_B^T A_B^{-1} A, & -c_B^T A_B^{-1} b \\ A_B^{-1} A, & A_B^{-1} b \end{pmatrix}$$

ein Simplextableau zur Basis  $B$  mit Zeilen  $i = 0, 1, \dots, m$  und Spalten  $j = 1, \dots, n+1$ .

□

Da ein Simplextableau stets eine volle  $(m, m)$ -Einheitsmatrix enthält, ist eine solche Darstellung für den Rechner redundant. Zum Kennenlernen der Simplexmethode und zur Übung der Rechentechnik ist sie allerdings sehr nützlich.

**(9.18) Update Formeln für ein Simplextableau.** Sei  $T_B$  ein Simplextableau zur Basis  $B$ . Ist  $q_s \in N$ , so daß  $B' = (p_1, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m)$  wieder eine zulässige Basis ist, so gelten für das Simplextableau  $T_{B'}$  folgende Transformationsregeln

$$(a) \quad \begin{aligned} (T_{B'})_{i.} &= (T_B)_{i.} - \frac{\bar{a}_{is}}{\bar{a}_{rs}} (T_B)_{r.}, \quad \text{für } i = 0, \dots, m, i \neq r, \\ (T_{B'})_{r.} &= \frac{1}{\bar{a}_{rs}} (T_B)_{r.} \end{aligned}$$

oder

$$(b) \quad T_{B'} = E' \cdot T_B$$

mit

$$E' = \begin{pmatrix} 1 & & & \eta'_0 & & & \\ & \ddots & & \vdots & & & \\ & & 1 & \eta'_{r-1} & & & \\ & & & \eta'_r & & & \\ & & & \eta'_{r+1} & 1 & & \\ & & & \vdots & & \ddots & \\ & & & \eta'_m & & & 1 \end{pmatrix}, \quad \begin{aligned} \eta'_i &= -\frac{(T_B)_{is}}{\bar{a}_{rs}}, \quad i \neq r \\ \eta'_r &= \frac{1}{\bar{a}_{rs}} \end{aligned}$$

(Die Transformationsformeln (9.18) sind im wesentlichen die gleichen wie in Satz (9.10), nur dass diesmal die Zielfunktion gleich mittransformiert wird.)

**Beweis :** Formel (a):

$$T_{B'} := \begin{pmatrix} c^T - c_{B'}^T A_{B'}^{-1} A, & -c_{B'}^T A_{B'}^{-1} b \\ A_{B'}^{-1} A, & A_{B'}^{-1} b \end{pmatrix},$$

d. h. mit Satz (9.10) gilt ( $A_{B'}^{-1} = EA_B^{-1}$ ):

$$T_{B'} := \begin{pmatrix} (c^T, 0) - c_{B'}^T E(A_B^{-1} A, A_B^{-1} b) \\ E(T_B)_{\{1, \dots, m\}} \end{pmatrix}.$$

Sei  $K := \{1, \dots, m\}$ . Aus  $(T_{B'})_K = E(T_B)_K$  ergeben sich die beiden unteren Transformationsformeln in (a). Sei  $z \in \mathbb{K}^m$  und  $y = A_B^{-1} z$ , dann gilt

$$c_{B'}^T A_{B'}^{-1} z = c_{B'}^T E A_B^{-1} z = c_{B'}^T E y = c_{B'}^T y + (c_{B'}^T \eta - c_{p_r}) y_r$$

(mit  $\eta$  aus Satz (9.10)). Mit (9.13) folgt nun  $c_{B'}^T \eta - c_{p_r} = \frac{\bar{c}_s}{\bar{a}_{rs}} = \frac{(T_B)_{0s}}{\bar{a}_{rs}}$ . Also gilt:

$$\begin{aligned} (c^T, 0) - c_{B'}^T A_{B'}^{-1} (A, b) &= (c^T, 0) - c_B^T A_B^{-1} (A, b) - (T_B)_{0s} \bar{a}_{rs} (A_B^{-1} (A, b))_r \\ &= (T_B)_{0\cdot} - \frac{(T_B)_{0s}}{\bar{a}_{rs}} (T_B)_r. \end{aligned}$$

Die Formeln (b) ergeben sich direkt aus (a). □

**(9.19) Beispiel.** Wir betrachten das lineare Programm

$$\begin{aligned} \max x_1 + 2x_2 \\ x_1 &\leq 4 \\ 2x_1 + x_2 &\leq 10 \\ -x_1 + x_2 &\leq 5 \\ x_1, x_2 &\geq 0, \end{aligned}$$

dessen Lösungsmenge in Abbildung 9.4 graphisch dargestellt ist.

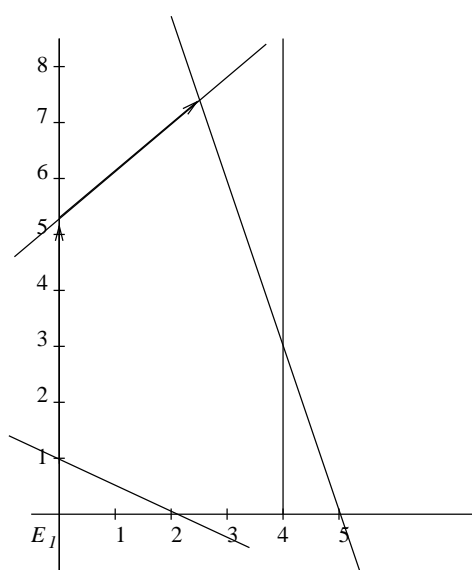


Abb. 9.4

Wir führen Schlupfvariablen  $x_3, x_4, x_5$  ein und transformieren unser LP in folgende Form.

$$\max (1, 2, 0, 0, 0) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0 & 1 & 0 & 0 \\ 2 & 1 & 0 & 1 & 0 \\ -1 & 1 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} 4 \\ 10 \\ 5 \end{pmatrix}$$

$$x_i \geq 0 \quad i = 1, \dots, 5.$$

Die Matrix  $A$  des Gleichungssystems des transformierten Systems hat 3 Zeilen und 5 Spalten und enthält eine Einheitsmatrix. Folglich hat sie vollen Zeilenrang.

Daraus folgt, dass die Spalten 3, 4, 5 eine Basis von  $A$  definieren, und diese ist zulässig, da die zugehörige Basislösung  $x_3 = 4, x_4 = 10, x_5 = 5, x_1 = 0, x_2 = 0$  nichtnegativ ist. Diese Basislösung entspricht übrigens der Ecke  $x_1 = 0, x_2 = 0$  des ursprünglichen Problems. Wir stellen nun das zur Anfangsbasis gehörige Tableau auf und führen den Simplexalgorithmus durch. Die Anfangsbasis ist gege-

ben durch  $B = (3, 4, 5)$ , und somit ist  $N = (1, 2)$ .

$$T_0 : \begin{array}{c|cccccc|c} & 1 & 2 & 3 & 4 & 5 & & \\ \hline & 1 & 2 & 0 & 0 & 0 & & 0 \\ 3 & 1 & 0 & 1 & 0 & 0 & & 4 \\ 4 & 2 & 1 & 0 & 1 & 0 & & 10 \\ 5 & -1 & \boxed{1} & 0 & 0 & 1 & & 5 \end{array} \quad \bar{a}_{rs} = \bar{a}_{32} = 1$$

Spaltenauswahl: Steilster Anstieg, d. h. der größte reduzierte Kostenkoeffizient wird gewählt. Das Pivotelement ist:  $\bar{a}_{rs} = \bar{a}_{32} = 1$ .

$$T_1 : \begin{array}{c|cccccc|c} & 3 & 0 & 0 & 0 & -2 & & -10 \\ \hline 3 & 1 & 0 & 1 & 0 & 0 & & 4 \\ 4 & \boxed{3} & 0 & 0 & 1 & -1 & & 5 \\ 2 & -1 & 1 & 0 & 0 & 1 & & 5 \end{array} \quad \bar{a}_{rs} = \bar{a}_{21} = 3$$

$$T_2 : \begin{array}{c|cccccc|c} & 0 & 0 & 0 & -1 & -1 & & -15 \\ \hline 3 & 0 & 0 & 1 & -\frac{1}{3} & \frac{1}{3} & & \frac{7}{3} \\ 1 & 1 & 0 & 0 & \frac{1}{3} & -\frac{1}{3} & & \frac{5}{3} \\ 2 & 0 & 1 & 0 & \frac{1}{3} & \frac{2}{3} & & \frac{20}{3} \end{array}$$

Das Tableau  $T_2$  ist optimal. Die optimale Lösung ist  $(\frac{5}{3}, \frac{20}{3}, \frac{7}{3}, 0, 0)$ , bzw.  $x_1 = \frac{5}{3}$ ,  $x_2 = \frac{20}{3}$ .

Wir sehen an diesem Beispiel, dass es lästig ist, immer die Einheitsmatrix im Tableau mitzuschleppen. Wenn wir uns die gegenwärtigen Basis- und Nichtbasisvariablen geeignet (z. B. am Rand des Tableaus) merken, können wir die Einheitsmatrix weglassen. Ein derartiges Tableau nennt man **verkürztes Tableau**. Wir schreiben die vorhergehende Rechnung noch einmal in dieser verkürzten Form auf.

$$\begin{array}{c|cc|cc} & 1 & 2 & & \\ \hline & 1 & 2 & & 0 \\ \hline 1 & 0 & & 4 & 3 \\ 2 & 1 & & 10 & 4 \\ -1 & 1 & & 5 & 5 \end{array}$$

$$\begin{array}{cc|cc}
 & 1 & 5 & & \\
 & 3 & -2 & -10 & \\
 \hline
 & 1 & 0 & 4 & 3 \\
 & 3 & -1 & 5 & 4 \\
 & -1 & 1 & 5 & 2 \\
 & 4 & 5 & & \\
 \hline
 & -1 & -1 & -15 & \\
 \hline
 & -\frac{1}{3} & \frac{1}{3} & \frac{7}{3} & 3 \\
 & \frac{1}{3} & -\frac{1}{3} & \frac{5}{3} & 1 \\
 & \frac{1}{3} & \frac{2}{3} & \frac{20}{3} & 2
 \end{array}$$

□

Ein weiteres Beispiel soll zur Einübung der Rechentechnik dienen.

**(9.20) Beispiel.** Wir betrachten das folgende lineare Programm.

$$\begin{aligned}
 \max \quad & x_1 + x_2 - x_3 \\
 & 3x_1 - 2x_2 \leq 3 \\
 & 2x_1 + x_3 \leq 4 \\
 & x_1 + 3x_2 - 2x_3 \leq 6 \\
 & x_1, x_2, x_3 \geq 0
 \end{aligned}$$

Wir führen Schlupfvariablen ein und erhalten das folgende zulässige Anfangstableau.

$$\begin{array}{cccccc|cc}
 & & 1 & 1 & -1 & 0 & 0 & 0 & 0 & & \\
 T_0 : & 4 & \boxed{4} & -2 & 0 & 1 & 0 & 0 & 3 & s = 1 & \\
 & 5 & 2 & 0 & 1 & 0 & 1 & 0 & 4 & r = 1 & \\
 & 6 & 1 & 3 & -2 & 0 & 0 & 1 & 6 & & \\
 \\
 & & 0 & \frac{5}{3} & -1 & -\frac{1}{3} & 0 & 0 & -1 & & \\
 T_1 : & 1 & 1 & -\frac{2}{3} & 0 & \frac{1}{3} & 0 & 0 & 1 & s = 2 & \\
 & 5 & 0 & \frac{4}{3} & 1 & -\frac{2}{3} & 1 & 0 & 2 & r = 3 & \\
 & 6 & 0 & \boxed{5} & -2 & -\frac{1}{3} & 0 & 1 & 5 & & 
 \end{array}$$

$$T_2 : \begin{array}{ccccccc|c} & & 0 & 0 & -\frac{1}{11} & -\frac{6}{33} & 0 & -\frac{5}{11} & -\frac{36}{11} \\ 1 & 1 & 0 & -\frac{4}{11} & \frac{9}{33} & 0 & \frac{2}{11} & & \frac{21}{11} \\ 5 & 0 & 0 & \frac{19}{11} & -\frac{18}{33} & 1 & -\frac{4}{11} & & \frac{2}{11} \\ 2 & 0 & 1 & -\frac{6}{11} & -\frac{1}{11} & 0 & \frac{3}{11} & & \frac{15}{11} \end{array}$$

Als Optimallösung ergibt sich somit:  $x_1 = \frac{21}{11}$ ,  $x_2 = \frac{15}{11}$ ,  $x_3 = 0$ , mit  $c^T x = \frac{36}{11}$ .  $\square$

Über die Grundversion der Simplexmethode können wir die folgende Aussage machen.

**(9.21) Satz.** Sei (9.1) ein lineares Programm, so dass  $P=(A, b) \neq \emptyset$  und alle zulässigen Basislösungen nicht entartet sind. Dann findet die Grundversion des Simplexalgorithmus (d. h. alle Schritte, bei denen eine nicht spezifizierte Auswahl besteht, werden in irgendeiner beliebigen Form ausgeführt) nach endlich vielen Schritten eine optimale Lösung oder stellt fest, dass das Problem unbeschränkt ist.

**Beweis :** Bei jedem Durchlauf des Simplexalgorithmus wird eine neue Basis erzeugt. Da jede Basislösung nicht entartet ist, hat die neue Basislösung nach (9.14) (c) einen strikt besseren Zielfunktionswert. Ist  $A_{B_1}, \dots, A_{B_i}, \dots$  die Folge der mit dem Simplexalgorithmus generierten Basen, dann folgt daraus, dass keine Basis mehrfach in dieser Folge auftreten kann. Da die Anzahl der verschiedenen Basen von  $A$  endlich ist, ist die Folge der Basen endlich. Ist  $A_{B_k}$  die letzte erzeugte Basis, so ist sie entweder optimal oder in Schritt (II.3) wurde Unbeschränktheit festgestellt.  $\square$

Besitzt (9.1) degenerierte Ecken, so kann es sein, dass die Grundversion des Simplexverfahrens irgendwann einmal nur noch zulässige Basen erzeugt, die alle zu derselben Ecke gehören und bei der die gleichen Matrizen  $\bar{A}$  immer wiederkehren. Man sagt, das Verfahren kreist oder **kreiselt**. In der Tat sind Beispiele konstruiert worden, bei denen dieses Phänomen auftreten kann, siehe (9.22).

**(9.22) Beispiel (Kreiseln).** Gegeben sei das folgende LP:

$$\begin{array}{r} \max \quad \frac{4}{5}x_1 - 18x_2 - x_3 - x_4 \\ \frac{16}{5}x_1 - 84x_2 - 12x_3 + 8x_4 \leq 0 \\ \frac{1}{5}x_1 - 5x_2 - \frac{2}{3}x_3 + \frac{1}{3}x_4 \leq 0 \\ x_1 \leq 1 \\ x_1, x_2, x_3, x_4 \geq 0 \end{array}$$

Das verkürzte Tableau zur Anfangsbasis  $B = (5, 6, 7)$  lautet:

$$\begin{array}{cccc|cc}
 & 1 & 2 & 3 & 4 & & \\
 \hline
 & \frac{4}{5} & -18 & -1 & -1 & 0 & \\
 \hline
 & \boxed{\frac{16}{5}} & -84 & -12 & 8 & 0 & 5 \\
 & \frac{1}{5} & -5 & -\frac{2}{3} & \frac{1}{3} & 0 & 6 \\
 & 1 & 0 & 0 & 0 & 1 & 7
 \end{array}$$

Wir führen nun den Simplexalgorithmus aus. Die Pivotelemente sind in den Tableaus gekennzeichnet.

$$\begin{array}{cccc|cc}
 5 & 2 & 3 & 4 & & \\
 \hline
 -\frac{1}{4} & 3 & 2 & -3 & 0 & \\
 \hline
 \frac{5}{16} & -\frac{105}{4} & -\frac{15}{4} & \frac{5}{2} & 0 & 1 \\
 -\frac{1}{16} & \boxed{\frac{1}{4}} & \frac{1}{12} & -\frac{1}{6} & 0 & 6 \\
 -\frac{5}{16} & \frac{105}{4} & \frac{15}{4} & -\frac{5}{2} & 1 & 7
 \end{array}
 \qquad
 \begin{array}{cccc|cc}
 5 & 6 & 3 & 4 & & \\
 \hline
 \frac{1}{2} & -12 & 1 & -1 & 0 & \\
 \hline
 -\frac{25}{4} & 105 & \boxed{5} & -15 & 0 & 1 \\
 -\frac{1}{4} & 4 & \frac{1}{3} & -\frac{2}{3} & 0 & 2 \\
 \frac{25}{4} & -105 & -5 & 15 & 1 & 7
 \end{array}$$

$$\begin{array}{cccc|cc}
 5 & 6 & 1 & 4 & & \\
 \hline
 \frac{7}{4} & -33 & -\frac{1}{5} & 2 & 0 & \\
 -\frac{5}{4} & 21 & \frac{1}{5} & -3 & 0 & 3 \\
 \frac{1}{6} & -3 & -\frac{1}{15} & \boxed{\frac{1}{3}} & 0 & 2 \\
 0 & 0 & 1 & 0 & 1 & 7
 \end{array}
 \qquad
 \begin{array}{cccc|cc}
 5 & 6 & 1 & 2 & & \\
 \hline
 \frac{3}{4} & -15 & \frac{1}{5} & -6 & 0 & \\
 \boxed{\frac{1}{4}} & -6 & -\frac{2}{5} & 9 & 0 & 3 \\
 \frac{1}{2} & -9 & -\frac{1}{5} & 3 & 0 & 4 \\
 0 & 0 & 1 & 0 & 1 & 7
 \end{array}$$

$$\begin{array}{cccc|cc}
 3 & 6 & 1 & 2 & & \\
 \hline
 -3 & 3 & \frac{7}{5} & -33 & 0 & \\
 4 & -24 & -\frac{8}{5} & 36 & 0 & 5 \\
 -2 & \boxed{3} & \frac{3}{5} & -15 & 0 & 4 \\
 0 & 0 & 1 & 0 & 1 & 7
 \end{array}
 \qquad
 \begin{array}{cccc|cc}
 3 & 4 & 1 & 2 & & \\
 \hline
 -1 & -1 & \frac{4}{5} & -18 & 0 & \\
 -12 & 8 & \frac{16}{5} & -84 & 0 & 5 \\
 -\frac{2}{3} & \frac{1}{3} & \frac{1}{5} & -5 & 0 & 6 \\
 0 & 0 & 1 & 0 & 1 & 7
 \end{array}$$

Das letzte Tableau ist bis auf Spaltenvertauschung das gleiche wie das erste. Alle Basen der oben generierten Folge gehören zur Ecke  $x^T = (0, 0, 0, 0)$  des Ausgangsproblems. Die bei der Berechnung des Beispiels benutzte Variante des Simplexverfahrens würde also nicht nach endlich vielen Schritten abbrechen.  $\square$

Wie Satz (9.21) zeigt, funktioniert die Grundversion der Simplexmethode, gleichgültig welche Arten von Auswahlregeln man in den Schritten (II.2) und (II.5) benutzt, wenn alle Basislösungen nicht entartet sind. Es ist daher sinnvoll zu fragen, ob man nicht das Ausgangsproblem so modifizieren (stören) kann, dass das neue Problem nicht entartet ist. Dass dies geht, zeigt der folgende Satz. Wie üblich bezeichnen wir das LP  $\max\{c^T x \mid Ax = b, x \geq 0\}$  mit (9.1). Sei für ein  $\varepsilon^T = (\varepsilon_1, \dots, \varepsilon_n)$ ,  $\varepsilon > 0$

$$\begin{aligned} \max c^T x \\ Ax = b + A\varepsilon \\ x \geq 0. \end{aligned}$$

**(9.23) Satz** (Perturbationsmethode).

- (a) Das lineare Programm (9.1) ist genau dann lösbar, wenn das LP (9.1 $\varepsilon$ ) lösbar ist für jedes  $\varepsilon > 0$ .
- (b) Es gibt ein  $\varepsilon_0 > 0$  derart, dass für alle  $0 < \varepsilon < \varepsilon_0$  jede zulässige Basis von (9.1 $\varepsilon$ ) zu einer nichtdegenerierten Basislösung von (9.1 $\varepsilon$ ) bestimmt und zum anderen eine zulässige Basis von (9.1) ist. Ferner bestimmt die zu einer optimalen Basislösung von (9.1 $\varepsilon$ ) gehörige Basis eine optimale Basislösung von (9.1).

**Beweis :** P. Kall, “Mathematische Methoden des Operations Research”, Teubner Studienbücher, Stuttgart, 1976, S. 52–53.  $\square$

Aufgrund von Satz (9.23) könnte also jedes LP durch Störung in ein nichtentartetes LP umgewandelt werden, und damit wäre die Konvergenz der Grundversion des Simplexalgorithmus gewährleistet. In der Praxis wird diese Methode jedoch kaum angewendet, da das  $\varepsilon_0$  nicht einfach zu bestimmen ist und der Ansatz mit irgendeiner Zahl  $\varepsilon > 0$  nicht zum Ziele führen muss. Ferner könnten durch eine geringe Störung des Problems durch Rundungsfehler numerische Schwierigkeiten auftreten, deren Auswirkungen nicht immer abschätzbar sind.

Wir werden im nächsten Kapitel andere Methoden zur Vermeidung des Kreisens kennenlernen. In der Praxis werden diese häufig jedoch gar nicht implementiert,

da bei praktischen Problemen trotz gelegentlicher Degeneration ein Kreisen sehr selten beobachtet wurde. Einer der Gründe dafür dürfte darin zu sehen sein, dass der Computer durch seine begrenzte Rechengenauigkeit sozusagen von allein eine Störungsmethode durchführt. So wird bei jedem Auf- oder Abrunden eine kleine Änderung des vorliegenden Programms vorgenommen, die sich positiv auf die Konvergenzeigenschaften des Simplexverfahrens auswirkt.

## 9.4 Die Phase I

Zur vollständigen Beschreibung der Grundversion der Simplexmethode fehlt noch eine Darstellung der Phase I von (9.15). Diese wollen wir nun nachholen. Ist das Ungleichungssystem  $Ax = b, x \geq 0$  gegeben, so können wir stets o. B. d. A.  $b \geq 0$  annehmen, denn Multiplikation einer Gleichung mit  $-1$  verändert das Polyeder nicht.

### (9.24) Phase I des Simplexverfahrens.

**Input:**  $A \in \mathbb{K}^{(m,n)}, b \in \mathbb{K}^m$  mit  $b \geq 0$ .

**Output:**

(a)  $P^=(A, b) = \emptyset$  oder

(b)  $P^=(A, b) = \{x\}$  oder

(c) Ausgabe von  $I \subseteq \{1, \dots, m\}$  und  $B = (p_1, \dots, p_k)$  mit folgenden Eigenschaften:

(1) Mit  $A' := A_I, b' := b_I$  gilt  $P^=(A', b') \neq \emptyset, \text{rang}(A') = |I| = k, k < n$  und  $P^=(A', b') = P^=(A, b)$ .

(2)  $A'_B$  ist eine zulässige Basis von  $A'$ .

Wir formulieren zunächst ein lineares ‘‘Hilfsprogramm’’. Sei  $D := (A, I)$  und betrachte das LP:

$$(9.25) \quad \begin{array}{ll} \max & \mathbb{1}^T Ax \\ D \begin{pmatrix} x \\ y \end{pmatrix} & = b \\ x, y & \geq 0 \end{array}$$

wobei  $x^T = (x_1, \dots, x_n)$ ,  $y^T = (y_{n+1}, \dots, y_{n+m})$  gesetzt wird. (9.25) erfüllt die Zusatzvoraussetzungen (9.2) (a), (b), (c), und mit  $B = (n+1, \dots, n+m)$  ist  $D_B$  eine zulässige Basis mit zugehöriger Basislösung  $x = 0$ ,  $y = b$ . Es gilt offenbar

$$\begin{array}{l} D \begin{pmatrix} x \\ y \end{pmatrix} = b \\ x, y \geq 0 \end{array} \iff \begin{array}{l} Ax + y = b \\ x, y \geq 0 \end{array}$$

daraus folgt  $\mathbb{1}^T Ax = \mathbb{1}^T b - \mathbb{1}^T y$ , d. h. (9.25) ist äquivalent zu  $\max\{\mathbb{1}^T b - \mathbb{1}^T y \mid Ax + y = b, x, y \geq 0\}$  bzw. zu

$$(9.26) \quad \begin{array}{ll} \mathbb{1}^T b - \min \mathbb{1}^T y & \\ Ax + y = & b \\ x, y \geq & 0. \end{array}$$

(9.26) bzw. (9.25) besagen also, dass die künstlich eingeführten Variablen möglichst klein gewählt werden sollen.

**(I.1)** Wende die Grundversion (9.15) des Simplexverfahrens auf (9.25) an. Die Matrix  $D$  hat vollen Zeilenrang und weniger Zeilen als Spalten,  $B = (n+1, \dots, n+m)$  definiert eine zulässige Basis, zu der die Matrix  $\bar{A}$  und die Vektoren  $\bar{b}$ ,  $\bar{c}$  und  $c_0$  trivial zu berechnen sind. Wir können also direkt mit diesen Daten mit Phase II beginnen.

Das LP (9.26) ist offenbar durch  $\mathbb{1}^T b - 0 = \mathbb{1}^T b$  nach oben beschränkt, also ist (9.25) durch  $\mathbb{1}^T y$  nach oben beschränkt. Der Dualitätssatz impliziert, dass (9.25) eine optimale Lösung besitzt. Somit endet das Simplexverfahren mit einer optimalen Basis  $D_B$  und zugehöriger Basislösung  $z = \begin{pmatrix} x \\ y \end{pmatrix}$ . Wir werden nun die optimale Lösung bzw. die optimale Basis analysieren, um aus Eigenschaften der Lösung bzw. Basis die Schlussfolgerungen (a), (b) oder (c) ziehen zu können.

**(I.2)** Falls  $\mathbb{1}^T Ax < \mathbb{1}^T b$ , so wird das Minimum in (9.26) nicht in  $y = 0$  angenommen. Hieraus folgt offenbar, dass  $P^=(A, b) = \emptyset$  gilt, und wir können das Verfahren mit Antwort (a) beenden, STOP!

Die folgenden Schritte behandeln den Fall  $\mathbb{1}^T Ax = \mathbb{1}^T b$ . D. h. es gilt  $\mathbb{1}^T Ax = \mathbb{1}^T b - \mathbb{1}^T y = \mathbb{1}^T b$  bzw.  $\mathbb{1}^T y = 0$ . Somit gilt  $y = 0$ , und  $x$  ist zulässig bzgl.  $Ax = b$ ,  $x \geq 0$ . Wir müssen nun noch evtl. Zeilen von  $A$  streichen, um die Rangbedingung zu erfüllen.

- (I.3)** Falls  $B \cap \{n+1, \dots, n+m\} = \emptyset$ , so befinden sich in der optimalen Basis keine künstlichen Variablen, d. h. es gilt  $D_B = A_B$ . Da  $D_B = A_B$  regulär ist, gilt  $\text{rang}(A) = m$ .

Falls  $N \cap \{1, \dots, n\} = \emptyset$ , so sind alle Nichtbasisvariablen künstliche Variablen, und wir haben  $m = n$ . Dann gilt  $x = A^{-1}b$  und  $P^=(A, b) = \{x\}$ . In diesem Falle können wir das Verfahren beenden mit Antwort (b), STOP! Der Vektor  $x$  ist die Optimallösung eines jeden LP über  $P^=(A, b)$ . Andernfalls ist  $m < n$  und  $A_B$  ein zulässige Basis von  $A$ . Wir setzen  $I = M$  und schließen das Verfahren mit Antwort (c) ab, STOP! Die beiden abschließenden Schritte sind für den Fall vorgesehen, dass sich noch künstliche Variablen in der Basis befinden, d. h. falls  $B \cap \{n+1, \dots, n+m\} \neq \emptyset$ . Wir versuchen zunächst, diese künstlichen Variablen aus der Basis zu entfernen. Da alle künstlichen Variablen  $y_{n+1}, \dots, y_{n+m}$  Null sind, ist die Basislösung  $z = (z_B, z_N)$ ,  $z_B = D_B^{-1}b$  degeneriert. Sei o. B. d. A.

$$z_B^T = (y_{p_1}, \dots, y_{p_t}, x_{p_{t+1}}, \dots, x_{p_m}),$$

d. h.  $B \cap \{n+1, \dots, n+m\} = \{p_1, \dots, p_t\}$ . Wenn wir also künstliche Variablen aus der Basis entfernen wollen, müssen wir prüfen, ob wir sie aus der Basis "hinauspivotisieren" können, d. h. ob in den Zeilen  $\bar{D}_i$  ( $i = 1, \dots, t$ ) ( $\bar{D} = (\bar{d}_{rs}) = D_B^{-1}D_N$ ) von Null verschiedene Pivotelemente vorhanden sind. Da wir ja nur Nichtbasisvariable gegen degenerierte Basisvariable austauschen, also Nullen gegen Nullen tauschen, können wir auch Pivotoperationen auf negativen Elementen zulassen.

- (I.4)** Falls es  $r \in \{1, \dots, t\}$  und  $s \in \{1, \dots, n\}$  gibt, so dass  $z_{q_s}$  keine künstliche Variable ist und  $\bar{d}_{rs} \neq 0$  gilt, dann

führe eine Pivotoperation mit dem Pivotelement  $\bar{d}_{rs}$  entsprechend Satz (9.10) durch. Wir erhalten dadurch eine neue Basis  $D_{B'}$ , mit einer künstlichen Variablen weniger, setzen  $B := B'$  und gehen zu (I.3).

- (I.5)** Falls  $\bar{d}_{ij} = 0 \forall i \in \{1, \dots, t\} \forall s \in \{1, \dots, n\}$ , für die  $z_{q_s}$  keine künstliche Variable ist, dann hat das Gleichungssystem  $z_B = D_B^{-1}b - D_B^{-1}D_N z_N$  bis auf Permutationen der Zeilen und Spalten folgende Form:

$$\begin{pmatrix} z_{p_1} \\ \vdots \\ z_{p_t} \end{pmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} - \left[ \begin{array}{c|c} 0 & * \end{array} \right] \begin{pmatrix} x_N \\ y_N \end{pmatrix}$$

und

$$\begin{pmatrix} z_{p_{t+1}} \\ \vdots \\ z_{p_m} \end{pmatrix} = \begin{bmatrix} * \\ * \\ * \end{bmatrix} - \begin{bmatrix} * & | & * \\ * & | & * \\ * & | & * \end{bmatrix} \begin{pmatrix} x_N \\ y_N \end{pmatrix}$$

(mit  $y_B = (z_{p_1}, \dots, z_{p_t})$ ,  $x_B = (z_{p_{t+1}}, \dots, z_{p_m})$ ,  $z_N = (x_N, y_N)$ ).

Setzen wir also die künstlichen Variablen Null, so ist der obere Teil des Gleichungssystems

$$y_B = 0 - 0x_N$$

unabhängig von der Belegung von  $x_N$  stets erfüllt, d. h. wir können alle Zeilen  $D_1, \dots, D_t$  streichen. Ist  $I = \{t+1, \dots, m\}$  und  $J = \{j \in \{1, \dots, n\} \mid z_{q_j} \text{ ist keine künstliche Variable}\}$  so gilt

$$x_B = (D_B^{-1}b)_I - \bar{D}_{IJ}x_N$$

und  $B = (p_{t+1}, \dots, p_m)$  ist eine zulässige Basis von  $A' := A_I$ . mit  $\text{rang}(A') = k := m - t$ .

Ist  $k = n$ , so gilt wieder  $P^=(A, b) = \{x\}$ , und wir enden mit Antwort (b) andernfalls gilt  $k < n$  und wir beenden das Verfahren mit Antwort (c), STOP!

□

Die Korrektheit des Phase I-Verfahrens (9.24) ist aufgrund der oben gemachten Bemerkungen offensichtlich. (9.24) zeigt auch, dass man bei der Lösung von linearen Programmen im wesentlichen mit der Grundversion des Simplexalgorithmus zur Lösung von Phase II auskommt. Es ist lediglich zusätzlich etwas Buchhaltung (Einführung künstlicher Variablen, Streichen von Spalten und Zeilen) zusätzlich durchzuführen. Außerdem sind u. U. im Schritt (I.4) Pivotoperationen auf negativen Pivotelementen erforderlich. Die Pivotformeln ändern sich dadurch aber nicht.

**(9.27) Bemerkung.** Liegen Probleme mit nicht vorzeichenbeschränkten Variablen vor, so wendet man die Transformationsregeln (2.4) aus Kapitel 2 an, um das Problem in ein LP mit vorzeichenbeschränkten Variablen zu transformieren.

□

Zur Auffindung einer Anfangslösung gibt es noch andere Varianten (z. B. die  $M$ -Methode), jedoch wird in kommerziellen Codes meist die Zweiphasenmethode verwendet.

**(9.28) Beispiel** (für Phase I+II). Gegeben sei das folgende lineare Programm in Standardform

$$\begin{aligned} \max x_1 - x_2 + 2x_3 \\ 2x_1 - 3x_2 + x_3 &= 3 \\ -x_1 + 2x_2 + x_3 &= -1 \\ 3x_1 - 5x_2 &= 4 \\ x_1, x_2, x_3 &\geq 0 \end{aligned}$$

Durch Multiplikation der zweiten Zeile mit  $-1$  machen wir die rechten Seiten nichtnegativ. Es gilt dann  $\mathbb{1}^T A = (6, -10, 0)$ ,  $\mathbb{1}^T b = 8$ . Wir stellen nun unser Simplextableau auf, wobei wir künstliche Variablen  $s_1, s_2, s_3$  einführen. Wir fügen dem Simplextableau eine neue oberste Zeile für die künstliche Zielfunktion hinzu.

$$T_0 : \begin{array}{cccccc|c} & x_1 & x_2 & x_3 & s_1 & s_2 & s_3 & & \\ & 6 & -10 & 0 & 0 & 0 & 0 & 0 & \leftarrow \text{künstliche Zielfunktion} \\ & 1 & -1 & 2 & 0 & 0 & 0 & 0 & \leftarrow \text{akt. Zielfunktion} \\ \hline s_1 & 2 & -3 & 1 & 1 & 0 & 0 & 3 & \\ s_2 & \boxed{1} & -2 & -1 & 0 & 1 & 0 & 1 & \\ s_3 & 3 & -5 & 0 & 0 & 0 & 1 & 4 & \end{array}$$

$$T_1 : \begin{array}{cccccc|c} & 0 & 2 & 6 & 0 & -6 & 0 & -6 & \\ & 0 & 1 & 3 & 0 & -1 & 0 & -1 & \\ \hline s_1 & 0 & 1 & \boxed{3} & 1 & -2 & 0 & 1 & \\ x_1 & 1 & -2 & -1 & 0 & 1 & 0 & 1 & \\ s_3 & 0 & 1 & 3 & 0 & -3 & 1 & 1 & \end{array}$$

$$T_2 : \begin{array}{ccccccc|c} & x_1 & x_2 & x_3 & s_1 & s_2 & s_3 & \\ \hline & 0 & 0 & 0 & -2 & -2 & 0 & -8 \\ & 0 & 0 & 0 & -1 & 1 & 0 & -2 \\ \hline x_3 & 0 & \frac{1}{3} & 1 & \frac{1}{3} & -\frac{2}{3} & 0 & \frac{1}{3} \\ x_1 & 1 & -\frac{5}{3} & 0 & \frac{1}{3} & \frac{1}{3} & 0 & \frac{4}{3} \\ s_3 & 0 & 0 & 0 & -1 & -1 & 1 & 0 \end{array}$$

Das Tableau  $T_2$  ist optimal bezüglich der künstlichen Zielfunktion. Es gilt  $\emptyset^T Ax = \emptyset^T b$ .  $B \cap \{n+1, \dots, n+m\} = \{n+3\} = \{p_3\} \neq \emptyset$  (d. h.  $s_3$  ist in der Basis). Schritt (I.4) wird nicht ausgeführt, da  $\bar{d}_{32} = 0$ , also dürfen wir entsprechend (I.5) die letzte Zeile und die Spalten 4, 5 und 6 streichen. Dies ergibt das folgende zulässige Tableau für das ursprüngliche Problem

$$\begin{array}{ccc|c} & x_1 & x_2 & x_3 \\ \hline & 0 & 0 & 0 \\ \hline x_3 & 0 & \frac{1}{3} & 1 \\ x_1 & 1 & -\frac{5}{3} & 0 \end{array} \quad \begin{array}{c} -2 \\ \frac{1}{3} \\ \frac{4}{3} \end{array}$$

Dieses Tableau ist optimal bezüglich der ursprünglichen Zielfunktion.  $x_1 = \frac{4}{3}$ ,  $x_2 = 0$ ,  $x_3 = \frac{1}{3}$ ,  $c^T x = 2$ .

# Kapitel 10

## Varianten der Simplex-Methode

Nachdem wir nun wissen, wie der Simplexalgorithmus im Prinzip funktioniert, wollen wir uns in diesem Kapitel überlegen, wie die noch offen gelassenen Schritte am besten ausgeführt und wie die verschiedenen Updates numerisch günstig implementiert werden können.

### 10.1 Der revidierte Simplexalgorithmus

Betrachtet man die Grundversion des Simplexalgorithmus (9.15), so stellt man fest, dass während einer Iteration i. a. nicht sämtliche Spalten der Matrix  $\bar{A}$  benötigt werden. Der revidierte Simplexalgorithmus nutzt diese Tatsache aus, indem er stets nur solche Spalten von  $\bar{A}$  berechnet, die im jeweiligen Schritt benötigt werden.

Alle numerischen Implementationen des Simplexverfahrens benutzen als Grundgerüst den revidierten Simplexalgorithmus. Wir wollen deshalb im folgenden den revidierten Simplexalgorithmus in Unterprogramme (Subroutines) zerlegen und später zeigen, wie diese Subroutines bei den jeweiligen numerischen Implementationsvarianten realisiert werden. Wir gehen immer davon aus, dass wir ein LP in Standardform (9.1) vorliegen haben. Alle Bezeichnungen werden, wie in Kapitel 9 festgelegt, verwendet.

**(10.1) Algorithmus (Revidierter Simplexalgorithmus).** Gegeben seien:  $A, A_B^{-1}, \bar{b}, c$  und  $B$ .

(1) **BTRAN** (*Backward Transformation*)

Berechnet:  $\pi^T := c_B^T A_B^{-1}$  (Schattenpreise)

(2) **PRICE** (Pivotspaltenauswahl)

Berechnet die reduzierten Kostenkoeffizienten:  $\bar{c}_j := (c_N^T)_j - \pi^T A_N e_j$ , für  $j = 1, \dots, n - m$  und wählt einen Index  $s$  mit  $\bar{c}_s > 0$ .

(1) **FTRAN** (Forward Transformation)

Aktualisiert die Pivotspalte:

$$\bar{d} := A_B^{-1} d = A_B^{-1} A_{\cdot q_s}.$$

(4) **CHUZR** (Pivotzeilenauswahl)

Berechnet:  $\lambda_0 = \min\{\frac{\bar{b}_i}{\bar{d}_i} \mid \bar{d}_i > 0, i = 1, \dots, m\}$  und wählt einen Index  $r \in \{i \in \{1, \dots, m\} \mid \frac{\bar{b}_i}{\bar{d}_i} = \lambda_0\}$ .

(4) **WRETA** (Updating der Basis)

Entferne das  $r$ -te Element von  $B$  und ersetze es durch  $q_s$ . Der neue Spaltenindexvektor heie  $B'$ .

Berechnet:  $A_{B'}^{-1}, \bar{b} = A_{B'}^{-1} b$ .

(4) Abbruchkriterium wie blich (siehe Schritt (II.1) von (9.15)).

□

Wir werden spter und in den bungen auf numerische Tricks zur Implementation der Updates eingehen. Es seien hier jedoch bereits einige Vorteile der revidierten Simplexmethode festgehalten:

- Wesentlich weniger Rechenaufwand, speziell bei Programmen mit erheblich mehr Variablen als Gleichungen.
- Die Ausgangsmatrix  $A$  kann auf externem Speicher bzw. bei dnn besetzten Matrizen durch spezielle Speichertechniken gehalten werden. Spalten von  $A$  werden je nach Bedarf generiert. (Sparse-Matrix-Techniques).
- Die Kontrolle ber die Rechengenauigkeit ist besser. Insbesondere kann bei Bedarf mit Hilfe eines Inversionsunterprogramms  $A_B^{-1}$  neu berechnet werden.
- Es gibt besondere Speichertechniken, die eine explizite Speicherung von  $A_B^{-1}$  vermeiden. Man merkt sich lediglich die Etavektoren und berechnet dann ber die Formel aus (9.10) aus der ersten Basisinversen durch Linksmultiplikation mit den Elementarmatrizen die gegenwrtige Basisinverse.

## 10.2 Spalten- und Zeilenauswahlregeln

Wir wollen uns nun damit beschäftigen, wie die noch unspezifizierten Schritte des Simplexverfahrens “vernünftig” ausgeführt werden können.

**(10.2) Varianten der PRICE-Routine (10.1) (2).** Sei  $S := \{j \in \{1, \dots, n\} \mid \bar{c}_j > 0\} \neq \emptyset$ . Folgende Regeln sind in der Literatur eingehend diskutiert worden und werden bei schlichten Implementationen des Simplexverfahrens angewendet.

- (1) **Kleinster-Index-Regel:** Wähle  $s = \min\{j \in S\}$
- (2) **Kleinster-Variablenindex-Regel:** Wähle  $s \in S$ , so dass  $q_s = \min\{q_j \in N \mid j \in S\}$ .
- (3) **Steilster-Anstieg-Regel:** Wähle  $s \in S$ , so dass  $\bar{c}_s = \max\{\bar{c}_j \mid j \in S\}$ .
- (4) **Größter-Fortschritt-Regel:** Berechne für jedes  $j \in S$  zunächst  $\lambda_0^j = \min\{\frac{\bar{b}_i}{\bar{a}_{ij}} \mid \bar{a}_{ij} > 0, i = 1, \dots, m\}$  und  $g_j = \bar{c}_j \lambda_0^j$ .  
Wähle  $s \in S$ , so dass  $g_s = \max\{g_j \mid j \in S\}$ .
- (5) **Varianten von (4):** Es gibt unzählige Varianten von (4), einige davon sind beschrieben in:

D. Goldfarb, J. K. Reid: “A practicable steepest-edge simplex algorithm”, *Mathematical Programming* 12 (1977), 361–371.

P. M. S. Harris: “Pivot selection methods for the Devex LP code”, *Mathematical Programming* 5 (1973), 1–28.

H. Crowder, J. M. Hattingh: “Partially normalized pivot selection in linear programming”, *Math. Progr. Study* 4 (1975), 12–25.

□

Die Regel (10.2) (1) ist die rechentechnisch einfachste. Man durchläuft den Vektor  $\bar{c}$ , sobald ein Index  $s$  mit  $\bar{c}_s > 0$  gefunden ist, hört man auf und geht zum nächsten Schritt. Die reduzierten Kosten müssen also nicht alle berechnet werden. Dadurch wird viel Aufwand gespart, aber aufgrund der simplen Wahl von  $s$  ist die Gesamtzahl der Pivotoperationen im Vergleich zu anderen Regeln recht hoch.

Ähnliches gilt für Regel (2). Hier müssen jedoch alle reduzierten Kostenkoeffizienten berechnet werden. Diese Regel ist jedoch aus einem theoretischen Grund, der noch diskutiert werden soll, interessant.

Die Regel (3) ist die Regel, die bei einfachen Implementationen (keine sophistifizierten Tricks) am häufigsten verwendet wird und bei Problemen bis zu mittleren Größenordnungen (jeweils bis zu 500 Variablen und Zeilen) recht gute Ergebnisse zeigt. Ihr liegt die Idee zu Grunde, dass diejenige Variable in die Basis genommen werden sollte, die pro Einheit den größten Zuwachs in der Zielfunktion bringt.

Es kann natürlich sein, dass die Nichtbasisvariable mit dem größten Zuwachs pro Einheit nur um sehr wenige Einheiten erhöht werden kann und dass ein Wechsel zu einer anderen Basis einen wesentlich größeren Gesamtfortschritt bringt. Hierzu ist Regel (4) geschaffen. Bei ihr wird für jeden möglichen Basiswechsel der tatsächliche Zuwachs  $g_j$  der Zielfunktion berechnet, und es wird der Basiswechsel vorgenommen, der den insgesamt größten Fortschritt bringt. Diese Verbesserung wird natürlich durch einen erheblichen rechnerischen Mehraufwand erkauft, bei dem es fraglich ist, ob er sich lohnt.

Aufgrund von Erfahrungen in der Praxis kann gesagt werden, dass sich die trivialen Regeln (1), (2) und (3) für kleinere bis mittlere Problemgrößen bewährt haben, jedoch für große LPs, Modifizierungen von (10.2) (4), wie sie in den Literaturangaben beschrieben sind, benutzt werden. Die trivialen Regeln führen im allgemeinen zu insgesamt mehr Pivotoperationen. Die komplizierten Regeln versuchen die Anzahl der Pivotoperationen minimal zu gestalten. Hierbei ist ein wesentlich höherer Rechenaufwand erforderlich (viele Spalten von  $\bar{A}$  sind zu generieren und für jede Spalte ist  $\lambda_0$  zu berechnen), der bei kleineren Problemen durchaus einem mehrfachen Basisaustausch entspricht und sich daher nicht auszahlt. Bei großen Problemen ist i. a. der Basiswechsel rechenaufwendiger als die komplizierte Auswahl der Pivotspalte. Hier zahlt sich dann die sorgfältige Auswahl der Pivotspalte bei der Gesamtrechenzeit aus.

Bei der Auswahl von Pivotzeilen gibt es nicht so viele interessante Regeln, aber hier kann durch geschickte Wahl die Endlichkeit des Simplexverfahrens erzwungen werden.

Im folgenden gehen wir davon aus, dass der Spaltenindex  $s$  durch eine Spaltenauswahlregel (10.2) bestimmt wurde und setzen  $R := \{i \in \{1 \dots, m\} \mid \frac{\bar{b}_i}{\bar{a}_{is}} = \lambda_0\}$ . Wir treffen zunächst eine Definition.

**(10.3) Definition.** Ein Vektor  $x^T = (x_1, \dots, x_n) \in \mathbb{K}^n$  heißt **lexikographisch positiv** (wir schreiben:  $x \succ 0$ ), wenn die erste von Null verschiedene Komponente positiv ist (d. h. ist  $i := \min \text{supp}(x)$ , dann ist  $x_i > 0$ ). Wir schreiben  $x \succ y$ , falls  $x - y \succ 0$  gilt.

□

**(10.4) Bemerkung.** “ $\succ$ ” definiert eine totale Ordnung im  $\mathbb{K}^n$ . Wir schreiben  $x \succeq y \iff x \succ y \vee x = y$  und bezeichnen mit  $\text{lex-min } S$  das lexikographische Minimum einer endlichen Menge  $S \subseteq \mathbb{K}^n$ .

□

**(10.5) 1. Lexikographische Zeilenauswahlregel.** Wähle  $r \in R$  so dass

$$\frac{1}{\bar{a}_{rs}}(A_B^{-1}A)_r. = \text{lex-min}\left\{\frac{1}{\bar{a}_{is}}(A_B^{-1}A)_i. \mid \bar{a}_{is} > 0, i \in \{1, \dots, m\}\right\}.$$

□

Mit Regel (10.5) wird jedes Simplexverfahren unabhängig von der Spaltenauswahlregel endlich.

**(10.6) (Endlichkeit des Simplexalgorithmus).** Wird im Simplexverfahren die 1. Lexikographische Regel (10.5) verwendet, so endet der Algorithmus nach endlich vielen Schritten, gleichgültig, welche Spaltenauswahlregel verwendet wird.

**Beweis :** Wir zeigen zuerst:

$$(a) \quad \frac{1}{\bar{a}_{rs}}(A_B^{-1}A)_r. \prec \frac{1}{\bar{a}_{is}}(A_B^{-1}A)_i. \quad \forall i \neq r,$$

das heißt, das lexikographische Minimum wird eindeutig angenommen. Seien  $r$  und  $r'$  Indizes mit

$$(*) \quad \frac{1}{\bar{a}_{rs}}(A_B^{-1}A)_r. = \frac{1}{\bar{a}_{r's}}(A_B^{-1}A)_{r'}.$$

und o. B. d. A. sei  $A = (A_B, A_N)$ . Dann gilt  $(A_B^{-1}A)_r. = (A_B^{-1})_r. A = (e_r^T, \bar{A}_r.)$  und  $(A_B^{-1}A)_{r'}. = (e_{r'}^T, \bar{A}_{r'}.)$  und  $(*)$  impliziert  $(e_r^T, \bar{A}_r.) = \frac{\bar{a}_{rs}}{\bar{a}_{r's}}(e_{r'}^T, \bar{A}_{r'}.)$ . Damit folgt  $e_r = \frac{\bar{a}_{rs}}{\bar{a}_{r's}}e_{r'}$  und somit  $r = r'$ . Also gilt (a).

Sei  $A_B$  nun eine zulässige Basis von  $A$  (z. B. die Startbasis (evtl. auch die der Phase I)) und sei  $T_B$  das zugehörige Simplextableau (vgl. (9.16)). O. B. d. A. sei

$$T_B = \begin{pmatrix} 0 & \bar{c} & -c_B^T A_B^{-1}b \\ I & \bar{A} & A_B^{-1}b \end{pmatrix}$$

(evtl. ist eine Permutation der Spalten der Ausgangsmatrix notwendig). Sei  $B' = (p_1, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m)$ . Dann gilt

$$(b) \quad (T_{B'})_r. \succ 0$$

denn  $(T_B)_{r.} \succ 0$  und  $\bar{a}_{rs} > 0$  (vgl. Formel (9.18) (a)). Da nach Voraussetzung  $(T_B)_{i.} \succ 0 \forall i = 1, \dots, m$ , folgt

$$(c) \quad (T_{B'})_{i.} \succ 0 \forall i \text{ mit } \bar{a}_{is} \leq 0$$

vgl. (9.18) (a).

Aus (a) folgt nun  $(T_B)_{i.} \succ \frac{\bar{a}_{is}}{\bar{a}_{rs}}(T_B)_{r.}$  für alle  $i \neq r$  mit  $\bar{a}_{is} > 0$ , d. h.

$$(d) \quad (T_{B'})_{i.} \succ 0 \forall i \text{ mit } \bar{a}_{is} > 0.$$

Somit sind bei Anwendung der 1. Lexikographischen Regel stets alle Zeilen  $(T_B)_{i.}$  ( $i \neq 0$ ) des Simplextableaus lexikographisch positiv, d. h. es gilt für 2 aufeinanderfolgende Basen  $B$  und  $B'$  (vgl. (9.18) (a))

$$(T_{B'})_{o.} - (T_B)_{o.} = -\frac{(T_B)_{os}}{\bar{a}_{rs}}(T_B)_{r.} \prec 0,$$

denn  $(T_B)_{os} > 0$ . Also sind alle Simplextableaus verschieden und da nur endlich viele zulässige Basislösungen und somit Simplextableaus existieren, folgt nun die Behauptung.  $\square$

**(10.7) 2. Lexikographische Zeilenauswahlregel.** Wähle  $r \in R$ , so dass

$$\frac{1}{\bar{a}_{rs}}(A_B^{-1})_{r.} = \text{lex-min} \left\{ \frac{1}{\bar{a}_{is}}(A_B^{-1})_{i.} \mid \bar{a}_{is} > 0, i \in \{1, \dots, m\} \right\}.$$

$\square$

Die Anwendung der 2. Lexikographischen Regel garantiert auch Endlichkeit des Simplexverfahrens, analog zu Satz (10.6) (siehe Kall oder Dantzig S. 269, diese folgt aus der Störungsmethode).

**(10.8) Weitere Zeilenauswahlregeln.**

- (1) **Kleinster-Index-Regel:**  $r := \min R$ .
- (2) **Kleinster-Variablenindex-Regel:** Wähle  $r \in R$ , so dass  $p_r = \min\{p_i \in B \mid i \in R\}$ .

Keine der beiden Regeln in (10.8) kann für sich allein Endlichkeit des Simplexverfahrens bewirken. Aber eine Kombination von (10.8) (2) und (10.2) (2) schafft es.

**(10.9) Bland-Regel:** R. Bland (“New finite pivoting rules for the simplex method”, *Mathematics of Operations Research* 2 (1977), 103–107) hat gezeigt, dass das Simplexverfahren auch endlich ist, wenn sowohl bei der Spaltenauswahl (PRICE-Routine) als auch bei der Zeilenauswahl die Kleinster-Variablen-Index Regel angewendet wird.

□

Von Avis und Chvátal (“Notes on Bland’s pivoting rule”, *Mathematical Programming Studies* 8 (1978), 24–34) sind Untersuchungen über die Praktikabilität der Bland-Regel gemacht worden. Dabei hat sich gezeigt, dass i. a. die Anzahl der Pivotoperationen bei Anwendung der Bland-Regel wesentlich höher liegt als z. B. bei der Steilster-Anstieg Regel (10.2) (3). Für Computerimplementierungen ist die Bland-Regel selbst bei hochdegenerierten Problemen nicht gut geeignet.

**(10.10) Worst-Case Verhalten des Simplexverfahrens.** Zu fast allen bekannten Pivotauswahlregeln (speziell zu allen, die hier genannt wurden) kennt man heute eine Klasse von Polyedern und Zielfunktionen, so dass der Simplexalgorithmus bei Anwendung einer zugehörigen Auswahlregel durch alle Ecken des Polyeders läuft. Da die Anzahl der Ecken dieser Polyeder exponentiell mit der Anzahl der Variablen wächst, sagt man, dass das Simplexverfahren ein exponentielles worst-case Verhalten besitzt.

### 10.3 Die Behandlung oberer Schranken

In linearen Programmen kommen häufig Beschränkungen der folgenden Form vor:

$$0 \leq x \leq u.$$

Da diese strukturell sehr einfach sind, kann man sie algorithmisch besser behandeln als allgemeine Ungleichungen. Normalerweise führen wir Ungleichungen durch Einfügung von Schlupfvariablen wie folgt in Gleichungen und Nichtnegativitätsbedingungen über:

$$x + \bar{x} = u, \quad x \geq 0, \quad \bar{x} \geq 0.$$

Ist jedoch der Wert von  $x$  (bzw.  $\bar{x}$ ) festgelegt, so ist der Wert der Komplementärvariablen  $\bar{x}$  (bzw.  $x$ ) bereits eindeutig bestimmt. Diesen Vorteil kann man nun so ausnutzen, dass man im Simplexverfahren nur eine der beiden Variablen mit-schleppt und die zusätzliche Gleichung völlig weglässt. Für jede Zeile oder Spalte muss jedoch festgehalten werden, ob sie  $x$  oder der Komplementärvariablen

$\bar{x} = u - x$  entspricht. Die Zeilenauswahlregeln sind ebenfalls etwas komplizierter, da eine neue in die “Basis” hineinzunehmende Variable nur maximal bis zu ihrer Beschränkung erhöht werden darf und die übrigen Basisvariablen ebenfalls ihre Schranken nicht überschreiten dürfen.

Wir wollen im folgenden die sogenannte **Obere-Schranken-Technik** zur Behandlung von linearen Programmen der Form

$$(10.11) \quad \begin{aligned} \max c^T x \\ Ax = b \\ 0 \leq x \leq u \end{aligned}$$

besprechen. Diese “upper-bound-technique” behandelt nur explizit nach oben beschränkte Variablen. Sind einige der Variablen nicht explizit nach oben beschränkt, so legen wir, um die weitere Diskussion und die Formeln zu vereinfachen, fest, dass ihre obere Schranke  $+\infty$  ist. Das heißt, die Komponenten von  $u$  haben Werte aus der Menge  $\mathbb{R}_+ \cup \{+\infty\}$ .

Weiter benötigen wir eine erweiterte Definition von Basis und Nichtbasis. Wir halten uns im Prinzip an Konvention (9.3), lassen jedoch zu, dass der Vektor der Indizes von Nichtbasisvariablen positive und negative Indizes enthalten kann. Durch das Vorzeichen wollen wir uns merken, ob eine Nichtbasisvariable die obere oder untere Schranke annimmt, und zwar legen wir fest, dass für eine Nichtbasisvariable  $q_s$  der Wert  $x_{q_s}$  Null ist, falls das Vorzeichen von  $q_s$  positiv ist, andernfalls ist der Wert von  $x_{q_s}$  die obere Schranke  $u_{q_s}$ . Um dies formeltechnisch einfach aufschreiben zu können, treffen wir die folgenden Vereinbarungen. Ist  $B = (p_1, \dots, p_m)$  ein Spaltenindexvektor, so dass  $A_B$  eine Basis ist und  $N = (q_1, \dots, q_{n-m})$  wie in (9.3) definiert, dann sei

$$(10.12) \quad \bar{N} := (\bar{q}_1, \dots, \bar{q}_{n-m}) \text{ mit } \bar{q}_i = q_i \text{ oder } \bar{q}_i = -q_i$$

$$(10.13) \quad \begin{aligned} N^+ &:= (q_i \mid \bar{q}_i > 0) \\ N^- &:= (q_i \mid \bar{q}_i < 0) \end{aligned}$$

Die in  $B$ ,  $N^+$  und  $N^-$  enthaltenen Indizes bilden also eine Partition der Spaltenindexmenge  $\{1, \dots, n\}$  von  $A$ . Bei einer Rechnerimplementation genügt es natürlich, sich  $\bar{N}$  zu merken, denn  $N = (|\bar{q}_1|, \dots, |\bar{q}_{n-m}|)$ . Ist  $A_B$  regulär, so nennen wir  $A_B$  eine **E-Basis** (erweiterte Basis) von  $A$  und  $\bar{N}$  eine **E-Nichtbasis** von  $A$ . Der Vektor  $x \in \mathbb{K}^n$  mit

$$(10.14) \quad \begin{aligned} x_{N^+} &= 0 \\ x_{N^-} &= u_{N^-} \\ x_B &= A_B^{-1}b - A_B^{-1}A_{N^-}u_{N^-} \end{aligned}$$

heißt **E-Basislösung** zur E-Basis  $A_B$ .

Eine  $\left\{ \begin{array}{l} \text{E-Basis } A_B \\ \text{E-Basislösung } x \end{array} \right\}$  heißt **zulässig**, wenn

$0 \leq x_B \leq u_B$  und heißt **nicht degeneriert (nicht entartet)**, falls  $0 < x_B < u_B$ , andernfalls **degeneriert (entartet)**. Gibt es keine oberen Schranken, so gilt offenbar  $N = \bar{N} = N^+$  und alle Formeln reduzieren sich auf den bereits behandelten Fall.

**(10.15) Satz.** Gegeben sei ein Polyeder  $P := \{x \in \mathbb{K}^n \mid Ax = b, 0 \leq x \leq u\}$  mit  $\text{rang}(A) = m$ . Ein Vektor  $x \in \mathbb{K}^n$  ist genau dann eine Ecke von  $P$ , wenn  $x$  zulässige E-Basislösung ist.

**Beweis :** Es gilt  $P = P(D, d)$  mit

$$D = \begin{pmatrix} A \\ -A \\ I \\ -I \end{pmatrix} \quad d = \begin{pmatrix} b \\ -b \\ u \\ 0 \end{pmatrix}.$$

Also folgt mit Satz (8.9):  $x$  Ecke von  $P \iff \text{rang}(D_{\text{eq}(\{x\})}) = n$ . Seien  $J = \text{eq}(\{x\})$  und  $J_1, J_2, J_3, J_4$  so gewählt, dass

$$D_J = \begin{array}{|c|cc|} \hline A_{J_1} & & \\ -A_{J_2} & * & \\ \hline 0 & I_{J_3} & 0 \\ \hline 0 & 0 & -I_{J_4} \\ \hline \end{array}$$

Gilt  $\text{rang}(D_J) = n$ , so besitzt  $A_{J_1 \cup J_2}$  vollen Rang und mit  $K := \{1, \dots, n\} \setminus (J_3 \cup J_4)$  ist  $A_B := A_{J_1 \cup J_2, K}$  regulär. Man verifiziert sofort, dass  $A_B$  zulässig ist, wenn  $N^- = J_3, N^+ = J_4$  gewählt wird. Die Rückrichtung ist nun evident.  $\square$

Wir kennzeichnen also im Weiteren Basen, Nichtbasen und Basislösung von (10.11) mit einem  $E$ , um sie von den in Kapitel 9 eingeführten Begriffen zu unterscheiden.

**(10.16) Algorithmus.** (Upper-Bound-Technique)

zur Lösung von linearen Programmen der Form (10.11).

Wie nehmen an, dass  $u \in (\mathbb{R}_+ \cup \{+\infty\})^n$  gilt und dass eine zulässige  $E$ -Basis  $A_B$  vorliegt. Wie beschreiben lediglich die Phase II. Es seien also gegeben: zulässige  $E$ -Basis  $A_B$ ,  $A_B^{-1}$ ,  $\bar{b} = A_B^{-1}b$ ,  $c$ ,  $B = (p_1, \dots, p_m)$  und  $\bar{N} = (\bar{q}_1, \dots, \bar{q}_{n-m})$  ( $N$ ,  $N^+$ ,  $N^-$  können daraus bestimmt werden),  $\bar{\bar{b}} = A_B^{-1}b - A_B^{-1}A_{N^-}u_{N^-}$ .

**(1) BTRAN:**

Berechne  $\pi^T := c_B^T A_B^{-1}$ .

**(2) PRICE:**

Berechne  $\bar{c}^T := c_N^T - \pi^T A_N$  und wähle ein  $s \in \{1, \dots, n-m\}$  mit

$$\bar{c}_s \begin{cases} > 0 & \text{falls } \bar{q}_s > 0 \\ < 0 & \text{falls } \bar{q}_s < 0 \end{cases}$$

mittels irgendeiner der Pivotspaltenauswahlregeln. Gibt es keinen solchen Index  $s$ , so ist die aktuelle  $E$ -Basislösung optimal.

