

Algorithmen und Berechnungskomplexität I

Skript WS 2017

Nach Aufzeichnungen von Elmar Langetepe, Heiko Röglin,
und Rolf Klein

Bonn, Dezember 2017

Inhaltsverzeichnis

1	Einführung	1
1.1	Algorithmen	1
1.1.1	Anwendungsgebiete	1
1.1.2	Maschinenmodell	3
1.1.3	Algorithmische Paradigmen	4
1.1.4	Analysemethode	4
1.2	Ein einfaches Beispiel: Insertionsort	6
2	Divide-and-Conquer	9
2.1	Sortieren von Zahlenketten, Mergesort	9
2.1.1	Laufzeitanalyse	10
2.1.2	Korrektheit	12
2.1.3	Untere Schranke von Sortieralgorithmen	12
2.2	Closest Pair von n Punkten	13
3	Lösen von Rekursionsgleichungen	17
3.1	Substitutionsmethode	17
3.1.1	Raten durch Iteration	18
3.1.2	Raten durch Rekursionsbaum	19
3.2	Typische Probleme	20
3.3	Das Mastertheorem	21
4	Dynamische Programmierung	25
4.1	Fibonacci Zahlen	25
4.2	Matrixmultiplikation	28
4.3	Rucksackproblem	30
5	Greedy Algorithmen	33
5.1	Rucksackproblem	33
5.2	Das optimale Bepacken von Behältern	35
5.3	Aktivitätenauswahl	38
5.4	Firefighting in Bäumen	39

5.4.1	Parametereffiziente Lösung	41
6	Das Sweepverfahren	45
6.0.2	Closest Pair von Punkten in der Ebene	46
7	Datenstrukturen	49
7.1	Stacks	50
7.2	Dynamische Listen	52
7.3	Bäume und Suchbäume	53
7.3.1	Binärbäume und Suchbäume	54
7.4	AVL-Bäume	59
7.4.1	Einfügen eines neuen Knoten	61
7.4.2	Entfernen eines Knoten	67
7.5	B-Bäume	70
7.6	Mittelwertbildung	72
7.6.1	Amortisierte Kosten: Berechnung der konvexen Hülle	72
7.6.2	Randomisierter Input: Bucketsort	73
7.6.3	Randomisierter Algorithmus: Bestimmung des Maximums	74
7.6.4	Randomisierter Algorithmus: Quicksort	75
7.7	Heaps	76
7.7.1	Heapsort	78
7.7.2	Die Prioritätenwarteschlange	79
7.8	Hashing	80
7.8.1	Kollisionsbehandlung mit verketteten Listen	81
7.9	Union-Find Datenstruktur	84
7.10	Datenstrukturen für Graphen	86
7.10.1	Adjazenzlisten	87
7.10.2	Nachbarschaftsmatrix	88
8	Elementare Graphalgorithmen	91
8.1	Minimale Spannbäume	91
8.1.1	Algorithmus von Kruskal	91
8.2	Kürzeste Wege	95
8.2.1	Single-Source Shortest Path Problem	97
8.2.2	All-Pairs Shortest Path Problem	103
8.3	Flussprobleme	105
8.3.1	Algorithmus von Ford und Fulkerson	108
8.3.2	Algorithmus von Edmonds und Karp	115
8.3.3	Der Heiratssatz	117

9	Lineare Programmierung	119
9.1	Grundlagen	121
9.1.1	Kanonische Form und Gleichungsform	121
9.1.2	Geometrische Interpretation	121
9.2	Anwendung des Simplex Algorithmus	125
9.2.1	Allgemeines Verfahren: Dictionary	127
9.2.2	Besonderheiten des Simplex-Algorithmus	129
9.3	Beispiel Anwendung: Integer LP	134
9.4	Algebraische Interpretation mit Lösungspolyeder	136
9.4.1	Anschauliches Beispiel der Basislösung	139
9.5	Simplex-Algorithmus in der Matrix-Variante	140
9.5.1	Formale Beschreibung	141
9.5.2	Berechnung der initialen Basislösung	143
9.5.3	Laufzeit	144
9.6	Komplexität von linearer Programmierung	144
	Literatur	145

Kapitel 1

Einführung

In dieser Vorlesung wird die Komplexität von Berechnungsproblemen untersucht und wir werden die dafür notwendigen theoretischen Konzepte entwickeln. Für viele Problemstellungen werden Algorithmen (Lösungspläne) vorgestellt und analysiert.

1.1 Algorithmen

Formal gesehen ist ein *Algorithmus* ein Lösungsplan, der für eine *Probleminstanz* eines definierten Berechnungsproblems die entsprechende *Ausgabe* liefert.

Beispiel: **Sortierproblem**

Eingabe: Eine Folge von n reellen Zahlen (a_1, a_2, \dots, a_n) .

Ausgabe: Eine Permutation π , so dass $(a_{\pi(1)}, a_{\pi(2)}, \dots, a_{\pi(n)})$ eine sortierte Folge der Eingabe ergibt, d.h., $a_{\pi(1)} \leq a_{\pi(2)} \leq \dots \leq a_{\pi(n)}$.

Konkrete Eingabefolgen werden *Instanzen* genannt.

Bezüglich der Algorithmen und der Problemstellung werden wir uns insgesamt zunächst mit drei Begriffen auseinandersetzen:

- **Korrektheit:** Wir wollen für die Algorithmen formal beweisen, dass zu jeder Instanz das richtige Ergebnis berechnet wird.
- **Laufzeitabschätzung:** Wie ist die Laufzeit des Algorithmus in Bezug auf die Eingabegröße n und wie effizient ist der Algorithmus im Vergleich zu anderen möglichen Algorithmen?
- **Speicherplatzabschätzung:** Wieviel Speicher in Bezug auf die Eingabegröße n benötigt der Algorithmus und wie effizient ist das im Vergleich zu anderen möglichen Algorithmen?

Die hier betrachteten Algorithmen sollen mithilfe eines Computers ihre Aufgabe erfüllen, deshalb werden wir uns an gängigen Programmiersprachen orientieren und dazu *Pseudocode* verwenden.

1.1.1 Anwendungsgebiete

Algorithmen spielen in verschiedenen Anwendungsgebieten eine große Rolle. Die Problemstellungen werden durch mathematische Modellierung so formuliert, dass sie durch einen Rechner

bearbeitet werden können.

- Human Genome Project: Datenanalyse, DNA-Vergleiche, Matchingalgorithmen, ...
- Internet: Routing, Datenhaltung, Suchmaschinen, Routenplaner, Verschlüsselung, ...
- Industrie: Beladen von Containern, Optimieren von Arbeitsabläufen, ...
- Infrastruktur: Optimieren von Netzen, Umwege minimieren, Kürzeste Wege berechnen/approximieren, ...
- ...

Ein Beispiel aus unserer aktuellen algorithmischen Forschung in Zusammenarbeit mit dem FKIE Wachtberg.

Wir betrachten ein Gebiet, in dem es einen Industrie-Unfall gegeben hat und das von einer Menge von Robotern beobachtet werden soll. Das Gebiet ist durch den Menschen nicht mehr betretbar. Die Roboter haben einen beschränkten Senderadius und sollen den Kontakt untereinander und zur Basisstation halten. Von der Basisstation aus sollen die Roboter zu vorher festgelegten Positionen gefahren werden. Die Positionen und die Wege der einzelnen Positionen zueinander konnten berechnet werden. Damit die Roboter den Kontakt untereinander nicht verlieren, ist es manchmal erforderlich, mehrere Roboter von Position A zu Position B zu bewegen. Für eine Bewegung von A nach B sind somit zum Beispiel mindestens 10 Roboter notwendig.

Die Frage lautet, wieviel Roboter brauchen wir, um von einer Basisstation alle gegebenen Positionen zu besetzen und bei der Bewegung zu gewährleisten, dass die Agenten alle in Kontakt bleiben?

Modellbildung durch einen Graphen $G = (V, E)$ mit ganzzahligen Knotengewichten w_v für $v \in V$ und ganzzahligen Kantengewichten w_e für $w \in E$, ein Startknoten v_s ; siehe Abbildung 1.1.

Folgende Bedingungen:

- Eine Kante e darf nur mit $k \geq w_e$ Agenten traversiert werden.
- Ein Knoten v darf nur mit $k \geq w_v$ Agenten besucht werden.
- Beim ersten Besuch eines Knotens v bleiben w_v Agenten am Knoten zurück.

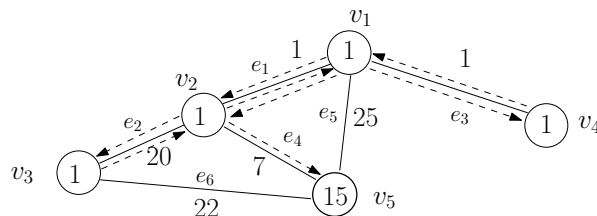


Abbildung 1.1: Falls die Agenten am Knoten v_1 starten, kann eine Tour mit einer minimalen Anzahl an Agenten wie folgt beschrieben werden. 23 Agenten in einer einzelnen Gruppe besuchen Knoten in der Reihenfolge v_1, v_2, v_3, v_4 and v_5 . Die Kanten e_5 and e_6 werden nicht besucht.

Frage: Wieviele Agenten werden benötigt, um alle Knoten v mit w_v Agenten zu füllen? Welchen Weg können die Agenten dabei zurücklegen? Lässt sich ein effizienter Algorithmus formulieren, der dieses Problem für beliebige Instanzen löst?

Die Abbildung zeigt ein Beispiel: Vom festgelegten Basisknoten v_1 aus, kommt man mit 23 Agenten aus. Beginnend am Knoten v_1 bleibt zunächst ein Agent zurück und 22 Agenten bewegen sich über die Kante e_1 nach v_2 . Hier wird ein Agent abgelegt und mit 21 Agenten geht es zum Knoten v_3 über die Kante e_2 . Nachdem wir dort einen Agenten abgelegt haben, können wir mit den verbleibenden 20 Agenten gerade noch über die Kante zurück. Hier stellt sich quasi heraus, dass es mit weniger als 23 gar nicht gehen kann!

Mit den 20 Agenten besuchen wir v_4 (über e_1 und e_3) und lassen einen Agenten dort. Dann geht es mit den verbleibenden über e_3 , e_1 und e_4 zu v_5 . Hier werden 15 Agenten benötigt und wir sind fertig.

Es stellt sich heraus, dass dieses Problem für allgemeine Graphen wirklich *schwer* zu lösen ist. Für einfache Graphen (Bäume) gibt es effiziente Algorithmen, durch die die optimale Lösung in best-möglicher Laufzeit (auch durch die Verwendung effizienter Datenstrukturen (Heaps)) ermittelt werden kann. Daraus resultiert auch eine Approximationslösung für allgemeine Graphen. Aus dem Graphen wird ein Baum mit guter Verbindungseigenschaft generiert und darauf der Baumalgorithmus angewendet.

Ein formaler Beweis zeigt jeweils die Güte der Approximation, die Optimalität des Baumalgorithmus und die Schwere des Problems im Allgemeinen. Dieses Beispiel enthält somit viele Aspekte, die in der Vorlesung vermittelt werden sollen.

1.1.2 Maschinenmodell

Die Laufzeitabschätzungen der Algorithmen soll unabhängig von der jeweils aktuellen Hardware-Rechnerausstattung (z.B. Taktfrequenz, Verarbeitungsbreite etc.) sein, deshalb zählen wir die Anzahl der ausgeführten *Elementaroperationen* in unserem jeweiligen Maschinenmodell. Der Einfachheit halber veranschlagen wir für jede Operation in etwa die gleichen Kosten.

Wir verwenden das Modell einer sogenannten **REAL RAM**, eine *Random Access Maschine*, die mit reellen Zahlen rechnet. In diesem Modell werden Anweisungen nacheinander (sequentiell) ausgeführt. Die REAL RAM Modell orientiert sich an realen Rechnern und hat folgende Spezifikation:

- Abzählbar unendlich viele Speicherzellen, die mit den natürlichen Zahlen adressiert werden.
- Jede Speicherzelle kann eine beliebige reelle oder natürliche Zahl enthalten. Direkter oder indirekter Zugriff ist möglich.
- Elementare Rechenoperationen: Addieren, Subtrahieren, Multiplizieren, Dividieren, Restbilden, Abrunden, Aufrunden.
- Elementare Relationen: kleiner gleich, größer gleich, gleich, \vee , \wedge .
- Kontrollierende Befehle: Verzweigung, Aufruf von Unterroutinen, Rückgabe.
- Datenbewegende Befehle: Laden, Speichern, Kopieren.

Alle angegebenen Operationen können in konstanter Zeit ausgeführt werden.

Wir nehmen an, dass Integer-Zahlen der Größe m mit $c \log m$ Bits für eine Konstante $c \geq 1$ dargestellt werden können. Reelle Zahlen werden gemäß des IEEE Standards im Rechner

durch Binärzahlen approximiert. Innerhalb des Darstellungsbereiches können arithmetische Operationen mit reellen Zahlen im Rechner sehr effizient durchgeführt werden.

1.1.3 Algorithmische Paradigmen

Algorithmen lassen sich häufig bezüglich ihrer Herangehensweise an das jeweilige Problem kategorisieren. Zunächst seien hier einige davon kurz erwähnt. Es geht hier im wesentlichen darum, wie das Problem gelöst wird.

- Brute-Force (Naive Methode)
- Inkrementelle Konstruktion (Schrittweise Vergrößerung der Eingabe)
- Divide-and-Conquer (Aufteilen und Zusammenfügen)
- Greedy (*gierig*, schnelle Verbesserungen)
- Dynamische Programmierung (Tabellarische Auflistung und Verwertung von Teillösungen)
- Sweep (geometrische Probleme, Dimensionsreduktion)
- ...

Die Methoden können auch kombiniert auftreten, zum Beispiel ein Divide-and-Conquer der im Merge-Schritt einen Sweep verwendet. Auf die einzelnen Paradigmen werden wir später gezielt für konkrete Problemstellungen eingehen.

1.1.4 Analysemethode

Sei P eine Probleminstanz eines Problems Π und sei A ein Algorithmus mit RAM-Anweisungen, der für jedes Problem $P \in \Pi$ eine Lösung liefert.

Die Anzahl der Elementaroperationen eines Algorithmus A bei einer Eingabegröße $|P| = n \in \mathbb{N}$ wird durch eine *Kostenfunktion* $T_A : \mathbb{N} \mapsto \mathbb{R}^+$ beschrieben. Die Instanz P könnte beispielsweise eine Menge von n Zahlen sein.

Die exakte Bestimmung dieser Funktion kann vernachlässigt werden, insbesondere, da auf unterschiedlichen Rechnern verschiedene Konstanten für die Elementaroperationen beachtet werden müssten. Also sind wir bei der Analyse nur an der Größenordnung der Funktion T_A in Abhängigkeit der Eingabegröße $|P| = n$ interessiert. Ebenso verhält es sich mit der Eingabegröße selbst, ob wir es mit n Punkten in der Ebene zu tun haben oder mit $2n$ Koordinaten, soll beispielsweise keine Rolle bei der Beschreibung der Größenordnung spielen.

Die Anzahl der Elementaroperationen eines Algorithmus mit Eingabegröße $n \in \mathbb{N}$ kann auch von der Zusammenstellung der Eingabe selbst abhängen. Deshalb können wir verschiedene Kennzahlen verwenden:

- Worst-Case: Maximal mögliche Anzahl an Elementaroperationen
- Average-Case: Mittlere Anzahl der Elementaroperationen, gemittelt über alle möglichen Eingabefolgen der entsprechenden Größenordnung

Die Analyse kleiner Eingabegrößen ist generell wenig aussagekräftig, da wir in diesem Fall die Summe aller Elementaroperationen unter einer großen Konstante subsumieren können. Deshalb untersuchen wir das *asymptotische Verhalten* der Kostenfunktion.

Das führt zu folgenden Notationen für Klassen von Funktionen:

Definition 1 (O -, Ω -, Θ -Notation)

$$O(f) = \{g \mid \text{es existiert ein } n_0 \geq 0 \text{ und ein } C > 0, \\ \text{so dass } g(n) \leq C f(n) \text{ für alle } n \geq n_0\}$$

$$g \in \Omega(f) \quad :\Leftrightarrow \quad f \in O(g) \\ :\Leftrightarrow \quad \text{es existiert ein } n_0 \geq 0 \text{ und ein } C > 0, \\ \text{so dass } f(n) \leq C g(n) \text{ für alle } n \geq n_0.$$

$$g \in \Theta(f) \quad :\Leftrightarrow \quad g \in O(f) \text{ und } g \in \Omega(f)$$

Wir erlauben also mit $n \geq n_0$ endliche viele Ausnahmestellen.

Die obige Notation wird in dem Sinne gebraucht, dass $f \in O(g)$ eine *obere Schranke* an die Funktion f durch g liefert, währenddessen $f \in \Omega(g)$, eine *untere Schranke* für f durch g darstellt.

Beispiele:

$3n+2 \in O(n)$, da für alle $n \geq 2$ und für $C = 4$ gilt: $3n+2 \leq 4n$. Desgleichen gilt $3n+2 \in \Omega(n)$, da für alle $n \geq 0$ bereits $n \leq 3n+2$ gilt. Somit gilt $3n+2 \in \Theta(n)$.

Für $g(n) = 2n^3 - 18n^2 - \sin(n) + 16n + 3$ führt zunächst die grobe Abschätzung

$$g(n) \leq 2n^3 + 16n + 3 \\ \leq (2 + 16 + 3)n^3 \text{ für alle } n \geq 1$$

zur Aussage $g \in O(n^3)$. Es gilt aber auch

$$g(n) \geq 2n^3 - 18n^2 \\ = n^3 + (n - 18)n^2 \\ \geq n^3 \text{ für alle } n \geq 15$$

und deshalb gilt hier ebenfalls $g \in \Omega(n^3)$ und $g \in \Theta(n^3)$. Wir sagen auch, dass die Funktion g in der *Größenordnung* n^3 wächst.

Eine weitere hilfreiche Notationkonvention ist, dass wir die obigen Symbole der Einfachheit halber bei einer asymptotischen Betrachtung auch in Formeln verwenden möchten. So drückt $2n^2 + 3n + 1 = 2n^2 + \Theta(n)$ aus, dass wir uns um die Funktion $3n+1$ in der Größenordnung $\Theta(n)$ nicht genau kümmern möchten. Diese Konvention kann hilfreich sein, wenn beispielsweise eine Laufzeitabschätzung *rekursiv* durch eine Formel $T(n) = 2T(\frac{n}{2}) + \Theta(n)$ beschrieben werden darf. Wahlweise kann hier auch $T(n) = 2T(\frac{n}{2}) + C \cdot n$ für ein nicht näher bestimmtes C verwendet werden.

Außerdem kann es sinnvoll sein, die Analyse der Laufzeiten *ausgabesensitiv* vorzunehmen, um die Ausgabekomplexität nicht der algorithmischen Komplexität anzulasten. Falls wir beispielsweise die Menge aller Schnittpunkte von n Liniensegmenten suchen, so ist klar, dass es

im worst-case $\Omega(n^2)$ viele Schnittpunkte geben kann, jeder Algorithmus also mindestens diese Laufzeit im worst-case benötigt. Eine ausgabesensitive Analyse könnte zu einem Ergebnis von $O((n+k)\log n)$ führen, wobei k die Anzahl der tatsächlichen Schnittpunkte angibt.

1.2 Ein einfaches Beispiel: Insertionsort

Die obigen Begriffe wollen wir zunächst auf das einfache Beispiel des Sortierens anwenden. Wir nehmen an, die n zu sortierenden Zahlen sind in einem *Array* A gespeichert.

Eingabe: n reelle Zahlen $A[1], A[2], \dots, A[n]$.

Ausgabe: Ein Array B der Zahlen aus A mit $B[1] \leq B[2] \leq \dots \leq B[n]$.

Beachte, dass wir hier eine etwas andere Ausgabe als zu Beginn erwarten (Permutation π). Die Aufgabenstellung ist aber äquivalent.

Eine *Brute-Force* Vorgehensweise wäre, das kleinste Element des Arrays zu ermitteln, dieses zu entfernen und mit dem restlichen Array fortzufahren.

Wir betrachten einen *inkrementellen* Ansatz. Sukzessive wird ein immer größerer Teil des Arrays A in B sortiert. Anders ausgedrückt, falls wir die ersten $j-1$ Elemente des Arrays A in B bereits sortiert haben, nehmen wir uns danach das j -te Element vor und sortieren es in B ein.

Algorithm 1 InsertionSort(A)

Array A (nichtleer) sortiert in Array B ausgeben

```

1:  $B[1] := A[1]$ ;
2: for  $j = 2$  to  $\text{length}(A)$  do
3:    $key := A[j]$ ;  $i := j - 1$ ;
4:   while  $i > 0$  und  $B[i] > key$  do
5:      $B[i + 1] := B[i]$ ;  $i := i - 1$ ;
6:   end while
7:    $B[i + 1] := key$ ;
8: end for
9: RETURN  $B$ 

```

Die Korrektheit des Algorithmus wird durch eine *Schleifeninvariante* gezeigt. Wir zeigen formal, dass für jeden Beginn der FOR-Schleife die folgende Invariante gilt: *Für den aktuellen Wert j gilt, dass B bis zum Index $j-1$ eine sortierte Folge der ersten $j-1$ Einträge von A ist.*

Formal gesehen, handelt es sich hier um eine Induktion, die bei einem bestimmten Parameter terminiert.

Induktionsanfang: $j = 2$. Die Eigenschaft ist erfüllt, durch die erste Anweisung.

Induktionsschritt: Wir nehmen an, dass die Aussage bereits für $j-1 \geq 2$ gilt. Jetzt werden in der WHILE-Schleife im Array B solange die Werte um einen Index nach hinten geschoben, bis das Element $key := A[j]$ kleiner gleich dem i -ten Element von B ist. Dann wird key als $(i+1)$ -tes Element in B eingetragen. Falls das bis $i = 0$ nicht eintritt, wird key als 1-tes Element von B eingetragen. Insgesamt ist B danach bis zum Index j ein sortiertes Array der ersten j -Elemente von A .

Terminierung: Falls $j = \text{length}(A) + 1$ wird, gilt zunächst die Invariante, dass für den aktuellen Wert j gilt, dass B bis zum Index $j - 1 = \text{length}(A)$ eine sortierte Folge der ersten $\text{length}(A)$ Einträge von A ist. Die FOR-Schleife bricht ab und das Ergebnis wird korrekt in B ausgegeben.

Streng genommen muss die WHILE-Schleife genauso analysiert werden. Welche Schleifeninvariante muss hier gewählt werden?

Für eine Laufzeitabschätzung analysieren wir die einzelnen Programmschritte. Dabei sei t_j die Anzahl der Durchläufe der WHILE-Schleife wenn die FOR-Schleife für j aufgerufen wird. Sei $\text{length}(A) = n$.

Es ergibt sich die folgende Tabelle:

Anweisungen in Algorithmus 1	Kosten	Häufigkeit
1:	c_0	1
3:	c_1	$n - 1$
5:	c_3	$\sum_{j=2}^n t_j$
7:	c_4	$n - 1$

Insgesamt ergibt sich nach unseren Konventionen eine Laufzeit $T(n) = c_3 \sum_{j=2}^n t_j + O(n)$.

Die Laufzeit hängt also von der Größe von t_j ab. Im schlechtesten Fall ist das Array A absteigend sortiert, alle Elemente aus B müssen verschoben werden und das j -te Element wird vorne eingefügt. In diesem Fall gilt $t_j = j$ und wir haben.

$$T(n) \leq c_3 \sum_{j=1}^n j + O(n) = c_3 \frac{n(n+1)}{2} + O(n) \in O(n^2)$$

Die worst-case Analyse ist im allgemeinen sinnvoll, weil ...

- ... wir eine Laufzeit für jede beliebige Eingabe garantieren.
- ... der worst-case in vielen Anwendungsfeldern relativ häufig vorkommen kann.
- ... die Analyse der mittleren Laufzeit nicht unbedingt bessere Ergebnisse erzielt.

Beispielsweise können wir bei der Analyse einer mittleren Laufzeit die Frage stellen, an welcher Stelle im Mittel das j -te Element eingeordnet werden muss. Im Mittel ist die Hälfte der Elemente größer als $A[j]$ und die Hälfte der Elemente kleiner als $A[j]$ und somit könnten wir $t_j = \frac{j}{2}$ betrachten. Das führt asymptotisch gesehen exakt zur gleichen Laufzeit.

Anders kann es aussehen, wenn wir aus der Eingabe durch *Würfeln* zufällig gleichverteilt die Elemente aus A ziehen.

Solche *probabilistische Analysen* werden wir später für sogenannte *randomisierte* Algorithmen durchführen. Im Gegensatz dazu haben wir hier zunächst einen *deterministischen* Algorithmus betrachtet. Der Lösungsplan ist von vorneherein vollständig festgelegt und nicht vom Zufall abhängig.

Kapitel 2

Divide-and-Conquer

In diesem Kapitel behandeln wird die Entwurfsmethode *Teile-und-Herrsche* und eine zugehörige allgemeine Analyse­methode. Gleichzeitig lernen wir dadurch das Prinzip der *Rekursion* kennen, das heißt wir rufen bestimmte Routinen unseres Algorithmus immer wieder auf, um Teilprobleme zu lösen. Am Ende werden die Lösungen der Teilprobleme zur Lösung des Gesamtproblems zusammengefügt bzw. zusammengemischt.

Das allgemeine Schema der Methode für ein Problem P der Größe n funktioniert wie folgt:

- | | |
|------------|--|
| 1. Divide | $\begin{cases} n > c & : & \text{Teile das Problem in } k \geq 2 \text{ Teilprobleme} \\ & & \text{der Größen } n_1, n_2, \dots, n_k, \text{ gehe zu 2.} \\ n \leq c & : & \text{Löse das Problem direkt} \end{cases}$ |
| 2. Conquer | Löse die k Teilprobleme auf dieselbe Art (rekursiv) |
| 3. Merge | Füge die k berechneten Teillösungen zu einer Gesamtlösung zusammen |

In diesem Kapitel werden wir zunächst das Beispiel des Sortierens von Zahlenketten betrachten und später ein allgemeines nützliches Theorem zur Lösung von Rekursionsgleichungen vorstellen.

2.1 Sortieren von Zahlenketten, Mergesort

Hier betrachten wir $k = 2$ und $n_1, n_2 \approx \frac{n}{2}$. Für eine konkrete Zahlenfolge ergibt sich das folgende Aufrufschema eines Divide-And-Conquer Algorithmus:

Eingabe	: (3, 2, 1, 7, 3, 4, 9, 2)
Divide I.	: (3, 2, 1, 7)(3, 4, 9, 2)
(Conquer)Divide II.	: (3, 2)(1, 7)(3, 4), (9, 2)
(Conquer)Divide III.	: (3)(2)(1)(7)(3)(4)(9)(2)
Merge III.	: (2,3)(1,7)(3,4)(2,9)
Merge II.	: (1,2,3,7)(2,3,4,9)
Merge I.	: (1,2,2,3,3,4,7,9)

Im Pseudocode und unter Verwendung von Arrays könnte eine Implementierung wie folgt aussehen. Wir verwenden ein nichtleeres Array A von Index 1 bis $n = \text{length}(A)$ und rufen zur vollständigen Sortierung $\text{MergeSort}(A, 1, n)$ auf.

Procedure 2 MergeSort(A, p, r)Array A wird zwischen Index p und r sortiert

```

1: if  $p < r$  then
2:    $q := \lfloor (p + r) / 2 \rfloor$ ;
3:   MergeSort( $A, p, q$ );
4:   MergeSort( $A, q + 1, r$ );
5:   Merge( $A, p, q, r$ );
6: end if

```

Die eigentliche Arbeit (bis auf die Anzahl der rekursiven Aufrufe) wird im Merge-Schritt vollzogen. Dafür kopieren wir die bereits von Index p bis q und von Index $q + 1$ bis r sortierten Abschnitte in zwei Arrays L und R und fügen diese dann wieder gemeinsam sortiert in A von Index p bis r ein. Für ein Abbruchkriterium des Mischens fügen wir einen Dummywert ∞ am Ende von L und R ein.

Danach werden die beiden Teilarrays gemischt, indem in einem gleichzeitigen Durchlauf jeweils das kleinste momentane *Kopfelement* von L und R in A eingefügt wird.

Procedure 3 Merge(A, p, q, r)Array A vorsortiert zwischen Index p und q und Index $q + 1$ und r wird zwischen p und r sortiert

```

1:  $n_1 := q - p + 1$ ;  $n_2 := r - q$ ;
2: for  $i = 1$  to  $n_1$  do
3:    $L[i] := A[p + i - 1]$ ;
4: end for
5: for  $j = 1$  to  $n_2$  do
6:    $R[j] := A[q + j]$ ;
7: end for // erzeuge Arrays  $L[1 \dots (n_1 + 1)]$  und  $R[1 \dots (n_2 + 1)]$ 
8:  $L[n_1 + 1] := \infty$ ;  $R[n_2 + 1] := \infty$ ;
9:  $i := 1$ ;  $j := 1$ ;
10: for  $k = p$  to  $r$  do
11:   if  $L[i] \leq R[j]$  then
12:      $A[k] := L[i]$ ;  $i := i + 1$ ;
13:   else
14:      $A[k] := R[j]$ ;  $j := j + 1$ ;
15:   end if
16: end for

```

2.1.1 Laufzeitanalyse

Für einen einzelnen Merge-Schritt entsteht offensichtlich ein Aufwand von $O(|L| + |R|)$; siehe Prozedur 3. In jedem Schritt der finalen FOR-Schleife wird ein Element von L oder R nicht mehr betrachtet. Deshalb kann die Schleife nicht mehr als $(|L| + |R|)$ Mal aufgerufen werden.

Zur Laufzeitanalyse verwenden wir eine Rekursionsgleichung. Eine klassische Vereinfachung ist dabei, dass wir annehmen, dass n eine Zeierpotenz ist, also $n = 2^k$. Wir können die Kostenfunktion dann für den obigen Algorithmus wie folgt beschreiben.

$$T(n) = \begin{cases} c & \text{falls } n = 1 \\ 2T\left(\frac{n}{2}\right) + cn & \text{falls } n > 1 \end{cases} \quad (2.1)$$

In der obigen Gleichung haben wir eine einzelne Konstante c verwendet. Das ist prinzipiell erlaubt, da wir für die Konstanten für das Problem der Größe 1 und den Merge-Schritt einfach die größte der beiden Konstanten wählen können. Die asymptotische Laufzeit ändert sich nicht. Beachten Sie, dass in den Kosten cn auch das Aufteilen der Aufgabenstellung in zwei Teilprobleme subsumiert ist. Hier musste lediglich ein Index in der Mitte in konstanter Zeit berechnet werden.

Die Rekursionsgleichung läßt sich nun leicht insbesondere auch wegen der Annahme $n = 2^k$ durch einen Rekursionsbaum abschätzen, siehe Abbildung 2.1. Die allerunterste Ebene des Baumes enthält exakt n Aufrufe für die Arrays der Größe 1. In jeder darüberliegenden Ebene werden jeweils zwei Mengen der darunterliegenden Ebene im Merge-Schritt vereinigt. Die Kosten sind jeweils *linear* zur Summe der Größen der darunterliegenden Arrays. Somit sind die Kosten auf jeder Ebene linear zur Gesamtzahl der Elemente überhaupt. Der Baum hat $k = \log_2 n$ viele Ebenen, da sich die Anzahl der Knoten von unten nach oben jeweils halbiert. Wir können also $T(n) = O(n \log n)$ und sogar $T(n) = \Theta(n \log n)$ folgern. Streng genommen

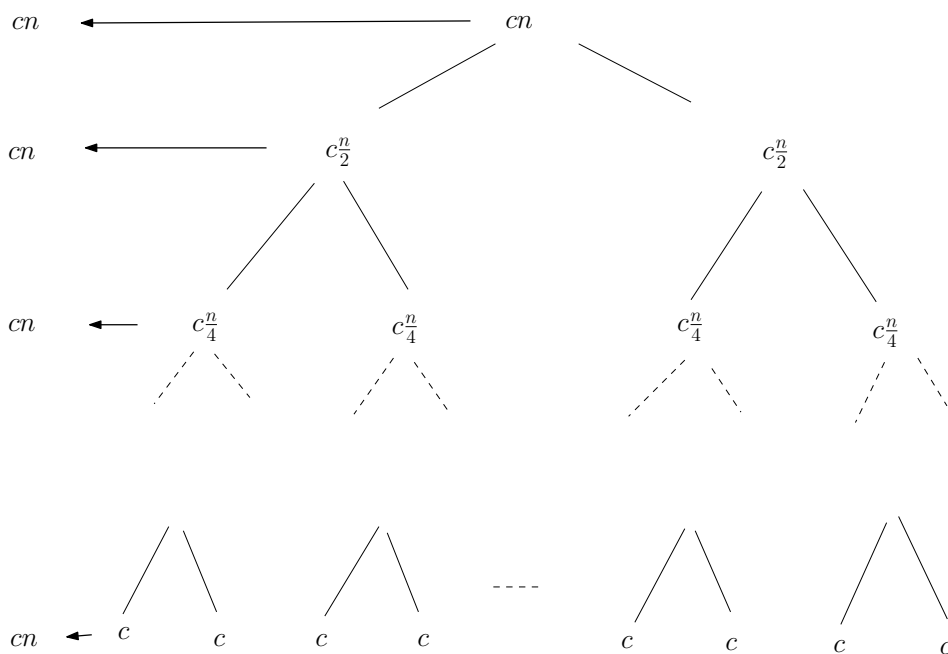


Abbildung 2.1: Auf jeder Ebene des Rekursionsbaumes fällt ein Aufwand von cn für die einzelnen Merge-Schritte an. Nach $k = \log_2 n$ Ebenen sind nur noch einzelne Elemente vorhanden.

hätten wir die Rekursionsgleichung (2.1) wie folgt formulieren müssen.

$$T(n) = \begin{cases} c & \text{falls } n = 1 \\ T(\lfloor \frac{n}{2} \rfloor) + T(\lceil \frac{n}{2} \rceil) + cn & \text{falls } n > 1 \end{cases} \quad (2.2)$$

In jedem Fall müssen wir uns beim Aufstellen und Lösen der Rekursionsformeln Gedanken darüber machen, ob die Einschränkung $n = 2^k$ nicht das asymptotischen Verhalten der Laufzeit beeinflusst. Im obigen Fall ist das offensichtlich nicht der Fall, der zugehörige Rekursionsbaum wäre zwar nicht so ausgeglichen, wie in Abbildung 2.1. Für jedes n könnten wir aber beispielsweise die nächste Zweierpotenz oberhalb von n nehmen, um eine obere Schranke zu erzielen, sowie die nächste Zweierpotenz unterhalb von n für eine untere Schranke der Laufzeit. Der Algorithmus hat dann entsprechend mehr oder entsprechend weniger Aufrufe zu absolvieren. Darüber müssen wir uns aber Gedanken machen.

Beispielsweise kann eine Rekursionsformel der Form

$$T(n) = \begin{cases} n & \text{falls } n \text{ gerade} \\ n^2 & \text{falls } n \text{ ungerade} \end{cases} \quad (2.3)$$

nicht mit dem Ansatz $n = 2^k$ betrachtet werden.

2.1.2 Korrektheit

Die Korrektheit des obigen Divide-and-Conquer Algorithmus ergibt sich direkt aus der Korrektheit der Merge Prozedur. Sobald diese gemäß ihrer Spezifikation funktioniert, erhalten wir die richtige Zwischenergebnisse rekursiv für weitere Aufrufe.

Zunächst ist klar, dass die Merge-Prozedur durch die ersten beiden FOR-Schleifen die Einträge des Arrays A von Index p bis Index q und von Index $q + 1$ bis r in die Kopien L und R der Größen $n_1 := q - p + 1$ $n_2 := r - q$ einträgt. Streng genommen müssen auch diese Behauptungen durch eine Schleifeninvariante bewiesen werden.

Wir beschränken uns hier auf die dritte FOR-Schleife und formulieren eine geeignete Schleifeninvariante.

Zu Beginn jeder Iteration der FOR-Schleife enthält das Array $A[p \dots k - 1]$ die $(k - p)$ kleinsten Einträge aus L und R in sortierter Reihenfolge. $L[i]$ und $L[j]$ sind die kleinsten Einträge der jeweiligen Arrays, die noch nicht in A zurückkopiert wurden.

Induktionsanfang: $k = p$. Das Array $A[p \dots k - 1]$ enthält die $k - p = 0$ kleinsten Elemente von L und R . $L[1]$ und $R[1]$ sind die kleinsten Elemente der Arrays, die noch nicht zurückkopiert wurden.

Induktionsschritt: Wir nehmen an, die Aussage gelte bereits für $k = p + l$ und wollen zeigen, dass die Aussage nach dem Durchlauf der FOR-Schleife für $k' = p + l + 1$ erhalten bleibt. Die Schleife wird mit $k = p + l$ ausgeführt. Es werden die aktuellen Köpfe der beiden Arrays L und R verglichen und das kleinste der beiden Elemente als $k = (p + l)$ -tes Element in A eingetragen. Also besteht danach $A[p \dots k' - 1]$ für $k' = p + l + 1$ aus den $(k' - p)$ kleinsten Einträge aus L und R in sortierter Reihenfolge. Je nachdem, ob $L[i] \leq R[j]$ oder $L[i] > R[j]$ gilt, wird i oder j um eins erhöht und die Aussage gilt auch für $k' = p + l + 1$ und das neue j oder i .

Terminierung: Beim Abbruch gilt $k = r + 1$ und $A[p \dots r]$ enthält die $r + 1 - p$ kleinsten Elemente von L und R in sortierter Reihenfolge. Zusammen enthalten L und R insgesamt $n_1 + n_2 + 2 = r - p + 3$ Elemente. Also muss sowohl i am Index $n_1 + 1$ und j am Index $n_2 + 1$ angekommen sein. Alle Werte außer den Dummywerten ∞ sind in A eingetragen worden.

2.1.3 Untere Schranke von Sortieralgorithmen

Mergesort kann also n Zahlen (a_1, a_2, \dots, a_n) in Zeit $O(n \log n)$ sortieren. Dabei greift der Algorithmus auf die Elemente der Eingabefolge nur durch *paarweise Vergleiche* zu (siehe Zeile 11 der Merge-Prozedur): Es wird getestet, ob zum Beispiel $a_i < a_j$ gilt, und der Ausgang dieses Tests entscheidet, was der Algorithmus als nächstes macht. "Gerechnet" wird mit den a_i aber nicht.

Mit solchen vergleichsbasierten Sortierverfahren kann man deswegen nicht nur Zahlen sortieren, sondern Objekte aus beliebigen vollständig geordneten Mengen, vorausgesetzt, es steht ein Größenvergleichstest zur Verfügung.

Interessanterweise ist Mergesort sogar optimal.

Theorem 2 *Jeder vergleichsbasierte Sortieralgorithmus benötigt im worst-case $\Omega(n \log n)$ viele Vergleiche, um n Objekte zu sortieren.*

Beweis. Sei A ein vergleichsbasierter Sortieralgorithmus. Wir nehmen der Einfachheit halber an, dass A deterministisch ist; das Theorem gilt aber auch für randomisierte Verfahren. Außerdem setzen wir voraus, dass die a_i paarweise verschieden sind. Weil A den Input (a_1, a_2, \dots, a_n) zunächst nicht kennt, wird stets derselbe Test $a_i < a_j$? als erster ausgeführt. In Abhängigkeit vom Ergebnis verzweigt der Algorithmus. Für alle Eingaben, bei denen $a_i < a_j$ gilt, wird jetzt stets derselbe Test $a_v < a_w$? als nächster ausgeführt. Ebenso kommt für alle Inputs mit $a_i > a_j$ als nächstes stets derselbe Test $a_r < a_s$? an die Reihe, und so fort. So entsteht der Entscheidungsbaum B_A von Algorithmus A , ein Binärbaum, der mindestens so viele Blätter enthält, wie es Permutationen gibt; denn ein korrekter Sortieralgorithmus darf für unterschiedlich permutierte Eingabefolgen nicht zum selben Ergebnis kommen. Also gilt

$$\begin{aligned} \text{Höhe}(B_A) &\geq \log_2(n!) \geq \log_2 \left(\left(\frac{n}{2} \right)^{\frac{n}{2}} \right) = \frac{n}{2} \log_2 \left(\frac{n}{2} \right) \\ &\geq \frac{1}{3} n \log_2(n) \end{aligned}$$

für alle $n \geq 8$. Dabei gilt die erste Ungleichung, weil es in $n!$ mindestens $n/2$ viele Faktoren gibt, die mindestens so groß sind wie $n/2$. Es gibt also im Baum B_A mindestens einen Pfad der Länge $\Omega(n \log n)$, und wenn die Eingabe so permutiert ist, wie es den Testergebnissen längs dieses Pfades entspricht, muss Algorithmus A alle diese Vergleiche ausführen. \square

Wenn die zu sortierenden Objekte aber Zahlen sind, mit denen man “rechnen” darf, gilt die untere Schranke von Theorem 2 nicht. Zum Beispiel geht der Algorithmus Radixsort folgendermaßen vor, um n k -stellige Zahlen zu sortieren: Zunächst werden leere 10 Listen L_0, L_1, \dots, L_9 angelegt, je eine für jede Ziffer. Nun werden abwechselnd k Verteilungs- und Sammelphasen ausgeführt. Man beginnt mit der letzten der k Stellen, durchläuft die Eingabereihenfolge und hängt die Zahl a_i an diejenige Liste an, die zur letzten Ziffer von a_i gehört. Anschließend werden die Zahlen aus den Listen wieder eingesammelt, in der Reihenfolge L_0, L_1, \dots, L_9 und unter Beibehalt der Reihenfolgen innerhalb der Listen. Mit der neuen Zahlenfolge verfährt man ebenso, verwendet aber nun die vorletzte Stelle, und so fort. Wenn man schließlich die Zahlen nach der ersten Stelle verteilt und wieder eingesammelt hat, sind sie sortiert (!).

Offensichtlich führt Radixsort nur $O(kn)$ viele Schritte aus. Wenn die Stellenzahl k als konstant betrachtet wird (was bei Standard-Zahlentypen meist der Fall ist), so ergibt sich die Laufzeit $O(n)$.

2.2 Closest Pair von n Punkten

Wir betrachten die folgende geometrische Aufgabenstellung. Für einer Menge S von n Punkten p_1, p_2, \dots, p_n in der Ebene, soll das Punktepaar mit kleinstem Euklidischen Abstand $d(\cdot, \cdot)$ berechnet werden.

Wir suchen also den Wert

$$d_S = \min_{1 \leq i, j \leq n, i \neq j} d(p_i, p_j).$$

Der Brute-Force Ansatz für dieses Problem vergleicht systematisch alle $O(n^2)$ vielen Distanzen und benötigt eine Laufzeit von $\Theta(n^2)$.

Für einen Divide-and-Conquer Algorithmus sortieren wir die Punkte zunächst nach X -Koordinaten und teilen die Punktmenge rekursiv in zwei gleichgroße Mengen anhand der X -Koordinaten auf.

1. Divide Teile das Problem in Teilmengen S_r und S_l gleicher Größe anhand der X -Koordinaten auf
2. Conquer Bestimme rekursiv d_{S_l} und d_{S_r}
3. Merge Berechne d_S durch die Verwendung von d_{S_l} und d_{S_r} .

Wir wollen im Merge-Schritt wiederum möglichst wenig Aufwand für die Bestimmung von d_S verwenden. Sei $d := \min\{d_{S_l}, d_{S_r}\}$. Wir brauchen für den Merge nur Punkte betrachten die von S_l zu Punkten aus S_r einen Abstand $< d$ haben. Nur die Punkte innerhalb des Streifens der Breite $2d$ symmetrisch entlang der *Splitvertikalen* können einen solchen Abstand $< d$ zueinander haben; siehe Abbildung 2.2. Im Merge-Schritt wird folgendes Verfahren angewendet:

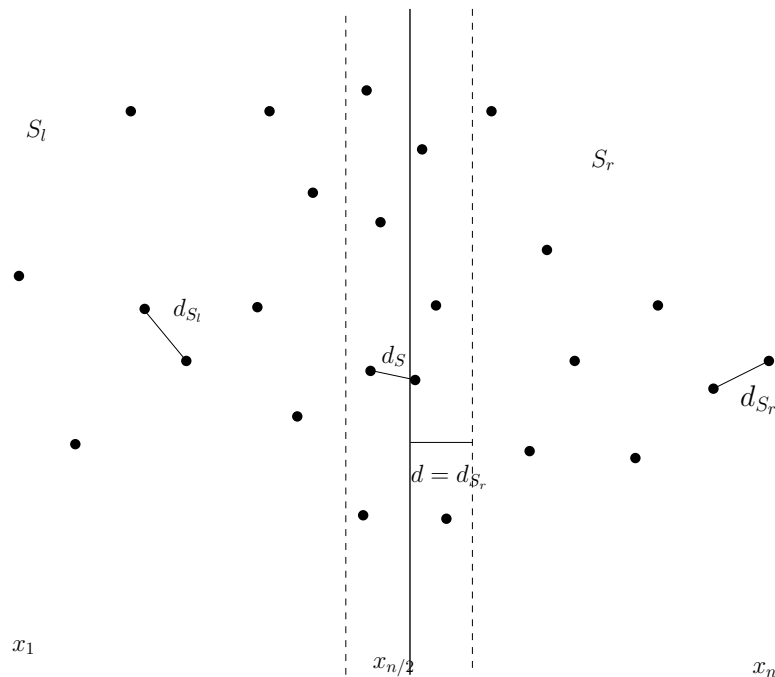


Abbildung 2.2: Für den Merge-Schritt müssen nur die Punkte der beiden Teilstreifen der Breite d entlang der Splitvertikalen betrachtet werden.

- Bestimme die Punkte in den zugehörigen rechten und linken Teilstreifen der Breite d und sortiere diese jeweils nach Y -Koordinaten
- Gehe die Punkte im linken Teilstreifen nach fallenden Y -Koordinaten durch
- Für jeden solchen Punkt p betrachte alle Punkte aus dem rechten Teilstreifen deren Y -Koordinaten im Bereich $[p_y - d, p_y + d]$ liegen
- Der zugehörige Punktbereich auf der rechten Seite wandert sukzessive nach unten

- Behauptung: Zu jedem Zeitpunkt sind nicht mehr als 10 Punkte im zugehörigen Bereich

Beweis. (Behauptung) Nehmen wir an, dass p und q ein dichtestes Punktepaar mit Abstand $< d$ sind. Dann können wir oBdA annehmen, dass p im linken und q im rechten Teilstreifen liegt. Der Punkt q kann nicht außerhalb des Rechtecks $R(p)$ der Breite d und Höhe $2d$ im rechten Teilstreifen liegen; siehe Abbildung 2.3. Außerdem haben alle Punkte in $R(p)$ einen Abstand von mindestens d zueinander. Dann sind die Kreisumgebungen $U_{\frac{d}{2}}(q)$ in $R(p)$ disjunkt. Für jeden Punkt q in $R(p)$ ist mindestens ein Viertel dieser Kreisfläche in $R(p)$ enthalten. $R(p)$ kann dann höchstens

$$\frac{2d^2}{\frac{\pi}{4} \left(\frac{d}{2}\right)^2} = \frac{32}{\pi} < 11$$

viele Punkte enthalten. □

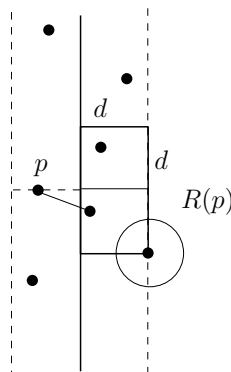


Abbildung 2.3: Für jeden Punkt p des linken Teilstreifens müssen nur konstant viele Punkte im rechten Teilstreifen getestet werden. Die zu testenden Punkte haben Y -Koordinaten, die im Bereich $[p_y - d, p_y + d]$ liegen.

Zur Analyse der Laufzeit betrachten wir nun die Einzelschritte des beschriebenen Algorithmus.

- In einem ersten Schritt werden die Punkte nach X -Koordinaten sortiert. Das kostet $O(n \log n)$ Zeit.
- Der Divide-Schritt kostet dann stets lineare Zeit.
- Für den Merge Schritt ergibt sich folgender Aufwand:
 1. Es müssen die Punkte mit X -Koordinaten aus den Teilstreifen ermittelt werden, das kostet $O(n)$ Zeit.
 2. Danach werden die Punkte in den jeweiligen Bereichen nach Y -Koordinate sortiert. Das kostet $O(n \log n)$ Zeit.
 3. In einem linearen Durchlauf wird für jeden Punkt p zu einer konstanten Anzahl (< 10) von Punkten q der Abstand ermittelt und der aktuell kleinste gespeichert. Insgesamt liegt der Aufwand in $O(n)$ für diesen Schritt.

Insgesamt erstellen wir daraus die folgende Rekursionsgleichung.

$$T(n) = \begin{cases} c & \text{falls } n = 1 \\ 2T\left(\frac{n}{2}\right) + c \cdot n \log n & \text{falls } n > 1 \end{cases} \quad (2.4)$$

Wie wir solche Rekursionsgleichungen im Allgemeinen lösen ist Gegenstand des nächsten Kapitels. Eine Übungsaufgabe besteht dann darin, die obige Gleichung abzuschätzen. Welche Vermutung haben Sie?

Kapitel 3

Lösen von Rekursionsgleichungen

Für die Analyse von Algorithmen, die das Prinzip der Rekursion verwenden, kann häufig eine Rekursionsgleichung aufgestellt werden.

Zunächst erwähnen wir einige Konventionen. Obwohl unsere Kostenfunktionen auf den natürlichen Zahlen definiert sind ignorieren wir diesen Sachverhalt. Deshalb können wir bei der Rekursionsformel für den Mergesort auch sogleich die Formel (2.1) verwenden und verzichten vollständig auf die Annahme $n = 2^k$.

Weiterhin liegen die Kosten für $T(n)$ für kleine n im Allgemeinen in $\Theta(1)$ und somit lassen wir diesen Teil der Rekursionsgleichung manchmal einfach weg. Die Bedingungen für kleine n werden auch Randbedingungen genannt. Die Randbedingungen spielen beim Induktionsbeweis für die Lösung einer Rekursionsgleichung eine Rolle. Insbesondere der Induktionsanfang ist davon betroffen.

Beispiele für Rekursionsgleichungen sind somit:

$$T(n) = T\left(\left\lceil \frac{n-1}{2} \right\rceil\right) + c \quad \text{Binäre Suche}$$

$$T(n) = 2T\left(\frac{n}{2}\right) + cn \quad \text{Mergesort}$$

$$T(n) = 7T\left(\frac{n}{2}\right) + cn \quad \text{Strassen's Matrixmultiplikation}$$

Die Frage lautet, wie wir diese Rekursionsgleichungen im Allgemeinen lösen können?

3.1 Substitutionsmethode

Die Substitutionsmethode besteht im Prinzip aus zwei Teilen.

1. Zunächst wird eine Lösung der Rekursionsgleichung *erraten*. Beispielsweise durch Erfahrung, durch Einsetzen einiger Werte oder durch die Analyse eines Rekursionsbaumes.
2. Dann wird die Lösung per *Induktion* verifiziert. Dabei können die verwendeten Konstanten präzisiert werden.

Beispiel:

$$T(n) = \begin{cases} 1 & \text{falls } n = 1 \\ 2T\left(\frac{n}{2}\right) + n & \text{falls } n > 1 \end{cases}$$

Wir vermuten, dass $T(n) \in O(n \log n)$ gilt und wollen $T(n) \leq c \cdot n \log n$ für geeignetes $c > 0$ per Induktion zeigen. Zunächst ist $T(1) = 1 > c \cdot 1 \log 1 = 0$. Deshalb muss hier $n_0 > 1$ gewählt werden. $T(2) = 4 \leq c \cdot 2 \log 2$ und $T(3) = 5 \leq c \cdot 3 \log 3$ für $c > 1$. Jetzt überprüfen wir die Aussage für $n_0 = 2$ per Induktion mit der *Substitutionsmethode*.

Induktionsanfang: $T(2) = 4 \leq c \cdot 2 \log 2$.

Induktionsschritt: Die Aussage gilt für alle Indizes von n_0 bis $n-1$. Zeige die Aussage auch für n .

$$\begin{aligned} T(n) &= 2T\left(\frac{n}{2}\right) + n \\ &\leq 2c \cdot \frac{n}{2} \log \frac{n}{2} + n && \text{(Induktionsannahme)} \\ &= c \cdot n \log n - c \cdot n \log 2 + n && \text{(Rechenregeln log)} \\ &\leq c \cdot n \log n + (1-c)n \\ &\leq c \cdot n \log n \text{ (für } c > 1) \end{aligned}$$

Als nächstes stellen wir zwei Methoden vor, wie die zu verifizierende Lösung erraten werden kann.

3.1.1 Raten durch Iteration

Wir betrachten die Rekursionsgleichung

$$T(n) = 3T\left(\frac{n}{4}\right) + n. \quad (3.1)$$

Für das Raten einer Lösung setzen wir die Gleichung sukzessive ein.

$$\begin{aligned} T(n) &= 3T\left(\frac{n}{4}\right) + n \\ &= 3\left(3T\left(\frac{n}{16}\right) + \frac{n}{4}\right) + n = 9T\left(\frac{n}{16}\right) + \left(\frac{3}{4}\right)^1 \cdot n + n && \text{(Einsetzen)} \\ &= 9\left(3T\left(\frac{n}{64}\right) + \frac{n}{16}\right) + 3\frac{n}{4} + n && \text{(Einsetzen)} \\ &= 27T\left(\frac{n}{64}\right) + \left(\frac{3}{4}\right)^2 \cdot n + \left(\frac{3}{4}\right)^1 \cdot n + \left(\frac{3}{4}\right)^0 \cdot n && \text{(Umformen)} \end{aligned}$$

Wenn wir $T(1) = c$ und $n = 4^k$, annehmen gelangen wir zur Vermutung

$$\begin{aligned} T(n) &= n \sum_{i=0}^{k-1} \left(\frac{3}{4}\right)^i + c \cdot 3^k \\ &= n \sum_{i=0}^{\log_4 n - 1} \left(\frac{3}{4}\right)^i + c \cdot n^{\log_4 3}. \end{aligned}$$

Hierbei wurde $3^{\log_4 n} = n^{\log_4 3}$ verwendet. Daraus wollen wir eine genaue Vermutung ableiten und diese danach beweisen.

$$\begin{aligned}
T(n) &= n \sum_{i=0}^{\log_4 n - 1} \left(\frac{3}{4}\right)^i + c \cdot n^{\log_4 3} \\
&< n \sum_{i=0}^{\infty} \left(\frac{3}{4}\right)^i + c \cdot n^{\log_4 3} && \text{(nach oben Abschätzen)} \\
&= n \cdot C + c \cdot n^{\log_4 3} && \text{(Potenzreihe } \sum_{i=1}^{\infty} X^i \text{ Konvergenzradius } |X| \leq 1) \\
&< C' \cdot n && (\log_4 3 \leq 1)
\end{aligned}$$

Also vermuten wir $T(n) \in O(n)$ und verifizieren die Aussage $T(n) \leq C' \cdot n$ durch Induktion mit Substitutionsmethode:

Induktionsanfang: $T(1) = c < C' \cdot 1$ gilt für $C' > c$.

Induktionsschritt: Die Aussage gilt für alle Indizes von n_0 bis $n-1$. Zeige die Aussage auch für n .

$$\begin{aligned}
T(n) &= 3T\left(\frac{n}{4}\right) + n \\
&\leq 3C' \cdot \frac{n}{4} + n && \text{(Induktionsannahme)} \\
&= C' \cdot \frac{3}{4} \cdot n + n && \text{(Umformung)} \\
&= \left(1 + C' \cdot \frac{3}{4}\right) n && \text{(Umformung)} \\
&\leq C' \cdot n && \text{(für } C' \geq 4).
\end{aligned}$$

3.1.2 Raten durch Rekursionsbaum

Die Idee bei der Verwendung eines Rekursionsbaumes ist, das obige iterative Einsetzen durch einen Baum darzustellen. Dabei werden die nicht-rekursiven Kosten und das Rekursionsargument notiert. Wir gehen wie folgt vor:

1. Jeder Knoten der Baumes stellt die Kosten des Teilproblems dar. Die Wurzel stellt die Kosten $T(n)$ dar und die Blätter die Kosten $T(1)$.
2. Wir summieren die Kosten auf jeder *Ebene* des Baumes.
3. Die Gesamtkosten ergeben sich durch die Summe der Kosten aller Ebenen.

Als Beispiel betrachten wir wiederum die Rekursionsgleichung (3.1) und nehmen $T(1) = c$ an. Damit der Baum *passend aufgeht* wählen wir $n = 4^k$.

Die Analyse über den Rekursionsbaum erfolgt durch

$$T(n) = \underbrace{\sum_{i=0}^{\log_4 n - 1}}_{\text{alle Ebenen}} \underbrace{\left(\frac{3}{4}\right)^i \cdot n}_{\text{Kosten pro Ebene}} + \underbrace{c \cdot n^{\log_4 3}}_{\text{Kosten der Blätter}} \quad (3.2)$$

Beachte: In der Praxis werden gelegentlich die Ratemethoden (insbesondere das iterative Einsetzen) bereits als *formaler Beweis* für eine Laufzeitanalyse verwendet. Auf die Induktion wird also tatsächlich gelegentlich verzichtet. Formal gesehen kann man das insbesondere unter

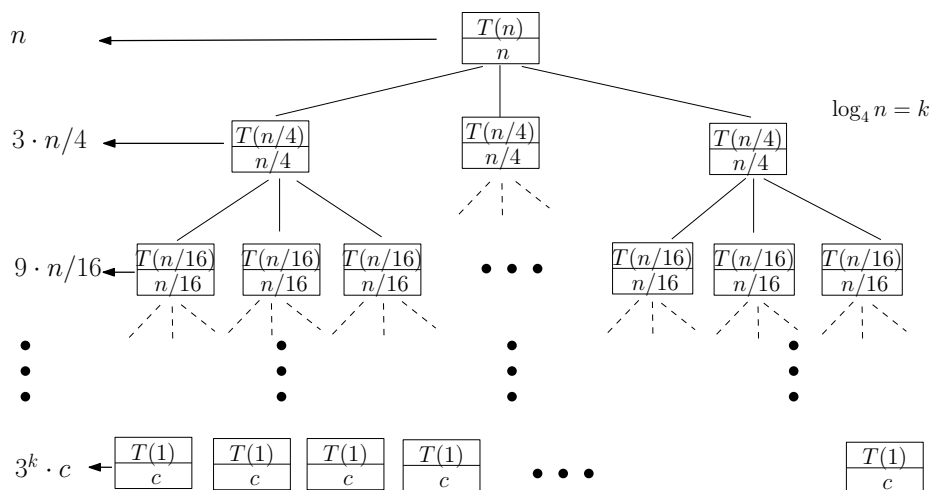


Abbildung 3.1: Der Rekursionsbaum zur Gleichung (3.1) hat $\log_4 n$ Ebenen. In jeder Ebene fallen konkrete Kosten (für das Aufteilen und das Kombinieren) an.

der Annahme, dass n eine bestimmte Potenz ist (zum Beispiel $n = 2^k$), gelten lassen. Das haben wir beispielsweise auch bei der Analyse des Mergesorts so gehandhabt. Die Aussage gilt dann aber streng genommen zunächst nur für exakt solche Potenzen.

3.2 Typische Probleme

Bei der Lösung von Rekursionsgleichungen können u.a. typischerweise die folgenden Probleme auftreten:

1. Der Induktionsbeweis ist aufgrund der Konstanten fehlerhaft.
2. Die Lösungsform ist zu ungenau.
3. Die Lösungsform erfordert eine Variablentransformation.

Wir geben hier Beispiele für die obigen Probleme an.

Korrekte Konstanten: Für

$$T(n) = 2T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + n$$

könnten wir

$$\begin{aligned} T(n) &\leq 2\left(c\left\lfloor \frac{n}{2} \right\rfloor\right) + n \\ &\leq c \cdot n + n \\ &= O(n) \quad \text{Falsch!!} \end{aligned}$$

folgern, da c eine Konstante ist. Der Fehler liegt darin, dass wir nicht exakt $T(n) \leq cn$ bewiesen haben. Wir müssen die Konstante genau angeben.

Lösungsform ist zu ungenau: Für

$$T(n) = \begin{cases} 1 & \text{falls } n = 1 \\ 2T\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + T\left(\left\lceil \frac{n}{2} \right\rceil\right) + 1 & \text{falls } n > 1 \end{cases}$$

vermuten wir $T(n) \leq cn$. Die Induktionsannahme ist für $c > 1$ erfüllt. Leider führt

$$\begin{aligned} T(n) &\leq c \left\lfloor \frac{n}{2} \right\rfloor + c \left\lceil \frac{n}{2} \right\rceil + 1 \\ &= cn + 1 > cn \end{aligned}$$

im Induktionsschluss durch Substitution nicht zum Ziel. Wir müssen hier $T(n) \leq cn - 1$ wählen und dann für die Induktionsannahme $c > 2$.

Lösung durch Variablentransformation:

Für

$$T(n) = \begin{cases} 1 & \text{falls } n = 1 \\ 2T(\lfloor \sqrt{n} \rfloor) + \log n & \text{falls } n > 1 \end{cases}$$

können wir eine geeignete Umformung vornehmen. Durch $m = \log n$ gelangen wir vereinfacht zu

$T(2^m) = T(2^{m/2}) + m$ und erhalten durch Umbenennung die Gleichung $S(m) = 2S(m/2) + m$. Da, wie bereits gezeigt $S(m) = O(m \log m)$ gilt erhalten wir

$$T(n) = T(2^m) = S(m) \in O(m \log m) \in O(\log n \log \log n).$$

3.3 Das Mastertheorem

Wir wollen nun ein allgemeines Theorem für die Lösung solcher Rekursionsgleichungen erstellen. In der Regel sieht das Problem folgendermaßen aus:

- Das Problem teilt sich in $a \geq 1$ Teilprobleme auf.
- Jedes der Teilprobleme hat eine Größe $\frac{n}{b}$ für $b > 1$.
- Die Kosten für das Aufteilen und das Kombinieren der Teillösungen sind durch $f(n)$ gegeben.

Die Kostenanalyse wird dann durch die Gleichung

$$T(n) = aT\left(\frac{n}{b}\right) + f(n) \tag{3.3}$$

beschrieben.

Theorem 3 Seien $a \geq 1$ und $b > 1$ Konstanten und $f(n)$ eine nichtnegative Funktion. Dann gilt für die Rekursionsgleichung (3.3):

1. Wenn $f(n) \in O(n^{\log_b a - \epsilon})$ für $\epsilon > 0$, dann gilt $T(n) \in \Theta(n^{\log_b a})$.
2. Wenn $f(n) \in \Theta(n^{\log_b a})$, dann gilt $T(n) \in \Theta(n^{\log_b a} \log n)$.
3. Wenn $f(n) \in \Omega(n^{\log_b a + \epsilon})$ für $\epsilon > 0$ und $a f(n/b) \leq c \cdot f(n)$ für $c < 1$ und für hinreichend große n , dann gilt $T(n) \in \Theta(f(n))$.

Die Lesart des obigen Theorems ist, dass wir hier einen Vergleich zwischen $f(n)$ und $n^{\log_b a}$ durchführen. Falls $f(n)$ asymptotisch kleiner ist als $n^{\log_b a}$ (Fall 1), wird $T(n)$ durch $n^{\log_b a}$ dominiert, falls $f(n)$ asymptotisch größer ist als $n^{\log_b a}$, wird $T(n)$ durch $f(n)$ dominiert (Fall 3). Wenn die beiden Funktionen die gleiche asymptotischen Größenordnung besitzen, gilt $T(n) = \Theta(f(n) \log n)$. Kleiner beziehungsweise größer ist hier als *polynomial* kleiner oder größer gedacht, also ein Unterschied in n^ϵ für ein $\epsilon > 0$.

Beispiel 1: Für $T(n) = 9 T(n/3) + n$ gilt $a = 9$ und $b = 3$ und $f(n) = n$, dann ist $n^{\log_b a} = n^{\log_3 9} \in \Theta(n^2)$. Fall 1 tritt ein und $T(n) \in \Theta(n^2)$.

Beispiel 2: Für $T(n) = T(2n/3) + 1$ gilt $a = 1$ und $b = \frac{3}{2}$ und $f(n) = 1$ und es ist $n^{\log_b a} = n^{\log_{3/2} 1} = n^0$ und somit $f(n) \in \Theta(n^{\log_b a})$. Fall 2 ist anwendbar und es gilt $T(n) \in \Theta(\log n)$.

Beispiel 3: Für $T(n) = 3 T(n/4) + n \log n$ gilt $a = 3$ und $b = 4$ und $f(n) = n \log n$ und es ist $n^{\log_b a} = n^{\log_4 3} = O(n^{0.8})$. Da nun $f(n) \in \Omega(n^{\log_4 3 + \epsilon})$ gilt und die *Regularitätsbedingung* $a f(n/b) \leq c f(n)$ für $c = \frac{3}{4}$ erfüllt ist, gilt $T(n) \in \Theta(n \log n)$. Es gilt $a f(n/b) = 3(n/4) \log(n/4) \leq 3/4 n \log n = 3/4 f(n)$.

Beispiel 4: Für $T(n) = 2 T(n/2) + n \log n$ gilt $a = 2$ und $b = 2$ und $f(n) = n \log n$ und $n^{\log_b a} = n^{\log_2 2} = n$. Nun ist aber $f(n)$ nicht polynomiell größer als $n^{\log_b a}$ da $\frac{f(n)}{n} = \log n$. Für jedes $\epsilon > 0$ ist n^ϵ asymptotisch größer als $\log n$. $f(n)$ ist zwar größer als $n^{\log_b a}$ aber nicht polynomiell größer. Das Theorem ist nicht anwendbar, eine Lücke zwischen Fall 2 und Fall 3. Das gleiche kann passieren, wenn $f(n)$ zwar kleiner ist als $n^{\log_b a}$ aber nicht polynomiell kleiner, das ist eine Lücke zwischen Fall 1 und Fall 2.

Wir beweisen das Mastertheorem in Teilen für Werte die durch Potenzen von b gegeben sind. Der allgemeine Beweis befindet sich zum Beispiel in [1]. Zunächst schätzen wir $T(n)$ für alle drei Fälle durch einen Rekursionsbaum ab. Danach werden für die Fälle die jeweiligen Eigenschaften der Funktion f benutzt.