**Begründung:** Offenbar gilt  $x \in \{y \in \mathbb{K}^n \mid Ay = b, 0 \leq y \leq u\} \iff x_B = A_B^{-1}b - A_B^{-1}A_{N^+}x_{N^+} - A_B^{-1}A_{N^-}x_{N^-}$  und  $0 \leq x_B, x_{N^+}, x_{N^-} \leq u$ . Also ist

$$\begin{aligned} c^T x &= c_B^T x_B + c_{N^+}^T x_{N^+} + c_{N^-}^T x_{N^-} \\ &= c_B^T A_B^{-1}b + (c_{N^+}^T - c_B^T A_B^{-1} A_{N^+})x_{N^+} + (c_{N^-}^T - c_B^T A_B^{-1} A_{N^-})x_{N^-} \\ &= \pi^T b + (c_{N^+}^T - \pi^T A_{N^+})x_{N^+} + (c_{N^-}^T - \pi^T A_{N^-})x_{N^-}. \end{aligned}$$

Da für die gegenwärtige  $E$ -Basislösung gilt  $x_{N^-} = u_{N^-}$ ,  $x_{N^+} = 0$ , kann der Zielfunktionswert nur verbessert werden, wenn für eine Nichtbasisvariable  $q_s$  entweder gilt  $\bar{q}_s > 0$  und  $\bar{c}_s = (c_N^T - \pi^T A_N)_s > 0$  oder  $\bar{q}_s < 0$  und  $\bar{c}_s < 0$ .

**(3) FTRAN:**

Berechne  $\bar{d} := A_B^{-1}A_{q_s} = \bar{A}_{.s}$ .

**(4) CHUZR:**

Setze  $\sigma = \text{sign}(\bar{q}_s)$  und berechne

$$\begin{aligned} \lambda_0 &:= \min \left\{ \frac{\bar{b}_i}{\sigma \bar{a}_{is}} \mid \sigma \bar{a}_{is} > 0, i = 1, \dots, m \right\} \\ \lambda_1 &:= \min \left\{ \frac{\bar{b}_i - u_{p_i}}{\sigma \bar{a}_{is}} \mid \sigma \bar{a}_{is} < 0, i = 1, \dots, m \right\} \\ \lambda_2 &:= u_{q_s} \end{aligned}$$

- (a) Ist keine der Zahlen  $\lambda_0, \lambda_1, \lambda_2$  endlich, so ist das Programm (10.11) unbeschränkt (das kann natürlich nur vorkommen, wenn  $x_{q_s}$  nicht explizit beschränkt ist).
- (b) Ist  $\lambda_0 = \min\{\lambda_0, \lambda_1, \lambda_2\}$ , so wähle einen Index

$$r \in \{i \in \{1, \dots, m\} \mid \frac{\bar{b}_i}{\sigma \bar{a}_{is}} = \lambda_0, \sigma \bar{a}_{is} > 0\}$$

*mittels irgendeiner Pivotzeilenauswahlregel.*

- (c) Ist  $\lambda_1 = \min\{\lambda_0, \lambda_1, \lambda_2\}$  so wähle

$$r \in \{i \in \{1, \dots, m\} \mid \frac{\bar{b}_i - u_{p_i}}{\sigma \bar{a}_{is}} = \lambda_1, \sigma \bar{a}_{is} < 0\}$$

*mittels irgendeiner Pivotzeilenauswahlregel.*

- (d) Ist  $\lambda_2 = \min\{\lambda_0, \lambda_1, \lambda_2\}$  so setze  $\bar{q}_s := -\bar{q}_s$ , berechne  $\bar{b}$  neu und gehe zurück zur PRICE-Routine (2).

##### (5) WRETA:

Setze  $B' = (p_1, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m)$ ,  $\bar{N} = (\bar{q}_1, \dots, \bar{q}_{s-1}, \bar{q}, \bar{q}_{s+1}, \dots, \bar{q}_{n-m})$ , wobei  $\bar{q} = p_r$ , falls  $\lambda_0 = \min\{\lambda_0, \lambda_1, \lambda_2\}$ , bzw.  $\bar{q} = -p_r$ , falls  $\lambda_1 = \min\{\lambda_0, \lambda_1, \lambda_2\}$ . Berechne  $A_{B'}^{-1}$ ,  $\bar{b} = A_{B'}^{-1}b - A_{B'}^{-1}A_{N'}u_{N'}$ .

**Begründung:** Die Pivotzeilenauswahl (CHUZR-Routine) dient dazu, die Zeilenauswahl so zu bestimmen, dass die transformierten Variablen nach erfolgter Pivotoperation wieder zulässig sind. Entsprechend der Transformationsregeln beim Pivotisieren mit dem Pivotelement  $\bar{a}_{rs}$  (vergleiche (9.10)) muss gewährleistet sein, dass nach Ausführung des Pivotschrittes für die neue  $E$ -Basislösung folgendes gilt:

$$\alpha) 0 \leq x'_{p_i} = \bar{b}_i - \frac{\bar{a}_{is}}{\bar{a}_{rs}} \bar{b}_r \leq u_{p_i}, \text{ für } i \neq r,$$

$$\beta) 0 \leq x'_{q_s} = \frac{1}{\bar{a}_{rs}} \bar{b}_r \leq u_{q_s},$$

$$\gamma) x'_{p_r} \in \{0, u_{p_r}\},$$

δ)  $x'_{q_i} \in \{0, u_{q_i}\}$ , für  $i \neq s$ .

Wollen wir nun den Wert der Variablen  $x_{q_s}$  um  $\lambda \geq 0$  ändern (d. h. erhöhen, falls  $\bar{q}_s > 0$ , bzw. um  $\lambda$  erniedrigen, falls  $\bar{q}_s < 0$ ), so können wir das so lange tun bis entweder

- eine Basisvariable den Wert ihrer oberen oder unteren Schranke annimmt oder
- die Nichtbasisvariable  $q_s$  den Wert ihrer anderen Schranke annimmt.

Aus letzterem folgt natürlich, dass  $\lambda \leq u_{q_s}$  gelten muss, ansonsten würden wir bei Erhöhung ( $\bar{q}_s > 0$ ) die obere Schranke bzw. bei Erniedrigung ( $\bar{q}_s < 0$ ) die untere Schranke für  $x_{q_s}$  verletzen. Dies erklärt die Bestimmung von  $\lambda_2$  in (4).

Ist  $\bar{q}_s > 0$ , so bewirkt eine Erhöhung von  $x_{q_s}$  um den Wert  $\lambda$ , dass  $x'_{p_i} = \bar{b}_i - \lambda \bar{a}_{i_s}$  gilt. Aus der Bedingung  $\alpha$ ) erhalten wir daher die folgenden Schranken für  $\lambda$ :

$$\begin{aligned} \lambda &\leq \frac{\bar{b}_i}{\bar{a}_{i_s}} \quad \text{für } \bar{a}_{i_s} > 0 \\ \lambda &\leq \frac{\bar{b}_i - u_{p_i}}{\bar{a}_{i_s}} \quad \text{für } \bar{a}_{i_s} < 0. \end{aligned}$$

Ist  $\bar{q}_s < 0$ , so bewirkt die Verminderung von  $x_{q_s}$  ( $= u_{q_s}$ ) um den Wert  $\lambda \geq 0$ , dass  $x'_{p_i} = \bar{b}_i + \lambda \bar{a}_{i_s}$  gilt. Aus  $\alpha$ ) erhalten wir somit:

$$\begin{aligned} \lambda &\leq -\frac{\bar{b}_i}{\bar{a}_{i_s}} \quad \text{für } \bar{a}_{i_s} < 0 \\ \lambda &\leq \frac{u_{p_i} - \bar{b}_i}{\bar{a}_{i_s}} \quad \text{für } \bar{a}_{i_s} > 0. \end{aligned}$$

Dies begründet die Definition von  $\lambda_0$  und  $\lambda_1$ . Sei nun  $\lambda_{\min} = \min\{\lambda_0, \lambda_1, \lambda_2\}$ . Gilt  $\lambda_{\min} = \lambda_2$ , so wird einfach der Wert einer Nichtbasisvariablen von der gegenwärtigen Schranke auf die entgegengesetzte Schranke umdefiniert. Die Basis ändert sich dadurch nicht, jedoch konnte aufgrund der Auswahl in PRICE eine Zielfunktionsverbesserung erreicht werden. Gilt  $\lambda_{\min} = \lambda_1$  oder  $\lambda_{\min} = \lambda_2$ , so führen wir einen üblichen Pivotschritt durch. Bei  $\lambda_{\min} = \lambda_1$  wird die neue Nichtbasisvariable  $x_{p_r}$  mit Null festgesetzt, da die untere Schranke von  $x_{p_r}$  die stärkste Einschränkung für die Veränderung von  $x_{q_s}$  darstellte. Bei  $\lambda_{\min} = \lambda_2$ , wird analog  $x_{p_r}$  eine Nichtbasisvariable, die den Wert ihrer oberen Schranke annimmt.  $\square$

Wir wollen nun anhand eines Beispiels die Ausführung von Algorithmus (10.16) in Tableautechnik demonstrieren. Dabei werden wir keine neuen Tableau-Update-Formeln definieren, sondern mit den bekannten Formeln für die Pivotschritte rechnen. Wir führen lediglich eine zusätzliche Spalte ein, in der jeweils aus  $\bar{b}$  und  $c_0$  der

Wert der gegenwärtigen  $E$ -Basislösung berechnet wird. Ist also  $c_0$  der gegenwärtige Wert der Basislösung und  $x_B = A_B^{-1}b = \bar{b}$  so hat die zusätzliche Spalte in der ersten Komponente den Eintrag  $\bar{c}_0 := -c_0 - \bar{c}_N - u_N$ , und die restlichen Komponenten berechnen sich durch

$$\bar{b} := \bar{b} - \bar{A}_N - u_N.$$

**(10.17) Beispiel.** Wir betrachten das folgende LP, dessen Lösungsmenge in Abbildung 10.1 dargestellt ist

$$\begin{array}{ll} \max & 2x_1 + x_2 \\ & -x_1 + x_2 \leq \frac{1}{2} \\ & x_1 + x_2 \leq 2 \\ & x_1 \leq \frac{3}{2} (= u_1) \\ & x_2 \leq 1 (= u_2) \\ & x_1, x_2 \geq 0. \end{array}$$

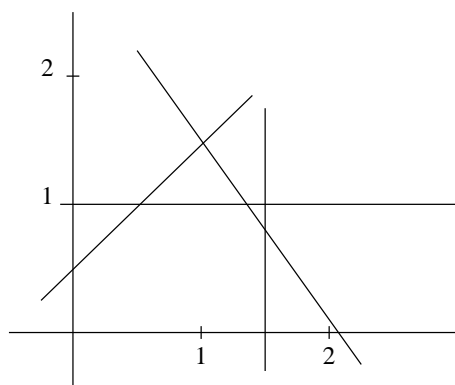


Abb. 10.1

Wir führen 2 Schlupfvariablen  $x_3, x_4$  ein (mit oberen Schranken  $+\infty$ ) und erhalten als Anfangstableau  $T_0$  (mit zusätzlicher Spalte  $\begin{pmatrix} \bar{c}_0 \\ \bar{b} \end{pmatrix}$ , die mit der Spalte  $\begin{pmatrix} -c_0 \\ \bar{b} \end{pmatrix}$  identisch ist, da keine Nichtbasisvariable von Null verschieden ist):

$$T_0 : \begin{array}{cccc|cc} & 1 & 2 & 3 & 4 & & \\ & 2 & 1 & 0 & 0 & 0 & 0 \\ \hline 3 & -1 & 1 & 1 & 0 & \frac{1}{2} & \frac{1}{2} \\ 4 & 1 & 1 & 0 & 1 & 2 & 2 \end{array} = \begin{array}{ccc|c|c} \bar{c} & 0 & & -c_0 & \bar{c}_0 \\ \hline \bar{A} & I & & \bar{b} & \bar{b} \end{array} \quad \begin{array}{l} B = (3, 4) \\ \bar{N} = (1, 2) \end{array}$$

Wir wählen die erste Spalte als Pivotspalte, d. h.  $\bar{q}_1 = 1$ . Schritt (4) CHUZR ergibt

$$\lambda_0 = \frac{\bar{b}_2}{\bar{a}_{21}} = \frac{2}{1} = 2, \quad \lambda_1 \text{ ist wegen } u_3 = \infty \text{ nicht definiert, } \lambda_2 = u_1 = \frac{3}{2},$$

d. h.  $\lambda_{\min} = \lambda_2$ . Wir führen also Schritt (4) (d) aus, setzen  $\bar{q}_1 := -\bar{q}_1$ , d. h.  $\bar{q}_1 = -1$ , und berechnen  $\bar{b}$  neu. Die neue letzte Spalte, d. h. die neue  $E$ -Basislösung zur Basis  $A_B$  erhalten wir aus der (normalen) Basislösung wegen  $N^- = (1)$  wie folgt:

$$\bar{b} = \bar{b} - \bar{A}_{\cdot 1} u_1 = \begin{pmatrix} \frac{1}{2} \\ 2 \end{pmatrix} - \frac{3}{2} \begin{pmatrix} -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ \frac{1}{2} \end{pmatrix},$$

und der Zielfunktionswert erhöht sich um  $\bar{c}_1 \cdot u_1 = 2 \cdot \frac{3}{2} = 3$ . Mithin erhalten wir als nächstes Tableau (mit  $\bar{q}_1 := -\bar{q}_1$ )

$$T_1 : \begin{array}{cccc|cc} & -1 & 2 & 3 & 4 & & \\ & 2 & 1 & 0 & 0 & 0 & -3 \\ \hline 3 & -1 & 1 & 1 & 0 & \frac{1}{2} & 2 \\ 4 & 1 & 1 & 0 & 1 & 2 & \frac{1}{2} \end{array} \quad \begin{array}{l} B = (3, 4) \\ \bar{N} = (-1, 2) \end{array}$$

Als Pivotspalte kommt nur die 2. Spalte ( $\bar{q}_2 = 2$ ) in Frage. CHUZR ergibt

$$\lambda_0 = \frac{\bar{b}_2}{\bar{a}_{22}} = \frac{1}{2}, \lambda_1 \text{ nicht definiert}, \lambda_2 = u_2 = 1.$$

Die Wahl der Austauschzeile in (4) (b) ergibt ein eindeutig bestimmtes  $r = 2$ . Wir führen dann (5) in Tableautechnik aus. Wir führen also einen Standardpivotschritt auf dem Element  $\bar{a}_{22} = 1$  durch und korrigieren anschließend  $\bar{b}$  und  $c_0$ , um  $\bar{b}$  und  $\bar{c}_0$  zu erhalten.

$$T_2 : \begin{array}{cccc|cc} & -1 & 2 & 3 & 4 & & \\ & 1 & 0 & 0 & -1 & -2 & -\frac{7}{2} \\ \hline 3 & -2 & 0 & 0 & -1 & -\frac{3}{2} & \frac{3}{2} \\ 2 & 1 & 1 & 0 & 1 & 2 & \frac{1}{2} \end{array} \quad \begin{array}{l} B = (3, 2) \\ \bar{N} = (-1, 4) \end{array}$$

Dieses Tableau ist optimal, mit  $x_1 = \frac{3}{2}$ ,  $x_2 = \frac{1}{2}$ . Zur Erläuterung der übrigen Fälle beginnen wir noch einmal mit Tableau  $T_0$ , wählen aber die 2. Spalte als Pivotspalte. CHUZR ergibt

$$\lambda_0 = \frac{\bar{b}_1}{\bar{a}_{12}} = \frac{1}{2}, \lambda_1 \text{ nicht definiert}, \lambda_2 = 1.$$

In diesem Falle machen wir einen Standardpivotschritt mit dem Pivotelement

$$\bar{a}_{rs} = \bar{a}_{12} = 1.$$

$$T_3 : \begin{array}{c} \\ \\ \\ \end{array} \begin{array}{cccc|cc} & 1 & 2 & 3 & 4 & & \\ & 3 & 0 & -1 & 0 & -\frac{1}{2} & -\frac{1}{2} \\ 2 & -1 & 1 & 1 & 0 & \frac{1}{2} & \frac{1}{2} \\ 4 & 2 & 0 & -1 & 1 & \frac{3}{2} & \frac{3}{2} \end{array} \quad \begin{array}{l} B = (2, 4) \\ \bar{N} = (1, 3) \end{array}$$

Die gegenwärtige Basislösung wird wie üblich in eine  $E$ -Basislösung transformiert. Als Pivotspalte kommt nur die erste in Frage. CHUZR ergibt

$$\lambda_0 = \frac{\bar{b}_2}{\bar{a}_{21}} = \frac{3}{4}, \quad \lambda_1 = \frac{\bar{b}_1 - u_2}{\bar{a}_{11}} = \frac{1}{2}, \quad \lambda_2 = u_1 = \frac{3}{2} \quad (\lambda_{\min} = \lambda_1)$$

Wir führen nun Schritt (5) des Algorithmus (10.6) durch einen Pivotschritt auf  $\bar{a}_{rs} = \bar{a}_{11} = -1$  aus und berechnen das gesamte Tableau neu. Anschließend müssen wir die rechte Spalte des Tableaus korrigieren. Da  $\lambda_{\min} = \lambda_1$  wird die neue Nichtbasisvariable  $x_2$  mit ihrer oberen Schranke  $u_2 = 1$  belegt (und  $\bar{q}_1 = -2$  gesetzt). Wir müssen daher vom Neuberechneten  $\bar{b} = (-\frac{1}{2}, \frac{5}{2})^T$  noch das  $u_2$ -fache der zum Index 2 gehörigen Spalte von  $\bar{A}$  also  $\bar{A}_{\cdot 1} = (-1, 2)^T$  subtrahieren. Vom Zielfunktionswert  $c_0$  subtrahieren wir  $u_2 \cdot \bar{c}_2 = 1 \cdot 3$ . Daraus ergibt sich

$$T_4 : \begin{array}{c} \\ \\ \\ \end{array} \begin{array}{cccc|cc} & 1 & -2 & 3 & 4 & & \\ & 0 & 3 & 2 & 0 & 1 & -2 \\ 1 & 1 & -1 & -1 & 0 & -\frac{1}{2} & \frac{1}{2} \\ 4 & 0 & 2 & 1 & 1 & \frac{5}{2} & \frac{1}{2} \end{array} \quad \begin{array}{l} B = (1, 4) \\ \bar{N} = (-2, 3) \end{array}$$

$$\bar{b} = \begin{pmatrix} -\frac{1}{2} \\ \frac{5}{2} \end{pmatrix} - 1 \begin{pmatrix} -1 \\ 2 \end{pmatrix}.$$

Als Pivotspalte müssen wir nun die 3. Spalte von  $T_4$ , also  $\bar{A}_{\cdot 2}$ , wählen.

$$\lambda_0 = \frac{\bar{b}_2}{\bar{a}_{22}} = \frac{1}{2}, \quad \lambda_1 = \frac{\bar{b}_1 - u_1}{\bar{a}_{12}} = 1, \quad \lambda_2 = u_1 = \frac{3}{2}.$$

Somit führen wir einen Standardpivotschritt auf  $\bar{a}_{22} = 1$  aus

$$T_5 : \begin{array}{c} \\ \\ \\ \end{array} \begin{array}{cccc|cc} & 1 & -2 & 3 & 4 & & \\ & 0 & -1 & 0 & -2 & -4 & -3 \\ 1 & 1 & 1 & 0 & 1 & 2 & 1 \\ 3 & 0 & 2 & 1 & 1 & \frac{5}{2} & \frac{1}{2} \end{array} \quad \begin{array}{l} B = (1, 3) \\ \bar{N} = (-2, 4) \end{array}$$

Als Pivotspalte kommt nur  $\bar{A}_{\cdot 1}$ , also die 2. Spalte von  $T_5$ , in Frage, d. h.  $q_s = q_1 = 2$ ,  $\bar{q}_1 = -2$ . Wir erhalten durch CHUZR

$$\lambda_0 \text{ undefiniert, } \lambda_1 = \frac{\bar{b}_1 - u_1}{\sigma \bar{a}_{11}} = \frac{1}{2}, \lambda_2 = u_2 = 1.$$

Unser Pivotelement ist somit  $\bar{a}_{11} = 1$ . Wie üblich führen wir einen Standardpivotschritt durch. Da wir die neue Nichtbasisvariable  $x_1$  mit ihrer oberen Schranke  $u_1 = \frac{3}{2}$  belegen, ziehen wir wieder vom Neuberechneten  $\bar{b}$  das  $u_1$ -fache der zu 1 gehörigen Spalte von  $\bar{A}$  (das ist  $\bar{A}_{\cdot 1} = (1, -2)^T$ ) ab, um  $\bar{b}$  zu erhalten.

$$T_6 : \quad \begin{array}{cc|cc|c|c} & -1 & 2 & 3 & 4 & & \\ & 1 & 0 & 0 & -1 & -2 & -\frac{7}{2} \\ \hline & 2 & 1 & 1 & 0 & 1 & 2 & \frac{1}{2} \\ & 3 & -2 & 0 & 1 & -1 & -\frac{3}{2} & \frac{3}{2} \end{array} \quad \begin{array}{l} B = (2, 3) \\ \bar{N} = (-1, 4) \end{array}$$

Wir haben wieder die Optimallösung erreicht. □

Analog zu oberen Schranken können natürlich auch untere Schranken behandelt werden. Ferner kann man auch allgemeine obere Schranken (generalized upper bounds, GUB), dies sind Restriktionen der folgenden Form:  $\sum x_i \leq a$ , auf ähnliche Weise berücksichtigen. Bei solchen Restriktionen können ebenfalls Pivotschemata entworfen werden, die wesentlich weniger Speicher- und Rechenaufwand erfordern als die Standardmethode. Weiterhin gibt es effiziente Techniken zur Behandlung von sogenannten variable upper bounds (VUB) der Form  $0 \leq x \leq y$ .

## 10.4 Das duale Simplexverfahren

Wir wollen nun eine Variante des Simplexverfahrens darstellen, die — wie wir noch sehen werden — bei bestimmten Problemstellungen von Vorteil ist. Sei  $A_B$  eine Basis von  $A$ , und betrachte das Paar dualer linearer Programme:

$$(P) \quad \begin{array}{l} \max c^T x \\ Ax = b \\ x \geq 0 \end{array} \quad (D) \quad \begin{array}{l} \min u^T b \\ u^T A \geq c^T \end{array}$$

bzw.

$$\begin{array}{l} \max c^T x \\ x_B = A_B^{-1}b - A_B^{-1}A_N x_N \\ x_N, x_B \geq 0 \end{array} \quad \begin{array}{l} \min u^T b \\ u^T A_B \geq c_B^T \\ u^T A_N \geq c_N^T \end{array}$$

**(10.18) Definition.** Die Basis  $A_B$  von  $A$  heißt **primal zulässig**, falls  $A_B^{-1}b \geq 0$  und **dual zulässig**, falls  $\bar{c}^T = c_N^T - c_B^T A_B^{-1} A_N \leq 0$ . Die zugehörigen Basislösungen  $x$  bzw.  $u$  mit  $x_B = A_B^{-1}b$ ,  $x_N = 0$  bzw.  $u^T = c_B^T A_B^{-1}$  heißen dann **primal** bzw. **dual zulässig**.

**(10.19) Satz.** Sei  $P = \{u \in \mathbb{K}^m \mid u^T A \geq c^T\}$ . Der Vektor  $u$  ist genau dann eine Ecke von  $P$ , wenn  $u$  eine dual zulässige Basislösung ist.

**Beweis :** Sei  $\bar{u}$  Ecke von  $P = P(-A^T, -c)$  und  $I = \text{eq}(\{\bar{u}\})$ . Mit Satz (8.9) folgt  $\text{rang}((-A^T)_I) = m$ , d. h. es existiert  $B \subseteq I$  mit  $A_B$  Basis von  $A$ , und es gilt  $\bar{u}^T A_B = c_B^T$ ,  $\bar{u}^T A_N \geq c_N^T$ . Also ist  $\bar{u}^T = c_B^T A_B^{-1}$  und  $c_B^T A_B^{-1} A_N \geq c_N^T$ , d. h.  $A_B$  ist dual zulässig. Ist umgekehrt  $A_B$  dual zulässig, so ist  $\bar{u} := c_B^T A_B^{-1}$  aufgrund von Satz (8.9) eine Ecke von  $P$ .  $\square$

**(10.20) Bemerkung.** Ist  $A_B$  eine Basis von  $A$ , und sind  $x$  bzw.  $u$  die zu  $A_B$  gehörenden (primalen bzw. dualen aber nicht notwendigerweise zulässigen) Basislösungen, so gilt:  $c^T x = u^T b$ .

**Beweis :**  $u^T b = c_B^T A_B^{-1} b = c_B^T x_B = c^T x$ .  $\square$

**(10.21) Folgerung.** Ist  $A_B$  eine Basis von  $A$ , so ist  $A_B$  optimal genau dann, wenn  $A_B$  primal und dual zulässig ist.

**Beweis :** (Dualitätssatz).

**(10.22) Folgerung.** Ist  $A_B$  eine optimale Basis für das Programm (P), dann ist  $c_B^T A_B^{-1}$  eine optimale Lösung des zu (P) dualen Programms (D).

$\square$

Der Vektor  $\pi := c_B^T A_B^{-1}$  (mit  $A_B$  dual zulässig) heißt der Vektor der **Schattenpreise** (vgl. ökonomische Interpretation der Dualität und Schritt (1) in (10.1)).

Wir wollen nun die dem dualen Simplexverfahren zugrunde liegende Idee entwickeln und bemerken, dass — im Prinzip — das duale Simplexverfahren das primale Simplexverfahren angewendet auf das duale Problem ist. Sei  $A_B$  eine Basis von  $A$ , so lässt sich (9.1) umformen in:

$$(10.23) \quad \begin{aligned} \max \quad & c_B^T A_B^{-1} b + \bar{c}^T x_N \\ & \bar{A} x_N + I x_B = A_B^{-1} b (= \bar{b}) \\ & x_B, x_N \geq 0 \end{aligned}$$

oder

$$(10.24) \quad \begin{aligned} c_B^T A_B^{-1} b + \max \bar{c}^T x_N \\ \bar{A} x_N &\leq \bar{b} \\ x_N &\geq 0 \end{aligned}$$

Das zu (10.24) duale Programm lautet (bei Weglassung des konstanten Terms  $c_B^T A_B^{-1} b$  der Zielfunktion):

$$(10.25) \quad \begin{aligned} \min u^T \bar{b} \\ \bar{A}^T u &\geq \bar{c} \\ u &\geq 0 \end{aligned}$$

Die Einführung von Schlupfvariablen  $y_N$  (und Variablenumbenennung  $y_B := u$ ) ergibt:

$$(10.26) \quad \begin{aligned} - \max -\bar{b}^T y_B \\ -\bar{A}^T y_B + I y_N &= -\bar{c} \\ y &\geq 0 \end{aligned}$$

Durch eine lange Kette von Umformungen ist es uns also gelungen, ein LP in Standardform (9.1) in ein anderes LP in Standardform (10.26) zu transformieren. Was haben wir gewonnen? Nicht viel, es sei denn, die Basis  $A_B$  ist **dual zulässig**. Denn in diesem Falle gilt, dass die rechte Seite von (10.26), also  $-\bar{c}$ , nichtnegativ ist. Also ist die Matrix  $I$  eine zulässige (Ausgangs-)basis für (10.26), und wir können auf (10.26) den Simplexalgorithmus (direkt mit Phase II beginnend) anwenden.

Eine besonders interessante Anwendung ergibt sich, wenn das ursprüngliche Problem die Form

$$\begin{aligned} \max c^T x \\ Ax &\leq b \\ x &\geq 0 \end{aligned}$$

hat, wobei  $c \leq 0$  und  $b \not\geq 0$  ist. Hier ist eine dual zulässige Basis direkt gegeben, während eine primal zulässige Basis erst mittels Phase I gefunden werden müsste.

### (10.27) Algorithmus (duale Simplexmethode).

**Input:**  $A' \in \mathbb{K}^{(m,n)}$ ,  $b' \in \mathbb{K}^m$ ,  $c \in \mathbb{K}^n$ .

**Output:** optimale Lösung  $x$  des LP  $\max\{c^T x \mid A'x = b', x \geq 0\}$ .

**Phase I:**

*Bestimmung eines Subsystems  $Ax = b, x \geq 0$  mit  $P^=(A, b) = P^=(A', b')$ , welches (9.2) erfüllt (falls möglich) und Bestimmung einer dual zulässigen Basis  $A_B$  von  $A$ . Berechne:  $\bar{A} = A_B^{-1}A_N, \bar{b} = A_B^{-1}b, \bar{c}^T = c_N^T - c_B^T A_B^{-1}A_N$ . (Dieser Teil des Algorithmus wird analog zur Phase I (9.24) der Grundversion des Simplexalgorithmus ausgeführt.)*

## Phase II: Optimierung

### (II.1) (Optimalitätsprüfung)

*Gilt  $\bar{b}_i \geq 0$  ( $i = 1, \dots, m$ ), so ist die gegenwärtige Basislösung optimal ( $A_B$  ist primal und dual zulässig). Setze  $x_B = \bar{b}$  und  $x_N = 0$ , andernfalls gehe zu (II.2).*

### (II.2) (Bestimmung der Pivotzeile)

*Wähle einen Index  $r$  mit  $\bar{b}_r < 0$ .*

### (II.3) (Prüfung auf Beschränktheit des Optimums)

*Gilt  $\bar{A}_r \geq 0$ , so hat das duale Programm keine endliche Lösung, also ist  $P^=(A, b) = \emptyset$ . STOP!*

### (II.4) Berechne $\lambda_0 := \min \left\{ \frac{\bar{c}_j}{\bar{a}_{rj}} \mid \bar{a}_{rj} < 0, j = 1, \dots, n \right\}$ .

### (II.5) (Bestimmung der Pivotspalte)

*Wähle Index  $s \in \{j \in \{1, \dots, n - m\} \mid \frac{\bar{c}_j}{\bar{a}_{rj}} = \lambda_0\}$ .*

### (II.6) (Pivotoperation)

*Setze  $A_{B'}^{-1} := EA_B^{-1}$  mit  $E$  aus Satz (9.10), und berechne alle notwendigen Parameter neu. Gehe zu (II.1).*

□

Der duale Simplexalgorithmus wurde 1954 von Lemke entwickelt. Wie wir gerade gezeigt haben, ist er (nach Transformation) jedoch nichts anderes als die Anwendung des primalen Simplexalgorithmus auf das duale Problem. Aus diesem Grunde hat sich lange Zeit niemand die Mühe gemacht, eine „echte“ Implementation des dualen Verfahrens vorzunehmen. Als in den 90er Jahren die Anwendung der ganzzahligen Optimierung immer wichtiger und Schnittebenenverfahren (die wir später erklären werden) als die zentralen Methoden zur Lösung ganzzahliger Optimierungsprobleme erkannt wurden, sind erstmals duale Methoden programmiert worden. Zur Überraschung vieler erwiesen sich diese dualen Codes als (in

der Praxis) schneller als die primalen, so dass sie heute bei praktischen Rechnungen dominieren, siehe Bixby (2002) für den experimentellen Nachweis und die heuristische Begründung.

**(10.28) Tableauform des dualen Simplexalgorithmus.** Wollen wir das duale Simplexverfahren auf das Problem (10.23) (mit dual zulässiger Basis  $A_B$ ) anwenden, so können wir das verkürzte Tableau

$$VT := \begin{array}{|c|c|} \hline -\bar{b}^T & -z_0 \\ \hline -\bar{A}^T & -\bar{c} \\ \hline \end{array}$$

verwenden und auf dieses Tableau die Update Formeln des verkürzten Simplextableaus (diese sind in Schritt (II.7) von (9.15) angegeben) anwenden. Jedoch zeigt eine einfache Überlegung, dass wir bei der Auswahl der Indizes  $r$  und  $s$  entsprechend (II.2) bzw. (II.5) die Updateformeln (II.7) aus (9.15) auch direkt auf die Matrix  $\bar{A}$  bzw.  $\bar{b}$  und  $\bar{c}$  anwenden können und zwar genau in der Form wie sie in (9.15) (II.7) aufgeschrieben wurden (um das neue Tableau  $VT'$  zu erhalten).

$$VT' := \begin{array}{|c|c|} \hline -(\bar{b}')^T & -z'_0 \\ \hline -(\bar{A}')^T & -\bar{c}' \\ \hline \end{array}$$

Wir führen also die Transponierung bzw. den Vorzeichenwechsel nicht explizit durch, sondern beachten den Vorzeichenwechsel bzw. die Transponierung implizit bei der Bestimmung der Pivotzeile bzw. Pivotspalte.

Dann könnten wir aber auch wieder das primale (verkürzte) Tableau verwenden und anhand dieses Tableaus sowohl den primalen als auch den dualen Simplexalgorithmus durchführen.  $\square$

**(10.29) Das Tucker-Tableau.** Für die nachfolgende spezielle Form linearer Programme erhalten wir primale und duale Optimallösungen direkt aus der Tableau-rechnung. Wir betrachten:

$$\begin{array}{ll} \max c^T x & \min u^T b \\ (P) \quad Ax \leq b & (D) \quad u^T A \geq c^T \\ x \geq 0 & u \geq 0 \end{array}$$

oder mit Schlupfvariablen

$$\begin{array}{rcl} \max z_0 & & \max z_0 \\ c^T x - z_0 = 0 & & u^T b + z_0 = 0 \\ Ax - b = -u & & -c^T + u^T A = x^T \\ x, u \geq 0 & & u, x \geq 0 \end{array}$$

oder in Tableauschreibweise

$c_1, c_2, \dots, c_n$	$z_0$
$A$	$b_1$
	$\vdots$
	$b_m$

- $N$  = Nichtbasisvariable des primalen Programms  
 = Basisvariable des dualen Programms  
 $B$  = Basisvariable des primalen Programms  
 = Nichtbasisvariable des dualen Programms

Ein solches Tucker-Tableau heißt **primal zulässig**, wenn  $b \geq 0$ , und **dual zulässig**, wenn  $c \leq 0$ .

Führen wir den primalen oder dualen Simplexalgorithmus bis zum Optimaltableau durch, so sind die jeweiligen Optimallösungen gegeben durch:

$$\left. \begin{array}{l} x_{B(i)} = \bar{b}_i \quad i = 1, \dots, m \\ x_{N(i)} = 0 \quad i = 1, \dots, n \end{array} \right\} \text{optimale Lösung des primalen Programms.}$$

$$\left. \begin{array}{l} u_{N(i)} = -\bar{c}_i \quad i = 1, \dots, n \\ u_{B(i)} = 0 \quad i = 1, \dots, m \end{array} \right\} \text{optimale Lösung des dualen Programms.}$$

□

Hat man eine Optimallösung von einem LP in Standardform ausgehend berechnet, so lässt sich die Optimallösung des dualen Programms nicht ohne Weiteres aus dem Tableau ablesen.

**(10.30) Beispiel (dualer Simplexalgorithmus).** Wir betrachten die folgenden zueinander dualen Programme (P) und (D). Die Lösungsmenge von (P) sei mit P bezeichnet, die von (D) mit D. P und D sind in Abbildung 10.2 dargestellt.

$$\begin{array}{rcl} \max -x_1 - 2x_2 & & \min -2y_3 - 3y_4 \\ (P) \quad \begin{array}{l} -x_1 - x_2 \leq -2 \quad (x_3) \\ -2x_1 - x_2 \leq -3 \quad (x_4) \\ x_1, x_2 \geq 0 \end{array} & (D) & \begin{array}{l} -y_3 - 2y_4 \geq -1 \quad (-y_1) \\ -y_3 - y_4 \geq -2 \quad (-y_2) \\ y_3, y_4 \geq 0 \end{array} \end{array}$$

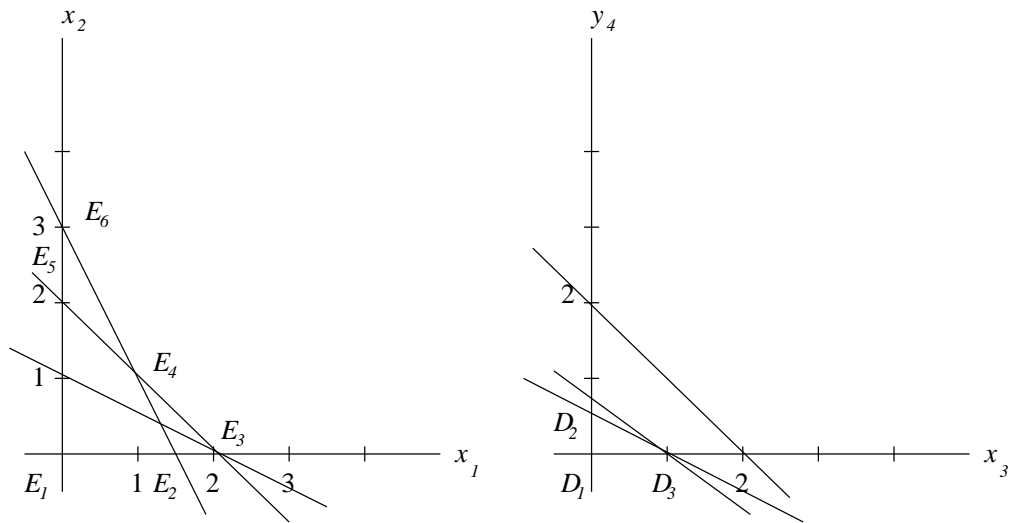


Abb. 10.2

Wir schreiben das Tucker-Tableau zu diesem Problem auf.

	1	2				
	-1	-2		0		
3	-1	-1		-2	$x_1 = 0$	$y_3 = 0$
4	-2	-1		-3	$x_2 = 0$	$y_4 = 0$
					$x_3 = -2$	$y_1 = 1$
					$x_4 = -3$	$y_2 = 2$

$N = (1, 2) =$  primale Nichtbasis, duale Basis  
 $B = (3, 4) =$  primale Basis, duale Nichtbasis

Die gegenwärtige Basis ist dual aber nicht primal zulässig. Wir machen einen Update mit dem Pivotelement  $\bar{a}_{rs} = \bar{q}_{21} = -2$ . Das neue Tableau hat folgende Form.

	4	2			
	$-\frac{1}{2}$	$-\frac{3}{2}$		$-\frac{3}{2}$	$N = (4, 2), B = (3, 1)$
3	$-\frac{1}{2}$	$-\frac{1}{2}$		$-\frac{1}{2}$	
1	$-\frac{1}{2}$	$\frac{1}{2}$		$\frac{3}{2}$	

$x_1 = \frac{3}{2}$	}	$= E_2$	$y_3 = 0$	}	$= D_2$
$x_2 = 0$			$y_4 = \frac{1}{2}$		

Es folgt ein Pivotschritt mit  $\bar{a}_{rs} = \bar{a}_{11} = -\frac{1}{2}$ .

$$\begin{array}{cc|c}
 & 3 & 2 \\
 & -1 & -1 & -2 \\
 4 & -2 & 1 & 1 \\
 1 & -1 & 1 & 2 \\
 \hline
 x_1 & = & 2 & \left. \vphantom{x_1} \right\} = E_3 \\
 x_2 & = & 0 & 
 \end{array}
 \quad
 N = (3, 2), \quad B = (4, 1)$$

$$\begin{array}{l}
 y_3 = 1 \\
 y_4 = 0
 \end{array}
 \left. \vphantom{y_3} \right\} = D_3$$

Dieses Tableau ist optimal. □

## 10.5 Zur Numerik des Simplexverfahrens

Dieser Abschnitt ist nur äußerst kursorisch und oberflächlich ausgearbeitet. Eine sorgfältige Behandlung des Themas würde eine gesamte Spezialvorlesung erfordern. Wir empfehlen dem Leser die Bemerkungen zur Numerik in

V. Chvátal, “Linear Programming”, Freeman, New York, 1983, (80 SK 870 C 564)

oder das Buch

M. Bastian, “Lineare Optimierung großer Systeme”, Athenäum/Hain/Skriptor/Hanstein, Königstein, 1980, (80 ST 130 B 326)

zu lesen, das einem Spezialaspekt dieses Themas gewidmet ist und diesen einigermaßen erschöpfend behandelt. Das Buch

B. A. Murtagh, “Advanced Linear Programming: Computation and Practice”, McGraw-Hill, New York, 1981, (80 QH 400 M 984).

ist ganz Rechen- und Implementierungstechniken der linearen Optimierung gewidmet. Generelle Methoden zur Behandlung spezieller (z. B. dünn besetzter oder triangulierter oder symmetrischer) Matrizen im Rahmen numerischer Verfahren (etwa Gauß-Elimination, Matrixinvertierung) finden sich in

S. Pissanetsky, “Sparse Matrix Technology”, Academic Press, London, 1984, (80 SK 220 P 673).

Generell wird bei Implementationen die revidierte Simplexmethode verwendet. Es gibt zwei grundsätzliche Varianten für die Reinverson: **Produktform der Inversen (PFI)**, **Eliminationsform der Inversen (EFI)**.

**(10.31) Produktform der Inversen.**  $B^{-1}$  liegt in Produktform vor:  
 $B^{-1} = E_k \cdot E_{k-1} \cdot \dots \cdot E_1$  mit  $E_i$  aus Satz (9.10).

**Prinzip:** Gauß-Jordan Verfahren.

*Freiheitsgrade:*

- a) Positionen der Pivotelemente,
- b) Reihenfolge der Pivots.
  - a) hat Einfluss auf die numerische Stabilität
  - b) hat Einfluss auf die Anzahl NNE (nicht Null Elemente) im Etafile (Speicherung der Spalten  $\eta$  von  $E_i$ ).

**Beispiel.**

$$B = \begin{pmatrix} 1 & 1 & 1 & 1 \\ -1 & 1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ -1 & 0 & 0 & 1 \end{pmatrix}$$

besitzt folgende Produktform-Darstellungen der Inversen:

- a) Pivolisieren auf der Hauptdiagonalen von “links oben” nach “rechts unten” ergibt

$$B^{-1} = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{pmatrix} \begin{pmatrix} 1 & -\frac{1}{3} & & \\ & 1 & -\frac{1}{3} & \\ & & \frac{2}{3} & \\ & & -\frac{1}{3} & 1 \end{pmatrix} \begin{pmatrix} 1 & -\frac{1}{2} & & \\ & \frac{1}{2} & & \\ & & 1 & \\ & & & 1 \end{pmatrix} \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{pmatrix}$$

Zu speichern sind 16 Werte (+ Positionen).

- b) Pivolisieren auf der Hauptdiagonalen von “rechts unten” nach “links oben” ergibt:

$$B^{-1} = \begin{pmatrix} \frac{1}{4} & & & \\ \frac{1}{4} & 1 & & \\ \frac{1}{4} & & 1 & \\ \frac{1}{4} & & & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 & & \\ & 1 & & \\ & & 0 & 1 \\ & & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 & & \\ & 1 & 0 & \\ & & 1 & \\ & & & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & & -1 \\ & 1 & 0 \\ & & 1 & 0 \\ & & & 1 \end{pmatrix}$$

Zu speichern sind 10 Werte (+ Positionen)

Folgende Gegebenheiten sollten bei der Pivotauswahl berücksichtigt werden:

**(10.32)** Es ist günstig, in einer Spalte mit wenigen NNE zu pivotisieren, da dann die Etavektoren dünn besetzt sind.

**(10.33)** Es ist günstig, in einer Zeile mit wenigen NNE zu pivotisieren, da dann in anderen Zeilen der umgeformten Restmatrix potentiell weniger neue NNE entstehen.

**(10.34)** Es ist aus Gründen der Stabilität nicht günstig, ein Pivotelement zu wählen, dessen Betrag sehr klein ist.

In den Inversionsroutinen, die bei den ersten Implementierungen des Simplexverfahrens benutzt wurden, beachtete man nur die numerische Stabilität, nahm sich die Spalten der Basismatrix in der Reihenfolge, in der sie abgespeichert waren, multiplizierte sie mit den bereits bestimmten Elementarmatrizen und pivotisierte auf der Komponente mit dem größten Absolutbetrag. Später nahm man eine Vorsortierung der Basisspalten in der Weise vor, so dass die dünn-besetzten zuerst bearbeitet wurden: dies war ohne großen Aufwand möglich, da aufgrund der spaltenweisen Speicherung der NNE die "column counts" (Anzahl NNE einer bestimmten Spalte) bekannt waren.

Heutzutage wird die Inversion üblicherweise in 2 Phasen zerlegt:

**(10.35) Boolesche Phase.** Entscheidung über Pivotposition!

**(10.36) Numerische Phase.** Rechentechnisch günstige Ausführung des Pivot-schrittes!

Die folgende Beobachtung über Dreiecksmatrizen ist für Implementationen nützlich.

**(10.37) Bemerkung.** Sei  $B$  eine untere  $(m, m)$ -Dreiecksmatrix mit  $b_{ii} \neq 0 \ i = 1, \dots, m$ . Dann ist durch

$$B^{-1} = E_m \cdot E_{m-1} \cdot \dots \cdot E_2 \cdot E_1$$





der NNE, und in der anderen werden entsprechend zeilenweise die Spaltenindices gespeichert. Darüber hinaus werden noch column und row counts (Anzahl der NNE in einer Spalte bzw. Zeile) festgehalten. Auf diese Weise wird ein rascher Zugriff auf Spalten oder Zeilen nach dem Kriterium der Besetzungsdichte ermöglicht.

Die Bestimmung der Pivotpositionen erfolgt in drei Schritten, wobei die beiden ersten Schritte der Identifizierung des “Bumps” in den oben besprochenen Triangularisierungsverfahren entsprechen:

1. Wähle Spalten mit row count 1; dadurch sind die zugehörigen Pivotzeilen ebenfalls eindeutig festgelegt.

Die Spalten- und Zeilenindices der NNE dieser Spalten werden aus den beiden Listen entfernt und die row und column counts angepasst. Dadurch entstehen möglicherweise wieder Zeilen mit row count 1, die anschließend ausgewählt werden.

2. Entsprechend wie im Schritt 1 werden dann Spalten mit column count 1 ausgewählt, Pivotspalten und -zeilen aus den Indices-Listen gestrichen und row sowie column counts angepasst.

Das Ergebnis der beiden ersten Schritte ist die Umordnung von Spalten der Ausgangsmatrix so, dass die vorbestimmten Pivotelemente auf der Hauptdiagonale liegen. Durch symmetrisches Vertauschen der Zeilen und Spalten gemäß der vorgesehenen Pivotfolge erhält man eine Matrix  $B$  der folgenden Gestalt:

$$B = \begin{array}{|c|c|c|} \hline & & \\ \hline & u^1 & u^3 \\ \hline L^1 & & N \\ \hline \end{array}$$

Der Nukleus  $N$  entspricht dem Bump beim Verfahren von Hellerman und Rarick und in der Tat könnte dies Verfahren auch zur Zerlegung des Nukleus benutzt werden.

3. Stattdessen beruht der Schritt 3 der Reinversionsroutine von MPSX/370 auf einer LU-Zerlegung des Nukleus in der Booleschen Matrix, wobei die folgenden einfachen Auswahlregeln angewandt werden:

- wähle unter den Spalten mit minimalem column count eine solche als Pivotspalte, die ein NNE in einer Zeile mit möglichst kleinem row count besitzt;
- wähle als Pivotzeile eine Zeile mit minimalem row count unter den Zeilen mit NNE in der Pivotspalte.

Ist eine Pivotposition bestimmt worden, so werden die beiden Indices-Listen aktualisiert; die Anpassung der row und column counts erfolgt entsprechend.

Durch besondere Maßnahmen wird versucht, das Durchsuchen von Listen sowie die Suche nach NNE in Vektoren einzuschränken (eine detaillierte Beschreibung der eingesetzten Methoden findet man bei Gustavson). So wird die Pivotwahl z. B. dadurch vereinfacht, dass zu jeder Spalte  $j$ , in der noch nicht pivotisiert wurde, nicht nur der column count sondern auch der minimale row count unter den Zeilen mit NNE in Spalte  $j$  mitgeführt wird. Spalten mit gleichem column count und minimalem row count werden verkettet. Zusätzlich werden die Listen der Zeilen- und Spaltenindices von NNE zu der vorbestimmten Pivotposition für die numerische Phase aufgehoben:

- die Liste der Spaltenindices in der Pivotzeile gibt die Spalten an, die aktualisiert werden müssen (durch Umspeicherung erhält man zu jeder Spalte die Adressen derjenigen Etavektoren, die zur Aktualisierung dieser Spalte herangezogen werden müssen);
- die Liste der Zeilenindices in der Pivotspalte entspricht der Besetzung mit NNE in der aktualisierten Pivotspalte, und ihre Speicherung verhindert, dass die gesamte Spalte nach NNE durchsucht werden muss.

Werden die Listen im Verlaufe zu umfangreich und damit auch die Aktualisierung zu aufwendig, so wird die Boolesche Matrix explizit als Matrix umgespeichert (ein Byte pro Element). Ist diese Matrix vollbesetzt oder erreicht sie eine vom Benutzer zu spezifizierende Dichte, so wird die Boolesche Phase abgebrochen.

#### (b) Die Numerische Phase

In der Numerischen Phase werden die Pivotspalten in der Reihenfolge, die in der Booleschen Phase vorbestimmt wurde, aktualisiert und die Etavektoren gebildet. Die Etavektoren von  $L^{-1}$  und  $U^{-1}$  werden dabei auf getrennten Files, dem  $L$ -File und dem  $U$ -File, abgespeichert. Die Etavektoren

zu  $L^1$  und zu  $U^1$  ergeben sich direkt aus den Ausgangsdaten und werden sofort extern gespeichert. Bei der anschließenden Zerlegung der Nukleusspalten hält man die neu entstehenden Etaspalten des  $L$ -Files zunächst im Hauptspeicher, sofern genug Speicherplatz vorhanden ist.

Sei  $d$  diejenige Basisspalte, aus der die Elementarmatrizen  $L_k$  und  $U_k$  berechnet werden sollen. Ausgangspunkt ist die Spalte

$$\hat{d} := L_{k-1} \dots L_2 \cdot L_1 \cdot d,$$

die man sich in einem Arbeitsbereich erzeugt (man beachte, dass die Etavektoren von  $L_1, \dots, L_j$  nicht benötigt werden, wenn  $L^1$  aus  $j$  Spalten besteht).

Diejenigen Komponenten von  $\hat{d}$ , die in Zeilen liegen, in denen bereits pivotisiert wurde, liefern nun den Etavektor von  $U_k$ ; aus den anderen ergibt sich der Etavektor von  $L_k$ .

Wegen der in der Booleschen Phase geleisteten Vorarbeiten sind hierbei nur sehr wenige Suchvorgänge erforderlich:

- die Pivotzeile ist bereits bekannt;
- diejenigen vorausgegangenen Pivots, die für die Berechnung von  $\hat{d}$  relevant sind, sind bekannt; das  $L$ -File muss also nicht durchsucht werden, sondern auf die benötigten  $L_j$  kann direkt zugegriffen werden;
- die Indices der NNE in der aktualisierten Spalte sind von vornherein bekannt; dies erleichtert nicht nur die Abspeicherung der Etavektoren sondern auch das Löschen des Arbeitsbereichs für die nachfolgenden Operationen.

Ist ein vorbestimmtes Pivotelement dem Betrag nach zu klein bzw. durch Differenzbildung tatsächlich gleich Null, so wird die betreffende Spalte zunächst zurückgestellt bis alle anderen vorbestimmten Pivots durchgeführt worden sind. Die Verarbeitung dieser Spalten ist wesentlich aufwendiger, da keine Informationen aus der Booleschen Phase verwendet werden können. Die Pivot-Wahl in den zurückgestellten Spalten sowie in denjenigen Spalten, für die in der Booleschen Phase kein Pivotelement bestimmt wurde, erfolgt nach dem Gesichtspunkt der numerischen Stabilität: man entscheidet sich für die Komponente mit dem größten Absolutbetrag in einer Zeile, in der noch nicht pivotisiert worden ist (Pivoting for Size).

## 10.6 Spezialliteratur zum Simplexalgorithmus

Die folgende Liste ist eine (unvollständige) Sammlung einiger wichtiger Paper zu dem in den Kapiteln 9 und 10 behandelten Themen.

D. Avis & V. Chvátal [1978]: *Notes on Bland's Pivoting Rule*, *Mathematical Programming Studies* 8 (1978) 24–34.

R. H. Bartels [1971]: *A Stabilization of the Simplex Method*, *Numerische Mathematik* 16 (1971) 414–434.

R. H. Bartels & G. H. Golub [1969]: *The Simplex Method of Linear Programming Using LU Decomposition*, *Communications of the Association for Computing Machinery* 12 (1969) 266–268. 275–278.

E. M. L. Beale & J. A. Tomlin [1970] *Special Facilities in a General Mathematical Programming System for Non-Convex Problems Using Sets of Variables*, in: J. Lawrence (ed.), “Proceedings of the fifth IFORS Conference”. Tavistock, London 1970, 447–454.

M. Bénichou, J. M. Gauthier, P. Girodet & G. Hentgès. G. Ribière & O. Vincent [1971]: *Experiments in Mixed Integer Linear Programming*, *Mathematical Programming* 1 (1971) 76–94.

M. Bénichou, J. M. Gauthier, G. Hentgès & G. Ribière [1977]: *The Efficient Solution of Large-Scale Linear Programming Problems — Some Algorithmic Techniques and Computational Results*, *Mathematical Programming* 13 (1977) 280–322.

R. G. Bland [1977]: *New Finite Pivoting Rules for the Simplex Method*, *Mathematics of Operations Research* 2 (1977) 103–107.

R. K. Brayton, F. G. Gustavson & R. A. Willoughby [1970]: *Some Results on Sparse Matrices*, *Mathematics of Computation* 24 (1970) 937–954.

H. Crowder & J. M. Hattingh [1975]: *Partially Normalized Pivot Selection in Linear Programming*, *Mathematical Programming Study* 4 (1975) 12–25.

A. R. Curtis & J. K. Reid [1972]: *On the automatic Scaling of Matrices for Gaussian Elimination*, *Journal of the Institute of Mathematics and its Applications* 10 (1972) 118–124.

G. B. Dantzig, R. P. Harvey, R. D. McKnight & S. S. Smith [1969]: *Sparse Matrix Techniques in Two Mathematical Programming Codes*, in: R. A. Willoughby (ed.), “Sparse Matrix Proceedings” RA-1, IBM Research Center, Yorktown Heights, New York, 1969, 85–99.

R. Fletcher & S. P. J. Matthews [1984]: *Stable Modification of Explicit LU Factors for Simplex Updates*, Mathematical Programming 30 (1984) 267–284.

J. J. H. Forrest & J. A. Tomlin [1972]: *Updated Triangular Factors of the Basis to Maintain Sparsity in the Product Form Simplex Method*, Mathematical Programming 2 (1972) 263–278.

A. M. Geoffrion & R. E. Marsten [1972]: *Integer Programming Algorithms: a Framework and State-of-the-Art Survey*, Management Science 18 (1972) 465–491.

D. Goldfarb & J. K. Reid [1977]: *A Practicable Steepest Edge Simplex Algorithm* Mathematical Programming 12 (1977) 361–371.

D. Goldfarb [1976]: *Using the Steepest Edge Simplex Algorithm to Solve Sparse Linear Programs*, in: J. Bunch and D. Rose (eds.), “Sparse Matrix Computations”, Academic Press, New York, 1976, 227–240.

F. G. Gustavson [1972]: *Some Basic Techniques for Solving Sparse Systems of Linear Equations*, in: D. J. Rose & R. A. Willoughby (eds.), “Sparse Matrices and Their Applications”, Plenum Press New York, 1972, 41–52.

P. M. J. Harris [1975]: *Pivot Selection Methods of the Devex LP Code*, Mathematical Programming Study 4 (1975) 30–57.

R. G. Jeroslow [1973]: *The Simplex Algorithm with the Pivot Rule of Maximizing Improvement Criterion*, Discrete Mathematics 4 (1973) 367–377.

V. Klee & G. J. Minty [1972]: *How Good is the Simplex Algorithm?*, in: O. Shisha (ed.), “Inequalities - III”, Academic Press, New York, 1972, 159–175.

H. Kuhn & R. E. Quandt [1963]: *An Experimental Study of the Simplex Method*, in: N. C. Metropolis et al. (eds.), “Experimental Arithmetic, High-Speed Computing and Mathematics”, Proceedings of Symposia on Applied Mathematics XV, American Mathematical Society, Providence, RI 1963, 107–124.

M. A. Saunders [1976]: *The Complexity of LU Updating in the Simplex Method*, in: R. S. Anderssen & R. P. Brent (eds.), “The Complexity of Computational Problem Solving”, Queensland University Press, Queensland, 1976, 214–230.

M. A. Saunders [1976]: *A Fast, Stable Implementation of the Simplex Method Using Bartels-Golub Updating*, in: J. R. Bunch & D. J. Rose (eds.), “Sparse Matrix Computations”, Academic Press New York, 1976, 213–226.

M. J. Todd [1984]: *Modifying the Forrest-Tomlin and Saunders Updates for Linear Programming Problems with Variable Upper Bounds*, Cornell University, School of OR/IE 1984, TR 617.

J. A. Tomlin [1975]: *On Scaling Linear Programming Problems*, Mathematical Programming Study (1975) 146–166.



# Kapitel 11

## Postoptimierung und parametrische Programme

Wir haben bisher die Theorie (und Praxis) des Simplexverfahrens entwickelt unter der Annahme, dass die Daten,  $A$ ,  $b$ , und  $c$  fest vorgegeben sind. Bei praktischen Problemen können jedoch häufig nicht alle der benötigten Zahlen mit Sicherheit angegeben werden, d. h. es kann Unsicherheit über die Beschränkungen  $b$  oder den zu erwartenden Gewinn  $c^T x$  geben. Man muss also die Tatsache mit in Betracht ziehen, dass die Wirklichkeit nicht exakt in das lineare Modell abgebildet wurde bzw. werden konnte. Darüber hinaus können im Nachhinein zusätzliche Variablen oder Restriktionen auftreten, die bei Aufstellung des Modells übersehen wurden oder nicht berücksichtigt werden konnten.

Wir werden uns nun überlegen, wie wir gewisse auftretende Änderungen behandeln können, wenn wir z. B. die optimale Lösung des ursprünglichen Programms gefunden haben. Wir wollen uns auf die folgenden Fälle beschränken:

1. Variation der rechten Seite  $b$ .
2. Variation der Zielfunktion  $c$ .
3. Änderung einer Nichtbasisspalte.
4. Hinzufügung einer neuen Variablen.
5. Hinzufügung einer neuen Restriktion.

Im Prinzip könnten wir natürlich versuchen, eine Theorie zu entwerfen, bei der

Schwankungen aller Daten berücksichtigt werden. Das ist jedoch außerordentlich kompliziert, und es gibt darüber kaum rechnerisch verwertbare Aussagen.

Wir gehen im weiteren davon aus, dass ein LP in Standardform (9.1) gegeben ist und dass wir eine optimale Basis  $A_B$  von  $A$  kennen.

## 11.1 Änderung der rechten Seite $b$

Wir wollen uns zunächst überlegen, welchen Einfluss eine Änderung der rechten Seite auf die Optimallösung bzw. den Wert der Optimallösung hat.

**(11.1) Änderung der Optimallösung.** Gegeben seien ein LP in Standardform (9.1)

$$\max\{c^T x \mid Ax = b, x \geq 0\}$$

und eine optimale Basis  $A_B$  von  $A$ . Die neue rechte Seite des LP sei  $b' := b + \Delta$ . Wir berechnen die neue Basislösung zur Basis  $A_B$ . Diese ergibt sich aus:

$$x'_B = A_B^{-1}(b + \Delta) = x_B + A_B^{-1}\Delta, \quad x'_N = 0.$$

Gilt  $x'_B \geq 0$ , so ist die neue Basislösung optimal, da sich ja an den reduzierten Kosten  $\bar{c}^T = c_N^T - c_B^T A_B^{-1} A_N$  nichts geändert hat, d. h. es gilt weiterhin  $\bar{c} \leq 0$ .

Gilt  $x'_B \not\geq 0$ , so ist die neue Basislösung primal nicht zulässig. Die Optimalitätsbedingung  $\bar{c} \leq 0$  ist jedoch weiterhin erfüllt. Das aber heißt nach (10.18), dass die Basis dual zulässig ist. Folglich haben wir mit  $A_B$  eine zulässige Basis für die duale Simplexmethode (10.27), und wir können direkt mit Phase II von (10.27) mit der Neuberechnung der Optimallösung beginnen.  $\square$

In diesem Zusammenhang sei auf eine wichtige Funktion, die die Änderung des Zielfunktionswertes eines LP bei Variationen der rechten Seite angibt, hingewiesen. Diese Funktion hat interessante Eigenschaften, von denen wir nachfolgend einige aufführen wollen. Die Funktion

$$L : \{b \in \mathbb{K}^m \mid P^=(A, b) \neq \emptyset\} \longrightarrow \mathbb{K} \cup \{\infty\}$$

definiert durch

$$(11.2) \quad L(b) := \sup\{c^T x \mid Ax = b, x \geq 0\}$$

heißt **Perturbationsfunktion** bzgl. des LPs

$$\begin{aligned} \max \quad & c^T x \\ \text{s.t.} \quad & Ax = b \\ & x \geq 0 \end{aligned}$$

Es bezeichne

$$\text{RP}(A) := \{b \in \mathbb{K}^m \mid P^=(A, b) \neq \emptyset\}$$

den Definitionsbereich von  $L$  (wegen  $0 \in \text{RP}(A)$  ist  $\text{RP}(A) \neq \emptyset$ ) und

$$\text{epi}(L) := \left\{ \begin{pmatrix} v \\ z \end{pmatrix} \in \mathbb{K}^{m+1} \mid v \in \text{RP}(A), z \in \mathbb{K}, L(v) \leq z \right\}$$

den Epigraphen von  $L$ .

**(11.3) Satz.** *Die Menge  $\text{epi}(L)$  ist ein polyedrischer Kegel.*

**Beweis:** Offenbar ist  $\text{epi}(L)$  eine Projektion des polyedrischen Kegels  $\{(x^T, v^T, z)^T \mid Ax - v = 0, x \geq 0, c^T x - z \leq 0\}$  auf die  $(v, z)$ -Koordinaten. Also ist nach (6.1)  $\text{epi}(L)$  ein Polyeder, das trivialerweise ein Kegel ist.  $\square$

**(11.4) Bemerkung.**

$$\exists b \in \text{RP}(A) \text{ mit } L(b) < \infty \iff \forall b \in \text{RP}(A) \text{ gilt } L(b) < \infty.$$

**Beweis:** Sei  $V := \{x \in \mathbb{K}^n \mid Ax = 0\}$ , dann gibt es zu jedem  $b \in \mathbb{K}^m$  einen Vektor  $x(b) \in \mathbb{K}^n$  mit  $P^=(A, b) = (V + x(b)) \cap \mathbb{K}_+^n$ . Folglich gilt für alle  $b \in \text{RP}(A)$ :  $L(b) = \infty \iff L(0) = \infty$ .  $\square$

**(11.5) Satz.** *Ist  $L \neq \infty$ , so gibt es Vektoren  $g^i \in \mathbb{K}^m$  ( $i = 1, \dots, k$ ), so daß*

$$L(v) = \max\{v^T g^i \mid i = 1, \dots, k\}$$

für alle  $v \in \text{RP}(A)$  gilt.