Lemma 4 Für Konstanten $a \geq 1$ und $b > 1$ und eine nichtnegative Funktion $f(n)$ sei

$$T(n) = \begin{cases} \Theta(1) & \text{falls } n = 1 \\ aT\left(\frac{n}{b}\right) + f(n) & \text{falls } n = b^i \end{cases} \quad (3.4)$$

über die Potenzen von b definiert. Dann gilt

$$T(n) = \Theta\left(n^{\log_b a}\right) + \sum_{j=0}^{\log_b n - 1} a^j f(n/b^j). \quad (3.5)$$

Beweis. Wir schätzen $T(n)$ durch die Rekursionsbaummethode exakt für $n = b^i$ ab. Der Baum ist *vollständig* und hat eine Höhe von $i = \log_b n$. Jeder Knoten in Ebene j beschreibt ein Teilproblem mit Zusatzkosten $f(n/b^j)$. Bis auf die Blätter besitzt jeder Knoten a Kinder in der nächsten Ebene. Für die Darstellung des Baumes nehmen wir an, dass a eine ganze Zahl ist, logisch betrachtet ist das aber nicht notwendig. Das bedeutet, dass in Ebene j insgesamt a^j Knoten liegen. Die Blätter auf Ebene $\log_b n$ haben Kosten $\Theta(1)$. Insgesamt gibt es $a^{\log_b n} = n^{\log_b a}$ Blätter, siehe Abbildung 3.2. \square

Die Gesamtkosten belaufen sich dann für die inneren Knoten auf

$$\sum_{j=0}^{\log_b n - 1} a^j f(n/b^j)$$

und für die Blätter auf $c \cdot n^{\log_b a}$. Jetzt schätzen wir die Summe aus dem Term (3.5) gemäß des Mastertheorems ab.

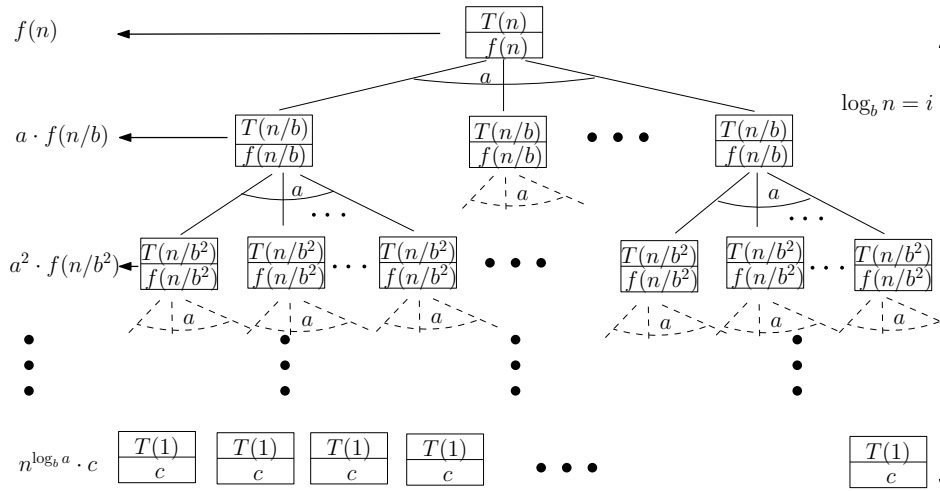


Abbildung 3.2: Der Rekursionsbaum zur Gleichung (3.4) hat $i = \log_b n$ Ebenen. In jeder Ebene fallen konkrete Kosten $a^j f(n/b^j)$ an.

Lemma 5 Seien $a \geq 1$ und $b > 1$ Konstanten und $f(n)$ eine nichtnegative Funktion. Dann gilt für die über Potenzen von b definierte Funktion

$$g(n) = \sum_{j=0}^{\log_b n - 1} a^j f(n/b^j). \tag{3.6}$$

1. Wenn $f(n) \in O(n^{\log_b a - \epsilon})$ für $\epsilon > 0$, dann gilt $g(n) \in O(n^{\log_b a})$.
2. Wenn $f(n) \in \Theta(n^{\log_b a})$, dann gilt $g(n) \in \Theta(n^{\log_b a} \log n)$.
3. Wenn $f(n) \in \Omega(n^{\log_b a + \epsilon})$ für $\epsilon > 0$ und $a f(n/b) \leq c f(n)$ für $c < 1$ und für hinreichend große n , dann gilt $g(n) \in \Theta(f(n))$.

Beweis. Wir beweisen hier den ersten Fall, die anderen Fälle werden ähnlich behandelt, siehe auch [1]. Wir haben $f(n) \in O(n^{\log_b a - \epsilon})$ und somit $f(n/b^j) \in O((n/b^j)^{\log_b a - \epsilon})$. Einsetzen in (3.6) ergibt

$$g(n) \in O\left(\sum_{j=0}^{\log_b n - 1} a^j \left(\frac{n}{b^j}\right)^{\log_b a - \epsilon}\right). \tag{3.7}$$

Wir vereinfachen den Term innerhalb der O -Notation und wenden die geometrische Reihe an.

$$\begin{aligned} \sum_{j=0}^{\log_b n - 1} a^j \left(\frac{n}{b^j}\right)^{\log_b a - \epsilon} &= n^{\log_b a - \epsilon} \sum_{j=0}^{\log_b n - 1} \left(\frac{ab^\epsilon}{b^{\log_b a}}\right)^j \\ &= n^{\log_b a - \epsilon} \sum_{j=0}^{\log_b n - 1} (b^\epsilon)^j \\ &= n^{\log_b a - \epsilon} \left(\frac{b^{\epsilon \log_b n} - 1}{b^\epsilon - 1}\right) \\ &= n^{\log_b a - \epsilon} \left(\frac{n^\epsilon - 1}{b^\epsilon - 1}\right) \end{aligned}$$

Da b und ϵ feste Konstanten sind kann der Term $\left(\frac{n^\epsilon - 1}{b^\epsilon - 1}\right)$ durch $O(n^\epsilon)$ abgeschätzt werden und es gilt: $g(n) \in O\left(n^{\log_b a}\right)$.

Die Fälle 2 und 3 sind noch etwas einfacher. □

Abschließend beweisen wir beispielhaft Nummer 1. in Theorem 3.

Beweis. (Beweis Theorem 3 Fall 1.) Aus Lemma 5 und Lemma 4 folgern wir

$$T(n) \in \Theta\left(n^{\log_b a}\right) + O\left(n^{\log_b a}\right) \in \Theta\left(n^{\log_b a}\right).$$

□

Wie bereits erwähnt können die vollständigen Beweise sehr gut in [1] nachgelesen werden.

Kapitel 4

Dynamische Programmierung

Das Prinzip der dynamischen Programmierung kann dann angewendet werden, wenn eine rekursive Problemzerlegung an vielen Stellen die gleiche Teillösung verlangt. In diesem Fall ist es ratsam, die Lösungen dieser Teilprobleme zu speichern und gegebenenfalls darauf zurückzugreifen.

Falls die Teillösungen immer wieder neu berechnet werden, kann der Rechenaufwand exponentiell steigen. Grundsätzlich besteht die dynamische Programmierung aus

- einem rekursiven Aufruf für Teilprobleme und aus
- einer Tabellierung von Zwischenergebnissen die
- wiederverwendet werden sollen.

Wir betrachten zunächst ein einfaches klassisches Beispiel, das diese wesentlichen Bestandteile bereits charakterisiert. Generell werden durch die dynamischen Programmierung meistens Optimierungsprobleme gelöst. In solche Fällen wird dann versucht, zu jedem Teilproblem eine optimale Lösung zu finden.

Dynamische Programmierung ist dann sinnvoll, wenn das sogenannte *Optimalitätsprinzip von Bellman* vorliegt:

Eine optimale Lösung hat stets eine Zerlegung in optimale Teilprobleme.

Ein solches Beispiel wird in Abschnitt 4.2 vorgestellt.

4.1 Fibonacci Zahlen

Wir betrachten die Fibonacci-Zahlen, die rekursiv wie folgt definiert werden.

$$\text{fib}(0) = 1 \quad (1)$$

$$\text{fib}(1) = 1 \quad (2)$$

$$\text{fib}(n) = \text{fib}(n-1) + \text{fib}(n-2) \quad (3)$$

Wollen wir nun beispielsweise $\text{fib}(5)$ berechnen, so ergibt sich folgender Rekursionsbaum für die Aufrufe von fib ; siehe Abbildung 4.1(I). Dieser Baum hat für $\text{fib}(n)$ eine Größe von $\Omega\left(\left(\frac{1+\sqrt{5}}{2}\right)^n\right)$, wie eine Übungsaufgabe zeigt.

Offensichtlich ist es ratsam, einige der Ergebnisse zu speichern, die rekursiven Abhängigkeiten können einfacher dargestellt werden, wie in Abbildung 4.1(II). Wir wollen zum Beispiel $\text{fib}(2)$

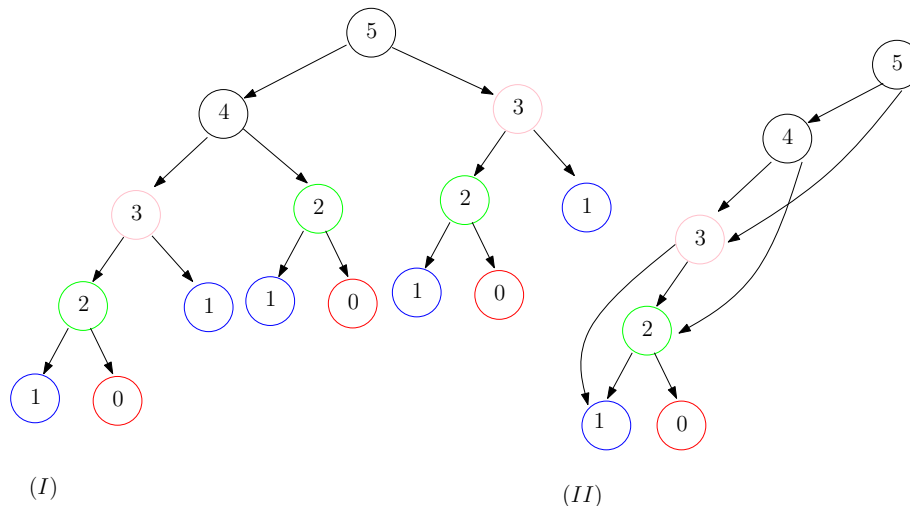


Abbildung 4.1: (I) Der vollständige Rekursionsbaum für fib(5) enthält einige Redundanzen. (II) Die rekursiven Abhängigkeiten können einfach dargestellt werden.

nicht dreimal wieder aus den Vorgängern neu berechnen, deshalb speichern wir die Werte ab. Es ist hier so, dass sich Teilprobleme überlappen und diese effizient gespeichert werden sollen. Der Abhängigkeitsgraph in Abbildung 4.1(II) hat nur eine Größe von $\Theta(n)$. Das Prinzip der Speicherung von Zwischenergebnissen wird auch *Memoisation* genannt.

Wenn wir nun die Abhängigkeiten nutzen wollen und von oben nach unten (Top-Down) die benötigten Werte fib(i) bestimmen wollen, könnten wir algorithmisch wie in Algorithmus 4 vorgehen, wobei FibMemo die memoisierende Prozedur 5 ist.

Algorithm 4 FibTopDown(n)

```

1:  $F[1] := 1; F[0] := 1;$ 
2: for  $i = 2$  to  $n$  do
3:    $F[i] := 0;$ 
4: end for
5: RETURN FibMemo( $n, F$ )

```

Hierbei ist FibMemo die folgende memoisierende Prozedur.

Procedure 5 FibMemo(n, F)

```

1: if  $F[n] > 0$  then
2:   RETURN  $F[n]$ 
3: end if
4:  $F[n] := \text{FibMemo}(n - 1, F) + \text{FibMemo}(n - 2, F);$ 
5: RETURN  $F[n]$ 

```

In Algorithmus 4 wird das Feld F zunächst mit Nullen initialisiert und beginnen dann den Aufruf der rekursiven Prozedur 5. Falls der Wert von $F[n]$ bekannt ist, geben wir diesen hier zurück, sonst wird er rekursiv gemäß der Rekursionsformel berechnet.

Laufzeitabschätzung: Wir benötigen im Algorithmus 4 $O(n)$ Zeit für die Initialisierung. Danach wird FibMemo(n, F) aufgerufen. Jeder Aufruf von FibMemo(m, F) in Prozedur 5 generiert höchstens einmal die Aufrufe FibMemo($m - 1, F$) und FibMemo($m - 2, F$). Für jedes

$l < n$ wird somit $\text{FibMemo}(l, F)$ höchstens zweimal aufgerufen. Ohne die rekursiven Aufrufe wird nur $O(1)$ Zeit in Prozedur 5 benötigt. Insgesamt gilt für die Laufzeit $T(n) \in O(n)$.

Wir machen nun die folgende Beobachtung. Obwohl im obigen Algorithmus $F[n]$ in einer Top-Down Vorgehensweise errechnet wird, werden die eigentlichen Werte von unten nach oben berechnet. Erst wenn die Rekursion bei $F[1]$ und $F[0]$ angekommen ist, werden neue Felder belegt und die Memoisation setzt ein. Wenn wir beispielsweise $F[5]$ berechnen, rufen wir sukzessive $\text{FibMemo}(5, F)$, $\text{FibMemo}(4, F)$, $\text{FibMemo}(3, F)$, $\text{FibMemo}(2, F)$ und $\text{FibMemo}(1, F)$ auf bis hier zum ersten Mal ein Wert $F[1]$ zurückgegeben wird. Dann wird wieder $\text{FibMemo}(0, F)$ aufgerufen und $F[2]$ belegt. Dann wird $\text{FibMemo}(1, F)$ nochmal aufgerufen und $F[3]$ belegt. Danach rufen wir $\text{FibMemo}(2, F)$ auf und belegen $F[4]$. Dann $\text{FibMemo}(3, F)$ und $F[5]$ wird belegt und ausgegeben.

Deshalb ist es offensichtlich sinnvoller gleich die Werte von unten nach oben (Bottom-Up) zu berechnen. Da für die Berechnung von $F[n]$ nur zwei Werte aus $F[0], \dots, F[n-1]$ benötigt werden, sparen wir außerdem Speicherplatz. Die obige Memoisation und die Top-Down vorgehensweise ist hier bezüglich des Speicherplatzes nicht effizient.

Insgesamt erhalten wir folgende effiziente dynamische Programmierungslösung.

Algorithm 6 $\text{FibDynProg}(n)$

```

1:  $F_{\text{letzt}} := 1; F_{\text{vorletzt}} := 1;$ 
2: for  $i = 2$  to  $n$  do
3:    $F := F_{\text{vorletzt}} + F_{\text{letzt}};$ 
4:    $F_{\text{vorletzt}} := F_{\text{letzt}}; F_{\text{letzt}} := F;$ 
5: end for
6: RETURN  $F$ 

```

- Es müssen insgesamt nur $\Theta(n)$ viele Berechnungen durchgeführt werden.
- Die Bottom-Up Berechnung stellt sicher, dass die benötigten Werte stets vorliegen.
- Der Speicherplatzbedarf liegt in $O(1)$.

Beachte: Manchmal kann eine Top-Down vorgehensweise auch sinnvoller sein, wenn zum Beispiel gar nicht alle Werte verwendet werden, die wir Bottom-Up berechnen würden. Die Prinzipien sind also nicht *dogmatisch* zu interpretieren.

Grundprinzip der dynamischen Programmierung:

- Formuliere das Problem insgesamt rekursiv. Stelle fest, dass sich Teillösungen überlappen und rekursiv voneinander abhängen.
- Stelle eine Rekursionsgleichung (Top-Down) auf.
- Löse das Problem (Bottom-Up) durch das Ausfüllen von Tabellen mit Teillösungen (Memoisation).
- Ziel: Die Zahl der Tabelleneinträge soll klein gehalten werden.
- Anwendung: Bei einer direkten Ausführung der Rekursion ist mit vielen Mehrfachberechnungen zu rechnen.

Die Korrektheit der dynamischen Programmierung ergibt sich in der Regel dadurch, dass die rekursive Problembeschreibung korrekt durchgeführt wurde. Die Laufzeitanalyse ist in

der Regel nicht schwer, da die Algorithmen typischerweise aus geschachtelten FOR-Schleifen bestehen. Das eigentliche Problem der dynamischen Programmierung ist es, eine geeignete rekursive Formulierung des Problems zu finden. Wir werden uns deshalb ein signifikantes Optimierungsbeispiel ansehen in dem auch das *Optimalitätsprinzip von Bellman* zum tragen kommt.

4.2 Matrixmultiplikation

Wir betrachten das Multiplizieren von reellwertigen Matrizen. Seien $A \in \mathbb{R}^{i \times j}$ und $B \in \mathbb{R}^{j \times k}$ Matrizen mit i Zeilen und j Spalten bzw. j Zeilen und k Spalten, dann beschreibt $A \cdot B$ die Multiplikation dieser beiden Matrizen. Insgesamt sind zur Berechnung von $A \cdot B = C \in \mathbb{R}^{i \times k}$ $i \cdot k \cdot j$ Floating-Point Berechnungen notwendig. Genauer: Die Matrixeinträge von C werden durch

$$c_{st} = \sum_{l=1}^j a_{sl} \cdot b_{lt}$$

für alle $s = 1, \dots, i$ und alle $t = 1, \dots, k$ berechnet, also $i \cdot j \cdot k$ Multiplikationen und Additionen.

Falls wir nun mehrere Matrizen multiplizieren wollen, also $E = A_1 \cdot A_2 \cdot \dots \cdot A_n$ dann müssen auch hier die entsprechenden Matrizen gemäß der Größen *zusammenpassen*. Die Anzahl der Spalten von A_i muss der Anzahl der Zeilen von A_{i+1} entsprechen für $i = 1, \dots, n-1$.

Die Matrixmultiplikation ist *assoziativ*, für drei Matrizen $A \in \mathbb{R}^{i \times j}$, $A \in \mathbb{R}^{j \times k}$ und $C \in \mathbb{R}^{k \times l}$ gilt also $(A \cdot B) \cdot C = A \cdot (B \cdot C) \in \mathbb{R}^{i \times l}$. Deshalb lassen sich auch allgemeine (zusammenpassende Ketten) beliebig Klammern. Diese Klammerung hat Einfluss auf die Anzahl der Floating-Point Operationen.

Beispiel: Für $A \in \mathbb{R}^{10 \times 50}$, $A \in \mathbb{R}^{50 \times 10}$ und $C \in \mathbb{R}^{10 \times 50}$, benötigt $(A \cdot B) \cdot C$ zuerst $10 \times 50 \times 10$ Operationen für $(A \cdot B)$ und danach $10 \times 10 \times 50$ Operationen für die Multiplikation des Ergebnisses mit C also insgesamt 10000 Operationen. Dahingegen benötigt $A \cdot (B \cdot C)$ zuerst $50 \times 10 \times 50$ Operationen für $(B \cdot C)$ und danach $10 \times 50 \times 50$ Operationen für die Multiplikation des Ergebnisses mit A also insgesamt 50000 Operationen. Die Klammerung $(A \cdot B) \cdot C$ ist also deutlich günstiger.

Allgemeine Problembeschreibung: Bestimme für die Berechnung von $E = A_1 \cdot A_2 \cdot \dots \cdot A_n$ mit Dimensionen $d_0 \times d_1, d_1 \times d_2, \dots, d_{n-1} \times d_n$ die optimale Klammerung für die minimale Anzahl an Operationen.

Dieses Problem wollen wir durch dynamische Programmierung lösen. Deshalb muss zunächst die Lösung rekursiv durch Teillösungen angegeben werden. Falls $M(n)$ alle Möglichkeiten beschreibt, wie man n Matrizen klammert, dann ist $M(n) = \sum_{k=1}^{n-1} M(k) \cdot M(n-k)$ mit $M(1) = 1$. Diese Rekursionsgleichung liegt in $\Omega(2^n)$. Es ist also nicht zweckmäßig, alle Möglichkeiten auszuprobieren.

Übungsaufgabe: Verifizieren Sie die Lösung $M(n) \in \Omega(2^n)$.

Wir wollen nun optimale Teillösungen wiederverwenden. Sei nun $K(l, k)$ der minimale Aufwand bzw. die optimale Klammerung für $A_l \cdot A_{l+1} \cdot \dots \cdot A_k$. Zum Beispiel zerlegt sich zu Beginn ein erster Multiplikationsschritt für die Stelle s , d.h. der minimale Aufwand für $(A_1 \cdot A_2 \cdot \dots \cdot A_s) \cdot (A_{s+1} \cdot A_{s+2} \cdot \dots \cdot A_n)$, rekursiv in die Kosten

$$K(1, s) + K(s, n) + d_0 \cdot d_s \cdot d_n \quad (4.1)$$

Das Ziel ist nun, das optimale s zu finden, und das wollen wir rekursiv beschreiben.

- Offensichtlich ist $K(l, l) = 0$ für $l = 1, \dots, n$

- Die Dimension einer Teilkette $A_l \cdot A_{l+1} \cdot \dots \cdot A_k$ ist $d_{l-1} \times d_k$.
- Die Kosten für ein Teilen bei s sind $K(l, s) + d_{l-1} \cdot d_s \cdot d_k + K(s + 1, k)$

Da wir $K(1, n)$ minimieren wollen erhalten wir insgesamt:

$$K(l, k) = \begin{cases} 0 & \text{falls } l = k \\ \min_{l \leq s \leq k-1} K(l, s) + d_{l-1} \cdot d_s \cdot d_k + K(s + 1, k) & \text{falls } l < k \end{cases} \quad (4.2)$$

Wir stellen fest, dass dynamische Programmierung geeignet ist:

- Wiederverwendung: Hier werden viele Teilprobleme mehrfach benutzt, zum Beispiel $K(1, 2)$ wird verwendet von $K(1, 3), K(1, 4), \dots, K(1, n)$.
- Beschränkte Anzahl: Für alle $1 \leq l < k \leq n$ gibt es genau ein Teilproblem, also insgesamt $\Theta(n^2)$ viele.
- Bearbeitungsreihenfolge: Falls alle Werte $K(l, s)$ und $K(s + 1, k)$ bekannt sind, kann $K(l, k)$ in Zeit $O(k - l) \in O(n)$ berechnet werden.

Die geschickte Auswertung lässt sich sehr gut durch eine Matrix beschreiben. Wir verwenden ein Array $K[1..n, 1..n]$.

Zuerst werden die Diagonalen eingetragen, dann berechnen wir sukzessive die Nebendiagonalen von rechts unten nach links oben, zu jedem Zeitpunkt sind für $K(l, k)$ alle Werte $K(l, s)$ und alle Werte $K(s + 1, k)$ für $l \leq s$ und $s + 1 \leq k$ berechnet worden; siehe Abbildung 4.2.

	1	2	3	4	5
1	0	[1, 2]	[1, 3]	??	
2		0	[2, 3]	[2, 4]	[2, 5]
3			0	[3, 4]	[3, 5]
4				0	[4, 5]
5					0

Abbildung 4.2: Die Nebendiagonalen können von unten nach oben berechnet werden. Nachdem zwei Nebendiagonalen bereits berechnet wurden beginnen wir die Berechnung der dritten Nebendiagonalen. Die benötigten Werte liegen bereits vor.

Korrektheit: Am Ende der Auswertung wird der Wert $[1, n]$ die minimale Anzahl an Operationen enthalten, wir haben den Wert rekursiv genauso festgelegt. Wir wollen uns natürlich auch die optimale Klammerung merken. Dazu müssen wir uns nur für jeden berechneten Matrixeintrag den jeweils optimalen Index $s[l, k]$ merken. Falls beispielsweise der optimale Index für die Berechnung von $[1, 5]$ bei $s = 2$ liegt, ist die Klammerung $(A_1 A_2)(A_3 A_4 A_5)$ rekursiv

die Beste und wir speichern den Wert $s[1, 5] = 2$. Die optimale Klammerung von $(A_1 A_2)$ zur Berechnung von $[1, 2]$ ist trivialerweise $S = 1$ also $s[1, 2] = 1$. Rekursiv haben wir auch die optimale Klammerung aus $[3, 5]$ gespeichert, diese kann $s = 3$ oder $s = 4$ sein, somit ist $s[3, 5] \in \{3, 4\}$. Insgesamt erhalten wir so die optimale Klammerung.

Laufzeitanalyse und Speicherplatz: Insgesamt wird Speicherplatz von $O(n^2)$ für beide Matrizen (Einträge $K[l, k]$ und $s[l, k]$) verwendet. Für die Berechnung eines Eintrages $K[l, k]$ und der Klammerung $s[l, k]$ wird ein Aufwand von $O(k - l) \in O(n)$ benötigt. Da wir $O(n^2)$ viele Aufrufe haben, ergibt sich ein Gesamtaufwand von $O(n^3)$.

Übungsaufgabe: Formulieren Sie einen Pseudocode, der den obigen Algorithmus beschreibt.

4.3 Rucksackproblem

Wir betrachten das Problem des optimalen Befüllens eines Rucksackes mit Gesamtgewichtskapazität G . Dazu sei eine Menge von n Gegenständen a_i mit Gewicht g_i und Wert w_i gegeben. Wir wollen den Rucksack optimal füllen.

Aufgabenstellung: Bestimme eine Teilmenge $\{a_{i_1}, a_{i_2}, \dots, a_{i_k}\}$ der Elemente aus $\{a_1, \dots, a_n\}$, so dass $\sum_{j=1}^k g_{i_j} \leq G$ gilt und $\sum_{j=1}^k w_{i_j}$ maximal ist. Bestimme also den Wert:

$$w_{\max} = \max_{\{a_{i_1}, a_{i_2}, \dots, a_{i_k}\} \subseteq \{a_1, \dots, a_n\}} \left\{ \sum_{j=1}^k w_{i_j} \mid \sum_{j=1}^k g_{i_j} \leq G \right\}.$$

Wenn wir dieses Problem rekursiv lösen wollen, können wir naiv einen Rekursionsbaum angeben, der in jeder Ebene i die Entscheidung für den i -ten Gegenstand berücksichtigt; siehe Abbildung 4.3. Dadurch werden alle möglichen Fälle beschrieben. Das ist aber nicht effizient, da der Baum eine exponentielle Größe.

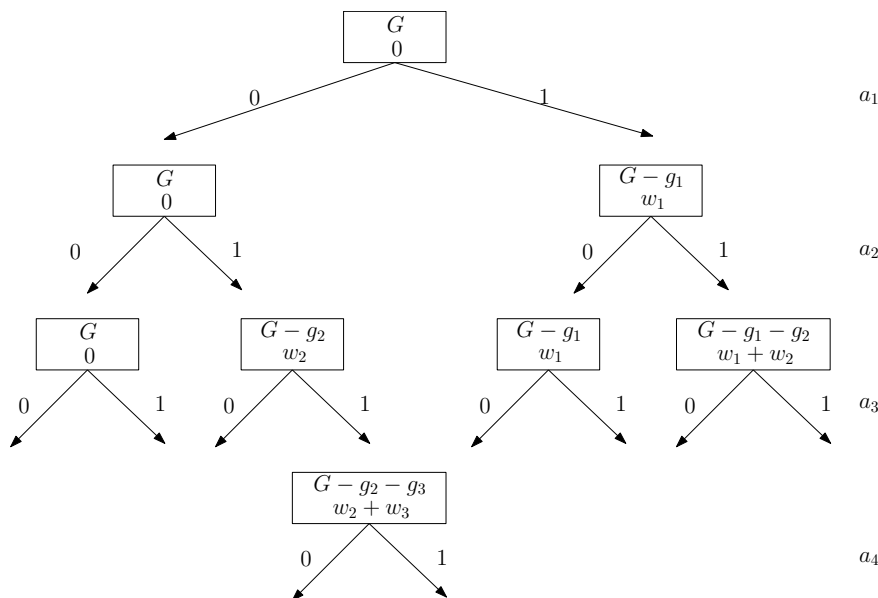


Abbildung 4.3: In Ebene i wird die Entscheidung für den i -ten Gegenstand getroffen. Der Baum hat exponentielle Größe.

Wir betrachten deshalb die folgende Vereinfachung und wollen Teilergebnisse verwenden. Sei dazu $W(i, j)$ der optimale Wert falls nur die ersten i Gegenstände betrachtet werden und der

Rucksack die Kapazität j hat. Am Ende soll $w_{\max} = W(n, G)$ berechnet werden.

Wir können $W(i, j)$ rekursiv wie folgt angeben. Die Frage stellt sich, ob bei Lösungen mit den ersten $i - 1$ Gegenständen der i -te Gegenstand hinzugenommen werden soll oder nicht. Falls nicht, bleibt der Wert bei $W(i - 1, j)$ stehen. Falls ja, kommt der Wert w_i hinzu. Allerdings müssen wir zur Berechnung von $W(i, j)$ dann rekursiv auf $W(i - 1, j - g_i)$ zurückgreifen, $W(i - 1, j - g_i) + w_i$ ist dann die Lösung. Der maximale Wert aus diesen beiden Fällen ergibt den Wert $W(i, j)$.

Falls j negativ wird, ist das Ergebnis ungültig. Das drücken wir durch $W(i, j) = -\infty$ aus. Falls i null wird, setzen wir den Wert $W(i, j)$ auf null, da kein Gegenstand verwendet wird.

Wir kommen zu folgender Rekursionsgleichung, die wir danach Bottom-Up lösen wollen.

$$W(i, j) = \begin{cases} 0 & \text{falls } i = 0 \\ -\infty & \text{falls } j < 0 \\ \max\{W(i - 1, j), W(i - 1, j - g_i) + w_i\} & \text{sonst} \end{cases} \quad (4.3)$$

Nun wollen wir wiederum die Werte $W(i, j)$ in einem Array $W[0..n, 0..G]$ berechnen und initialisieren $W[0, j]$ und $W[i, 0]$ mit 0. Aus (4.3) ergibt sich offensichtlich, dass wir die Werte am besten Zeilenweise von links nach rechts und mit steigender Zeilenzahl berechnen. Zur Berechnung von $W(i, j)$ muss $W(i - 1, j)$ und $W(i - 1, j - g_i)$ vorliegen.

Beispiel: Für $n = 5$ notieren $a_i = (g_i, w_i)$ mit $a_1 = (6, 4)$, $a_2 = (1, 2)$, $a_3 = (2, 3)$, $a_4 = (5, 8)$ und $a_5 = (9, 10)$ mit $G = 12$.

Wir beschreiben einen Teilausschnitt der zugehörigen Matrix während der Berechnung. Zur Berechnung von $W(2, 7)$ wird $W(1, 7) = 4$ (ohne a_2) und $W(1, 7 - g_2) + w_2 = W(1, 6) + 2 = 4 + 2 = 6$ (mit a_2) verglichen. $W(1, 7)$ ist also gleich 6. Zeilenweise wird die Matrix von links nach rechts gefüllt.

		0	1	2	3	4	5	6	7	8	9	10	11	12
0)	0	0	0	0	0	0	0	0	0	0	0	0	0
1		0	0	0	0	0	0	4	4	4	4	4	4	4
2		0	2	2	2	2	2	4	$W(2, 7)?$	x	x	x	x	x
3		0	x	x	x	x	x	x	x	x	x	x	x	x
4		0	x	x	x	x	x	x	x	x	x	x	x	x
5		0	x	x	x	x	x	x	x	x	x	x	x	x

Auch für diesen Fall wollen wir uns am Ende die optimale Lösung generieren können. Dazu merken wir uns in einer zweiten Matrix $L[0..n, 0..G]$, woraus $W(i, j)$ entstanden ist. Zum Beispiel merken wir uns für $W(1, 7)$ das a_2 benutzt wird und der Plan aus $L(1, 6)$. Diese Information wird in $L(2, 7)$ gespeichert, zum Beispiel durch $L[2, 7] := (a_2, [1, 6])$. In $L(1, 6)$ steht dann wiederum, dass a_1 benutzt wird.

Korrektheit: Offensichtlich läßt sich nach Berechnung von $W(n, G)$ und $L(n, G)$ der Lösungsplan rekursiv angeben und das optimale Packen ist berechnet worden. Die Lösung ist rekursiv genauso formuliert worden.

Wir geben hier nochmal explizit die Formulierung im Pseudocode in Algorithmus 7 an und analysieren dann die Korrektheit und die Laufzeit.

Laufzeit und Speicherplatz: In Algorithmus 7 werden nach der Initialisierung der Werte $W[i, 0]$ und $W[0, j]$ zwei *FOR*-Schleifen geschachtelt, wobei die Äußere von 1 bis n läuft und die Innere von 1 bis G . Im Inneren ist der Aufwand $\Theta(1)$. Also ist der Aufwand insgesamt in $O(n \cdot G)$ und für das Array benötigen wir $O(n \cdot G)$ Speicherplatz.