**Beweis:** Mit Satz (11.3) ist  $K := \text{epi}(L)$  ein polyedrischer Kegel. Also gibt es nach Bemerkung (2.6) eine  $(r, m+1)$ -Matrix  $B$ , so dass  $K = P(B, 0)$  gilt. Da  $L(0) \geq 0$ , ist  $\begin{pmatrix} 0 \\ z \end{pmatrix} \notin K$  für alle  $z < 0$ . Also kann die letzte Spalte von  $B$  kein Nullvektor sein, und wir können o. B. d. A. annehmen, dass  $B$  die Form

$$B = \begin{pmatrix} H & 0 \\ G & -h \end{pmatrix}$$

mit  $h \neq 0$  besitzt. Da  $L \neq \infty$ , folgt mit (11.4)  $L(0) < \infty$ , und daraus folgt direkt  $L(0) = 0$ . Also wissen wir, dass  $\begin{pmatrix} 0 \\ 1 \end{pmatrix} \in K$ . Dies impliziert  $h > 0$ . Sei o. B. d. A.  $h = 1$ . Dann gilt:  $\text{epi}(L) = \{(v^T, z)^T \mid Hv \leq 0, Gv \leq z\}$ , und es ist  $L(v) = z \iff G_i v = z$  für ein  $i$ . Bezeichnen  $g^i$  ( $i = 1, \dots, k$ ) die Zeilen von  $G$ , so gilt

$$L(v) = \max\{v^T g^i \mid i = 1, \dots, k\}.$$

□

**(11.6) Folgerung.** Die Perturbationsfunktion  $L$  ist stückweise linear und konkav.

□

Speziell folgt daraus, dass sich der Wert der Zielfunktion eines linearen Programms stetig mit einer Variation der rechten Seite des LP ändert. Die stetige Änderung lässt sich durch (11.5) explizit angeben. Sie ist “fast überall” linear, bis auf einige “Knicke”.

## 11.2 Änderungen der Zielfunktion $c$

Wir gehen wieder davon aus, dass wir eine Optimalbasis  $A_B$  von (9.1) kennen und dass sich die Zielfunktion  $c$  um  $\Delta$  ändert. Da wir eine Basis  $A_B$  gegeben haben, können wir die neuen Kosten der Basislösung ausrechnen. Es gilt

$$\begin{aligned} (c^T + \Delta^T)x &= (c_B^T + \Delta_B^T)A_B^{-1}b + (c_N^T + \Delta_N^T - (c_B^T + \Delta_B^T)A_B^{-1}A_N)x_N \\ &= c_B^T \bar{b} + \bar{c}x_N + \Delta_B^T \bar{b} + \underbrace{(\Delta_N^T - \Delta_B^T \bar{A})}_{=: \bar{\Delta}} x_N. \end{aligned}$$

- Wird keiner der Koeffizienten von  $c_B^T$  geändert (ist also  $\Delta_B = 0$ ), ändert sich wegen  $x_N = 0$  der Zielfunktionswert der Basislösung zur Basis  $A_B$  nicht.
- Sind die neuen reduzierten Kosten  $\bar{c} + \bar{\Delta}$  nicht positiv, so ist das Optimalitätskriterium (9.13) (b) weiterhin erfüllt, d. h.  $A_B$  bleibt die optimale Basis, jedoch ändert sich u. U. der Wert der zugehörigen Basislösung wenn  $\Delta_B \neq 0$ .
- Ist einer der Koeffizienten  $\bar{c}_j + \bar{\Delta}_j$  positiv, so wenden wir Phase II des primalen Simplexalgorithmus mit (zulässiger) Anfangsbasislösung  $A_B$  auf das neue Problem an.

(Zur Berechnung von  $\bar{\Delta}$  genügt die Kenntnis von  $\bar{A}$  oder  $A_B^{-1}$  und  $A_N$ .)

An den obigen Berechnungen kann man sehen, wie man eine Schar von linearen Programmen, bei denen entweder nur  $b$  oder nur  $c$  einer Änderung unterworfen wird, lösen kann.

### 11.3 Änderungen des Spaltenvektors $A_{\cdot j}$ , $j = N(s)$

Wollen wir die  $s$ -te Nichtbasisspalte (d. h. die  $j$ -te Spalte von  $A$ ) um  $\Delta$  ändern und kennen wir die Basisinverse  $A_B^{-1}$ , so können wir die neue  $s$ -te Spalte von  $\bar{A}$  berechnen durch

$$\bar{A}_{\cdot s}' = A_B^{-1}(A_{\cdot j} + \Delta) = \bar{A}_{\cdot s} + A_B^{-1}\Delta.$$

Die Basislösung  $\bar{x}$  zu  $A_B$  bleibt weiterhin primal zulässig, da diese Änderung keinen Einfluss auf  $\bar{x}_B = \bar{b}$ ,  $\bar{x}_N = 0$  hat. Jedoch bleibt  $\bar{x}$  i. a. nicht optimal, da sich der  $s$ -te Koeffizient der reduzierten Kosten ( $j = N(s)$ ) wie folgt ändert

$$\begin{aligned} \bar{c}_s' &= c_j - c_B^T(\bar{A}_{\cdot s} + A_B^{-1}\Delta) \\ &= c_j - c_B^T\bar{A}_{\cdot s} - c_B^T A_B^{-1}\Delta \\ &= \bar{c}_s - c_B^T A_B^{-1}\Delta = \bar{c}_s - u^T\Delta, \end{aligned}$$

wobei  $u$  eine optimale duale Lösung ist.

Die Optimalität der gegenwärtigen Basislösung  $\bar{x}$  bleibt genau dann erhalten, wenn  $\bar{c}_s' \leq u^T\Delta$  gilt. Andernfalls können wir die neue Optimallösung mit Phase II des primalen Simplexalgorithmus berechnen, wobei wir von der zulässigen Basis  $A_B$  ausgehen (die revidierte Methode ist natürlich von Vorteil).

(Bei Änderung einer Basisspalte kann man kaum Vorhersagen machen. Es muss neu invertiert werden. Dabei kann sowohl primale als auch duale Zulässigkeit verloren gehen.)

**(11.7) Beispiel.** Wir betrachten nochmals unser Beispiel (10.30)

$$\begin{aligned} \max \quad & -x_1 - 2x_2 \\ & -x_1 - x_2 \leq -2 \\ & -2x_1 - x_2 \leq -3 \\ & x_1, x_2 \geq 0 \end{aligned}$$

Das optimale (Tucker) Tableau ist:

-1	-1	-2
-2	1	1
-1	1	2

duale Lösung  $y_3 = 1, y_4 = 0$  primale Lösung  $x_1 = 2, x_2 = 0$

Die Spalte  $A_{.2}$  ist eine Nichtbasisspalte. Wir untersuchen, um wieviel diese Spalte geändert werden kann, ohne Optimalität zu verlieren.

$$A_{.2} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \quad \Delta = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}, \quad \bar{c}_2 = -1, \quad u = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\text{also } \bar{c}'_2 = -1 - (1, 0) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = -1 - \alpha$$

$$\bar{c}'_2 \leq 0 \iff -1 - \alpha \leq 0 \iff \alpha \geq -1.$$

Die obige Basislösung ist also für alle linearen Programme der Form

$$\begin{aligned} \max \quad & -x_1 - 2x_2 \\ -x_1 + (\alpha - 1)x_2 & \leq -2 \\ -2x_1 + (\beta - 1)x_2 & \leq -3 \\ x_1, x_2 & \geq 0 \end{aligned}$$

optimal, wenn nur  $\alpha \geq -1$  gilt. Zur geometrischen Interpretation siehe Abbildung 11.1.  $\square$

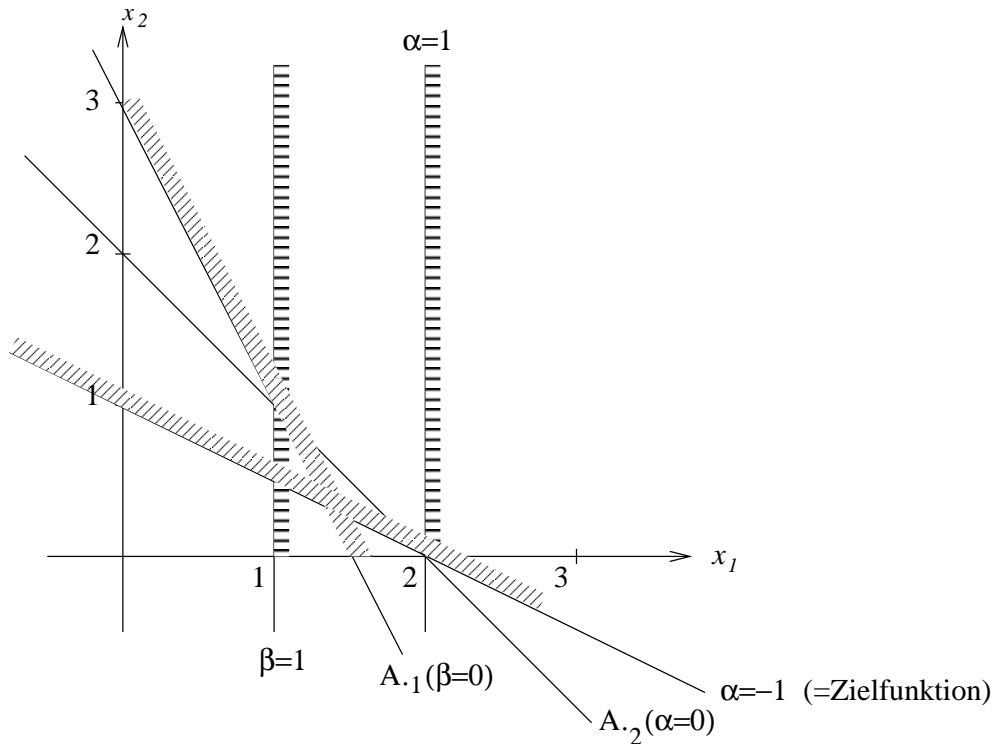


Abb. 11.1

## 11.4 Hinzufügung einer neuen Variablen

Sei  $A_{n+1}$  die neue Spalte von  $A$  und  $c_{n+1}$  der Zielfunktionswert der neuen Variablen  $x_{n+1}$ . Wir verlängern den Vektor  $N$  der Nichtbasisindizes um eine Komponente und setzen  $N(n-m+1) := n+1$ . Wir berechnen

$$\bar{c}_{n-m+1} := c_{n+1} - \underbrace{c_B^T A_B^{-1}}_{u^T} A_{n+1}.$$

Gilt  $\bar{c}_{n-m+1} \leq 0$ , so bleibt aufgrund von (9.13) (b) die alte Lösung optimal. Andernfalls führen wir mit der Anfangsbasis  $A_B$  den primalen Simplexalgorithmus aus. Dabei ist

$$\bar{A}_{n-m+1} = A_B^{-1} A_{n+1}.$$

Im ersten Schritt wird die neue Variable in die Basis aufgenommen.

## 11.5 Hinzufügung einer neuen Restriktion

Zu den bestehenden Nebenbedingungen fügen wir die neue Restriktion

$$A_{m+1,\cdot}x = \sum_{i=1}^n a_{m+1,i}x_i = b_{m+1}$$

hinzu.

1. Fall: Die Basislösung  $\bar{x}$  zu  $A_B$  erfüllt  $A_{m+1,\cdot}\bar{x} = b_{m+1}$ .

In diesem Fall ist  $\bar{x}$  auch die optimale Lösung des erweiterten Problems. Ist die neue Zeile linear abhängig von den übrigen Zeilen von  $A$ , so ist die neu hinzugefügte Gleichung irrelevant für das Problem und kann gestrichen werden. Ist die neue Zeile linear unabhängig von den übrigen, so ist  $\bar{x}$  eine entartete Basislösung. Eine der Nichtbasisvariablen wird mit Wert 0 in die Basis aufgenommen und die Basis entsprechend erweitert.

2. Fall: Die Basislösung  $\bar{x}$  erfüllt die neue Gleichung nicht.

Wir führen eine neue Schlupfvariable  $x_{n+1} \geq 0$  (mit dem Zielfunktionswert  $c_{n+1} = 0$ ) ein und verändern die neue Gleichung in

$$\sum_{i=1}^n a_{m+1,i}x_i \pm x_{n+1} = b_{m+1},$$

wobei wir ein  $+$  Zeichen schreiben, falls  $A_{m+1,\cdot}\bar{x} > b_{m+1}$  gilt, andernfalls ziehen wir  $x_{n+1}$  ab. Wir verlängern  $B$  um eine Komponente und setzen  $B(m+1) := n+1$ , d. h. wir nehmen  $x_{n+1}$  mit dem Wert  $\bar{b}_{m+1} = -|b_{m+1} - A_{m+1,\cdot}\bar{x}|$  in die Basis auf. Da  $\bar{b}_{m+1}$  negativ ist, ist die neue Basis primal nicht zulässig; wir haben jedoch keine Änderung an der Zielfunktion vorgenommen, d. h. die reduzierten Kosten sind weiterhin nicht positiv. Also ist die gegenwärtige neue Basis dual zulässig. Die neue Basis hat die Form

$$\begin{pmatrix} A_B & 0 \\ A_{m+1,B} & 1 \end{pmatrix}.$$

Die Basisinverse lässt sich wie folgt schreiben:

$$\begin{pmatrix} A_B^{-1} & 0 \\ -A_{m+1,B}A_B^{-1} & 1 \end{pmatrix}.$$

Die neue  $(m + 1)$ -te Zeile von  $\bar{A}$  lautet:

$$-A_{m+1,B}\bar{A} + A_{m+1,N} = \bar{A}_{m+1,\cdot}.$$

Wir führen nun einen dualen Pivotschritt auf der  $(m + 1)$ -ten Zeile von  $\bar{A}$  durch, wobei wir versuchen, die Variable  $x_{n+1}$  aus der Basis zu entfernen.

Ergibt der duale Beschränktheitstest ( $\bar{A}_{m+1,\cdot} \geq 0$  (10.27) (II.3)), dass das duale Problem unbeschränkt ist, so ist das neue primale Problem unzulässig, d. h. die Hyperebene  $A_{m+1,\cdot}x = b_{m+1}$  hat einen leeren Schnitt mit der Menge der zulässigen Lösungen des alten Problems, und wir können den Algorithmus beenden.

Andernfalls führen wir Phase II des dualen Simplexalgorithmus (10.27) aus. Endet der Algorithmus mit einer primal und dual zulässigen Lösung, so gibt es eine optimale primale Lösung  $x^*$  mit  $x_{n+1}^* = 0$ . Diese kann wie folgt konstruiert werden, falls  $x_{n+1}^*$  nicht bereits Null ist:

Sei  $x'$  die primale zulässige Basislösung nach Beendigung des dualen Verfahrens und  $\bar{x}$  die Anfangsbasislösung bei Beginn des dualen Programms, d. h.  $\bar{x}_{n+1} = \bar{b}_{m+1} < 0$ . Setzen wir

$$\lambda := \frac{x'_{n+1}}{x'_{n+1} - \bar{x}_{n+1}},$$

so gilt  $\lambda > 0$ , da  $x'_{n+1} > 0$  und  $\bar{x}_{n+1} < 0$ . Sei  $x^* := \lambda\bar{x} + (1 - \lambda)x'$ , dann gilt

$$\begin{aligned} x_i^* &\geq 0 \text{ für } i = 1, \dots, n, \text{ da } \bar{x}_i \geq 0, x'_i \geq 0 \\ x_{n+1}^* &= \frac{x'_{n+1}}{x'_{n+1} - \bar{x}_{n+1}}\bar{x}_{n+1} + x'_{n+1} - \frac{x'_{n+1}}{x'_{n+1} - \bar{x}_{n+1}}x'_{n+1} \\ &= \frac{1}{x'_{n+1} - \bar{x}_{n+1}}(x'_{n+1}\bar{x}_{n+1} - x'_{n+1}x'_{n+1} + x'_{n+1}x'_{n+1} - \bar{x}_{n+1}x'_{n+1}) \\ &= 0. \end{aligned}$$

Also gilt  $x^* \geq 0$  und ferner

$$\begin{aligned} c^T x^* &= \lambda \underbrace{c^T \bar{x}}_{\geq c^T x'} + (1 - \lambda)c^T x' \\ &\geq c^T x'. \end{aligned}$$

Daraus folgen die Zulässigkeit und die Optimalität von  $x^*$ .



# Kapitel 12

## Die Ellipsoidmethode

In (10.10) haben wir angemerkt, dass man Beispiele von Polyedern und Zielfunktionen konstruieren kann, so dass der Simplexalgorithmus alle Ecken der Polyeder durchläuft. Die in den Übungen besprochene Klasse von Klee & Minty-Beispielen, bei denen dieser Fall auftritt, ist definiert durch  $n$  Variable und  $2n$  Ungleichungen. Die zugehörigen Polyeder haben  $2^n$  Ecken. Es ist offensichtlich, dass  $2^n$  Iterationen des Simplexalgorithmus zu nicht tolerierbaren Rechenzeiten führen. Man kann zwar zeigen (Borgwardt “The Average Number of Pivot Steps Required by the Simplex-Method is Polynomial”, Zeitschrift für Operations Research 26 (1982) 157–177), dass das Laufzeitverhalten des Simplexalgorithmus im Durchschnitt gut ist, aber dennoch bleibt die Frage, ob es nicht möglich ist, eine generelle “niedrige” Schranke für die Anzahl der Iterationenschritte des Simplexalgorithmus (mit einer speziellen Zeilen- und Spaltenauswahlregel) zu finden. Dieses Problem ist bis heute ungelöst!

Im Jahre 1979 hat L. G. Khachiyan (“A polynomial algorithm in linear programming”, Doklady Akademia Nauk SSSR 244 (1979) 1093–1096, engl. Übersetzung in Soviet Mathematics Doklady 20 (1979) 191–194) gezeigt, dass lineare Programme in “polynomialer Zeit” gelöst werden können. Das Verfahren, das er dazu angegeben hat, die sogenannte Ellipsoidmethode, arbeitet allerdings völlig anders als der Simplexalgorithmus und scheint in der Praxis im Laufzeitverhalten dem Simplexalgorithmus “im Durchschnitt” weit unterlegen zu sein. Theoretisch beweisbar ist jedoch, dass es keine Beispiele linearer Programme (z. B. vom Klee & Minty-Typ) gibt, die die Ellipsoidmethode zu exorbitanten Laufzeiten zwingen.

Die dieser Methode zugrundeliegenden Ideen benutzen ganz andere geometrische Überlegungen, als wir sie bisher gemacht haben. Außerdem sind zu einer korrekten Implementation so viele numerische Details zu klären, dass der zeitliche

Rahmen dieser Vorlesung durch eine vollständige Diskussion aller Probleme gesprengt würde. Wir wollen daher nur einen Überblick über die Methode geben und nur einige Beweise ausführen.

## 12.1 Polynomiale Algorithmen

Um das Laufzeitverhalten eines Algorithmus exakt beschreiben zu können, ist es notwendig Maschinenmodelle, Kodierungsvorschriften, Algorithmusbeschreibungen etc. exakt einzuführen. Dies wird z. B. in einer Vorlesung, die das Thema „Komplexitätstheorie“ ausführlich behandelt, gemacht. Hier wollen wir die wichtigsten Konzepte dieser Theorie informell für den Spezialfall der linearen Programmierung einführen. Eine mathematisch exakte Darstellung der Theorie findet man z. B. in dem Buch

M. R. Garey & D. S. Johnson, “Computers and Intractability: A Guide to the Theory of  $\mathcal{NP}$ -Completeness”, Freeman, San Francisco, 1979, (80 ST 230 G 229).

Zunächst müssen wir unser Konzept, dass Zahlen aus dem Körper  $\mathbb{K}$  gewählt werden können, aufgeben. Jede Zahl, die wir in unserem Berechnungsmodell betrachten, muss endlich kodierbar sein. Es gibt aber kein Kodierungsschema, so dass z. B. alle reellen Zahlen durch endlich viele Symbole beschrieben werden könnten. Wir beschränken uns daher auf rationale Zahlen. Haben wir eine Ungleichung

$$a^T x \leq \alpha$$

mit  $a \in \mathbb{Q}^n$ ,  $\alpha \in \mathbb{Q}$ , so können wir auf einfache Weise das kleinste gemeinsame Vielfache  $p$  der Nenner der Komponenten von  $a$  und des Nenners von  $\alpha$  bestimmen. Die Ungleichung

$$pa^T x \leq p\alpha$$

hat offenbar die gleiche Lösungsmenge wie  $a^T x \leq \alpha$ , während alle Koeffizienten dieser Ungleichung ganzzahlig sind. Der Fall rationaler Daten lässt sich also direkt auf den Fall ganzzahliger Daten reduzieren. Es ist zwar häufig einfacher unter der Voraussetzung ganzzahliger Daten zu rechnen, und fast alle Aufsätze zur Ellipsoidmethode machen diese Annahme, wir wollen aber dennoch für dieses Kapitel voraussetzen:

**(12.1) Annahme.** *Alle Daten linearer Programme sind rational, d. h. für jedes LP der Form  $\max c^T x$ ,  $Ax \leq b$  gilt  $c \in \mathbb{Q}^n$ ,  $A \in \mathbb{Q}^{(m,n)}$ ,  $b \in \mathbb{Q}^m$ .*  $\square$

Nun müssen wir die Frage klären, wie Daten “gegeben” sind. Da unsere Computer üblicherweise mit Binärcodes arbeiten, wollen wir annehmen, dass alle vorkom-

menden ganzen Zahlen binär kodiert sind. Die binäre Darstellung einer ganzen Zahl  $n$  benötigt  $\lceil \log_2(|n| + 1) \rceil$  Stellen (Bits) und eine Stelle für das Vorzeichen. Wir nennen daher

$$(12.2) \quad \langle n \rangle := \lceil \log_2(|n| + 1) \rceil + 1$$

die **Kodierungslänge** von  $n$ . Jede rationale Zahl  $r$  hat eine Darstellung  $r = \frac{p}{q}$  mit  $p$  und  $q$  teilerfremd und  $q > 0$ . Wir nehmen immer an, dass rationale Zahlen in dieser Form dargestellt sind. Also können wir sagen, dass die Kodierungslänge einer rationalen Zahl  $r = \frac{p}{q}$  gegeben ist durch

$$\langle r \rangle := \langle p \rangle + \langle q \rangle.$$

Die **Kodierungslänge einer Matrix**  $A = (a_{ij}) \in \mathbb{Q}^{(m,n)}$  (oder analog eines Vektors) ist gegeben durch

$$(12.3) \quad \langle A \rangle := \sum_{i=1}^m \sum_{j=1}^n \langle a_{ij} \rangle.$$

Daraus folgt, daß die Kodierungslänge eines linearen Programms der Form  $\max c^T x$ ,  $Ax \leq b$  gegeben ist durch

$$\langle c \rangle + \langle A \rangle + \langle b \rangle.$$

Um über Laufzeiten sprechen zu können, müssen wir uns überlegen, wie wir die Laufzeit eines Algorithmus messen wollen. Wir legen fest, dass wir dazu zunächst die

*Anzahl der elementaren Rechenschritte*

zählen wollen, wobei elementare Rechenschritte die arithmetischen Operationen

*Addition, Subtraktion, Multiplikation, Division, Vergleich*

von ganzen Zahlen oder rationalen Zahlen sind. Dies reicht aber noch nicht aus, denn wie jeder weiß, dauert z. B. die Multiplikation großer ganzer Zahlen länger als die Multiplikation kleiner ganzer Zahlen. Wir müssen also die Zahlengrößen in Betracht ziehen und jede elementare Rechenoperation mit der Kodierungslänge der dabei involvierten Zahlen multiplizieren. Für unsere Zwecke ist folgende Definition angemessen.

**(12.4) Definition.** Die **Laufzeit eines Algorithmus**  $\mathcal{A}$  zur Lösung eines Problems  $\Pi$  (kurz  $L_{\mathcal{A}}(\Pi)$ ) ist die Anzahl der elementaren Rechenschritte, die während der Ausführung des Algorithmus durchgeführt wurden, multipliziert mit der Kodierungslänge der (bezüglich ihrer Kodierungslänge) größten Zahl, die während der Ausführung des Algorithmus aufgetreten ist.  $\square$

Die beiden folgenden Begriffe sind entscheidend für das Verständnis des Weiteren.

**(12.5) Definition.** Sei  $\mathcal{A}$  ein Algorithmus zur Lösung einer Klasse von Problemen (z. B. die Klasse aller linearen Optimierungsprobleme). Die Elemente (Probleme) dieser Klasse (z. B. spezielle lineare Programme) bezeichnen wir mit  $\Pi$ .

(a) Die Funktion  $f : \mathbb{N} \rightarrow \mathbb{N}$  definiert durch

$$f_{\mathcal{A}}(n) := \max\{L_{\mathcal{A}}(\Pi) \mid \text{die Kodierungslänge der Daten zur Beschreibung von } \Pi \text{ ist höchstens } n\}$$

heißt **Laufzeitfunktion** von  $\mathcal{A}$ .

(b) Der Algorithmus  $\mathcal{A}$  hat eine **polynomiale Laufzeit** (kurz:  $\mathcal{A}$  ist ein **polynomialer Algorithmus**), wenn es ein Polynom  $p : \mathbb{N} \rightarrow \mathbb{N}$  gibt, so daß

$$f_{\mathcal{A}}(n) \leq p(n) \quad \forall n \in \mathbb{N} \text{ gilt.}$$

□

Polynomiale Algorithmen sind also solche Verfahren, deren Schrittzahl multipliziert mit der Kodierungslänge der größten auftretenden Zahl durch ein Polynom in der Kodierungslänge der Problemdaten beschränkt werden kann. Für die Klasse der linearen Programme der Form  $\max c^T x, Ax \leq b$  muss also die Laufzeit eines Algorithmus durch ein Polynom in  $\langle A \rangle + \langle b \rangle + \langle c \rangle$  beschränkt werden können, wenn er polynomial sein soll. Wie die Klee & Minty-Beispiele zeigen, ist der Simplexalgorithmus kein polynomialer Algorithmus zur Lösung linearer Programme!

## 12.2 Reduktionen

Die Ellipsoidmethode ist kein Optimierungsverfahren, sondern lediglich eine Methode, die in einem gegebenen Polyeder einen Punkt findet, falls ein solcher existiert. Wir müssen daher das allgemeine lineare Optimierungsproblem auf diesen Fall zurückführen. Aus Kapitel 2 (allgemeine Transformationsregeln) wissen wir, dass jedes lineare Programm in der Form

$$(12.6) \quad \begin{array}{ll} \max & c^T x \\ & Ax \leq b \\ & x \geq 0 \end{array}$$

geschrieben werden kann. Das zu (12.6) duale Programm ist

$$(12.7) \quad \begin{array}{ll} \min & b^T y \\ & A^T y \geq c \\ & y \geq 0. \end{array}$$

Aus dem Dualitätssatz (5.14) wissen wir, dass (12.6) und (12.7) genau dann optimale Lösungen haben, deren Werte gleich sind, wenn beide Programme zulässige Lösungen haben. Daraus folgt, dass jedes Element  $\begin{pmatrix} x \\ y \end{pmatrix}$  des Polyeders  $P$ , definiert durch

$$(12.8) \quad \begin{pmatrix} -c^T & b^T \\ A & 0 \\ 0 & -A^T \\ -I & 0 \\ 0 & -I \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \leq \begin{pmatrix} 0 \\ b \\ -c \\ 0 \\ 0 \end{pmatrix}$$

eine Optimallösung  $x$  von (12.6) und eine Optimallösung  $y$  von (12.7) bestimmt. Dies impliziert, dass es zur Lösung von (12.6) genügt, einen Punkt in (12.8) zu finden.

Der obige „Trick“, das Primal- und das Dualprogramm zusammenzufassen, bläht natürlich das zu bearbeitende System ungeheuer auf. Eine andere Methode der Reduktion ist die binäre Suche, die wir nun schildern. Zur korrekten Darstellung benötigen wir jedoch noch einige (auch für andere Zwecke) wichtige Hilfssätze. Zunächst wollen wir einige Parameter durch die Kodierungslänge abschätzen. Um den ersten Hilfssatz zu beweisen, benötigen wir die bekannte Hadamard-Ungleichung, die besagt, dass das Volumen eines Parallelepipeds in  $\mathbb{R}^n$  mit Kantenlängen  $d_1, \dots, d_n$  nicht größer ist als das Volumen des Würfels mit Kantenlängen  $d_1, \dots, d_n$ . Bekanntlich ist das Volumen eines Parallelepipeds, das durch die Vektoren  $a_1, \dots, a_n$  aufgespannt wird, nichts anderes als der Absolutbetrag der Determinante der Matrix  $A$  mit Spalten  $A_{.1} = a_1, \dots, A_{.n} = a_n$ . Die **Hadamard-Ungleichung** lautet also

$$(12.9) \quad |\det A| \leq \prod_{j=1}^n \|A_{.j}\|_2.$$

**(12.10) Lemma.**

- (a) Für jede Zahl  $r \in \mathbb{Q}$  gilt:  $|r| \leq 2^{\langle r \rangle - 1} - 1$ .
- (b) Für je zwei rationale Zahlen  $r, s$  gilt:  $\langle rs \rangle \leq \langle r \rangle + \langle s \rangle$ .
- (c) Für jeden Vektor  $x \in \mathbb{Q}^n$  gilt:  $\|x\|_2 \leq \|x\|_1 \leq 2^{\langle x \rangle - n} - 1$ .
- (d) Für jede Matrix  $A \in \mathbb{Q}^{(n,n)}$  gilt:  $|\det(A)| \leq 2^{\langle A \rangle - n^2} - 1$ .

**Beweis :**

(a) und (b) folgen direkt aus der Definition.

(c) Sei  $x = (x_1, \dots, x_n)^T \in \mathbb{Q}^n$ , dann gilt nach (a)

$$\begin{aligned} 1 + \|x\|_1 &= 1 + \sum_{i=1}^n |x_i| \leq \prod_{i=1}^n (1 + |x_i|) \\ &\leq \prod_{i=1}^n 2^{\langle x_i \rangle - 1} \\ &= 2^{\langle x \rangle - n}. \end{aligned}$$

Die Ungleichung  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \leq \sum_{i=1}^n |x_i| = \|x\|_1$  ist trivial.

(d) Aus der Hadamard-Ungleichung (12.9) und (c) folgt

$$\begin{aligned} 1 + |\det(A)| &\leq 1 + \prod_{j=1}^n \|A_{\cdot j}\|_2 \leq \prod_{j=1}^n (1 + \|A_{\cdot j}\|_2) \\ &\leq \prod_{j=1}^n 2^{\langle A_{\cdot j} \rangle - n} = 2^{\langle A \rangle - n^2}. \end{aligned}$$

□

Diesen Hilfssatz können wir zur Abschätzung der „Größe“ der Ecken von Polyedern wie folgt benutzen.

**(12.11) Satz.** Für jede Ecke  $v = (v_1, \dots, v_n)^T$  eines Polyeders  $P$  der Form  $P(A, b)$ ,  $P^=(A, b)$  oder  $P = \{x \in \mathbb{R}^n \mid Ax \leq b, x \geq 0\}$ ,  $A \in \mathbb{Q}^{(m,n)}$ ,  $b \in \mathbb{Q}^m$ , gilt

(a) Der Absolutbetrag des Zählers von  $v_i$  ist höchstens  $2^{\langle A \rangle + \langle b \rangle - n^2}$ , der Absolutbetrag des Nenners von  $v_i$  ist höchstens  $2^{\langle A \rangle - n^2}$ ,  $i = 1, \dots, n$ .

(b)  $|v_i| \leq 2^{2\langle A \rangle + \langle b \rangle - 2n^2}$ ,  $i = 1, \dots, n$ .

(c) Falls  $A \in \mathbb{Z}^{(m,n)}$ , so gilt  $|v_i| \leq 2^{\langle A \rangle + \langle b \rangle - n^2}$   $i = 1, \dots, n$ .

**Beweis :** Ist  $v$  Ecke von  $P$ , so gibt es nach Satz (8.9) eine reguläre Untermatrix und einen entsprechenden Untervektor des Restriktionensystems von  $P$ , so dass  $v$  die eindeutig bestimmte Lösung des hierdurch bestimmten Gleichungssystems ist. Wir führen den Fall  $P = \{x \mid Ax \leq b, x \geq 0\}$  vor. Die beiden anderen Fälle verlaufen analog.

Es gibt also eine  $(n, n)$ -Matrix  $D$ , deren Zeilen Zeilen von  $A$  oder negative Einheitsvektoren sind und einen entsprechenden Vektor  $d$ , dessen Komponenten Elemente von  $b$  oder Nullen sind, so dass  $v$  die eindeutige Lösung von  $Dx = d$  ist. Nach Cramer's Regel gilt dann für  $i = 1, \dots, n$

$$v_i = \frac{\det D_i}{\det D}$$

wobei  $D_i$  aus  $D$  dadurch entsteht, dass die  $i$ -te Spalte von  $D$  durch den Vektor  $d$  ersetzt wird. Hat  $D$  Zeilen, die negative Einheitsvektoren sind, so bezeichnen wir mit  $\overline{D}$  die Matrix, die durch Streichen dieser Zeilen und der Spalten, in denen sich die Elemente  $-1$  befinden, entsteht. Aufgrund des Determinantenentwicklungssatzes gilt  $|\det D| = |\det \overline{D}|$ , und ferner ist  $\overline{D}$  eine Untermatrix von  $A$ . Daraus folgt mit (12.10) (d)

$$|\det D| = |\det \overline{D}| \leq 2^{\langle \overline{D} \rangle - n^2} \leq 2^{\langle A \rangle - n^2}.$$

Analog folgt

$$|\det D_i| \leq 2^{\langle A \rangle + \langle b \rangle - n^2}.$$

Damit ist (a) bewiesen.

Ist  $|\det D| \geq 1$ , dies ist z. B. dann der Fall, wenn  $A \in \mathbb{Z}^{(m,n)}$  gilt, so erhalten wir mit (a)

$$|v_i| \leq |\det D_i| \leq 2^{\langle A \rangle + \langle b \rangle - n^2}.$$

Hieraus folgt (c).

Ist  $|\det D| < 1$ , so müssen wir  $|\det D|$  nach unten abschätzen. Sei  $q = \prod_{i,j} q_{ij}$  das Produkt der Nenner der Elemente  $d_{ij} = \frac{p_{ij}}{q_{ij}}$  von  $D$ . Dann gilt  $|\det D| = \frac{q|\det D|}{q}$ , und sowohl  $q$  als auch  $q|\det D|$  sind ganze Zahlen. Aus (12.10) (a), (b) folgt

$$q = \prod_{i,j} q_{ij} \leq 2^{\langle \Pi q_{ij} \rangle} \leq 2^{\sum \langle q_{ij} \rangle} \leq 2^{\langle D \rangle - n^2} \leq 2^{\langle A \rangle - n^2}.$$

Daraus ergibt sich

$$|v_i| \leq \frac{|\det D_i|}{|\det D|} = \frac{|\det D_i|q}{q|\det D|} \leq |\det D_i|q \leq 2^{2\langle A \rangle + \langle b \rangle - 2n^2}.$$

Damit ist auch (b) bewiesen.  $\square$

Da nach Satz (8.12) alle Polyeder  $P$  der Form  $P = \{x \mid Ax \leq b, x \geq 0\}$  spitz sind und nach Folgerung (8.15) lineare Programme über spitzen Polyedern optimale Ecklösungen haben (falls sie überhaupt Optimallösungen haben) folgt:

**(12.12) Satz.** *Das lineare Programm (12.6) hat eine Optimallösung genau dann, wenn die beiden linearen Programme*

$$(12.13) \quad \begin{aligned} \max c^T x \\ Ax &\leq b \\ x &\geq 0 \\ x_i &\leq 2^{2\langle A \rangle + \langle b \rangle - 2n^2}, \quad i = 1, \dots, n \end{aligned}$$

$$(12.14) \quad \begin{aligned} \max c^T x \\ Ax &\leq b \\ x &\geq 0 \\ x_i &\leq 2^{2\langle A \rangle + \langle b \rangle - 2n^2} + 1, \quad i = 1, \dots, n \end{aligned}$$

*eine Optimallösung haben und die Werte der Optimallösungen übereinstimmen. Die Zielfunktionswerte von (12.13) und (12.14) stimmen genau dann nicht überein, wenn (12.6) unbeschränkt ist. (12.6) ist genau dann nicht lösbar, wenn (12.13) oder (12.14) nicht lösbar ist.  $\square$*

Wir haben damit das lineare Programmierungsproblem (12.6) über einem (allgemeinen) Polyeder auf die Lösung zweier linearer Programme über Polytopen reduziert. Wir müssen also im Prinzip nur zeigen, wie man LP's über Polytopen löst.

**(12.15) Satz.** *Ist  $P \neq \emptyset$  ein Polytop der Form  $P(A, b)$ ,  $P^=(A, b)$  oder  $P = \{x \mid Ax \leq b, x \geq 0\}$  mit  $A \in \mathbb{Q}^{(m,n)}$ ,  $b \in \mathbb{Q}^m$ , so kann für  $c \in \mathbb{Q}^n$  das lineare Programm  $\max c^T x, x \in P$  nur Optimalwerte in der endlichen Menge*

$$S := \left\{ \frac{p}{q} \in \mathbb{Q} \mid |p| \leq n 2^{\langle A \rangle + \langle b \rangle + 2\langle c \rangle - n^2 - n}, 1 \leq q \leq 2^{\langle A \rangle + \langle c \rangle - n^2 - n} \right\}$$

*annehmen.*

**Beweis :** Da alle Optimalwerte Zielfunktionswerte von Ecken von  $P$  sind, brauchen wir nur die Werte  $c^T v$ ,  $v$  Ecke von  $P$ , abzuschätzen. Sei also  $v = (v_1, \dots, v_n)^T$  eine Ecke von  $P$ . Wie im Beweis von (12.11) können wir feststellen, dass die Komponenten  $v_i$  von  $v$  eine Darstellung der Form

$$v_i = \frac{\det D_i}{\det D} =: \frac{p_i}{d}$$

haben, wobei  $D$  eine reguläre Untermatrix von  $A$  bzw.  $\begin{pmatrix} A \\ I \end{pmatrix}$  ist und  $D_i$  aus  $D$  dadurch hervorgeht, dass die  $i$ -te Spalte von  $D$  durch einen Untervektor der rechten Seite ersetzt wird. Satz (12.11) zeigt  $d \leq 2^{\langle A \rangle - n^2}$ ,  $p_i \leq 2^{\langle A \rangle + \langle b \rangle - n^2}$ . Sei nun

$c = (\frac{s_1}{t_1}, \dots, \frac{s_n}{t_n})^T \in \mathbb{Q}^n$  eine teilerfremde Darstellung der Zielfunktion, so gilt

$$c^T v = \sum_{i=1}^n \frac{s_i p_i}{t_i d} = \frac{1}{dt} \sum_{i=1}^n s_i p_i \bar{t}_i \quad \text{mit } t := t_1 \cdot \dots \cdot t_n, \bar{t}_i := \frac{t}{t_i}.$$

Aus (12.10) (a) folgt  $t = t_1 \cdot \dots \cdot t_n \leq 2^{(t_1)-1} \cdot \dots \cdot 2^{(t_n)-1} = 2^{\sum (t_i)-1} \leq 2^{(c)-n}$  und somit

$$q := dt \leq 2^{(A)+(c)-n^2-n}.$$

Analog folgt

$$p := \sum_{i=1}^n s_i p_i \bar{t}_i \leq \sum_{i=1}^n 2^{(s_i)-1} 2^{(A)+(b)-n^2} 2^{(c)-n} \leq n 2^{(A)+(b)+2(c)-n^2-n}.$$

□

Satz (12.15) gibt uns nun die Möglichkeit, das Verfahren der binären Suche zur Lösung von  $\max c^T x, x \in P$  anzuwenden. Diese Methode funktioniert bezüglich der in Satz (12.15) definierten Menge  $S$  wie folgt.

### (12.16) Binäre Suche.

- (1) Wähle ein Element  $s \in S$ , so dass für  $S' := \{t \in S \mid t < s\}$  und  $S'' := \{t \in S \mid t \geq s\}$  gilt

$$|S'| \leq |S''| \leq |S'| + 1.$$

- (2) Überprüfe, ob das Polyeder

$$P_s = \{x \mid Ax \leq b, x \geq 0, c^T x \geq s\}$$

nicht leer ist.

- (3) Ist  $P_s$  leer, so setze  $S := S'$ , andernfalls setze  $S := S''$ .

- (4) Ist  $|S| = 1$ , so gilt für  $s \in S$ :  $s = \max\{c^T x \mid Ax \leq b, x \geq 0\}$ , und jeder Punkt in  $P_s$  ist eine Optimallösung von  $\max c^T x, x \in P$ . Andernfalls gehe zu (1).

□

Die Korrektheit dieses Verfahrens ist offensichtlich. Ist die Binärsuche auch effizient? Da in jedem Schritt die Kardinalität  $|S|$  von  $S$  (fast) halbiert wird, ist klar, dass höchstens

$$(12.17) \quad \overline{N} := \lceil \log_2(|S| + 1) \rceil$$

Aufspaltungen von  $S$  in zwei Teile notwendig sind, um eine einelementige Menge zu erhalten. Also muss Schritt (2) von (12.16)  $\overline{N}$ -mal ausgeführt werden. Nach (12.15) gilt

$$|S| \leq 2n2^{2\langle A \rangle + \langle b \rangle + 3\langle c \rangle - 2n^2 - 2n} + 1,$$

und daraus folgt, dass Test (2) von (12.16) höchstens

$$(12.18) \quad \overline{N} := 2\langle A \rangle + \langle b \rangle + 3\langle c \rangle - 2n^2 - 2n + \log_2 n + 2$$

mal durchgeführt wurde. Ist Test (2) in einer Zeit ausführbar, die polynomial in  $\langle A \rangle + \langle b \rangle + \langle c \rangle$  ist, so ist auch die binäre Suche ein polynomialer Algorithmus, da ein polynomialer Algorithmus für (2) nur  $\overline{N}$ -mal also nur polynomial oft ausgeführt werden muss.

**(12.19) Folgerung.** *Es gibt einen polynomialen Algorithmus zur Lösung linearer Programme genau dann, wenn es einen Algorithmus gibt, der in polynomialer Zeit entscheidet, ob ein Polytop  $P$  leer ist oder nicht und der, falls  $P \neq \emptyset$ , einen Punkt in  $P$  findet.*

□

Damit haben wir das lineare Programmierungsproblem reduziert auf die Frage der Lösbarkeit von Ungleichungssystemen, deren Lösungsmenge beschränkt ist.

## 12.3 Beschreibung der Ellipsoidmethode

In diesem Abschnitt setzen wir  $n \geq 2$  voraus. Wir wollen zunächst einige Eigenschaften von Ellipsoiden beschreiben. Wir erinnern daran, dass Ellipsoide volldimensionale konvexe Körper im  $\mathbb{R}^n$  sind, die eine besonders einfache Darstellung haben. Und zwar ist eine Menge  $E \subseteq \mathbb{R}^n$  ein **Ellipsoid (mit Zentrum  $a$ )** genau dann, wenn es einen Vektor  $a \in \mathbb{R}^n$  und eine (symmetrische) positiv definite Matrix  $A$  gibt, so dass gilt

$$(12.20) \quad E = E(A, a) := \{x \in \mathbb{R}^n \mid (x - a)^T A^{-1}(x - a) \leq 1\}.$$

Aus der linearen Algebra wissen wir, dass für symmetrische Matrizen  $A \in \mathbb{R}^{(n,n)}$  die folgenden Aussagen äquivalent sind:

- (i)  $A$  ist positiv definit.
- (ii)  $A^{-1}$  ist positiv definit.
- (12.21) (iii) Alle Eigenwerte von  $A$  sind positive reelle Zahlen.
- (iv)  $\det(A_{II}) > 0$  für alle Mengen  $I = \{1, \dots, i\}, i = 1, \dots, n$ .
- (v)  $A = B^T B$  für eine reguläre Matrix  $B \in \mathbb{R}^{(n,n)}$ .

Das Ellipsoid  $E(A, a)$  ist durch die (ebenfalls positiv definite) Inverse  $A^{-1}$  von  $A$  definiert. Das erscheint zunächst seltsam, liegt aber daran, dass sich viele Eigenschaften von  $E(A, a)$  aus algebraischen Eigenschaften von  $A$  ableiten lassen. Z. B. ist der Durchmesser (d. h. die Länge der längsten Achse) von  $E(A, a)$  gleich  $2\sqrt{\Lambda}$ , wobei  $\Lambda$  der größte Eigenwert von  $A$  ist. Die längsten Achsen sind durch die Eigenvektoren, die zu  $\Lambda$  gehören, gegeben. Die Symmetrieachsen von  $E(A, a)$  entsprechen ebenfalls den Eigenvektoren von  $A$ . In Abbildung 12.1 ist das Ellipsoid  $E(A, 0)$  dargestellt mit

$$A = \begin{pmatrix} 16 & 0 \\ 0 & 4 \end{pmatrix}.$$

Die Eigenwerte von  $A$  sind  $\Lambda = 16$  und  $\lambda = 4$  mit den zugehörigen Eigenvektoren  $e_1 = (1, 0)^T$  und  $e_2 = (0, 1)^T$ . Der Durchmesser von  $E(A, 0)$  ist  $2\sqrt{\Lambda} = 8$ .

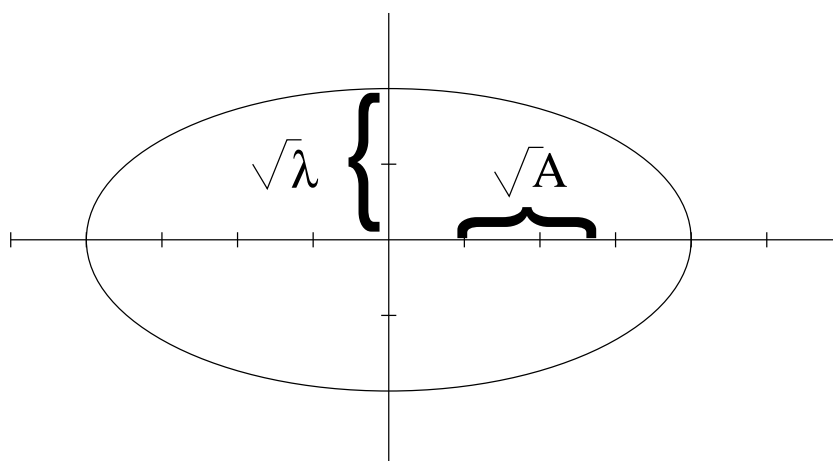


Abb. 12.1

Zu jeder positiv definiten Matrix  $A$  gibt es eine eindeutig bestimmte positive definite Matrix, die mit  $A^{\frac{1}{2}}$  bezeichnet und **Wurzel von  $A$**  genannt wird, mit

$$A = A^{\frac{1}{2}} A^{\frac{1}{2}}.$$

Die Einheitskugel  $S(0, 1) \subseteq \mathbb{R}^n$  (um den Nullpunkt mit Radius 1) ist das Ellipsoid  $E(I, 0)$ . Man rechnet leicht nach, dass gilt:

$$(12.22) \quad E(A, a) = A^{\frac{1}{2}} S(0, 1) + a.$$

Damit sei die Aufzählung von Eigenschaften von Ellipsoiden beendet. Wir wollen nun die geometrische Idee, die hinter der Ellipsoidmethode steckt, erläutern.

**(12.23) Geometrische Beschreibung der Ellipsoidmethode.** Gegeben sei ein Polytop  $P$ . Wir wollen einen Punkt in  $P$  finden oder beweisen, dass  $P$  leer ist.

- (1) Konstruiere ein Ellipsoid  $E_0 = E(A_0, a_0)$ , das  $P$  enthält. Setze  $k := 0$ .
- (2) Ist das gegenwärtige Ellipsoid “zu klein”, so brich ab mit der Antwort “ $P$  ist leer”.
- (3) Teste ob der Mittelpunkt  $a_k$  von  $E_k$  in  $P$  enthalten ist.
- (4) Gilt  $a_k \in P$ , dann haben wir unser Ziel erreicht, STOP.
- (5) Gilt  $a_k \notin P$ , dann gibt es eine  $P$  definierende Ungleichung, sagen wir

$$c^T x \leq \gamma,$$

die von  $a_k$  verletzt wird, d. h.  $c^T a_k > \gamma$ . Mit

$$E'_k := E_k \cap \{x \mid c^T x \leq c^T a_k\}$$

bezeichnen wir das Halbellipsoid von  $E_k$ , das  $P$  enthält. Wir konstruieren nun das Ellipsoid kleinsten Volumens, das  $E'_k$  enthält und nennen es  $E_{k+1}$ .

- (6) Setze  $k := k + 1$  und gehe zu (2).

□

Das Prinzip des Verfahrens ist klar. Man zäunt das Polytop  $P$  durch ein Ellipsoid  $E_k$  ein. Ist das Zentrum von  $E_k$  nicht in  $P$ , so sucht man sich ein kleineres Ellipsoid  $E_{k+1}$ , das  $P$  enthält und fährt so fort. Die Probleme, die zu klären sind, sind die folgenden:

- Wie findet man ein Ellipsoid, das  $P$  enthält?
- Wie kann man  $E_{k+1}$  aus  $E_k$  konstruieren?
- Wann kann man abbrechen, d. h. was heißt “zu klein”?
- Wieviele Iterationen des Verfahrens sind durchzuführen?

Die nachfolgenden Sätze beantworten diese Fragen, wobei wir nur einige der Beweise, die zum Nachweis der Polynomialität des Verfahrens notwendig sind, angeben wollen.

Unser Anfangsellipsoid soll eine Kugel mit dem Nullpunkt als Zentrum sein. Enthalten die Restriktionen, die das Polytop  $P$  definieren, explizite obere und untere Schranken für die Variablen, sagen wir

$$l_i \leq x_i \leq u_i, \quad i = 1, \dots, n$$

so sei

$$(12.24) \quad R := \sqrt{\sum_{i=1}^n (\max\{|u_i|, |l_i|\})^2}.$$

Dann gilt trivialerweise  $P \subseteq S(0, R) = E(R^2 I, 0)$ . Andernfalls kann man zeigen

**(12.25) Lemma.** Sei  $P$  ein Polyeder der Form  $P(A, b)$ ,  $P^=(A, b)$  oder  $\{x \in \mathbb{R}^n \mid Ax \leq b, x \geq 0\}$  mit  $A \in \mathbb{Q}^{(m,n)}$ ,  $b \in \mathbb{Q}^m$ , dann gilt

(a) Alle Ecken von  $P$  sind in der Kugel  $S(0, R)$  enthalten mit

$$R := \sqrt{n} 2^{2\langle A \rangle + \langle b \rangle - 2n^2}.$$

(b) Ist  $P$  ein Polytop, so gilt  $P \subseteq S(0, R) = E(R^2 I, 0)$ .

**Beweis :** Nach Satz (12.11) gilt für jede Ecke  $v^T = (v_1, \dots, v_n)$  von  $P$

$$|v_i| \leq 2^{2\langle A \rangle + \langle b \rangle - 2n^2} \quad (i = 1, \dots, n),$$

und daraus folgt für die euklidische Norm von  $v$ :

$$\|v\| = \sqrt{\sum_{i=1}^n v_i^2} \leq \sqrt{n \max\{v_i^2\}} \leq \sqrt{n} 2^{2\langle A \rangle + \langle b \rangle - 2n^2}.$$

Also ist jede Ecke von  $P$  in  $S(0, R)$  enthalten. Ist insbesondere  $P$  ein Polytop, so folgt daraus  $P \subseteq S(0, R)$ .  $\square$

Damit haben wir durch (12.24) oder (12.25) ein Anfangsellipsoid  $E_0$  gefunden, mit dem wir die Ellipsoidmethode beginnen können. Die Konstruktion von  $E_{k+1}$  aus  $E_k$  geschieht wie folgt.

**(12.26) Satz.** Sei  $E_k = E(A_k, a_k) \subseteq \mathbb{R}^n$  ein Ellipsoid,  $c \in \mathbb{R}^n \setminus \{0\}$  und  $E'_k := \{x \in \mathbb{R}^n \mid c^T x \leq c^T a_k\} \cap E_k$ . Setze

$$(12.27) \quad d := \frac{1}{\sqrt{c^T A_k c}} A_k c,$$

$$(12.28) \quad a_{k+1} := a_k - \frac{1}{n+1} d,$$

$$(12.29) \quad A_{k+1} := \frac{n^2}{n^2-1} \left( A_k - \frac{2}{n+1} d d^T \right),$$

dann ist  $A_{k+1}$  positiv definit und  $E_{k+1} := E(A_{k+1}, a_{k+1})$  ist das eindeutig bestimmte Ellipsoid minimalen Volumens, das  $E'_k$  enthält.

**Beweis:** Siehe R. G. Bland, D. Goldfarb & M. J. Todd, “The Ellipsoid Method: A Survey”, Operations Research 29 (1981) 1039–1091 oder M. Grötschel, L. Lovász & A. Schrijver (1988)  $\square$

Wir wollen kurz den geometrischen Gehalt der Schritte in Satz (12.26) erläutern.

Ist  $c \in \mathbb{R}^n$ , und wollen wir das Maximum oder Minimum von  $c^T x$  über  $E_k$  finden, so kann man dies explizit wie folgt angeben. Für den in (12.27) definierten Vektor  $d$  und

$$z_{\max} := a_k + d, \quad z_{\min} := a_k - d$$

gilt

$$\begin{aligned} c^T z_{\max} &= \max\{c^T x \mid x \in E_k\} = c^T a_k + \sqrt{c^T A_k c}, \\ c^T z_{\min} &= \min\{c^T x \mid x \in E_k\} = c^T a_k - \sqrt{c^T A_k c}. \end{aligned}$$

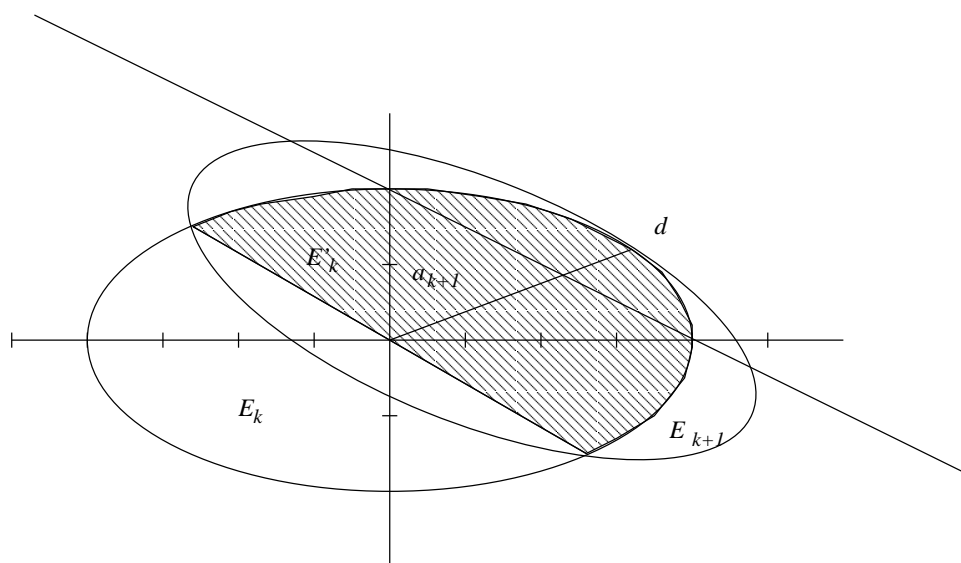
Diese Beziehung kann man leicht — z. B. aus der Cauchy-Schwarz-Ungleichung — ableiten. Daraus folgt, dass der Mittelpunkt  $a_{k+1}$  des neuen Ellipsoids  $E_{k+1}$  auf dem Geradenstück zwischen  $a_k$  und  $z_{\min}$  liegt. Die Länge dieses Geradenstücks ist  $\|d\|$ , und  $a_{k+1}$  erreicht man von  $a_k$  aus, indem man einen Schritt der Länge  $\frac{1}{n+1}\|d\|$  in Richtung  $-d$  macht. Der Durchschnitt des Randes des Ellipsoids  $E_{k+1}$  mit dem Rand von  $E'_k$  wird gebildet durch den Punkt  $z_{\min}$  und  $E''_k := \{x \mid (x - a_k)^T A_k (x - a_k) = 1\} \cap \{x \mid c^T x = c^T a_k\}$ .  $E''_k$  ist der Rand eines “ $(n-1)$ -dimensionalen Ellipsoids” im  $\mathbb{R}^n$ .

In Abbildung 12.2 sind das Ellipsoid  $E_k$  aus Abbildung 12.1 und das Ellipsoid  $E_{k+1}$ , definiert durch Satz (12.26) bezüglich des Vektors  $c^T = (-1, -2)$  dargestellt.  $E'_k$  ist gestrichelt eingezeichnet. Als Anfangsdaten haben wir daher: Die verletzte Ungleichung ist  $-x_1 - 2x_2 \leq -4$ . Die zugehörige Gerade  $\{x \mid -x_1 - 2x_2 = -4\}$  ist ebenfalls gezeichnet.

$$a_k = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad A_k = \begin{pmatrix} 16 & 0 \\ 0 & 4 \end{pmatrix}, \quad c = \begin{pmatrix} -1 \\ -2 \end{pmatrix}.$$

Die Formeln (12.27), (12.28), (12.29) ergeben:

$$\begin{aligned} d &= \frac{-1}{\sqrt{2}} \binom{4}{2}, \\ a_{k+1} &= a_k - \frac{1}{3}d = \frac{1}{3\sqrt{2}} \binom{4}{2} \approx \begin{pmatrix} 0.9428 \\ 0.4714 \end{pmatrix}, \\ A_{k+1} &= \frac{n^2}{n^2-1} \left( A_k - \frac{2}{n+1} dd^T \right) = \frac{4}{9} \begin{pmatrix} 32 & -8 \\ -8 & 8 \end{pmatrix}, \\ E_{k+1} &= E(A_{k+1}, a_{k+1}). \end{aligned}$$



**Abb. 12.2**

Das Stopkriterium der Ellipsoidmethode beruht auf einem Volumenargument. Nach Konstruktion ist klar, dass das Volumen von  $E_{k+1}$  (bezeichnet mit  $\text{vol}(E_{k+1})$ ) kleiner ist als das von  $E_k$ . Man kann den Volumenschrumpfungsfaktor explizit berechnen. Er hängt nur von der Dimension  $n$  des Raumes  $\mathbb{R}^n$  und nicht etwa von Update-Vektor  $c$  ab.

**(12.30) Lemma.**

$$\frac{\text{vol}(E_{k+1})}{\text{vol}(E_k)} = \left( \left( \frac{n}{n+1} \right)^{n+1} \left( \frac{n}{n-1} \right)^{n-1} \right)^{\frac{1}{2}} \leq e^{\frac{-1}{2n}} < 1.$$

**Beweis :** Siehe Grötschel, Lovász & Schrijver (1988), Lemma (3.1.34).  $\square$

Wir sehen also, dass mit den Formeln aus Satz (12.26) eine Folge von Ellipsoiden konstruiert werden kann, so dass jedes Ellipsoid  $E_{k+1}$  das Halbellipsoid  $E'_k$  und

somit das Polytop  $P$  enthält und dass die Volumina der Ellipsoide schrumpfen. Die Folge der Volumina konvergiert gegen Null, muss also einmal das Volumen von  $P$ , falls  $P$  ein positives Volumen hat, unterschreiten. Daher muss nach endlich vielen Schritten der Mittelpunkt eines der Ellipsoide in  $P$  sein, falls  $P \neq \emptyset$ . Wir wollen nun ausrechnen, nach wievielen Schritten, dies der Fall ist.

**(12.31) Lemma.** Seien  $P = P(A, b) \subseteq \mathbb{R}^n$ ,  $\overset{\circ}{P} = \{x \in \mathbb{R}^n \mid Ax < b\}$ , und  $R := \sqrt{n}2^{2\langle A \rangle + \langle b \rangle - 2n^2}$ .

(a) Entweder gilt  $\overset{\circ}{P} = \emptyset$  oder

$$\text{vol}(\overset{\circ}{P} \cap S(0, R)) \geq 2^{-(n+1)(\langle A \rangle + \langle b \rangle - n^2)}.$$

(b) Ist  $P$  volldimensional (d. h.  $\dim P = n$ ), dann gilt

$$\text{vol}(P \cap S(0, R)) = \text{vol}(\overset{\circ}{P} \cap S(0, R)) > 0.$$

□

Lemma (12.31) zusammen mit Lemma (12.25) sind geometrisch und algorithmisch interessant. Die beiden Hilfssätze implizieren folgendes.

- Wenn ein Polyeder  $P$  nicht leer ist, dann gibt es Elemente des Polyeders, die “nah” beim Nullvektor liegen, d. h. in  $S(0, R)$  enthalten sind.
- Es reicht aus zu überprüfen, ob  $P \cap S(0, R)$  leer ist oder nicht. Daraus kann man schließen, dass  $P$  leer ist oder nicht.
- Will man einen Konvergenzbeweis über Volumen führen, so sollte man strikte statt normale Ungleichungssysteme betrachten. Denn ein striktes Ungleichungssystem ist entweder unlösbar oder seine Lösungsmenge hat ein positives (relativ großes) Volumen.

Damit können wir die Ellipsoidmethode formulieren.

**(12.32) Die Ellipsoidmethode.**

**Input:**  $A \in \mathbb{Q}^{(m,n)}$ ,  $b \in \mathbb{Q}^m$ .

**Output:** Ein Vektor  $x \in \mathbb{Q}^n$  mit  $Ax \leq b$  oder die Feststellung, dass  $\{x \in \mathbb{R}^n \mid Ax < b\}$  leer ist.

**(1) Initialisierung:** Setze

$$\begin{aligned}
A_0 &:= R^2 I, \text{ mit } R = \sqrt{n} 2^{2\langle A \rangle + \langle b \rangle - 2n^2} \text{ (oder } R \text{ durch andere Vorinformationen} \\
&\quad \text{kleiner gewählt),} \\
a_0 &:= 0, \\
k &:= 0, \\
N &:= 2n((3n+1)\langle A \rangle + (2n+1)\langle b \rangle - n^3).
\end{aligned}$$

(Das Anfangsellipsoid ist  $E_0 := E(A_0, a_0)$ .)

**(2) Abbruchkriterium:**

(2.a) Gilt  $k = N$ , dann hat  $Ax < b$  keine Lösung, STOP!

(2.b) Gilt  $Aa_k \leq b$ , dann ist eine Lösung gefunden, STOP!

(2.c) Andernfalls sei  $c^T$  irgendeine Zeile von  $A$  derart, dass der Mittelpunkt  $a_k$  von  $E_k$  die entsprechende Ungleichung verletzt.

**(3) Update:** Setze

$$(3.a) \quad a_{k+1} := a_k - \frac{1}{n+1} \frac{1}{\sqrt{c^T A_k c}} A_k c$$

$$\begin{aligned}
(3.b) \quad A_{k+1} &:= \frac{n^2}{n^2-1} \left( A_k - \frac{2}{n+1} \frac{1}{c^T A_k c} A_k c c^T A_k^T \right) \\
(E_{k+1} &:= E(A_{k+1}, a_{k+1}) \text{ ist das neue Ellipsoid.)} \\
k &:= k + 1.
\end{aligned}$$

Gehe zu (2). □

**(12.33) Satz.** Die Ellipsoidmethode arbeitet korrekt.

**Beweis :** Gibt es ein  $k \leq N$ , so dass  $Aa_k \leq b$ , so ist offenbar ein Vektor aus  $P(A, b)$  gefunden. Bricht die Ellipsoidmethode in Schritt (2.a) ab, so müssen wir zeigen, dass  $\overset{\circ}{P} := \{x \mid Ax < b\}$  kein Element besitzt.

Angenommen  $\overset{\circ}{P} \neq \emptyset$ . Sei  $P'$  der Durchschnitt von  $\overset{\circ}{P}$  mit  $E_0$ . Dann gilt nach Lemma (12.31), dass das Volumen von  $P'$  mindestens  $2^{-(n+1)(\langle A \rangle + \langle b \rangle - n^2)}$  beträgt. Ist  $a_k$  nicht in  $P(A, b)$ ,  $0 \leq k < N$ , so wird in (2.c) eine verletzte Ungleichung gefunden, sagen wir  $c^T x \leq \gamma$ . Wegen  $\gamma < c^T a_k$  enthält das Halbellipsoid

$$E'_k := E_k \cap \{x \mid c^T x \leq c^T a_k\}$$

die Menge  $P'$ . Das durch die Formeln (3.a), (3.b) konstruierte neue Ellipsoid  $E_{k+1}$  ist nach Satz (12.26) das volumenmäßig kleinste Ellipsoid, das  $E'_k$  enthält. Wegen  $P' \subseteq E'_k$  gilt natürlich  $P' \subseteq E_{k+1}$ . Daraus folgt

$$P' \subseteq E_k, \quad 0 \leq k \leq N.$$

Das Volumen des Anfangsellipsoids  $E_0$  kann man wie folgt berechnen:

$$\text{vol}(E_0) = \sqrt{\det(R^2 I)} V_n = R^n V_n,$$

wobei  $V_n$  das Volumen der Einheitskugel im  $\mathbb{R}^n$  ist. Wir machen nun eine sehr grobe Abschätzung. Die Einheitskugel ist im Würfel  $W = \{x \in \mathbb{R}^n \mid |x_i| \leq 1, i = 1, \dots, n\}$  enthalten. Das Volumen von  $W$  ist offensichtlich  $2^n$ . Daraus folgt

$$\begin{aligned} \text{vol}(E_0) &< R^n 2^n = 2^{n(2\langle A \rangle + \langle b \rangle - 2n^2 + \log \sqrt{n+1})} \\ &< 2^{n(2\langle A \rangle + \langle b \rangle - n^2)}, \end{aligned}$$

In jedem Schritt der Ellipsoidmethode schrumpft nach (12.30) das Volumen um mindestens den Faktor  $e^{\frac{1}{2n}}$ . Aus der Abschätzung von  $\text{vol}(E_0)$  und der in (1) von (12.33) angegebenen Formel für  $N$  erhalten wir somit

$$\text{vol}(E_N) \leq e^{\frac{-N}{2n}} \text{vol}(E_0) < 2^{-(n+1)(\langle A \rangle + \langle b \rangle - n^2)}.$$

Also gilt  $\text{vol}(E_N) < \text{vol}(P')$  und  $P' \subseteq E_N$ . Dies ist ein Widerspruch. Hieraus folgt, dass  $P'$  und somit  $\{x \mid Ax < b\}$  leer sind, wenn die Ellipsoidmethode in (2.a) abbricht.  $\square$

Aus (12.31) (b) folgt nun unmittelbar

**(12.34) Folgerung.** *Ist  $P = P(A, b) \subseteq \mathbb{R}^n$  ein Polyeder, von dem wir wissen, dass es entweder volldimensional oder leer ist, dann findet die Ellipsoidmethode entweder einen Punkt in  $P$  oder beweist, dass  $P$  leer ist.*

$\square$

Nun kann man natürlich einem durch ein Ungleichungssystem gegebenen Polyeder nicht unmittelbar ansehen, ob es volldimensional ist oder nicht. Ferner können gerade diejenigen Polyeder, die durch Transformation linearer Programme entstehen (siehe (12.8), hier gilt immer  $c^T x = b^T y$ ), nicht volldimensional sein, so dass die Ellipsoidmethode zu keiner befriedigenden Antwort führt. Diesen Defekt kann man jedoch durch die folgende Beobachtung reparieren.

**(12.35) Satz.** *Seien  $A \in \mathbb{Q}^{(m,n)}$  und  $b \in \mathbb{Q}^m$ , dann hat das Ungleichungssystem*

$$Ax \leq b$$

hat genau dann eine Lösung, wenn das strikte Ungleichungssystem

$$Ax < b + 2^{-2\langle A \rangle - \langle b \rangle} \mathbb{1}$$

eine Lösung hat. Ferner kann man aus einer Lösung des strikten Ungleichungssystems in polynomialer Zeit eine Lösung von  $Ax \leq b$  konstruieren.  $\square$

Damit ist die Beschreibung der Ellipsoidmethode (bis auf die Abschätzung der Rechenzeit) vollständig. Wollen wir entscheiden, ob ein Polyeder  $P(A, b)$  einen Punkt enthält, können wir die Methode (12.32) auf das Ungleichungssystem  $Ax \leq b$  anwenden. Finden wir einen Punkt in  $P(A, b)$ , dann sind wir fertig. Andernfalls wissen wir, dass  $P(A, b)$  nicht volldimensional ist. In diesem Falle können wir auf Satz (12.35) zurückgreifen und starten die Ellipsoidmethode neu, und zwar z. B. mit dem ganzzahligen Ungleichungssystem

$$(12.36) \quad 2^{2\langle A \rangle + \langle b \rangle} Ax \leq 2^{2\langle A \rangle + \langle b \rangle} b + \mathbb{1}.$$

Entscheidet die Ellipsoidmethode, dass das zu (12.36) gehörige strikte Ungleichungssystem keine Lösung hat (Abbruch in Schritt (2.a)), so können wir aus (12.35) folgern, dass  $P(A, b)$  leer ist. Andernfalls findet die Ellipsoidmethode einen Vektor  $x'$ , der (12.36) erfüllt. Gilt  $x' \in P(A, b)$ , haben wir das gewünschte gefunden, falls nicht, kann man mit (einfachen Methoden der linearen Algebra) aus  $x'$  einen Punkt  $x \in P(A, b)$  konstruieren, siehe (12.35).

Zur Lösung linearer Programme kann man die in Abschnitt 12.2 besprochenen Reduktionen benutzen. Entweder man fasst das lineare Programm (12.6) und das dazu duale (12.7) zu (12.8) zusammen und sucht wie oben angegeben im nicht volldimensionalen Polyeder (12.8) einen Punkt, oder man wendet  $\bar{N}$ -mal, siehe (12.18), die Ellipsoidmethode für niederdimensionale Polyeder des Typs  $P_s$  aus (12.16) (2) im Rahmen eines binären Suchverfahrens (12.16) an.

In jedem Falle ist das Gesamtverfahren polynomial, wenn die Ellipsoidmethode (12.32) polynomial ist.

## 12.4 Laufzeit der Ellipsoidmethode

Wir wollen nun die Laufzeit der Ellipsoidmethode untersuchen und auf einige bisher verschwiegene Probleme bei der Ausführung von (12.32) aufmerksam machen. Offenbar ist die maximale Iterationszahl

$$N = 2n((3n + 1)\langle A \rangle + (2n + 1)\langle b \rangle - n^3)$$

polynomial in der Kodierungslänge von  $A$  und  $b$ . Also ist das Verfahren (12.32) genau dann polynomial, wenn jede Iteration in polynomialer Zeit ausgeführt werden kann.

Bei der Initialisierung (1) besteht kein Problem. Test (2.a) ist trivial, und die Schritte (2.b) und (2.c) führen wir dadurch aus, dass wir das Zentrum  $a_k$  in die Ungleichungen einsetzen und überprüfen, ob die Ungleichungen erfüllt sind oder nicht. Die Anzahl der hierzu benötigten elementaren Rechenschritte ist linear in  $\langle A \rangle$  und  $\langle b \rangle$ . Sie ist also polynomial, wenn die Kodierungslänge des Vektors  $a_{k+1}$  polynomial ist.

An dieser Stelle beginnen die Schwierigkeiten. In der Update-Formel (3.a) muss eine Wurzel berechnet werden. I. a. werden also hier irrationale Zahlen auftreten, die natürlich nicht exakt berechnet werden können. Die (möglicherweise) irrationale Zahl  $\sqrt{c^T A_k c}$  muss daher zu einer rationalen Zahl gerundet werden. Dadurch wird geometrisch bewirkt, dass der Mittelpunkt des Ellipsoids  $E_{k+1}$  ein wenig verschoben wird. Mit Sicherheit enthält das so verschobene Ellipsoid nicht mehr die Menge  $E'_k$  (siehe Satz (12.26)) und möglicherweise ist auch  $P(A, b)$  nicht mehr in diesem Ellipsoid enthalten. Also bricht unser gesamter Beweis der Korrektheit des Verfahrens zusammen.

Ferner wird beim Update (3.b) durch möglicherweise große Zahlen geteilt, und es ist nicht a priori klar, dass die Kodierungslänge der Elemente von  $A_{k+1}$  bei wiederholter Anwendung von (3.b) polynomial in  $\langle A \rangle + \langle b \rangle$  bleibt. Also müssen auch die Einträge in  $A_{k+1}$  gerundet werden. Dies kann zu folgenden Problemen führen. Die gerundete Matrix, sagen wir  $A_{k+1}^*$ , ist nicht mehr positiv definit und das Verfahren wird sinnlos. Oder  $A_{k+1}^*$  bleibt positiv definit, aber durch die Rundung hat sich die Form des zugehörigen Ellipsoids, sagen wir  $E_{k+1}^*$  so geändert, dass  $E_{k+1}^*$  das Polyeder  $P(A, b)$  nicht mehr enthält.

Alle diese Klippen kann man mit einem einfachen Trick umschiffen, dessen Korrektheitsbeweis allerdings recht aufwendig ist. Die geometrische Idee hinter diesem Trick ist die folgende. Man nehme die binäre Darstellung der Komponenten des in (3.a) berechneten Vektors und der in (3.b) berechneten Matrix und runde nach  $p$  Stellen hinter dem Binärkomma. Dadurch ändert man die Lage des Mittelpunkts und die Form des Ellipsoids ein wenig. Nun bläst man das Ellipsoid ein bisschen auf, d. h. man multipliziert  $A_{k+1}$  mit einem Faktor  $\zeta > 1$  und zwar so, dass die Menge  $E'_k$  beweisbar in dem aufgeblasenen Ellipsoid enthalten ist. Durch die Vergrößerung des Ellipsoids wird natürlich die in (12.30) bestimmte Schrumpfrate verschlechtert, was bedeutet, dass man insgesamt mehr Iterationen, sagen wir  $N'$ , durchführen muss. Daraus folgt, dass der Rundungsparameter  $p$  und der Aufblasparameter  $\zeta$  aufeinander so abgestimmt sein müssen, dass alle gerun-

deten und mit  $\zeta$  multiplizierten Matrizen  $A_k$ ,  $1 \leq k \leq N'$  und alle gerundeten Mittelpunkte  $a_k$ ,  $1 \leq k \leq N'$  polynomial in  $\langle A \rangle + \langle b \rangle$  berechnet werden können, so dass alle  $A_k$  positiv definit sind,  $P \cap S(0, R)$  in allen Ellipsoiden  $E_k$  enthalten ist und die Iterationszahl  $N'$  ebenfalls polynomial in  $\langle A \rangle + \langle b \rangle$  ist. Dies kann in der Tat durch Anwendung von Abschätzungstechniken aus der Analysis und linearen Algebra realisiert werden, siehe Grötschel, Lovász & Schrijver (1988).

Die Ellipsoidmethode (mit Rundungsmodifikation) ist also ein polynomialer Algorithmus, der entscheidet, ob ein Polyeder leer ist oder nicht. Mit Hilfe der im Vorhergehenden beschriebenen Reduktion kann sie dazu benutzt werden, lineare Programme in polynomialer Zeit zu lösen.

Wie sich der Leser denken kann, gibt es eine ganze Reihe von Varianten der Ellipsoidmethode, die dazu dienen sollen, schnellere Konvergenz zu erzwingen. Derartige Varianten sind z. B. in Bland, Goldfarb & Todd (1982), Grötschel, Lovász & Schrijver (1988) und M. Akgül, "Topics in Relaxation and Ellipsoidal Methods", Pitman, Boston, 1984, beschrieben. Aus Platz- und Zeitgründen können wir hier nicht weiter darauf eingehen.

## 12.5 Ein Beispiel

Wir wollen hier den Ablauf der Ellipsoidmethode anhand eines Beispiels im  $\mathbb{R}^2$  geometrisch veranschaulichen. Wir starten mit dem folgenden Polyeder  $P(A, b) \subseteq \mathbb{R}^2$  definiert durch

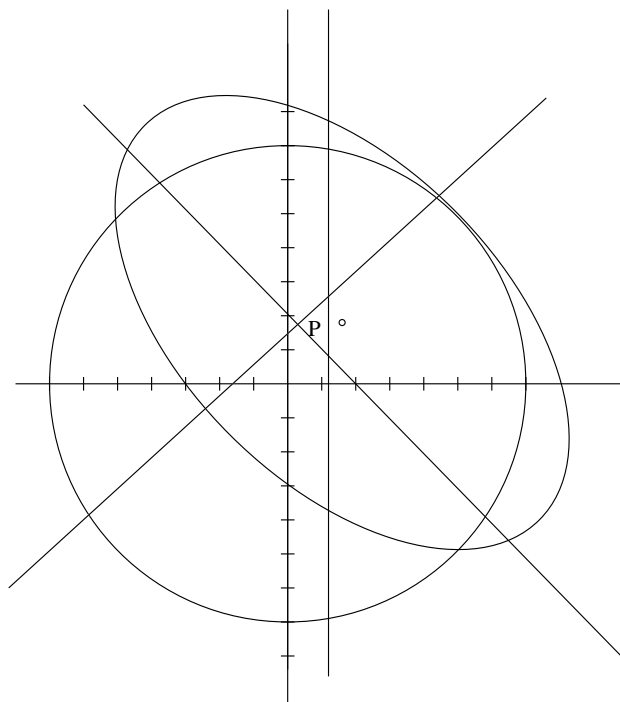
$$\begin{aligned} -x_1 - x_2 &\leq -2 \\ 3x_1 &\leq 4 \\ -2x_1 + 2x_2 &\leq 3 \end{aligned}$$

Dieses Polyeder  $P$  ist in der Kugel (Kreis) um den Nullpunkt mit Radius 7 enthalten. Diese Kugel soll unser Anfangsellipsoid  $E_0 = E(A_0, 0)$  sein. Wir führen mit dieser Initialisierung die Schritte (2) und (3) des Ellipsoidverfahrens (12.32) durch. Wir rechnen natürlich nicht mit der eigentlich erforderlichen Genauigkeit sondern benutzen die vorhandene Maschinenpräzision. In unserem Fall ist das Programm in PASCAL geschrieben, und die Rechnungen werden in REAL-Arithmetik durchgeführt, d. h. wir benutzen ein Mantisse von 3 byte und einen Exponenten von 1 byte zur Zahlendarstellung. Das Verfahren führt insgesamt 7 Iterationen durch und endet mit einem zulässigen Punkt. In den nachfolgenden Abbildungen 12.3, ..., 12.9 sind (im Maßstab 1:2) jeweils die Ellipsoide  $E_k$  und  $E_{k+1}$  mit ihren Mittelpunkten  $a_k$  und  $a_{k+1}$ ,  $k = 0, \dots, 6$  aufgezeichnet. Man sieht also, wie sich die Form und Lage eines Ellipsoids in einem Iterationsschritt ändert.

Das Startellipsoid ist gegeben durch

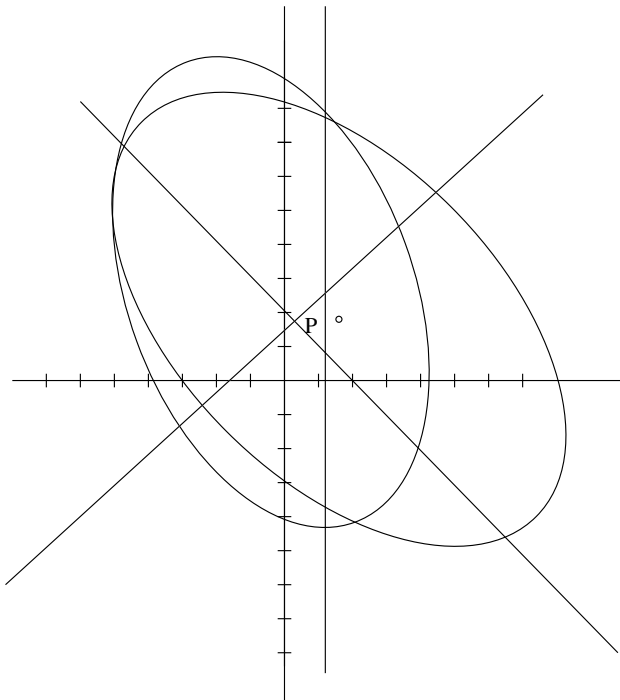
$$a_0^T = (0, 0), \quad A = \begin{pmatrix} 49 & 0 \\ 0 & 49 \end{pmatrix}.$$

Neben jeder Abbildung sind das neue Zentrum  $a_{k+1}$  und die neue Matrix  $A_{k+1}$  mit  $E_{k+1} = E(A_{k+1}, a_{k+1})$  angegeben. Jede Abbildung enthält die Begrenzungsgeraden des Polytops  $P$ , so dass man auf einfache Weise sehen kann, welche Ungleichungen durch den neuen Mittelpunkt verletzt sind. Die Abbildung 12.10 enthält eine ‘‘Gesamtschau’’ des Verfahrens. Hier sind alle während des Verfahrens generierten Ellipsoide (im Maßstab 1:1) gleichzeitig dargestellt. Man sieht hier recht gut, wie die Mittelpunkte ‘‘herumhüpfen’’ und auf den ersten Blick keine ‘‘Ordnung’’ in der Folge der Mittelpunkte zu erkennen ist.



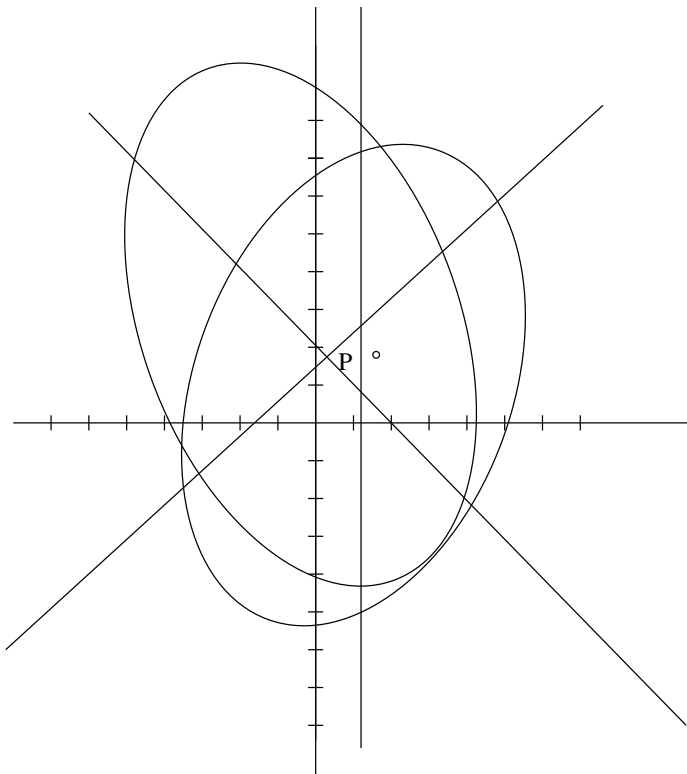
$$\begin{aligned} a_0 &= (0.0, 0.0) \\ A_0 &= \begin{pmatrix} 49.000 & 0.0 \\ 0.0 & 49.000 \end{pmatrix} \\ a_1 &= (1.6499, 1.6499) \\ A_1 &= \begin{pmatrix} 43.5555 & -21.7777 \\ -21.7777 & 43.5555 \end{pmatrix} \end{aligned}$$

**Abb. 12.3**



$$a_2 = (-0.5499, 2.7499)$$
$$A_2 = \begin{pmatrix} 19.3580 & -9.6790 \\ -9.6790 & 48.3951 \end{pmatrix}$$

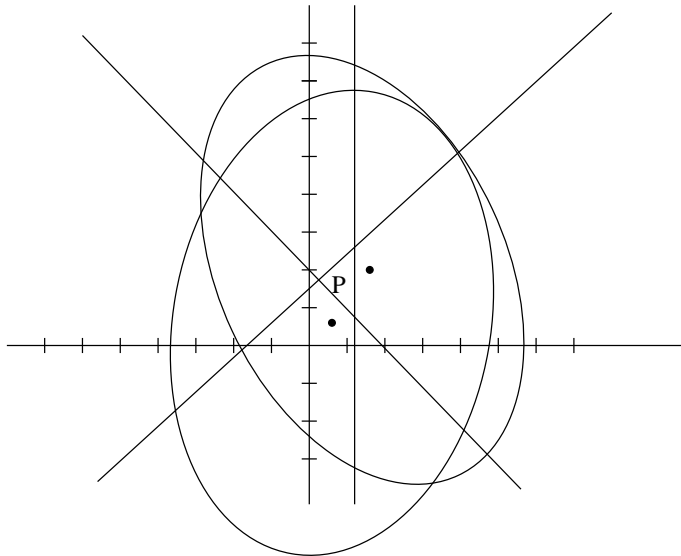
Abb. 12.4



$$a_3 = (0.4871, 0.6758)$$

$$A_3 = \begin{pmatrix} 17.2071 & 4.3018 \\ 4.3018 & 30.1125 \end{pmatrix}$$

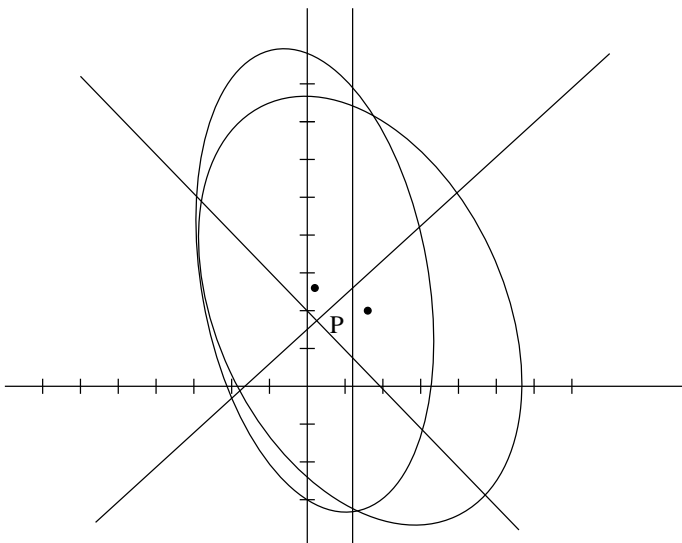
Abb. 12.5



$$a_4 = (1.4458, 2.2098)$$

$$A_4 = \begin{pmatrix} 15.5894 & -6.0299 \\ -6.0299 & 21.351 \end{pmatrix}$$

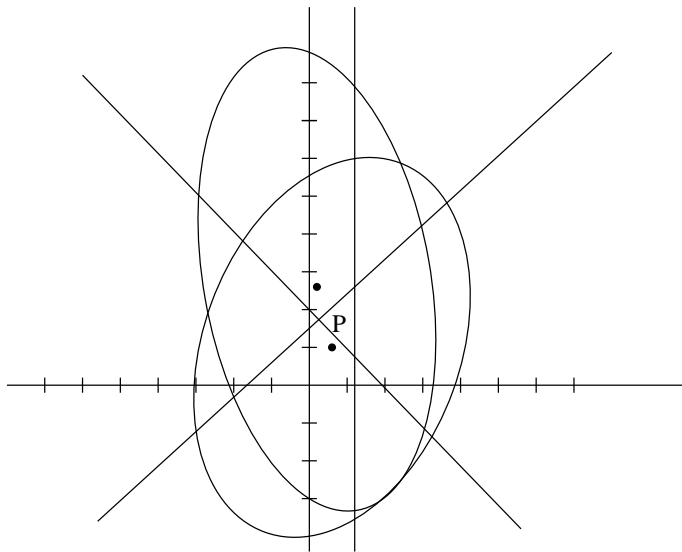
Abb. 12.6



$$a_5 = (0.1297, 2.7188)$$

$$A_5 = \begin{pmatrix} 6.9286 & -2.6799 \\ -2.6799 & 26.3603 \end{pmatrix}$$

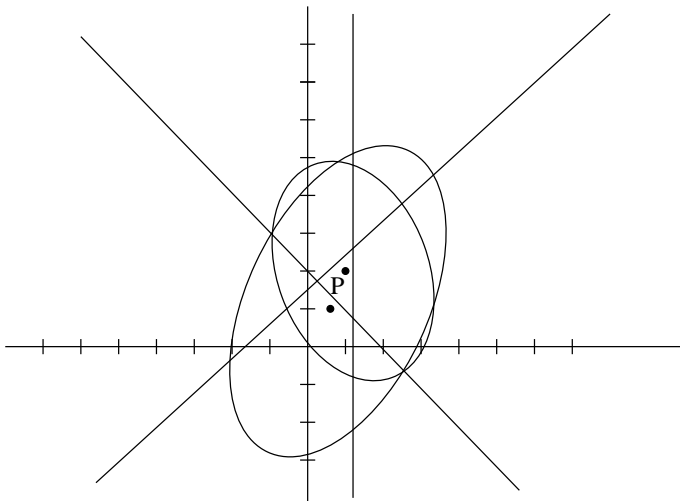
Abb. 12.7



$$a_6 = (0.6449, 1-1618)$$

$$A_6 = \begin{pmatrix} 7.1148 & 2.8443 \\ 2.8443 & 15.7511 \end{pmatrix}$$

Abb. 12.8



$$a_7 = (1.2661, 2.317)$$

$$A_7 = \begin{pmatrix} 6.3988 & -1.9726 \\ -1.9726 & 10.2372 \end{pmatrix}$$

Abb. 12.9



# Kapitel 13

## Der Karmarkar-Algorithmus und Innere-Punkte-Methoden

Im Jahre 1984 hat N. Karmarkar (AT&T Bell Laboratories) einen neuen Algorithmus zur Lösung linearer Programme entwickelt. Karmarkar's Aufsatz "A New Polynomial Time Algorithm for Linear Programming" ist in der Zeitschrift *Combinatorica*, Vol. 4, Nr. 4, 1984, 373–395, veröffentlicht. Im Herbst 1984 hat Karmarkar auf Vorträgen mitgeteilt, dass sein Verfahren nicht nur theoretisch polynomial ist (wie auch die in Kapitel 12 beschriebene Ellipsoidmethode), sondern in der Praxis (Implementation bei AT& Bell Laboratories) erheblich rascher als das Simplexverfahren arbeitet. Karmarkar behauptete, dass die Implementation seines Algorithmus etwa 50-mal schneller ist als MPSX, ein von der Firma IBM angebotenes und auf dem Simplexalgorithmus basierendes Softwarepaket zur Lösung linearer Programme, das zur damaligen Zeit als eines der besten LP-Pakete galt.

Da bei großen Industriefirmen sehr große LP's täglich in nicht unerheblicher Zahl gelöst werden, war die Ankündigung Karmarkars auf sehr großes Interesse bei Praktikern gestoßen. Ja, sein Verfahren hat sogar Aufsehen in der Presse verursacht. Zum Beispiel haben New York Times, Science, Time, Der Spiegel (falsch — wie bei Artikeln über Mathematik üblich), Die Zeit, Süddeutsche Zeitung (sehr ordentlich) über Karmarkars Verfahren berichtet, wie immer natürlich unter dem Hauptaspekt, dass dadurch Millionen von Dollars gespart werden können.

Im Nachfolgenden wollen wir die Grundidee des Algorithmus von Karmarkar darstellen. Neu- und Uminterpretation dieser Idee haben zu einer stürmischen Entwicklung geführt, bei der insbesondere in den 80er und 90er Jahren hunderte Artikel geschrieben wurden. Die hierbei entworfenen Algorithmen, genannt „Interior Point“ oder „Barrier“ Methoden haben zu Codes geführt, die dem Simplex-

Algorithmus ebenbürtig sind. In der Vorlesung habe ich über das prinzipielle Vorgehen der Innere-Punkte-Verfahren, ihre Vorzüge und Nachteile einen Überblick gegeben, den ich jedoch nicht schriftlich ausgearbeitet habe. Ich verweise hierzu auf die Bücher:

Roos, Terlaky, and Vial (2006), Interior point methods for linear optimization (2nd edition), Springer, und Wright (1997), Primal-dual interior-point methods, SIAM.

## 13.1 13.1 Reduktionen

Wie der Simplexalgorithmus und die Ellipsoidmethode, so ist auch der Karmarkar-Algorithmus auf einen speziellen Problemtyp zugeschnitten und nicht auf allgemeine lineare Programme anwendbar. Wie üblich müssen wir daher zunächst zeigen, dass ein allgemeines LP so transformiert werden kann, dass es mit Karmarkars Algorithmus gelöst werden kann. Die Grundversion des Karmarkar-Algorithmus (und nur die wollen wir hier beschreiben) ist ein Verfahren zur Entscheidung, ob ein Polyeder einer recht speziellen Form einen Punkt enthält oder nicht. Das Problem, das Karmarkar betrachtet ist das folgende.

**(13.1) Problem.** Gegeben seien eine Matrix  $A \in \mathbb{Q}^{(m,n)}$  und ein Vektor  $c \in \mathbb{Q}^n$ , entscheide, ob das System

$$Ax = 0, \quad \mathbf{1}^T x = 1, \quad x \geq 0, \quad c^T x \leq 0$$

eine Lösung hat, und falls das so ist, finde eine. □

Um den Karmarkar-Algorithmus starten zu können, ist die Kenntnis eines zulässigen Startvektors im relativ Inneren von  $\{x \mid Ax = 0, \mathbf{1}^T x = 1, x \geq 0\}$  notwendig. Wir wollen sogar noch mehr fordern, und zwar soll gelten

$$(13.2) \quad \bar{x} := \frac{1}{n} \mathbf{1} \quad \text{erfüllt} \quad A\bar{x} = 0.$$

Trivialerweise erfüllt  $\bar{x}$  die Bedingungen  $\mathbf{1}^T \bar{x} = 1, \bar{x} \geq 0$ . Gilt  $c^T \bar{x} \leq 0$ , ist man natürlich fertig. Die eigentliche Aufgabe besteht also darin, aus  $\bar{x}$  einen Punkt zu konstruieren, der alle Bedingungen von (13.1) erfüllt.

Wir werden nun, wie bei der Darstellung der Ellipsoidmethode, ein allgemeines LP in ein (bzw. mehrere) Probleme des Typs (13.1) umformen.

Wir wissen bereits, dass jedes LP in ein LP in Standardform (9.1) transformiert

werden kann. Gegeben sei also das folgende Problem

$$(13.3) \quad \begin{aligned} \max w^T x \\ Bx &= b \\ x &\geq 0 \end{aligned}$$

mit  $w \in \mathbb{Q}^n$ ,  $B \in \mathbb{Q}^{(m,n)}$ ,  $b \in \mathbb{Q}^m$ . Wie in Abschnitt 12.2 gezeigt gibt es zwei prinzipielle Reduktionsmöglichkeiten. Man fasst (13.3) und das dazu duale Problem wie in (12.8) zusammen, oder man führt die in (12.16) beschriebene Binärsuche durch. Wir stellen hier kurz die Reduktion über Binärsuche dar. Daraus ergibt sich auch, wie man die Kombination des primalen mit dem dualen Problem behandeln kann.

Genau wie in (12.16) angegeben (und mit den in den davor angegebenen Sätzen vorgelegten Begründungen) führen wir eine binäre Suche durch über der Menge

$$S := \left\{ \frac{p}{q} \in \mathbb{Q} \mid |p| \leq n2^{\langle B \rangle + \langle b \rangle + 2\langle w \rangle - n^2 - n}, 1 \leq q \leq 2^{\langle B \rangle + \langle w \rangle - n^2 - n} \right\}$$

der möglichen optimalen Zielfunktionswerte (falls (13.3) eine endliche Optimallösung hat). Zur Bestimmung einer optimalen Lösung von (13.3) sind nach (12.18) höchstens

$$\bar{N} := 2\langle B \rangle + \langle b \rangle + 3\langle c \rangle - 2n^2 - 2n + \log_2 n + 2$$

Aufrufe eines Algorithmus nötig, der für ein  $s \in S$  entscheidet, ob

$$(13.4) \quad \begin{aligned} w^T x &\leq s \\ Bx &= b \\ x &\geq 0 \end{aligned}$$

eine Lösung hat. Genau dann, wenn wir zeigen können, dass (13.4) in polynomialer Zeit gelöst werden kann, haben wir also ein polynomiales Verfahren zur Lösung von (13.3). Wir wissen aus Satz (12.11), dass wir alle Variablen durch  $2^{2\langle B \rangle + \langle b \rangle - 2n^2}$  beschränken können. Führen wir für die Ungleichung in (13.4) eine Schlupfvariable  $x_{n+1}$  ein, so gilt  $0 \leq x_{n+1} \leq |S| \leq n2^{\langle B \rangle + \langle b \rangle + 2\langle w \rangle - n^2 - n}$ . Setzen wir also

$$M := n(2^{2\langle B \rangle + \langle b \rangle - 2n^2} + 2^{\langle B \rangle + \langle b \rangle + 2\langle w \rangle - n^2 - n}),$$

so ist (13.4) genau dann lösbar, wenn

$$(13.5) \quad \begin{aligned} \begin{pmatrix} w^T & 1 \\ B & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_{n+1} \end{pmatrix} &= \begin{pmatrix} s \\ b \end{pmatrix} \\ \sum_{i=1}^{n+1} x_i &\leq M \\ x_i &\geq 0, \quad i = 1, \dots, n+1 \end{aligned}$$

lösbar ist. Führen wir eine weitere Schlupfvariable  $x_{n+2}$  ein (d. h. es gilt nun  $x \in \mathbb{R}^{n+2}$ ) und skalieren wir die Variablen um, so ist (13.5) genau dann lösbar, wenn

$$(13.6) \quad \begin{aligned} Cx &:= \begin{pmatrix} w^T & 1 & 0 \\ B & 0 & 0 \end{pmatrix} x = \frac{1}{M} \begin{pmatrix} s \\ b \end{pmatrix} =: \begin{pmatrix} \bar{c}_0 \\ \vdots \\ \bar{c}_m \end{pmatrix} \\ \mathbb{1}^T x &= 1 \\ x_i &\geq 0, \quad i = 1, \dots, n+2 \end{aligned}$$

lösbar ist. Von der  $i$ -ten Zeile der Matrix  $C$  ( $i = 0, \dots, m$ ) ziehen wir nun das  $\bar{c}_i$ -fache der letzten Gleichung  $\mathbb{1}^T x = 1$  ab. Dadurch machen wir die ersten  $m+1$  Gleichungen homogen und erreichen, dass aus dem Ungleichungs- und Gleichungssystem (13.6) ein äquivalentes System der folgenden Form wird:

$$(13.7) \quad \begin{aligned} Dx &= 0 \\ \mathbb{1}^T x &= 1 \\ x &\geq 0 \end{aligned}$$

Um die Voraussetzung (13.2) herzustellen, führen wir eine weitere Schlupfvariable  $x_{n+3}$  wie folgt ein. Wir betrachten

$$(13.8) \quad \begin{aligned} Dx - D\mathbb{1}x_{n+3} &= 0 \\ \mathbb{1}^T x + x_{n+3} &= 1 \\ x_i &\geq 0, \quad i = 1, \dots, n+3. \end{aligned}$$

Offenbar hat (13.7) genau dann eine Lösung, wenn (13.8) eine Lösung mit  $x_{n+3} = 0$  hat. Setzen wir

$$A := (D, -D\mathbb{1}) \in \mathbb{Q}^{(m+1, n+3)}, \quad c := (0, \dots, 0, 1)^T \in \mathbb{Q}^{n+3}$$

und betrachten wir  $x$  als Vektor im  $\mathbb{R}^{n+3}$ , so hat also (13.7) genau dann eine Lösung, wenn

$$(13.9) \quad Ax = 0, \quad \mathbb{1}^T x = 1, \quad x \geq 0, \quad c^T x \leq 0$$

eine Lösung hat. Ferner erfüllt nach Konstruktion der Vektor

$$\bar{x}^T := \left( \frac{1}{n+3}, \dots, \frac{1}{n+3} \right) \in \mathbb{Q}^{n+3}$$

bezüglich des Systems (13.9) die Forderung (13.2).

Folglich kann jedes lineare Program durch polynomial viele Aufrufe eines Algorithmus zur Lösung von (13.1) (mit Zusatzvoraussetzung (13.2)) gelöst werden.

## 13.2 Die Grundversion des Karmarkar-Algorithmus

Wir wollen nun das Karmarkar-Verfahren zur Lösung von (13.1) unter der Zusatzvoraussetzung (13.2) beschreiben. Wie wir am Ende des letzten Abschnittes gesehen haben, können wir jedes LP in eine Folge von Problemen der Form (13.1) transformieren, und wir können sogar das Zentrum des Simplex  $\{x \in \mathbb{R}^n \mid \mathbb{1}^T x = 1, x \geq 0\}$  als Startpunkt wählen. Im weiteren betrachten wir also das folgende

**(13.10) Problem.** Gegeben seien eine Matrix  $A \in \mathbb{Q}^{(m,n)}$  und ein Vektor  $c \in \mathbb{Q}^n$ . Es seien

$$\begin{aligned} E &:= \{x \in \mathbb{R}^n \mid \mathbb{1}^T x = 1\} && \text{(Hyperebene)} \\ \Sigma &:= E \cap \mathbb{R}_+^n && \text{(Simplex)} \\ \Omega &:= \{x \in \mathbb{R}^n \mid Ax = 0\} && \text{(linearer Raum)}. \end{aligned}$$

Ferner wird vorausgesetzt, dass der Vektor  $\bar{x} := \frac{1}{n}\mathbb{1}$  in  $\Sigma \cap \Omega$  enthalten ist. Gesucht ist ein Vektor  $x \in \mathbb{Q}^n$  mit

$$x \in P := \Sigma \cap \Omega, \quad c^T x \leq 0.$$

□

Problem (13.10) kann sicherlich dadurch gelöst werden, daß eine Optimallösung des folgenden Programms bestimmt wird.

$$(13.11) \quad \begin{array}{ll} \min c^T x & \\ Ax = 0 & \text{bzw.} \quad \min c^T x \\ \mathbb{1}^T x = 1 & x \in \Omega \cap E \cap \mathbb{R}_+^n \\ x \geq 0 & \end{array}$$

Offenbar ist der optimale Zielfunktionswert von (13.11) genau dann größer als Null, wenn (13.10) keine Lösung hat. Unser Ziel wird nun sein (13.11) statt (13.10) zu lösen, wobei wir voraussetzen, dass  $\frac{1}{n}\mathbb{1}$  für (13.11) zulässig ist.

Es erscheint natürlich reichlich umständlich, ein lineares Programm, sagen wir der Form (13.3), durch Binärsuche auf eine Folge von Zulässigkeitsproblemen (13.10) zu reduzieren, die dann wieder durch spezielle lineare Programme der Form (13.11) gelöst werden. Diesen Umweg kann man durch die sogenannte Sliding Objective Function Technique begründen. Die dabei notwendigen zusätzlichen Überlegungen tragen jedoch nicht zu einem besseren Verständnis der Grundversion des Verfahrens bei, um die es hier geht. Ziel dieses Kapitels ist nicht die Darstellung einer besonders effizienten Version des Karmarkar-Algorithmus sondern dessen prinzipielle Idee.

## Geometrische Beschreibung eines Iterationsschrittes

Zur Lösung von (13.11) hat Karmarkar ein Verfahren entworfen, das wie fast alle Optimierungsverfahren (insbesondere die Methoden der nichtlinearen Optimierung) auf der folgenden simplen Idee beruht. Angenommen man befinde sich an einem zulässigen Punkt, sagen wir  $x^k$ , dann sucht man eine Richtung (also einen Vektor  $d \in \mathbb{R}^n$ ), bezüglich der die Zielfunktion verbessert werden kann. Daraufhin bestimmt man eine Schrittweite  $\rho$ , so dass man von  $x^k$  zum nächsten Punkt  $x^{k+1} := x^k + \rho d$  gelangt. Der nächste Punkt  $x^{k+1}$  soll natürlich auch zulässig sein und einen “wesentlich” besseren Zielfunktionswert haben.

Die Essenz eines jeden solchen Verfahrens steckt natürlich in der Wahl der Richtung und der Schrittweite. Bei derartigen Verfahren tritt häufig die folgende Situation ein. Man ist in der Lage, eine sehr gute Richtung zu bestimmen (d. h. die Zielfunktion wird in Richtung  $d$  stark verbessert), aber man kann in Richtung  $d$  nur einen sehr kleinen Schritt ausführen, wenn man die zulässige Menge nicht verlassen will. Trotz guter Richtung kommt man also im Bezug auf eine tatsächliche Verbesserung kaum vorwärts und erhält u. U. global schlechtes Konvergenzverhalten. Man muss sich also bemühen, einen guten Kompromiss zwischen “Qualität der Richtung” und “mögliche Schrittweite” zu finden, um insgesamt gute Fortschritte zu machen.

Ist man — wie im vorliegenden Fall — im relativen Inneren der Menge  $P$ , aber nah am Rand und geht man z. B. in Richtung des Normalenvektors der Zielfunktion, so kann man sehr schnell an den Rand von  $P$  gelangen, ohne wirklich weiter gekommen zu sein, siehe Abbildung 13.1.

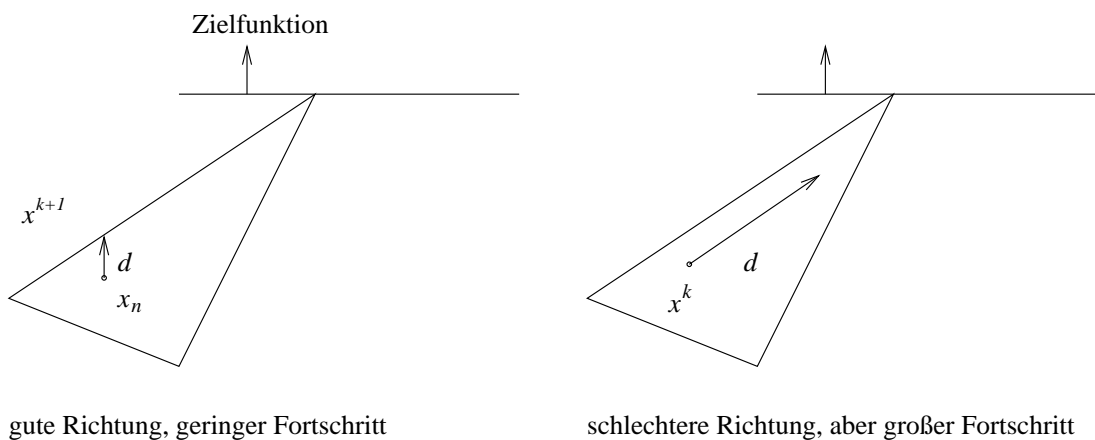


Abb. 13.1

Karmarkars Idee zur Lösung bzw. Umgehung dieser Schwierigkeit ist die folgende. Er führt eine projektive Transformation aus, die den Simplex  $\Sigma$  auf sich selbst, den affinen Teilraum  $\Omega$  auf einen anderen affinen Teilraum  $\Omega'$  abbildet und den relativ inneren Punkt  $x^k$  auf das Zentrum  $\frac{1}{n}\mathbb{1}$  von  $\Sigma$  wirft.  $P$  wird dabei auf ein neues Polyeder  $P_k$  abgebildet. Offenbar kann man von  $\frac{1}{n}\mathbb{1}$  aus recht große Schritte in alle zulässigen Richtungen machen, ohne sofort den zulässigen Bereich  $P_k$  zu verlassen. Man bestimmt so durch Festlegung einer Richtung und einer Schrittlänge von  $\frac{1}{n}\mathbb{1}$  ausgehend einen Punkt  $y^{k+1} \in P_k$  und transformiert diesen zurück, um den nächsten (zulässigen) Iterationspunkt  $x^{k+1} \in P$  zu erhalten.

Zur Bestimmung der Richtung macht Karmarkar folgende Überlegung. Ideal wäre es natürlich, direkt das Optimum des transformierten Problems zu bestimmen. Aber die lineare Zielfunktion des ursprünglichen Problems wird durch die projektive Transformation in eine nichtlineare Abbildung übergeführt, so dass dies nicht so einfach zu bewerkstelligen ist. Dennoch kann man diese transformierte Zielfunktion auf natürliche Weise linearisieren. Mit  $\bar{c}^T x$  sei die neue (linearisierte) Zielfunktion bezeichnet. Man hat nun ein neues Problem: Es soll eine lineare Zielfunktion über dem Durchschnitt eines Simplex mit einem affinen Raum optimiert werden, wir erhalten

$$(13.12) \quad \begin{aligned} \min \bar{c}^T x \\ x \in \Omega' \cap E \cap \mathbb{R}_+^n \end{aligned}$$

Dieses Problem ist offenbar wiederum vom Typ unseres Ausgangsproblems (13.11). Aber diese neue Aufgabe ist ja nur ein Hilfsproblem! Möglicherweise ist daher eine gute Approximation der Optimallösung von (13.12) ausreichend für das, was wir bezwecken. Durch die Lösung einer Vereinfachung dieses Problems (13.12) könnten u. U. ein Richtungsvektor und eine Schrittlänge gefunden werden, die von hinreichend guter Qualität bezüglich globaler Konvergenz des Verfahrens sind. Die Vereinfachung, die wir betrachten wollen, geschieht dadurch, dass die bei der Durchschnittsbildung beteiligte Menge  $\mathbb{R}_+^n$  durch die größte Kugel, sagen wir  $K$ , mit Zentrum  $\frac{1}{n}\mathbb{1}$ , die im Simplex  $\Sigma$  enthalten ist ersetzt wird. Statt (13.12) wird also die Aufgabe

$$(13.13) \quad \begin{aligned} \min \bar{c}^T x \\ x \in \Omega' \cap E \cap K \end{aligned}$$

betrachtet. Man optimiere eine lineare Zielfunktion über dem Durchschnitt einer Kugel  $K$  mit einem affinen Raum  $\Omega' \cap E$ . Dieses Problem ist ohne Einschaltung eines Algorithmus durch eine explizite Formel lösbar. Durch die Optimallösung von (13.13) werden die gesuchte Richtung und die Schrittlänge explizit geliefert.

Eine Modifikation ist jedoch noch nötig. Das Optimum über  $\Omega' \cap E \cap K$  ist i. a. irrational und muss gerundet werden. Dadurch ist es möglich, dass  $y^{k+1}$  nicht mehr im relativen Inneren von  $P_k$  bzw. der nächste (gerundete) Iterationspunkt  $x^{k+1}$  nicht mehr im relativ Inneren von  $P$  liegt. Statt  $K$  wählt man daher eine kleinere Kugel  $K'$  mit dem gleichen Zentrum, so dass auch nach Rundung und Rücktransformation garantiert ist, dass der nächste Iterationspunkt ein relativ innerer Punkt ist.

## Analytische Beschreibung eines Iterationsschrittes

Die oben gegebene anschauliche Darstellung wollen wir nun explizit vorführen. Wir starten also mit Problem (13.11). Wir setzen (wie beim Simplexverfahren) voraus, dass  $A$  vollen Zeilenrang hat. Außerdem kennen wir mit  $x^k$  einen relativ inneren Punkt von  $P$ . Ist  $D = \text{diag}(x^k)$  die  $(n, n)$ -Diagonalmatrix, die die Komponenten  $x_1^k, \dots, x_n^k$  von  $x^k$  auf der Hauptdiagonalen enthält, so ist durch

$$(13.14) \quad T_k(x) := \frac{1}{\mathbb{1}^T D^{-1} x} D^{-1} x$$

eine projektive Transformation definiert. ( $T_k(x)$  ist natürlich nur dann endlich, wenn  $x \notin \bar{H} = \{y \in \mathbb{R}^n \mid \mathbb{1}^T D^{-1} y = 0\}$  gilt. Wir werden beim Rechnen mit  $T_k$  diese Einschränkung nicht weiter erwähnen und gehen davon aus, dass klar ist, wie die Formeln, die  $T_k$  enthalten, zu interpretieren sind.) Die projektive Transformation  $T_k$  hat, wie leicht zu sehen ist, folgende Eigenschaften:

### (13.15) Eigenschaften von $T_k$ .

- (a)  $x \geq 0 \implies T_k(x) \geq 0$ .
- (b)  $\mathbb{1}^T x = 1 \implies \mathbb{1}^T T_k(x) = 1$ .
- (c)  $T_k^{-1}(y) = \frac{1}{\mathbb{1}^T D y} D y$  für alle  $y \in \Sigma$ .
- (d)  $T_k(\Sigma) = \Sigma$ .
- (e)  $P_k := T_k(P) = T_k(\Sigma \cap \Omega) = \Sigma \cap \Omega_k$  wobei  $\Omega_k = \{y \in \mathbb{R}^n \mid A D y = 0\}$ .
- (f)  $T_k(x^k) = \frac{1}{\mathbb{1}^T \mathbb{1}} D^{-1} x^k = \frac{1}{n} \mathbb{1} \in P_k$ .

□

Die Transformation  $T_k$  ist also eine bijektive Abbildung  $P \rightarrow P_k$  und  $\Sigma \rightarrow \Sigma$ , die den Punkt  $x^k$  in das Zentrum  $\frac{1}{n}\mathbb{1}$  von  $\Sigma$  abbildet. Aus (13.15) folgt

$$(13.16) \quad \begin{array}{llll} \min c^T x & = & \min c^T x & = & \min c^T T_k^{-1}(y) & = & \min \frac{1}{\mathbb{1}^T D y} c^T D y \\ Ax = 0 & & x \in P & & y \in P_k = T_k(P) & & ADy = 0 \\ \mathbb{1}^T x = 1 & & & & & & \mathbb{1}^T y = 1 \\ x \geq 0 & & & & & & y \geq 0 \end{array}$$

Das letzte Programm in (13.16) ist das in der geometrisch-anschaulichen Erläuterung weiter oben genannte Programm mit nichtlinearer Zielfunktion  $\frac{1}{\mathbb{1}^T D y} D y$ . Wir linearisieren diese durch Weglassen des Nenners wie folgt:

$$(13.17) \quad \bar{c}^T := c^T D$$

und erhalten das lineare (Hilfs)-Programm

$$(13.18) \quad \begin{array}{ll} \min \bar{c}^T y & \\ ADy = 0 & \\ \mathbb{1}^T y = 1 & \\ y \geq 0 & \end{array}$$

Wir vereinfachen (13.18) dadurch, dass wir die Nichtnegativitätsbedingung in (13.18) durch die Bedingung ersetzen, dass  $y$  in einer Kugel um  $\frac{1}{n}\mathbb{1}$  liegt. Die größte Kugel mit diesem Zentrum, die man dem Simplex  $\Sigma$  einbeschreiben kann, hat den Radius  $\frac{1}{\sqrt{n(n-1)}}$ . Wie bereits erwähnt, müssen wir aus technischen Gründen eine kleinere Kugel wählen. Es wird sich zeigen, dass man mit der Hälfte dieses optimalen Radius auskommt. Wir betrachten also das Programm

$$(13.19) \quad \begin{array}{ll} \min \bar{c}^T y & \\ ADy = 0 & \\ \mathbb{1}^T y = 1 & \\ \|y - \frac{1}{n}\mathbb{1}\| \leq \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} & \end{array}$$

Das Minimum von (13.19) kann explizit bestimmt werden.

**(13.20) Lemma.** Die Optimallösung von (13.19) ist der Vektor

$$y^{k+1} := \frac{1}{n}\mathbb{1} - \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} \frac{(I - DA^T(AD^2A^T)^{-1}AD - \frac{1}{n}\mathbb{1}\mathbb{1}^T)Dc}{\|(I - DA^T(AD^2A^T)^{-1}AD - \frac{1}{n}\mathbb{1}\mathbb{1}^T)Dc\|}.$$

**Beweis:** Zunächst projizieren wir  $\bar{c}$  orthogonal auf den linearen Raum  $L := \{x \in \mathbb{R}^n \mid ADx = 0, \mathbb{1}^T x = 0\}$ . Setzen wir

$$B = \begin{pmatrix} AD \\ \mathbb{1}^T \end{pmatrix},$$

so ist die Projektion von  $\bar{c}$  auf  $L$  gegeben durch

$$\bar{\bar{c}} = (I - B^T(BB^T)^{-1}B)\bar{c},$$

denn offenbar gilt  $B\bar{c} = 0$ ,  $(\bar{c} - \bar{\bar{c}})^T \bar{c} = 0$ . Beachten wir, dass  $AD\mathbb{1} = Ax^k = 0$  gilt, so folgt durch einfaches Ausrechnen

$$\begin{aligned} BB^T &= \begin{pmatrix} AD^2 A^T & AD\mathbb{1} \\ (AD\mathbb{1})^T & \mathbb{1}^T \mathbb{1} \end{pmatrix} = \begin{pmatrix} AD^2 A^T & 0 \\ 0 & n \end{pmatrix} \\ (BB^T)^{-1} &= \begin{pmatrix} (AD^2 A^T)^{-1} & 0 \\ 0 & \frac{1}{n} \end{pmatrix} \\ B^T(BB^T)^{-1}B &= DA^T(AD^2 A^T)^{-1}AD + \frac{1}{n}\mathbb{1}\mathbb{1}^T, \end{aligned}$$

und somit

$$\bar{\bar{c}} = (I - DA^T(AD^2 A^T)^{-1}AD - \frac{1}{n}\mathbb{1}\mathbb{1}^T)\bar{c}.$$

Für  $y \in L$  gilt nach Konstruktion  $\bar{c}^T y = \bar{\bar{c}}^T y$  und somit ist das Minimum von (13.19) gleich dem Minimum der Funktion  $\bar{\bar{c}}^T y$  unter den Nebenbedingungen von (13.19). Aufgrund unserer Konstruktion ist dieses Minimum gleich dem Minimum von

$$\begin{aligned} \min \bar{\bar{c}}^T y \\ \|y - \frac{1}{n}\mathbb{1}\| \leq \frac{1}{2} \frac{1}{\sqrt{n(n-1)}}, \end{aligned}$$

also dem Minimum einer linearen Zielfunktion über einer Kugel. Offenbar wird das Minimum hier durch den Vektor angenommen, den man durch einen Schritt vom Mittelpunkt  $\frac{1}{n}\mathbb{1}$  aus in Richtung  $-\bar{\bar{c}}$  mit der Länge des Kugelradius erhält, also durch

$$y^{k+1} = \frac{1}{n}\mathbb{1} - \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} \frac{1}{\|\bar{\bar{c}}\|} \bar{\bar{c}}.$$

Hieraus folgt die Behauptung. □

Damit haben wir ein Minimum von (13.19) analytisch bestimmen können und unser vereinfachtes Hilfsproblem gelöst. Den nächsten Iterationspunkt erhält man durch Anwendung der inversen projektiven Transformation  $T_k^{-1}$  auf den in (13.20) bestimmten Vektor  $y^{k+1}$ :

$$(13.21) \quad x^{k+1} := T_k^{-1}(y^{k+1}) = \frac{1}{\mathbb{1}^T D y^{k+1}} D y^{k+1}.$$

Dem Hilfsprogramm (13.19) kann man übrigens noch eine andere geometrische Interpretation geben. Setzen wir

$$z := Dy \quad \text{bzw.} \quad y := D^{-1}z$$

so kann man (13.19) wie folgt schreiben

$$(13.22) \quad \begin{aligned} \min c^T z \\ Az &= 0 \\ \mathbb{1}^T D^{-1}z &= 1 \\ z &\in E\left(\frac{1}{4n(n-1)}D^2, \frac{1}{n}x^k\right) \end{aligned}$$

Denn wegen (13.17) gilt  $\bar{c}^T y = c^T z$  und aus  $D^{-1}x^k = \mathbb{1}$  folgt:

$$\begin{aligned} \|y - \frac{1}{n}\mathbb{1}\| \leq \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} &\iff \|D^{-1}z - \frac{1}{n}D^{-1}x^k\| \leq \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} \\ &\iff \|D^{-1}(z - \frac{1}{n}x^k)\| \leq \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} \\ &\iff (z - \frac{1}{n}x^k)^T D^{-2}(z - \frac{1}{n}x^k) \leq \frac{1}{4n(n-1)} \\ &\iff (z - \frac{1}{n}x^k)(4n(n-1))D^{-2}(z - \frac{1}{n}x^k) \leq 1 \\ &\stackrel{(12.20)}{\iff} z \in E\left(\frac{1}{4n(n-1)}D^2, \frac{1}{n}x^k\right) \end{aligned}$$

Gehen wir vom Ursprungsproblem (13.11) aus, so bedeutet dies also, dass wir in unserem Hilfsprogramm die Originalzielfunktion  $c$  und den linearen Raum  $\Omega$  unverändert lassen. Wir ersetzen die Hyperebene  $E$  durch  $E_k = \{z \mid \mathbb{1}^T D^{-1}z = 1\}$  und die Nichtnegativitätsbedingung  $x \in \mathbb{R}_+^n$  durch das Ellipsoid  $E(\frac{1}{4n(n-1)}D^2, \frac{1}{n}x^k)$  und optimieren hierüber. Aus einer Optimallösung, sagen wir  $z^{k+1}$ , von (13.22) erhalten wir eine Optimallösung von (13.19) durch

$$y^{k+1} = D^{-1}z^{k+1}.$$

Wie im Beweis von (13.20) kann man zeigen, dass (13.22) durch eine direkte Formel gelöst werden kann (das folgt natürlich auch aus (13.20) durch die obige Gleichung), und zwar gilt:

$$(13.23) \quad z^{k+1} = \frac{1}{n}x^k - \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} \frac{D(I - DA^T(AD^2A^T)^{-1}AD - \frac{1}{n}\mathbb{1}\mathbb{1}^T)Dc}{\|(I - DA^T(AD^2A^T)^{-1}AD - \frac{1}{n}\mathbb{1}\mathbb{1}^T)Dc\|}.$$

Hieraus ergibt sich über (13.21) eine neue (einfache) Formel zur Berechnung des nächsten Iterationspunktes.

$$(13.24) \quad x^{k+1} := \frac{1}{\mathbb{1}^T Dy^{k+1}} Dy^{k+1} = \frac{1}{\mathbb{1}^T z^{k+1}} z^{k+1}.$$

Damit haben wir die wesentlichen Bestandteile eines Iterationsschrittes des Karmarkar-Algorithmus beschrieben und können ihn zusammenfassend darstellen, wobei noch das Abbruchkriterium zu begründen ist.

**(13.25) Der Karmarkar-Algorithmus.**

**Input:**  $A \in \mathbb{Q}^{(m,n)}$  und  $c \in \mathbb{Q}^n$ . Zusätzlich wird vorausgesetzt, dass  $\frac{1}{n}A\mathbb{1} = 0$  und  $c^T\mathbb{1} > 0$  gilt.

**Output:** Ein Vektor  $x$  mit  $Ax = 0$ ,  $\mathbb{1}^T x = 1$ ,  $x \geq 0$  und  $c^T x \leq 0$  oder die Feststellung, dass kein derartiger Vektor existiert.

(1) **Initialisierung.** Setze

$$\begin{aligned} x^0 &:= \frac{1}{n}\mathbb{1} \\ k &:= 0 \\ N &:= 3n(\langle A \rangle + 2\langle c \rangle - n) \end{aligned}$$

(2) **Abbruchkriterium.**

(2.a) Gilt  $k = N$ , dann hat  $Ax = 0$ ,  $\mathbb{1}^T x = 1$ ,  $x \geq 0$ ,  $c^T x \leq 0$  keine Lösung, **STOP!**

(2.b) Gilt  $c^T x^k \leq 2^{-(A)-\langle c \rangle}$ , dann ist eine Lösung gefunden. Falls  $c^T x^k \leq 0$ , dann ist  $x^k$  eine Lösung, andernfalls kann wie bei der Ellipsoidmethode (Satz (12.34)) aus  $x^k$  ein Vektor  $\bar{x}$  konstruiert werden mit  $c^T \bar{x} \leq 0$ ,  $A\bar{x} = 0$ ,  $\mathbb{1}^T \bar{x} = 1$ ,  $\bar{x} \geq 0$ , **STOP!**

**Update.**

(3) (3.a)  $D := \text{diag}(x^k)$

(3.b)  $\bar{c} := (I - DA^T(AD^2A^T)^{-1}AD - \frac{1}{n}\mathbb{1}\mathbb{1}^T)Dc$  (siehe Beweis von (13.20))

(3.c)  $y^{k+1} := \frac{1}{n}\mathbb{1} - \frac{1}{2} \frac{1}{\sqrt{n(n-1)} \|\bar{c}\|} \bar{c}$

(3.d)  $x^{k+1} := \frac{1}{\mathbb{1}^T D y^{k+1}} D y^{k+1}$

(3.e)  $k := k + 1$

Gehe zu (2).