Algorithm 7 RucksackDynProg($a_i = (g_i, w_i), i = 1, \dots, n$, Gesamtgew.: G)

```
1: for  $i = 0$  to  $G$  do
2:    $W[0, i] := 0$ ;
3: end for
4: for  $j = 0$  to  $n$  do
5:    $W[j, 0] := 0$ ;
6: end for
7: for  $i = 1$  to  $n$  do
8:   for  $j = 1$  to  $G$  do
9:     if  $(j - g_i) < 0$  then
10:       $W[i, j] := W[i - 1, j]$ ; //  $a_i$  passt gar nicht rein
11:    else
12:       $W[i, j] := \max\{W[i - 1, j], W[i - 1, j - g_i] + w_i\}$ ;
13:    end if
14:  end for
15: end for
16: RETURN  $W[n, G]$ 
```

Übungsaufgabe: Füllen Sie die Matrizen $W[0..5, 0..12]$ und $L[0..5, 0..12]$ für das obige Beispiel. Formulieren Sie die algorithmische Lösung im Pseudocode auch für die Belegung der Matrix L zur Rekonstruktion der optimalen Lösung.

Bemerkungen: Der obige Algorithmus ist für ganzzahlige Werte entworfen worden und die Laufzeit hängt polynomiell von der Größe der Rucksacks ab. Für nicht-ganzzahlige Werte ist das Problem deutlich schwerer zu lösen. Es ist nach heutiger Erkenntnis nicht zu erwarten, dass ein polynomieller Algorithmus für das allgemeine Rucksackproblem mit nicht-ganzzahligen Werten existiert.

Kapitel 5

Greedy Algorithmen

Neben den bereits behandelten klassischen Entwurfsmethoden der Algorithmik gibt es eine weitere Methode, die sich auf die intuitive Idee stützt, dass lokal optimale Berechnungsschritte auch global optimal sind. Es geht also auch in diesem Kapitel um Optimierungsprobleme. Der lokal optimale Berechnungsschritt soll dabei möglichst nicht wieder rückgängig gemacht werden. Außerdem werden keine Alternativen ausprobiert. In jedem Schritt wird der nächste lokal optimale Schritt berechnet und tatsächlich ausgeführt. In diesem Sinne nennt man diese Vorgehensweise auch *iterativ*. Die Kosten für einen *Iterationsschritt* sollen klein gehalten werden. Das ist meistens möglich, da wir zunächst gar nicht auf die Gesamtlösung achten. Zur Bestimmung des nächsten Schrittes wird allerdings häufig eine Liste von Kandidaten für den nächsten Schritt vorgehalten. Diese Liste wird nach dem Schritt aktualisiert. Greedy-Verfahren sind dennoch in der Regel einfach zu beschreiben und benötigen wenig Rechenaufwand.

Dieses *Greedy*-Prinzip wird leider nicht immer die optimale Lösung erzielen, drei Szenarien sind denkbar:

1. Wir erhalten die optimale Lösung.
2. Wir erhalten nicht die optimale Lösung, aber wir können zumindest beweisen, dass die Lösung im Vergleich zum Optimum nicht beliebig schlecht ist.
3. Die Lösung ist im Vergleich zum Optimum beliebig schlecht.

Zusammengefasst ergibt sich folgendes Grundprinzip der Greedy-Vorgehensweise:

- Lokal optimale Schritte werden iterativ ausgeführt.
- Eine Kandidatenliste wird dafür bereithalten und ggf. aktualisiert.
- Es werden keine Alternativen ausprobiert oder Zwischenergebnisse tabellarisch abgespeichert.
- Die Kosten für den nächsten Schritt sind nicht sehr hoch.
- Die Lösung kann je nach Aufgabenstellung optimal, beliebig schlecht oder approximativ sein.

5.1 Rucksackproblem

Das bereits bekannte Rucksackproblem aus dem letzten Abschnitt lässt sich leicht in Greedy-Manier formulieren. Lokal optimal erscheint dabei, dass wir den Gegenstand mit dem jewei-

ligen besten Preis/Gewichtverhältnis als erstes Einfügen.

Beispiel: Für $n = 5$ notieren $a_i = (g_i, w_i)$ mit $a_1 = (6, 4)$, $a_2 = (1, 2)$, $a_3 = (2, 3)$, $a_4 = (5, 8)$ und $a_5 = (9, 10)$ mit $G = 12$. Die optimale Lösung ist hier offensichtlich (a_5, a_2, a_3) mit einem Wert von 15.

Zur Erinnerung, die Gegenstände haben ein Gewicht von g_i und einen Wert von w_i . Deshalb bestimmen wir zunächst eine Kandidatenliste gemäß des Verhältnisses $v_i = \frac{w_i}{g_i}$. Wir haben $v_1 = \frac{4}{6} = \frac{2}{3}$, $v_2 = \frac{2}{1} = 2$, $v_3 = \frac{3}{2} = 1.5$, $v_4 = \frac{8}{5}$ und $v_5 = \frac{10}{9}$. Das ergibt folgende Sortierung $v_2 > v_4 > v_3 > v_5 > v_1$. Ein Greedy Algorithmus wählt sukzessive den aktuell besten Kandidaten und revidiert diese Entscheidung nicht. Für den Algorithmus 8 nehmen wir zur einfachen Beschreibung an, dass die Gegenstände nach dem Preis/Gewichtsverhältnis sortiert vorliegen, also $v_1 \geq v_2 \geq \dots \geq v_n$.

Algorithm 8 RucksackGreedy($v_i = \frac{w_i}{g_i}$ nach Größe sortiert, Gewicht: G)

```

1: for  $i = 1$  to  $n$  do
2:   if  $g_i \leq G$  then
3:      $\lambda_i := 1$ ;  $G := G - g_i$ ;
4:   else
5:      $\lambda_i := 0$ ;
6:   end if
7: end for
8: RETURN  $\sum_{i=1}^n \lambda_i w_i$ 

```

Mit einer entsprechenden Umbenennung bedeutet das in unserem Beispiel, der Algorithmus sukzessive a_2 und a_4 und a_3 auswählt und dann die Kapazität $12 - 1 - 5 - 2 = 4$ verbleibt. Nun kann a_5 und a_1 nicht mehr verwendet werden. Insgesamt ergibt sich ein Gesamtwert von $2 + 8 + 3 = 13$. Der Greedy-Algorithmus ist also nicht optimal.

Schlimmer noch, wir können ein ganz einfaches Beispiel angeben, indem der Greedy-Algorithmus beliebig schlecht wird.

Worst-Case Beispiel: Sei $a_1 = (1, 1)$ und $a_2 = (G, G - 1)$. Dann ist $v_1 = 1 > \frac{G-1}{G} = v_2$. Der Greedy Algorithmus legt $\lambda_1 = 1$ und $\lambda_2 = 0$ fest und das Ergebnis ist 1. Optimal wäre $\lambda_1 = 0$ und $\lambda_2 = 1$ mit Ergebnis $G - 1$. Da G beliebig groß werden kann, wird das Verhältnis zwischen der Greedy-Lösung und der optimalen Lösung $\frac{G}{1}$ beliebig groß.

Wir haben jetzt ein Beispiel gesehen, in dem eine Greedy-Vorgehensweise nicht empfehlenswert ist, da wir keine Garantie für die Güte abgeben können. Wir möchten eine obere Schranke für das Verhältnis von $Lösung_mit_Greedy(P)$ und $Optimale_Lösung(P)$ bei einem Optimierungsproblem für alle Probleme $P \in \Pi$ für eine Problemklasse Π angeben. Das Verhältnis soll durch eine Konstante $C \geq 1$ angegeben werden, die besagt, dass die Greedy-Lösung höchstens C mal schlechter ist als die optimale Lösung. Deshalb betrachten wir verschiedene Quotienten wie folgt:

Maximierungsproblem: Gilt

$$\forall P \in \Pi \quad \frac{Optimale_Lösung(P)}{Lösung_mit_Greedy(P)} \leq C,$$

für eine Konstante C , so sagen wir, dass die Greedy-Lösung eine C -Approximation liefert.

Minimierungsproblem: Gilt

$$\forall P \in \Pi \quad \frac{Lösung_mit_Greedy(P)}{Optimale_Lösung(P)} \leq C,$$

für eine Konstante C , so sagen wir, dass die Greedy-Lösung eine C -Approximation liefert.

5.2 Das optimale Bepacken von Behältern

Das Problem des optimalen Bepackens von mehreren Behältern gleicher Gewichtskapazität G mit verschiedenen Paketen (Binpacking) gestaltet sich wie folgt. Die Pakete p haben keinen Preis sondern lediglich verschiedene Gewichte g . Insgesamt haben wir zunächst n Behälter b_j für $j = 1, \dots, n$ mit Kapazität G und n Gegenstände p_i mit Gewichten g_i für $i = 1, \dots, n$. Wir können $g_i \leq G$ annehmen, sonst brauchen wir einen Gegenstand gar nicht zu betrachten.

Ziel der Aufgabenstellung ist es, die Anzahl der benötigten Behälter zu minimieren. Im Prinzip bestimmen wir eine Zuordnung $Z : \{1, \dots, n\} \mapsto \{1, \dots, k\}$ die jedem p_i einen Behälter b_j über $Z(i) = j$ zuweist. Wir suchen das minimale k , für das es eine gültige Zuordnung Z gibt. Gültig heißt, dass die Summe aller g_i mit $Z(i) = j$ nicht größer als G ist.

Aufgabenstellung: Minimiere k , so dass eine Zuordnung $Z : \{1, \dots, n\} \mapsto \{1, \dots, k\}$ existiert mit

$$\sum_{Z(i)=j} g_i \leq G \text{ für alle } i \in \{1, \dots, n\}. \quad (5.1)$$

Für eine Greedystrategie benutzen wir hier keine besondere Kandidatenliste, wir nehmen die Pakete so wie sie gegeben sind in der Reihenfolge a_1, a_2, \dots, a_n . Die Annahme ist wirklich sinnvoll, da wir uns vorstellen können, dass die Pakete so vom Band kommen oder so geliefert werden und zugeordnet werden müssen.

Wir berücksichtigen nicht, welche Pakete in der Zukunft auftreten. Die einzige Entscheidung die wir treffen, ist also, in welchen Behälter b_j legen wir das Paket a_i nachdem bereits die Pakete a_1, a_2, \dots, a_{i-1} einsortiert wurden. Der maximale Index j soll klein gehalten werden. Naheliegend sind die folgenden beiden Vorgehensweisen. Angenommen, wir haben bereits a_1, a_2, \dots, a_{i-1} Pakete zugeordnet. Platziere nun a_i .

First-Fit: Unter allen Behältern wähle den ersten aus, in den a_i noch hineinpasst. Genauer: Finde kleinstes j , so dass $g_i + \sum_{1 \leq l \leq (i-1), Z(l)=j} g_l \leq G$ gilt.

Best-Fit: Packe a_i in die aktuell vollste Kiste, in die es noch hinein passt. Genauer: Finde ein j , so dass $g_i + \sum_{1 \leq l \leq (i-1), Z(l)=j} g_l \leq G$ gilt und dabei $\sum_{1 \leq l \leq (i-1), Z(l)=j} g_l$ maximal ist.

Beispiel: Betrachten wir Behälter der Kapazität $G = 10$ und 7 Gegenstände der Größen $g_1 = 2, g_2 = 5, g_3 = 4, g_4 = 7, g_5 = 1, g_6 = 3$ und $g_7 = 8$ dann belegt First-Fit die Behälter wie in Abbildung 5.1 und benötigt 4 Behälter. Optimal wäre es gewesen, drei Behälter jeweils mit $(g_1, g_7), (g_2, g_3, g_5)$ und (g_4, g_6) zu belegen. Beide Algorithmen lassen sich mit einer Laufzeit von $O(n^2)$ realisieren. Beispielsweise wird bei First-Fit für alle Behälter j der aktuelle Belegungswert gespeichert und wir überprüfen sukzessive von $j = 1, \dots, n$ in welchen Behälter der Gegenstand a_i noch hineinpasst. Dieses Verfahren wird wiederum sukzessive für $i = 1, \dots, n$ für die Gegenstände a_i durchgeführt. Insgesamt werden zwei FOR-Schleifen geschachtelt, die von $j = 1, \dots, n$ respektive von $i = 1, \dots, n$ laufen. Ähnlich kann bei Best-Fit vorgegangen werden.

Wir wollen nun zeigen, dass beide Algorithmen eine 2-Approximation der optimalen Anzahl der benötigten Behälter liefern. Wir gehen auch hier der Einfachheit halber davon aus, dass die Gegenstände ganzzahliges Gewicht haben.

Lemma 6 *Für das Binpacking Problem liefern die Algorithmen First-Fit und Best-Fit eine 2-Approximation für die optimale (minimale) Anzahl an benötigten Behältern.*

Beweis. Die Beweisidee stützt sich auf folgende Beobachtung. Nach Anwendung des jeweiligen Algorithmus gilt: Für je zwei *benutzte* Behälter gilt, dass das Gesamtgewicht der enthaltenen Gegenstände mindestens $G + 1$ ist.

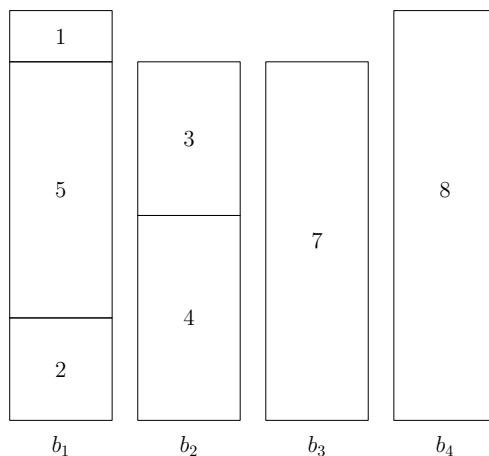


Abbildung 5.1: First-Fit für $g_1 = 2$, $g_2 = 5$, $g_3 = 4$, $g_4 = 7$, $g_5 = 1$, $g_6 = 3$ und $g_7 = 8$ benötigt einen Behälter zu viel.

Zum Beweis dieser Aussage führen wir einen Widerspruchsbeweis. Falls die obige Bedingung bei einer Belegung durch First-Fit bzw. Best-Fit nicht eintritt, ist die Belegung gar nicht durch den jeweiligen Algorithmus entstanden, siehe auch das Beispiel in 5.2.

Angenommen die obige Bedingung stimmt nach der Ausführung einer der beiden Algorithmen nicht, dann gibt es zwei Behälter b_l und b_k mit $l < k$, so dass die Gegenstände von b_k vollständig in b_l gepasst hätten und umgekehrt.

First-Fit hätte einen Gegenstand a_i aus b_k in ein b_j mit $j < k$ einfügen müssen, egal zu welchem Zeitpunkt a_i aufgetreten ist. Zumindest b_l mit $l < k$ war ein Kandidat dafür. Die Belegung ist somit sicher nicht durch First-Fit entstanden. Ein Widerspruch.

Best-Fit kann ebenfalls nicht ausgeführt worden sein. Zwischen den Alternativen b_l und b_k benutzt Best-Fit stets den Behälter mit dem höchsten Füllstand. Falls zuerst b_l das erste Paket von beiden Behältern überhaupt erhält, kann der Algorithmus die weiteren Pakete unter dem Best-Fit Aspekt nur noch diesem Behälter (oder einem anderen Behälter $\neq b_k$) zuordnen. Der Behälter b_k erhält jedenfalls am Ende gar kein Paket. Ein Widerspruch zu der Annahme, dass überhaupt beide Behälter benutzt wurden.

Jetzt zeigen wir unter dieser Annahme, die 2-Approximation. Sei dazu g das kleinste Gewicht eines Behälters nach der Ausführung des jeweiligen Algorithmus und sei k_{opt} die optimale Anzahl an Behältern und k die Anzahl der jeweils verwendeten Behälter.

Fall 1, $g \geq \frac{G}{2}$: Die Gesamtbeladung, GB, ist $\geq k \times \frac{G}{2}$. Für k_{opt} gilt stets $k_{\text{opt}} \geq \frac{\text{GB}}{G}$. Daraus ergibt sich:

$$k_{\text{opt}} \geq \frac{k \times G/2}{G} = \frac{k}{2}.$$

Fall 2, $g < \frac{G}{2}$, $k = 1$: Dann gilt $k_{\text{opt}} = k$.

Fall 3, $g < \frac{G}{2}$, $k > 1$: Gemäß der obigen Aussage kann es nur einen einzigen Behälter mit Füllmenge $g < \frac{G}{2}$ geben. Dann sind $(k - 1)$ Behälter mindestens bis $\frac{G}{2}$ gefüllt. Die Gesamtbeladung, GB, ist echt größer als $(k - 1) \times \frac{G}{2}$. Für k_{opt} gilt dann $k_{\text{opt}} \geq \frac{\text{GB}}{G}$.

Daraus ergibt sich: $k_{\text{opt}} > \frac{(k-1) \times G/2}{G}$, woraus

$$k_{\text{opt}} \geq \frac{k}{2}$$

wegen der Ganzzahligkeit gefolgert werden kann.

□

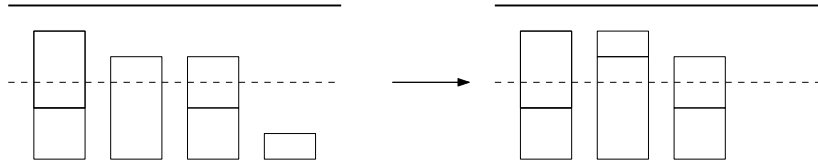


Abbildung 5.2: Nach Anwendung von First-Fit und Best-Fit haben je zwei Behälter ein Gesamtgewicht von $G + 1$.

Beide Algorithmen haben also eine beweisbare Güte. Sie sind in jedem Fall nicht schlechter als zweimal das Optimum.

Bemerkungen:

- Der obige Beweis läßt sich auch für nicht-ganzzahlige Gewichte führen. Die Hilfsaussage im Beweis von Lemma 6 lautet dann: *Für je zwei benutzte Behälter gilt, dass das Gesamtgewicht der enthaltenen Gegenstände mindestens $G + \epsilon$ für ein $\epsilon > 0$ ist.* Übungsaufgabe: Führen Sie den Beweis insgesamt für nicht-ganzzahlige Gewichte.
- Für beide Algorithmen läßt sich durch genauere Analyse ein Approximationsfaktor von $C \leq \frac{17}{10}$ beweisen.

Wir wollen abschließend auch noch eine untere Schranke für den Approximationsfaktor der beiden Algorithmen angeben. Das heißt, wir konstruieren eine Sequenz, die zu einem bestimmten Faktor führt. Wir wollen dabei auch zeigen, dass das Ergebnis für Problemstellungen gilt, die beliebig viele Behälter benötigen.

Lemma 7 *Für das Binpacking Problem und die Algorithmen First-Fit und Best-Fit gibt es Beispiel-Sequenzen, die beliebig lang sein können und für die die Lösung aus den beiden Algorithmen niemals besser sein kann als $\frac{5}{3}$ mal die optimale Lösung.*

Beweis. Sei dazu $n = 18m$ und $G = 131$. Gegenstände mit Gewicht 20 plus Gewicht 45 plus Gewicht 66 füllen genau einen Behälter.

Insgesamt betrachten wir $18m$ Gegenstände, wobei in der Sequenz der Gegenstände

1. ... die ersten $6m$ ein Gewicht $g_i = 20$ für $i = 1, \dots, 6m$,
2. ... die nächsten $6m$ ein Gewicht von $g_i = 45$ für $i = 6m + 1, \dots, 12m$,
3. ... und die letzten $6m$ ein Gewicht von $g_i = 66$ für $i = 12m + 1, \dots, 18m$

besitzen. Wir haben also drei Sequenzblöcke 1., 2. und 3. und exakt $6m$ Behälter reichen aus.

Wir zeigen, dass für diese Sequenz First-Fit und Best-Fit jeweils $10m$ Behälter befüllen. Beide Algorithmen füllen für den ersten Sequenzblock sukzessive die ersten m Behälter mit

Gewichten $6 \times 20 = 120$. Jeder Block hat eine Restkapazität von 11. Danach werden für den zweiten Block $3m$ Behälter mit jeweils 2×45 Gewicht und Restkapazität von 41. Danach werden $6m$ Behälter mit Gewicht 66 gefüllt und Restkapazität jeweils 65.

Insgesamt ergibt sich in beiden Fällen ein Quotient von $\frac{10m}{m} = 10$. \square

Wenn die Sequenzen aus dem obigen Beweis in absteigend sortierter Reihenfolge vorgelegen hätten, dann hätten beide Algorithmen die optimale Belegung mit $6m$ Behältern erzielt.

5.3 Aktivitätenauswahl

Wir betrachten das Problem, dass eine Ressource für verschiedenen zeitlich beschränkte Aktivitäten benutzt werden soll. Die Aktivitäten können die Resource nicht zeitgleich benutzen. Es sollen aber soviel wie möglich Aktivitäten eingeplant werden. Beispielsweise soll die Belegung eines Konzertsaals für verschiedene Veranstaltungen geplant werden. Allgemein gibt es eine Menge von n Aktivitäten $a_i = (s_i, t_i)$ die jeweils durch einen Starttermin, s_i , und einen Endtermin, t_i , mit $s_i < t_i$ bestimmt sind.

Zwei Aktivitäten a_i und a_j sind kompatibel, wenn die Intervalle $[s_i, t_i)$ und $[s_j, t_j)$ nicht überlappen.

Aufgabenstellung Aktivitäten-Auswahl: Für eine Menge von n Aktivitäten $a_i = (s_i, t_i)$ für $i = 1, \dots, n$ finde die maximale Menge gegenseitig kompatibler Aktivitäten. Mit maximal ist hier die Anzahl der einplanbaren Aktivitäten gemeint.

Zunächst ist klar, dass wir die Aktivitäten zwischen dem kleinsten Startzeitpunkt $S = \min\{s_i | i = 1, \dots, n\}$ und dem größten Endzeit $T = \max\{t_i | i = 1, \dots, n\}$ einplanen müssen. Eine naheliegende Greedy-Idee ist, die bislang nicht verwendete Zeit zu maximieren. Deshalb sortieren wir die Aktivitäten nach den Endzeitpunkten t_i . Nehmen wir an, dass wir die Folge a_i bereits so vorliegen haben, das heißt $t_1 \leq t_2 \leq \dots \leq t_n$.

Greedy-Strategie: Wähle stets die Aktivität, die den frühesten Endzeitpunkt hat und noch legal eingeplant werden kann.

Beispiel:

i	1	2	3	4	5	6	7	8	9	10	11
s_i	1	3	0	5	3	5	6	8	8	2	12
t_i	4	5	6	7	8	9	10	11	12	13	14

Hier sind zum Beispiel $\{a_3, a_9, a_{11}\}$ kompatibel, $\{a_1, a_4, a_8, a_{11}\}$ ist eine maximale kompatible Teilmenge, die durch die Greedy Strategie ausgewählt wird. Auch $\{a_2, a_4, a_9, a_{11}\}$ ist optimal bezüglich der Anzahl der einplanbaren Aktivitäten.

Algorithmisch beschreiben wir die Greedy-Strategie ganz einfach wie in Algorithmus 9.

Algorithm 9 AktivitätenGreedy($a_i = (s_i, t_i)$ nach t_i sortiert)

```

1:  $A_1 := \{a_1\}$ ; last := 1;
2: for  $i = 2$  to  $n$  do
3:   if  $s_i \geq t_{\text{last}}$  then
4:      $A_i := A_{i-1} \cup \{a_i\}$ ; last :=  $i$ ;
5:   else
6:      $A_i := A_{i-1}$ ;
7:   end if
8: end for
9: RETURN  $A_n$ 

```

Laufzeit: Es müssen lediglich vorab die Endzeitpunkte sortiert werden. Das kann in $O(n \log n)$ durchgeführt werden. Danach wird die Liste in $O(n)$ Zeit in der FOR-Schleife abgearbeitet.

Korrektheit: Der Algorithmus ist einfach und liefert eine Lösung, das ist ein Grundprinzip der Greedy-Algorithmen. Die Analyse, dass die Lösung optimal ist, kann manchmal sehr schwer sein. Im Folgenden versuchen wir die Optimalität formal zu beweisen.

Theorem 8 Für das Aktivitäten-Auswahl-Problem berechnet der Greedy-Algorithmus 9 eine optimale Lösung.

Beweis. Wir zeigen per Induktion, dass stets eine optimale Lösung A^* existiert, so dass für jedes i gilt, dass $A^* \cap \{a_1, a_2, \dots, a_i\} = A_i$ ist. Sei $A^* = \{a_{i_1}, a_{i_2}, \dots, a_{i_k}\}$ mit $i_1 < i_2 < \dots < i_k$.

Induktionsanfang: Für $A_1 = \{a_1\}$ gilt, dass a_{i_1} nicht vor a_1 endet. Also ist auch $A^{*'} = \{a_1, a_{i_2}, \dots, a_{i_k}\}$ eine optimale Lösung mit der gewünschten Eigenschaft.

Induktionsschritt: Nehmen wir an, die Aussage gilt bereits für $1, \dots, i-1$.

Fall 1: Falls a_i nicht kompatibel zu A_{i-1} ist, ist a_i nicht in A_i enthalten. Da es eine optimale Lösung A^* gibt, die alle Elemente aus A_{i-1} enthält, ist diese Lösung ebenfalls nicht kompatibel zu a_i (enthält also a_i nicht) und die Aussage gilt.

Fall 2: Falls a_i kompatibel zu A_{i-1} ist und somit in A_i vorhanden ist, sei nun $B := A^* \setminus A_{i-1} = \{b_1, b_2, \dots, b_j\}$, wobei die b_i 's nach den Endzeitpunkten sortiert sind. A^* erfülle dabei die Bedingung für A_{i-1} , das heißt, A^* ist eine optimale Lösung, so dass $A^* \cap \{a_1, a_2, \dots, a_{i-1}\} = A_{i-1}$ gilt. Nun gilt, dass b_1 nicht vor a_i enden kann und außerdem kompatibel zu A_{i-1} ist. Dann können wir b_1 durch a_i ersetzen. Es liegt eine Situation wie in Abbildung 5.3 vor. Somit ist auch $A^* = A_{i-1} \cup \{a_i\} \cup \{b_2, \dots, b_j\}$ eine optimale Lösung mit der angegebenen Bedingung für A_i .

Für $i = n$ erhalten wir somit durch A_n eine optimale Lösung, da es eine optimale Lösung A^* gibt, die genau aus den Elementen von A_n besteht. \square

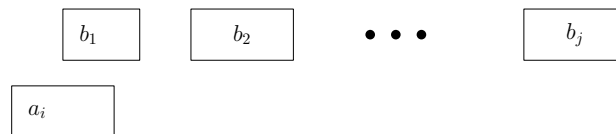


Abbildung 5.3: Nach Annahme gilt $A^* \cap \{a_1, a_2, \dots, a_{i-1}\} = A_{i-1}$. Falls a_i kompatibel mit A_{i-1} ist, kann b_1 in $B := A^* \setminus A_{i-1} = \{b_1, b_2, \dots, b_j\}$ durch a_i ersetzt werden.

Insgesamt haben wir drei Greedy Beispiele kennengelernt, in denen entweder eine optimale Lösung, eine gute Approximation oder eine beliebig schlechte Lösung erzielt wird.

5.4 Firefighting in Bäumen

Wir betrachten abschließend noch ein weiteres Problem, das im allgemeinen eine schweres Problem darstellt aber durch einen

Für einen Baum T mit Wurzel r betrachten wir das folgende dynamische Modell. An der Wurzel r des Baumes bricht ein *Feuer* aus, das sich in einem Zeitschritt auf die benachbarten Knoten ausdehnen kann. In jedem weiteren Zeitschritt breitet sich das Feuer von jedem brennenden Knoten auf nicht-brennende Nachbarknoten aus.

Wir verweisen hier kurz bereits auf das Kapitel 7 und Abschnitt 7.3. Jeder Knoten v des Baumes definiert einen eigenen Teilbaum T_v . Die *Tiefe* eines Knotens v ist die Länge des Weges von der Wurzel r zum Knoten v , also die Anzahl der Kanten, die dazwischen liegen.

Begleitend darf ein Feuerbekämpfer *vor* jedem Ausbreitungsschritt genau einen Knoten außerhalb des *Feuers* blockieren, so dass das Feuer sich in diesen Teil des Baumes nicht weiter ausbreiten kann. Der Agent soll nun eine Strategie angeben, die die maximal mögliche Anzahl an Knoten rettet.

Aufgabe: Finde für einen gegebenen Baum T mit ausgewiesener Wurzel r , diejenige Strategie, die in der obigen Konfiguration die maximale Anzahl an Knoten sichert.

Eine Greedy-Strategie würde in diesem Fall stets einen unverbrannten Knoten v blockieren, dessen Teilbaum T_v momentan die maximale Anzahl an Knoten aufweist. Diese Strategie ist niemals besser als eine 2-Approximation, wie bereits die folgende Abbildung beweist.

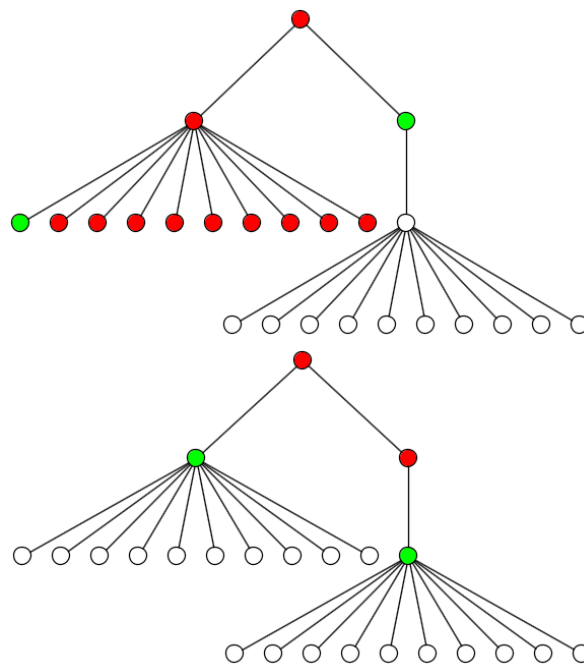


Abbildung 5.4: Stets den maximalen Teilbaum zu blockieren liefert eine 2-Approximation. Die Schranke wird angenommen.

Sei $n = 2k + 4$ die Anzahl der Knoten. In einer generalisierten Form des Baumes aus Abbildung 5.4 hat der rechte Teilbaum $k + 2$ Knoten und der linke Teilbaum $k + 1$ Knoten. Folglich sichert Greedy den rechten Teilbaum und kann im nächsten Schritt nur noch einen Knoten sichern. Somit bleiben $k + 3$ Knoten unverbrannt. Optimal wäre es gewesen, den linken Teilbaum der Größe $k + 1$ zuerst zu sichern und dann einen weiteren Teilbaum der Größe $k + 1$. Das Verhältnis $\frac{2k+2}{k+3}$ nähert sich beliebig nah an den Faktor 2 an. Wir können aber auch zeigen, dass Greedy niemals schlechter ist.

Theorem 9 *Für die Feuerbekämpfung in Bäumen gibt es Beispiele, in denen die Greedy-Strategie (asymptotisch) nicht besser ist als eine 2-Approximation. Die Greedy-Strategie ist aber niemals schlechter als eine 2-Approximation.*

Beweis. Den ersten Teil haben wir bereits oben gesehen. Für den zweiten Teil der Aussage betrachten wir folgende Argumentation. Wir teilen die Greedy-Schritte in zwei Kategorien A

und B ein, in A liegen alle Schritte, die mindestens so gut sind, wie die optimale Lösung in diesem Schritt und in B liegen alle Greedy-Schritte, die schlechter sind.

Jeder Schritt einer Strategie schneidet einen Teilbaum T ab. In diesem Teilbaum wird danach kein Knoten mehr gesichert. Außerdem können wir annehmen, dass wir die Teilbäume bei jeder Strategie stets direkt am Feuer sichern, sonst können wir jede Strategie leicht verbessern. Zudem wird somit im i -ten Schritt ein Knoten mit Distanz i zur Wurzel gesichert, solange das noch möglich ist.

Sei nun S_G die gesamte Anzahl der Knoten, die Greedy sichert und seien opt_A und opt_B die Anzahl der Knoten, die die optimale Strategie in den Kategorien A und B sichert. Wir wollen $2S_G \geq \text{opt}_A + \text{opt}_B$ zeigen.

Greedy schützt im ersten Schritt mindestens soviele Knoten wie die optimale Strategie. Das ist ja gerade das Konzept von Greedy. Somit liegt dieser Schritt in A . Falls nun irgendwann ein Schritt der optimalen Strategie einen Knoten v gewählt hat, der besser ist als Greedy, dann hat die Greedy-Strategie aber vorher irgendwann einen Vorgänger von v gewählt. Der Knoten v wäre ja sonst eine Alternative für Greedy gewesen und kann noch nicht unter Feuer stehen. Das heißt, das es für alle gesicherten Knoten aus opt_A und aus opt_B auch entsprechende gesicherte Knoten in S_G gibt. \square

Das gerade beschriebene Problem ist ebenfalls ein schweres Problem, wie wir im zweiten Teil der Vorlesung feststellen werden. Die im obigen Beweis festgestellten strukturellen Eigenschaften helfen aber, das Problem einigermaßen effizient optimal zu lösen. Dazu eignet sich die dynamische Programmierung.

5.4.1 Parametereffiziente Lösung

Wir wollen eine effiziente Lösung vorstellen, die effizient für eine beschränkte Anzahl von Agenten ist. Die Lösung basiert hier wiederum auf dynamische Programmierung, siehe Kapitel 4. Wir werden hier ein Beispiel sehen, in dem die Reihenfolge zur Memoisation (Bottom-Up) keine große Rolle spielt, da im wesentlichen alle erlaubten Konfigurationen betrachtet werden.

Zunächst betrachten wir das Problem Top-Down. Wir hatten bereits festgestellt: In Schritt j wird ein Knoten in Abstand j von der Wurzel gewählt, dessen Teilbaum gesichert wird. Nehmen wir an, dass wir k Agenten zur Verfügung haben. Falls im ersten Schritt ein Teilbaum (oder Knoten) in Abstand (*Tiefe*) 1 von der Wurzel ausgesucht und mit einem Agenten blockiert wurde, können wir diesen Teilbaum streichen und haben dann ein neues Teilproblem. Im verbleibenden Baum dürfen noch $k - 1$ Agenten verwendet werden, aber in Abstand 1 von der Wurzel (bzw. Knoten der Tiefe 1) darf dabei nichts mehr blockiert werden. In jeder Tiefe wird also genau ein Knoten blockiert.

Aus diesem Grund verwenden wir im allgemeinen einen Binärvektor $X \in \{0, 1\}^k$ mit k -Einträgen, der uns für ein Teilproblem sagt, in welchen Tiefen wir aktuell (rekursiv) noch Agenten setzen dürfen. $X(j) = 1$ bedeutet, dass ein Agent in Ebene j gesetzt werden darf, $X(j) = 0$ bedeutet, dass kein Agent in Ebene j gesetzt werden darf. Wenn wir nun in Schritt j in *Tiefe* j einen Knoten v auswählen und den Teilbaum T_v des Knoten dadurch beschützen, dann werden $|T_v|$ (Größe der Teilbaums) viele Knoten geschützt.

Durch X^v beschreiben wir dann einen Vektor, der den j -ten Eintrag von X auf 0 setzt, falls v die Tiefe j hat. Wenn wir das Problem rekursiv an einem Knoten v lösen wollen, müssen wir die Alternativen, dass der Teilbaum blockiert wird oder nicht blockiert wird betrachten. Zudem müssen wir beachten, dass beim Blockieren des Knoten v , der Pfad von der Wurzel zu v nicht blockiert werden darf. Außerdem müssen wir uns anschauen, wo wir nach der jeweiligen Bearbeitung des Knoten v weitermachen.

Bearbeitungsreihenfolge: Für eine Bearbeitungsreihenfolge betrachten wir einen sogenannten Preorder-Durchlauf des Baumes. Beginnend mit Ziffer 1 von der Wurzel vergeben wir aufsteigende numerische Labels an die Knoten und gehen dabei zunächst soweit in die Tiefe, bis es nicht mehr weitergeht. Danach gehen wir den Pfad zurück (Backtracking), bis wir wieder ungelabelte Nachbarn in größerer Tiefe entdecken. Das wird rekursiv durchgeführt, bis wir am Ende wieder an der Wurzel angekommen sind. Die Nummern werden fortlaufend vergeben. Eine solche Vorgehensweise nennt man auch Depth-First-Search (DFS) Durchlauf. Preorder bedeutet, dass bei diesem Durchlauf der nächste erstmals besuchte Knoten die nächste Nummer erhält. Die Abbildung 5.5 zeigt eine Preorder eines Baumes. Für jeden Knoten v ist $\text{pre}(v)$ der Preorder-Wert des Knoten v .

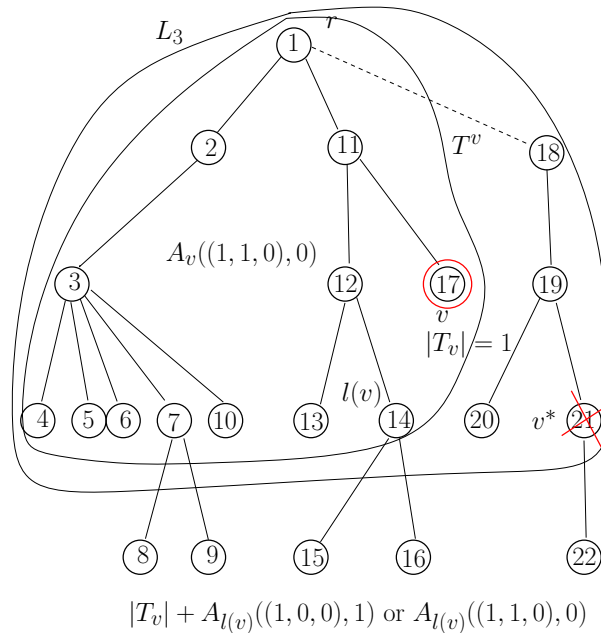


Abbildung 5.5: Die Preorder eines Baumes T und die Menge L_k der Knoten der Tiefe $\leq k$. Die Berechnung von $A_v((1, 1, 0), 0)$ bedeutet, dass wir eine Strategie im Baum T^v suchen, die Knoten in Tiefe 1 und 2 blockiert. Falls ein Knoten entlang des Pfades von r zu v blockiert werden soll, dann mit Tiefe größer als 0. Falls Knoten v blockiert wird, müssen wir $|T_v| + A_{l(v)}((1, 0, 0), 1)$ berechnen, wobei der zweite Parameter 1 sagt, dass ein Vorgänger von v nicht mehr blockiert werden soll. Falls Knoten v nicht blockiert wird, können wir $A_{l(v)}((1, 1, 0), 0)$ berechnen. Das Maximum dieser Werte ist dann die beste Wahl.

Wir wollen für k Agenten den Algorithmus beschreiben. Sei nun L_k die Menge aller Knoten von T mit Tiefe $\leq k$ und sei T^v der Teilbaum von T der alle Knoten von L_k (und deren Teilbäume) enthält mit Preorder kleiner-gleich $\text{pre}(v)$; siehe Abbildung 5.5. Sei $l(v)$ der Knoten mit der Preorder $\text{pre}(v) - 1$.

Wir wollen das Maximum, $A_v(X)$, der beschützbaaren Knoten in T^v bestimmen, falls wir eine Strategie mit k Agenten in T^v starten, die gemäß des Vektors X Knoten beschützt. Offensichtlich ist dann der optimale Wert $A_{v^*}(1^k)$ für den Knoten v^* mit der höchsten Preorder in L_k .

Betrachte Abbildung 5.5: Falls Knoten v in der Tiefe j geschützt wird, wird der Teilbaum T_v der Größe $|T_v|$ zur Lösung hinzugefügt. Der Teilbaum wird dann in Gedanken abgeschnitten, in dem wir mit dem Baum $T^{l(v)}$ fortfahren. Im Vektor X muss dann der j -te Eintrag auf 0 gesetzt werden. Der Vektor X^v ist der entsprechende Vektor. Also müssen wir $A_{l(v)}(X^v) + |T_v|$

berechnen. Falls Knoten v in der Tiefe j nicht geschützt wird, können wir mit $A_{l(v)}(X)$ fortfahren.

Ein kleines Problem besteht darin, dass wir auf jedem Pfad nur einen Knoten blockieren wollen und damit die Berechnung sogar vereinfachen können. Dazu merken wir uns für den aktuellen Knoten v , ab welcher Tiefe auf dem Weg zu v wir denn überhaupt Knoten blockieren dürfen. $A_v(X, i)$ besagt dann, dass wir hier aktuell auf dem Weg von r nach v nur Knoten nach Tiefe i blockieren dürfen. Im obigen Fall hätten wir $A_{l(v)}(X^v, \text{Tiefe}(v) - 1)$ aufrufen müssen, wenn wir die Variante betrachten, dass v blockiert wird.

Wir nehmen an, dass wir im allgemeinen rekursiv am Knoten v den Wert $A_v(X, i)$ bestimmen wollen und der obige erste Fall (Blockieren von v und Aufruf $A_{l(v)}(X^v, \text{Tiefe}(v) - 1)$) kommt dann nur vor, falls $X(\text{Tiefe}(v)) = 1$ und $\text{Tiefe}(v) > i$ gilt. Wenn wir zum Beispiel auf einem Pfad von der Wurzel zu einem Blatt v einen Knoten in Tiefe i blockieren, sollen die Vorgänger in diesem Zweig im rekursiven Aufruf nicht mehr blockiert werden dürfen. Das verringert die Anzahl der rekursiven Aufrufe.

Wenn der andere Fall (v wird nicht blockiert) betrachtet wird, ändern wir wie gehabt X nicht, und wechseln zu $A_{l(v)}$ für den Preorder-Vorgänger $l(v)$ von v in L_k . Der direkte Vorgänger von v der Tiefe $(\text{Tiefe}(v) - 1)$ liegt stets auf dem Pfad von r zu $l(v)$. Falls es nur erlaubt war, einen Knoten auf dem Pfad von r nach v ab der Tiefe i zu setzen, müssen wir das jetzt nur noch mit $\text{Tiefe}(v) - 1$ vergleichen. Das heißt, $A_{l(v)}(X, \min(\text{Tiefe}(v) - 1, i))$ wird berechnet.

Die allgemeine **Rekursionsgleichung** der dynamischen Programmierung lautet somit:

$$A_v(X, i) = \max \left\{ \begin{array}{l} A_{l(v)}(X, \min(\text{Tiefe}(v) - 1, i)) \\ [X(\text{Tiefe}(v)) = 1 \text{ und } \text{Tiefe}(v) > i] \cdot (|T_v| + A_{l(v)}(X^v, \text{Tiefe}(v) - 1)) \end{array} \right\}$$

wobei $[\phi]$ gleich 1 ist, wenn der Ausdruck ϕ gilt und 0 sonst. Bei der Bottom-Up Berechnung dieser Formel, berechnen wir rekursiv für fallende Preorder-Werte die entsprechenden Werte, deshalb muss jede dieser Berechnungen mit der Wurzel r enden. Mit $A_r(X, 0) := 0$ beenden wir also jeden rekursiven Durchlauf. In welcher Reihenfolge die Auswertungen vorgenommen werden ist uns dabei egal. Wir müssen sowieso alle Alternativen überprüfen. Eine geschickte tabellarische Auflistung der Aufrufe entfällt.

Theorem 10 *Die optimale Protektionsstrategie für k Agenten in einem Baum T der Größe n kann in Zeit $O(n2^k k)$ berechnet werden.*

Beweis.

Gemäß unserer Interpretation gibt für $v \in L_k$ der Wert $A_v(X, i)$ rekursiv die maximale Anzahl beschützbarer Knoten einer Strategie im Baum T^v an. Die Strategie setzt dabei die Agenten nach den Einträgen in X und setzt Agenten auf dem Pfad von r nach v nur *nach* Tiefe i .

Dann suchen wir nun lediglich nach $A_{v^*}(1^k, 0)$, wobei v^* der am weitesten rechts liegende Knoten der Tiefe k in T ist.

Nach der obigen Beschreibung ist unsere Prozedur korrekt und terminiert und wir müssen lediglich feststellen, wieviele Aufrufe es insgesamt geben kann. In L_k traversieren wir die Knoten von rechts nach links nur einmal. Ein Aufruf $A_v(X, i)$ kann aber für jedes sinnvolle i und jeden Vektor X vorkommen, also haben wir maximal $n \times 2^k \times k$ Einträge $A_v(X, i)$, wobei n für v , 2^k für X und k für i steht. \square

Für kleine Werte von k haben wir somit eine lineare Laufzeit in der Anzahl der Knoten und Kanten des Baumes. Man spricht hier auch von Algorithmen die *fixed parameter tractable* sind,

da wir es mit einem insgesamt *schweren* Problem zu tun haben und die Laufzeit $f(k) \cdot p(n)$ für ein Polynom p . Hier ist n die Eingabegröße des Problems und k der Parameter.

Kapitel 6

Das Sweepverfahren

In den vorherigen Kapiteln haben wir an verschiedenen Stellen ein Meta-Verfahren angewendet, bei dem wir Listen durchlaufen und dabei Invarianten aufrecht erhalten haben, die uns am Ende des Durchlaufs das gewünschte Ergebnis liefern.

Beispielsweise haben wir beim Mergesort (Prozedur 3) in Kapitel 2 zwei sortierte(!) Listen (in Arrays gespeichert) parallel durchlaufen und dabei stets die aktuellen Listenköpfe der Größe nach verglichen. Sukzessive wurde das größte Element ausgewählt und der *Zeiger* auf den Listenkopf für diese Liste nach vorne bewegt. Im Ergebnisarray stand dann am Ende die vollständig sortierte Liste von Punkten. Die Laufzeit war einfach durch die Längen der Listen beschränkt, in jedem Schritt wurden zwei Elemente verglichen. Man kann sagen, dass wir beiden Listen von *links* nach *rechts ausgefegt* haben und lokal geringen Aufwand hatten. Ein solches Verfahren bezeichnen wir deshalb als *Sweep*.

Die intuitive Idee des Sweepverfahrens stammt eigentlich aus der Betrachtung geometrischer Eingaben bezüglich des *Ausfegens*. Ein gutes einfaches Beispiel zeigt der Greedy-Algorithmus 9 aus Abschnitt 5.3 zur Berechnung einer maximalen Aktivitätenauswahl. Zunächst wurden hier die Aktivitäten nach den Endzeitpunkten sortiert (Ereignisstruktur/Eventstruktur). Dann wurden die *Segmente* (Zeitintervalle) gemäß dieser Ordnung durchlaufen. Wir können uns das geometrisch so vorstellen, wie in Abbildung 6.1. Die Segmente liegen untereinander und wir betrachten die Elemente in einem Sweep mit einer *Sweepline* von links nach rechts. Zu jedem Zeitpunkt merken wir uns die momentan aktive aktuelle Aktivität als *Sweep-Status-Struktur* (SSS). Die Sweepline rückt zum nächsten *Event* vor und überprüft in konstanter Zeit, ob das nächste Segment mit dem aktuellen in Konflikt steht. Ist das nicht der Fall, wird ein neues Segment aktiv und auch als Teil des Ergebnisses ausgegeben.

Korrektheit: Die allgemeine Korrektheit eines Sweepverfahrens ergibt sich durch das Aufrechterhalten der Invariante. Das ist die Eigenschaft der Sweep-Status-Struktur in Zusammenhang mit den bisher berichteten Ergebnissen.

Laufzeit: Die Laufzeit des Verfahrens ergibt sich zunächst aus der Vorbereitung (Preprocessing) als Berechnung der Bearbeitungsreihenfolge und der beginnenden SSS. Danach wird der Aufwand der Events mit der Anzahl der Events multipliziert.

Geometrisch betrachtet haben wir es mit einer Reduktion der Dimension zu tun. Eigentlich ist zum Beispiel das obige Problem ein eindimensionales Problem bezüglich des Überlappens von Intervallen. Wir machen daraus entlang der Sweepline sukzessive ein 0-dimensionales Problem, weil wir in der Nähe der Sweepline nur konstanten Aufwand haben.

Wir fassen den Sweep als Verfahren wie folgt zusammen:

- Ausfegen der Ebene bzw. Bearbeiten der Ereignis entlang einer Sweepline

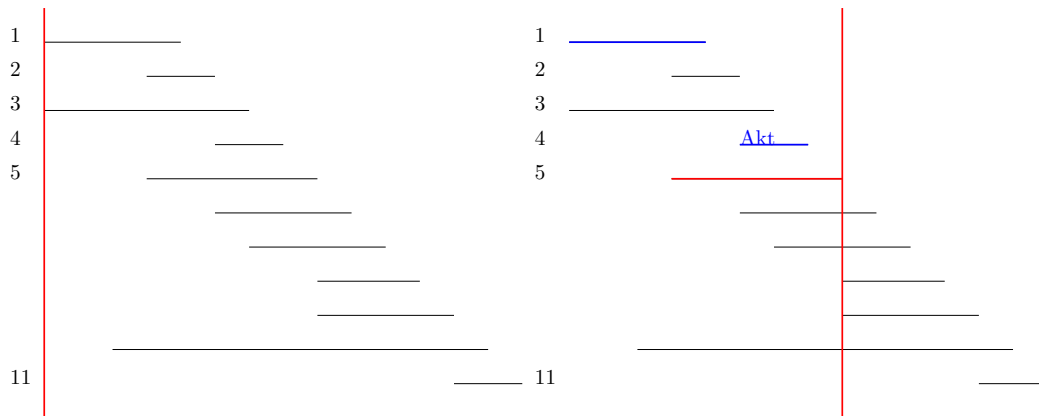


Abbildung 6.1: Die geometrische Umsetzung des Greedy-Algorithmus 9 zur Berechnung einer maximalen Anzahl an Aktivitäten. Zur Abarbeitung der Intervalle wurden diese nach den Endzeitpunkten sortiert. Wenn die Sweepline vorrückt, wird stets das aktuelle Segment mit dem aktuell aktiven Segment auf Überlappung überprüft. Eventuell wird das Ergebnis angepasst und die nächste Aktivität berichtet.

- Komplexitätsreduktion: Zum Beispiel von 2-dimensional auf 1-dimensional
- Sortieren der Ereignisse und gemäß der Sortierung bearbeiten
- SSS (Sweep-Status-Struktur): Invariante in der Nähe der Sweepline
- ES (Ereignisstruktur): Haltepunkte der Sweepline
- Ereignisverarbeitung: Aktualisierung der SSS und ggf. auch der Ereignisse
- Laufzeitanalyse: ($\#$ Ereignisse) \times Kosten(Verab)) plus die Kosten für die Berechnung der Ereignisstruktur und der Initialisierung der Sweep-Status-Struktur
- Korrektheit des Ergebnisses: Invariante ist erfüllt und es wird korrekt berichtet
- Einfache Beispiele: Punkte auf einer Linie, Closest Pair oder die maximale Anzahl nicht-überlappender Intervalle

Im obigen Beispiel sind die Events die Endpunkte der Intervalle und die SSS besteht aus dem aktiven Segment. Falls ein Segment erreicht wird, das sich mit dem aktiven Segment nicht überlappt wird die SSS angepasst und das Ergebnis berichtet.

6.0.2 Closest Pair von Punkten in der Ebene

Wir wollen das Sweepverfahren und seine zentralen Eigenschaften an einem zweidimensionalen Problem erläutern. Wir betrachten wie in Abschnitt 2.2 eine Menge S von n Punkten p_1, p_2, \dots, p_n in der Ebene, es soll das Punktepaar mit kleinstem Euklidischen Abstand $d(\cdot, \cdot)$ mittels Sweep berechnet werden. Wir suchen den Wert

$$d_S = \min_{1 \leq i, j \leq n, i \neq j} d(p_i, p_j).$$

Dazu sortieren wir die Punkte zunächst nach X -Koordinaten (Preprocessing) und bearbeiten die Punkte nach steigenden X -Koordinaten (Ereignisstruktur). Sei p_1, p_2, \dots, p_n die sortierte

Folge. Nachdem wir den Abstand $\text{MinSoFar} := d(p_1, p_2)$ zwischen den ersten beiden Punkten betrachten wir einen vertikalen Streifen $S(p_2)$ mit Breite MinSoFar und rechtem Rand p_2 .

Die Abbildung 6.2 zeigt die allgemeine Situation für einen Punkt p . Alle Punkte im Streifen $S(p)$ und auch links davon haben keinen kleineren Abstand als MinSoFar . Nun bewegt sich die Sweepline vorwärts, dabei können zwei Events auftreten:

- Der linke Rand des Streifens trifft zuerst auf einen Punkt r . Dieser Punkt verläßt den Streifen; Abbildung 6.2 ii).
- Der rechte Rand trifft auf einen neuen Punkt q ; Abbildung 6.2 iii). q betritt den Streifen, $S(q)$, und mit allen Punkten $w \in R(p)$ (maximal 11 viele) wird nun der Abstand getestet. Falls ein Abstand kleiner als MinSoFar ermittelt wird, dann schrumpft der Streifen zunächst. In diesem Fall könnten links weitere Punkte den Streifen verlassen. Zum Beispiel r' in Abbildung 6.2 iii).

Im Prinzip verwenden wir hier parallel zwei Sweeplines, die rechte und die linke Kante des momentanen Streifens. Ob zunächst die linke Kante den nächsten Punkt erreicht oder die rechte, können wir durch die Sortierung lokal entscheiden.

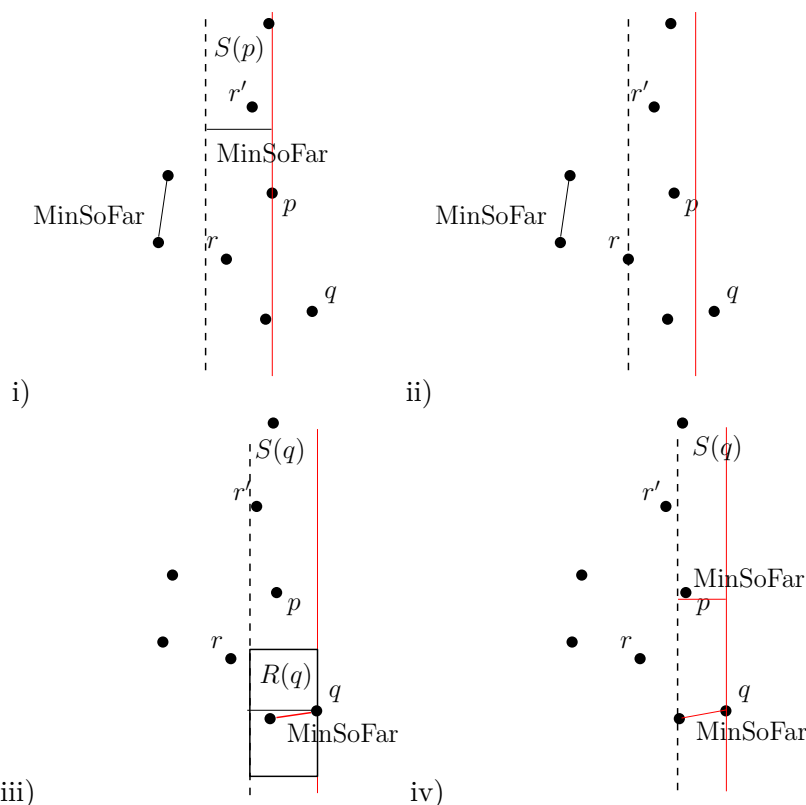


Abbildung 6.2: Beim Sweep wird in der SSS der Streifen der Breite MinSoFar aufrecht erhalten. Dabei verlassen Punkte den Streifen auf der linken Seite ii) oder neue Punkte betreten den Streifen iii). Im Fall iii) muss mit dem neuen Punkt getestet werden, ob sich MinSoFar verringert hat. Falls das der Fall ist, könnten zunächst weitere Punkte links den Streifen verlassen, wie r' in iv).

Offensichtlich betritt jeder Punkt maximal einmal den Streifen und verläßt einmal den Streifen. Nach der Sortierung ist das Anpassen der Streifen also in linearer Zeit möglich. Bei der

Aktualisierung des aktuell kleinsten Abstands wissen wir, dass wir nicht mit allen Punkten in $S(q)$ einen Test machen müssen sondern nur mit den Punkten aus $R(q)$; siehe Abbildung 6.2iii). Die Schwierigkeit besteht einzig darin, diese Punkte im Streifen $S(q)$ schnell zu finden. Falls wir alle durchgehen müssen, kann eine Laufzeit von $O(n^2)$ entstehen.

Was wir benötigen, ist eine effiziente Art der Speicherung der Punkte in $S(q)$, so dass wir eine *Bereichsanfrage* $R(q)$ für die Punkte in $S(q)$ machen können. Diese Datenstruktur muss sowohl das *Einfügen* als auch das *Entfernen* von Punkten effizient erlauben. Falls die Datenstruktur für eine Bereichsanfrage, das Einfügen und das Entfernen eine Laufzeit von jeweils $f(n)$ für n Punkte benötigt, haben wir eine Gesamtlaufzeit von $O(n \cdot f(n))$. Effiziente Datenstrukturen und der Gebrauch für Algorithmen behandeln wir im nächsten Abschnitt.

Kapitel 7

Datenstrukturen

Nachdem wir einige algorithmische Standardtechniken betrachtet haben, wollen wir uns nun mit der geschickten Repräsentation von Daten im Rechner beschäftigen und stellen ein paar einfache Datenstrukturen vor.

Die einfache Datenstruktur des Arrays haben wir bereits kennengelernt, für ein Array $A[1..n]$ können wir die Werte $A[i]$ in $O(1)$ Zeit auslesen oder ändern. Arrays stehen in allen gängigen Programmiersprachen zur Verfügung, über die interne Realisierung müssen wir uns zunächst keine Gedanken machen. Wie bereits gesehen, lassen sich die Arrays auch in mehreren Dimensionen verwenden.

Im Abschnitt 2 haben wir gesehen, wie ein Array A mit n Werten sortiert werden kann. Den direkten Zugriff über $A[i]$ kann man zum Beispiel nutzen, um effizient festzustellen, ob ein Element x im Array $A[1..n]$ vorhanden ist. Falls wir die Einträge eines sortierten Arrays A sukzessive mit x vergleichen, liegt die Laufzeit dieser Vorgehensweise in $O(n)$.

Das Prinzip der *binären Suche* erlaubt eine effizientere Lösung. Dabei wird der Wert x jeweils mit dem Wert des Arrays an der Stelle $\lfloor \frac{p+r}{2} \rfloor$ verglichen und je nach Ergebnis wird in den Teilbereichen weitergesucht. Algorithmus 10 beschreibt dieses Verfahren.

Algorithm 10 BinarySearch(A, p, r, x)

Ist x ein Element des aufsteigend sortierten Arrays A zwischen Index p und r ?

```
1: if  $p < r$  then
2:    $q := \lfloor (p + r) / 2 \rfloor$ ;
3:   if  $x < A[q]$  then
4:     BinarySearch( $A, p, q, x$ );
5:   else if  $x > A[q]$  then
6:     BinarySearch( $A, q + 1, r$ );
7:   else
8:     RETURN  $A[q] = x$ 
9:   end if
10: else if  $p = r$  then
11:   if  $A[p] = x$  then
12:     RETURN  $A[p] = x$ 
13:   else
14:     RETURN  $x \notin A$ 
15:   end if
16: end if
```

Laufzeitabschätzung: Dieser rekursive Algorithmus wird mit $\text{BinarySearch}(A, 1, n, x)$ aufgerufen und testet für $n = 1$ in konstanter Zeit, ob $A[1] = x$ ist. Sonst wird rekursiv ein Array der Größe $\lceil \frac{n}{2} \rceil$ betrachtet und konstanter Aufwand c für die konstante Anzahl an zusätzlichen Operationen verwendet.

Die zugehörige Rekursionsgleichung lautet $T(n) = T(\lceil \frac{n}{2} \rceil) + c$ und $T(n)$ liegt beweisbar in $O(\log n)$.

Korrektheit: Falls das Array nur einen Eintrag hat, testen wir, ob dieser dem Element x entspricht. Ansonsten betrachten wir das mittlere Element $q := \lfloor (p + r)/2 \rfloor$ des aufsteigend sortierten Arrays $A[p..r]$, falls $x < A[q]$ suchen wir im linken Teil, sonst im rechten Teil des Arrays. Falls beides nicht zutrifft, muss $x = A[q]$ gelten und wir geben das Element aus.

7.1 Stacks

Die oben beschriebenen Arrays lassen sich auch für die Beschreibung komplizierterer Strukturen verwenden. Wir betrachten hier einen sogenannten *Stack* auch Stapel oder Keller genannt.

Der Stack ist intern realisiert durch ein unendliches Array $S[1..\infty]$ mit einem ausgewiesenen Kopfelement $S[\text{Top}]$ für die Integer-Variable Top . Dabei ist $S[1], S[2], \dots, S[\text{Top}]$ der aktuelle Stack und $S[\text{Top}]$ das Kopfelement, der Stack ist leer, falls $\text{Top} = 0$ gilt.

Extern stehen dem Benutzer folgende Operationen zur Verfügung.

- $\text{Push}(x)$, ein neues Kopfelement x wird auf den Stack gelegt. Intern wird $\text{Top} := \text{Top} + 1$ und $S[\text{Top}] := x$ gesetzt
- $\text{Pop}(x)$, gibt das aktuelle Top-Element in x zurück und löscht es aus dem Stack. Intern wird $x := S[\text{Top}]$ und $\text{Top} := \text{Top} - 1$ gesetzt. Falls Top bereits 0 ist, wird berichtet, dass der Stack leer ist, extern beispielsweise über den Wert Nil.
- $\text{Top}(x)$ gibt das aktuelle Kopfelement in x zurück. Intern über $S[\text{Top}]$.

Wir wollen auch hier ein Beispiel für die effiziente Verwendung von Stacks angeben. Nehmen wir an, die Entwicklung verschiedener Aktien verläuft über einen bestimmten Zeitraum $I = [0, d]$ jeweils linear, das heißt beschrieben durch eine Funktion $f_i(x) = a_i x + b_i$ für $x \in I$. Nun möchte ein Analyst für jedes einzelne $x \in I$ schnell ermitteln, welche Aktie den größten Wert erzielt.

Naiv müsste man für den jeweiligen Wert x alle n Funktionswerte $f_i(x)$ miteinander vergleichen und könnte dann den maximalen Wert ermitteln.

Besser wäre es, gleich die sogenannte *obere Kontur* aller Funktionen explizit anzugeben. Zu jedem Zeitpunkt beschreibt eine der Funktionen f_i die momentan maximalen Werte. Wenn wir die Übergänge der jeweiligen *maximalen* Funktionen speichern könnten, können wir für die jeweiligen Teil-Intervalle schnell den optimalen Wert angeben.

Aufgabenstellung: Wir wollen die Funktion

$$f_{\max} : x \mapsto \max\{f_1(x), f_2(x), \dots, f_n(x)\}$$

explizit durch Intervalle und zugehörige Teil-Funktionen beschreiben. f_{\max} bezeichnen wir als die *obere Kontur* von f_1, f_2, \dots, f_n über dem Intervall $I = [0, d]$, siehe auch Abbildung 7.1.

Der Einfachheit halber nehmen wir an, dass die Werte b_i paarweise verschieden sind. Zunächst sortieren wir die Funktionen an der Stelle $x = 0$ nach den Y -Koordinaten b_i . Nehmen wir nun an, das o.B.d.A. $b_1 > b_2 > b_3 > \dots > b_n$ gilt.

Wir verwenden einen Stack, und der Stack soll am Ende sukzessive die maximalen Funktionen und die Übergänge enthalten.

Zunächst ist f_1 maximal und wir legen $(f_1, (0, b_1))$ auf den Stack als Kopfelement. Das ist die momentane obere Kontur der betrachteten Funktionen. Nun nehmen wir f_2 und betrachten den aktuellen Stack. Falls die beiden Funktionen einen Schnittpunkt $s(f_1, f_2) = (x_{f_1, f_2}, y_{f_1, f_2})$ besitzen, wechselt ab diesem Punkt das Maximum von f_1 zu f_2 . Wir legen somit f_2 mit $\text{Push}((f_2, s(f_1, f_2)))$ auf den Stack, der Stack repräsentiert die aktuelle obere Kontur. Das unterste Element f_1 startet bei $(0, b_1)$ und wird durch das nächste Element f_2 bei $s(f_1, f_2)$ abgewechselt.

Nun betrachten wir f_3 . Falls f_3 keinen Schnittpunkt mit f_2 bildet, kann f_3 ausgelassen werden, die Funktion liegt vollständig unterhalb von f_2 . Falls jedoch f_3 mit f_2 einen Schnittpunkt $s(f_2, f_3) = (x_{f_2, f_3}, y_{f_2, f_3})$ entscheidet die Lage des Schnittpunktes über die aktuelle Kontur der beteiligten Funktionen.

Fall 1.: Liegt der Punkt $s(f_2, f_3)$ mit seiner X -Koordinate x_{f_2, f_3} links vom X -Wert x_{f_1, f_2} des Startpunktes $s(f_1, f_2)$ des aktuellen Top-Elementes f_2 , so kann f_2 offensichtlich insgesamt gar nicht an der oberen Kontur teilnehmen. In diesem Fall wird f_2 durch $\text{Pop}(x)$ vom Stack genommen. Der gleiche Test wird mit dem nächsten Top-Element (hier f_1) wiederholt.

Fall 2.: Der Punkt $s(f_2, f_3)$ liegt mit seiner X -Koordinate x_{f_2, f_3} rechts vom X -Wert x_{f_1, f_2} des Startpunktes $s(f_1, f_2)$ des aktuellen Top-Elementes f_2 . Dann wird mit $\text{Push}((f_3, s(f_2, f_3)))$ der nächste Abschnitt der Kontur auf den Stack *gepusht* und der aktuelle Stack beschreibt tatsächlich die obere Kontur der beteiligten Funktionen mit den jeweiligen Übergängen.

In Abbildung 7.1 sehen wir, wie dieses Prinzip im Allgemeinen funktioniert. Nachdem bereits sukzessive $(f_1, (0, b_1))$, $(f_2, s(f_1, f_2))$, $(f_3, s(f_2, f_3))$, $(f_4, s(f_3, f_4))$ auf den Stack *gepusht* wurden, wird nun f_5 betrachtet. Wir betrachten jeweils die Lage des Schnittpunktes mit dem Top-Element. Da die Schnittpunkte $s(f_4, f_5)$ und $s(f_3, f_5)$ jeweils links von den Punkten $s(f_3, f_4)$ respektive $s(f_2, f_3)$ liegen, werden f_4 und f_3 nacheinander vom Stack *gepoppt*. Dann wird mit dem dann aktuellen Top-Element $(f_2, s(f_1, f_2))$ der Schnittpunkttest mit f_5 durchgeführt. Der Test der Schnittpunkte $s(f_2, f_5)$ und $s(f_1, f_2)$ ergibt, dass $s(f_2, f_5)$ rechts von $s(f_1, f_2)$ liegt und somit $(f_5, s(f_2, f_5))$ auf den Stack *gepusht* werden muss, $(f_2, s(f_1, f_2))$ bleibt erhalten. Der Stack gibt die aktuelle Kontur wieder.