□

Die Geschwindigkeit des Algorithmus hängt natürlich nur von der Implementierung des Schrittes (3) ab. Wir haben hier die aus Lemma (13.20) gewonnene Formel (13.21) für den nächsten Iterationspunkt  $x^{k+1}$  gewählt. Man kann  $x^{k+1}$  natürlich auch über (13.23), (13.24) bestimmen. Wesentlich ist, dass man eine Update-Formel findet, in der möglichst wenig arithmetische Operationen auftreten.

Es ist offensichtlich, dass die Hauptarbeit von (3) im Schritt (3.b) liegt. Führt man ihn kanonisch aus, so werden  $O(n^3)$  Rechenschritte benötigt. Man beachte, dass sich in jedem Schritt nur die Diagonalmatrix  $D$  ändert. Diese Tatsache kann man, wie Karmarkar gezeigt hat, ausnutzen, um (3) in  $O(n^{2.5})$  Rechenschritten zu erledigen. Die Laufzeit beträgt dann insgesamt  $O(n^{2.5}(\langle A \rangle + \langle c \rangle))$ .

Wie bei der Ellipsoidmethode können bei der Ausführung des Karmarkar-Algorithmus irrationale Zahlen (durch Wurzelziehen in (3.c)) auftreten. Wir sind also gezwungen, alle Rechnungen approximativ auszuführen. Die dadurch auftretenden Rundungsfehler akkumulieren sich natürlich. Wir müssen also unsere Rundungsvorschriften so einrichten, dass beim Abbruch des Verfahrens eine korrekte Schlussfolgerung gezogen werden kann. Auch das kann man wie bei der Ellipsoidmethode erledigen. Dies sind (auf Abschätzungstechniken der linearen Algebra und Analysis beruhende) Tricks, mit deren Hilfe Verfahren (13.25) in einen polynomialen Algorithmus abgewandelt werden kann. Da insgesamt höchstens  $N$  Iterationen durchgeführt werden, folgt

**(13.26) Satz.** *Die Laufzeit des (geeignet modifizierten) Karmarkar-Algorithmus (13.25) zur Lösung von Problemen des Typs (13.10) ist  $O(n^{3.5}(\langle A \rangle + \langle c \rangle)^2)$ .*  $\square$

Wir wollen im Weiteren die Rundungsfehler vernachlässigen und annehmen, dass wir in perfekter Arithmetik arbeiten. Es bleibt noch zu zeigen, dass Algorithmus (13.25) mit einem korrekten Ergebnis endet.

Zunächst überlegen wir uns, dass alle Punkte  $x^k$  im relativ Inneren von  $\Omega \cap E \cap \mathbb{R}_+^n$  enthalten sind.

**(13.27) Lemma.** *Für alle  $k \in \{0, 1, \dots, N\}$  gilt  $Ax^k = 0$ ,  $x^k > 0$ ,  $\mathbb{1}^T x^k = 1$ .*

**Beweis :** Durch Induktion über  $k$ ! Für  $k = 0$  gilt die Behauptung nach Voraussetzung. Für  $k + 1$  ist  $\bar{c}$  die Projektion von  $c$  auf  $\{x \mid ADx = 0, \mathbb{1}^T x = 0\}$ , siehe Beweis von (13.20). Daraus folgt  $AD\bar{c} = 0$ ,  $\mathbb{1}^T \bar{c} = 0$ . Mithin ergibt sich

$$\begin{aligned} (\mathbb{1}^T Dy^{k+1})Ax^{k+1} &= ADy^{k+1} = AD\left(\frac{1}{n}\mathbb{1} - \frac{1}{2} \frac{1}{\sqrt{n(n-1)}} \frac{1}{\|\bar{c}\|} \bar{c}\right) \\ AD\left(\frac{1}{n}\mathbb{1}\right) &= Ax^k = 0. \end{aligned}$$

Daraus folgt  $Ax^{k+1} = 0$ .

Um  $x^k > 0$  zu zeigen, stellen wir zunächst fest, dass  $K := \{y \mid \|y - \frac{1}{n}\mathbb{1}\| \leq \frac{1}{2} \frac{1}{\sqrt{n(n-1)}}\} \subseteq \mathbb{R}_+^n$ . Denn gibt es ein  $\bar{y} \in K$  mit einer nichtpositiven Komponente, sagen wir  $\bar{y}_1 \leq 0$ , so gilt  $\sum_{i=2}^n (\bar{y}_i - \frac{1}{n})^2 \leq \frac{1}{4n(n-1)} - (\bar{y}_1 - \frac{1}{n})^2 \leq \frac{1}{4n(n-1)} - \frac{1}{n^2} =$

$\frac{-3n+4}{4n^2(n-1)} < 0$ , ein Widerspruch. Daraus folgt  $y^{k+1} > 0$ , und (13.21) ergibt direkt  $x^{k+1} > 0$ .

Aus (13.21) folgt ebenfalls sofort, dass  $\mathbb{1}^T x^{k+1} = 1$  gilt. □

Entscheidend für die Abschätzung der Anzahl der Iterationsschritte ist die folgende Beobachtung.

**(13.28) Lemma.** *Gibt es ein  $x \geq 0$  mit  $Ax = 0$ ,  $\mathbb{1}^T x = 1$  und  $c^T x \leq 0$  dann gilt für alle  $k$*

$$\frac{(c^T x^{k+1})^n}{\prod_{i=1}^n x_i^{k+1}} \leq \frac{2 (c^T x^k)^n}{e \prod_{i=1}^n x_i^k}.$$

□

Wir wollen uns nun überlegen, wieviele Iterationen im Verfahren (13.25) höchstens durchgeführt werden müssen. Der erste Iterationspunkt ist der Vektor

$$x^0 = \frac{1}{n} \mathbb{1}.$$

Daraus folgt mit einer groben Abschätzung aus (12.10) (a)

$$c^T x^0 \leq \frac{1}{n} \sum_{i=1}^n |c_i| \leq \frac{1}{n} \sum_{i=1}^n 2^{(c_i)-1} < \frac{1}{n} 2^{(c)-n}.$$

Aus Lemma (13.28) ergibt sich (unter den dort gemachten Voraussetzungen)

$$\frac{(c^T x^k)^n}{\prod_{i=1}^n x_i^k} \leq \left(\frac{2}{e}\right)^k \frac{(c^T x^0)^n}{\prod_{i=1}^n x_i^0}$$

und somit wegen  $\prod_{i=1}^n x_i^k \leq 1$  und  $\prod_{i=1}^n x_i^0 = \left(\frac{1}{n}\right)^n$

$$(13.29) \quad \begin{aligned} c^T x^k &< \left(\frac{2}{e}\right)^{\frac{k}{n}} \left(\frac{\prod_{i=1}^n x_i^k}{\prod_{i=1}^n x_i^0}\right)^{\frac{1}{n}} \frac{1}{n} 2^{(c)-n} \\ &\leq \left(\frac{2}{e}\right)^{\frac{k}{n}} 2^{(c)-n}. \end{aligned}$$

In jedem Schritt schrumpft also der Zielfunktionswert um den nur von der Dimension  $n$  abhängigen Faktor  $\sqrt[n]{\frac{2}{e}}$ . Setzen wir

$$(13.30) \quad N := 3n(\langle A \rangle + 2n\langle c \rangle - n)$$

dann gilt

**(13.31) Satz.** *Es gibt ein  $x \geq 0$  mit  $Ax = 0$ ,  $\mathbb{1}^T x = 1$ ,  $c^T x \leq 0$  genau dann, wenn es ein  $k \in \{0, \dots, N\}$  gibt mit*

$$c^T x^k < 2^{-\langle A \rangle - \langle c \rangle}.$$

**Beweis :** “ $\implies$ ” Angenommen, es gibt ein  $x \in P = \Sigma \cap \Omega$  mit  $c^T x \leq 0$ , dann sind die Voraussetzungen von Lemma (13.28) erfüllt und folglich gilt mit (13.29) und (13.30)

$$c^T x^N < \left(\frac{2}{e}\right)^{\frac{N}{n}} 2^{\langle c \rangle - n} < 2^{-\langle A \rangle - \langle c \rangle},$$

wobei wir die Tatsache ausgenutzt haben, dass  $3 > \frac{-1}{\log(\frac{2}{e})}$  gilt.

“ $\impliedby$ ” Angenommen  $c^T x^k < 2^{-\langle A \rangle - \langle c \rangle}$  gilt für ein  $k \in \{0, \dots, N\}$ . Da  $P \neq \emptyset$ , ein Polytop ist, wird  $\min\{c^T x \mid x \in P\}$  in einer Ecke von  $P$  angenommen. Nach Satz (12.15) ist der Optimalwert dieses Programms eine rationale Zahl  $\frac{p}{q}$  mit

$$1 \leq q \leq 2^{\langle A \rangle + \langle \mathbb{1} \rangle + \langle c \rangle - n^2 - n} < 2^{\langle A \rangle + \langle c \rangle}.$$

Aus  $\frac{p}{q} \leq c^T x^k < 2^{-\langle A \rangle - \langle c \rangle}$  folgt dann  $\frac{p}{q} \leq 0$ . □



# Kapitel 14

## Duale Heuristiken, Relaxierungen

In den Kapiteln 9 – 12 haben wir uns damit beschäftigt, wie man (gute) zulässige Lösungen für ein kombinatorisches Optimierungsproblem findet.

Mit den in diesen Kapiteln vorgestellten Methoden berechnet man obere Schranken (bzw. untere Schranken) für Minimierungsprobleme (bzw. Maximierungsprobleme). Kapitel 13 brachte eine ganz andere Betrachtungsweise: probabilistische Analyse von Verfahren bzw. Optimalwerten. Nun wollen wir uns Sicherheit verschaffen, und zwar dadurch, dass wir Methoden angeben, die die Zielfunktionswerte von der “anderen Seite” beschränken.

Dieses Thema ist meiner Meinung nach genau so wichtig, wie der Entwurf von Primalheuristiken. Denn nur durch die gleichzeitige Kenntnis guter unterer und oberer Schranken kann man sich ein Bild von der Qualität von Lösungen machen. Leider wird der Berechnung derartiger Schranken in der Literatur bisher wenig Beachtung geschenkt, in der Praxis scheint diese Möglichkeit weithin unbekannt zu sein.

Wie bereits gesagt, wollen wir Verfahren, die für ein Maximierungsproblem eine obere Schranke bzw. für ein Minimierungsproblem eine untere Schranke für den optimalen Zielfunktionswert bestimmen, *duale Heuristiken* nennen. Duale Heuristiken liefern i. a. keine zulässigen Lösungen, ermöglichen aber — zusammen mit primalen Heuristiken — häufig gute Abschätzungen des optimalen Zielfunktionswertes.

Wir haben bereits mehrere duale Heuristiken kennengelernt, ohne dies explizit zu sagen. Denn bei jeder primalen Heuristik, für die wir eine Gütegarantie bewiesen haben, muss ten wir eine Abschätzung des Optimalwertes finden. Diese Abschätzung kann man als duale Heuristik ansehen. Wir lassen einige dieser Me-

thoden direkt Revue passieren.

Beginnen wir mit der LIST-Heuristik (9.3) für das Parallel-Shop-Problem (9.1). Hier stehen  $m$  identische Maschinen parallel zur Verfügung, und  $n$  Aufgaben mit Bearbeitungsdauern  $t_1, \dots, t_n$  sind gegeben. Offenbar muss insgesamt  $\sum_{i=1}^n t_i$  Bearbeitungszeit zur Verfügung gestellt werden. Verteilen wir diese gleichmäßig auf alle  $m$  Maschinen, so gilt für die optimale Lösung  $c_{\text{opt}}$

$$c_{\text{opt}} \geq \frac{1}{m} \sum_{i=1}^n t_i.$$

Diese triviale Abschätzung haben wir in (9.6) benutzt. Wir können sie natürlich als duale Heuristik auffassen.

Bemerkung (12.8) liefert die folgende Dualheuristik für das 0/1-Rucksackproblem:

**(14.1) Dualheuristik für das 0/1-Knapsackproblem.**

**Input:**  $c_1, \dots, c_n \in \mathbb{N}$ ,  $a_1, \dots, a_n \in \mathbb{N}$ ,  $b \in \mathbb{N}$ .

**Output:** Obere Schranke für den Optimalwert  $c_{\text{opt}}$  von

$$\max\{c^T x \mid a^T x \leq b, x \in \{0, 1\}^n\}$$

1. MIN :=  $+\infty$
2. DO  $k = 0$  TO  $n - 1$ : MIN :=  $\min\{\text{MIN}, \sum_{j=1}^k c_j + \frac{c_{k+1}}{a_{k+1}}(b - \sum_{j=1}^k a_j)\}$ .
3. Gib MIN aus.

□

Für die Christofides-Heuristik (10.1) (g) haben wir in Satz (10.4) gezeigt, dass sie ein  $\frac{1}{2}$ -approximatives Verfahren für das symmetrische TSP mit Dreiecksungleichung ist. Ist  $T_{\text{CH}}$  eine durch die Christofides-Heuristik gefundene Tour, so gilt daher für jede optimale Tour  $T_{\text{opt}}$  die Abschätzung  $c(T_{\text{opt}}) \geq \frac{2}{3}c(T_{\text{CH}})$ . Also kann man aus der (primalen) Christofides-Heuristik eine duale Heuristik für das euklidische symmetrische TSP machen.

## 14.1 Zwei Relaxierungstechniken

Derzeit werden im wesentlichen zwei grundsätzliche Relaxierungstechniken für kombinatorische Optimierungsprobleme benutzt. Diese wollen wir hier kurz beschreiben und an Beispielen erläutern. Eine Methode, Relaxierungen zu verbessern, wird in Abschnitt 14.2 vorgeführt.

Generell ist eine Relaxierung eine “Einbettung” eines Problems in ein größeres. Uns interessieren die beiden folgenden Relaxierungen.

**(14.2) Definition.** Gegeben sei ein kombinatorisches Optimierungsproblem  $(E, \mathcal{I}, c)$ .

- (a) Jedes kombinatorische Optimierungsproblem  $(E', \mathcal{I}', c')$ , in das  $(E, \mathcal{I}, c)$  eingebettet werden kann, nennen wir eine **kombinatorische Relaxierung** von  $(E, \mathcal{I}, c)$ . “Einbettung” heißt hier, es gibt eine injektive Abbildung  $\varphi : E \rightarrow E'$  mit  $\varphi(I) \in \mathcal{I}'$  für alle  $I \in \mathcal{I}$  und  $c(I) = c'(\varphi(I))$  für alle  $I \in \mathcal{I}$ .
- (b) Jedes lineare Programm  $\max$  (oder  $\min$ )  $c^T x$ ,  $Ax \leq b$ ,  $x \geq 0$  mit der Eigenschaft  $\{x \in \mathbb{K}^E \mid Ax \leq b, x \geq 0, x \text{ ganzzahlig}\}$  ist die Menge der Inzidenzvektoren von  $\mathcal{I}$  heißt **LP-Relaxierung** von  $(E, \mathcal{I}, c)$ .

□

Was können wir nun mit Relaxierungen anfangen? Angenommen  $(E, \mathcal{I}, c)$  ist ein kombinatorisches Maximierungsproblem. Aus der Einbettungsvorschrift folgt direkt

$$(14.3) \quad \max\{c'(I') \mid I' \in \mathcal{I}'\} \geq \max\{c(I) \mid I \in \mathcal{I}\}$$

für jede kombinatorische Relaxierung von  $(E, \mathcal{I}, c)$  und analog

$$(14.4) \quad \max\{c^T x \mid Ax \leq b, x \geq 0\} \geq \max\{c(I) \mid I \in \mathcal{I}\}$$

für jede LP-Relaxierung von  $(E, \mathcal{I}, c)$ . Entsprechendes gilt für Minimierungsprobleme. Daraus ergibt sich, daß jeder Algorithmus zur Lösung einer kombinatorischen oder LP-Relaxierung eines kombinatorischen Optimierungsproblems als Dualheuristik angesehen werden kann. Der Erfolg derartiger Dualheuristiken hängt natürlich ganz entscheidend von der Wahl der Relaxierung ab. Bei kombinatorischen Relaxierung ist das Erreichen von zwei Zielen wünschenswert:

- (14.4) (a)  $|\mathcal{I}'| - |\mathcal{I}|$  soll möglichst klein sein, d. h.  $\mathcal{I}'$  soll nur einige wenige Elemente mehr als  $\mathcal{I}$  enthalten. Dadurch besteht die begründete Hoffnung, daß die Optimallösung über  $\mathcal{I}'$  bereits in  $\mathcal{I}$  liegt bzw. die Optimalwerte der beiden Probleme nicht stark voneinander abweichen.

(b)  $(E', \mathcal{I}', c')$  ist effizient (polynomial) lösbar.

Bei LP-Relaxierungen sollten das Polyeder  $\{x \mid Ax \leq b, x \geq 0\}$  möglichst “nahe” an  $\text{conv}\{\chi^I \mid I \in \mathcal{I}\}$  liegen, und ferner sollte das LP  $\max c^T x, Ax \leq b, x \geq 0$  möglichst effizient lösbar sein.

Dem Thema LP-Relaxierungen und der Lösung derartiger linearer Programme werden wir die Kapitel 18 und 19 widmen und dort genauer auf die derzeit gegebenen Möglichkeiten eingehen, kombinatorische Optimierungsprobleme mit Hilfe linearer Optimierungsverfahren zu lösen.

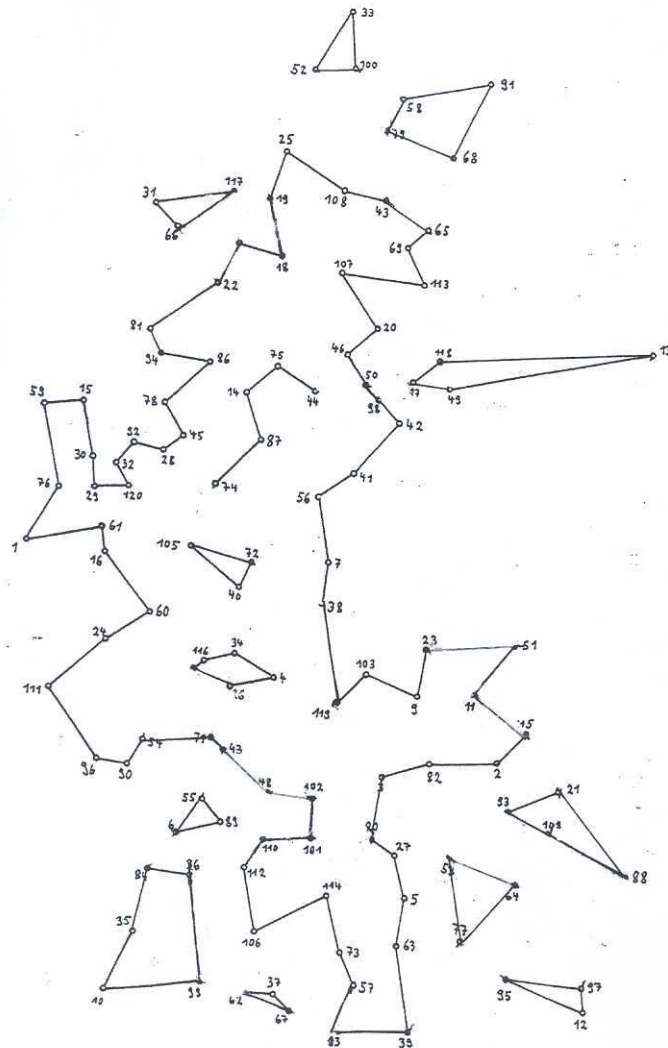
Wir geben einige Beispiele für kombinatorische und LP-Relaxierungen an.

**(14.5) 2-Matching Relaxierung des symmetrischen TSP.** Ein perfektes 2-Matching (kurz **2-Matching**) ist eine Kantenmenge  $M$  eines Graphen  $G$ , so dass jeder Knoten auf genau zwei Kanten liegt. 2-Matchings sind also Vereinigungen von knotendisjunkten Kreisen, so dass jeder Knoten auf einem Kreis liegt. Offenbar ist jede Tour ein 2-Matching. Daraus folgt für jede Zielfunktion  $c : \mathbb{K}^E \rightarrow \mathbb{K}$

$$\min\{c(M) \mid M \text{ 2-Matching}\} \leq \min\{c(T) \mid T \text{ Tour}\}.$$

Minimale 2-Matchings kann man mit einem Algorithmus von Edmonds in polynomialer Zeit bestimmen. Jedoch sind bisher noch keine wirklich effizienten 2-Matching Codes entwickelt worden, so dass die Nützlichkeit der 2-Matching-Relaxierung des symmetrischen Travelling-Salesman-Problems noch nicht intensiv untersucht wurde.  $\square$

Um ein Beispiel für die Abweichung des Wertes einer 2-Matching Lösung von der minimalen Tourlänge zu geben, habe ich mit LP-Techniken (Schnittebenenverfahren, siehe Kapitel 19) ein minimales 2-Matching für das 120-Städte Deutschland-Problem, dessen kürzeste Rundreise in Abbildung 2.3 dargestellt ist, berechnet. Dieses minimale 2-Matching ist in Abbildung 14.1 zu sehen. Die Länge dieses 2-Matchings beträgt 6694 km, während die optimale Tourlänge 6942 km beträgt. Das minimale 2-Matching ist also 3,57 % kürzer als die Optimaltour.



2-Matching für das 120-Städte Deutschland-Problem, Länge: 6694 km

Abb. 14.1

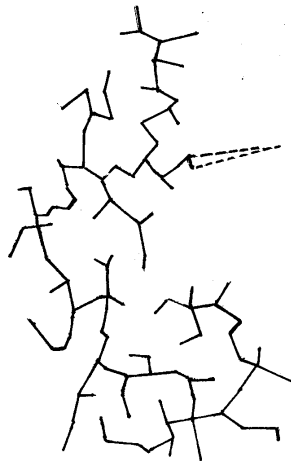
**(14.6) 1-Baum-Relaxierung des symmetrischen TSP.** Sei  $G = (V, E)$  ein zweifach zusammenhängender Graph und  $v$  ein Knoten von  $G$ . Jede Kantenmenge  $B$ , die aus einem aufspannenden Baum von  $G - v$  und zwei Kanten, die mit  $v$  inzidieren, besteht, heißt **1-Baum** (von  $G$ ). Die Menge aller 1-Bäume von  $G$  ist das Basissystem eines Matroids auf  $E$ , siehe Kapitel 5. Man kann also einen minimalen 1-Baum von  $G$  mit dem Greedy-Algorithmus (5.19) bestimmen. Konkret berechnet man zunächst mit einem der Algorithmen aus Abschnitt 4.1 einen minimalen aufspannenden Baum von  $G - v$ , dazu fügt man zwei Kanten, die mit  $v$  inzidieren, deren Gewichtsumme minimal ist.

Jede Tour des  $K_n$  ist ein 1-Baum des  $K_n$ , denn sie besteht aus einem hamiltonschen Weg von  $K_n - v$  und den beiden Kanten, die  $v$  mit den Endknoten des hamiltonschen Weges verbinden. Daher gilt für jede Zielfunktion  $c \in \mathbb{K}^E$

$$\min\{c(B) \mid B \text{ 1-Baum}\} \leq \min\{c(T) \mid T \text{ Tour}\}.$$

□

Wie wir oben angemerkt haben, sind 1-Bäume sehr einfach zu berechnen. Für das 120-Städte Deutschland-Problem ist ein minimaler 1-Baum ( $v = \text{Knoten 13} = \text{Berlin}$ ) in Abbildung 14.2 dargestellt. Seine Länge beträgt 6 025 km. Er ist also 13,21 % kürzer als die Optimaltour. I. a. ist die 1-Baum-Relaxierung erheblich schlechter als die 2-Matching-Relaxierung, aber sie ist viel einfacher zu berechnen. In Abschnitt 14.2 werden wir sehen, wie man die 1-Baum-Relaxierung erheblich verbessern kann.



minimaler 1-Baum: Länge 6025 km

**Abb. 14.2**

**(14.7) Zuordnungsrelaxierung des asymmetrischen TSP.** Gegeben sei ein vollständiger Digraph (ohne Schleifen)  $D = (V, A)$  mit Bogengewichten  $c_a \in \mathbb{K}$  für alle  $a \in A$ . Gesucht ist eine Bogenmenge  $Z$  minimalen Gewichts, so dass jeder Knoten jeweils Anfangsknoten und Endknoten genau eines Bogens aus  $Z$  ist. Jede Zuordnung ist also die Vereinigung knotendisjunkter gerichteter Kreise, so dass jeder Knoten auf genau einem Kreis liegt. Offenbar ist jede Tour eine Zuordnung. □

Diese Version des Zuordnungsproblems ist ein Spezialfall des Problems, einen kostenminimalen Fluss in einem Digraphen zu finden. Für dieses Zuordnungsproblem gibt es sehr effiziente Algorithmen (worst-case-Laufzeit  $O(n^3)$ , praktisch

schneller). Es ist die “beliebteste” Relaxierung des asymmetrischen Travelling-Salesman-Problems und wird mit gutem Erfolg in Branch & Bound Verfahren für das asymmetrische Travelling-Salesman-Problems eingesetzt, siehe Kapitel 15.

Die Zuordnungsrelaxierung des asymmetrischen Travelling-Salesman-Problems nimmt eine Zwitterstellung zwischen den beiden Relaxierungstypen ein. Wie man leicht sieht ist das folgende ganzzahlige Programm eine 0/1-Formulierung des asymmetrischen Travelling-Salesman-Problems

$$(14.8) \quad \min \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n c_{ij} x_{ij}$$

$$\begin{array}{ll} \text{(i)} & \sum_{i=1, i \neq j}^n x_{ij} = 1 \quad \text{für alle } j \in V \\ \text{(ii)} & \sum_{j=1, j \neq i}^n x_{ij} = 1 \quad \text{für alle } i \in V \\ \text{(iii)} & \sum_{i \in W} \sum_{j \in W, i \neq j} x_{ij} \leq |W| - 1 \quad \text{für alle } W \subseteq V, 3 \leq |W| \leq n - 1 \\ \text{(iv)} & x_{ij} \geq 0 \quad \text{für alle } (i, j) \in A \\ \text{(v)} & x_{ij} \text{ ganzzahlig} \quad \text{für alle } (i, j) \in A \end{array}$$

Die Ungleichungen (iii) nennt man Kurzzzyklus-Bedingungen. Lässt man diese aus (14.8) weg, so erhält man eine Formulierung des Zuordnungsproblems aus (14.7) als ganzzahliges Programm. Aus den Resultaten von Kapitel 8 folgt, dass in diesem Programm die Ganzzahligkeitsbedingung überflüssig ist, da alle Basislösungen des Systems (i), (iii), (iv) ganzzahlig sind. Das Zuordnungsproblem ist also auch eine gewisse LP-Relaxierung des asymmetrischen Travelling-Salesman-Problems.

Damit wollen wir kurz auf weitere LP-Relaxierungen eingehen, weiterführende Studien dieses Themas folgen in den Kapiteln 18 und 19.

**(14.9) LP-Relaxierungen des Stabile-Mengen-Problems** (siehe (2.14)). Sei  $G = (V, E)$  ein Graph mit Knotengewichten  $c_v$  für alle  $v \in V$ . Gesucht ist eine stabile Menge maximalen Gewichts. Da keine stabile Knotenmenge  $S \subseteq V$  eine Kante  $uv$  enthalten kann, sieht man sofort, dass das LP

$$(14.10) \quad \max \sum_{v \in V} c_v x_v$$

$$\begin{array}{ll} \text{(i)} & x_u + x_v \leq 1 \quad \text{für alle } uv \in E \\ \text{(ii)} & 0 \leq x_v \leq 1 \quad \text{für alle } v \in V \\ \text{(iii)} & x_v \text{ ganzzahlig} \quad \text{für alle } v \in V \end{array}$$

eine 0/1-Formulierung des Stabile-Mengen-Problems ist. Das LP, das aus (14.10) durch Weglassen der Ganzzahligkeitsbedingung (iii) entsteht, ist somit eine LP-Relaxierung des Stabile-Mengen-Problems. Man kann beweisen, dass alle Basislösungen dieses LP's genau dann ganzzahlig sind, wenn  $G$  ein bipartiter Graph ist. Diese LP-Relaxierung ist einfach (mit Spezialalgorithmen) lösbar, liefert allerdings i. a. keine allzu guten oberen Schranken für das Maximalgewicht einer stabilen Menge.  $\square$

Man kann diese LP-Formulierung (14.10) auf vielfältige Weise verschärfen. Umfangreiche Untersuchungen dieser Art sind in Grötschel, Lovász & Schrijver (1986) zu finden.

**(14.11) LP-Relaxierung des Knotenfärbungsproblems (2.15).** Gegeben sei ein Graph  $G = (V, E)$  mit Knotengewichten  $b_v \in \mathbb{Z}_+$ ,  $v \in V$ . Gesucht ist eine Folge stabiler Mengen (eine **Färbung**)  $S_1, \dots, S_t$ , so dass jeder Knoten in mindestens  $b_v$  dieser Mengen liegt und  $t$  minimal ist. Jeder stabilen Menge  $S \subseteq V$  ordnen wir eine ganzzahlige Variable  $\lambda_S$  zu. (Die Anzahl der stabilen Mengen eines Graphen ist i. a. exponentiell in  $|V|$ !) Die Variable  $\lambda_S$  gibt an, wie oft die stabile Menge  $S$  in einer Färbung  $S_1, \dots, S_t$  vorkommt. Ist  $\mathcal{S}$  die Menge der stabilen Teilmengen von  $V$ , so ist folglich

$$(14.12) \quad \begin{array}{ll} \min \sum_{S \in \mathcal{S}} \lambda_S & \\ \text{(i)} \quad \sum_{S \in \mathcal{S}, v \in S} \lambda_S \geq b_v & \text{für alle } v \in V, \\ \text{(ii)} \quad \lambda_S \geq 0 & \text{für alle } S \in \mathcal{S}, \\ \text{(iii)} \quad \lambda_S \text{ ganzzahlig} & \text{für alle } S \in \mathcal{S} \end{array}$$

eine Formulierung des Knotenfärbungsproblems als ganzzahliges lineares Programm (in exponentiell vielen Variablen). Lässt man die Ganzzahligkeitsbedingung (iii) in (14.12) weg, so erhält man eine LP-Relaxierung des Färbungsproblems.

Es ist natürlich unklar, wie man ein derartig großes LP lösen kann. Wir wollen daher auf die LP-Dualität zurückgreifen und weiter relaxieren. Das zu diesem LP duale lineare Programm lautet

$$(14.13) \quad \begin{array}{ll} \max \sum_{v \in V} b_v x_v & \\ \text{(i)} \quad \sum_{v \in S} x_v \leq 1 & \text{für alle } S \subseteq \mathcal{S} \\ \text{(ii)} \quad x_v \geq 0 & \text{für alle } v \in V \end{array}$$

Das Programm (14.13) hat  $|V|$  Variable und  $|\mathcal{S}| + |V|$ , also i. a. exponentiell viele Nebenbedingungen. Aus der Dualitätstheorie wissen wir, dass jede Lösung von (14.13) einen Wert hat, der nicht größer ist als der Wert jeder Lösung der LP-Relaxierung von (14.12), also insbesondere liefert jede zulässige Lösung von (14.13) eine untere Schranke von (14.12) und somit eine untere Schranke für den Wert einer Knotenfärbung.

Man kann zeigen, dass Probleme des Typs (14.13)  $\mathcal{NP}$ -schwer sind, aber da jede Lösung von (14.13) eine untere Schranke für (14.12) liefert, kann man versuchen mit Heuristiken "gute" Lösungen von (14.13) zu finden.

Offenbar ist jeder Inzidenzvektor einer Clique von  $G$  eine (ganzzahlige) Lösung von (14.13). Also liefert jede Heuristik für das Cliquesproblem (2.14) eine untere Schranke für (14.12).  $\square$

## 14.2 Lagrange-Relaxierungen

Wir wollen nun eine Methode vorstellen, mit der man gegebene Relaxierungen verbessern kann. Wir werden die Idee zunächst an der 1-Baum-Relaxierung des symmetrischen TSPs erklären (hierbei ist sie auch von Held & Karp (1970) entwickelt worden), und sie dann in voller Allgemeinheit formulieren. Wir beginnen mit einer simplen Beobachtung.

**(14.14) Lemma.** Gegeben sei ein symmetrisches TSP mit Entfernungen  $c_{ij}$ ,  $1 \leq i < j \leq n$ . Seien  $\lambda_i$ ,  $i = 1, \dots, n$ , Knotengewichte,  $\lambda := \sum_{i=1}^n \lambda_i$  und

$$c'_{ij} := c_{ij} + \lambda_i + \lambda_j, \quad 1 \leq i < j \leq n.$$

Für das modifizierte TSP mit Entfernungen  $c'_{ij}$  gilt:

$$\begin{aligned} c'(T) &= c(T) + 2 \sum_{i=1}^n \lambda_i = c(T) + \lambda \quad \text{für alle Touren } T, \\ c'(B) &= c(B) + \sum_{i=1}^n d_B(i) \lambda_i, \quad \text{für alle 1-Bäume } B, \end{aligned}$$

wobei  $d_B(i)$  den Grad des Knotens  $i$  im 1-Baum  $B$  bezeichnet. Das heißt, die optimalen Touren bezüglich  $c$  bleiben auch optimal bezüglich  $c'$ , während ein optimaler 1-Baum bezüglich  $c$  nicht unbedingt optimal bezüglich  $c'$  ist.

**Beweis :** Trivial. □

**(14.15) Folgerung.** Gegeben sei ein symmetrisches TSP mit Entfernungen  $c_{ij}$ ,  $1 \leq i < j \leq n$ . Dann ist

$$f^* := \max_{\lambda \in \mathbb{R}^n} \left( \min_{B \text{ 1-Baum}} \left( c(B) + \sum_{i=1}^n (d_B(i) - 2) \lambda_i \right) \right)$$

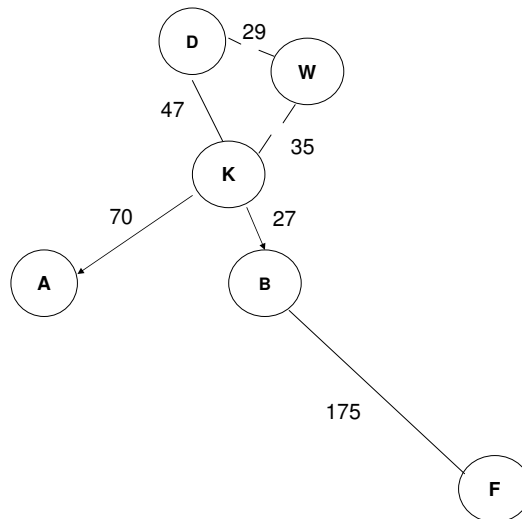
eine untere Schranke für die Länge der optimalen Tour.

**Beweis :** Jede Tour ist ein 1-Baum, und es ist  $d_T(i) = 2$  für jeden Knoten  $i$  in jeder Tour  $T$ . Also gilt für jede Tour  $T$

$$c(T) + \sum_{i=1}^n (d_T(i) - 2) \lambda_i = c(T),$$

d. h. für jedes  $\lambda \in \mathbb{K}^n$  ist der Wert des minimalen 1-Baumes höchstens so groß wie die Länge der kürzesten Tour. □

**(14.16) Beispiel.** Wir betrachten erneut das 6-Städte TSP aus Beispiel (10.2). Die Entfernungen sind in Tabelle 10.1 angegeben. Wir eliminieren Wuppertal und erhalten den in Abbildung 14.3 dargestellten minimalen 1-Baum mit Länge 403 km.



**Abb. 14.3**

Wir führen als Knotengewichte  $\lambda_K = 15$ ,  $\lambda_F = -120$ ,  $\lambda_A = -5$ ,  $\lambda_D = 0$ ,  $\lambda_B = 0$  ein und erhalten die in Tabelle 14.1 wiedergegebene Entfernungsmatrix.

		-5	0	0	-120	15	0
		A	B	D	F	K	W
-5	A	—	85	80	134	80	116
0	B	85	—	77	55	43	84
0	D	80	77	—	112	62	29
-120	F	134	55	112	—	84	116
15	K	80	42	62	84	—	70
0	W	116	84	29	116	70	—

Tabelle 14.1

Der kürzeste 1-Baum bezüglich der neuen Entfernungen hat die Länge 478 km und ist in Abbildung 14.4 dargestellt.

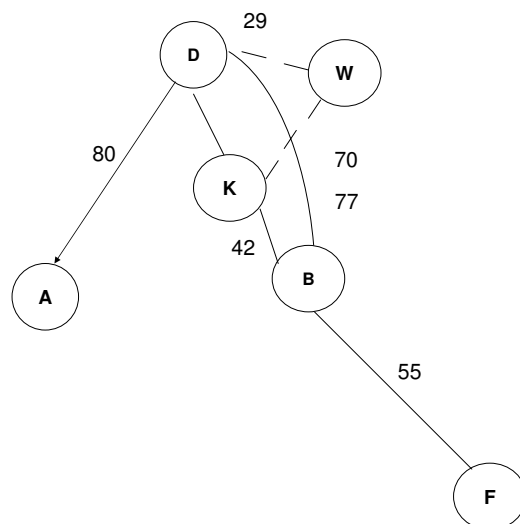


Abb. 14.4

Damit haben wir also eine bessere untere Schranke für die optimale Tour gefunden.

□

Wir wissen zwar, wie wir 1-Bäume schnell berechnen können, aber wie soll das in (14.15) formulierte freie Maximierungsproblem (über alle  $\lambda \in \mathbb{K}^n$ ) gelöst werden? Wie finden wir  $f^*$ ?

Auf folgende Weise können wir eine Funktion  $f : \mathbb{K}^n \rightarrow \mathbb{K}$  definieren:

$$(14.17) \quad f(\lambda) := \min_{B1\text{-Baum}} \left( c(B) + \sum_{i=1}^n (d_B(i) - 2)\lambda_i \right),$$

d. h.  $f$  ist eine Funktion, deren Auswertung die Berechnung der Optimallösung eines Optimierungsproblems erfordert. Und das Problem,  $f^*$  zu finden, ist damit das unbeschränkte Optimierungsproblem

$$\max_{\lambda \in \mathbb{K}^n} f(\lambda).$$

Um diese Situation zu analysieren, verallgemeinern wir die vorliegende Fragestellung:

Wir betrachten ein ( $\mathcal{NP}$ -schwieriges) ganzzahliges lineares Programm der folgen-

den Form

$$\begin{aligned} c^* &:= \min c^T x \\ Ax &= b \\ Dx &\leq d \\ x &\geq 0 \\ x &\text{ ganzzahlig,} \end{aligned}$$

wobei  $D$  eine  $(m, n)$ -Matrix sei.

Wir setzen

$$\begin{aligned} P &:= \{x \in \mathbb{R}^n \mid Ax = b, Dx \leq d, x \geq 0, x \text{ ganzzahlig}\}, \\ Q &:= \{x \in \mathbb{R}^n \mid Dx \leq d, x \geq 0, x \text{ ganzzahlig}\}. \end{aligned}$$

**(14.20) Definition.** Die Funktion  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  definiert durch

$$f(\lambda) := \min\{c^T x + \lambda^T (Ax - b) \mid x \in Q\}$$

heißt **Lagrange-Funktion** von (14.18), der Vektor  $\lambda \in \mathbb{R}^m$  heißt **Vektor der Lagrange-Multiplikatoren**. Das Optimierungsproblem

$$\max_{\lambda \in \mathbb{R}^m} f(\lambda)$$

heißt **Lagrange-Relaxierung** von (14.18).

□

Die in (14.17) definierte Funktion ist die Lagrange-Funktion des symmetrischen TSP bezüglich der 1-Baum Relaxierung. Denn, wie wir wissen, ist das folgende Programm eine ganzzahlige Formulierung des TSP:

(14.21)

$$\begin{aligned} &\min c^T x \\ \text{(i)} & \quad x(\delta(v)) = 2 \quad \text{für alle } v \in V \\ \text{(ii)} & \quad x(E(W)) \leq |W| - 1 \quad \text{für alle } W \subseteq V \setminus \{1\}, 3 \leq |W| \leq n - 3 \\ \text{(iii)} & \quad x(\delta(1)) \leq 2 \\ \text{(iv)} & \quad -x(\delta(1)) \leq -2 \\ \text{(v)} & \quad x \leq \emptyset \\ \text{(vi)} & \quad x \geq 0 \\ \text{(vii)} & \quad x \text{ ganzzahlig.} \end{aligned}$$

Schreiben wir das Gleichungssystem (i) von (14.21) als  $Ax = b$  und die Ungleichungen (ii), ..., (v) als  $Dx \leq d$ , so ist die in (14.17) definierte Funktion genau die Lagrange-Funktion von (14.21).

**(14.22) Satz.** Seien  $Q$  und  $P$  (siehe (14.19)) nicht leer und  $Q$  endlich. Dann gilt:

(a) Die Lagrange-Funktion  $f(\lambda) := \min_{x \in Q} (c^T x + \lambda^T (Ax - b))$  ist konkav, stückweise linear und nach oben beschränkt.

(b) Sei  $\lambda^* \in \mathbb{R}^n$  mit

$$f(\lambda^*) = \max_{\lambda \in \mathbb{R}^n} f(\lambda)$$

eine Optimallösung der Lagrange-Relaxierung von (14.18), so gilt

$$f(\lambda^*) \leq c^*,$$

d. h.  $f(\lambda^*)$  ist eine untere Schranke für den Optimalwert von (14.18).

**Beweis :** (a) Wir können für jeden Vektor  $x \in Q$  eine affine Funktion  $g_x$  in  $\lambda$  wie folgt definieren:

$$g_x(\lambda) := c^T x + \lambda^T (Ax - b).$$

Da  $Q$  endlich ist, ist  $f$  als Minimum endlich vieler affiner Funktionen darstellbar, d. h.

$$f(\lambda) = \min_{x \in Q} g_x(\lambda).$$

Wir zeigen

1)  **$f$  ist konkav.** Es ist zu beweisen:  $f(\alpha\lambda + (1 - \alpha)\mu) \geq \alpha f(\lambda) + (1 - \alpha)f(\mu)$  für alle  $\alpha \geq 0$  und für alle  $\lambda, \mu \in \mathbb{R}^n$ .

$$\begin{aligned} f(\alpha\lambda + (1 - \alpha)\mu) &= \min_{x \in Q} g_x(\alpha\lambda + (1 - \alpha)\mu) \\ &= \min_{x \in Q} (\alpha g_x(\lambda) + (1 - \alpha)g_x(\mu)) \\ &\geq \alpha \min_{x \in Q} g_x(\lambda) + (1 - \alpha) \min_{x \in Q} g_x(\mu) \\ &= \alpha f(\lambda) + (1 - \alpha)f(\mu). \end{aligned}$$

2)  **$f$  ist stückweise linear.** Sei für jedes  $x \in Q$ :  $L_x := \{\lambda \in \mathbb{R}^n \mid f(\lambda) = g_x(\lambda)\}$ . Es gilt natürlich  $\bigcup_{x \in Q} L_x = \mathbb{R}^n$ . Wenn wir zeigen können, dass  $L_x$  konvex ist, dann gilt  $f(\lambda) = g_x(\lambda) \forall \lambda \in L_x$ . Somit ist  $f$  linear auf  $L_x$ , d. h. stückweise linear auf  $\mathbb{R}^n$ .  $L_{x_0}$  konvex  $\iff \alpha\lambda + (1 - \alpha)\mu \in L_{x_0}$  für alle  $1 \geq \alpha \geq 0$  und für alle  $\lambda, \mu \in L_{x_0}$ .

$$\begin{aligned} g_{x_0}(\alpha\lambda + (1 - \alpha)\mu) &\geq f(\alpha\lambda + (1 - \alpha)\mu) \geq \alpha f(\lambda) + (1 - \alpha)f(\mu) = \\ &\alpha g_{x_0}(\lambda) + (1 - \alpha)g_{x_0}(\mu) = g_{x_0}(\alpha\lambda + (1 - \alpha)\mu). \end{aligned}$$

Daraus folgt,  $L_{x_0}$  ist konvex, und somit ist  $f$  stückweise linear.

3)  $f$  ist nach oben beschränkt. Dies ist offensichtlich.

(b) ist offensichtlich. □

Es besteht nun das Problem, das Maximum der Lagrange-Funktion  $f$  zu bestimmen. Die Standardverfahren der nichtlinearen Optimierung (die wir in der Vorlesung "Optimierungsmethoden II" behandeln werden) kann man hier nicht anwenden, da die Funktion  $f$  i. a. nicht differenzierbar ist. Eine Verallgemeinerung des Gradientenkonzepts aus der Analysis führt jedoch zu einer Lösungsmöglichkeit.

**(14.23) Definition.** Sei  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  eine konkave Funktion, dann heißt ein Vektor  $u \in \mathbb{R}^m$  ein **Subgradient** an  $f$  im Punkte  $x_0 \in \mathbb{R}^m$ , falls gilt

$$f(x) - f(x_0) \leq u^T(x - x_0) \quad \text{für alle } x \in \mathbb{R}^m.$$

Die Menge

$$\partial f(x_0) := \{u \in \mathbb{R}^m \mid u \text{ ist Subgradient an } f \text{ in } x_0\}$$

heißt **Subdifferential** von  $f$  in  $x_0$ . □

**(14.24) Bemerkung.** Ist die konkave Funktion  $f$  differenzierbar in  $x_0$ , so gilt

$$\partial f(x_0) = \{\nabla f(x_0)\},$$

wobei  $\nabla f(x_0)$  den Gradienten von  $f$  an der Stelle  $x_0$  bezeichnet. □

Für die speziellen Lagrange-Funktionen, die wir hier betrachten, können wir die Subdifferential einfach berechnen.

**(14.25) Satz.** Sei  $Q$  nicht leer und endlich und  $f(\lambda) := \min_{x \in Q} (c^T x + \lambda^T (Ax - b))$ , so gilt folgendes: Setzen wir für  $\lambda_0 \in \mathbb{R}^m$ ,  $L_0 := \{x_0 \in \mathbb{R}^m \mid f(\lambda_0) = c^T x_0 + \lambda_0^T (Ax_0 - b)\}$ , so ist

$$\partial f(\lambda_0) = \text{conv}\{(Ax_0 - b) \mid x_0 \in L_0\}.$$

□

Aus Zeitgründen können wir auf den Beweis von (14.25) hier nicht eingehen. Satz (14.25) hat einen einfachen geometrischen Gehalt. Für jedes  $\lambda_0 \in \mathbb{R}^m$  ist  $\partial f(\lambda_0)$  ein Polytop. Dieses Polytop erhält man wie folgt. Man nehme alle Punkte  $x \in \mathbb{R}^m$ , die das Minimum des Optimierungsproblems  $\min\{c^T x + \lambda_0^T (Ax - b) \mid x \in Q\}$  realisieren, also alle Optimalwerte bezüglich des Vektors der Lagrange-Multiplikation  $\lambda_0$ . Dann bilde man die konvexe Hülle all dieser Punkte und erhält das Subdifferential  $\partial f(\lambda_0)$ .

**(14.26) Beispiel.** Im Falle des symmetrischen TSPs gilt für einen Vektor von Knotengewichten

$$\lambda := (\lambda_1, \dots, \lambda_n)$$

$$\partial f(\lambda) = \text{conv} \left\{ \begin{pmatrix} d_B(1) - 2 \\ d_B(2) - 2 \\ \vdots \\ d_B(n) - 2 \end{pmatrix} \mid B \text{ optimaler 1-Baum bezüglich } c' \right\}.$$

Hat man einen optimalen 1-Baum  $B$  bezüglich  $\lambda$ , so kann man also einen Subgradienten der Lagrange-Funktion  $f$  direkt aus den Knotengraden von  $B$  ablesen.  $\square$

Der folgende Satz verallgemeinert die uns aus der Schule bekannte Tatsache, dass ein Punkt  $x^*$  eine differenzierbare konkave Funktion  $f$  maximiert, wenn die Ableitung von  $f$  in  $x^*$  Null ist.

**(14.27) Satz.** Ein Vektor  $x^* \in \mathbb{R}^m$  maximiert die Funktion  $f$  genau dann, wenn  $0 \in \partial f(x^*)$ .

**Beweis :** Falls  $0 \in \partial f(x^*)$ , dann ergibt die Subgradientenungleichung der Definition:

$$f(x) - f(x^*) \leq 0^T(x - x^*) \implies f(x) \leq f(x^*) \text{ für alle } x \in \mathbb{R}^m.$$

Falls  $f(x^*)$  maximal ist, so ist  $f(x) \leq f(x^*)$  für alle  $x \in \mathbb{R}^m$ , also  $f(x) - f(x^*) \leq 0 = 0^T(x - x^*)$ . Daraus folgt  $0 \in \partial f(x^*)$ ,  $\square$

Zur Lösung der Optimierungsaufgabe aus (14.20) ist die folgende Methode vorgeschlagen worden:

**(14.28) Subgradientenverfahren.**

**Input:** konkave Funktion  $f : \mathbb{R}^n \rightarrow \mathbb{R}$

**Gesucht:**  $\max\{f(x) \mid x \in \mathbb{R}^n\}$

1. Wähle  $\lambda_0 \in \mathbb{R}^m$  beliebig, setze  $j := 0$ .
2. Berechne  $f(\lambda_j)$  und  $\partial f(\lambda_j)$ .
3. Ist ein (noch zu definierendes) STOP-Kriterium erfüllt?

Falls ja, so ist eine “gute” untere Schranke für das Maximum oder das Maximum selbst gefunden.  $\rightarrow$  STOP.

Falls nein, gehe zu 4.

4. Wähle neues  $\lambda_{j+1}$ , setze z. B.

$$\lambda_{j+1} := \lambda_j + t_j u_j / \|u_j\|,$$

wobei  $t_j \in \mathbb{R}$  eine Schrittlänge ist und  $u_j \in \partial f(\lambda_j)$ .

5. Setze  $j := j + 1$  und gehe zu 2.

□

Bezüglich des Verfahrens (14.28) kann man Folgendes beweisen.

**(14.29) Satz (Polyak, 1967).** Falls die konkave Funktion  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  ein Maximum besitzt und falls die Folge  $(t_j)$  der Schrittlängen die Bedingungen

(a)  $t_j > 0$  für alle  $j \in \mathbb{N}$ ,

(b)  $\lim_{j \rightarrow \infty} t_j = 0$ ,

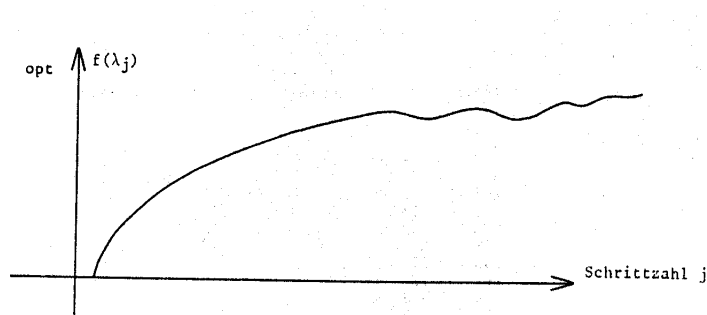
(c)  $\sum_{j=0}^{\infty} t_j = \infty$

erfüllt, so gilt

$$\lim_{j \rightarrow \infty} f(\lambda_j) = \max_{\lambda \in \mathbb{R}^m} f(\lambda).$$

□

Bei konkreten Rechnungen zeigt sich allerdings im allgemeinen, dass die (nach (14.29) theoretisch existierende) Konvergenz praktisch nicht erreicht wird. Abbildung 14.5 zeigt einen typischen Verlauf des Wertes  $f(\lambda_j)$ , abhängig von der Schrittzahl  $j$ .



**Grafik/Abb. 14.5**

Mit anderen Schrittlängenparametern kann man praktisch i. a. ein schnelleres Anwachsen der Funktionswerte  $f(\lambda_j)$  erreichen, jedoch theoretisch nicht unbedingt eine konvergente Methode. Da es sich ja hier “nur” um eine Heuristik zur Bestimmung unterer Schranken handelt, ist der Konvergenzaspekt gleichgültig. Man bricht einfach nach einer bestimmten Rechenzeit oder einem STOP-Kriterium ab und benutzt den bisher besten Zielfunktionswert als untere Schranke.

Aus empirischer Erfahrung kann man festhalten, dass die Schrittlängen  $t_j$  problemabhängig zu wählen sind und dass sie i. a. durch experimentelle Berechnung bestimmt werden müssen. Ein Beispiel für eine derartige Schrittlängenwahl ist in Schritt 5 des Algorithmus aus Beispiel (14.30) zu finden.

**(14.30) Beispiel.** Bei praktischen Versuchen bezüglich des symmetrischen TSP hat sich folgendes Vorgehen (modifiziertes Subgradientenverfahren) recht gut bewährt:

**Input:** symmetrisches TSP mit Zielfunktion  $c$ ,  $c_{ij} \geq 0$ ,  $1 \leq i < j \leq n$ .

**Output:** untere Schranke  $L$  für den Optimalwert des TSP.

1. Initialisierung:

$U$  sei eine obere Schranke für die kürzeste Tour (mit Heuristik gefunden).

$L := 0$  (gegenwärtig beste untere Schranke),

$\lambda^T := (0, 0, \dots, 0)$  (Anfangsknotengewichte),

$0 < \alpha \leq 2$  (i. a.  $\alpha := 2$ ) (Skalierungsfaktor für die Schrittlänge),

$\varepsilon_1, \varepsilon_2, \varepsilon_3 > 0$  numerische Genauigkeitsschranken, Zählparameter  $k$ .

2. Berechne den kürzesten 1-Baum  $B$  bezüglich  $c^\lambda$  (d. h.  $c_{ij}^\lambda := c_{ij} + \lambda_i + \lambda_j$ ). Falls  $B$  eine Tour ist  $\rightarrow$  STOP.

3. Berechne einen Subgradienten durch  $u_i := d_B(i) - 2$  für  $i = 1, \dots, n$ .

4. Setze  $L := \max\{L, c(B) + \lambda^T u\}$  (gegenwärtig beste untere Schranke).

4.1 Gilt  $U - L < \varepsilon_1 \rightarrow$  STOP (gewünschte Güte erreicht).

4.2 Hat sich  $L$  in den letzten  $k$  Schritten nicht wenigstens um  $\varepsilon_2$  verbessert  $\rightarrow$  STOP. (Wir geben auf, weil das Verfahren stagniert. Weiterrechnen wäre nur Zeitverschwendung).

5. Setze

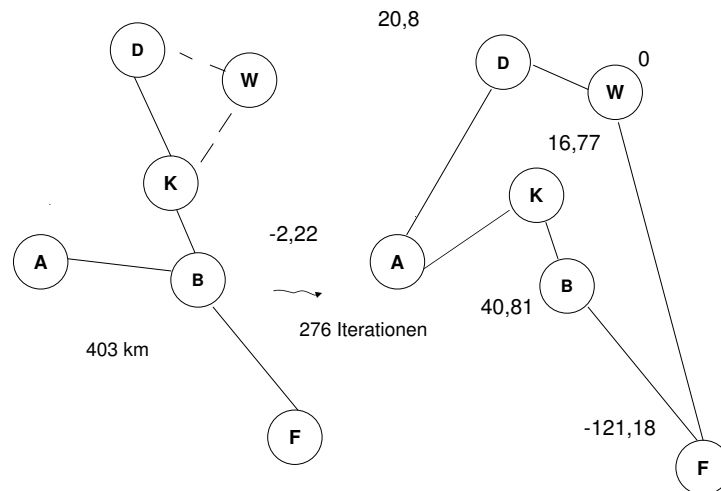
$$t := \alpha \frac{U - L}{\sum_{i=1}^n u_i^2} \quad (\text{Schrittlänge}).$$

Ist  $t < \varepsilon_3 \rightarrow \text{STOP}$ . (Die Schrittlänge ist so klein geworden, dass numerisch nichts mehr passieren wird. Wir geben auf.)

6. Setze  $\lambda := \lambda + tu$  und gehe zu 2.

ENDE.

Wir haben zu Testzwecken  $\alpha := 2$ ,  $k := 100$  und  $\varepsilon_1, \varepsilon_2, \varepsilon_3 := 10^{-3}$  gesetzt und für unser 6-Städte Beispiel (10.2) das obige Programm ausgeführt.



**Grafik/Abb. 14.6**

Das Subgradientenverfahren hat nach 276 Iterationen (d. h. 276 1-Baum Berechnungen) eine Tour geliefert. Damit ist diese Tour die optimale Lösung unseres 6-Städte Problems. Die optimalen Knotengewichte der 6 Städte sind in Abbildung 14.6 verzeichnet. (Diese Tour wurde auch durch einige unserer Primalheuristiken aus Kapitel 10 gefunden.)  $\square$

Es ist zu bemerken, dass die Optimallösung von

$$\max_{\lambda \in \mathbb{R}^n} f(\lambda)$$

im Falle des TSP (und nicht nur hier) nicht unbedingt eine Tour liefern muss, sondern es können durchaus “gaps” auftreten.

Trotzdem hat sich der Subgradientenansatz über die Lagrange-Relaxierung speziell für Probleme mittlerer Größenordnung bewährt. Verfahren dieser Art liefern i. a. sehr gute untere (bzw. bei Maximierungsproblemen obere) Schranken und kombiniert mit Branch & Bound (siehe Kapitel 15) gehören sie vielfach zu den besten Algorithmen, die wir für schwere kombinatorische Optimierung kennen.

Zum Ende dieses Abschnittes möchten wir noch zwei TSP-Beispiele vorführen und auf einige (numerische) Eigenschaften des Subgradientenverfahrens zur Lösung der Lagrange-Relaxierung hinweisen.

Die Werte (und damit unteren Schranken), die Verfahren (14.28) liefert, hängen natürlich stark von den gewählten STOP-Kriterien ab. Je nachdem, wie man in dem Spezialverfahren (14.30) für das TSP die Parameter  $k, \varepsilon_1, \varepsilon_2, \varepsilon_3, \alpha$  setzt, wird man natürlich andere Resultate erhalten. Außerdem muss man bei der 1-Baum Relaxierung einen Knoten  $v$  spezifizieren. Obwohl theoretisch die Wahl von  $v$  irrelevant ist (für alle  $v$  ist  $\max\{f(\lambda) \mid \lambda \in \mathbb{R}^n\}$  gleich), zeigt sich doch bei praktischen Rechnungen, dass unterschiedliche untere Schranken, abhängig von der Wahl, auftreten können. Wir demonstrieren das am folgenden Beispiel, von dem mir U Zaw Win berichtet hat.

**(14.31) 14-Städte Burma-Problem.** Burma besteht aus 14 Bundesländern. Zur Erinnerung an den Freiheitskampf gegen England und zur Festigung der nationalen Einheit wird in jedem Jahr die Nationalflagge auf eine Rundreise durch die Hauptstädte der 14 Bundesländer geschickt. Diese Hauptstädte sind mit ihren geographischen Koordinaten in Tabelle 14.2 aufgelistet. Die Rundreise beginnt jeweils am 30. Januar eines jeden Jahres in der Bundeshauptstadt Rangoon, führt dann durch die übrigen 13 Städte und kehrt zum Nationalfeiertag am 12. Februar nach Rangoon zurück. Der Transport der Flagge erfolgt — je nach Verkehrsverbindungen — zu Fuß (Mandalay nach Sagaing), mit dem Schiff, dem Auto oder dem Flugzeug. Wir gehen davon aus, dass die Entfernung zwischen den Städten durch die Luftliniendistanzen auf der Erde gegeben sind (man berechnet diese mit Hilfe des Programms (2.20) aus den Daten der Tabelle 14.2. Die so bestimmte Entfernungsmatrix findet sich in Tabelle 14.3.) Wie sollte die Rundreise der Nationalflagge ausgelegt werden, so dass der Reiseweg möglichst kurz ist?  $\square$

**14-Städte in Burma (Geografische Koordinaten)**

14			
5			
1.0			
1	Rangoon	16.47	96.10
2	Bassein	16.47	94.44
3	Sittwe	20.09	92.54
4	Haka	22.39	93.37
5	Myitkyina	25.23	97.24
6	Mandalay	22.00	96.05
7	Taunggyi	20.47	97.02
8	Pegu	17.20	96.29
9	Moulmein	16.30	97.38
10	Tavoy	14.05	98.12
11	Pa-an	16.53	97.38
12	Sagain	21.52	95.59
13	Loikaw	19.41	97.13
14	Magwe	20.09	94.55

**Tabelle 14.2****14-Städte in Burma (Distanzen auf Erdoberfläche)**

14												
1												
1.0												
153	510	706	966	581	455	70	160	372	157	567	342	398
422	664	997	598	507	197	311	479	310	581	417	376	
289	744	390	437	491	645	880	618	374	455	211		
491	265	410	664	804	1070	768	259	499	310			
400	514	902	990	1261	947	418	635	636				
168	522	634	910	593	18	284	239					
389	482	757	439	163	124	232						
154	406	133	508	273	355							
276	43	623	358	498								
318	898	633	761									
582	315	464										
275	221											
247												

**Tabelle 14.3**

Für das Burma-Problem (14.31) haben wir eine Rundreise mit Hilfe des Farthest-Insert-Algorithmus (10.1) (d) und dem Anfangskreis 7, 11, 14 berechnet. Diese Rundreise hat die Länge 3 322 km und ist in Abbildung 14.7 (a) dargestellt. Um die Qualität dieser Lösung abzuschätzen, haben wir untere Schranken mit Hilfe des Subgradientenverfahrens (14.30) berechnet.

Wählen wir als Knoten  $v$  des 1-Baumes den Knoten 1 (Rangoon), so erhalten wir den in Abbildung 14.7 (b) gezeigten 1-Baum der Länge 3 313.58398 km. Die "optimalen" Knotengewichte sind in Abbildung 14.6 (b) an den Knoten eingetragen.