Korrektheit: Wir wissen, dass sich zwei Geraden nur einmal schneiden können. Induktiv nehmen wir an, dass der aktuelle Stack die obere Kontur für die ersten f_1, f_2, \dots, f_{k-1} Geraden wiedergibt. Für das erste Element f_1 ist das der Fall.

Das aktuelle Top-Element $(f_j, s(f_i, f_j))$ des Stacks ist das letzte Stück der Kontur und wird mit der nächsten Geraden f_k verglichen.

Falls der Schnittpunkt $s(f_j, f_k)$ links von $s(f_i, f_j)$ liegt, kann f_j nicht zur oberen Kontur gehören. Rechts von $s(f_j, f_k)$ dominiert f_k die Gerade f_j , links von $s(f_i, f_j)$ dominiert die Gerade f_i die Gerade f_j . Die Push-Operation ist also in diesem Fall korrekt und das nächste Top-Element des Stacks wird mit f_k verglichen.

Falls der Schnittpunkt $s(f_j, f_k)$ rechts von $s(f_i, f_j)$ liegt, wird korrekterweise $(f_k, s(f_j, f_k))$ ans Ende der Kontur und in den Stack eingefügt.

Zu jedem Zeitpunkt ist die Kontur korrekt, am Ende gibt der Stack die Kontur mit den jeweiligen Intervallen korrekt wieder.

Laufzeitabschätzung: Gemäß der Beschreibung betrachten wir sukzessive die Geraden f_1, f_2, \dots, f_n und vergleichen diese mit dem jeweiligen Top-Element. Jede Gerade kann nur einmal auf den Stack *gepusht* werden wird aber auch nur einmal vom Stack *gepoppt*. Insgesamt haben wir somit nur $O(n)$ viele Push- und Pop-Operationen, das gleiche gilt dann für die Top-

verändern. Es gibt dazu (programmiersprachenspezifische) Befehle zum *allokieren* von neuen Speicherelementen mit Datenobjekt und zugehörigem Zeiger. Außerdem können die Zeiger entsprechend umgelenkt werden.

Soll in die obige Liste nach dem Wert 33 ein neues Element 17 eingefügt werden, so wird zunächst (programmiersprachenspezifisch) eine freie Adresse $A = 6$ für ein neues Datenobjekt mit noch nicht spezifiziertem Wert W_A und noch nicht spezifiziertem Zeiger Z_A erzeugt. Danach wird $W_A := 17$ eingetragen. Der Zeiger vom Element 33 wird auf den neuen Eintrag 6 gesetzt und der Zeiger Z_A verweist auf 10, also auf das nachfolgende Element 14. In Java läßt sich das beispielsweise durch *Listenobjekte* mit der entsprechenden Ausgestaltung realisieren. Für ein neues Element wird ein neues Listenobjekt angelegt.

Der Vorteil der Listen liegt in der einfachen dynamischen Veränderung, für das Einfügen und Löschen von Listenelementen ist jeweils nur konstanter Aufwand notwendig. Will man bei Arrays den Zusammenhang erhalten, dann sind Kopiervorgänge mit linearer Laufzeit für das Einfügen und Löschen notwendig.

7.3 Bäume und Suchbäume

Eine natürliche Verallgemeinerung von Listen sind Bäume. Formal können Bäume wie folgt über die Knoten rekursiv definiert. Sei dazu die Menge $V = \{v_1, v_2, \dots\}$ eine unendliche Menge von Knoten.

1. Jeder Knoten $v_i \in V$ ist ein Baum, v_i ist die Wurzel des Baumes.
2. Falls T_1, T_2, \dots, T_m mit $m \geq 1$ Bäume mit paarweise disjunkter Knotenmenge sind, dann ist auch $T = (v, T_1, T_2, \dots, T_m)$ ein Baum mit Wurzel v und T_i ist der i -te *direkte Teilbaum* von T . Die Zahl m bezeichnet den *Grad* des Knoten v .

Grafisch werden Bäume dadurch dargestellt, dass für jeden Knoten v , ein Knoten angelegt wird und die Verbindungen zu den direkten Teilbäumen durch Kanten repräsentiert werden; siehe Abbildung 7.3. Genauer handelt es sich hierbei um eine *grafische Realisation* des Baumes. Im Folgenden werden wir die grafische Realisation und die obige formale Definition als identisch betrachten. Der i -te *direkte Teilbaum* soll von links nach rechts betrachtet an der i -ten Stelle per Kante angefügt werden.

Klassische Bezeichnungen: Für einen Baum T mit Wurzel v und direkten Teilbäumen T_1, T_2, \dots, T_m mit Wurzeln w_1, w_2, \dots, w_m heißt w_i der i -te *Sohn von v* und v der *Vater von w_i* . Kurzschreibweisen $v = \text{Wurzel}(T)$, $v = \text{Vater}(w_i)$, $w_i = \text{Sohn}(v, i)$. Ein Knoten v ist *Nachfolger* eines Knoten w innerhalb eines Baumes T falls es eine Folge von Knoten (v_1, v_2, \dots, v_n) mit $v_1 = v$ und $v_n = w$, so dass $v_{i+1} = \text{Vater}(v_i)$ für $i = 1, 2, \dots, n-1$ gilt. Ein Knoten vom Grad 0 ist ein *Blatt*. Ein Knoten vom Grad > 0 ist ein *innerer Knoten*.

Die Tiefe eines Knoten $v \in T$ wird rekursiv wie folgt definiert.

$$\text{Tiefe}(v, T) = \begin{cases} 0 & \text{falls } v = \text{Wurzel}(T) \\ 1 + \text{Tiefe}(v, T_i) & \text{falls } T_i \text{ direkter Teilbaum} \\ & \text{von } T \text{ mit } v \in T_i \end{cases} \quad (7.1)$$

Die Höhe eines Baumes T wird durch

$$\text{Höhe}(T) = \max\{\text{Tiefe}(b, T) \mid b \text{ ist Blatt von } T\} \quad (7.2)$$

festgelegt

Die Abbildung 7.3 zeigt ein Beispiel eines Baumes der Höhe 4. Die Blätter v_{12} und v_{13} haben die maximale Tiefe 4.

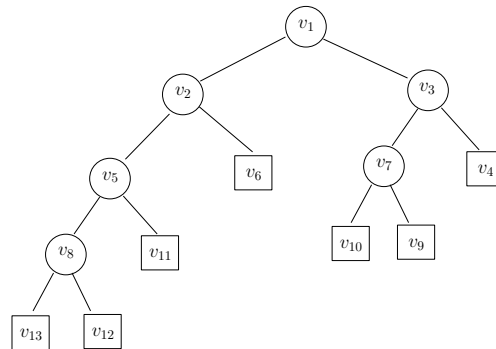


Abbildung 7.3: Die grafische Darstellung eines Baumes der Höhe 4.

Interne Realisation: Intern kann ein Baum analog zur verketteten Liste durch Zeiger realisiert werden. Ein Datenobjekt *Knoten* enthält dabei zum Beispiel einen Datensatz, einen Schlüssel (siehe unten) und eine Liste von Zeigern auf die Knoten der jeweiligen Söhne und ggf. einen Zeiger auf den Vaterknoten. Methoden für den Knoten erlauben es, auf den i -ten Sohn zuzugreifen, die Daten auszugeben oder die Zeiger zu verändern. Blätter werden als spezielle Knoten realisiert (zum Beispiel durch Nil-Zeiger, wenn keine Information in den Blättern gebraucht wird).

7.3.1 Binärbäume und Suchbäume

Bäume deren innere Knoten alle den Grad 2 haben werden Binärbäume genannt. Jeder innere Knoten hat dann einen eindeutigen linken und einen eindeutigen rechten direkten Teilbaum gemäß der eindeutigen grafischen Realisation. In den Knoten sollen außerdem Schlüssel abgespeichert werden. Über diese Schlüssel soll dann auf bestimmte Daten zugegriffen werden können. Je nachdem, ob sich die Schlüssel in den inneren Knoten oder nur in den Blättern befinden, gibt es die folgende allgemeine Unterteilung:

Art der Informationsspeicherung:

Suchbaum: Die Schlüssel befinden sich nur an den inneren Knoten, die Blätter sind leer (klassisch).

Blattsuchbaum: Die Schlüssel befinden sich nur in den Blättern, in den inneren Knoten gibt es nur *Wegweiser* zu den Knoten. Blattsuchbäume erlauben beispielsweise effiziente Bereichsanfragen.

Abbildung 7.4 zeigt einen Suchbaum (I) und einen Blattsuchbaum (II) für die Schlüssel (1, 3, 5, 7, 12, 15).

Im Folgenden werden wir zunächst Suchbäume betrachten, in den Blättern gibt es also keine Informationen. Die inneren Knoten, die Blätter als Nachfolger haben verweisen intern auf Nil. Insgesamt ergeben sich nun Möglichkeiten, wie man bei einem Durchlauf (Besuch aller Knoten) durch den Baum die Ausgabe der Schlüssel oder Daten (oder die Besuchsreihenfolge der Knoten) einfach beschreiben kann.

Präorder: Besuche die Wurzel (gebe die Daten aus), besuche dann den linken Teilbaum rekursiv und danach den rechten Teilbaum rekursiv, kurz *wLR*.

Postorder: Besuche zuerst den linken Teilbaum rekursiv und danach den rechten Teilbaum rekursiv und danach die Wurzel (gebe die Daten aus), kurz *LRw*.

Lexikographische Ordnung (In-Order): Besuche zuerst den linken Teilbaum rekursiv und danach die Wurzel (gebe die Daten aus) und dann den rechten Teilbaum rekursiv, kurz *LwR*.

Wird beispielsweise der Baum in Abbildung 7.4(I) nach der lexikographischen Ordnung durchlaufen, erhalten wir die Schlüssel in der Reihenfolge (1, 3, 5, 7, 12, 15). Die Präorder-Ausgabe ergibt: (7, 5, 3, 1, 15, 12).

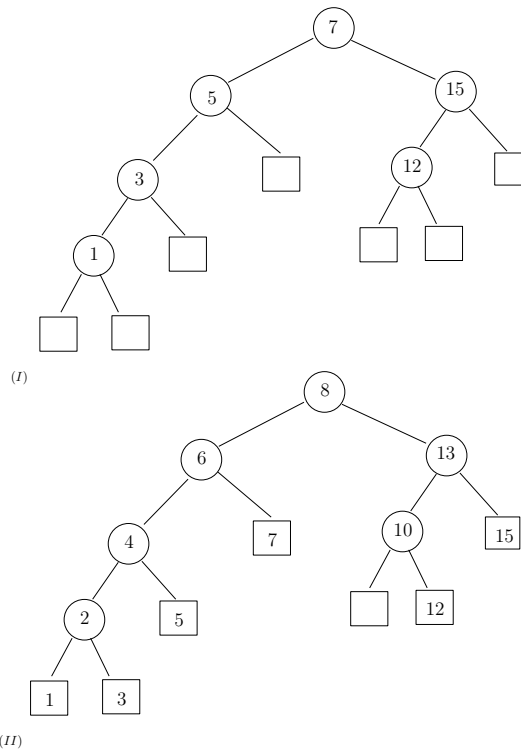


Abbildung 7.4: (I) Der Baum speichert Daten(Schlüssel) in den inneren Knoten, der Durchlauf nach *lexikographischer Ordnung LwR* gibt die Schlüssel sortiert in der Reihenfolge (1, 3, 5, 7, 12, 15) aus. (II) Ein Blattsuchbaum für die gleichen Schlüssel. Innere Knoten enthalten Wegweiser zu den Blättern.

Für einen Knoten v eines Baumes bezeichne $\text{Schlüssel}(v)$ den zugehörigen Schlüssel. Ein binärer Suchbaum hat die sogenannte *Suchbaumeigenschaft* falls für die Knoten und Schlüssel des Baumes gilt:

- Falls v ein Knoten des binären Suchbaumes ist und w ein Knoten im linken Teilbaum von v ist, dann gilt $\text{Schlüssel}(w) \leq \text{Schlüssel}(v)$.
- Falls v ein Knoten des binären Suchbaumes ist und w ein Knoten im rechten Teilbaum von v ist, dann gilt $\text{Schlüssel}(w) > \text{Schlüssel}(v)$.

In den Blattsuchbäumen erfüllen die Wegweiser zusammen mit den Schlüsseln der Blätter die Suchbaumeigenschaft. Suchbäume sollen es ermöglichen, möglichst effizient auf Daten über die Schlüssel zuzugreifen. Datenstrukturen sollen es insgesamt ermöglichen einfache Operationen effizient durchzuführen.

Ein klassisches Beispiel ist das **Wörterbuch-Problem** (Dictionary-Problem). Für eine Menge von Wörtern(Daten) wird dabei wie bereits oben erwähnt für jedes Wort ein Zahlenschlüssel verwendet, über den auf das Wort zugegriffen werden soll. Typische Operationen für dieses Problem sind:

Suche: Suche das Datenobjekt zu einem vorgegebenen Schlüssel x .

Einfügen: Füge ein neues Datenobjekt mit dem Schlüssel x ein.

Entfernen: Entferne ein Datenobjekt mit Schlüssel x .

Berichte Min/Max: Gebe das Datenobjekt mit kleinsten/größten Schlüssel aus.

Nächster Nachbar: Für einen Schlüssel m , finde ein Datenobjekt mit größten (kleinsten) Schlüssel $k \leq m$ ($k \geq m$).

Bereichsanfrage: Für ein Intervall $[l, r]$ finde alle Datenobjekte mit Schlüssel $k \in [l, r]$.

Im Folgenden beschreiben wir die Operationen jeweils für die Knoten eines Baumes mit zugehörigen Schlüssel. Dazu nehmen wir an, dass wir für einen Knoten v auf den linken Sohn bzw. rechten Sohn durch $\text{links}(v)$ respektive $\text{rechts}(v)$ zugreifen können. Blätter haben keine Nachfolger und können als Blatt identifiziert werden Die Operation Suchen kann nun wie folgt implementiert werden.

Algorithm 11 Search(T, k)

Suche im Baum T den Knoten v mit Schlüssel k oder gebe Blatt zurück, wenn kein solcher Knoten vorhanden ist

```

1:  $v := \text{Wurzel}(T)$ ;
2: while  $v \neq \text{Blatt}$  and  $k \neq \text{Schlüssel}(v)$  do
3:   if  $k < \text{Schlüssel}(v)$  then
4:      $v := \text{links}(v)$ ;
5:   else
6:      $v := \text{rechts}(v)$ ;
7:   end if
8: end while
9: RETURN  $v$ 

```

Laufzeit: In jedem Schritt der Ausführung der WHILE-Schleife steigen wir definitiv eine Stufe tiefer in den Baum ab, Die Laufzeit der Suchoperation hängt somit von der Höhe des Baumes ab. Für $h = \text{Höhe}(T)$ liegt die Laufzeit in $O(h)$. Für einen Baum mit n Knoten kann die Laufzeit sogar in $\Omega(n)$ liegen, falls der Baum entsprechend entartet ist.

Korrektheit: Wenn der Baum die Suchbaumeigenschaft erfüllt, dann zweigen wir in der WHILE-Schleife jeweils in den richtigen Teilbaum ab. Wenn wir auf ein Blatt stoßen ist $v = \text{Blatt}$ und die Schleife bricht ab. Falls der Schlüssel existiert gilt $k = \text{Schlüssel}(v)$, die Schleife bricht ab und v wird zurückgegeben.

Bei den Operationen Einfügen und Entfernen soll natürlich die Suchbaumeigenschaft erhalten bleiben. Wir stellen hier das Einfügen vor. Zunächst wurde dabei bereits ein Knoten v mit $\text{Schlüssel}(v) = k$ als Datenobjekt mit Verweisen $\text{links}(v) = \text{rechts}(v) = \text{Blatt}$ angelegt.

Korrektheit und Laufzeit: Auch hier müssen wir wieder sukzessive tiefer in den Baum absteigen, um das Element an der richtigen Stelle einzufügen. Der neue Knoten v ersetzt ein Blatt des Baumes T . In der Variablen y merken wir uns den inneren Knoten an dem wir jeweils verzeigen, über marker merken wir uns außerdem die Verzweigungsrichtung. Wenn wir

Algorithm 12 Insert(T, v)

 Füge v mit neuem Schlüssel(v) = k sortiert in nichtleeren Baum T ein

```

1:  $w := \text{Wurzel}(T)$ ;
2:  $y := \text{Blatt}$ ;
3: while  $w \neq \text{Blatt}$  do
4:    $y := w$ 
5:   if  $k < \text{Schlüssel}(w)$  then
6:      $w := \text{links}(w)$ ;  $\text{marker} := \text{links}$ ;
7:   else
8:      $w := \text{rechts}(w)$ ;  $\text{marker} := \text{rechts}$ ;
9:   end if
10: end while
11: if  $\text{marker} := \text{links}$  then
12:    $\text{links}(y) := v$ ;
13: else if  $\text{marker} := \text{rechts}$  then
14:    $\text{rechts}(y) := v$ ;
15: end if

```

am richtigen Blatt angekommen sind, können wir über y und die Verzweigungsrichtung den neuen Knoten einfügen. Die Laufzeit liegt wiederum in $O(h)$ wobei $h = \text{Höhe}(T)$ gilt.

Das Löschen eines Knoten v eines binären Suchbaumes ist nicht ganz so einfach und hängt davon ab, welche Söhne der jeweilige Knoten besitzt. Der Knoten wird wiederum über seinen Schlüssel in $O(h)$ Zeit gefunden. Wir unterscheiden folgende Fälle.

Entfernen eines Knoten v :

- Falls der Knoten v nur Blätter als Söhne hat, entfernen wir ihn einfach. Zum Beispiel beim Entfernen des Knoten mit Schlüssel 1 in der Abbildung 7.4(I).
- Falls der Knoten v ein Blatt als Sohn hat und einen Teilbaum T_1 mit Wurzel w , ersetzt der Knoten w den Knoten v , der Knoten v wird ausgeschnitten. Zum Beispiel beim Entfernen des Knoten mit Schlüssel 5 in der Abbildung 7.4(I).
- Falls der Knoten v zwei echte Teilbäume T_1 und T_2 als Unterbäume hat, wählen wir den Knoten w mit dem Nachfolgerschlüssel des Schlüssels von v aus. Dieser hat entweder zwei Blätter aber mindestens ein Blatt als Sohn. Der Nachfolger w wird ausgeschnitten und durch seinen Nachfolger wie in den ersten beiden Fällen ersetzt. Der Knoten v wird nun allerdings durch w ersetzt. Zum Beispiel beim Entfernen des Knoten mit Schlüssel 7 in der Abbildung 7.4(I), kann der Knoten mit Schlüssel 12 den Knoten mit Schlüssel 7 ersetzen.

In jedem der beschriebenen Fälle haben wir im von der Struktur her einen Teilbaum durch seinen linken oder rechten Unterbaum ersetzt, der andere Teilbaum war dabei stets ein Blatt. Wir mussten dabei den Knoten des zugehörigen Schlüssels finden und ggf. den Knoten des Nachfolgerschlüssels. Die Operation kann in $O(h)$ ausgeführt werden.

Übungsaufgabe: Falls ein Knoten v mit Schlüssel k in einem Suchbaum zwei echte Unterbäume hat, dann hat der Knoten mit dem k nachfolgenden Schlüssel mindestens ein Blatt als Unterbaum.

Für Binärbäume lassen sich viele interessante Eigenschaften ableiten bzw. beweisen. Exemplarisch beweisen wir hier, dass für einen Binärbaum mit n Blättern die Höhe in $\Omega(\log n)$

liegt. Außerdem hängt die Anzahl der inneren Knoten fest von der Anzahl der Blätter ab und umgekehrt.

Lemma 11

1. Ein Binärbaum mit n inneren Knoten hat $n + 1$ Blätter.
2. Ein Binärbaum mit n Blättern hat $n - 1$ innere Knoten.
3. In jedem Binärbaum gibt es maximal 2^i Knoten der Tiefe i .
4. Die Höhe eines Binärbaumes mit insgesamt n Knoten liegt in $\Omega(\log n)$.

Beweis. 1.: Zunächst beweisen wir induktiv die Aussage über die Anzahl der Blätter bei n inneren Knoten-

Induktionsanfang: $n = 0$. Ein Baum mit keinem inneren Knoten hat genau ein Blatt.

Induktionsschritt: Sei w die Wurzel eines Binärbaumes mit insgesamt $n \geq 1$ inneren Knoten. Dann teilt sich der Baum bei w in zwei Teilbäume T_1 und T_2 auf, die beide zusammen $n_1 + n_2 = n - 1$ innere Knoten besitzen. Nach Induktionsannahme gilt: T_1 hat genau $n_1 + 1$ und T_2 hat genau $n_2 + 1$ Blätter, insgesamt hat T genau $n_1 + n_2 + 2 = n + 1$ Blätter.

2.: Der Beweis kann analog zum obigen Beweis geführt werden, eine leichte Übungsaufgabe. Man kann aber auch einfach $n := n - 1$ setzen.

3.: Dass es in jedem Binärbaum nur maximal 2^i Knoten der Tiefe i geben kann, lässt sich leicht induktiv zeigen bzw. folgt direkt aus der Definition der Tiefe und ist eine leichte Übungsaufgabe.

4.: Nun betrachten wir die Höhe eines Binärbaumes T_n mit n Knoten. Wir zeigen direkt, dass $\text{Höhe}(T_n) \geq \lceil \log(n + 1) \rceil - 1$ für alle $n \geq 1$ gilt.

Die Höhe des Baumes ist die maximale Tiefe und deshalb wollen wir einen Binärbaum erstellen, dessen Knoten eine kleinstmögliche Tiefe haben. Nach obiger Aussage kann der Baum maximal $\sum_{i=0}^l 2^i = 2^{l+1} - 1$ Knoten der Tiefe $\leq l$ besitzen. Somit muss für n Knoten eine gewisse Mindesttiefe gelten, um alle Knoten unterzubringen.

Sei nun $k = \min_l \{l | 2^{l+1} - 1 \geq n\}$, offensichtlich ist k eine kleinstmögliche Höhe eines Baumes für n Knoten. Dann gilt offensichtlich $k \geq \lceil \log(n + 1) \rceil - 1$ denn aus $2^{k+1} - 1 \geq n$ folgt $2^{k+1} \geq n + 1$ und wenn wir darauf den Logarithmus anwenden, dann folgt wegen der Ganzzahligkeit von $k + 1$, dass $k + 1 \geq \lceil \log(n + 1) \rceil$ gilt. Kein Baum mit geringerer Höhe kann alle n Knoten unterbringen und die Aussage gilt. \square

Vorbereitend auf das nächste Kapitel wollen wir zeigen, dass wir in einem Suchbaum Knoten lokal umorganisieren können, ohne die Suchbaumeigenschaft zu verletzen. Die sogenannten Rotationen wollen wir dann im nächsten Kapitel für die Einhaltung von Höhenbalancen nutzen.

Linksrotation eines Knoten y mit Knoten x : Sei T ein Teilbaum mit folgender Eigenschaft. Die Wurzel x von T ist ein innerer Knoten mit linken Teilbaum A und rechten Teilbaum mit Wurzel y und seien B und C die Teilbäume (links und rechts) von y wie in der Abbildung 7.5. Dann können die Knoten y und x so nach *links* rotiert werden, dass y die Wurzel von T bildet mit rechten Teilbaum C . Der Knoten x wird linker Teilbaum von y mit linken und rechten Teilbäumen A und B .

Rechtsrotation eines Knoten x mit Knoten y : Symmetrisch zur Linksrotation, siehe Abbildung 7.5.

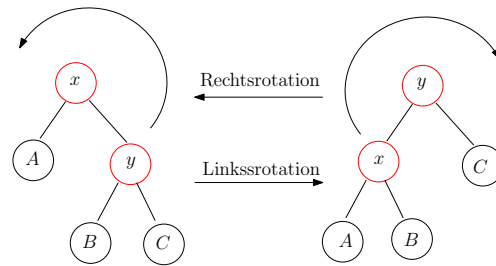


Abbildung 7.5: Bei einer Links- bzw. Rechtsrotation bleibt die Suchbaumeigenschaft erhalten.

Bei einer Linksrotation werden die Tiefen der Knoten in C um eins verringert, die Tiefen in B bleiben erhalten und die Tiefen in A werden um eins erhöht.

Bei einer Rechtsrotation werden die Tiefen der Knoten in A um eins verringert, die Tiefen in B bleiben erhalten und die Tiefen in C werden um eins erhöht.

Lemma 12 *Die Linksrotation erhält die Suchbaumeigenschaft des Teilbaumes T und des Gesamtbaumes. Der In-Order-Durchlauf gibt die gleiche Reihenfolge aus. Die Rotation kann mit konstanten Aufwand implementiert werden.*

Beweis. Dass die Rotation mit konstanten Aufwand implementiert werden kann, ergibt sich daraus, dass bei einer Implementierung mit Zeigern nur konstant viele Zeiger für die Rotation verändert werden müssen.

Die Schlüssel in A sind kleiner gleich Schlüssel $[x]$ und die Schlüssel in B und C sind größer gleich Schlüssel $[x]$. Es gilt Schlüssel $[x] \leq$ Schlüssel $[y]$. Schlüssel $[y]$ ist größer gleich den Schlüssel in B aber kleiner gleich den Schlüssel in C . Somit bleibt die Suchbaumeigenschaft insgesamt erhalten.

Der In-Order-Durchlauf beginnt jeweils mit einem In-Order-Durchlauf in A gibt dann x aus, führt einen In-Order-Durchlauf von B durch, gibt dann y aus und endet mit einem In-Order-Durchlauf von C . \square

Naheliegend ist folgende Eigenschaft. Übungsaufgabe: Ein binärer Baum erfüllt genau dann die Suchbaumeigenschaft, wenn der In-Order-Durchlauf die Schlüssel in aufsteigend sortierter Reihenfolge ausgibt.

Im Allgemeinen werden auch Bäume mit einem Grad von ≤ 2 für jeden inneren Knoten als Binärbäume bezeichnet. In diesem Kontext wird ein Binärbaum mit Grad exakt 2 für jeden inneren Knoten als *vollständig* bezeichnet. Insbesondere weil für die Suchbäume manchmal auch die eigentlichen Blätter weggelassen werden, findet man in der Literatur auch allgemeine Binärbäume. Wir werden stets die (leeren) Blätter mit notieren und somit vollständige Binärbäume betrachten.

7.4 AVL-Bäume

Wie wir bereits im letzten Abschnitt gesehen haben, ist die Höhe des Baumes entscheidend für elementare Operationen auf der Datenstruktur eines Binärbaumes. Das Lemma 11 zeigt uns allerdings bereits, dass eine bessere Laufzeit als $\Omega(\log n)$ für das Einfügen, Suchen und Entfernen nicht zu erwarten ist.

Ein binärer Suchbaum heißt AVL-Baum, falls für *jeden* Knoten v gilt, dass die Höhe der beiden Teilbäume von v sich nur um Eins unterscheiden.

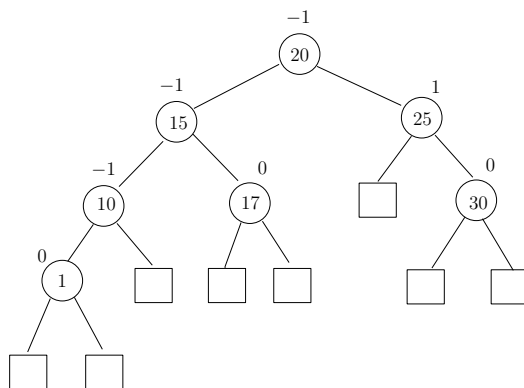


Abbildung 7.6: An jedem Knoten eines AVL-Baumes betragen die Höhenunterschieden der Teilbäume maximal 1.

Für jeden Knoten v gibt es einen sogenannten *Balancefaktor*.

$$\text{bal}(v) := \text{Höhe}(\text{rechts}(v)) - \text{Höhe}(\text{links}(v)) \quad (7.3)$$

In AVL-Bäumen (AVL nach Adelson, Velskij und Landis) gilt für alle Knoten $\text{bal}(v) \in \{-1, 0, 1\}$ ein Beispiel zeigt die Abbildung 7.6. Wir wollen nun die Höhe eines AVL-Baumes ermitteln. Dazu geben wir eine Höhe h vor und stellen fest, wieviele Knoten $n(h)$ ein AVL-Baum der Höhe h mindestens hat. Dann können wir umgekehrt feststellen welche maximale Höhe ein AVL-Baum für n Knoten aufweisen kann.

Ein AVL-Baum der Höhe h mit minimaler Knotenanzahl $n(h)$ heißt auch *minimaler AVL-Baum der Höhe h* . Wir zählen die Blätter $b(h)$ eines solchen Baumes und können nach Lemma 11 dann die Anzahl der Knoten ermitteln. Sei T_h ein minimaler AVL-Baum mit $b(h)$ Blättern. Betrachte Abbildung 7.7. Offensichtlich gilt:

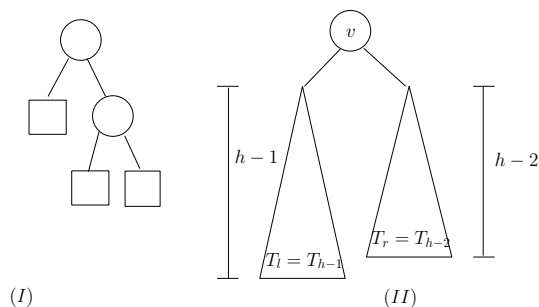


Abbildung 7.7: (I) Der minimale AVL-Baum der Höhe 2 hat $b(2) = 3$ Blätter. (II) Der minimale AVL-Baum der Höhe $h \geq 2$ hat $b(h) = b(h-1) + b(h-2)$ Blätter.

- Ein minimaler AVL-Baum der Höhe 0 hat ein Blatt, $b(0) = 1$.
- Ein minimaler AVL-Baum der Höhe 1 hat zwei Blätter, $b(1) = 2$.
- Ein minimaler AVL-Baum der Höhe 2 hat drei Blätter, $b(2) = 3$.
- Sei T_h ein minimaler AVL-Baum der Höhe h mit Wurzel v und seien T_l und T_r die Teilbäume von v . Dann hat zumindest einer der beiden Teilbäume eine Höhe von $h-1$ und der andere muss mindestens eine Höhe von $h-2$ haben. Da wir die Anzahl der

Bätter minimieren wollen, wählen wir o.B.d.A. $T_l = T_{h-1}$ und $T_r = T_{h-2}$ und für $h \geq 2$ gilt somit:

$$b(h) = b(h-1) + b(h-2) \quad (7.4)$$

Theorem 13 Die Höhe eines AVL-Baumes mit n Blättern (und $n-1$ inneren Knoten) beträgt $O(\log n)$.

Beweis. Sei $F_0 = 0$, $F_1 = 1$ und $F_2 = 1$ und $F_{h+2} = F_{h+1} + F_h$ für $h \geq 2$ dann ist offensichtlich $F_{h+2} = b(h)$ für $h \geq 0$. Die Folge F_h ist die bereits behandelte Fibonacci-Folge.

Wir zeigen zunächst, dass $F_{i+2} \geq 2F_i$ für $i \geq 1$ gilt. Das gilt für $i = 1$ und für $i \geq 2$ wenden wir direkt die Definition der Folge an, also $F_{i+2} = F_i + F_{i+1} = F_i + F_i + F_{i-1} \geq 2F_i$.

Mit $F_1 = F_2 = 1$ folgt nun $F_{h+2} \geq 2^{h/2}$ für $h \geq 1$. Das gilt zunächst für $h = 1$ und dann folgt induktiv, dass $F_h \geq 2^{(h-2)/2}$ für $h \geq 3$ bereits gilt. Dann ist $F_{h+2} \geq 2F_h \geq 2 \cdot 2^{(h-2)/2} = 2^{h/2}$.

Sei nun $b(h) = n$. Dann ist $n = F_{h+2} \geq 2^{h/2}$ und somit $h \leq 2 \log n$. \square

Als nächstes wollen wir die Operationen Einfügen und Entfernen für AVL-Bäume betrachten. Diese Operationen sollen

- ... *strukturerhaltend* durchgeführt werden, die AVL-Eigenschaft soll erhalten bleiben.
- ... *kostengünstig* durchgeführt werden, der Durchlauf in $O(\log n)$ und eventuelles Ausbalancieren in konstanter Zeit durch Umhängen von Zeigern.

7.4.1 Einfügen eines neuen Knoten

Dazu betrachten wir das Beispiel aus Abbildung 7.6. Im Baum notieren wir den Wert $\text{bal}(v)$ für jeden Knoten. Wir merken uns an jedem Knoten die Höhe des linken und rechten Teilbaumes, so können wir schnell den Wert $\text{bal}(v)$ bestimmen bzw. neu berechnen. Wir wollen einen Knoten z mit Schlüssel 5 einfügen.