Wählen wir als Knoten  $v$  den Knoten 8 (Pegu), so erhalten wir einen 1-Baum der Länge 3 316.98 km. Er ist in Abbildung 14.7 (c) dargestellt. Die "optimalen" Knotengewichte sind dort ebenfalls verzeichnet.

Wir sehen also, dass bei unterschiedlicher Wahl des Knotens  $v$  in der 1-Baum-Relaxierung aufgrund unserer heuristischen Stopregeln verschiedene untere Schranken berechnet werden können. Da unser Burma-Problem ganzzahlige Daten hat, wissen wir, dass die kürzeste Rundreise nicht kürzer als 3 317 km sein kann. Der maximale absolute Fehler der von uns gefundenen Rundreise beläuft sich somit auf 5 km, dies entspricht einem maximalen relativen Fehler von etwa 0,15 %.

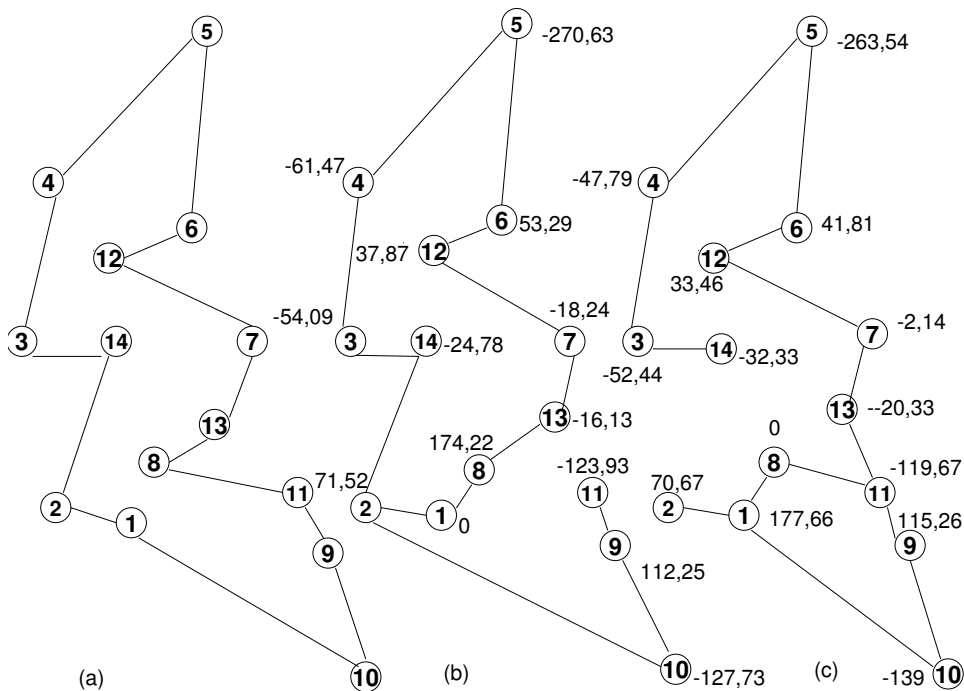
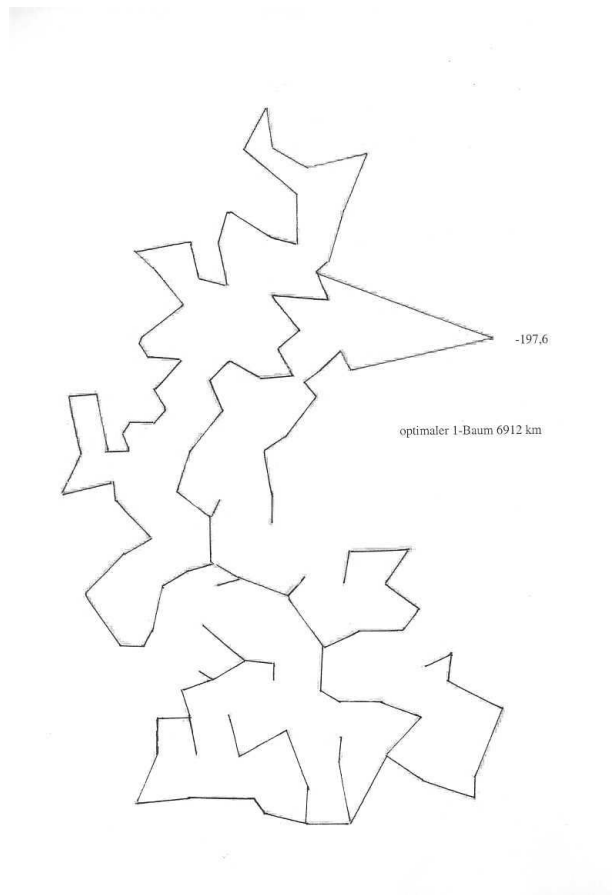


Abb. 14.7

Schließlich wollen wir noch demonstrieren, dass nicht nur bei den bisher betrachteten

kleinen “Spielbeispielen” durch die Lagrange-Relaxierung gute Resultate erzielt werden, sondern auch bei großen. Erinnern wir uns an das 120-Städte Deutschland-Problem mit Optimallösung 6 942 km, siehe Abbildung 2.3. Die (einfache) 1-Baum-Relaxierung, siehe Abbildung 14.2, brachte eine (nicht allzu gute) untere Schranke von 6 025 km (Knoten  $v = 13 = \text{Berlin}$ ).

Wir haben das Subgradientenverfahren aus (14.30) auf diese 1-Baum-Relaxierung angesetzt und laufen lassen, bis eines der STOP-Kriterien erfüllt war. Das Ergebnis ist in Abbildung 14.8 zu sehen. (An einigen Knoten sind die Lagrange-Multiplikatoren der letzten (“optimalen”) Lösung verzeichnet). Dieser 1-Baum hat eine Länge von 6912 km. Die Abweichung dieses Wertes von der Länge der Optimaltour beträgt also nur noch 0,43 %.



**Abb. 14.8**

# Kapitel 15

## Branch & Bound-Verfahren

Wir wollen uns nun mit ganzzahliger und gemischt-ganzzahliger Optimierung beschäftigen und werden eine Verfahrensklasse angeben, mit der man derartige Probleme exakt lösen kann. Sie basiert auf der Strategie, alle möglichen Lösungen auf geschickte Weise zu enumerieren und durch Benutzung von Primal- und Dualheuristiken den Rechenaufwand zu beschränken.

Die generelle Idee dieser sogenannten Branch & Bound-Verfahren ist ganz einfach. Wir betrachten ein Minimierungsproblem. Man benutze zunächst eine Primal- und eine Dualheuristik (Berechnung von Bounds), sind die Werte der gefundenen Lösung gleich, ist man fertig. Falls nicht, wird die Lösungsmenge partitioniert (Branching), und die Heuristiken werden auf die Teilmengen angewandt. Ist für eine (oder mehrere) dieser Teilmengen die für sie berechnete untere Schranke nicht kleiner als die beste bisher gefundene obere Schranke, braucht man die Lösungen in dieser Teilmenge nicht mehr zu betrachten (man erspart sich also die weitere Enumeration). Ist die untere Schranke kleiner als die beste gegenwärtige obere Schranke, muss man die Teilmenge weiter zerteilen. Man fährt so lange mit der Zerteilung fort, bis für alle Lösungsteilmengen die untere Schranke mindestens so groß ist wie die (global) beste obere Schranke.

Techniken dieser Art erscheinen in der Literatur unter einer Vielzahl von Namen wie z. B.

Branch & Bound-Verfahren,

Verfahren der implizierten Enumeration,

Divide-and-Conquer-Methoden,

Backtracking-Methoden,

Zerlegungs-Verfahren etc.

Durch geeignete Darstellung kann man zeigen, dass jede dieser Techniken ein Spezialfall der übrigen ist, d. h. im Prinzip tun diese Verfahren alle das gleiche, sie unterscheiden sich nur bezüglich der jeweils angewendeten Strategie bzw. Philosophie. Manche Autoren machen sehr feine Unterschiede zwischen diesen Methoden, wir wollen jedoch alle Verfahren, die Enumerationstechniken benutzen, **Branch & Bound-Algorithmen** nennen.

Das Prinzip der Branch & Bound-Verfahren ist trivial, theoretisch sind sie nicht sonderlich interessant, aber es scheint so, dass man ohne sie einfach nicht auskommt. Bei den meisten in der Praxis erfolgreichen Methoden zur Lösung schwieriger kombinatorischer oder ganzzahliger Optimierungsprobleme wird irgendwann einmal Branch & Bound eingesetzt.

Noch eine Warnung! So simpel die Idee ist, so schwierig ist es i. a. Branch & Bound-Methoden effizient zu implementieren. Eine Hauptschwierigkeit besteht darin, „gute“ Zerlegungstechniken und einfache Datenstrukturen zu erfinden, die die effiziente Abarbeitung und Durchsuchung der einzelnen Teilmengen ermöglichen. Ganz allgemein kann man diesen Algorithmentyp formalisiert wie folgt darstellen.

### (15.1) Allgemeines Branch & Bound-Verfahren.

**Input:** Eine Menge  $L$  von zulässigen Lösungen (implizit gegeben) und eine Zielfunktion  $c : L \rightarrow \mathbb{K}$ , die jedem  $S \in L$  einen Wert  $c(S)$  zuordnet.

**Output:** Lösung von

$$(P) \quad \max_{S \in L} c(S).$$

**Annahme:** Wir nehmen an, dass es ein „Universum“  $\bar{L}$  gibt mit  $L \subseteq \bar{L}$ , dass  $c(S)$  wohldefiniert ist für alle  $S \in \bar{L}$  und dass das relaxierte Problem (RP)  $\max_{S \in \bar{L}} c(S)$  „einfach“ lösbar ist. Ferner soll die Lösung von (RP) in „nicht-trivialem“ Zusammenhang mit der Lösung von (P) stehen. Wir nehmen also an, dass wir P relaxieren können und (durch die Lösung von (RP)) eine effiziente Dualheuristik zur Verfügung steht.

1. Berechne eine obere Schranke  $U$  für (P). Man wende z. B. eine Primalheuristik an oder setze  $U := -\infty$ .
2. Setze  $\mathcal{K} = \{L\}$  ( $\mathcal{K}$  ist die Menge der Kandidatenmengen)
3. Falls  $\mathcal{K} = \emptyset \rightarrow$  STOP (die beste gegenwärtige Lösung ist eine Optimallösung von (P)). Andernfalls wähle eine Menge  $M \in \mathcal{K}$ . (Diesen Schritt nennt man **Branching** oder **Verzweigung**. Man muss sich hier genau überlegen, welche der noch nicht bearbeiteten Kandidatenmengen untersucht werden soll.)
4. **Bounding:** Wähle eine „geeignete“ Menge  $\bar{M} \subseteq \bar{L}$  mit  $M \subseteq \bar{M}$ , so dass das relaxierte Problem

$$(RPM) \quad \max_{S \in \bar{M}} c(S)$$

„einfach“ zu lösen ist. (Das auf  $M$  eingeschränkte Ursprungsproblem wird also relaxiert. Wegen  $M \subseteq \overline{M}$  bildet der Wert  $c(S^*)$  der Optimallösung  $S^*$  von (RPM) eine obere Schranke (**Bound**) für den Wert der Optimallösung von

$$(PM) \quad \max_{S \in M} c(S)$$

5. Ist  $c(S^*) \leq U$ , so hat keine Lösung aus  $\overline{M}$  und somit auch keine aus  $M$  einen besseren Wert als die beste bereits bekannte Lösung von (P), d. h. die Lösungen aus  $M$  brauchen nicht weiter betrachtet zu werden. Entferne also  $M$  aus  $\mathcal{K}$  und gehe zu 3.
6. Ist  $S^* \in L$  und  $c(S^*) > U$ , so ist eine zulässige Lösung für (P) gefunden, deren Wert besser als  $U$  ist. Da die beste Lösung aus  $M$  gefunden ist, braucht  $M$  nicht mehr weiter untersucht zu werden.

Entferne  $M$  aus  $\mathcal{K}$ , setze  $U := c(S^*)$  und gehe zu 3.

7. **Separation** (Zerlegung): Es gilt also  $c(S^*) > U$  und  $S^* \notin L$ . In diesem Falle war die Berechnung der unteren Schranke nutzlos, da wir keine verwertbare Aussage machen können. Wir zerlegen daher  $M$  in „geeignete“ kleinere Mengen  $M_1, M_2, \dots, M_k$ , entfernen  $M$  aus  $\mathcal{K}$ , fügen  $M_1, M_2, \dots, M_k$  zu  $\mathcal{K}$  hinzu und gehen zu 3.

□

Falls  $M$  in den Schritten 5. oder 6. verworfen wird, so sagt man, dass  $M$  ausgelotet (fathomed) worden ist.

Ist man in Schritt 7 angelangt, so ist es häufig nützlich zu versuchen, aus  $S^*$  eine zulässige Lösung  $S$  von  $M$  zu konstruieren, um damit (durch eine Primalheuristik) die obere Schranke verbessern zu können.

Branch & Bound-Verfahren unterscheiden sich vor allem durch die Wahl verschiedener Relaxierungen  $\overline{L}$ , durch unterschiedliche Zerlegungsstrategien in 7. und durch die Auswahl der nächsten Kandidatenmenge in 3. Mit jedem Branch & Bound-Verfahren kann man einen sogenannten **Branch & Bound-Baum** assoziieren, der über die Verzweigungen Auskunft gibt. Ein derartiger Verzweigungsbaum ist in Abbildung 15.1 gezeigt.

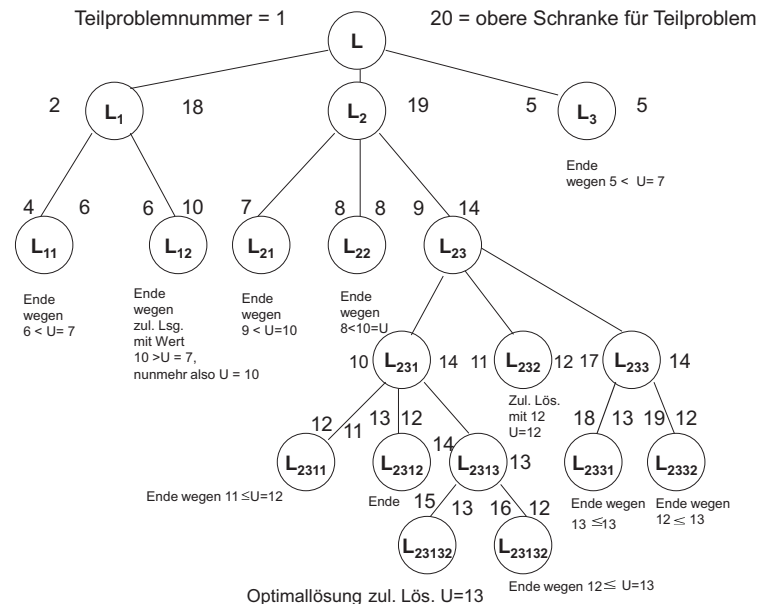


Abb. 15.1

Bei der Lösung des Teilproblems 1 wird eine obere Schranke mit Wert 20 gefunden und durch eine Heuristik eine zulässige Lösung von  $(P)$  mit Wert 7. Links vom jeweiligen Knoten ist die Reihenfolge vermerkt, in welcher die Teilprobleme gelöst wurden, rechts neben den Knoten die jeweils berechnete obere Schranke. Bei der Lösung des 6. Teilproblems wird eine zulässige Lösung mit Wert 10 gefunden. Das 15. Teilproblem liefert die Optimallösung.

Das Hauptproblem bei der Implementierung von Branch & Bound-Verfahren besteht darin, dass man sich alle Kandidatenmengen merken muss. Dies muss auf geschickte Weise geschehen, um nicht zu viel suchen zu müssen und nicht zu viel Speicherplatz zu vergeuden. Außerdem muss man dafür Sorge tragen, dass der Baum nicht zu groß wird und der Speicherplatz ausreicht. Man muss also in vielerlei Hinsicht Kompromisse schließen, um ein praktikables Verfahren zu entwerfen.

Das Branch & Bound-Verfahren ist so organisiert, dass in jedem Schritt folgendes gilt

*Eine zulässige Lösung  $S$  gehört entweder zu einer bereits verworfenen Menge oder sie ist in genau einer der Kandidatenmengen enthalten.*

Falls also die Menge der zulässigen Lösungen endlich ist und jedes Subproblem (RPM) in endlicher Zeit gelöst wird, bricht das Verfahren nach endlich vielen Schritten ab, wenn wir die (sinnvolle) Konvention treffen, dass einelementige Mengen  $M$  nicht relaxiert werden.

Wir wollen uns nun überlegen, wie wir die Branch & Bound-Technik bei ganzzahligen Programmierungsproblemen einsetzen können.

**(15.2) Das Branch & Bound-Verfahren von Dakin für gemischt-ganzzahlige Programme.****Input:** Gemischt-ganzzahliges Programm mit rationalen Daten

$$(MIP^=) = (MIP_0^=) \quad \begin{array}{l} \max c^T x \\ Ax = b \\ x \geq 0, \\ x_i \text{ ganzzahlig für alle } i \in N_1 \end{array}$$

$$P_0 := \{x \mid Ax = b, x \geq 0, x_i \text{ ganzzahlig für alle } i \in N_1\}$$

1. Setze:

$$\begin{array}{ll} \text{untere Schranke} & U := -\infty \\ \text{Kandidatenmengen} & \mathcal{K} := \{P_0\} \\ \text{Zählindex} & k := 0 \\ \text{gegenwärtig beste Lösung} & \bar{x} \text{ (am Anfang leer)} \end{array}$$

2. Falls  $\mathcal{K} = \emptyset$ , dann STOP.Ergebnis: Falls  $U = -\infty$ , so hat  $(MIP^=)$  keine Lösung.Falls  $U > -\infty$ , so ist  $U$  der Optimalwert und  $\bar{x}$  eine Optimallösung.3. **Branching:** Wähle eine Menge  $P_j \in \mathcal{K}$ .4. **Bounding:** Löse das durch Weglassen der Ganzzahligkeitsbedingung entstehende lineare Programm  $(LMIP_j^=) \max c^T x, x \in LP_j$ .5. Ist  $LP_j = \emptyset$  oder  $(LMIP_j^=)$  unbeschränkt, so sei  $c^* = -\infty$ , andernfalls sei  $x^*$  die Optimallösung von  $(LMIP_j^=)$  und  $c^*$  der Optimalwert.6. Ist  $c^* \leq U$ , so entferne  $P_j$  aus  $\mathcal{K}$ .7. Ist  $c^* > U$  und  $x_i^* \in \mathbb{Z} \forall i \in N_1$ , so setze  $U := c^*$ ,  $\bar{x} := x^*$  und entferne  $P_j$  aus  $\mathcal{K}$ .8. **Separation:** Im Falle  $c^* > U$  und  $x_i^* \notin \mathbb{Z}$  für einige  $i \in N_1$  wähle ein  $i \in N_1$  mit  $x_i^* \notin \mathbb{Z}$ . Entferne  $P_j$  aus  $\mathcal{K}$  und setze

$$\begin{array}{ll} P_{k+1} & := P_j \cap \{x \mid x_i \leq \lfloor x_i^* \rfloor\}, \\ P_{k+2} & := P_j \cap \{x \mid x_i \geq \lceil x_i^* \rceil\}, \\ k & := k + 2 \end{array}$$

und gehe zu 2.

□

**(15.3) Bemerkung.** Die Daten des gemischt-ganzzahligen Programms ( $MIP^=$ ) seien — wie in (15.2) verlangt — rational.

- (a) Ist ( $LMIP^=$ ) unzulässig oder unbeschränkt, so bricht das Verfahren von Dakin ab und zeigt an, dass ( $MIP^=$ ) keine Lösung oder keine endliche Optimallösung besitzt.
- (b) Hat ( $LMIP^=$ ) ein endliches Optimum, und ist  $P_0 \neq \emptyset$ , so findet das Verfahren von Dakin nach endlich vielen Schritten eine Optimallösung von ( $MIP^=$ ), da alle ganzzahligen Werte der  $x_i, i \in N_1$ , durchsucht werden.
- (c) Hat ( $LMIP^=$ ) ein endliches Optimum, und ist  $P_0 = \emptyset$ , so kann (theoretisch) durch Einführung einer unteren Schranke für den Zielfunktionswert ein endlicher Abbruch erzwungen werden.

Folglich ist das Dakin-Verfahren ein endliches Verfahren zur Lösung gemischt-ganzzahliger Programme.  $\square$

**(15.4) Bemerkung.** Das ursprüngliche lineare Programm (LP-Relaxierung von ( $MIP^=$ )) wird bei den verschiedenen Verzweigungsschritten jeweils um sehr einfache Ungleichungen, nämlich um obere bzw. untere Schranken der Variablen erweitert. Für diesen Spezialfall kennen wir bereits eine Variante des Simplexalgorithmus (dualer Simplexalgorithmus mit oberen Schranken, (upper-bounding-technique), siehe “Optimierungsmethoden I”, Kapitel 10), die es erlaubt, die oberen Schranken durch Transformationen im Tableau implizit zu berücksichtigen, ohne die Ungleichungen explizit hinzuzufügen.  $\square$

**(15.5) Beispiel.** Wir betrachten das folgende rein-ganzzahlige Problem

$$\begin{aligned} \max & -7x_1 - 3x_2 - 4x_3 \\ & x_1 + 2x_2 + 3x_3 - x_4 = 8 \\ & 3x_1 + x_2 + x_3 - x_5 = 5 \\ & x_i \geq 0, \quad i = 1, \dots, 5 \\ & x_i \text{ ganzzahlig, } i = 1, \dots, 5 \end{aligned}$$

$P_0$  sei die Lösungsmenge dieses Programms. Eine optimale Lösung des zugehörigen linearen Programms ist gegeben durch  $x_3 = x_4 = x_5 = 0$  und

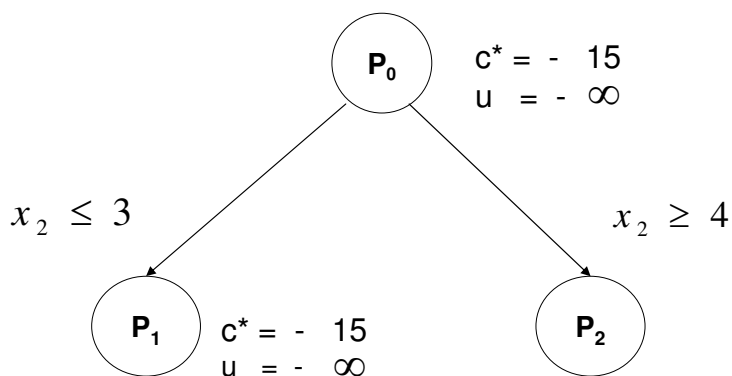
$$\begin{aligned} x_1 &= \frac{2}{5} \\ x_2 &= \frac{19}{5} \\ c^* &= -\frac{71}{5}. \end{aligned}$$

Wir merken uns bei jedem Programm den Wert der Optimallösung und den gegenwärtigen Wert der unteren Schranke  $U$ . Da das Problem rein ganzzahlig ist, muss der Zielfunktionswert rein ganzzahlig sein, wir können also den Wert des LP-Optimums abrunden

und erhalten damit eine untere Schranke für alle ganzzahligen Lösungen des gegenwärtig gelösten LP. Für  $P_0$  erhalten wir

$$\begin{aligned}c^* &= -15 \\U &= -\infty\end{aligned}$$

Zur Verzweigung wählen wir die Variable  $x_2$  und fügen einerseits  $x_2 \leq 3$  und andererseits  $x_2 \geq 4$  zu den Nebenbedingungen von  $P_0$  hinzu. Wir erhalten auf diese Weise neue Kandidatenmengen  $P_1, P_2$ .  $P_0$  wird aus  $\mathcal{K}$  entfernt. Wir merken uns, was wir getan haben bildlich in dem in Abbildung 15.2 gezeigten Branch & Bound-Baum.



**Abb. 15.2**

Unsere Kandidatenmenge  $\mathcal{K}$  enthält nunmehr die Mengen  $P_1, P_2$ . Wir entscheiden uns zur Lösung von  $\max c^T x, x \in P_1$ . Die Optimallösung der LP-Relaxierung dieses Problems ist die folgende:  $x_4 = x_5 = 0$  und

$$\begin{aligned}x_1 &= \frac{1}{2} \\x_2 &= 2 \\x_3 &= \frac{1}{2} \\c^* &= -\frac{29}{2} \quad (\text{Abrunden ergibt } c^* = -15).\end{aligned}$$

Wir tragen dieses Resultat in den Branch & Bound-Baum 15.2 ein. Da die obige Lösung nicht ganzzahlig ist, müssen wir  $P_1$  weiter zerlegen. Wir wählen dazu die Variable  $x_1$ , wobei wir zu den Restriktionen von  $P_1$  einerseits  $x_1 \geq 1$ , andererseits  $x_1 \leq 0$  hinzufügen können. Wir entfernen  $P_1$  aus  $\mathcal{K}$  und fügen die neuen Kandidatenmengen  $P_3$  und  $P_4$  zu  $\mathcal{K}$  hinzu, siehe Abbildung 15.3.

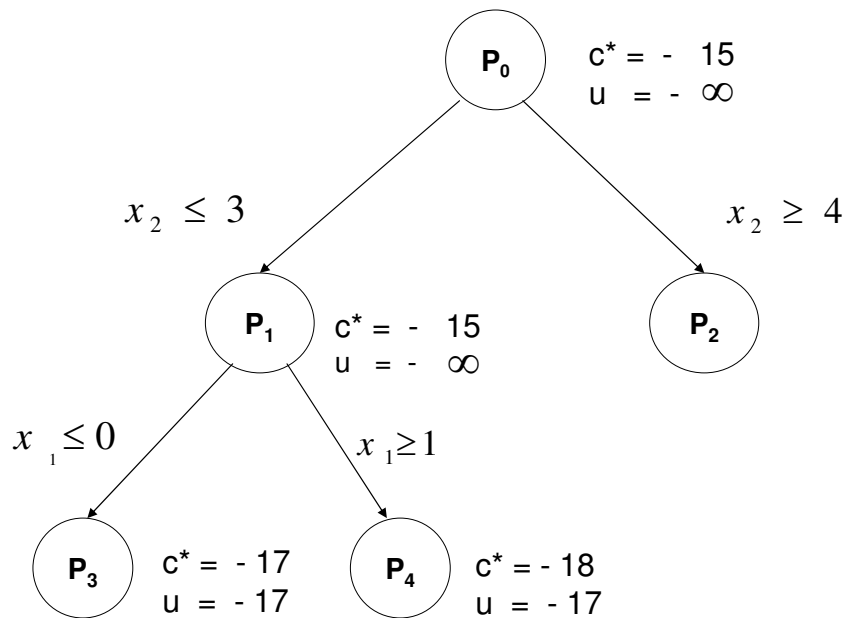


Abb. 15.3

Die Kandidatenmenge enthält nun  $P_2, P_3$  und  $P_4$ . Wir lösen  $P_3$  und erhalten  $x_1 = 0$ ,  $x_5 = 0$  und

$$\begin{aligned} x_2 &= 3 \\ x_3 &= 2 \\ x_4 &= 4 \\ c^* &= -17. \end{aligned}$$

Die optimale Lösung der LP-Relaxierung von  $P_3$  ist also ganzzahlig. Wir haben die untere Schranke verbessert und setzen  $U := -17$ , außerdem ist die Menge  $P_3$  vollständig erledigt, wir können sie aus  $\mathcal{K}$  eliminieren.

Es verbleiben  $P_2$  und  $P_4$  in  $\mathcal{K}$ . Wir entscheiden uns zur Bearbeitung von  $P_4$ . Das Optimum der LP-Relaxierung von  $\max c^T x, x \in P_4$  ist:  $x_4 = 0$  und

$$\begin{aligned} x_1 &= 1 \\ x_2 &= 3 \\ x_3 &= \frac{1}{3} \\ x_5 &= \frac{4}{3} \\ c^* &= -\frac{52}{3} \quad (\text{Abrunden ergibt } c^* = -18). \end{aligned}$$

Also ist der Wert des LP-Optimums kleiner als der Wert der gegenwärtig besten ganzzahligen Lösung.  $P_4$  ist damit auch vollständig erledigt. Die oben erzielten Resultate sind auch in Abbildung 15.3 eingetragen.

In der Kandidatenmenge ist nur noch  $P_2$ . Die Optimallösung der LP-Relaxierung von  $\max c^T x, x \in P_2$  ist  $x_3 = x_5 = 0$  und

$$\begin{aligned} x_1 &= \frac{1}{3} \\ x_2 &= 4 \\ x_4 &= \frac{1}{3} \\ c^* &= -\frac{43}{3} \quad (\text{Abrunden ergibt } c^* = -15). \end{aligned}$$

Die Lösung ist weder ganzzahlig noch unterschreitet der Optimalwert die gegenwärtig beste Schranke. Wir verzweigen bezüglich  $x_1$  und fügen einerseits  $x_1 \leq 0$  und andererseits  $x_1 \geq 1$  zu  $P_2$  hinzu. Wir erhalten neue Mengen  $P_5, P_6$ , die wir zu  $\mathcal{K}$  addieren.  $P_2$  wird aus  $\mathcal{K}$  entfernt. Der gegenwärtige Branch & Bound-Baum ist in Abbildung 15.4 gezeigt.

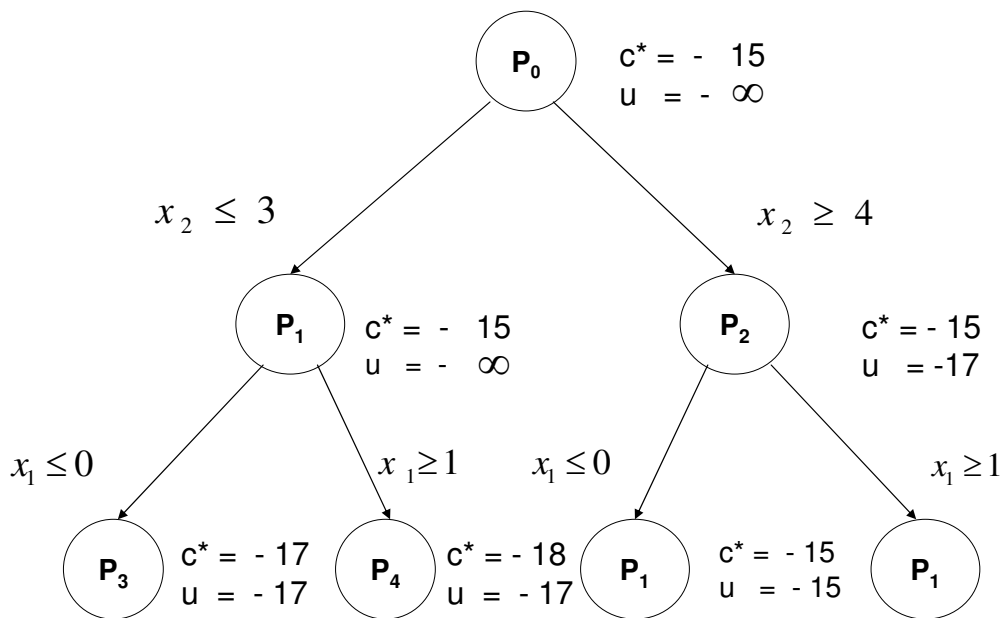


Abb. 15.4

Wir lösen die LP-Relaxierung von  $\max c^T x, x \in P_5$  und erhalten

$$\begin{aligned} x_2 &= 5 \\ x_4 &= 2 & x_i = 0 \text{ sonst} \\ c^* &= -15. \end{aligned}$$

Diese Lösung ist ganzzahlig, und wir sind mit unserem Algorithmus fertig, da  $P_6$  keine bessere ganzzahlige Lösung enthalten kann, denn eine ganzzahlige Lösung in  $P_2$  hat höchstens den Wert  $-15$ . Damit ist eine optimale Lösung des Ausgangsproblems gefunden.  $\square$

Ich glaube, dass damit das Prinzip der Branch & Bound-Verfahren klar ist. Wichtig ist, eine Relaxierung zu wählen, die schnell lösbar ist und gute Schranken liefert, also eine gute duale Heuristik auszuwählen. Im Dakin-Verfahren haben wir einfach die kanonische LP-Relaxierung gewählt.

Wir wollen nun noch ein einfaches Beispiel eines Branch & Bound-Verfahrens angeben und zwar eines für das asymmetrische TSP.

Der Großvater aller Branch & Bound-Algorithmen ist das Verfahren von Little, Murty, Sweeney & Karel zur Lösung asymmetrischer Travelling-Salesman-Probleme, das 1963 veröffentlicht wurde und in dem der Name Branch & Bound-Algorithmus geprägt wurde. Weil dieses Verfahren so einfach ist, wollen wir es hier kurz besprechen, jedoch darauf hinweisen, dass es zur Lösung großer Probleme nicht besonders geeignet ist.

### (15.6) Das Verfahren von Little et al für das asymmetrische Travelling-Salesman-Problem.

**Input:** asymmetrisches TSP auf  $n$  Städten gegeben durch Distanzmatrix  $c = (c_{ij})$  mit  $c_{ii} = \infty$   $i = 1, \dots, n$  und  $c_{ij} \geq 0$ .

**Output:** Optimallösung des gegebenen TSPs.

**Verfahrensweise:** Zeilen- und Spaltenreduktion der Matrix  $C$  und sukzessive Erzeugung von Teilproblemen.

Ein Teilproblem ist definiert durch

$$P(EINS, NULL, I, J, C, L, U)$$

wobei

- $EINS$  = die Menge der auf 1 fixierten Bögen
- $NULL$  = die Menge der auf 0 fixierten Bögen
- $C$  = die Distanzmatrix des Teilproblems
- $I$  = die Zeilenindizes von  $C$
- $J$  = die Spaltenindizes von  $C$
- $L$  = die untere Schranke für das Teilproblem
- $U$  = die obere Schranke für das Globalproblem

## 1. Definiere das Anfangsproblem

$$P(\emptyset, \emptyset, \{1, \dots, n\}, \{1, \dots, n\}, C, 0, +\infty),$$

wobei  $C$  die Ausgangsmatrix ist, und setze dieses Problem auf die Problemliste.

2. Sind alle Teilprobleme gelöst  $\rightarrow$  STOP. Andernfalls wähle ein ungelöstes Teilproblem  $P(\text{EINS}, \text{NULL}, I, J, C, L, U)$  aus der Problemliste.

**Regel:** Wähle das Problem mit kleinster unterer Schranke oder das Problem, bei dem EINS am meisten Elemente enthält. Bei (a) erhält man hoffentlich eine optimale Tour, bei (b) hat man das Problem schnell abgearbeitet.

3. **Bounding:**a) **Zeilenreduktion:**

FOR all  $i \in I$  DO:

Berechne das Minimum der  $i$ -ten Zeile von  $C$

$$c_{ij_0} = \min_{j \in J} c_{ij}$$

Setze  $c_{ij} := c_{ij} - c_{ij_0} \forall j \in J$

und  $L := L + c_{ij_0}$

END;

b) **Spaltenreduktion**

FOR all  $j \in J$  DO:

Berechne das Minimum der  $j$ -ten Spalte von  $C$

$$c_{i_0j} = \min_{i \in I} c_{ij}$$

Setze  $c_{ij} := c_{ij} - c_{i_0j} \forall i \in I$  und  $L := L + c_{i_0j}$

END;

c) Ist  $L \geq U$ , so entferne das gegenwärtige Teilproblem aus der Problemliste und gehe zu 2.d) Definiere den Nulldigraphen  $G_0 = (V, A)$  mit  $A := \text{EINS} \cup \{(i, j) \in I \times J \mid c_{ij} = 0\}$ e) Enthält  $G_0$  keine Tour so gehe zu 4.f) Enthält  $G_0$  eine Tour, so hat die Tour die Länge  $L$ .

$f_1$ ) Entferne alle Teilprobleme aus der Problemliste, deren untere Schranke größer gleich  $L$  ist.

$f_2$ ) Setze in allen noch nicht gelösten Teilproblemen  $U := L$ .

$f_3$ ) Gehe zu 2.

(i. a. wird der Nulldigraph  $G_0$  nicht berechnet sondern so lange enumeriert bis nur noch (2, 2)-Probleme übriggeblieben sind.)

**4. Branching:**

a) Wähle nach einem Plausibilitätskriterium einen Bogen  $(i, j)$ , der 0 oder 1 gesetzt wird.

z. B. Definiere  $u(i, j) := \min\{c_{ik} \mid k \in J \setminus \{j\}\} + \min\{c_{kj} \mid k \in I \setminus \{i\}\} - c_{ij}$  d. h.  $u(i, j)$  sind die Zusatzkosten (Umwegkosten), die entstehen, wenn wir von  $i$  nach  $j$  gehen, ohne  $(i, j)$  zu benutzen.

**Branching-Regel:** Wähle  $(i, j) \in I \times J$ , so dass

$$c_{ij} = 0 \text{ und} \\ u(i, j) = \max\{u(p, q) \mid c_{pq} = 0\}$$

(Wähle also einen Bogen, der zu möglichst hohen Umwegkosten führt)

b) Definiere die neuen Teilprobleme

b<sub>1</sub>)  $P(\text{EINS} \cup \{(i, j)\}, \text{NULL}, I \setminus \{i\}, J \setminus \{j\}, C', L, U)$

wobei  $C'$  aus  $C$  durch Streichen der Zeile  $i$  und der Spalte  $j$  entsteht und  $c'_{ji} := \infty$  (zur Vermeidung des Kurzzyklus  $\langle i, j \rangle$ ).

b<sub>2</sub>)  $P(\text{EINS}, \text{NULL} \cup \{(i, j)\}, I, J, C'', L, U)$ ,

wobei  $C''$  aus  $C$  entsteht durch  $c_{ij} := \infty$ .

Füge diese Teilprobleme zur Problemliste hinzu und gehe zu 2.

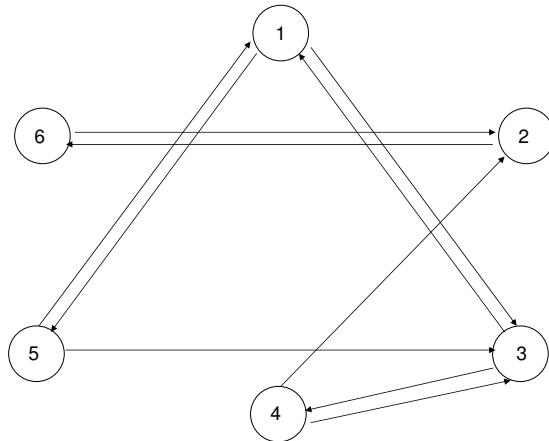
□

Das Hauptproblem bei (15.6) ist — wie immer — die Buchhaltung der Teilprobleme. Die Anzahl der Touren in 4. b<sub>1</sub>) ist kleiner als die in 4. b<sub>2</sub>), die Wahl von  $(i, j)$  wird so getroffen, dass “hoffentlich” in 4. b<sub>1</sub>) eine optimale Tour ist.

**(15.7) Beispiel.** Wir beginnen mit  $P = (\emptyset, \emptyset, \{1, \dots, 6\}, \{1, \dots, 6\}, C, 0, +\infty)$

		Zeilenmin												
							↓							
Anfangsmatrix	∞	5	1	2	1	6	1	∞	4	0	2	0	6	
	6	∞	6	3	7	2	2	4	∞	4	1	5	0	
	1	4	∞	1	2	5	1	→	0	3	∞	0	1	4
	4	3	3	∞	5	4	3		1	0	0	∞	2	1
	1	5	1	2	∞	5	1		0	4	0	1	∞	4
	6	2	6	4	5	∞	2		4	0	4	2	3	∞
							-----							
							10	Spaltenmin	0	0	0	0	0	0

Die Zeilenreduktion (Schritt 3. a)) ergibt insgesamt einen Wert von 10, die Spaltenreduktion erbringt nichts, da nach der Zeilenreduktion bereits jede Spalte eine Null enthält. Die untere Schranke ist somit  $L = 10$ . Wir können nun den Nulldigraphen  $G_0$  definieren. Er ist in Abbildung 15.5 dargestellt (alle Bögen mit dem Wert 0 sind eingezeichnet.).  $G_0$  enthält keine Tour



**Nulldigraph**

**Abb. 15.5**

Wir verwenden die in Schritt 4. a) angegebene Branchingregel und wählen den Bogen  $(i, j) = (2, 6)$  mit  $u(2, 6) = 2$ . Das Ausgangsproblem ist nun erledigt. Wir definieren zwei neue Teilprobleme.

**Erstes neues Problem:**  $P = (\{(2, 6)\}, \emptyset, \{1, 3, 4, 5, 6\}, \{1, 2, 3, 4, 5\}, C', 10, \infty)$

		1	2	3	4	5			1	2	3	4	5
	1	$\infty$	4	0	2	0	0	$\infty$	4	0	2	0	0
	3	0	3	$\infty$	0	1	0	0	3	$\infty$	0	1	0
$C'$	=	4	1	0	0	$\infty$	2	0	1	0	0	$\infty$	2
		5	0	4	0	1	$\infty$	0	0	4	0	1	$\infty$
		6	4	$\infty$	4	2	3	2	2	$\infty$	2	0	1

Die untere Schranke (nach Zeilen- und Spaltenreduktion) für dieses Teilproblem ist  $L = 12$ .

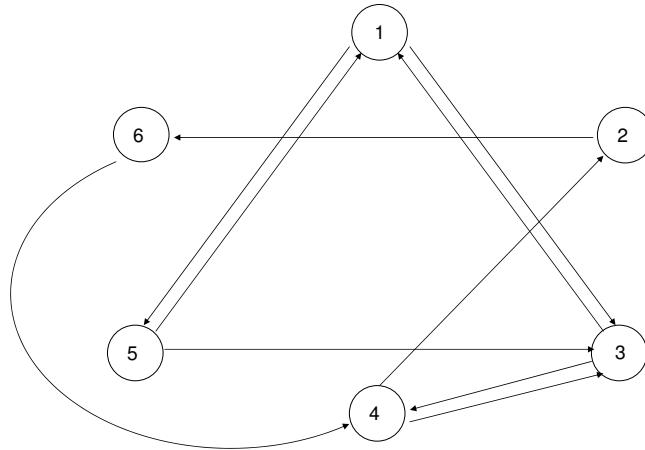


Abb. 15.6

Zweites neues Problem:  $P = (\emptyset, \{(2, 6)\}, \{1, 2, 3, 4, 5, 6\}, \{1, 2, 3, 4, 5, 6\}, C'', 10, \infty)$

$$C'' = \begin{array}{cccccc|cccc}
 \infty & 4 & 0 & 1 & 0 & 5 & 0 & \infty & 4 & 0 & 1 & 0 & 4 \\
 4 & \infty & 4 & 1 & 5 & \infty & 1 & 3 & \infty & 3 & 0 & 4 & \infty \\
 0 & 3 & \infty & 0 & 1 & 4 & 0 & 0 & 3 & \infty & 0 & 1 & 3 \\
 1 & 0 & 0 & \infty & 2 & 1 & 0 & 1 & 0 & 0 & \infty & 2 & 0 \\
 0 & 4 & 0 & 1 & \infty & 4 & 0 & 0 & 4 & 0 & 1 & \infty & 3 \\
 4 & 0 & 4 & 2 & 3 & \infty & 0 & 4 & 0 & 4 & 2 & 3 & \infty \\
 0 & 0 & 0 & 0 & 0 & 1 & u = 12 & & & & & & 
 \end{array} \longrightarrow \begin{array}{cccccc}
 \infty & 4 & 0 & 1 & 0 & 4 \\
 3 & \infty & 3 & 0 & 4 & \infty \\
 0 & 3 & \infty & 0 & 1 & 3 \\
 1 & 0 & 0 & \infty & 2 & 0 \\
 0 & 4 & 0 & 1 & \infty & 3 \\
 4 & 0 & 4 & 2 & 3 & \infty
 \end{array}$$

Die untere Schranke für das zweite Problem ist ebenfalls 12. Das Verfahren wird auf diese Weise fortgeführt. Insgesamt erhält man den in Abbildung 15.7 dargestellten **Branch & Bound**-Baum. □

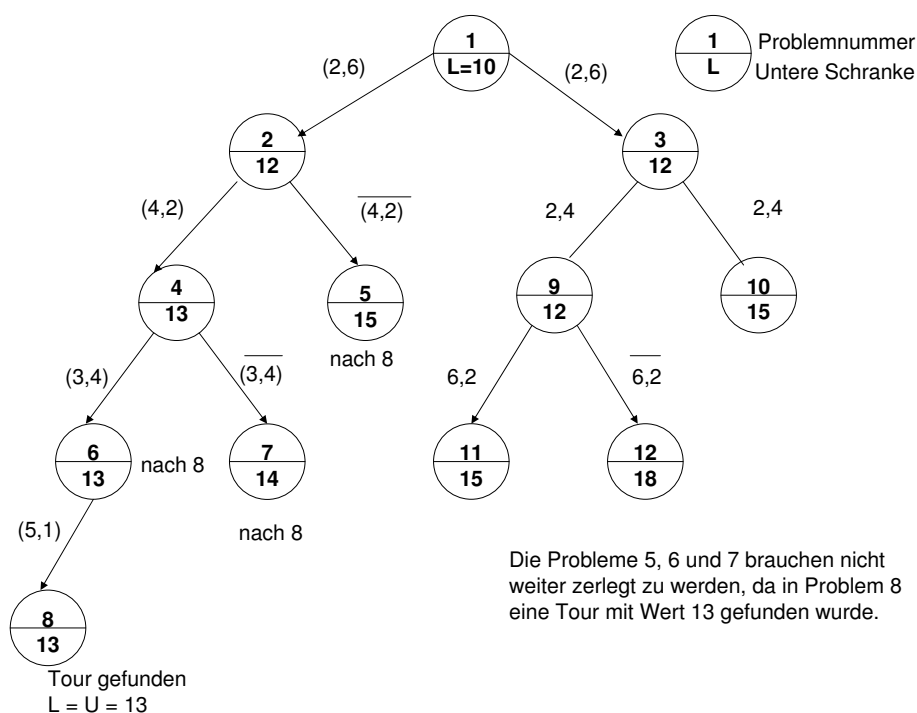


Abb. 15.7

Das Verfahren (15.6) ist sehr simpel und (daher) in der Praxis nicht sonderlich effizient. Wesentlich bessere Ergebnisse erzielt man bei Benutzung der Zuordnungsrelaxierung (14.7). Der Leser sollte nun in der Lage sein, ein auf der Zuordnungsrelaxierung basierendes Branch & Bound-Verfahren selbst zu entwerfen.

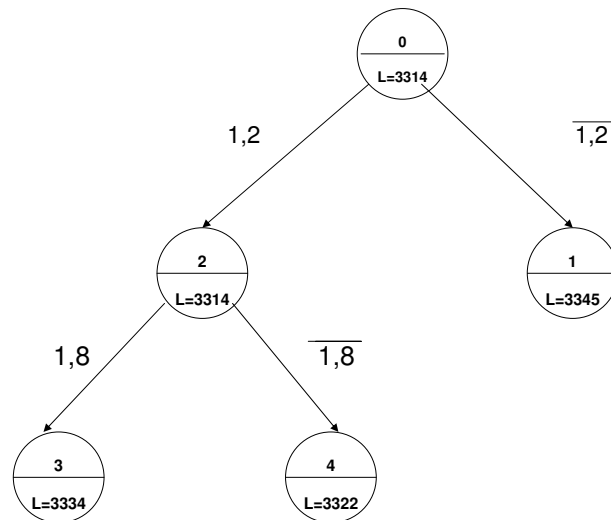
Recht gute Erfolge sind für das symmetrische TSP bei Benutzung der 1-Baum Lagrange-Relaxierung (siehe (14.30)) zu verzeichnen. Mit einem derartigen Branch & Bound-Verfahren sind Travelling-Salesman-Probleme mit bis zu 150 Städten gelöst worden. Wir deuten hier kurz an, wie man das 14-Städte Burma-Problem (14.31) auf diese Weise lösen kann.

Wir haben als Knoten  $v$  der 1-Baum-Relaxierung den Knoten  $v = 1$  (Rangoon) gewählt. Das Subgradientenverfahren (14.30) liefert den in Abbildung 14.6 (b) gezeigten "optimalen" 1-Baum mit Wert 3 313.58 km. Damit wissen wir, dass die kürzeste Rundreise (aufgrund der Ganzzahligkeit der Daten) nicht länger als  $L := 3\,314$  km sein kann.

Wir entscheiden uns bezüglich der Kante 1,2 zu verzweigen. Das heißt, wir definieren ein erstes Teilproblem, bei dem alle Rundreisen, die Kante 1,2 nicht enthalten, und ein zweites, bei dem alle Rundreisen die Kante 1,2 enthalten. Statt nun die kürzeste Rundreise für jedes dieser zwei Teilprobleme zu berechnen, bestimmen wir mit Verfahren (14.30) neue "optimale" 1-Bäume (unter den gerade festgelegten Nebenbedingungen). (14.30) liefert, dass der optimale 1-Baum von Burma, der die Kante 1,2 nicht enthält, mindestens die Länge 3 345 km hat. Da wir eine Rundreise der Länge 3 322 km kennen (siehe Abbildung

14.6 (a)), kann keine in Teilproblem 1 enthaltene Rundreise optimal sein. Also können wir folgern, dass jede kürzeste Rundreise durch Burma die Kante 1,2 enthalten muss. Da der (global) optimale 1-Baum die Kante 1,2 enthält, brauchen wir keine Neuberechnung des 1-Baumes für Teilproblem 2 vorzunehmen, sondern verzweigen erneut.

Wir zerlegen nun Teilproblem 2 in Teilproblem 3 (alle diejenigen Rundreisen, die die Kante 1,8 enthalten) und Teilproblem 4 (alle diejenigen Rundreisen, die die Kante 1,8 nicht enthalten). Man beachte, dass bei Teilproblem 3 durch die Festlegung, dass alle Touren die Kanten 1,8 und 1,2 enthalten müssen, impliziert wird, dass die Rundreisen weder Kanten  $1, i$  ( $i \neq 1, 2, 8$ ) noch die Kante 2,8 enthalten können. Diese implizierten Restriktionen benutzt man natürlich bei der Berechnung des optimalen 1-Baumes. Für Teilproblem 3 bestimmt man also einen optimalen 1-Baum unter allen (auch den implizierten) Restriktionen und erhält durch (14.30) einen 1-Baum mit Wert 3 335 km. Teilproblem 3 kann also keine optimale Tour enthalten. Für Teilproblem 4 erhält man einen 1-Baum mit Wert 3 321.9 km. Aufgrund der Ganzzahligkeit der Daten und der Tatsache, dass man eine Tour der Länge 3 322 km kennt, kann man nun schließen, dass diese Tour optimal ist. Der Branch & Bound-Baum dieses Vorgehens ist in Abbildung 15.8 dargestellt.



**Abb. 15.8**

Insgesamt sind also nur 4 Aufrufe des Subgradientenverfahrens notwendig gewesen, um die Optimalität der durch FARTHEST INSERT gefundenen Burma-Rundreise nachzuweisen.

# Kapitel 16

## Theorie der ganzzahligen und gemischt-ganzzahligen Optimierung

Wir möchten hier auf Kapitel 1 der Vorlesung zurückverweisen. In diesem Kapitel sind die verschiedenen Typen linearer ganzzahliger und gemischt-ganzzahliger Programme formuliert und erklärt worden. Fast die gesamte bisherige Vorlesung hat sich mit Spezialfällen dieser Probleme beschäftigt. In Kapitel 15 haben wir mit dem Branch & Bound-Verfahren einen allgemeinen Ansatz zur Lösung linearer ganzzahliger Programme dargestellt. Nun soll eine zweite generelle Methode zur Lösung derartiger Probleme vorgestellt werden: das sogenannte Schnittebenenverfahren. Hierzu benötigen wir jedoch ein wenig Theorie, die in diesem Kapitel zusammengestellt ist.

Wir werden uns hauptsächlich auf ganzzahlige Programme konzentrieren und gemischt-ganzzahlige nebenbei (in Anmerkungen) behandeln.

### 16.1 Einführung

In der Vorlesung haben wir uns bereits intensiv mit der Polyedertheorie beschäftigt. Zur Erinnerung: **Polyeder** sind Teilmengen des  $\mathbb{K}^n$  der folgenden Form

$$P(A, b) = \{x \in \mathbb{K}^n \mid Ax \leq b\},$$

d. h. Polyeder sind Durchschnitte von endlich vielen abgeschlossenen Halbräumen. Wir haben gezeigt, dass Polyeder auch eine andere Darstellung besitzen. Es gibt nämlich endliche Mengen  $V \subseteq \mathbb{K}^n$  und  $E \subseteq \mathbb{K}^n$ , so dass gilt

$$P(A, b) = \text{conv}(V) + \text{cone}(E).$$

Eine **Seitenfläche** eines Polyeders  $P$  ist eine Teilmenge  $F \subseteq P$  mit der Eigenschaft: Es existiert eine Ungleichung  $c^T x \leq c_0$  mit  $P \subseteq \{x \in \mathbb{R}^n \mid c^T x \leq c_0\}$  und  $F = \{x \in P \mid c^T x = c_0\}$ . Wir haben gezeigt, dass zwei Typen von Seitenflächen eine besondere Rolle spielen, nämlich die Ecken und die Facetten, d. h. die minimalen und die maximalen echten Seitenflächen (im Sinne der Mengeninklusion).

Zur Notationsvereinfachung führen wir nun einige neue Bezeichnungen ein.

**(16.1) Definition**

(a) Seien  $A \in \mathbb{K}^{(m,n)}$  und  $b \in \mathbb{K}^m$ , dann setzen wir

$$IP(A, b) := \{x \in \mathbb{Z}^n \mid Ax \leq b\}$$

und

$$IP^=(A, b) := \{x \in \mathbb{Z}^n \mid Ax = b, x \geq 0\}.$$

(b) Seien  $A \in \mathbb{K}^{(m,n_1)}$ ,  $B \in \mathbb{K}^{(m,n_2)}$  und  $b \in \mathbb{K}^m$ , dann setzen wir

$$MIP(A, B, b) := \left\{ \begin{pmatrix} x \\ y \end{pmatrix} \in \mathbb{K}^{n_1+n_2} \mid Ax + By \leq b, x \in \mathbb{Z}^{n_1} \right\}$$

und

$$MIP^=(A, B, b) := \left\{ \begin{pmatrix} x \\ y \end{pmatrix} \in \mathbb{K}^{n_1+n_2} \mid Ax + By = b, x, y \geq 0, x \in \mathbb{Z}^{n_1} \right\}.$$

□

$IP(A, b)$  ist also die Menge aller ganzzahligen Vektoren im Polyeder  $P(A, b)$ ,  $IP^=(A, b)$  ist die Menge aller ganzzahligen Vektoren im Polyeder  $P^=(A, b)$ . Analoges gilt für  $MIP(A, B, b)$  und  $MIP^=(A, B, b)$ . Im Weiteren gehen wir davon aus, dass alle auftretenden ganzzahligen linearen Programme entweder in der Form

$$(16.2) \quad \max_{x \in IP(A, b)} c^T x \quad \text{oder} \quad (16.3) \quad \max_{x \in IP^=(A, b)} c^T x$$

und alle gemischt-ganzzahligen Programme entweder in der Form

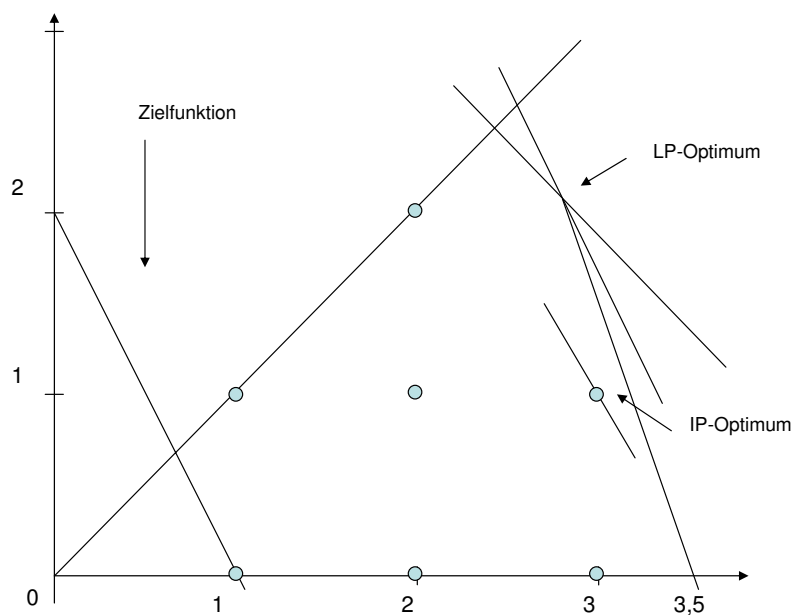
$$(16.4) \quad \max_{(x^T, y^T)^T \in MIP(A, B, b)} c^T x + d^T y \quad \text{oder} \quad (16.5) \quad \max_{(x^T, y^T)^T \in MIP^=(A, B, b)} c^T x + d^T y$$

gegeben sind. Statt des Wortungetüms „ganzzahliges lineares Programm“ schreiben wir in Zukunft häufig **IP** oder **GLP** und statt „gemischt-ganzzahliges lineares Programm“ **MIP** oder **GGLP**. Um die Unterschiede zwischen linearen, ganzzahligen und gemischt-ganzzahligen Programmen und zwischen den Mengen  $P(A, b)$ ,  $IP(A, b)$ ,  $MIP(A, B, b)$  zu verdeutlichen, betrachten wir zunächst ein Beispiel.

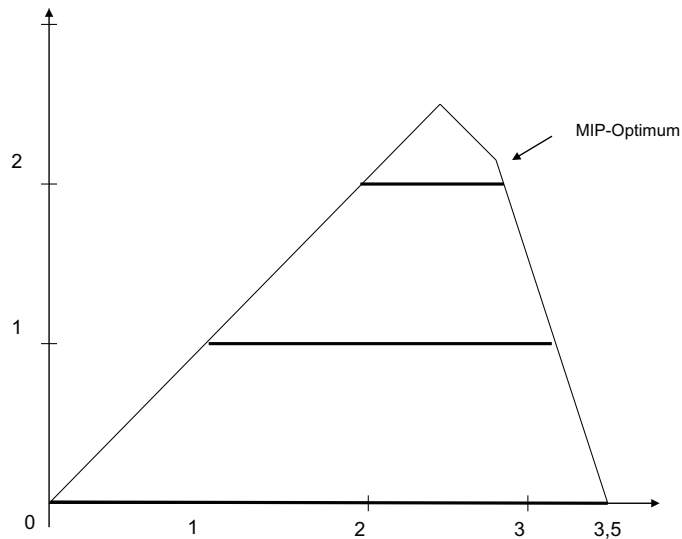
**(16.6) Beispiel.**

- (a) Gegeben sei das folgende IP, dessen Lösungsmenge  $IP(A, b)$  in Abbildung 16.1 dargestellt ist. Die Menge  $IP(A, b)$  enthält 8 Vektoren, die als dicke Punkte gekennzeichnet sind.

$$\begin{aligned}
 \max \quad & 2x_1 + x_2 \\
 & x_1 + x_2 \leq 5 \\
 & -x_1 + x_2 \leq 0 \\
 & 6x_1 + 2x_2 \leq 21 \\
 & x_1, x_2 \geq 0 \quad \text{und ganzzahlig.}
 \end{aligned}$$

**Abb. 16.1**

- (b) Wir streichen im obigen IP die Ganzzahligkeitsforderung für  $x_1$  und erhalten so ein MIP.

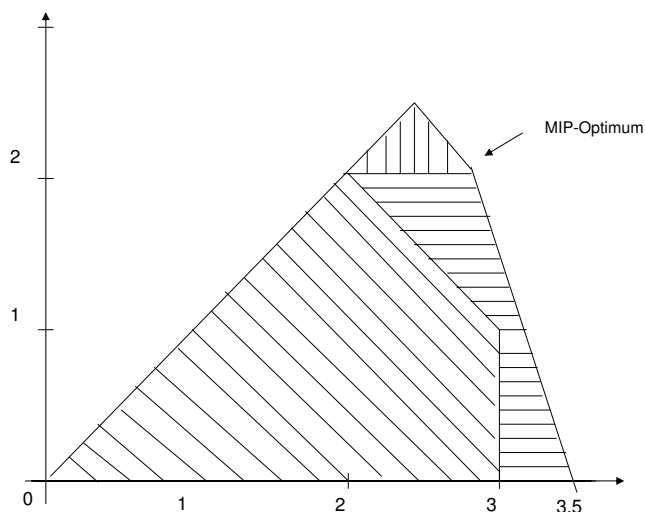


**Abb. 16.2**

Die Menge der zulässigen Lösungen dieses MIP besteht aus den drei dick gezeichneten Streifen in Abbildung 16.2.

Aus den beiden obigen Abbildungen ist klar, dass die Optimallösungen der LP-Relaxierung des ganzzahligen Programms (a), des IPs aus (a) und des MIPs aus (b) verschieden sind.  $\square$

Das Ziel dieses und des nächsten Kapitels ist die Anwendung polyedertheoretischer Methoden auf IPs und MIPs. Die Mengen  $IP(A, b)$ ,  $IP^=(A, b)$ ,  $MIP(A, B, b)$  und  $MIP^=(A, B, b)$  sind nach Definition Teilmengen von Polyedern. Wenn man die konvexe Hülle dieser Mengen bildet, so könnte man (intuitiv) meinen, dass dann auch wieder ein Polyeder entsteht. Abbildung 16.3 zeigt z. B. das durch die Ungleichungen aus (16.6) (a) (ohne Ganzzahligkeitsbedingung) definierte Polyeder  $P(A, b)$ , die konvexe Hülle der Lösungen von (b) und die konvexe Hülle der Lösungen von (a). Alle diese Mengen sind



Polyeder.

**Abb. 16.3**

Wäre dies immer so, so könnten wir u. U. alle bereits diskutierten Methoden zur Lösung linearer Programme (Simplexverfahren, Ellipsoidmethode, Karmarkar-Algorithmus) benutzen, um ganzzahlige Probleme anzugehen. Jedoch ist i. a. die konvexe Hülle von IP's oder MIP's kein Polyeder, wie das nachfolgende Beispiel zeigt.

**(16.7) Beispiel.** Wir betrachten das Programm

$$\begin{aligned} \max \quad & p \\ 0 \leq p & \leq \sqrt{2} \\ p & \in \mathbb{Q}. \end{aligned}$$

Bekanntlich ist  $\sqrt{2}$  keine rationale Zahl. Jede irrationale Zahl kann aber durch rationale Zahlen beliebig genau (z. B. von unten) angenähert werden. Daher hat (\*) kein Maximum. Wir verwandeln nun (\*) in ein ganzzahliges Programm, indem wir die rationale Zahl  $p$  als Bruch  $\frac{y}{x}$  mit nichtnegativen ganzen Zahlen  $x, y$  schreiben. Wir erhalten

$$\begin{aligned} \max \quad & -\sqrt{2}x + y \\ -\sqrt{2}x + y & \leq 0 \\ x & \geq 1 \\ y & \geq 0 \\ x, y & \text{ ganzzahlig.} \end{aligned}$$

Das Programm (\*\*) hat die folgenden Eigenschaften:

- (a) Es gibt zulässige Lösungen.
- (b) Der Zielfunktionswert ist nach oben durch 0 beschränkt, da 0 das Optimum des zugehörigen (reellen) LPs ist.
- (c) Es gibt keine ganzzahlige Optimallösung (sonst wäre  $\sqrt{2}$  rational).

Bezeichnen wir mit  $S$  die Lösungsmenge von (\*\*), so folgt aus (a), (b), (c):

$\text{conv}(S)$  ist kein Polyeder.

Denn für jeden Zielfunktionsvektor  $c$  gilt trivialerweise  $\max\{c^T z \mid z \in S\} = \max\{c^T z \mid z \in \text{conv}(S)\}$ . Da für  $c^T = (-\sqrt{2}, 1)$  das Problem  $\max\{c^T z \mid z \in S\}$  keine Optimallösung hat, hat auch  $\max\{c^T z \mid z \in \text{conv}(S)\}$  keine. Wäre  $\text{conv}(S)$  ein Polyeder, so müsste, da die Zielfunktion  $-\sqrt{2}x + y$  über  $S$  und somit über  $\text{conv}(S)$  beschränkt ist und da  $\text{conv}(S) \neq \emptyset$  gilt, das lineare Programm  $\max\{c^T z \mid z \in \text{conv}(S)\}$  eine Optimallösung haben.  $\square$

Die Tatsache, dass im obigen Beispiel  $\text{conv}(S)$  kein Polyeder ist, hat zwei Gründe:

Die Daten des Programms sind irrational.

$S$  ist unbeschränkt.

Wir werden zeigen, dass man unter der Annahme, dass ein Programm eine der beiden Eigenschaften nicht hat, das gewünschte Resultat beweisen kann.

Dabei ist die Aussage, dass die konvexe Hülle der ganzzahligen Punkte eines durch rationale Ungleichungen definierten Polyeders ein Polyeder ist, nicht trivial. Sie wird im nächsten Abschnitt behandelt.

**(16.8) Satz.** *Ist  $B \subseteq \mathbb{K}^n$  eine beschränkte Menge, dann ist  $\text{conv}\{x \in B \mid x \text{ ganzzahlig}\}$  ein Polytop. Ferner gibt es eine ganzzahlige  $(m, n)$ -Matrix  $D$  und  $d \in \mathbb{Z}^m$  mit  $\text{conv}(\{x \in B \mid \{x \in \mathbb{Z}^n\}) = P(D, d)$ .*

**Beweis :** Ist  $B$  beschränkt, so ist  $IB := \{x \in B \mid x \text{ ganzzahlig}\}$  natürlich endlich. Wir wissen, dass die konvexe Hülle einer endlichen Menge ein Polytop ist. Da  $IB \subseteq \mathbb{Q}^n$ , ist  $\text{conv}(IB)$  ein Polytop in  $\mathbb{Q}^n$ . Daraus folgt, dass  $\text{conv}(IB)$  eine Darstellung der Form  $P(D', d')$  hat mit  $D', d'$  rational. Durch Multiplikation der Zeilen von  $D'x \leq d'$  mit geeigneten ganzen Zahlen lässt sich ein System  $Dx \leq d$  mit den gewünschten Eigenschaften konstruieren.  $\square$

Aus (16.8) folgt unmittelbar, dass für jede Matrix  $A \in \mathbb{K}^{(m,n)}$  und jeden Vektor  $b \in \mathbb{K}^m$ , mit  $P(A, b)$  beschränkt, die Menge  $\text{conv}(IP(A, b))$  ein Polytop (mit einer ganzzahligen Darstellung) ist.

## 16.2 Ganzzahlige Punkte in rationalen Polyedern

Aufgrund der Vorbemerkungen in Abschnitt 16.1 wollen wir uns nun auf Probleme beschränken, bei denen alle Daten rational sind. Zur Abkürzung nennen wir ein Polyeder  $P \subseteq \mathbb{K}^n$  **rational**, wenn es eine Matrix  $A \in \mathbb{Q}^{(m,n)}$  und einen Vektor  $b \in \mathbb{Q}^m$  gibt mit  $P = P(A, b)$ . Wir wissen, dass ein Polyeder genau dann rational ist, wenn es endliche Mengen  $V \subseteq \mathbb{Q}^n$  und  $E \subseteq \mathbb{Q}^n$  gibt mit  $P = \text{conv}(V) + \text{cone}(E)$ . Der Einfachheit halber werden wir uns auch auf rationale Wertebereiche beschränken und nur Polyeder in  $\mathbb{Q}^n$  betrachten.

Ist  $P$  eine beliebige Teilmenge des reellen oder rationalen Vektorraums, so setzen wir

$$P_I := \text{conv}\{x \in P \mid x \text{ ganzzahlig}\}.$$

Ist also zum Beispiel  $P = P(A, b)$  ein Polyeder, so bezeichnet  $IP(A, b)$  die ganzzahligen Punkte in  $P$  und  $P_I$  die konvexe Hülle von  $IP(A, b)$ . Wir zeigen nun, dass  $P_I$  wiederum ein Polyeder ist. Wir wissen dies bereits aus (16.8), falls  $P$  beschränkt ist. Ferner gilt offenbar

**(16.9) Satz.** *Ist  $C$  ein rationaler Kegel, so gilt*

$$C = C_I.$$

**Beweis :** Ein rationaler Kegel  $C$  hat die Form  $C = \text{cone}(E)$ ,  $E \subseteq \mathbb{Q}^n$  endlich. Durch Multiplikation der Vektoren aus  $E$  mit geeigneten ganzen Zahlen erhält man eine Menge  $E' \subseteq \mathbb{Z}^n$  mit  $C = \text{cone}(E')$ . Jeder ganzzahlige Vektor in  $C$  ist also konische Kombination der Vektoren aus  $E'$ . Daraus folgt die Behauptung.  $\square$

Damit können wir nun zeigen:

**(16.10) Satz.** *Ist  $P$  ein rationales Polyeder, dann ist  $P_I$  ebenfalls ein (rationales) Polyeder. Ferner gilt  $\text{rec}(P) = \text{rec}(P_I)$ .*

**Beweis :**  $P$  hat eine Darstellung der Form  $P = Q + C$ , wobei  $Q$  ein Polytop und  $C = \text{rec}(P)$  ein Kegel ist. Nach Satz (16.9) können wir annehmen, dass  $C = \text{cone}\{y_1, \dots, y_s\}$  gilt mit  $y_i \in \mathbb{Z}^n$ ,  $i = 1, \dots, s$ . Wir setzen

$$B := \left\{ \sum_{i=1}^s \mu_i y_i \mid 0 \leq \mu_i \leq 1, i = 1, \dots, s \right\}$$

und behaupten:

$$P_I = (Q + B)_I + C.$$

Aus dieser Behauptung folgt der Satz, denn  $Q + B$  ist beschränkt, also ist  $(Q + B)_I$  ein Polytop; und falls  $P_I \neq \emptyset$ , dann ist  $C$  der Rezeptionskegel von  $P_I$ .

Wir beweisen zunächst  $P_I \subseteq (Q + B)_I + C$ . Es sei  $p$  ein ganzzahliger Punkt in  $P$ . Dann gibt es Vektoren  $q \in Q$  und  $c \in C$  mit  $p = q + c$ . Der Vektor  $c$  hat eine Darstellung der Form  $c = \sum_{i=1}^s \mu_i y_i$ ,  $\mu_i \geq 0$ . Setzen wir  $c' := \sum_{i=1}^s \lfloor \mu_i \rfloor y_i$  und  $b := c - c'$ , dann ist  $c'$  ganzzahlig; und weil  $p$  und  $c'$  ganzzahlig sind, ist auch  $q + b = p - c'$  ganzzahlig. Ferner gilt nach Definition  $b \in B$ . Daraus folgt  $q + b \in (Q + B)_I$  und somit  $p = (q + b) + c' \in (Q + B)_I + C$ .

Aus  $(Q + B)_I + C \subseteq P_I + C = P_I + C_I = (P + C)_I = P_I$  folgt die umgekehrte Inklusion.  $\square$

Analog kann man zeigen

**(16.11) Satz.** Sind  $A \in \mathbb{Q}^{(m, n_1)}$ ,  $B \in \mathbb{Q}^{m \times n_2}$ ,  $b \in \mathbb{Q}^m$ , dann ist  $\text{conv}\{x \in MIP(A, B, b)\}$  ein rationales Polyeder.  $\square$

Wir haben schon mehrfach angemerkt, dass für  $c \in \mathbb{Q}^n$ ,  $S \subseteq \mathbb{Q}^n$  folgendes gilt:

$$\max\{c^T x \mid x \in S\} = \max\{c^T x \mid x \in \text{conv}(S)\}.$$

Die obigen Sätze zeigen also, dass wir IP's oder MIP's in lineare Programme transformieren können. Haben wir damit eine Lösungsmethode für ganzzahlige und gemischt-ganzzahlige Programme gefunden? Nicht ganz, denn sind IP's oder MIP's durch (16.2), (16.3), (16.4) oder (16.5) gegeben, so ist überhaupt nicht klar, wie man die Ungleichungen finden kann, die aus den ganzzahligen oder gemischt-ganzzahligen Programmen lineare Programme machen. Wir wissen nun zwar, dass es theoretisch geht, aber ob es praktisch durchführbar ist, ist a priori nicht zu sehen.

Bevor wir dieser Frage nachgehen, wollen wir untersuchen, ob man aus Informationen über die Kodierungslängen der Ungleichungen eines Systems  $Ax \leq b$  Abschätzungen über die Kodierungslängen der Facetten von  $P(A, b)_I$  ableiten kann. Zunächst beweisen wir einen Hilfssatz.

**(16.12) Lemma.** Sei  $P \subseteq \mathbb{Q}^n$  ein Polyeder.

- (a) Hat  $P$  eine Darstellung der Form  $P = P(A, b)$ , so daß die Kodierungslänge jeder Ungleichung des Systems  $Ax \leq b$  höchstens  $\varphi$  ist, dann gibt es endliche Mengen  $V \subseteq \mathbb{Q}^n$  und  $E \subseteq \mathbb{Q}^n$  mit  $P = \text{conv}(V) + \text{cone}(E)$ , so daß die Vektoren aus  $V \cup E$  höchstens die Kodierungslänge  $4n^2\varphi$  haben.
- (b) Hat  $P$  eine Darstellung der Form  $P = \text{conv}(V) + \text{cone}(E)$ , so daß die Kodierungslänge jedes Vektors aus  $V \cup E$  höchstens  $\nu$  ist, dann gibt es ein Ungleichungssystem  $Ax \leq b$  mit  $P = P(A, b)$ , so daß jede Ungleichung des Systems höchstens die Kodierungslänge  $3n^2\nu$  hat.

**Beweis :** (a) Angenommen die Voraussetzungen von (a) gelten, dann gibt es Ungleichungen  $a_i^T x \leq b_i$ ,  $a_i \in \mathbb{Q}^n$ ,  $b_i \in \mathbb{Q}$ ,  $i = 1, \dots, m$ , die  $P$  definieren, mit  $\langle a_i \rangle + \langle b_i \rangle \leq \varphi$  für  $i = 1, \dots, m$ . Aus "Optimierungsmethoden I" wissen wir, dass es endliche Mengen  $V, E \subseteq \mathbb{Q}^n$  gibt mit  $P = \text{conv}(V) + \text{cone}(E)$ , so dass die Komponenten der Vektoren aus  $V \cup E$  entweder 0 oder Quotienten von zwei Subdeterminanten der folgenden Matrix sind

$$C = \begin{pmatrix} a_1^T & b_1 \\ \vdots & \vdots \\ a_m^T & b_m \end{pmatrix}.$$

Sei  $D$  eine quadratische  $(k, k)$ -Untermatrix von  $C$ , dann folgt aus einfachen Abschätzungen (Hausaufgabe!)

$$\langle \det D \rangle \leq 2\langle D \rangle - k^2.$$

Da nach Voraussetzung  $\langle D \rangle \leq k\varphi$  gilt, erhalten wir

$$\langle \det D \rangle \leq 2\langle D \rangle - k^2 \leq 2n\varphi.$$

Die Zähler und Nenner der Komponenten der Vektoren aus  $V \cup E$  haben also höchstens die Kodierungslänge  $2n\varphi$ , somit hat jede Komponente höchstens die Kodierungslänge  $4n\varphi$  und damit jeder dieser Vektoren höchstens die Kodierungslänge  $4n^2\varphi$ .

(b) beweist man ähnlich. (Hausaufgabe!) □

**(16.13) Satz.** Sei  $P = P(A, b) \subseteq \mathbb{Q}^n$  ein rationales Polyeder, so dass die Kodierungslänge jeder Ungleichung des Systems  $Ax \leq b$  höchstens  $\varphi$  ist. Dann gibt es ein Ungleichungssystem  $Dx \leq d$  mit  $P_I = P(D, d)$ , so dass die Kodierungslänge jeder Ungleichung des Systems höchstens  $15n^6\varphi$  ist.

**Beweis :** Ist  $P_I = \emptyset$ , so ist die Behauptung trivial. Wir können also annehmen, dass  $P_I \neq \emptyset$  gilt.

Nach Lemma (16.12) gibt es eine Darstellung von  $P$  der Form

$$P = \text{conv}\{v_1, \dots, v_t\} + \text{cone}\{y_1, \dots, y_s\},$$

so dass jeder der Vektoren  $v_1, \dots, v_t, y_1, \dots, y_s$  höchstens die Kodierungslänge  $4n^2\varphi$  hat. Wir setzen

$$Q := P \cap \{x \in \mathbb{Q}^n \mid -(n+1)2^{4n^3\varphi} \leq x_i \leq (n+1)2^{4n^3\varphi}, i = 1, \dots, n\}$$

und zeigen

$$(*) \quad P_I = Q_I + \text{conv}\{y_1, \dots, y_s\}.$$

Die Inklusion  $\supseteq$  in (\*) gilt trivialerweise. Um die umgekehrte Inklusion zu zeigen, wählen wir einen beliebigen (ganzzahligen) Vektor  $x \in P$ . Für  $j = 1, \dots, s$  sei  $y_j'$  derjenige

ganzahlige Vektor, der aus  $y_j$  dadurch entsteht, dass man  $y_j$  mit dem k. g. V. der Nenner der Komponenten von  $y_j$  multipliziert. Die Kodierungslänge von  $y'_j$  ist somit höchstens  $4n^3\varphi$ . Aufgrund des Satzes von Caratheodory können wir  $x$  in folgender Form darstellen

$$x = \lambda_1 v_1 + \dots + \lambda_t v_t + \mu_1 y'_1 + \dots + \mu_s y'_s,$$

wobei  $\lambda_1, \dots, \lambda_t, \mu_1, \dots, \mu_s \geq 0$ ,  $\lambda_1 + \dots + \lambda_t = 1$  und höchstens  $n$  der rationalen Zahlen  $\mu_j$  von Null verschieden sind. Wir setzen nun

$$\begin{aligned} x' &:= \lambda_1 v_1 + \dots + \lambda_t v_t + (\mu_1 - \lfloor \mu_1 \rfloor) y'_1 + \dots + (\mu_s - \lfloor \mu_s \rfloor) y'_s \\ x'' &:= \lfloor \mu_1 \rfloor y'_1 + \dots + \lfloor \mu_s \rfloor y'_s. \end{aligned}$$

Da  $x$  und  $x''$  ganzzahlig sind, ist auch  $x' = x - x''$  ganzzahlig. Für jeden der Vektoren  $v_i$  und  $y'_j$  hat jede seiner Komponenten einen Betrag, der nicht größer als  $2^{4n^3}\varphi$  ist. Da höchstens  $n$  der Zahlen  $\mu_1, \dots, \mu_s$  von Null verschieden sind, gilt  $\lambda_1 + \dots + \lambda_t + (\mu_1 - \lfloor \mu_1 \rfloor) + \dots + (\mu_s - \lfloor \mu_s \rfloor) \leq n + 1$ ; also hat jede Komponente von  $x'$  einen Absolutbetrag, der nicht größer als  $(n+1)2^{4n^3}\varphi$  ist. Mithin gilt  $x' \in Q$ . Da  $x'' \in \text{cone}\{y'_1, \dots, y'_s\}$ , gehört  $x = x' + x''$  zu  $Q_I + \text{cone}\{y_1, \dots, y_s\}$ . Daraus folgt die Gleichheit in (\*).

Sei  $Z$  die Menge der ganzzahligen Vektoren in  $Q$ . Nach Definition hat jeder Vektor in  $Z$  höchstens die Kodierungslänge  $5n^4\varphi$ . Aufgrund von (\*) gilt

$$P_I = \text{conv}(Z) + \text{cone}\{y_1, \dots, y_s\}.$$

In dieser Darstellung von  $P_I$  als Summe eines Polytops und eines Kegels hat jeder Vektor höchstens die Kodierungslänge  $5n^4\varphi$ ; folglich können wir aus Lemma (16.12) (b) schließen, dass  $P_I$  durch ein Ungleichungssystem  $Dx \leq d$  dargestellt werden kann, so dass jede Ungleichung dieses Systems höchstens die Kodierungslänge  $15n^6\varphi$  hat.  $\square$

Im obigen Beweis haben wir eine Aussage mitbewiesen, die für weitere Anwendungen interessant ist.

**(16.14) Folgerung.** Sei  $P = P(A, b)$  ein rationales Polyeder, so dass die Kodierungslänge jeder Ungleichung des Systems  $Ax \leq b$  höchstens  $\varphi$  ist. Falls  $P$  einen ganzzahligen Vektor enthält, dann gibt es in  $P$  einen ganzzahligen Vektor, der im Würfel

$$\{x \in \mathbb{Q}^n \mid -(n+1)2^{4n^3}\varphi \leq x_i \leq (n+1)2^{4n^3}\varphi\}$$

enthalten ist, und es gibt in  $P$  einen ganzzahligen Vektor der Kodierungslänge höchstens  $5n^4\varphi$ .  $\square$

Aus dieser Folgerung ergibt sich eine wichtige Komplexitätstheoretische Konsequenz.

**(16.15) Folgerung.** Das Entscheidungsproblem "Hat ein gegebenes rationales Ungleichungssystem  $Ax \leq b$  eine ganzzahlige Lösung?" ist in der Klasse  $\mathcal{NP}$ .

**Beweis :** Um (nichtdeterministisch) zu zeigen, dass eine ganzzahlige Lösung von  $Ax \leq b$  existiert, können wir einfach eine Lösung  $y$  raten. Die Korrektheit der geratenen Lösung  $y$  kann aber nur dann (durch Substitution von  $y$  in das Ungleichungssystem) in polynomialer Zeit überprüft werden, wenn die Kodierungslänge von  $y$  polynomial in  $\langle A \rangle + \langle b \rangle$  ist. Aus (16.14) folgt, dass — wenn  $Ax \leq b$  überhaupt eine ganzzahlige Lösung hat — eine ganzzahlige Lösung  $y$  mit  $\langle y \rangle \leq 5n^4 \langle A, b \rangle$  existiert. Also kann man tatsächlich eine “kleine” ganzzahlige Lösung von  $Ax \leq b$  raten, wenn das Entscheidungsproblem eine ja-Antwort besitzt.  $\square$

Ein weiteres interessantes Korollar von (16.13) ist die folgende Beobachtung.

**(16.16) Folgerung.** Sei  $P \subseteq \mathbb{Q}^n$  ein Polyeder, für das es eine Beschreibung  $P = P(A, b)$  gibt, so dass jede der Ungleichungen des Systems höchstens die Kodierungslänge  $\varphi$  hat, und sei  $c \in \mathbb{Q}^n$ . Falls  $\max\{c^T x \mid x \in P, x \text{ ganzzahlig}\}$  endlich ist, dann gibt es eine Optimallösung, deren Kodierungslänge höchstens  $5n^4 \varphi$  ist, und der Absolutbetrag des Optimalwertes dieses IP's ist höchstens  $2^{\langle c \rangle + 5n^4 \varphi - 2n}$ .

**Beweis :** Aus dem Beweis von Satz (16.13) folgt, dass das Maximum von  $\max\{c^T x \mid x \in P, x \text{ ganzzahlig}\}$  in einem der Vektoren  $z \in Z$  (das sind die ganzzahligen Vektoren aus  $Q$ ) angenommen wird. Aus  $\langle z \rangle \leq 5n^4 \varphi$  folgt die erste Behauptung. Die zweite folgt mit Cauchy-Schwarz und Lemma (12.10) (c) aus

$$|c^T x| \leq \|c\|_2 \|z\|_2 \leq (2^{\langle c \rangle - n} - 1)(2^{\langle z \rangle - n} - 1) \leq 2^{\langle c \rangle + 5n^4 \varphi - 2n}.$$

$\square$

Folgerung (16.16) bzw. (16.14) impliziert, dass für ein ganzzahliges Programm  $\max c^T x$ ,  $x \in IP(A, b)$  mit  $c \in \mathbb{Z}^n$  die Menge der möglichen Optimallösungen im Würfel

$$\{x \in \mathbb{Z}^n \mid -(n+1)2^{4n^3 \langle A, b \rangle} \leq x_i \leq (n+1)2^{4n^3 \langle A, b \rangle}\}$$

liegt und dass der Optimalwert dieses IP's, wenn er existiert, eine ganze Zahl  $z^*$  im folgenden Intervall ist

$$-2^{\langle c \rangle + 5n^4 \langle A, b \rangle - 2n} \leq z^* \leq 2^{\langle c \rangle + 5n^4 \langle A, b \rangle - 2n}.$$

Daraus folgt direkt, dass IP's durch binäre Suche gelöst werden können, wir haben dazu “lediglich” die Entscheidungsfrage

$$\text{“Gibt es einen Vektor } x \in IP(A, b) \text{ mit } c^T x \geq z^* \text{?”}$$

für die  $z^*$  aus dem obigen Intervall zu lösen. Leider ist jedoch das letztere Entscheidungsproblem  $\mathcal{NP}$ -vollständig.

Analoge Resultate wie die oben für ganzzahlige Programme bewiesenen, lassen sich auch für gemischt-ganzzahlige Programme ableiten. Aus ihnen folgt die Endlichkeit des Dakin-Verfahrens (15.2).

### 16.3 Schnittebentheorie

Wir haben im vorhergehenden Abschnitt gezeigt, dass für jedes rationale Polyeder  $P$  die Menge  $P_I$ , die konvexe Hülle der ganzzahligen Punkte in  $P$ , ein Polyeder ist ( $P_I$  nennt man häufig **das zu  $P$  gehörige ganzzahlige Polyeder**). Ferner haben wir bewiesen, dass es für  $P = P(A, b)$  eine Darstellung  $P_I = P(D, d)$  gibt, so dass die Ungleichungen des Systems  $Dx \leq d$  „kleine“ Koeffizienten haben. Wir wollen nun der Frage nachgehen, ob man  $Dx \leq d$  algorithmisch bestimmen kann und ob man die Anzahl der Ungleichungen von  $Dx \leq d$  durch ein Polynom in der Anzahl der Ungleichungen von  $Ax \leq b$  beschränken kann.

Die erste Frage hat eine positive Antwort, die zweite nicht. Selbst dann, wenn ein kombinatorisches Optimierungsproblem in polynomialer Zeit gelöst werden kann, kann die Anzahl der Ungleichungen, die notwendig sind, um die konvexe Hülle der zulässigen Punkte zu beschreiben, exponentiell in den Daten des Problems sein. Wir geben hierzu ein Beispiel.

Ein **Matching**  $M$  ist eine Teilmenge der Kantenmenge eines Graphen  $G = (V, E)$ , so dass jeder Knoten  $v \in V$  auf höchstens einer Kante von  $M$  liegt.

Das Matching-Problem ist die Aufgabe, in einem Graphen  $G = (V, E)$  mit Kantengewichten  $c_e \in \mathbb{Q}$  für alle  $e \in E$  ein Matching maximalen Gewichts zu finden. Aus der Definition können wir eine Formulierung des Matching-Problems als ganzzahliges Programm ableiten. Diese lautet

$$(16.17) \quad \begin{aligned} \max \quad & c^T x \\ x(\delta(v)) := \sum_{e \in \delta(v)} x_e \leq 1 & \quad \forall v \in V \\ x_e \geq 0 & \quad \forall e \in E \\ x_e \in \{0, 1\} & \quad \forall e \in E. \end{aligned}$$

Hierbei ist  $\delta(v)$  die Menge aller Kanten, die den Knoten  $v$  enthalten. Bezeichnen wir mit  $P(G)$  die Menge der zulässigen Lösungen von (16.17) ohne Ganzzahligkeitsbedingung und mit  $\text{MATCH}(G)$  die konvexe Hülle aller Inzidenzvektoren von Matchings in  $G$ , so gilt offenbar

$$P(G)_I = \text{MATCH}(G).$$

Man kann zeigen, dass

$$P(G) = \text{MATCH}(G)$$

genau dann gilt, wenn  $G$  bipartit ist und keine isolierten Knoten hat (Folgerung aus der *totalen Unimodularität* der Matrix-Restriktionen in 16.17, falls  $G$  bipartit ist). Man kann ferner zeigen (Edmonds (1965)), dass  $\text{MATCH}(G)$  durch das folgende System von Ungleichungen beschrieben werden kann:

$$(16.18) \quad \begin{aligned} x(\delta(v)) &\leq 1 & \forall v \in V, \\ x(E(W)) &\leq \frac{1}{2}(|W| - 1) & \forall W \subseteq V, |W| \geq 3 \text{ und ungerade,} \\ x_e &\geq 0 & \forall e \in E. \end{aligned}$$

Hierbei bezeichnet  $E(W)$  die Menge aller Kanten, bei denen beide Endknoten in  $W$  liegen.

Falls  $G$  der vollständige Graph  $K_n$ ,  $n \geq 3$  ist, ist das System (16.18) nicht redundant, d. h. jede der in (16.18) angegebenen Ungleichungen definiert eine Facette von  $\text{MATCH}(K_n)$ . Somit hat  $\text{MATCH}(K_n)$  rund  $2^{n-1}$  Facetten; daraus folgt, dass die Speicherkomplexität jeder vollständigen Beschreibung von  $\text{MATCH}(K_n)$  exponentiell in der Kodierungslänge von  $K_n$  bzw. (16.17) ist. Folglich gibt es (für das Matching-Polyeder und somit generell) keinen Algorithmus, der aus einem Ungleichungssystem  $Ax \leq b$  eine vollständige Beschreibung von  $P(A, b)_I$  der Form  $P(A, b)_I = \{x \mid Dx \leq d\}$  konstruiert und dessen Laufzeit polynomial in  $\langle A \rangle + \langle b \rangle$  ist.

Trotz dieses negativen Resultates ist es — wie wir noch sehen werden — sinnvoll, sich Gedanken zu machen, wie man aus  $Ax \leq b$  ein System  $Dx \leq d$  mit  $P(A, b)_I = P(D, d)$  algorithmisch konstruieren kann. Dies wird mit sogenannten Schnittebenenmethoden gemacht. Hier wollen wir den theoretischen Hintergrund dieser Verfahren erläutern.

**(16.19) Definition.** Es seien  $A \in \mathbb{Q}^{(m,n)}$  und  $b \in \mathbb{Q}^m$ . Mit  $S$  bezeichnen wir das System der linearen Ungleichungen  $Ax \leq b$ , und wir setzen  $P := P(A, b)$ .

- (a) Wir sagen, dass eine Ungleichung  $d^T x \leq d_0$ ,  $d \in \mathbb{Z}^n$ , zum **elementaren Abschluss** von  $S$  gehört, wenn es einen Vektor  $\lambda \in \mathbb{Q}^m$ ,  $\lambda \geq 0$  gibt mit

$$\begin{aligned} \lambda^T A &= d^T \\ \text{und } \lfloor \lambda^T b \rfloor &\leq d_0. \end{aligned}$$

- (b) Wir setzen  $e^0(S) := S$  und bezeichnen die Menge aller Ungleichungen, die zum elementaren Abschluss von  $S$  gehören mit  $e^1(S)$ .

- (c) Für  $k > 1$  definieren wir  $e^k(S) := e^1(S \cup e^{k-1}(S))$ .

- (d)  $\text{cl}(S) := \bigcup_{k=1}^{\infty} e^k(S)$  heißt der **Abschluss** von  $S$ .

- (e) Wir setzen  $P^0 = P$  und  $P^k := \{x \in \mathbb{Q}^n \mid x \text{ erfüllt alle Ungleichungen des Systems } e^k(S)\}$ .

□

Der (einfache) Hintergrund der obigen Begriffsbildung ist der folgende. Für jeden Vektor  $\lambda \in \mathbb{Q}^m$ ,  $\lambda \geq 0$  wird die Ungleichung

$$(\lambda^T A)x \leq \lambda^T b$$

von allen Punkten aus  $P(A, b)$  erfüllt. Ist  $\lambda^T A$  ein ganzzahliger Vektor, so ist für jeden Vektor  $x \in IP(A, b)$  der Wert  $(\lambda^T A)x$  eine ganze Zahl. Folglich erfüllt in diesem Falle jeder Vektor  $x \in P(A, b)_I$  die Ungleichung

$$(\lambda^T A)x \leq \lfloor \lambda^T b \rfloor.$$

Diese Ungleichung (also eine Ungleichung des elementaren Abschlusses) ist somit gültig für das Polyeder  $P(A, b)_I$ . Durch Induktion folgt, dass jede Ungleichung aus dem Abschluss  $\text{cl}(S)$  von  $S$  gültig bezüglich  $P(A, b)_I$  ist. Es ist relativ erstaunlich, dass  $\text{cl}(S)$  alle Ungleichungen enthält, die Facetten von  $P(A, b)_I$  definieren. Der folgende Satz fasst die Eigenschaften der Abschlussoperationen zusammen.

**(16.20) Satz.**  $P = P(A, b) \subseteq \mathbb{Q}^n$  sei ein rationales Polyeder, dann gilt:

- (a) Ist  $c^T x \leq c_0$  mit  $c \in \mathbb{Z}^n$ ,  $c_0 \in \mathbb{Z}$  gültig bezüglich  $P_I$ , dann ist  $c^T x \leq c_0$  ein Element von  $\text{cl}(Ax \leq b)$ .
- (b) Es gibt eine ganze Zahl  $k \geq 0$ , so dass  $\text{cl}(Ax \leq b) = e^k(Ax \leq b)$  gilt und dass endlich viele der Ungleichungen aus  $e^k(Ax \leq b)$  das Polyeder  $P_I$  definieren, d. h. es gibt ein  $k \geq 0$  mit

$$P_I = P^k.$$

- (c) Sei  $c \in \mathbb{Z}^n$ , dann gibt es eine ganze Zahl  $k \geq 0$ , so dass für das ganzzahlige Programm

$$(*) \quad \max\{c^T x \mid x \in IP(A, b)\}$$

und das lineare Programm

$$(**) \quad \max\{c^T x \mid x \in P^k\}$$

*gilt: (\*) hat keine bzw. keine endliche Lösung genau dann, wenn (\*\*) keine bzw. keine endliche hat. Sind (\*) und (\*\*) endlich lösbar, so stimmen die Maxima überein.*

**Beweis :** Siehe V. Chvátal, "Edmonds' polytopes and a hierarchy of combinatorial problems", Discrete Mathematics 4 (1973) 305–337 für beschränkte Polyeder und A. Schrijver, "On cutting planes", Annals of Discrete Mathematics 9 (1980) 291–296 für den allgemeinen Fall.  $\square$

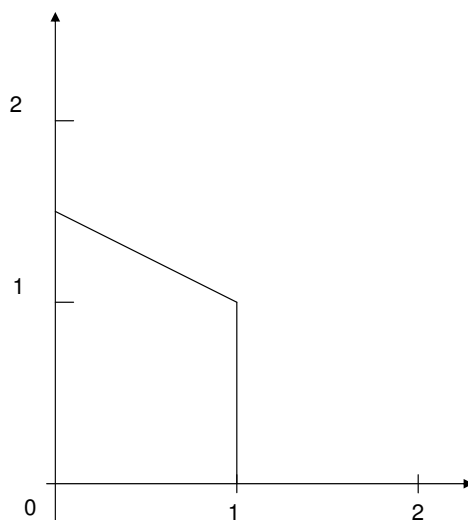
Man kann anhand von Beispielen zeigen, dass das  $k$  in Satz (16.20) (b) bzw. (c) beliebig groß sein kann, siehe Chvátal (1973). Mit einigem Recht kann man sagen, dass ein Problem mit "großem  $k$ " schwieriger ist als eines mit „kleinem“. Die Zahl  $k$  misst in einem gewissen Sinne die „Nähe zur Ganzzahligkeit“ eines linearen Programmierungsproblems.

Trotzdem bleibt festzuhalten, dass man nach Satz (16.20) aus den Daten eines Polytopen  $P(A, b)$  prinzipiell das definierende Ungleichungssystem des zugehörigen ganzzahligen Polytopen  $P_I = \text{conv}(IP(A, b))$  konstruieren kann. Wir werden im nächsten Kapitel sehen, dass dieses Verfahren sogar algorithmisch ausgewertet werden kann und zu einem endlichen Algorithmus führt. Zum Abschluss dieses Abschnitts betrachten wir noch ein Beispiel.

**(16.21) Beispiel.** (a) Gegeben sei das folgende Ungleichungssystem  $Ax \leq b$

$$\begin{aligned} (1) \quad & x_1 + 2x_2 \leq 3 \\ (2) \quad & x_1 \leq 1 \\ (3) \quad & -x_1 \leq 0 \\ (4) \quad & -x_2 \leq 0, \end{aligned}$$

dessen Lösungsmenge  $P(A, b)$  in Abbildung 16.4 dargestellt ist.



**Abb. 16.4**

Es gilt  $IP(A, b) = \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}$  und damit

$$P_I = \text{conv}(IP(A, b)) = \{x \in \mathbb{R}^2 \mid x_1 \leq 1, x_2 \leq 1, x_1 \geq 0, x_2 \geq 0\}.$$

Zum Beispiel gilt  $x_1 + x_2 \leq 2 \in e^1(S)$ , denn

$$(1) + (2) \implies 2x_1 + 2x_2 \leq 4 \implies x_1 + x_2 \leq 2.$$

Zur Beschreibung von  $\text{conv}(IP(A, b))$  fehlt noch die Ungleichung  $x_2 \leq 1$ , aber auch diese ist in  $e^1(S)$  enthalten, denn

$$(1) + (3) \implies 2x_2 \leq 3 \implies x_2 \leq \frac{3}{2} \implies x_2 \leq 1 = \lfloor \frac{3}{2} \rfloor.$$

Damit haben wir gezeigt  $\text{cl}(S) = e^1(S)$ , d. h.

$$P_I = \text{conv}(IP(A, b)) = \bigcap \{x \in \mathbb{R}^2 \mid a^T x \leq a_0 \in e^1(S)\}.$$

(b) Hausaufgabe!

Gegeben sei das folgende Ungleichungssystem  $S$

$$\begin{aligned} (1) \quad & -kx + y \leq 0 \\ (2) \quad & kx + y \leq k \\ (3) \quad & -y \leq 0, \end{aligned}$$

dessen Lösungsmenge in Abbildung 16.5 zu sehen ist.

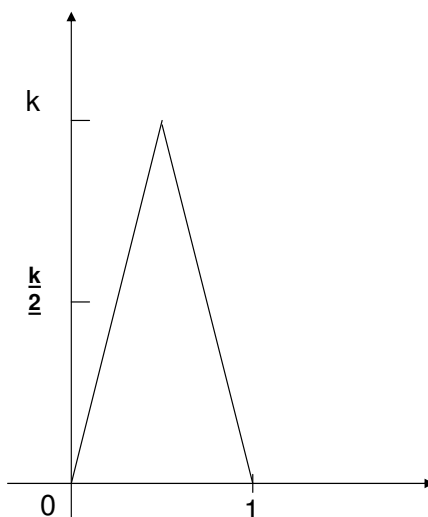


Abb. 16.5

Die ganzzahligen Lösungen von  $S$  sind die Vektoren  $\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ . Zeigen Sie

i)

$$\left. \begin{aligned} (4) \quad & y \leq \lfloor \frac{k}{2} \rfloor \\ (5) \quad & -y \leq 0 \\ (6) \quad & x \leq 1 \end{aligned} \right\} \text{ gehören zu } e^1(S)$$

ii) Die Ungleichung (7)  $y \leq 0$  gehört nicht zu  $e^1(S)$ , falls  $k \geq 2$  (offenbar gehört (7) zu  $\text{cl}(S)$ ).  $\square$

**(16.22) Hausaufgabe.** Zeigen Sie, dass für jeden Graphen  $G$  gilt

$$\text{MATCH}(G) = P(G)^1,$$

d. h. das Ungleichungssystem (16.18) gehört zum elementaren Abschluss des Ungleichungssystems (16.17).

Wir stellen fest, dass  $x(E(w)) \leq \frac{1}{2}(|w| - 1)$   $w \subseteq V$ ,  $|w|$  gerade ist im (ersten) elementaren Schluss.

$$\begin{aligned} \sum_{v \in W} x(\delta(v)) &= \sum_{e \in E(w)} x_e + \sum_{e \in \delta(w)} x_e \\ &\leq |w| \\ \Rightarrow \sum_{e \in E(w)} x_e &\leq \frac{1}{2} \left( \sum_{e \in E(w)} 2x_e + \sum_{e \in \delta(w)} x_e \right) \leq \frac{|w|}{2} \\ \Rightarrow x(E(w)) &\leq \frac{1}{2}(|w| - 1), \end{aligned}$$

wenn  $|w|$  gerade ist.

□



# Kapitel 17

## Allgemeine Schnittebenenverfahren der ganzzahligen Programmierung

In Kapitel 16 haben wir gesehen, dass die konvexe Hülle der Lösungen eines ganzzahligen oder gemischt-ganzzahligen Programms nicht unbedingt ein Polyeder sein muss (Beispiel (16.7)). Unter der Voraussetzung der Rationalität der Daten oder Beschränktheit der Lösungsmenge kann jedoch gezeigt werden (Sätze (16.8), (16.10)), dass diese konvexe Hülle immer ein Polyeder ist. Da wir zur Lösung ganzzahliger und gemischt-ganzzahliger Programme das Simplexverfahren bzw. andere Verfahren der linearen Programmierung einsetzen wollen, werden wir voraussetzen müssen, dass alle Daten rational sind. Für Probleme aus der Praxis ist dies immer der Fall, da auf einem Rechner sowieso nur rationale Zahlen repräsentiert werden können. Haben wir ein Ungleichungssystem  $Ax \leq b$  mit rationalen Daten, so kann jedes Matrixelement als Bruch mit positivem Nenner geschrieben werden, das gleiche gilt für den Vektor  $b$ . Daher können wir für jede Ungleichung  $A_i \cdot x \leq b_i$  das kleinste gemeinsame Vielfache  $k_i$  der Nenner bestimmen. Die Ungleichung

$$k_i A_i \cdot x \leq k_i b_i$$

definiert denselben Halbraum und hat nur ganzzahlige Elemente. Auf diese Weise erhalten wir ein ganzzahliges System  $Dx \leq d$ , so dass  $P(A, b) = P(D, d)$  gilt. Daraus folgt, dass es keine Einschränkung der Allgemeinheit ist, wenn wir statt Rationalität der Daten Ganzzahligkeit der Daten fordern.

**Im weiteren Verlauf dieses Kapitels seien alle Daten von  $Ax \leq b$ , bzw.  $Ax = b$ , etc. ganzzahlig.**

Wir wissen also, dass bei ganzzahligen Daten die konvexe Hülle der Lösungen eines IP oder MIP ein Polyeder ist und dass der Optimalwert eines ganzzahligen oder gemischt-ganzzahligen Programms mit dem Optimalwert des linearen Programms über der kon-

vexen Hülle übereinstimmt. Ferner haben wir in Satz (16.20) gesehen, dass ein lineares System zur Beschreibung der konvexen Hülle in einem gewissen Sinne “effektiv“ konstruierbar ist. Wir werden uns nun damit beschäftigen, wie diese polyedertheoretischen Erkenntnisse mit Hilfe von Schnittebenenverfahren algorithmisch ausgewertet werden können.

## 17.1 Ein Schnittebenenverfahren für ganzzahlige Programme

Schnittebenenverfahren arbeiten im Prinzip wie folgt:

### (17.1) Allgemeines Schnittebenenverfahren für ganzzahlige Programme

**Input:**  $A \in \mathbb{Z}^{(m,n)}$ ,  $b \in \mathbb{Z}^m$ ,  $c \in \mathbb{Z}^n$ .

**Output:** Optimale Lösung des ganzzahligen Programms

$$\begin{aligned} \max \quad & c^T x \\ & Ax = b \\ & x \geq 0 \\ & x \text{ ganzzahlig} \end{aligned}$$

1. Löse das zu (16.3) gehörige lineare Programm (LP), das durch Weglassen der Ganzzahligkeitsbedingung entsteht. Hat (LP) keine endliche Optimallösung, so ist (16.3) unbeschränkt, bzw.  $IP^=(A, b) = \emptyset$ . Falls (LP) beschränkt ist, so sei  $x^*$  eine Optimallösung von (LP).
2. Ist  $x^*$  ganzzahlig  $\rightarrow$  STOP (gib  $x^*$  aus). Andernfalls bestimme eine Ungleichung  $d^T x \leq d_0$  mit folgenden Eigenschaften (**Schnittebene**)

$$\begin{aligned} d^T x^* &> d_0 \\ d^T x &\leq d_0 \quad \forall x \in IP^=(A, b) \end{aligned}$$

3. Füge die Ungleichung  $d^T x \leq d_0$  unter Hinzufügung einer Schlupfvariablen zum gegenwärtigen System hinzu. Nenne dieses neue System (16.3) und gehe zu 1.

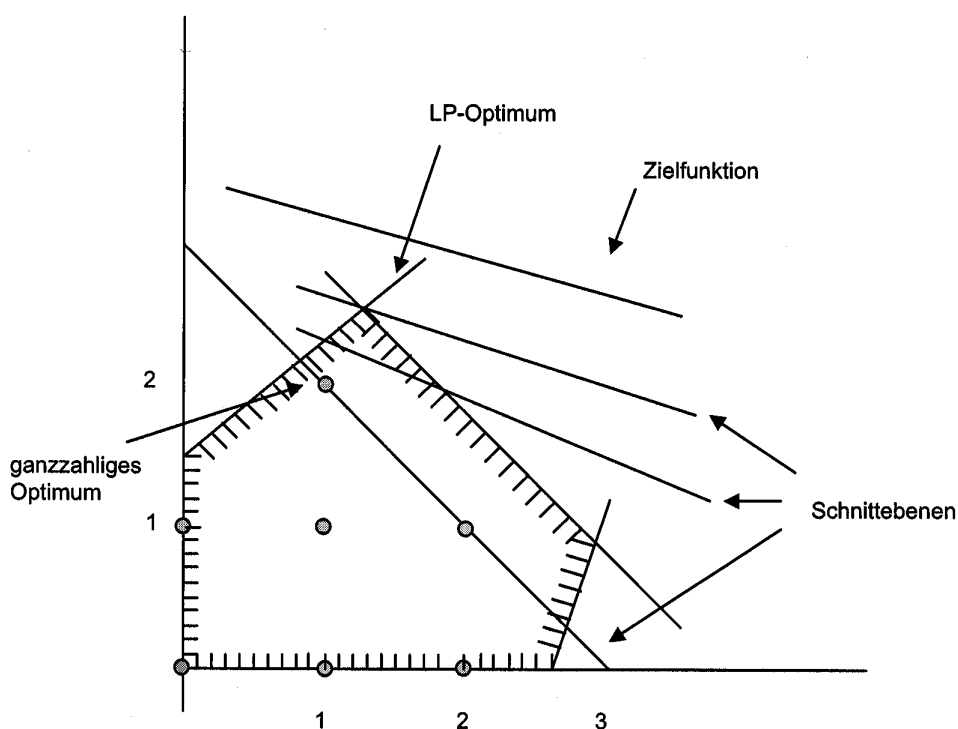
□

Für ganzzahlige Programme der Form (16.2) und für gemischt-ganzzahlige Programme verlaufen Schnittebenenverfahren analog.

Schnittebenenverfahren unterscheiden sich nur durch die Art Bestimmung der Schnittebenen — dafür gibt es verschiedene Techniken — und durch die Art der Lösung des erweiterten linearen Programms. Im Allgemeinen wird Schritt 1 bei der ersten Ausführung mit

dem primalen Simplexverfahren gelöst und nach jeder Ausführung von Schritt 3 mit Hilfe der dualen Simplexmethode.

Der Name Schnittebene kommt daher, dass die gegenwärtige Optimallösung abgeschnitten wird, siehe Abbildung 14.1.



**Abb. 17.1**

Satz (16.20) besagt, dass mit Hilfe des Abschlussoperators  $cl$  im Prinzip Schnittebenen konstruiert werden können, so dass irgendwann die konvexe Hülle der ganzzahligen Punkte erreicht ist. Die Frage ist nun, ob eine (möglichst einfache) Regel angegeben werden kann, die Schnittebenen erzeugt und endliche Konvergenz garantiert. Bis zum Jahre 1958 sind verschiedene Regeln konstruiert worden, die aber alle nicht funktionierten. Gomory gelang dann 1958 die Entwicklung eines endlichen Schnittebenenverfahrens für (rein-)ganzzahlige Programme. Später hat er weitere Techniken zur Schnittebenenerzeugung entworfen, die auch im Falle gemischt-ganzzahliger Programme arbeiten.

Bevor Gomory seine Methode veröffentlichte war von Dantzig (s. Garfinkel-Nemhauser (1972) S. 166) der folgende Schnitt vorgeschlagen worden. Ist eine optimale Lösung des

zu (16.3) gehörigen LP nicht ganzzahlig, so muss mindestens eine der Nichtbasisvariablen  $N$  von Null verschieden sein, wegen der Ganzzahligkeit also mindestens 1, daraus folgt, dass die Ungleichung

$$\sum_{j \in N} x_j \geq 1$$

von allen ganzzahligen Punkten aber nicht von der gegenwärtigen Optimallösung erfüllt wird. Sie liefert also eine Schnittebene. Gomory und Hu haben 1963 ein Beispiel konstruiert, für das ein Verfahren mit diesen Schnittebenen nicht endlich terminiert.

Wir geben nun eine Klasse „guter Schnittebenen“ an.

**(17.2) Lemma.** Gegeben sei ein ganzzahliges Programm  $(IP^=)$   $\max c^T x$ ,  $Ax = b$ ,  $x \geq 0$  und ganzzahlig mit ganzzahligen Daten. Sei  $A_B$  eine Basis der LP-Relaxierung  $(LP^=)$ , und sei  $h \neq 0$  eine ganze Zahl, dann sind

$$(a) \sum_{j \in N} (\lfloor h\bar{a}_{ij} \rfloor - h\bar{a}_{ij})x_j \leq \lfloor h\bar{b}_i \rfloor - h\bar{b}_i, \quad \forall i \in B$$

$$(b) \sum_{j \in N} (\lfloor h\bar{c}_j \rfloor - h\bar{c}_j)x_j \leq \lfloor hc^* \rfloor - hc^*,$$

Schnittebenen, genannt **fundamentale Schnittebenen**, falls  $h\bar{b}_i$  bzw.  $hc^*$  nicht ganzzahlig sind. Dabei seien:  $c^*$  der Wert der Basislösung zu  $A_B$  von  $(LP^=)$ , und (wie üblich)  $B$  die Indizes der Basisvariablen,  $N$  Indizes der Nichtbasisvariablen und

$$\bar{A} = (\bar{a}_{ij}) = A_B^{-1}A_N, \quad \bar{b} = A_B^{-1}b, \quad \bar{c}^T = (c_N^T - c_B^T\bar{A}) \text{ und } \tilde{c} = -\bar{c}.$$

**Beweis :** Jedes System  $Ax = b$  hat bezüglich einer Basis  $A_B$  die Darstellung  $x_B = A_B^{-1}b - A_B^{-1}A_Nx_N = \bar{b} - \bar{A}x_N$ . (Die Basislösung bezüglich  $A_B$  ist gegeben durch  $x_B = \bar{b}$ ,  $x_N = 0$ .) Es gilt also

$$x_i = \bar{b}_i - \bar{A}_i x_N \quad \forall i \in B$$

bzw.

$$(1) \quad x_i + \sum_{j \in N} \bar{a}_{ij}x_j = \bar{b}_i \quad \forall i \in B.$$

Multiplizieren mit  $h \in \mathbb{Z}$ ,  $h \neq 0$ , ergibt

$$(2) \quad hx_i + \sum_{j \in N} h\bar{a}_{ij}x_j = h\bar{b}_i.$$

Da bei Zulässigkeit  $x_k \geq 0 \forall k$  gelten muss, können wir die Koeffizienten auf der linken Seite von (2) abrunden und erhalten so aus der Gleichung (2) eine gültige Ungleichung

$$(3) \quad hx_i + \sum_{j \in N} \lfloor h\bar{a}_{ij} \rfloor x_j \leq h\bar{b}_i.$$

Für ganzzahlige Lösungen ergibt die linke Seite von (3) einen ganzzahligen Wert. Wenn wir somit die rechte Seite ebenfalls abrunden, bleibt die neu entstehende Ungleichung für alle  $x \in IP^=(A, b)$  gültig. Wir erhalten also die für  $IP^=(A, b)$  gültige Ungleichung

$$(4) \quad hx_i + \sum_{j \in N} \lfloor h\bar{a}_{ij} \rfloor x_j \leq \lfloor h\bar{b}_i \rfloor.$$

Ziehen wir nun (2) von (4) ab, so ergibt sich die gültige Ungleichung

$$(5) \quad \sum_{j \in N} (\lfloor h\bar{a}_{ij} \rfloor - ha_{ij}) x_j \leq \lfloor h\bar{b}_i \rfloor - h\bar{b}_i \quad \forall i \in B.$$

Ist nun die gegenwärtige Basislösung nicht ganzzahlig, so gibt es ein  $i \in B$  mit  $\bar{b}_i \notin \mathbb{Z}$ . Ist  $h\bar{b}_i \notin \mathbb{Z}$ , so besagt die Ungleichung (5), dass die linke Seite kleiner oder gleich  $\lfloor h\bar{b}_i \rfloor - h\bar{b}_i < 0$  sein muss. Da  $x_N = 0$  für die Basislösung gilt, erfüllt die gegenwärtige Basislösung die Ungleichung (5) nicht, wird also abgeschnitten.

Der Zielfunktionswert  $c_0 := c^T x$  bezüglich der Basis  $A_B$  lässt sich wie folgt berechnen:

$$c_0 = c^T x = c_B^T x_B + (c_N^T - c_B^T A_B^{-1} A_N) x_N = c_B^T \bar{b} + \bar{c}^T x_N,$$

wobei  $\bar{c}$  die reduzierten Kosten sind. Setzen wir  $\tilde{c} = -\bar{c}$ , so ergibt dies

$$c_0 + \tilde{c}^T x_N = c_B^T \bar{b} =: c^*.$$

Multiplizieren mit  $h \in \mathbb{Z}$  ergibt

$$(6) \quad hc_0 + \sum_{j \in N} h\tilde{c}_j x_j = hc^*.$$

Für zulässige Punkte  $x$  gilt  $x_j \geq 0$ , also liefert Abrunden

$$(7) \quad hc_0 + \sum_{j \in N} \lfloor h\tilde{c}_j \rfloor x_j \leq hc^*.$$

Ist  $x$  ganzzahlig, so ist  $c_0$  ganzzahlig, weil  $c \in \mathbb{Z}^n$  ist; also stehen auf der linken Seite von (7) nur ganze Zahlen, folglich kann die rechte Seite von (7) ebenfalls abgerundet werden, ohne die Gültigkeit für Punkte aus  $IP^=(A, b)$  zu verlieren. Es ergibt sich

$$(8) \quad hc_0 + \sum_{j \in N} \lfloor h\tilde{c}_j \rfloor x_j \leq \lfloor hc^* \rfloor.$$

Ziehen wir (6) von (8) ab, so erhalten wir

$$(9) \quad hc_0 + \sum_{j \in N} \lfloor h\tilde{c}_j \rfloor x_j \leq hc^*.$$

was zu zeigen war. □

Es ist klar, dass ein ganzzahliges Programm mit ganzzahligen Daten als Optimalwert eine ganze Zahl haben muss. Haben wir das zu  $(IP^=)$  gehörige lineare Programm  $(LP^=)$  gelöst, und hat dieses keinen ganzzahligen Optimalwert, so können wir mittels einer Schnittebene des Typs (b) die gegenwärtige Basislösung abschneiden. Ist der Optimalwert des LP ganzzahlig, aber die optimale Basislösung nicht ganzzahlig, so gibt es wegen  $x_N = 0$  eine Basisvariable  $x_i$  mit  $x_i = \bar{b}_i \notin \mathbb{Z}$ . Dann können wir mit einer Ungleichung des Typs (a) wiederum die gegenwärtige Basislösung abschneiden.

Die obigen Überlegungen zeigen also, dass zu jeder nichtganzzahligen Basislösung eines LP eine Schnittebene existiert, die bezüglich aller ganzzahligen Punkte zulässig ist, die die betrachtete Basislösung abschneidet, und die — was wichtig ist — sehr einfach aus den LP-Daten abgeleitet werden kann. Im Prinzip kann dieses Verfahren unendlich lange dauern, wir werden aber zeigen, dass bereits für den Spezialfall  $h = 1$  unter gewissen Voraussetzungen endliche Konvergenz gezeigt werden kann.

**(17.3) Folgerung.** Die Voraussetzungen seien wie in Lemma (17.2). Wir setzen

$$\left. \begin{aligned} f_{ij} &:= \bar{a}_{ij} - \lfloor \bar{a}_{ij} \rfloor & i \in B, j \in N \\ f_{0j} &:= -\bar{c}_j - \lfloor -\bar{c}_j \rfloor & j \in N \\ f_{i0} &:= \bar{b}_i - \lfloor \bar{b}_i \rfloor & i \in B \\ f_{00} &:= c^* - \lfloor c^* \rfloor. \end{aligned} \right\} \quad \begin{array}{l} \text{(diese rationalen Zahlen nennt} \\ \text{man gebrochene Teile)} \end{array}$$

Falls  $f_{i0} \neq 0$  für  $i = 0$  oder  $i \in B$ , dann ist der **erste Gomory-Schnitt**

$$\sum_{j \in N} -f_{ij} x_j \leq -f_{i0}$$

eine Ungleichung, die für alle Punkte aus  $IP^=(A, b)$  zulässig ist und die die gegenwärtige Basislösung abschneidet.  $\square$

Wir stellen nun einen Algorithmus zusammen, der auf der Basis des Simplex-Verfahrens ganzzahlige Programmierungsprobleme löst. In einer ersten Phase wird mit dem primalen Simplexalgorithmus eine optimale Lösung des zugehörigen LP produziert, dann werden Schnittebenen abgeleitet – und zwar erste Gomory-Schnitte – um die das Restriktionssystem erweitert wird. Mit dem dualen Simplexverfahren wird anschließend eine Reoptimierung vorgenommen. Der Algorithmus ist sehr ausführlich dargestellt.

#### (17.4) Erster Gomory-Algorithmus.

**Input:**  $A \in \mathbb{Z}^{(m,n)}, b \in \mathbb{Z}^m, c \in \mathbb{Z}^n$ .

**Output:** Lösung des ganzzahligen Programms  $(IP^=)$

$$\begin{aligned} \max \quad & c^T x \\ & Ax = b \\ & x \geq 0, \quad x \text{ ganzzahlig} \end{aligned}$$

Zur einfacheren Darstellung des Algorithmus führen wir folgende Notation ein:

$$\begin{aligned} a_{00} &:= c_0 && \text{(gegenwärtiger Zielfunktionswert)} \\ a_{0j} &:= -c_j && j = 1, \dots, n \text{ (negativer Wert der Zielfunktionskoeffizienten, um die Schnitte aus der Zielfunktion einfacher ableiten zu können)} \\ a_{i0} &:= b_i && i = 1, \dots, m \text{ (rechte Seite)} \end{aligned}$$

Unser  $(IP^=)$  lautet also in der neuen Notationsform

$$\begin{aligned} \max a_{00} = & - \sum_{j=1}^n a_{0j} x_j \\ & \sum_{j=1}^n a_{ij} x_j = a_{i0} \quad \text{für } i = 1, \dots, m \\ & x_j \geq 0 \quad \text{und ganzzahlig } j = 1, \dots, n \end{aligned}$$

$(LP^=)$  ist das obige Programm ohne Ganzzahligkeitsbedingungen.

1. In einer Phase 1 wird die Zulässigkeit des zugehörigen  $(LP^=)$  geprüft.
  - 1.1 Falls  $(LP^=)$  unzulässig ist, ist natürlich auch  $(IP^=)$  unzulässig.  
→ STOP.
  - 1.2 Ist  $(LP^=)$  zulässig, so endet die Phase 1 mit einer zulässigen Basis. Hierbei können u. U. einige linear abhängige Zeilen von  $Ax = b$  gestrichen worden sein. Wir würden dann mit dem reduzierten System beginnen. Der Einfachheit halber nehmen wir an, dass  $A$  vollen Rang  $m$  hat. Wir erhalten also

$$B = (p_1, p_2, \dots, p_m), \quad N = (q_1, q_2, \dots, q_{n-m})$$

und somit eine Basismatrix  $A_B$  von  $A$

2. Wir berechnen

$$\begin{aligned} \bar{A} &= A_B^{-1} A_N && (m, n-m)\text{-Matrix} \\ \bar{a}_{\cdot 0} &= A_B^{-1} b = A_B^{-1} a_{\cdot 0} && m\text{-Vektor} \\ \bar{a}_{0\cdot} &= -(c_N^T - c_B^T A_B^{-1} A_N) && (n-m)\text{-Vektor} \\ &= a_{0,N}^T - a_{0,B}^T \bar{A} \\ \bar{a}_{00} &= c_B^T A_B^{-1} b = -a_{0,B}^T \bar{a}_{\cdot 0}. \end{aligned}$$

**Primaler Simplex-Algorithmus** (Voraussetzung: zulässige Basis bekannt).

3. (Optimalitätsprüfung)
 

Gilt  $\bar{a}_{0j} \geq 0$  für  $j = 1, \dots, n-m$ , so ist die gegenwärtige Basislösung optimal. Gehe zu Schritt 8. Andernfalls gehe zu Schritt 4.

4. (Bestimmung der Austauschspalte)  
Wähle einen Index  $s \in \{1, \dots, n - m\}$ , so dass  $\bar{a}_{0s} < 0$  gilt. (Hierzu haben wir mehrere Möglichkeiten kennengelernt; z. B., wähle  $s$ , so dass  $\bar{a}_{0s}$  minimal ist.)
5. (Prüfung auf Beschränktheit des Optimums)  
Gilt  $\bar{a}_{is} \leq 0$  für  $i = 1, \dots, m$ , so ist das lineare Programm nach “Optimierungsmethoden I”, Satz 9.14 (a) unbeschränkt, also ist entweder ( $IP^=$ ) unbeschränkt oder hat keine zulässige Lösung.  $\rightarrow$  STOP.
6. (Bestimmung der Austauschzeile)

$$6.1 \text{ Berechne } \lambda_0 := \min \left\{ \frac{\bar{a}_{i0}}{\bar{a}_{is}} \mid \bar{a}_{is} > 0, i = 1, \dots, m \right\}$$

6.2 Wähle einen Index  $r \in \{1, \dots, m\}$ , so dass gilt

$$\frac{\bar{a}_{r0}}{\bar{a}_{rs}} = \lambda_0.$$

(Hierzu haben wir in Kapitel 10 von “Optimierungsmethoden I” mehrere Möglichkeiten kennengelernt, die die Endlichkeit des Verfahrens garantieren.)

7. (Basisaustausch, Pivotoperationen)

7.1 Setze

$$\begin{aligned} B' &:= (p_1, p_2, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m) \\ N' &:= (q_1, q_2, \dots, q_{s-1}, p_r, q_{s+1}, \dots, q_{n-m}) \end{aligned}$$

7.2 (Neuberechnung der erweiterten Matrix  $\bar{A}$ )

Führe folgende Pivotoperationen durch

$$7.2.1 \quad \bar{\bar{a}}_{rs} := \frac{1}{\bar{a}_{rs}}$$

$$7.2.2 \quad \bar{\bar{a}}_{rj} := \frac{\bar{a}_{rj}}{\bar{a}_{rs}}, j = 0, 1, \dots, n - m; j \neq s.$$

$$7.2.3 \quad \bar{\bar{a}}_{is} := -\frac{\bar{a}_{is}}{\bar{a}_{rs}}, i = 0, 1, \dots, m; i \neq r.$$

$$7.2.4 \quad \bar{\bar{a}}_{ij} := \bar{a}_{ij} - \frac{\bar{a}_{is}\bar{a}_{rj}}{\bar{a}_{rs}}, i = 0, 1, \dots, m; i \neq r, j = 0, 1, \dots, n - m; j \neq s.$$

7.3 Setze  $B := B'$ ,  $N := N'$ , und  $\bar{a}_{ij} := \bar{\bar{a}}_{ij}$   $i = 0, \dots, m; j = 0, \dots, n - m$  und gehe zu Schritt 3.

**Ende des primalen Simplex-Algorithmus**  
**Bestimmung der Schnittebenen, Optimalitätstest**

## 8. (Prüfung auf Ganzzahligkeit der Lösung)

Gilt  $\bar{a}_{i0} \in \mathbb{Z}$  für  $i = 0, 1, \dots, m$ , so ist die optimale Lösung des gegenwärtigen LP ganzzahlig und folglich auch eine optimale Lösung von  $(IP^=)$ .