Zunächst fügen wir dabei **abwärtslaufend** wie in der Prozedur 12 den Knoten ein; siehe Abbildung 7.8. Danach werden **aufwärtslaufend** die $\text{bal}(v)$ -Werte geändert bis zum ersten Mal eine Disbalance auftritt. Die neuen $\text{bal}(v)$ -Werte werden durch die veränderten Höhen sukzessive berechnet. Am Knoten y mit Schlüssel 10 tritt zum ersten Mal eine Disbalance auf. Die Teilbäume des Knoten sollen umgeordnet werden. Eine einfache Rechtsrotation am Knoten x mit y führt hier leider nicht zum Erfolg, der neu eingefügte Knoten z (Baum B aus Abbildung 7.5) würde seine Höhe nicht verlieren und am neuen Knoten x wird dann lediglich der Wert $\text{bal}(x) = +2$ notiert.

Deshalb führen wir hier die folgende *Doppelrotation* aus:

1. Zunächst eine Linksrotation des Knoten z mit x .
2. Danach eine Rechtsrotation des Knoten z mit y .

Zunächst verliert dabei der Knoten x an Höhe, das wird aber durch die zweite Rotation wieder ausgeglichen. In Abbildung 7.9 sehen wir die Zwischenergebnisse der Rotationen. Der zugehörige Teilbaum ist danach ausbalanciert. Abschließend stellt sich die Frage ob gegebenenfalls noch die $\text{bal}(v)$ Einträge bis Wurzel aktualisiert müssen. Wir haben den Teilbaum des Knoten y ausbalanciert aber er hat seine Höhe behalten, somit ändern sich die Werte nicht. Die Abbildung 7.10 zeigt das Ergebnis nach dem Einfügen.

Das Einfügen gestaltet sich insgesamt in den folgenden Schritten:

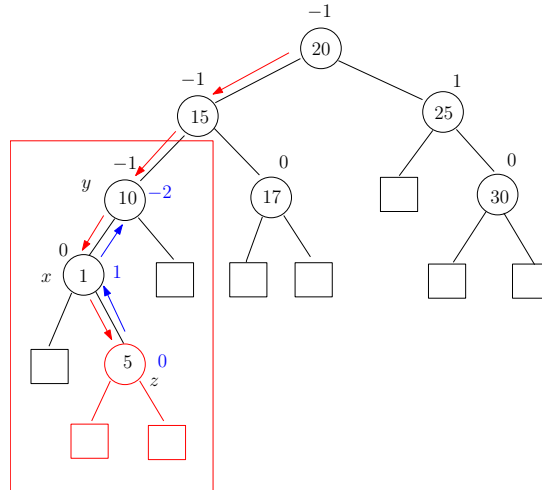


Abbildung 7.8: Das Einfügen eines Knoten z mit Schlüssel 5 in den AVL-Baum aus Abbildung 7.6 durch einen **Durchlauf** zum entsprechenden Blatt. Danach werden **rückwärtslaufend** die $bal(v)$ -Werte geändert, bis zum ersten Mal eine Disbalance auftritt. Der Teilbaum in der eingezeichneten Box muss ausbalanciert werden.

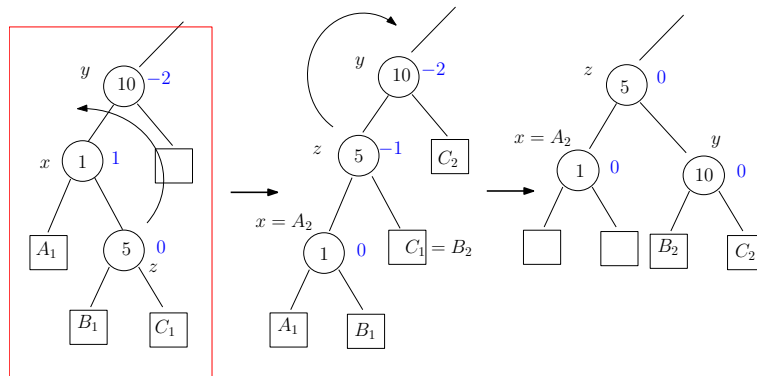


Abbildung 7.9: Zuerst wird eine Linksrotation am Knoten z mit x vorgenommen, danach eine Rechtsrotation am Knoten z mit y . Der Teilbaum ist danach balanciert.

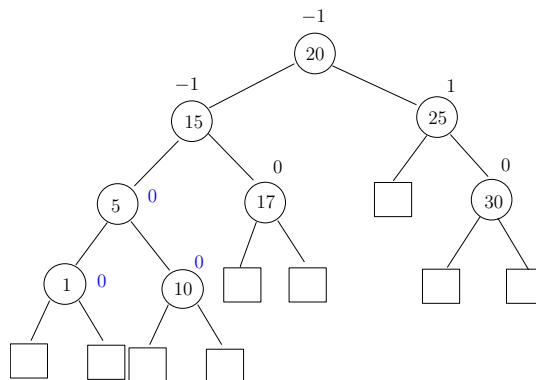


Abbildung 7.10: Der ausbalancierte Baum nach dem Einfügen des Knoten z des Schlüssels 5. Die $bal(v)$ -Werte alle anderen Knoten bleiben gleich, da die Höhe des Teilbaumes erhalten bleibt.

1. Vorwärtsdurchlauf gemäß der Schlüssel und Eintrag des neuen Knoten. Laufzeit: $O(\log n)$
2. Rückwärtsdurchlauf zum Anpassen der $\text{bal}(v)$ -Werte bis die erste Disbalance an einem Knoten v auftritt. Laufzeit: $O(\log n)$
3. Ausbalancierung des zugehörigen Teilbaums durch Rotation. Laufzeit: $O(1)$, inklusive umsetzen der Zeiger

Kategorisierung der Rotationen:

Die Doppelrotation (erst links, dann rechts) im obigen Beispiel war notwendig, weil der Weg vom Knoten y der ersten Disbalance zum neu eingefügten Knoten z zuerst links und danach rechts abgebogen ist. Der neue Knoten liegt also vom Knoten y aus gesehen im rechten Unterbaum des linken Unterbaumes. Dieser rechte Unterbaum muss an Höhe verlieren.

Falls der besagte Weg zuerst rechts und danach links abbiegt, wird die Doppelrotation in der Reihenfolge rechts/links ausgeführt. Der neue Knoten liegt dann also vom Knoten y aus gesehen im linken Unterbaum des rechten Unterbaumes.

Manchmal reicht auch eine einfache Rechts- oder Linksrotation aus. Das ist der Fall, wenn der Weg vom Knoten y der ersten Disbalance zum neu eingefügten Knoten z nur rechts oder nur links abbiegt.

Falls wir beispielsweise in den neuen Baum aus Abbildung 7.10 einen weiteren Knoten z mit Schlüssel 32 einfügen, so erhalten wir wie in Abbildung 7.11 die erste Disbalance am Knoten y und führen eine einfache Linksrotation des Knoten x mit y aus. Das Ergebnis ist in Abbildung 7.12 zu sehen.

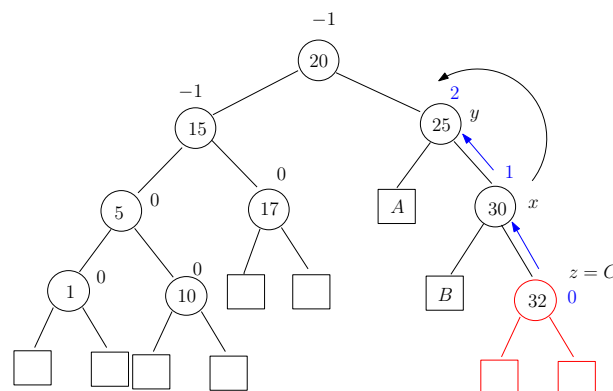


Abbildung 7.11: Nach dem Einfügen des Knoten z mit Schlüssel 32 reicht eine einfache Linksrotation.

Nach diesen Beispielen beschreiben wir die Doppelrotation (links/rechts) und die einfache Rotation (links) schematisch. Die Doppelrotation (rechts/links) und die einfache Rotation (rechts) verlaufen analog.

Doppelrotation (links/rechts): Nach dem Einfügen eines Knoten z und dem Ermitteln des untersten Disbalance-Knoten y läuft der Weg zum Knoten z in den rechten Unterbaum des linken Unterbaumes, siehe Abbildung 7.13. Zunächst wird eine Linksrotation des Knoten w mit x ausgeführt, siehe das Ergebnis in Abbildung 7.14. Danach wird eine Rechtsrotation des Knoten w mit y ausgeführt, siehe das Ergebnis in Abbildung 7.14. Der Baum ist ausgeglichen und das wäre auch der Fall, wenn der Knoten z in B eingefügt worden wäre, siehe z' in Abbildung 7.14.

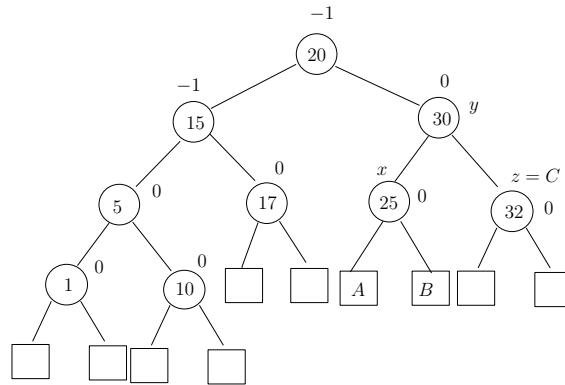


Abbildung 7.12: Die Höhe des ausbalancierten Baumes ändert sich nicht, alle anderen $bal(v)$ -Werte bleiben gleich.

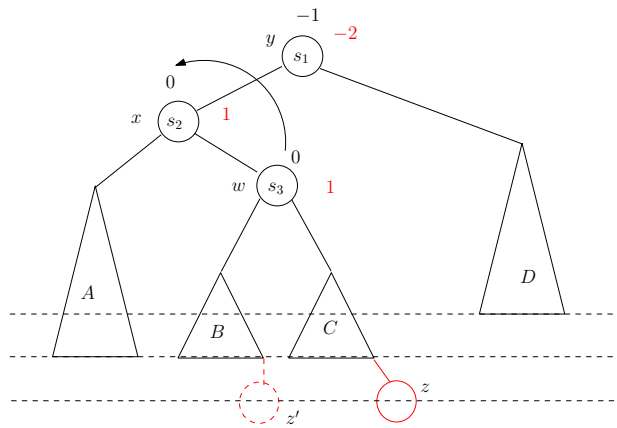


Abbildung 7.13: Die Situation vor einer (links/rechts) Doppelrotation. Der Knoten z könnte auch in B eingefügt worden sein.

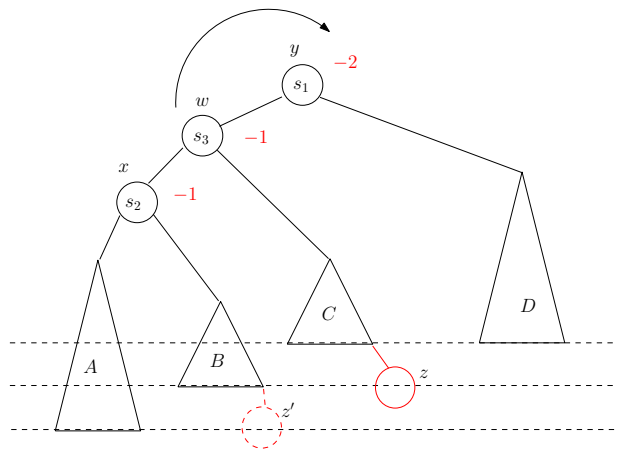


Abbildung 7.14: Die Situation vor einer (links/rechts) Doppelrotation. Der Knoten z könnte auch in B eingefügt worden sein.

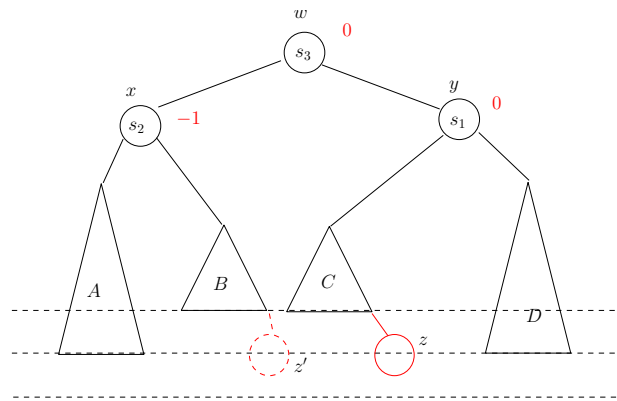


Abbildung 7.15: Nach der zweiten Rotation ist der Baum ausgeglichen. Falls der Knoten z in den Teilbaum B eingefügt worden wäre, wäre der Baum ebenfalls ausgeglichen. In jedem Fall ist die Höhe des Baumes wieder identisch zur Ausgangshöhe.

Einfache Rotation (links): Nach dem Einfügen eines Knoten z und dem Ermitteln des untersten Disbalance-Knoten y läuft der Weg zum Knoten z in den rechten Unterbaum des rechten Unterbaumes, siehe Abbildung 7.16.

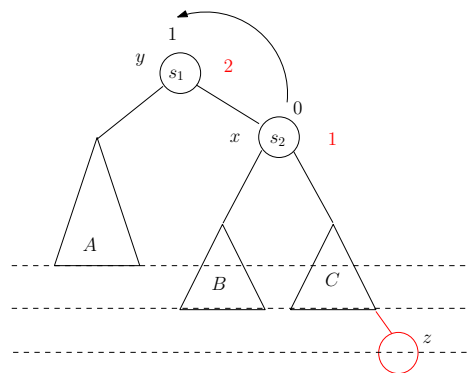


Abbildung 7.16: Die Situation vor einer einfachen Rotation (links). Der neue Knoten z liegt im rechten Unterbaum des rechten Unterbaumes von y .

Wir beschreiben nun in der Prozedur 13, wie die Zeiger bei einer Linksrotation verändert werden müssen.

Eine Rechtsrotation kann analog ausgeführt werden. Die Doppelrotation (links/rechts) die für einen Disbalance-Knoten y ausgeführt wird, kann dann durch folgende Prozedur beschrieben werden.

Übungsaufgabe: Beschreiben Sie die Schemata für eine Doppelrotation (rechts/links) und eine einfache Rotation (rechts) und die zugehörigen Zeigerprozeduren.

Korrektheit: Von der Logik her können Disbalancen am Knoten y durch das Einfügen eines Knoten z nur durch eine der vier Einfügefälle von z (links/rechts, rechts/links, links/links, rechts/rechts) entstehen. Wird der Knoten z beispielsweise direkt nur links unter y eingefügt dann kann $\text{rechts}(y)$ nur aus einem Blatt oder einem Baum der Höhe 1 bestanden haben. In jedem Fall entsteht bei y gar keine Disbalance. Der andere Fall geht analog.

Außerdem wird durch das Beseitigen der Disbalance beim Knoten y stets der gesamte Baum ausgeglichen. Falls beispielsweise die Doppelrotation (links/rechts) durchgeführt wurde, muss

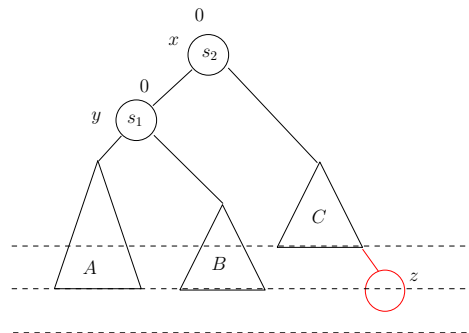


Abbildung 7.17: Die Situation nach einer einfachen Rotation (links). Der Teilbaum C hat an Höhe verloren, der Teilbaum A an Höhe gewonnen, der Baum ist ausgeglichen. Außerdem ist die Höhe des Baumes wieder identisch zur Ausgangshöhe.

Prozedur 13 Linksrotation(T,x)

Linksrotation des Knoten x mit $y = \text{Vater}(x)$

```

1:  $y := \text{Vater}(x)$ ;  $w := \text{Vater}(y)$ ;
2: if  $w \neq \text{Nil}$  and  $y = \text{LinkerSohn}(w)$  then
3:   marker := links;
4: else
5:   marker := rechts;
6: end if
7:  $\text{rechts}(y) := \text{links}(x)$ ;
8:  $\text{links}(x) := y$ ;
9: if marker := links then
10:   $\text{links}(w) := x$ ;
11: else
12:   $\text{rechts}(w) := x$ ;
13: end if

```

Prozedur 14 DoppelrotationLR(T,y)

Doppelrotation (links/rechts) am Disbalance-Knoten y

```

1:  $x := \text{links}(y)$ ;
2:  $w := \text{rechts}(x)$ ;
3: Linksrotation( $T, w$ );
4: Rechtsrotation( $T, w$ );

```

der Knoten y vor dem Einfügen den Wert -1 gehabt haben. Nach der Rotation hat der Teilbaum mit neuer Wurzel w zwar den bal-Wert 0 aber die Höhe des Teilbaumes bleibt erhalten. Das heißt, alle anderen $\text{bal}(v)$ -Werte bleiben erhalten. Der Baum ist insgesamt ausgeglichen. Die anderen Fälle gelten analog.

Insgesamt führen wir ein korrektes Einfügen eines Knoten z in Laufzeit $O(\log n)$ durch.

7.4.2 Entfernen eines Knoten

Beim Entfernen von Knoten kommen analoge Operationen zur Anwendung. Wie bereits in Abschnitt 7.3.1 bemerkt, wird stets ein Knoten v durch seinen linken oder rechten Teilbaum w ersetzt, der andere Teilbaum war dabei ein Blatt. Wenn der Baum nach dieser Operation am Knoten w unbalanciert ist, muss der Baum rebalanciert werden.

Wir geben zunächst ein Beispiel an und beschreiben dann die allgemeine Vorgehensweise. Wenn in Abbildung 7.12 der Knoten mit Schlüssel 15 und zwei echten Teilbäumen entfernt werden soll, so wird der Knoten zuerst durch den Knoten mit nächstgrößtem Schlüssel 17 ersetzt und der vormalige Knoten mit Schlüssel 17 entfernt. Dieser hatte mindestens ein Blatt, in unserem Fall sogar zwei Blätter. Der Teilbaum wurde entfernt und in diesem Fall durch ein Blatt ersetzt. Auf dem Weg zur Wurzel werden dann die $\text{bal}(v)$ -Werte angepasst, bis zum ersten Mal eine Disbalance entdeckt wird. Nun müssen wiederum Rotationen ausgeführt werden, um die Disbalance zu beseitigen. Offensichtlich reicht hier eine einfache Rechtsrotation am Knoten mit Schlüssel 5 aus, siehe Abbildung 7.18. Es kann ebenfalls sein, dass

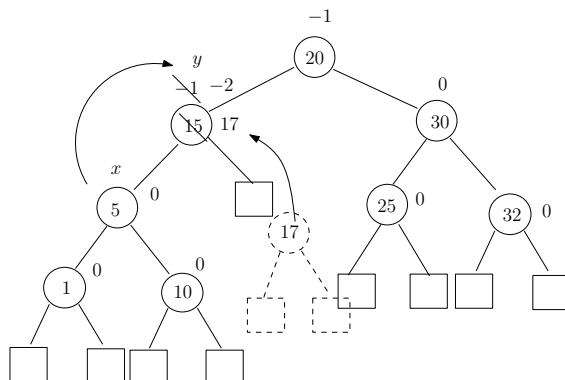


Abbildung 7.18: Die Situation vor einer einfachen Rotation nach dem Entfernen des Knoten mit Schlüssel 15 . Der ehemalige Knoten des Schlüssels 17 wird entfernt.

Doppelrotationen notwendig sind. Wir beschreiben kurz das allgemeine Schema. Im Gegensatz zur Beschreibung des Einfügens eines Knoten beschreiben wir hier die einfache Rotation (rechts) und die Doppelrotation (links/rechts) die durch eine -2 Disbalance entstehen. Analoge Operationen einfache Rotation (links) und Doppelrotation (rechts/links) entstehen bei $+2$ Disbalancen.

Angenommen die Disbalance tritt am Knoten y zuerst auf. Falls die Disbalance im Knoten y durch einen -2 Wert gegeben ist, wurde ein rechter Teilbaum C von y angehoben. Für den linken Nachfolger x ergeben sich zwei Fälle. Falls der bisherige $\text{bal}(x)$ -Wert 0 oder -1 war, reicht eine einfache Rechtsrotation am Knoten x mit Knoten y aus, siehe Abbildung 7.20. Der Teilbäume A von x wird um eins angehoben, der Teilbaum C um eins gesenkt. Der Teilbaum B respektive B' behält seine Höhe. Es kann auch vorkommen (Fall B'), dass die Höhe des Baumes insgesamt um Eins sinkt.

Falls allerdings der $\text{bal}(x)$ Wert $+1$ war, hat der rechte Nachfolger w von x zwei Teilbäume B

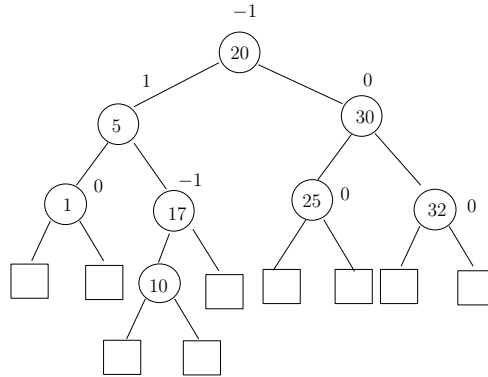


Abbildung 7.19: Nach der Rotation ist der Baum ausgeglichen.

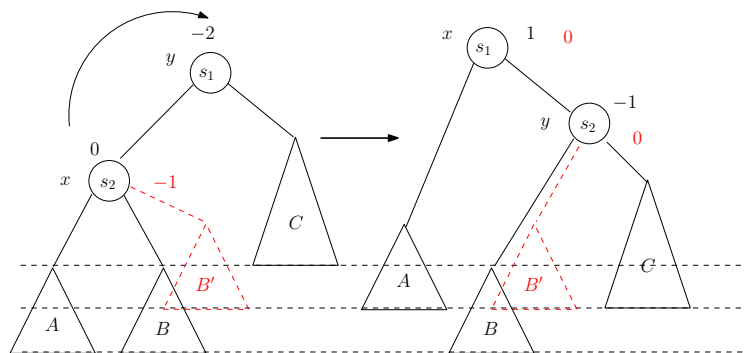


Abbildung 7.20: Die Situation vor einer einfachen Rotation. Der Disbalance-Knoten y ($\text{bal}(y) = -2$) hat einen linken Nachfolger x mit $\text{bal}(x) \in \{0, -1\}$. Dann reicht eine einfache Rechtsrotation des Knoten x mit y aus.

und C , von denen mindestens einer zusammen mit w eine höhere Tiefe als der linke Nachfolger A von x aufweist. Eine einfache Rechtsrotation von x mit y wird die Höhe des Baumes w erhalten, die Höhe von A aber absenken, somit entsteht eine $+2$ Disbalance.

Hier ist somit wiederum eine Doppelrotation notwendig. Zunächst eine Linksrotation des Knoten w mit x und danach eine Rechtsrotation des Knoten w mit y . Das vollständige Ergebnis ist in Abbildung 7.21 zu sehen. Beachte, dass C und B auch die Rollen tauschen können. Das heißt, dass entweder B oder C die Disbalance am Knoten y zusammen mit D verursacht. In jedem Fall wird der Baum rebalanciert, hat aber insgesamt in jedem Fall die Höhe um Eins verringert. Übungsaufgabe: Beschreiben Sie die Schemata für die Rotationen nach dem

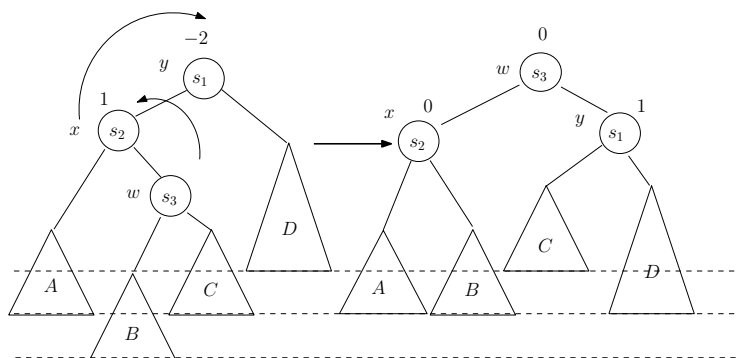


Abbildung 7.21: Die Situation vor einer Doppelrotation (links/rechts). Der Disbalance-Knoten y ($\text{bal}(y) = -2$) hat einen linken Nachfolger x mit $\text{bal}(x) = 1$. Dann ist eine Doppelrotation (zuerst Linksrotation Knoten w mit x und danach Rechtsrotation Knoten w mit y) notwendig.

Entfernen, falls sich die Disbalance durch einen ersten -2 Wert an einem Knoten y ergibt und somit ein rechter Teilbaum von y an Höhe verloren hat. Beschreiben Sie die zwei Unterfälle, dass der linke Nachfolger x von y den Wert $\text{bal}(x) = -1$ bzw. $\text{bal}(x) \in \{0, 1\}$ aufweist.

Korrektheit: Für den ersten Disbalance Knoten y gibt es nur die beiden Möglichkeiten $\text{bal}(y) \in \{2, -2\}$, wobei die Disbalance nur durch die jeweils beschriebenen Fälle aufgetreten sein konnten. Entweder wurde die Höhe eines rechten oder eines linken Teilbaumes von y reduziert. Für den jeweils anderen Nachfolger x von y ergeben sich dann die zwei möglichen Unterfälle. Insgesamt wird zur Rebalancierung von y entweder eine einfache Rotation (links oder rechts) oder eine Doppelrotation (links/rechts oder rechts/links) ausgeführt. Danach ist der Knoten y rebalanciert. Wir sehen aber, dass der Teilbaum von y (im Vergleich zur Situation vor dem Entfernen) insgesamt immer um 1 an Höhe verliert, deshalb kann es sein, dass die verringerte Höhe von y weitere $\text{bal}(v)$ -Werte auf dem Weg zur Wurzel verändert.

Das ist nicht weiter tragisch. Wir wenden die Argumentation sukzessive wieder an. Für den nächsten Disbalance-Knoten auf dem Weg zur Wurzel führen wir die gleichen Operationen durch. Spätestens wenn wir an der Wurzel angekommen sind, ist der gesamte Baum ausgeglichen.

Laufzeit: Für das Entfernen wird zunächst ein Aufwand in $O(\log n)$ benötigt. Für die Rebalancierung laufen wir aufsteigend zu den nächsten zu rebalancierenden Knoten. Der Weg zur Wurzel hat eine Länge von $O(\log n)$. Jede einzelne Balancierungsoperation kostet $O(1)$ Zeit. Insgesamt wird nur $O(\log n)$ Zeit benötigt.

7.5 B-Bäume

Wenn die Datenmenge so groß wird, dass sie nicht mehr in den schnellen Hauptspeicher passt, muss sie in langsamere Speicher ausgelagert werden. Hier gilt das Modell der REAL RAM nur noch bedingt, denn die Laufzeit von Algorithmen wird jetzt von der *Anzahl der Zugriffe auf das langsame Speichermedium* dominiert.

Zum Beispiel muss bei externen Festplatten ein Schreib-/Lesekopf auf die richtige Spur bewegt werden; das dauert zwar nur einige Millisekunden, aber doch fast 10^3 mal so lang wie ein Hauptspeicherzugriff. Dafür liefert ein Festplattenzugriff nicht nur einen einzelnen Wert, sondern eine Datenseite (Datenblock) mit Volumen 512 Byte bis 1 KByte.

Diesem Umstand trägt das Modell der *B-Bäume* Rechnung.

Definition 14 Ein *B-Baum der Ordnung k* ist durch folgende Eigenschaften bestimmt:

1. Alle Blätter haben dieselbe Tiefe.
2. Alle Knoten außer der Wurzel haben mindestens k und höchstens $2k - 1$ viele Söhne.
3. Die Wurzel hat mindestens 2 und höchstens $2k - 1$ viele Söhne.

Wenn h die Höhe eines solchen B-Baums bezeichnet und n die Anzahl seiner Blätter, so gilt

$$2 \cdot k^{h-1} \leq n \leq (2k - 1)^h;$$

die untere Schranke wird erreicht, wenn die Wurzel nur 2 und alle anderen internen Knoten nur k Söhne haben, die obere, wenn alle Knoten maximal viele, nämlich $2k-1$, Kinder besitzen. Daraus folgt

$$\log_{2k-1} n \leq h \leq 1 + \log_k \frac{n}{2} \quad (7.5)$$

Abbildung 7.22 zeigt, wie ein Knoten eines B-Suchbaums aufgebaut ist.

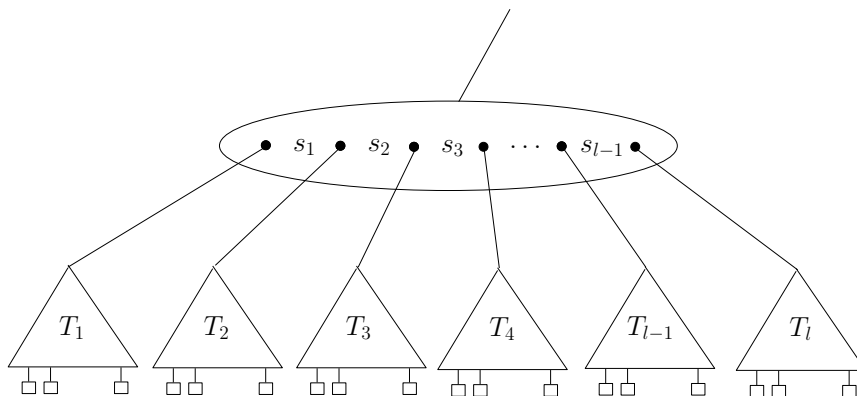


Abbildung 7.22: Im Knoten eines B-Baums wechseln Zeiger auf Teilbäume und Schlüssel einander ab. Alle Schlüssel in T_i sind kleiner als s_i , und s_i ist kleiner als die Schlüssel in T_{i+1} .

Schlüssel stehen nur in den inneren Knoten, zwischen den Zeigern auf die Teilbäume. Sie sind aufsteigend sortiert: Für alle i gilt

$$\text{Schlüssel von } T_i < s_i < \text{Schlüssel von } T_{i+1}.$$

Der Parameter k wird meist so gewählt, dass ein maximal voller Knoten mit $2k - 1$ Zeigern und $2k - 2$ Schlüsseln gerade auf eine Datenseite passt. Wegen der unteren Schranke k für die Anzahl der Söhne ist dann jede Datenseite (außer der Wurzel) zumindest halb voll.

Beim *Suchen* eines Schlüssels x startet man im Wurzelknoten und durchläuft von links nach rechts die dort gespeicherten Schlüssel s_1, \dots, s_{l-1} , bis man zum letzten Schlüssel s_i mit $s_i \leq x$ gelangt. Gilt $s_i = x$, so ist man fertig. Andernfalls setzt man die Suche im Teilbaum T_{i+1} fort.

Um einen Schlüssel x *einzufragen*, wird x zunächst gesucht. Kommt x bereits vor, so wird eine Meldung generiert. Andernfalls endet die Suche erfolglos in einem Blatt b , das in seinem Vaterknoten v hinter einem Schlüssel $s_i < x$ angehängt ist. Wir fügen x als neuen Schlüssel hinter s_i in die Liste ein und ein neues Blatt b' hinter x . Wenn Knoten v jetzt immer noch $\leq 2k - 1$ viele Söhne hat, ist man fertig. Wenn nicht, hat v jetzt genau $2k$ viele Söhne und wird in zwei gleich große Knoten v_1 und v_2 aufgeteilt; dabei wandert der k -te Schlüssel (nach Einfügen von x) in den Vater v' von v , siehe Abbildung 7.23. Möglicherweise ist jetzt auch Knoten v' überfüllt und muss geteilt werden. Das kann sich bis zur Wurzel w fortpflanzen; in dem Fall wird w in zwei Knoten aufgeteilt und eine neue Wurzel darübersetzt, die genau diese zwei Söhne hat.

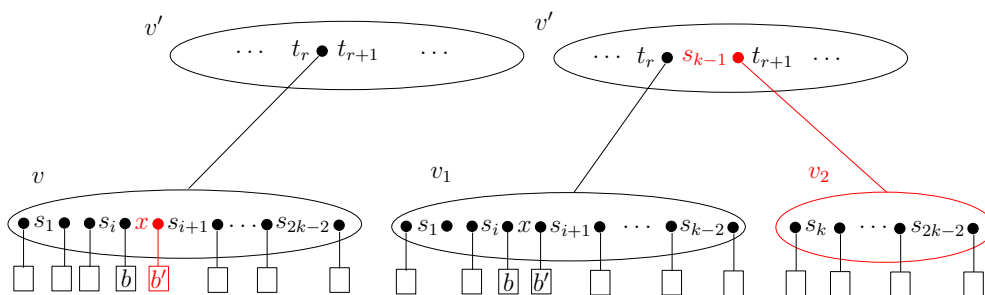


Abbildung 7.23: Knoten v wird geteilt, wenn er nach Einfügen von Schlüssel x und einem neuen Blatt b' nun $2k$ viele Söhne hat.