Ausgabe: Zielfunktionswert  $\bar{a}_{00}$

Optimallösung

$$\begin{aligned} x_i &= 0 & \forall i \in N \\ x_i &= \bar{a}_{i0} & \forall i \in B \end{aligned}$$

STOP.

Andernfalls gehe zu Schritt 9.

## 9. (Ableitung einer neuen Restriktion, Schnittebenengenerierung)

Wähle ein  $i \in \{0, 1, \dots, m\}$  so dass  $\bar{a}_{i0} \notin \mathbb{Z}$  (hierfür kann man sich mehrere Strategien überlegen), und setze für  $j = 0, 1, \dots, n - m$

$$\bar{a}_{m+1,j} := \bar{a}_{ij} - \lfloor \bar{a}_{ij} \rfloor.$$

Füge die Zeile

$$\bar{a}_{m+1,0}, \bar{a}_{m+1,1}, \dots, \bar{a}_{m+1,n-m}$$

als  $(m + 1)$ -ste Zeile zum gegenwärtigen Tableau hinzu.

## 10. (Basiserweiterung)

Nimm die zur neu hinzugefügten Ungleichung gehörige Schlupfvariable  $x_{n+1}$  mit dem Wert  $\bar{a}_{m+1,0}$  in die Basis auf, d. h. setze  $B = (p_1, p_2, \dots, p_m, n+1)$  und setze

$$m := m + 1, \quad n := n + 1.$$

Gehe zu Schritt 11.

**Duales Simplex-Verfahren** (Voraussetzung: dual zulässige Basis bekannt)

In Abschnitt 10.4 von "Optimierungsmethoden I" haben wir das duale Simplex-Verfahren (siehe 10.27) und das Vorgehen bei Hinzufügung einer neuen Restriktion, falls eine Basislösung des linearen Programms bekannt ist, behandelt (siehe 11.5). Das zum (primalen) linearen Programm

$$(LP^=) \max c^T x, Ax = b, x \geq 0$$

duale Programm lautet  $\min u^T b, u^T A \geq c$ . Die zu einer Basis  $A_B$  von  $A$  gehörige Basislösung ist  $x_B = A_B^{-1}b, x_N = 0$ . Falls  $x_B \geq 0$  ist, heißt diese Basislösung primal zulässig, da sie eine zulässige Lösung von  $(LP^=)$  darstellt. Eine Basislösung heißt dual zulässig, wenn der Vektor  $c_B^T A_B^{-1}$  eine zulässige Lösung des dualen Programms ist. Eine Basislösung ist also genau dann dual zulässig, wenn  $c_B^T A_B^{-1} A \geq c^T$  gilt, bzw. (1)  $c_B^T A_B^{-1} A_B \geq c_B^T$  und (2)  $c_B^T A_B^{-1} A_N \geq c_N^T$  gilt. Bedingung (1) ist natürlich immer erfüllt, Bedingung (2) ist erfüllt, wenn  $\bar{c}^T := c_N^T - c_B^T A_B^{-1} A_N \leq 0$  ist. Dies ist aber nach (9.13) (b) "Optimierungsmethoden I" gerade die primale Optimalitätsbedingung (die reduzierten Kosten sind

nicht positiv). Es folgt, dass eine Basislösung primal und dual zulässig ist genau dann, wenn sie optimal ist.

Erweitern wir unser Restriktionssystem um eine Schnittebene, so können wir durch Aufnahme der Schlupfvariablen in die Basis sofort eine neue Basis bzw. Basislösung angeben. Da der Wert der neuen Basisvariablen nach Konstruktion negativ ist, ist die Basislösung nicht primal zulässig. Sie ist jedoch dual zulässig, da sich die reduzierten Kosten nicht geändert haben. Wir haben somit eine Startbasis zur Ausführung des dualen Simplexverfahrens.

11. (Optimalitätstest, d. h. Test auf primale Zulässigkeit)

Gilt  $\bar{a}_{i0} \geq 0$ ,  $i = 1, \dots, m$  so ist eine optimale Lösung des gegenwärtigen linearen Programms gefunden, gehe zu Schritt 8.

Andernfalls gehe zu Schritt 12.

12. (Bestimmung der in die duale Basis eintretenden Variablen)

Wähle einen Index  $r \in \{1, \dots, m\}$  mit  $\bar{a}_{r0} < 0$ . (Hierzu gibt es wiederum mehrere Strategien, z. B.  $\bar{a}_{r0}$  minimal.)

13. (Prüfung auf Beschränktheit des Optimums)

Gilt  $\bar{a}_{rj} \geq 0$ , für  $j = 1, \dots, n-m$ , so hat das duale Programm keine endliche Optimallösung, d. h. das primale Programm hat überhaupt keine Lösung, also gilt  $IP^=(A, b) = \emptyset$ .  $\rightarrow$  STOP.

14. (Bestimmung der aus der dualen Basis austretenden Variablen)

14.1 Berechne  $\lambda_0 := \max \left\{ \frac{\bar{a}_{0j}}{\bar{a}_{rj}} \mid \bar{a}_{rj} < 0, j = 1, \dots, n-m \right\}$ .

14.2 Wähle einen Index  $s \in \{1, \dots, n-m\}$ , so dass gilt

$$\frac{\bar{a}_{0s}}{\bar{a}_{rs}} = \lambda_0$$

(mehrere Strategien möglich, falls  $s$  nicht eindeutig bestimmt ist).

15. (Basisaustausch, Pivotoperation)

15.1 Setze

$$\begin{aligned} B' &:= (p_1, p_2, \dots, p_{r-1}, q_s, p_{r+1}, \dots, p_m) \\ N' &:= (q_1, q_2, \dots, q_{s-1}, p_r, q_{s+1}, \dots, q_{n-m}). \end{aligned}$$

15.2 (Neuberechnung des Tableaus)

$$15.2.1 \quad \bar{\bar{a}}_{rs} := \frac{1}{\bar{a}_{rs}}$$

$$15.2.2 \quad \bar{\bar{a}}_{rj} := \frac{\bar{a}_{rj}}{\bar{a}_{rs}}, j = 0, 1, \dots, n-m; j \neq s.$$

$$15.2.3 \quad \bar{\bar{a}}_{is} := -\frac{\bar{a}_{is}}{\bar{a}_{rs}}, i = 0, 1, \dots, m; i \neq r.$$

$$15.2.4 \quad \bar{\bar{a}}_{ij} := \bar{a}_{ij} - \frac{\bar{a}_{is}\bar{a}_{rj}}{\bar{a}_{rs}}, i = 0, 1, \dots, m; i \neq r, j = 0, 1, \dots, n-m; j \neq s.$$

15.3 Setze  $B := B'$ ,  $N := N'$ , und

$$\bar{a}_{ij} := \bar{\bar{a}}_{ij} \quad \text{für } i = 0, \dots, m; j = 0, \dots, n - m$$

und gehe zu Schritt 11.

**Ende des dualen Simplex-Algorithmus.**

**Ende des ersten Gomory-Verfahrens.** □

**(17.5) Beispiel.** Wir lösen das folgende ganzzahlige Programm mit dem ersten Gomory-Algorithmus. Die Lösungsmenge ist in Abbildung 17.2 dargestellt.

$$\begin{array}{ll} \max & x_1 + 2x_2 \\ x_3 \longrightarrow & (1) \quad x_1 \leq 4 \\ x_4 \longrightarrow & (2) \quad 2x_1 + x_2 \leq 10 \\ x_5 \longrightarrow & (3) \quad -x_1 + x_2 \leq 5 \\ & x_1, x_2 \geq 0, \\ & x_1, x_2 \text{ ganzzahlig} \end{array}$$

↑  
Schlupfvariable

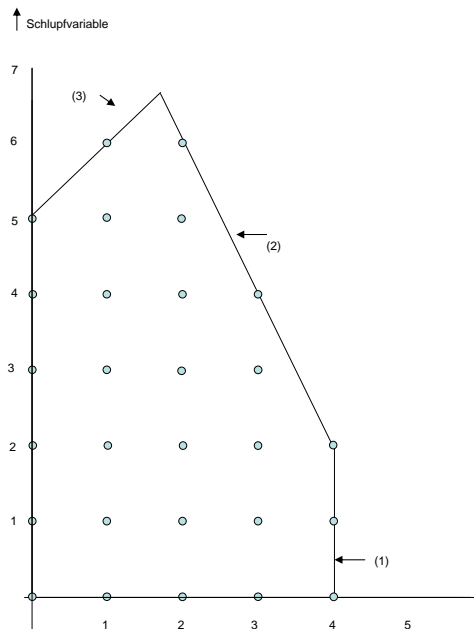


Abb. 17.2

(a) Anfangsbasis = Schlupfvariable,  $B = (3, 4, 5)$ ,  $N = (1, 2)$ .

1. Tableau

			1	2	$\leftarrow N$	
	0		-1	-2		
3	4	1	0			$s = 2$
4	10	2	1			$r = 3$
5	5	-1	1			
$\uparrow$						
$B$						

Eine Pivotoperation ergibt das 2. Tableau.

			1	5	
	10		-3	2	
3	4	1	0		$s = 1$
4	5	3	-1		$r = 2$
2	5	-1	1		

3. Tableau

			4	5	
	15		1	1	Optimalitätstest erfüllt
3	$\frac{7}{3}$	$-\frac{1}{3}$	$\frac{1}{3}$		
1	$\frac{5}{3}$	$\frac{1}{3}$	$-\frac{1}{3}$		
2	$\frac{20}{3}$	$\frac{1}{3}$	$\frac{2}{3}$		Optimallösung $x_1 = \frac{5}{3}, x_2 = \frac{20}{3}$ $c^* = 15$

Ende des primalen Simplexverfahrens

b) **Schnittebenenbestimmung**

Außer der Zielfunktionszeile können alle anderen Zeilen zur Schnittebenenberechnung benutzt werden, da  $\bar{a}_{i0} \notin \mathbb{Z}$ ,  $i = 1, 2, 3$ . Wir bestimmen alle drei Schnittebenen (Kandidaten für die neue vierte Zeile des Tableaus).

$$1. \text{ Zeile: } \bar{a}_{40} = -\left(\frac{7}{3} - \frac{6}{3}\right) = -\frac{1}{3}, \bar{a}_{41} = -\left(\frac{1}{3} - \frac{-3}{3}\right) = -\frac{2}{3}, \bar{a}_{42} = -\left(\frac{1}{3} - 0\right) = -\frac{1}{3},$$

Daraus resultiert die Schnittebene

$$-\frac{2}{3}x_4 - \frac{1}{3}x_5 \leq -\frac{1}{3}.$$

In den ursprünglichen Variablen — wenn wir  $x_4 = 10 - 2x_1 - x_2$  und  $x_5 = 5 + x_1 - x_2$  substituieren — ergibt dies:

$$\begin{aligned} -\frac{2}{3}(10 - 2x_1 - x_2) - \frac{1}{3}(5 + x_1 - x_2) &\leq -\frac{1}{3} \\ -20 + 4x_1 + 2x_2 - 5 - x_1 + x_2 &\leq -1 \\ 3x_1 + 3x_2 &\leq 24 \end{aligned}$$

Unsere Schnittebene ist also äquivalent zur Ungleichung:

$$x_1 + x_2 \leq 8.$$

2. Zeile:  $\bar{a}_{40} = -\frac{2}{3}$ ,  $\bar{a}_{41} = -\frac{1}{3}$ ,  $\bar{a}_{42} = -\frac{2}{3}$

$$-\frac{1}{3}x_4 - \frac{2}{3}x_5 \leq -\frac{2}{3}$$

bzw. in ursprünglichen Variablen

$$\begin{aligned} -\frac{1}{3}(10 - 2x_1 - x_2) - \frac{2}{3}(5 + x_1 - x_2) &\leq -\frac{2}{3} \\ -10 + 2x_1 + x_2 - 10 - 2x_1 + 2x_2 &\leq -2 \\ 3x_2 &\leq 18 \end{aligned}$$

das heißt, die Schnittebene ist äquivalent zu:

$$x_2 \leq 6.$$

3. Zeile:  $\bar{a}_{40} = -\frac{2}{3}$ ,  $\bar{a}_{41} = -\frac{1}{3}$ ,  $\bar{a}_{42} = -\frac{2}{3}$  also

$$-\frac{1}{3}x_4 - \frac{2}{3}x_5 \leq -\frac{2}{3}.$$

Damit ergibt sich dieselbe Schnittebene wie bei der Ableitung aus der zweiten Zeile.

Die beiden auf diese Weise gefundenen Schnittebenen sind in Abbildung 17.3 dargestellt.

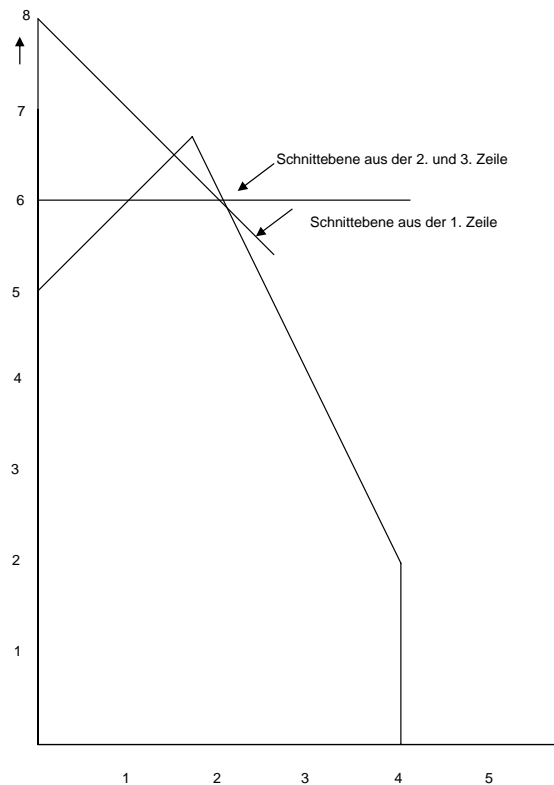


Abb. 17.3

An dem obigen Beispiel kann man sehen, dass die Gomory-Schnitte verschieden “tief” sein können. Der Schnitt aus der zweiten Zeile ist tiefer, weil der Durchschnitt von  $P$  mit dem zugehörigen Halbraum echt im Durchschnitt von  $P$  mit dem aus der ersten Zeile abgeleiteten Halbraum enthalten ist. Die Schnittebene aus der zweiten Zeile ergibt zusammen mit den ursprünglichen Ungleichungen sogar die konvexe Hülle der ganzzahligen Punkte. Daraus folgt, dass bei Hinzufügung dieses Schnittes das ganzzahlige Optimum erreicht wird.

Bei normalem Durchlauf des Algorithmus hätten wir wahrscheinlich den ersten Schnitt gewählt, deshalb fügen wir diesen zusammen mit der neuen Schlupfvariablen  $x_6$  zu unserem Tableau hinzu.

Neues Tableau

		4	5
	15	1	1
3	$\frac{7}{3}$	$-\frac{1}{3}$	$\frac{1}{3}$
1	$\frac{5}{3}$	$\frac{1}{3}$	$-\frac{1}{3}$
2	$\frac{20}{3}$	$\frac{1}{3}$	$\frac{2}{3}$
6	$-\frac{1}{3}$	$-\frac{2}{3}$	$-\frac{1}{3}$

← neue Zeile,  
Schlupfvariable  $x_6$  in der Basis

c) Duales Simplexverfahren angewendet auf das obige Tableau

$$\begin{aligned} r &= 4 \\ \lambda_0 &= -\frac{3}{2} \\ s &= 1 \end{aligned}$$

Die Pivotoperation ergibt:

		6	5
	$\frac{29}{2}$	$\frac{3}{2}$	$\frac{1}{2}$
3	$\frac{5}{2}$	$-\frac{1}{2}$	$\frac{1}{2}$
1	$\frac{3}{2}$	$\frac{1}{2}$	$-\frac{1}{2}$
2	$\frac{13}{2}$	$\frac{1}{2}$	$\frac{1}{2}$
4	$\frac{1}{2}$	$-\frac{3}{2}$	$\frac{1}{2}$

dieses Tableau ist optimal,  
jedoch nicht ganzzahlig

d) Ableitung eines neuen Schnittes aus der Zielfunktion  $\bar{a}_{50} = -\frac{1}{2}$ ,  $\bar{a}_{51} = -\frac{1}{2}$ ,  $\bar{a}_{52} = -\frac{1}{2}$ . Daraus ergibt sich

$$-\frac{1}{2}x_5 - \frac{1}{2}x_6 \leq -\frac{1}{2}.$$

Wegen  $x_5 = 5 + x_1 - x_2$  und  $x_6 = 8 - x_1 - x_2$  sieht diese Ungleichung in den ursprünglichen Variablen wie folgt aus

$$x_2 \leq 6.$$

Diese Ungleichung hätten wir bereits nach Lösung des Anfangs-LP mit dem primalen Verfahren aus der 2. oder 3. Zeile erhalten können.

Neues Tableau

		6	5	
	$\frac{29}{2}$	$\frac{3}{2}$	$\frac{1}{2}$	
3	$\frac{5}{2}$	$-\frac{1}{2}$	$\frac{1}{2}$	
1	$\frac{3}{2}$	$\frac{1}{2}$	$-\frac{1}{2}$	
2	$\frac{13}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	
4	$\frac{1}{2}$	$-\frac{3}{2}$	$\frac{1}{2}$	
7	$-\frac{1}{2}$	$-\frac{1}{2}$	<span style="border: 1px solid black; padding: 2px;"><math>-\frac{1}{2}</math></span>	← neue Zeile, Schlupfvariable $x_7$

e) erneute Anwendung des dualen Simplexverfahrens

		6	7	
	14	1	1	
3	2	-1	1	$r = 5$
1	2	1	-1	$s = 2$
2	6	0	1	
4	0	-2	1	
5	1	1	-2	
	↑			

entartete Basislösung, wegen  $x_4 = 0$ 

Optimallösung gefunden:

$$\begin{aligned} x_1 &= 2 \\ x_2 &= 6 \\ c^* &= 14 \end{aligned}$$

Der Gesamttablauf des Verfahrens ist in Abbildung 17.4 zusammengefasst.

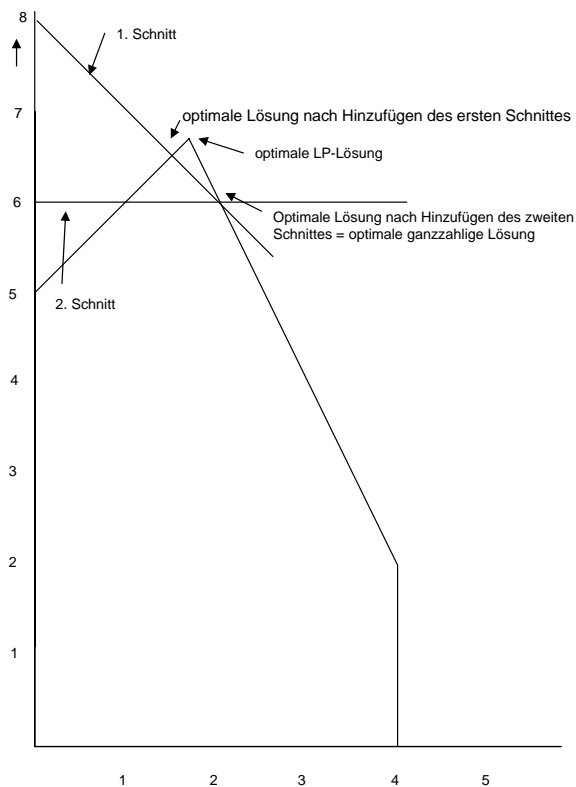


Abb. 17.4

**(17.6) Beobachtungen zum ersten Gomory-Algorithmus.**

- Zeigt das primale Simplexverfahren in Schritt 5 die Unbeschränktheit der LP-Relaxierung an, so ist  $(IP^=)$  nicht sinnvoll lösbar, d. h. entweder unbeschränkt oder  $IP^=(A, b) = \emptyset$ .
- Wird in Schritt 13 die Unbeschränktheit des dualen Programms zum gegenwärtigen LP (dieses ist das um Schnittebenen erweiterte ursprüngliche LP) festgestellt, so ist  $IP^=(A, b) = \emptyset$ .
- Obwohl an zwei Stellen des Algorithmus die Unzulässigkeit von  $(IP^=)$  festgestellt werden kann, gibt es keine Garantie dafür, dass der Algorithmus in dieser Form die Unzulässigkeit auch beweist. Wir können diesen Mangel mit Folgerung (16.14) (oder (16.16)) reparieren und einen Zusatzschritt formulieren.

- (d) Falls das primale oder duale Simplex-Verfahren eine nicht ganzzahlige Lösung liefert, so kann das Schnittebenen-Unterprogramm theoretisch immer einen Gomory-Schnitt finden, der diese Lösung abschneidet. Wir schreiben hier “theoretisch”, da bei praktischen Rechnungen mit reeller floating-point-Arithmetik aufgrund von Rundfehlern und Toleranzgrößen die Entscheidung, ob eine Größe ganzzahlig ist oder nicht, schwierig zu fällen ist.

□

Bevor wir uns dem Beweis der Korrektheit von Algorithmus (17.4) zuwenden, formulieren wir einige Zusatzregeln — die allerdings nur von theoretischem Interesse sind.

### (17.7) Sonderregeln zum ersten Gomory-Algorithmus.

- (a) Ein Vektor  $a$  heißt **lexikographisch größer** als ein Vektor  $b$ , wenn entweder  $a_1 > b_1$  gilt oder ein  $k$  existiert mit  $1 < k \leq n$ ,  $a_i = b_i$  für  $i = 1, \dots, k-1$  und  $a_k > b_k$ . Die Spalten des Simplextableaus (inklusive der 0-ten Zeile) bezeichnen wir wie immer mit  $\bar{A}_{.j}$ , d. h.

$$\bar{A}_{.j} = \begin{pmatrix} \bar{a}_{0j} \\ \bar{a}_{1j} \\ \vdots \\ \bar{a}_{mj} \end{pmatrix} \quad j = 1, \dots, n - m.$$

Wir gestalten unser Verfahren so, dass alle Spalten  $\bar{A}_{.j}$ ,  $j = 1, \dots, n - m$  nach Ende des primalen Simplexverfahrens lexikographisch positiv sind und weiterhin bleiben.

Ist die optimale LP-Lösung nicht **dual entartet**, d. h. sind alle (negativen) reduzierten Kosten  $\bar{a}_{0j}$  positiv, so ist die lexikographische Positivität offenbar gegeben. Da aufgrund der Optimalität  $\bar{a}_{0j} \geq 0$   $j = 1, \dots, n - m$  gilt, kann kein  $\bar{a}_{0j} < 0$  sein. Probleme treten auf, wenn einige  $\bar{a}_{0j} = 0$  sind. Durch Umstellen der Zeilen des Tableaus kann u. U. lexikographische Positivität erreicht werden. Falls dies nicht geht, dann fügen wir eine Ungleichung hinzu, die die lexikographische Positivität gewährleistet und keine Ecke von  $\text{conv}(IP=(A, b))$  abschneidet. Diese Ungleichung kann ganzzahlige Punkte abschneiden, lässt aber mindestens eine der ganzzahligen Optimallösungen, falls überhaupt welche existieren, übrig. Eine Ungleichung dieses Typs ist z. B.

$$\sum_{j \in N} x_j \leq (n - m)(n + 1)2^{4n^3 \langle A, b \rangle}$$

was aus (16.14) folgt.

Diese Ungleichung fügen wir dann als erste Zeile unterhalb der Zielfunktion in unser Tableau ein, wobei eine neue Schlupfvariable addiert werden muss. Die

Schlupfvariable geht mit dem Wert der rechten Seite in die Basis. Nunmehr ist unser Tableau lexikographisch positiv, da  $\bar{a}_{0j} \geq 0$  und  $\bar{a}_{1j} = 1$  für  $j = 1, \dots, n - m$ .

- (b) Zur Beibehaltung der lexikographischen Positivität muss eine besondere Spaltenauswahlregel benutzt werden.

14.2 Sei  $S = \{j \in \{1, \dots, n - m\} \mid \frac{\bar{a}_{0j}}{\bar{a}_{rj}} = \lambda_0, \bar{a}_{rj} < 0\}$ .  
Wähle denjenigen Index  $s \in S$ , so dass  $\frac{1}{\bar{a}_{rs}} \bar{A}_{\cdot s}$  der lexikographisch größte der Vektoren  $\frac{1}{\bar{a}_{rj}} \bar{A}_{\cdot j}$ ,  $j \in S$  ist.

**Hausaufgabe:** Ist  $r$  eine in Schritt 12 gewählte Pivotzeile und  $s$  eine mit obiger Regel bestimmte Pivotspalte, dann ist das Tableau nach der Pivotoperation weiterhin lexikographisch positiv.

- (c) Wir müssen noch gewährleisten, dass der Algorithmus erkennt, dass  $(IP^=)$  unbeschränkt oder  $IP^=(A, b)$  leer ist. Hierzu gibt es zwei Möglichkeiten. Aus den Folgerungen (16.14), (16.16) kann man ableiten, dass der Absolutwert des Optimalwertes nicht größer als  $2^{\langle c \rangle + 5n^4 \langle A, b \rangle - 2n}$  ist bzw. dass die Menge der Optimallösung von  $(IP^=)$  im Würfel

$$W := \{x \in \mathbb{Z}^n \mid -(n+1)2^{4n^3 \langle A, b \rangle} \leq x_i \leq (n+1)2^{4n^3 \langle A, b \rangle}\}$$

liegt. Wir können daher entweder die Restriktion von  $W$  zu unserem ursprünglich ganzzahligen Programm hinzufügen (auf diese Weise haben wir allgemeine ganzzahlige Programme auf beschränkte ganzzahlige Programme reduziert) oder wir bauen in Schritt 8 die folgende Abfrage ein:

**Zusatz zu Schritt 8**  
Gilt für den gegenwärtigen Zielfunktionswert  $\bar{a}_{00}$  :  $\bar{a}_{00} < -2^{\langle c \rangle + 5n^4 \langle A, b \rangle - 2n}$ ,  
so besitzt das ganzzahlige Programm keine zulässige Lösung.  $\rightarrow$  STOP.

Es muss darauf hingewiesen werden, dass dieses Abbruchkriterium, die Entdeckung der Unzulässigkeit des  $(IP^=)$  nur theoretisch gewährleistet, denn die obige Schranke ist i. a. ungeheuer groß und kaum auf dem Rechner darstellbar. Der Gomory-Algorithmus dürfte selbst bei sehr kleinen Beispielen Tausende von Jahren brauchen, um mit diesem Kriterium abzurechnen.

- (d) Ferner müssen wir noch gewährleisten, dass die Schnittebenen in einer bestimmten Weise ausgewählt werden, und zwar leiten wir den ersten Gomory-Schnitt immer aus der am weitesten oben stehenden Zeile des Tableaus ab

**Zusatz zu Schritt 9**  
Wähle den kleinsten Zeilenindex  $i \in \{0, 1, \dots, m\}$  mit  $\bar{a}_{i0} \notin \mathbb{Z}$ .

- e) Wir wollen nun noch eine Regel einführen, die garantiert, dass die Zahl unserer Zeilen, d. h. die Zahl der Restriktionen nicht beliebig groß werden kann. Angenommen nach Ende des dualen Simplexverfahrens ist eine Schlupfvariable  $x_{pi}$ ,

die zu einer Gomory-Schnittebene gehört, Basisvariable. Dann hat  $x_{pi}$  den Wert  $\bar{a}_{i0} \geq 0$ . Ist  $x_{pi} > 0$ , so bedeutet das, dass die Schnittebene bezüglich der gegenwärtigen Lösung nicht bindend ist. Streichen wir die Schnittebene aus unserem gegenwärtigen Restriktionensystem und eliminieren wir die Schlupfvariable, so ist unsere gegenwärtige Lösung eine optimale Basislösung des neuen Systems. Ist  $x_{pi} = 0$ , so ist die Basislösung entartet, zwar liegt die gegenwärtige Lösung auf der Schnittebene, jedoch können wir die Schnittebene weglassen, ohne die Optimalität zu verlieren. Das Streichen einer Zeile des Tableaus kann jedoch dazu führen, dass das Tableau nicht mehr lexikographisch positiv ist.

Man kann jedoch zeigen (siehe Garfinkel & Nemhauser (1972), S. 163–164): Wird bei einer Pivotoperation eine Variable, die Schlupfvariable einer Schnittebene ist, in die Basis aufgenommen, so können die Pivotzeilen nach der Pivotoperation gestrichen und die Schlupfvariable eliminiert werden, so dass das reduzierte Tableau weiterhin lexikographisch positiv ist.

**Zusatz zu Schritt 9**

15.2.5 Ist die neue Basisvariable die Schlupfvariable einer Schnittebene, d. h. ist ihr Index größer als  $n$ , so streichen wir die  $r$ -te Zeile des Tableaus und das  $r$ -te Element von  $B'$ . Anschließend setzen wir  $m = m - 1$ ,  $n = n - 1$  und numerieren die Zeilen mit Index  $\geq r + 1$  und die Schlupfvariablen neu.

Diese Zusatzregel garantiert uns, dass nur ursprüngliche Strukturvariable in der Basis sind, also kann unser Tableau niemals mehr als  $n$  Zeilen (plus Zielfunktion) enthalten.  $\square$

Wir werden anschließend zeigen, dass mit diesen Zusatzregeln Konvergenz des Gomory-Verfahrens bewiesen werden kann. Bei echten Implementationen wird jedoch kaum die lexikographische Spaltenauswahlregel (b) benutzt, da andere Regeln viel schneller zu einer Optimallösung des dualen Simplexverfahrens führen. Konvergiert das Gomory-Verfahren mit den anderen Regeln nicht, so konvergiert es mit der lexikographischen Regel in der Praxis auch nicht. Bei den Schnittebenen nimmt man meistens solche, die einen besonders "tiefen" Schnitt ermöglichen und nicht die erste mögliche. In der Praxis ist das Weglassen von Restriktionen ebenfalls umstritten, wobei hier keine "guten" Strategien angegeben werden können.

**(17.8) Satz.** Sei  $(IP^=)$  ein ganzzahliges Programm der Form (16.3) mit ganzzahligen Daten, dann liefert der erste Gomory-Algorithmus (14.4) mit den Zusatzregeln (14.7) nach endlich vielen Schritten eine ganzzahlige Lösung, oder er weist nach endlich vielen Schritten nach, dass entweder  $(IP^=)$  unbeschränkt oder  $(IP^=(A, b))$  leer ist.

**Beweis :** Der Beweis verläuft wie folgt: Wir wissen, dass bei jedem Schritt des dualen Simplexalgorithmus der Zielfunktionswert abnimmt oder gleich bleibt. Also kann nach einem Durchlauf des dualen Simplexalgorithmus der Zielfunktionswert nicht größer geworden sein. Wir zeigen nun zuerst, dass bei Ausführung des Gomory-Verfahrens der Zielfunktionswert nach endlich vielen Schritten einen festen ganzzahligen Wert annimmt und sich der Zielfunktionswert von da an nicht mehr ändert. Dann zeigen wir, dass nach weiteren endlich vielen Schritten das Tableauelement  $\bar{a}_{10}$  einen festen ganzzahligen Wert annimmt und diesen Wert fortan beibehält. Wir führen diese Argumentation weiter und zeigen, dass das gleiche für  $\bar{a}_{20}, \bar{a}_{30}, \dots$  gilt. Aufgrund der "Streichregel" (Zusatz zu Schritt 15) wissen wir, dass unser Tableau niemals mehr als  $n$  Zeilen enthalten kann. Da nach endlich vielen Schritten  $\bar{a}_{10}$  ganzzahlig ist und der Wert sich nicht mehr verändert, da nach weiteren endlich vielen Schritten das gleiche für  $\bar{a}_{20}$  gilt usw., ist nach endlich vielen Schritten das letzte Tableauelement  $\bar{a}_{m0}$  ( $m \leq n$ ) ganzzahlig. Damit ist dann eine ganzzahlige Lösung des  $(IP^=)$  gefunden.

Liefert der Gomory-Algorithmus nach endlich vielen Durchläufen die Unbeschränktheit oder die Unzulässigkeit von  $(IP^=)$ , so ist der zweite Teil unserer Behauptung gezeigt.

Damit der Gomory-Algorithmus nicht nach endlich vielen Schritten abbricht, müssen also das primale LP und anschließend jedes duale LP ein endliches Optimum haben. Aufgrund

der Zusatzregel zu Schritt 8 kann der Zielfunktionswert zum dualen Programm niemals kleiner werden als  $-2^{(c)+5n^4(A,b)-2n}$ , daraus folgt, dass die Folge  $(\bar{a}_{00}^k)_{k \in \mathbb{N}}$  der Zielfunktionswerte nach dem  $k$ -ten Durchlauf des dualen Programms eine nach unten beschränkte, monoton fallende Folge ist. Diese Folge konvergiert also.

Wir untersuchen nun, wie sich der Algorithmus verhält, wenn der Optimalwert  $\bar{a}_{00}^k$  des  $k$ -ten dualen Programms nicht ganzzahlig ist.

Die Zusatzregel (d) zu Schritt 9 besagt, dass die neue Schnittebene aus der Zielfunktionszeile abgeleitet werden muss. Fügen wir die Schnittebene zum Tableau hinzu, so wird die neue Schlupfvariable mit dem Wert  $\bar{a}_{m+1,0} = \lfloor a_{00}^k \rfloor - a_{00}^k < 0$  Basisvariable. Da das Tableau nach Ende des dualen Verfahrens primal und dual zulässig ist, d. h.  $\bar{a}_{i0} \geq 0$ ,  $i = 1, \dots, m$  und  $\bar{a}_{0j} \geq 0$ ,  $j = 1, \dots, n$ , ist die neue Basisvariable die einzige mit negativem Wert. Folglich wird die neu hinzugefügte Zeile  $m + 1$  in Schritt 12 als Pivotzeile gewählt. Sei  $s$  die in Schritt 14 gewählte Pivotspalte, dann gilt  $\bar{a}_{m+1,s} < 0$ . Weiterhin impliziert  $\bar{a}_{0s} \geq 0$ , dass  $\bar{a}_{0s} \geq \bar{a}_{0s} - \lfloor \bar{a}_{0s} \rfloor = f_{0s} = -\bar{a}_{m+1,s}$  gilt (siehe Schritt 9). Daraus folgt  $\frac{\bar{a}_{0s}}{-\bar{a}_{m+1,s}} \geq 1$ . Die Pivotoperation in Schritt 15 liefert nunmehr:

$$\begin{aligned} \bar{\bar{a}}_{00} &= \bar{a}_{00}^k - \frac{\bar{a}_{0s} \bar{a}_{m+1,0}}{\bar{a}_{m+1,s}} = \bar{a}_{00}^k + \frac{\bar{a}_{0s}}{-\bar{a}_{m+1,s}} \bar{a}_{m+1,0} \leq \bar{a}_{00}^k + \bar{a}_{m+1,0} \\ &= \bar{a}_{00}^k + \lfloor \bar{a}_{00}^k \rfloor - \bar{a}_{00}^k = \lfloor \bar{a}_{00}^k \rfloor. \end{aligned}$$

Also ergibt die erste Pivotoperation nach Hinzufügung der aus der Zielfunktion abgeleiteten Schnittebene:  $\bar{\bar{a}}_{00} \leq \lfloor \bar{a}_{00}^k \rfloor$ . Da im weiteren Verlauf bis zur nächsten Optimallösung  $\bar{\bar{a}}_{00}^{k+1}$  der Zielfunktionswert nicht größer werden kann, gilt nach Ende des  $(k + 1)$ -ten Durchlaufs des dualen Programms

$$\bar{\bar{a}}_{00}^{k+1} \leq \lfloor \bar{a}_{00}^k \rfloor.$$

Nähme der Zielfunktionswert  $\bar{\bar{a}}_{00}^k$  nach Beendigung des dualen Verfahrens unendlich oft einen nicht ganzzahligen Wert an, so wäre nach der geraden abgeleiteten Beziehung die Folge  $(\bar{\bar{a}}_{00}^k)$  nach unten unbeschränkt. Aufgrund des Zusatzes zu Schritt 8 würde dann aber nach endlich vielen Schritten ein Abbruch erfolgen. Also ist die Folge  $(\bar{\bar{a}}_{00}^k)$  nur endlich oft nicht ganzzahlig. Da sie nach unten beschränkt ist, nimmt sie nach endlich vielen Schritten einen ganzzahligen Wert an und behält diesen Wert von da an bei, d. h. die Folge  $(\bar{\bar{a}}_{00}^k)$  ist nach endlich vielen Gliedern konstant und zwar mit einem ganzzahligen Wert.

Wir wollen uns nun überlegen, dass  $\bar{\bar{a}}_{10}$  nach endlich vielen Schritten einen festen ganzzahligen Wert annimmt und diesen Wert von da an beibehält. Aufgrund der obigen Überlegungen können wir voraussetzen, dass  $\bar{\bar{a}}_{00}$  nach dem  $k_0$ -ten Durchlauf des dualen Simplexalgorithmus seinen endgültigen ganzzahligen Wert erreicht hat. In der Optimallösung des  $k_0$ -ten dualen Programms hat das Tableauelement  $\bar{a}_{10}$  einen nicht-negativen Wert (Schritt 11).

Wir zeigen nun, dass nach Beendigung des  $k_0$ -ten Durchlaufs bei jeder weiteren Pivotoperation der Wert von  $\bar{a}_{10}$  nicht-negativ bleibt, nicht größer wird und dass— falls die Schnittebene aus der ersten Zeile abgeleitet wird — nach der ersten anschließenden Pivotoperation  $\bar{a}_{10} \leq \lfloor \bar{a}_{10} \rfloor$  gilt.

Nehmen wir an, dass in den Schritten 12 und 14 ein Pivotelement  $\bar{a}_{rs} < 0$  ausgewählt wurde. Angenommen  $\bar{a}_{0s} > 0$ , dann folgt aus der Pivotformel 15.2.4 (wegen  $\bar{a}_{r0} < 0$ ,  $\bar{a}_{rs} < 0$ , dass der Zielfunktionswert  $\bar{a}_{00}$  abnimmt. Nach Voraussetzung ist aber  $\bar{a}_{00}$  konstant, also gilt  $\bar{a}_{0s} = 0$ . Folglich muss aufgrund der lexikographischen Positivität des Tableaus  $\bar{a}_{1s} \geq 0$  gelten. Daraus ergibt sich:  $r \neq 1$ . Also ist die erste Zeile nicht die Pivotzeile. Daraus folgt, dass nach dem  $k_0$ -ten Durchlauf des dualen Programms als Pivotspalten nur Spalten  $\bar{A}_j$  in Frage kommen mit  $\bar{a}_{0j} = 0$  und dass die erste Zeile  $\bar{A}_1$  niemals Pivotzeile ist. Wird bei einer Pivotoperation der Wert  $\bar{a}_{10}$  neu berechnet, so kommt wegen  $r \neq 1$  die Formel 15.2.4 in Frage also

$$\bar{a}_{10} = \bar{a}_{10} - \frac{\bar{a}_{1s}\bar{a}_{r0}}{\bar{a}_{rs}}.$$

Aufgrund der Regel zur Auswahl von Pivotspalte und -zeile gilt  $\bar{a}_{r0} < 0$ ,  $\bar{a}_{rs} < 0$ . Aus  $\bar{a}_{1s} \geq 0$  folgt, dass  $\bar{a}_{10} \leq \bar{a}_{10}$  gilt. Mithin wird das Tableauelement  $\bar{a}_{10}$  bei jeder Pivotoperation kleiner oder bleibt gleich. Würde  $\bar{a}_{10}$  irgendwann einmal negativ, so könnte niemals mehr primale Zulässigkeit (d. h.  $\bar{a}_{i0} \geq 0$   $i = 1, \dots, m$ ) erreicht werden, da — wie wir oben gezeigt haben — die erste Zeile niemals Pivotzeile wird. Weil aber das duale Simplexverfahren nach endlich vielen Schritten aufhört, muss  $\bar{a}_{10}$  immer nicht-negativ bleiben.

Daraus folgt, dass  $\bar{a}_{10}$  nach dem  $k_0$ -ten Durchlauf des dualen Programms bei jeder dualen Pivotoperation nicht größer wird, aber nicht-negativ bleibt. Also ist die Folge  $(\bar{a}_{10}^k)_{k \geq k_0}$  monoton fallend und nach unten beschränkt.

Ist  $\bar{a}_{10}$  in der Optimallösung eines dualen Programms nicht ganzzahlig, so wird aufgrund der Zusatzregel zu Schritt 9 die Schnittebene aus der ersten Zeile abgeleitet. Ersetzen wir in der obigen Diskussion der Zielfunktion den Zeilenindex 0 durch den Zeilenindex 1, so erhalten wir analog, dass nach dem ersten anschließenden Pivotschritt gilt:  $\bar{a}_{10} \leq \lfloor \bar{a}_{10} \rfloor$ . Wäre also  $\bar{a}_{10}$  unendlich oft nicht ganzzahlig, so wäre die Folge  $(\bar{a}_{10}^k)_{k \geq k_0}$  nach unten unbeschränkt. Da aber Null eine untere Schranke ist, muss  $(\bar{a}_{10}^k)_{k \geq k_0}$  nach endlich vielen Schritten, sagen wir  $k_1$ , einen ganzzahligen Wert annehmen und diesen Wert von da an beibehalten.

Um zu zeigen, dass auch  $\bar{a}_{20}$  nach endlich vielen weiteren Schritten ganzzahlig wird, können wir voraussetzen, dass  $\bar{a}_{00}$  und  $\bar{a}_{10}$  bereits ihren endgültigen ganzzahligen Wert angenommen haben. Dann verläuft der Beweis völlig analog wie bei  $\bar{a}_{10}$ . Auf diese Weise weisen wir nach, dass  $\bar{a}_{00}, \bar{a}_{10}, \dots, \bar{a}_{m0}$  nach endlich vielen Schritten einen ganzzahligen Wert annehmen. Damit ist der Beweis erledigt.  $\square$

Die Ganzzahligkeit aller Daten ist beim Gomory-Algorithmus unbedingt notwendig.

(17.9) **Beispiel.** Wir lösen das folgende IP mit dem ersten Gomory-Algorithmus:

$$\begin{aligned} \max \quad & 3x_1 + x_2 \\ & 17x_1 + 11x_2 \leq 86.5 \\ & x_1 + 2x_2 \leq 10.2 \\ & x_1 \leq 3.87 \\ & x_1, x_2 \geq 0 \text{ und ganzzahlig.} \end{aligned}$$

**LP-Lösung:**

		1	2				5	2
	0	- 3	- 1			11.61	3	- 1
3	86.50	17	11		3	20.71	-17	11
4	10.20	1	2	→	4	6.33	- 1	2
5	3.87	1	0		1	13.87	1	0

		5	3				5	6
	$\frac{14842}{1100}$	$\frac{16}{11}$	$\frac{1}{11}$			13	1	1
2	$-\frac{2071}{1100}$	$-\frac{17}{11}$	$\frac{1}{11}$		2	$\frac{1529}{1100}$	-2	1
4	$-\frac{2821}{1100}$	$\frac{23}{11}$	$-\frac{2}{11}$	→	4	$\frac{3905}{1100}$	3	-2
1	13.87	1	0		1	3.87	1	0
	$-\frac{542}{1100}$	$-\frac{5}{11}$	$-\frac{1}{11}$		3	5.42	5	-11
						$-f_{i0}$	0	0

Aus allen Zeilen des letzten Tableaus kann man einen Gomory-Schnitt ableiten, der die Form

$$0 \cdot x_5 + 0 \cdot x_6 \leq -f_{i0}$$

hat, mit  $f_{i0} > 0$ . Wäre dieser Schnitt gültig, hieße dies, dass unser Problem keine ganzzahlige Lösung hat, denn diese Ungleichung kann nicht erfüllt werden. Also liegt der Fehler in der Nichtganzzahligkeit der rechten Seite.  $\square$

Wir wissen ja bereits, dass das ganzzahlige Programmierungsproblem zur Klasse der  $\mathcal{NP}$ -vollständigen Probleme gehört. Man sollte also nicht erwarten, dass der Gomory-Algorithmus in polynomialer Zeit läuft. Selbst wenn wir statt des Simplexalgorithmus ein polynomiales Verfahren (Ellipsoidmethode, Karmarkar-Algorithmus) zur Lösung der jeweils auftretenden LP's wählen würden, wäre das gesamte Verfahren nicht polynomial, da die Zahl der Schnitte, die man betrachten muss, nicht durch ein Polynom beschränkt werden kann. Dies verdeutlicht das folgende ganzzahlige Programm.

**(17.10) Beispiel.** Für eine vorgegebene positive ganze Zahl  $k$  sei

$$\begin{aligned} \max \quad & -x_2 \\ & -kx_1 + x_2 \leq 0 \\ & kx_1 + x_2 \leq k \\ & x_1, x_2 \geq 0 \quad \text{und ganzzahlig,} \end{aligned}$$

ein ganzzahliges Programm (vergleiche (16.21) (b)). Dieses Programm hat nur zwei ganzzahlige Lösungen,  $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$  und  $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$ , die beide optimal sind.

Wendet man den Gomory-Algorithmus auf dieses IP an, so wird man feststellen, dass bis zur Auffindung einer Optimallösung  $\lceil \frac{3k}{2} \rceil - 1$  Gomory-Schnitte benötigt werden. Die Inputlänge des obigen LP kann durch  $8 \langle k \rangle$  abgeschätzt werden. Die Anzahl der Gomory-Schnitte kann also nicht durch ein Polynom in  $\log k$  beschränkt werden. Daraus folgt, dass die Anzahl der Iterationen des Gomory-Algorithmus nicht polynomial in den Ausgangsdaten ist.  $\square$

## 17.2 Ein Schnittebenenverfahren für gemischt-ganzzahlige Programme

Die vorhergehenden Überlegungen zur Lösung ganzzahliger Programme mit Hilfe von Schnittebenen-Verfahren lassen sich recht einfach auf den gemischt-ganzzahligen Fall übertragen. Hierzu müssen wir uns nur überlegen, wie wir eine geeignete Schnittebene ableiten können. Gegeben sei also ein gemischt-ganzzahliges Programm

$$\begin{aligned} \max \quad & c^T x + d^T y \\ & Ax + By = b \\ & x, y \geq 0, \quad x \text{ ganzzahlig.} \end{aligned}$$

Um Bezeichnungen einfacher zu machen, wollen wir dieses Programm in folgender Form schreiben:

$$\begin{aligned} \max \quad & c^T x \\ (MIP^=) \quad & Ax = b \\ & x_i \geq 0 \text{ und ganzzahlig, } i = 1, \dots, n_1 \leq n \end{aligned}$$

d. h. wir nehmen an, dass die Variablen  $x_1, \dots, x_{n_1}$  mit den niedrigen Indizes ganzzahlig sind und die übrigen rational. Wir setzen  $N_1 := \{1, \dots, n_1\}$ .

Wir gehen wieder davon aus, dass wir die LP-Relaxierung ( $LMIP^=$ ) von ( $MIP^=$ ) betrachten und dort an einer Basislösung angelangt sind, für die wir eine Schnittebene ableiten wollen, falls die Variablen  $x_i, i \in N_1$  nicht ganzzahlig sind.

**(17.11) Lemma.** Gegeben sei ein gemischt-ganzzahliges Programm

$$(MIP^=) \quad \max c^T x, Ax = b, x \geq 0, x_i \text{ ganzzahlig für } i \in N_1 = \{1, \dots, n_1\}$$

mit ganzzahligen Daten. Sei  $A_B$  eine Basis der LP-Relaxierung ( $LMIP^=$ ), dann gilt:

Ist  $x_{B(i)}$  eine Basisvariable mit  $B(i) \in N_1$  und ist  $\bar{b}_i = \bar{a}_{i0} \notin \mathbb{Z}$ , so ist

$$(*) \quad \sum_{j \in N_{\mathbb{Z}}} -f_{ij}x_j - \sum_{j \in N_{\mathbb{Q}}^+} \bar{a}_{ij}x_j + \sum_{j \in N_{\mathbb{Q}}^-} \frac{f_{i0}\bar{a}_{ij}}{1-f_{i0}}x_j \leq -f_{i0}$$

eine Schnittebene, die die gegenwärtige Basislösung abschneidet, falls  $f_{i0} = \bar{a}_{i0} - [\bar{a}_{i0}] = b_i - [\bar{b}_i] \neq 0$  gilt.

Hierbei sei  $B$  die Indexmenge der Basisvariablen,  $N$  die Indexmenge der Nichtbasisvariablen,  $N_{\mathbb{Z}}$  die Indexmenge der Nichtbasisvariablen, die der Ganzzahligkeitsforderung unterliegen, und  $N_{\mathbb{Q}} := N \setminus N_{\mathbb{Z}}$ . Wie üblich ist  $f_{ij} = \bar{a}_{ij} - [\bar{a}_{ij}]$  der gebrochene Teil, und es ist  $N_{\mathbb{Q}}^- := \{j \in N_{\mathbb{Q}} \mid \bar{a}_{ij} \leq 0\}$ ,  $N_{\mathbb{Q}}^+ := \{j \in N_{\mathbb{Q}} \mid \bar{a}_{ij} > 0\}$ .

**Beweis :** O. B. d. A. sei  $B = (1, 2, \dots, m)$ ,  $N = (m+1, \dots, n)$ . Wie üblich starten wir mit der Darstellung:

$$(1) \quad x_i = \bar{b}_i - \sum_{j \in N} \bar{a}_{ij}x_j = \bar{b}_i - \sum_{j \in N_{\mathbb{Q}}} \bar{a}_{ij}x_j - \sum_{j \in N_{\mathbb{Z}}} \bar{a}_{ij}x_j.$$

Daraus folgt

$$(2) \quad \sum_{j \in N_{\mathbb{Z}}} f_{ij}x_j + \sum_{j \in N_{\mathbb{Q}}} \bar{a}_{ij}x_j - f_{i0} = [\bar{a}_{i0}] - \sum_{j \in N_{\mathbb{Z}}} [\bar{a}_{ij}]x_j - x_i.$$

Die rechte Seite von (2) ist für alle zulässigen Lösungen von ( $MIP^=$ ) ganzzahlig, also gilt entweder

$$(3) \quad \sum_{j \in N_{\mathbb{Z}}} f_{ij}x_j + \sum_{j \in N_{\mathbb{Q}}} \bar{a}_{ij}x_j - f_{i0} \geq 0$$

oder

$$(4) \quad \sum_{j \in N_{\mathbb{Z}}} f_{ij}x_j + \sum_{j \in N_{\mathbb{Q}}} \bar{a}_{ij}x_j - f_{i0} \leq -1.$$

Falls (3) gilt, so gilt auch

$$(5) \quad \sum_{j \in N_{\mathbb{Z}}} f_{ij}x_j + \sum_{j \in N_{\mathbb{Q}}^+} \bar{a}_{ij}x_j \geq f_{i0}.$$

Falls (4) gilt, so gilt (da für alle  $j$  gilt:  $f_{ij} \geq 0$ )

$$(6) \quad \sum_{j \in N_{\mathbb{Q}}^-} \bar{a}_{ij}x_j \leq f_{i0} - 1.$$

Multiplizieren wir (6) mit  $\frac{f_{i0}}{f_{i0} - 1} (< 0)$ , so erhalten wir

$$(7) \quad -\sum_{j \in N_{\mathbb{Q}}^-} \frac{f_{i0} \bar{a}_{ij}}{1 - f_{i0}} x_j \geq f_{i0}.$$

Die linken Seiten von (7) und (5) sind beide nicht-negativ; da (5) oder (7) gelten muss, erhalten wir durch Aufaddieren der linken Seiten

$$(8) \quad \sum_{j \in N_{\mathbb{Z}}} f_{ij} x_j + \sum_{j \in N_{\mathbb{Q}}^+} \bar{a}_{ij} x_j - \sum_{j \in N_{\mathbb{Q}}^-} \frac{f_{i0} \bar{a}_{ij}}{1 - f_{i0}} x_j \geq f_{i0},$$

also ist die angegebene Ungleichung für alle zulässigen Lösungen gültig. Da für die gegenwärtige Basislösung  $x_N = 0$  gilt, erfüllt diese Basislösung die Ungleichung nicht.  $\square$

Durch eine genaue Analyse des obigen Beweises kann man den Schnitt (\*) noch etwas verbessern.

**(17.12) Lemma.** Die Voraussetzungen seien wie in Lemma (17.11). Falls es ein  $j \in N_{\mathbb{Z}}$  gibt, mit  $f_{ij} > f_{i0}$ , so kann die obige Ungleichung zu folgender Ungleichung verschärft werden:

$$(**) \quad \sum_{j \in N_{\mathbb{Z}}^+} -f_{ij} x_j - \sum_{j \in N_{\mathbb{Z}}^-} \frac{f_{i0}(1 - f_{ij})}{1 - f_{i0}} x_j - \sum_{j \in N_{\mathbb{Q}}^+} \bar{a}_{ij} x_j + \sum_{j \in N_{\mathbb{Q}}^-} \frac{f_{i0} \bar{a}_{ij}}{1 - f_{i0}} x_j \leq -f_{i0},$$

wobei  $N_{\mathbb{Z}}^+ = \{j \in N_{\mathbb{Z}} \mid f_{ij} \leq f_{i0}\}$ ,  $N_{\mathbb{Z}}^- = \{j \in N_{\mathbb{Z}} \mid f_{ij} > f_{i0}\}$ .

**Beweis:** Da nach Annahme  $f_{i0} > 0$ , gilt  $f_{ij} > 0 \forall j \in N_{\mathbb{Z}}^-$ . Subtrahieren wir  $\sum_{j \in N_{\mathbb{Z}}^-} x_j$  von beiden Seiten von (2), so erhalten wir:

$$(9) \quad \sum_{j \in N_{\mathbb{Z}}^+} + \sum_{j \in N_{\mathbb{Z}}^-} (f_{ij} - 1)x_j + \sum_{j \in N_{\mathbb{Q}}} \bar{a}_{ij} x_j - f_{i0} = \lfloor \bar{a}_{i0} \rfloor - \sum_{j \in N_{\mathbb{Z}}^+} \lfloor \bar{a}_{ij} \rfloor x_j - \sum_{j \in N_{\mathbb{Z}}^-} \lceil \bar{a}_{ij} \rceil x_j - x_i. \text{ Da die rechte Seite wiederum ganzzahlig ist, erhalten wir, dass entweder}$$

$$(10) \quad \sum_{j \in N_{\mathbb{Z}}^+} f_{ij} x_j + \sum_{j \in N_{\mathbb{Q}}^-} \bar{a}_{ij} x_j \geq f_{i0}$$

oder

$$(11) \quad -\frac{f_{i0}}{1 - f_{i0}} (\sum_{j \in N_{\mathbb{Z}}^-} (f_{ij} - 1)x_j + \sum_{j \in N_{\mathbb{Q}}^-} \bar{a}_{ij} x_j) \geq f_{i0} \text{ gilt.}$$

Also muss die Summe der beiden linken Seiten größer oder gleich  $f_{i0}$  sein. Diese Summe ist gerade die behauptete Ungleichung.

Die Koeffizienten von (\*) sind die gleichen wie die von (\*\*), außer im Falle  $j \in N_{\mathbb{Z}}^-$ . Aus  $f_{ij} > f_{i0}$  für  $j \in N_{\mathbb{Z}}^-$  folgt

$$f_{ij} - f_{i0} f_{ij} > f_{i0} - f_{i0} f_{ij}$$

und damit

$$f_{ij} > \frac{f_{i0}(1 - f_{ij})}{1 - f_{i0}}.$$

Daraus folgt, dass (\*\*) schärfer ist als (\*).  $\square$

**(17.13) Gomory's gemischt-ganzzahliger Algorithmus.****Input:**  $A \in \mathbb{Z}^{(m,n)}$ ,  $b \in \mathbb{Z}^m$ ,  $c \in \mathbb{Z}^n$ , Indexmenge  $N_1$ .**Output:** Lösung des gemischt-ganzzahligen Programms

$$(MIP^=) \quad \begin{array}{ll} \max & c^T x \\ & Ax = b \\ & x \geq 0 \\ & x_i \text{ ganzzahlig für } i \in N_1 \end{array}$$

Setze wie in (17.4):  $\bar{a}_{i0} := b_i$ ,  $i = 1, \dots, m$ ;  $\bar{a}_{0j} := -c_j$ ,  $j = 1, \dots, n$ ;  $\bar{a}_{00} :=$  Zielfunktionswert.

1. Löse das zugehörige lineare Programm ( $LMIP^=$ ) mit dem primalen Simplexalgorithmus.
2. Gilt  $\bar{a}_{i0} \in \mathbb{Z}$  für alle Basisvariablen  $x_{B(i)}$  mit  $B(i) \in N_1$ , so ist eine optimale gemischt-ganzzahlige Lösung von  $(MIP^=)$  gefunden.
3. Wähle ein  $i \in \{0, 1, \dots, m\}$  mit  $\bar{a}_{i0} \notin \mathbb{Z}$ , das zu einer ganzzahligen Variablen gehört, und setze

$$\begin{array}{ll} f_{ij} & := \bar{a}_{ij} - \lfloor \bar{a}_{ij} \rfloor \\ f_{00} & := \bar{a}_{00} - \lfloor \bar{a}_{00} \rfloor \end{array} \quad \text{für alle } j \in N \text{ (Nichtbasisvariable) mit } j \in N_1$$

$$\bar{a}_{m+1,j} := \begin{cases} -f_{ij}, & \text{falls } j \in N_1 \text{ und } f_{ij} \leq f_{i0} \\ -\frac{f_{i0}(1-f_{ij})}{1-f_{i0}}, & \text{falls } j \in N_1 \text{ und } f_{ij} > f_{i0} \\ -\bar{a}_{ij}, & \text{falls } j \in \{1, \dots, n\} \setminus N_1 \text{ und } \bar{a}_{ij} > 0 \\ \frac{f_{i0}\bar{a}_{ij}}{1-f_{i0}}, & \text{falls } j \in \{1, \dots, n\} \setminus N_1 \text{ und } \bar{a}_{ij} \leq 0 \end{cases}$$

Füge die Zeile

$$\bar{a}_{m+1,0}, \bar{a}_{m+1,1}, \dots, \bar{a}_{m+1,n-m}$$

als  $(m+1)$ -te Zeile zum gegenwärtigen Tableau hinzu. Setze  $m := m+1$ ,  $n := n+1$  und erweitere die Basis um die Schlupfvariable  $n+1$ .

4. Löse das erweiterte Programm mit dem dualen Simplexverfahren und gehe anschließend zu 2. □

Für das oben angegebene Verfahren für gemischt-ganzzahlige Programme kann man nun analog zu Satz (17.8) zeigen:

**(17.14) Satz.** Falls die Zusatzregeln (17.7) zum ersten Gomory-Verfahren analog beim gemischt-ganzzahligen Gomory-Verfahren angewendet werden und unter der zusätzlichen Voraussetzung, dass der Wert der Zielfunktion in der Optimallösung ganzzahlig sein muss, dann liefert der gemischt-ganzzahlige Algorithmus von Gomory nach endlich vielen Schritten eine Optimallösung oder zeigt nach endlich vielen Schritten an, dass entweder ( $MIP=$ ) unbeschränkt oder  $MIP=(A, B, b)$  leer ist.

**Beweis :** völlig analog zu (17.7). □

Die beiden Schnittebenentypen, die wir bisher kennengelernt haben, sind bei weitem nicht alle, die in Schnittebenenverfahren eingesetzt werden. Die sechziger Jahre brachten eine Fülle interessanter Studien zu Schnittebenenverfahren. Die neu erfundenen Schnittebenen erhielten schöne Namen wie

Intersection-Cuts,

Diamond-Cuts,

Polaroid-Cuts,

Cuts from the Hypercube,

Outer-Polar-Cuts

(siehe hierzu Garfinkel & Nemhauser (1972)), jedoch änderten die neuen Techniken nichts an der bis Mitte der 90er Jahre empirisch beobachteten Tatsache, dass Schnittebenenverfahren dieser allgemeinen Art relativ ineffiziente Methoden zur Lösung ganzzahliger und gemischt-ganzzahliger Programme sind. Neuimplementierungen von Schnittebenenverfahren in kommerziellen MIP-Codes Ende der 90er Jahre haben zu einer Änderung der „herrschenden Meinung“ geführt. Die (praktisch effiziente) Einbindung von Schnittebenen in Branch & Bound-Codes hat zu einer enormen Beschleunigung dieser Verfahren geführt, siehe Bixby-Vortrag am 20.01.2010.

Die beiden Verfahren, die wir hier behandelt haben, werden auch zu den dualen Verfahren gezählt und zwar, weil in jedem Schritt eine dual zulässige Lösung behalten wird und das Ziel ist, eine primal zulässige und somit dann eine optimale Lösung zu finden.

Da unsere Tableaus auch gebrochene Werte enthalten und aus den gebrochenen rechten Seiten und den zugehörigen Tableau-Zeilen Schnitte abgeleitet werden, nennt man die Verfahren (17.4) bzw. (17.13) in der Literatur auch

**Dual Fractional Integer Programming Algorithm** bzw.  
**Dual Fractional Mixed Integer Programming Algorithm**

Es gibt zu diesem Ansatz einige Varianten, und zwar kann man versuchen, mit einer primal zulässigen Lösung zu starten und immer dann, wenn im nächsten Simplexschritt eine nicht-ganzzahlige Lösung generiert wird, einen Schritt hinzuzufügen und auf diese Weise immer primale Zulässigkeit zu erhalten. Am Ende folgt dann aus der dualen Zulässigkeit die Optimalität. Verfahren dieses Typs nennt man **primale Schnittebenenverfahren**.

Ein großes Problem, das bei den "Fractional-Verfahren" auftritt, sind die Rundfehler, die häufig zum Abbruch ohne Optimalität führen. Um diese unangenehmen Erscheinungen zu vermeiden, hat Gomory ein Verfahren entwickelt, das mit rein ganzzahliger Arithmetik arbeitet und deshalb **All-Integer Verfahren** genannt wird. Hierzu sind ebenfalls Varianten angegeben worden, so dass wir **Dual-All-Integer** und **Primal-All-Integer Integer Programming Verfahren** kennen. Aus Zeitgründen können wir auf diese Varianten nicht eingehen. In der Praxis werden sie nicht verwendet.