Auch beim *Entfernen* eines Schlüssels x wird dieser zunächst gesucht. Falls x nicht vorkommt, wird eine Meldung generiert. Andernfalls steht x in einem Knoten v hinter dem Zeiger auf einen Teilbaum T_j . Im ersten Fall sind T_j und alle anderen Söhne von v Blätter. Dann werden T_j und x einfach entfernt. Dabei kann es passieren, dass v nur noch $k - 1$ Söhne bleiben; wir überlegen uns gleich, was dann zu tun ist. Im zweiten Fall hat v echte Teilbäume als Söhne. Dann sei s der größte Schlüssel in T_j , enthalten in einem Knoten w . Rechts von s steht in w nur noch der Zeiger auf ein Blatt b ; siehe Abbildung 7.24. Schlüssel s ist der Vorgänger von x im Baum. Darum ersetzen wir im Knoten v den Schlüssel x durch s und entfernen dann im Knoten w den Schlüssel s und das Blatt b (so, als hätten wir von Anfang an s entfernen wollen, wie im ersten Fall). Auch hier kann es passieren, dass w nur $k - 1$ Söhne behält.

Dann wird w mit einem Nachbarknoten u mit $l \geq k$ Söhnen vereinigt; das ist gerade die Umkehrung der Knotenteilung beim Einfügen. Falls $l > k$ war, können wir den Vereinigungsknoten gleich wieder in zwei Knoten mit jeweils $\geq k$ vielen Söhnen teilen und sind nach dieser Umverteilung von Schlüsseln fertig. Andernfalls hat der Vater w' von w und u durch die Vereinigung einen Sohn verloren und könnte jetzt selbst unterbelegt sein. Dann muss w' mit einem Nachbarknoten vereinigt werden, und so fort. Sollte die Wurzel nur noch einen Sohn haben, wird sie entfernt, und der Sohn wird zur neuen Wurzel.

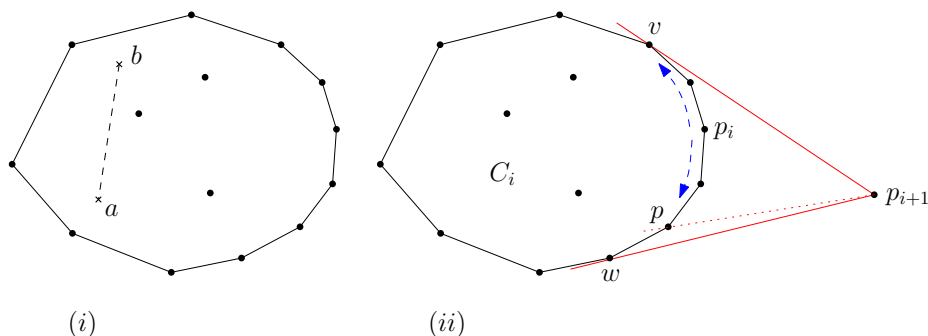


Abbildung 7.25: (i) Die konvexe Hülle von Punkten in der Ebene. (ii) Punkt p_{i+1} kommt hinzu.

Zu Beginn werden die n Punkte nach aufsteigenden X -Koordinaten sortiert; der Einfachheit halber wird angenommen, dass diese Koordinaten paarweise verschieden sind und dass außerdem keine drei Punkte auf einer Geraden liegen. Jetzt verfährt man inkrementell. Die konvexe Hülle von p_1, p_2, p_3 ist einfach das Dreieck mit diesen Eckpunkten. Angenommen, die konvexe Hülle C_i von p_1, \dots, p_i ist schon konstruiert, und p_{i+1} soll nun als weiterer Punkt hinzugenommen werden. Wegen der Sortierung liegt C_i links von der Senkrechten durch p_i , und p_{i+1} liegt rechts davon; siehe Abbildung 7.25 (ii). Deshalb kann p_{i+1} den Punkt p_i sehen. Wir starten bei p_i und durchlaufen den Rand von C_i nach oben und unten, bis die beiden Tangentialpunkte v und w gefunden sind. Unterwegs muss man bei jedem Eckpunkt p testen, ob die Halbgerade von p_{i+1} durch p die konvexe Hülle C_i berührt oder schneidet; das geht in konstanter Zeit, weil man nur die beiden Kanten von C_i betrachten muss, die sich in p treffen.

Wenn v und w gefunden sind, ergibt sich die neue konvexe Hülle C_{i+1} , indem man die Segmente $\overline{vp_{i+1}}$ und $\overline{wp_{i+1}}$ hinzufügt und das Randstück von C_i zwischen v und w , welches p_i enthält, entfernt.

Diese Aktion kann Kosten in $O(n)$ verursachen, wenn zwischen v und w viele Eckpunkte p liegen, die besucht werden müssen. Aber alle diese Eckpunkte werden anschließend entfernt und kommen nie wieder vor, weil sie jetzt innerhalb von C_{i+1} liegen. Also betragen die Gesamtkosten aller Einfügeoperationen $O(n)$, weil es ja nur n Punkte gibt, und eine einzelne Operation braucht im Mittel nur konstante Zeit. Damit haben wir folgendes Ergebnis:

Theorem 16 *Die konvexe Hülle von n Punkten in der Ebene lässt sich nach Sortieren der Punkte in Zeit $O(n)$ konstruieren, insgesamt also in Zeit $O(n \log n)$.*

Der Speicherplatzbedarf ist linear, wenn man die konvexen Hüllen als verkettete Listen verwaltet. Beide Schranken sind optimal. Der oben vorgestellte Algorithmus ist sowohl ein Beispiel für *Inkrementelle Konstruktion* als auch für *Sweep*, je nach Betrachtungsweise.

7.6.2 Randomisierter Input: Bucketsort

Bei manchen Algorithmen hängt die Laufzeit stark von der Eingabe ab. Dann kann es sinnvoll sein, die *mittlere Laufzeit* zu betrachten. Dazu muss man wissen, mit welcher Wahrscheinlichkeit jede mögliche Eingabe an die Reihe kommt.

Betrachten wir als Beispiel den Algorithmus *Bucketsort*, einen einfachen Verwandten vom Radixsort aus Kapitel 2. Um n reelle Zahlen x aus dem Intervall $[0, 1)$ zu sortieren, werden zunächst alle Eingabewerte in denselben Behälter (Bucket) gelegt, die im gleichen Intervall $[\frac{i}{n}, \frac{i+1}{n})$ liegen, für $i = 0, 1, 2, \dots, n-1$. Nun werden die Zahlen eines jeden Behälters einzeln

sortiert, zum Beispiel mit einem quadratischen Verfahren wie Insertionsort. Schließlich werden die sortierten Behälterinhalte von links nach rechts ausgegeben.

Offensichtlich braucht Bucketsort Zeit $\Omega(n^2)$, wenn alle n Eingabewerte im selben Intervall $[\frac{i}{n}, \frac{i+1}{n})$ liegen, denn dann wird die gesamte Arbeit von der quadratischen Prozedur Insertionsort geleistet. Allgemein hat Bucketsort die Laufzeit

$$T(n) = dn + c \sum_{i=0}^{n-1} B_i^2,$$

wobei B_i die Anzahl der Eingabewerte im Intervall $[\frac{i}{n}, \frac{i+1}{n})$ bedeutet; dabei misst dn den Aufwand für die Aufteilung in Buckets.

Die Summe der Quadrate der B_i entspricht genau der Anzahl aller Paare (x, y) von Eingabewerten, bei denen x und y im selben Intervall liegen. Das sind alle Paare (x, x) , und außerdem gewisse Paare (x, y) mit $x \neq y$. Wenn man nun annehmen darf, dass die Eingabewerte gleichverteilt und unabhängig aus dem Intervall $[0, 1)$ gezogen werden, dann liegt zu gegebenem x der Wert y genau mit Wahrscheinlichkeit $\frac{1}{n}$ im gleichen Intervall wie x , denn die Intervalle haben alle die gleiche Länge $\frac{1}{n}$, und die Ziehung von y hängt nicht von der Ziehung von x ab. Für den Erwartungswert der Laufzeit ergibt sich demnach

$$\begin{aligned} E(T(n)) &= dn + c \left(\sum_x 1 + \sum_x \sum_{y \neq x} \frac{1}{n} \right) \\ &= dn + c \left(n + n(n-1) \frac{1}{n} \right) \\ &< (d+2c)n \end{aligned}$$

aus der Linearität des Erwartungswerts. Also gilt:

Theorem 17 *Bei gleichverteilten und unabhängig gezogenen Eingabewerten aus $[0, 1)$ hat Bucketsort die mittlere Laufzeit $O(n)$.*

Auch bei Bucketsort finden, wie schon bei Radixsort, nicht nur Schlüsselvergleiche statt, denn die Eingabewerte x müssen ja ihren Intervallen zugeordnet werden. Um die Bedingung $\frac{i}{n} \leq x < \frac{i+1}{n}$, oder äquivalent: $i \leq nx < i+1$, testen zu können, wird die Rundungsfunktion $\lfloor nx \rfloor$ benötigt. Daher widerspricht Theorem 17 nicht der unteren Schranke in Theorem 2.

In manchen Fällen ist man nicht daran interessiert, die Laufzeit über alle möglichen Eingaben zu mitteln, sondern nur in der Umgebung besonders "schwieriger" Inputs; so kann man Aufschluss darüber gewinnen, ob es sich um singuläre Ausreißer handelt.

7.6.3 Randomisierter Algorithmus: Bestimmung des Maximums

In den meisten Fällen ist nicht bekannt, welcher Wahrscheinlichkeitsverteilung die Eingaben unterliegen; dann muss ein Algorithmus mit den konkreten Eingabewerten arbeiten, die er bekommt. Trotzdem kann er Zufall ins Spiel bringen, indem er würfelt, welche von mehreren möglichen Aktionen als nächste ausgeführt werden soll.

Wir beginnen mit einem ganz einfachen Beispiel. Um das Maximum von n Schlüsseln a_i zu bestimmen, verfährt man wie folgt: Eine Variable \max wird zunächst auf a_1 gesetzt. Dann wird für $i = 2$ bis n getestet, ob $\max < a_i$ gilt, und gegebenenfalls das bisherige Maximum durch $\max := a_i$ aktualisiert. Das erfordert $n-1$ Tests und schlimmstenfalls $n-1$ Aktualisierungen, wenn die Schlüssel aufsteigend sortiert sind.

Angenommen, ein Test kostet Zeit T , eine Aktualisierung aber Zeit $A \gg T$, weil dabei Daten bewegt werden müssen. Dann würde man gern die Anzahl der Aktualisierungen begrenzen. Ein naheliegender Ansatz besteht darin, die gegebenen n Schlüssel zunächst in eine zufällige Reihenfolge zu bringen und dann obigen Algorithmus zu starten.

Wie wahrscheinlich ist es nun, dass die Aktualisierung $\max := a_i$ ausgeführt werden muss? Das geschieht genau dann, wenn a_i der größte Schlüssel in $A_i := \{a_1, a_2, \dots, a_{i-1}, a_i\}$ ist. Die Wahrscheinlichkeit hierfür beträgt $\frac{1}{i}$, denn nach der anfänglichen Zufallspermutation ist a_i eine beliebige Zahl in der i -elementigen Menge A_i . Also hat der randomisierte Algorithmus die zu erwartenden Kosten

$$E(T(n)) = (n-1) \cdot T + \sum_{i=2}^n A \cdot \frac{1}{i} \quad (7.6)$$

$$< (n-1) \cdot T + A \cdot \int_1^n \frac{1}{x} dx \quad (7.7)$$

$$< n \cdot T + \ln(n) \cdot A, \quad (7.8)$$

was deutlich besser ist als der Worst Case $n \cdot (T + A)$. Die Abschätzung der Summe durch das Integral ergibt sich durch Vergleich der Flächen unter den Funktionen.

7.6.4 Randomisierter Algorithmus: Quicksort

Eine einfache und weit verbreitete Variante des Sortierens mit Schlüsselvergleichen funktioniert folgendermaßen: Gegeben ist ein Array A mit n Schlüsseln a_1, \dots, a_n , die sortiert werden sollen. Eine dieser Zahlen wird als sogenanntes *Pivotelement* p ausgewählt. Dafür gibt es verschiedene Möglichkeiten; zum Beispiel könnte man einfach $p := a_n$ setzen. Jetzt werden alle $a_i < p$ in einer Folge L zusammengefasst, und alle $a_j > p$ in einer Folge R ; das geht in Zeit $O(n)$. Nun wendet man Quicksort rekursiv auf L und R an, schreibt die sortierten Teilfolgen hintereinander, dazwischen p , und ist fertig. (In der Literatur wird beschrieben, wie man direkt im Array A die Folge L links von p anordnet und R rechts von p , ohne zusätzlichen Speicherplatz zu benutzen.)

Die Effizienz von Quicksort hängt von der Größe von L und R ab und damit von der Wahl des Pivotelements p . Im besten Fall sind L und R stets gleich groß. Dann ergibt sich für die Laufzeit die Rekursion $T(n) = 2T(\frac{n}{2}) + cn$, woraus $T(n) \in O(n \log n)$ folgt; vergleiche Kapitel 3. Tatsächlich könnte man in Zeit $O(n)$ den sogenannten *Median* von a_1, \dots, a_n bestimmen, also das Element, welches in der Mitte der sortierten Folge liegt, und damit $|L| \approx |R|$ garantieren. Der Median-Algorithmus ist aber selbst rekursiv, und die Konstante in der $O(n)$ -Abschätzung recht hoch.

Der schlimmste Fall für Quicksort tritt ein, wenn die Eingabefolge bereits aufsteigend sortiert ist. Wir müssen dann $O(n)$ Zeit aufwenden, um festzustellen, dass die ersten $n-1$ Elemente kleiner als $p = a_n$ sind. Dann enthält L diese $n-1$ Elemente, und R ist leer. In diesem Fall braucht Quicksort quadratische Laufzeit.

Es liegt deshalb nahe, das Pivotelement p *zufällig* in a_1, \dots, a_n auszuwählen. Dann ist jede Position von p in der sortierten Folge gleich wahrscheinlich, und für die mittlere Laufzeit ergibt sich die Rekursion

$$E(T(n)) = c(n+1) + \frac{1}{n} \sum_{i=1}^n (E(T(i-1)) + E(T(n-i))) \quad (7.9)$$

$$= c(n+1) + \frac{2}{n} \sum_{i=0}^{n-1} E(T(i)), \quad (7.10)$$

denn bei beiden Summanden läuft der Index von 0 bis $n-1$. Zur Abkürzung sei $e_i := E(T(i))$. In der Rekursion (7.10) können wir die Summation durch Differenzenbildung eliminieren und bekommen die Gleichungen

$$\begin{aligned} e_{n+1}(n+1) - e_n n &= c \left((n+2)(n+1) - (n+1)n \right) + 2 \sum_{i=0}^n e_i - 2 \sum_{i=0}^{n-1} e_i \\ &= 2c(n+1) + 2e_n \\ (n+1)e_{n+1} &= 2c(n+1) + (n+2)e_n \\ \frac{e_{n+1}}{n+2} &= \frac{2c}{n+2} + \frac{e_n}{n+1} \end{aligned}$$

Aus der letzten Zeile ergibt sich durch iteriertes Einsetzen

$$\begin{aligned} \frac{e_{n+1}}{n+2} &= 2c \left(\frac{1}{n+2} + \frac{1}{n+1} + \dots + \frac{1}{2} \right) + \frac{e_0}{1} \\ &< 2c \ln(n+2) \end{aligned}$$

wie in (7.8), wobei $e_0 = 0$ ist. Damit haben wir folgendes Ergebnis:

Theorem 18 *Die mittlere Laufzeit von Quicksort, angesetzt auf eine beliebige Eingabe von n Schlüsseln, liegt in $O(n \log n)$.*

Quicksort und Mergesort sind beides Algorithmen vom Typ Divide-and-Conquer. Quicksort leistet die Hauptarbeit bei der Aufteilung in Teilmengen, Mergesort beim Zusammenfügen der Ergebnisse.

Im Kapitel 7.8 werden wir Hash-Verfahren untersuchen, bei denen sowohl randomisierte Eingaben als auch amortisierte Kosten eine Rolle spielen.

7.7 Heaps

In diesem Kapitel lernen wir eine Datenstruktur kennen, die sehr ähnlich zu den bisher behandelten Bäumen ist. Sie erbt die Vorzüge eines balancierten Baumes, muss jedoch nicht die volle Suchbaumeigenschaft erfüllen. Als Anwendungen der Datenstruktur betrachten wir zwei Beispiele. Zum einen den Sortieralgorithmus Heapsort, zum anderen die Prioritätenwarteschlange.

Ein Max-Heap ist ein fast vollständiger Binärbaum, welcher folgende Eigenschaften erfüllt:

- Der Baum ist, mit Ausnahme der untersten Ebene, vollständig besetzt. Die unterste Ebene ist von links nach rechts bis zu einem gewissen Punkt vollständig besetzt.
- Der Wurzelknoten enthält den größten Schlüsselwert. Alle anderen im Max-Heap gespeicherten Schlüssel sind nicht größer. Der linke und rechte Teilbaum der Wurzel ist jeweils ein Max-Heap.

Ganz ähnlich lässt sich auch ein Min-Heap definieren, in welchem der Elternknoten stets einen Wert kleiner oder gleich dem Schlüsselwert der Kinder enthält. In Abbildung 7.26a wird ein Beispiel für einen Max-Heap gegeben. Ein Heap kann sehr einfach in einem Array A abgelegt werden (s. Abbildung 7.26a). Dies hat unter anderem den Vorteil, dass für einen Wert in Zelle i leicht der Elternknoten bzw. das linke und rechte Kind bestimmt werden kann:

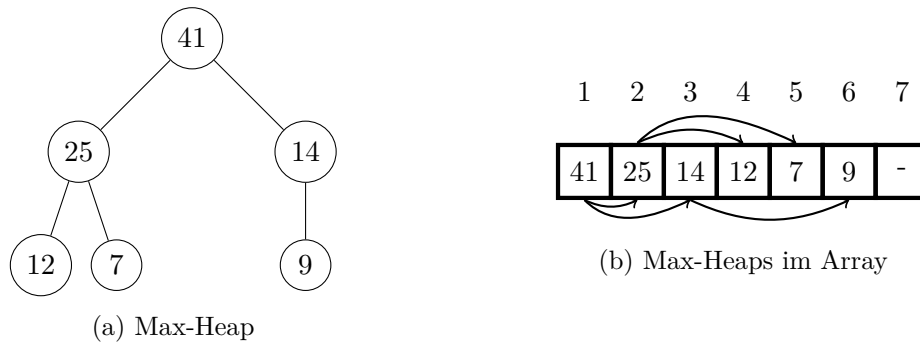


Abbildung 7.26: Beispiel eines Max-Heaps für die Schlüsselwerte 7,9,12,14,25,41 mit Einbettung im Array.

```

parent( $i$ ) : return  $\lfloor i/2 \rfloor$ ;
links( $i$ )   : return  $2i$ ;
rechts( $i$ )  : return  $2i + 1$ ;

```

Im Folgenden werden wir uns nur mehr mit Max-Heaps beschäftigen. Die Bezeichnung Heap steht daher stets für einen Max-Heap.

Bevor wir uns damit befassen, wie ein Heap gebaut werden kann und wie Schlüssel eingefügt bzw. gelöscht werden können, betrachten wir den grundlegenden Algorithmus 15. Er stellt die Heap-Eigenschaft für einen Knoten i her, wenn sowohl der linke als auch der rechte Teilbaum von i bereits die Heap-Eigenschaft erfüllen.

Prozedur 15 Max-Heapify(A, i)

```

1:  $l := \text{links}(i)$ ,  $r := \text{rechts}(i)$ ;
2: if  $l \leq A.\text{heap-size}$  und  $A[l] > A[i]$  then
3:    $max := l$ ;
4: else
5:    $max := i$ ;
6: end if
7: if  $r \leq A.\text{heap-size}$  und  $A[r] > A[max]$  then
8:    $max := r$ ;
9: end if
10: if  $max \neq i$  then
11:   vertausche  $A[i]$  und  $A[max]$ ;
12:   Max-Heapify( $A, max$ );
13: end if

```

Heap-Eigenschaft herstellen: Ist im Knoten i die Heap-Eigenschaft nicht erfüllt, so vertauscht Algorithmus 15 den Schlüssel von i mit dem Schlüssel des Kindes, das den maximalen Schlüsselwert enthält. Aus diesem Grund ist die Heap-Eigenschaft nun mehr höchstens in dem Teilbaum nicht mehr erfüllt, dessen Wurzel nun den Schlüssel von i enthält. Um dieses Problem zu lösen, ruft der Algorithmus sich rekursiv auf diesem Teilbaum auf. Der Schlüssel des Knotens i sinkt somit solange ab, bis die Heap-Eigenschaft im kompletten Baum erfüllt ist. Da der Binärbaum vollständig ist, hat er Höhe $O(\log n)$, wobei n die Anzahl der gespeicherten Knoten ist. Wie wir im Algorithmus sehen, so wird die Prozedur höchstens einmal rekursiv aufgerufen. Die Laufzeit von Max-Heapify hängt wegen der linearen Rekursion im Wesentlichen von der Höhe des Heaps ab, ist also $O(\log n)$.

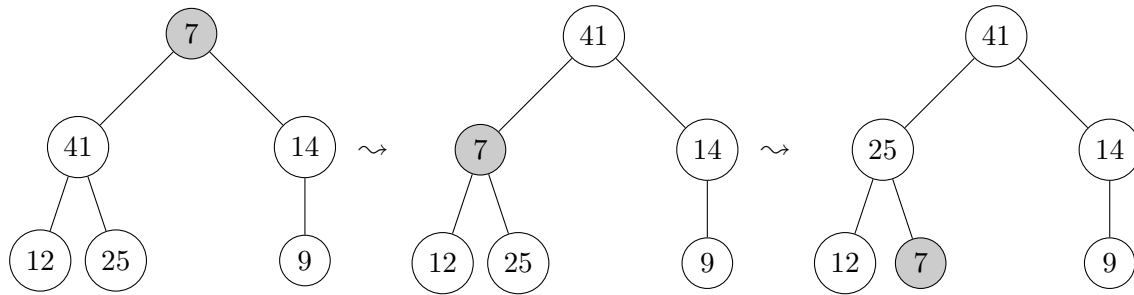


Abbildung 7.27: Beispiel für einen Aufruf von $\text{Max-Heapify}(A, 1)$. Die Heap-Eigenschaft ist im linken und rechten Teilbaum der Wurzel bereits erfüllt. Der Schlüssel der Wurzel wird mit dem größten Schlüssel der Kinder vertauscht und im Anschluss der Algorithmus rekursiv auf dem Teilbaum des linken Kindes aufgerufen, also $\text{Max-Heapify}(A, 2)$ aufgerufen.

Heaps bauen: Aus dem Algorithmus zum Herstellen der Heap-Eigenschaft können wir leicht einen Algorithmus zum Bau eines Heaps in einem Array A herleiten:

Prozedur 16 $\text{Build-Heap}(A)$

- 1: $A.\text{heap-size} = A.\text{length}$;
 - 2: **for** $i = \lfloor A.\text{length}/2 \rfloor$ **to** 1 **do**
 - 3: $\text{Max-Heapify}(A, i)$;
 - 4: **end for**
-

Die Blätter erfüllen bereits die Heap-Eigenschaft. Daher können wir für alle Ebenen nacheinander, von der Blattebene bis zur Wurzelebene die Heap-Eigenschaft mit Aufrufen von Max-Heapify herstellen. Insgesamt liegt die Bauzeit somit in $O(n \log n)$. In den Übungen werden wir durch eine genauere Analyse zeigen, dass diese Bauzeit tatsächlich sogar linear ist.

7.7.1 Heapsort

Die Heap-Datenstruktur können wir auch verwenden, um Zahlenketten effizient zu sortieren. Dies veranschaulicht außerdem ein weiteres Entwurfsparadigma für Algorithmen, nämlich die Verwendung einer geeigneten Datenstruktur. Der resultierende Algorithmus Heapsort verbindet die Vorteile der bereits behandelten Algorithmen *Insertionsort* (Abschnitt 1.2) und *Mergesort* (Abschnitt 2.1). Wie auch Insertionsort benötigt er zusätzlich zur Eingabe nur konstant viel Speicherplatz. Solche Sortierverfahren werden auch als *in place* Verfahren bezeichnet. Darüber hinaus hat der Algorithmus eine worst-case Laufzeit von $O(n \log n)$, genau wie Mergesort. Er ist damit in puncto Laufzeit und Speicherplatz optimal.

Prozedur 17 $\text{Heapsort}(A)$

- 1: Baue Max-Heap für Array A ;
 - 2: **for** $i = A.\text{length}$ **to** 2 **do**
 - 3: vertausche $A[1]$ und $A[i]$;
 - 4: $A.\text{heap-size} = A.\text{heap-size} - 1$;
 - 5: $\text{Max-Heapify}(A, 1)$;
 - 6: **end for**
-

Zunächst baut Algorithmus 17 über dem Array A einen Max-Heap auf. Aufgrund der Heap-

Eigenschaft, steht der größte Wert anschließend in der Wurzel des Heaps, also $A[1]$. Solange der Heap noch mindestens 2 Werte enthält, vertauschen wir den größten Wert in der Wurzel mit dem Wert, welcher in der untersten Ebene am weitesten rechts steht. Diesen Knoten können wir durch Reduktion der Größe des Heaps um 1 einfach herauslöschen. Der Rest des Arrays, welcher nicht vom Heap belegt ist, ist bereits sortiert. Im linken und rechten Teilbaum der Wurzel ist dann noch immer die Heap-Eigenschaft erfüllt, nicht jedoch in der Wurzel. Um dies zu beheben, wird schließlich Algorithmus 15 auf dem Wurzelknoten aufgerufen.

7.7.2 Die Prioritätenwarteschlange

Als zweites Anwendungsbeispiel sehen wir im folgenden Abschnitt, wie wir Heaps zum Bau einer Prioritätenwarteschlange verwenden können. Eine Prioritätenwarteschlange ist eine Datenstruktur über einer Menge S von Elementen. Jedem Element x aus der Menge S ist zusätzlich ein Schlüsselwert (Priorität) zugeordnet. Wie auch bei den Heaps, unterscheiden wir zwei Typen, je nachdem, ob das Warten mit auf- oder absteigender Priorität erfolgt. Eine Max-Prioritätenwarteschlange unterstützt die folgenden Anfragen:

- $\text{MAX}(S)$: gibt das Element mit der höchsten Priorität (größter Schlüssel) zurück;
- $\text{EXTRACT-MAX}(S)$: entfernt das Element mit der höchsten Priorität (dem größten Schlüssel) und gibt es zurück;
- $\text{INSERT}(S, x)$: fügt das Element x zur Warteschlange S hinzu;
- $\text{INCREASE-KEY}(S, x, k)$: erhöht die Priorität des Elements x auf den Wert k (wobei $k > \text{Schlüssel}(x)$ gelten muss);

Diese Operationen erlauben es, Prioritätenwarteschlangen vielfältig einzusetzen. Beispielsweise bei der Abarbeitung von Rechenprozessen bzw. Rechenjobs auf mehreren CPUs oder Rechnern. Außerdem wird die Datenstruktur im Algorithmus von Dijkstra benötigt, wie wir später sehen werden.

Befassen wir uns nun damit, wie diese Operationen implementiert werden können. Sehr einfach ist es, das Element mit der höchsten Priorität in konstanter Zeit auszugeben:

Prozedur 18 $\text{Max}(A)$

1: RETURN $A[1]$;

Auch wie wir dieses Element aus dem Heap entfernen können, haben wir bereits in Sortieralgorithmus 17 gesehen. Wir ersetzen den Wurzelknoten durch den Knoten, der in der untersten Ebene am weitesten rechts steht. Anschließend lassen wir den Wurzelknoten absinken, bis die Heap-Eigenschaft erfüllt ist:

Prozedur 19 $\text{Extract-Max}(A)$

1: **if** $A.\text{heap-size} < 1$ **then**
 2: **error** "Heap ist leer"
 3: **end if**
 4: $max = A[1]$;
 5: $A[1] = A[A.\text{heap-size}]$;
 6: $A.\text{heap-size} = A.\text{heap-size} - 1$;
 7: $\text{Max-Heapify}(A, 1)$; RETURN max ;

Die Laufzeit von Extract-Max hängt im Wesentlichen vom Aufruf von Max-Heapify ab, da zusätzlich nur konstante Zeit benötigt wird. Somit liegt die Laufzeit von Extract-Max in $O(\log n)$.

Erhöhen wir den Schlüssel eines Knotens, ist die Heap-Eigenschaft nicht mehr unbedingt erfüllt. In diesem Fall können wir jedoch nicht die Heapify-Prozedur benutzen, da diese einen Knoten im Heap absinken lässt. Wir benötigen hingegen eine Prozedur, die den entsprechenden Knoten so lange im Heap aufsteigen lässt, bis die Heap-Eigenschaft erfüllt ist:

Prozedur 20 Increase-Key(A, i, key)

```

1: if  $key < A[i]$  then
2:   error "Aktueller Schlüssel ist kleiner als  $key$ "
3: end if
4:  $A[i] = key$ ;
5: while  $i > 1$  und  $A[\text{parent}(i)] < A[i]$  do
6:   vertausche  $A[i]$  und  $A[\text{parent}(i)]$ ;
7:    $i = \text{parent}(i)$ ;
8: end while

```

Der Beweis der Korrektheit der Prozedur ist eine gute Übung. Die Laufzeit hängt im Wesentlichen von der **while**-Schleife ab, in jedem Schleifendurchlauf sinkt jedoch der Abstand des Knotens i zur Wurzel des Heaps. Die Durchläufe der Schleife sind somit durch die Höhe des Heaps ($O(\log n)$) beschränkt.

Diese Prozedur, die einen Knoten im Heap aufsteigen lässt, können wir nun auch verwenden, um einen neuen Knoten einzufügen. Dazu legen wir zunächst einen Knoten mit dem Schlüsselwert $-\infty$ an. Anschließend erhöhen wir den Schlüssel dieses Knotens auf den gewünschten Wert und lassen ihn im Heap aufsteigen, bis die Heap-Eigenschaft erfüllt ist:

Prozedur 21 Insert(A, key)

```

1:  $A.\text{heap-size} = A.\text{heap-size} + 1$ ;
2:  $A[A.\text{heap-size}] = -\infty$ ;
3: Increase-Key( $A, A.\text{heap-size}, key$ );

```

Ein neuer Schlüssel kann also in $O(\log n)$ Zeit eingefügt werden. Insgesamt können wir also mit einem Heap jede der Operationen der Prioritätenwarteschlange in $O(\log n)$ Zeit ausführen.

7.8 Hashing

Wir haben gesehen, dass sich balancierte Bäume sehr gut dazu eignen, sortierte Schlüssel zu speichern. Für einen gegebenen Schlüssel findet sich so der zugehörige Datensatz in $O(\log n)$ Zeit. Das Einfügen und Entfernen eines neuen Objektes mit neuem Schlüssel hat die gleiche Zeitkomplexität.

Wie man sich leicht überlegt, gelten die gleichen Laufzeiten auch für die *mittlere* Laufzeit dieser Operationen. Wenn wir also eine Folge von n Operationen betrachten und $K(i)$ die Kosten der i -ten Operation beschreibt, so erhalten wir durch

$$\frac{1}{n} \sum_{i=1}^n K(i)$$

die mittleren Kosten der Operationen. Es kann sich dabei um verschiedene Operationen handeln.

Für das Einfügen und Entfernen liegen die mittleren Kosten in jedem Fall in $\Theta(\log n)$, da bei jedem Einfügen und bei jedem Entfernen die Mindestdiefe $\Omega(\log n)$ aus dem Beweis von Lemma 11 zu Durchlaufen ist und andererseits die Kosten nach den Betrachtungen des vorherigen Kapitels stets in $O(\log n)$ liegen können.

Für Datenbanken mit sehr großen Datenmengen verzichtet man beim erwähnten **Wörterbuch-Problem** auf die sortierte Speicherung der Schlüssel. Stattdessen sollen die Operationen (Exakte) Suche, Einfügen und Entfernen im Mittel nur konstante Zeit verursachen.

In solchen Fälle kommt das sogenannte *Hashing* zur Anwendung. Gegeben ist also eine Menge von Objekten (Daten) die über einen eindeutigen Schlüssel (ganze Zahl) identifizierbar sind und es wird eine Struktur zur Speicherung gesucht, die mindestens die Operationen

- (Exakte) Suche eines Objektes mit Schlüssel x
- Einfügen eines Objektes mit Schlüssel x
- Entfernen eines Objektes mit Schlüssel x

effizient (im Mittel konstant) ausgeführt werden können. Dafür wird das Schlüsseluniversum partitioniert. Die Position des Datenelementes im Speicher ergibt sich durch eine Berechnung aus dem Schlüssel.

Formal steht uns für die Verwaltung einer ungeordneten Menge S von Schlüsseln

1. ... ein Feld $T[0..m-1]$ für m Behälter $T[0], T[1], \dots, T[m-1]$ die sogenannte *Hashtabelle* und
2. ... eine *Hash-Funktion* $h : U \rightarrow [0..m-1]$ die das Universum der Schlüssel auf das Feld T abbildet

zur Verfügung. Ziel ist es, für eine Schlüsselmenge $S \subseteq U$ für jedes Element $x \in S$ das zugehörige Datenobjekt in $T[h(x)]$ zu speichern. Im Folgenden lassen wir der Übersichtlichkeit halber die Datenobjekte weg und speichern lediglich die Schlüssel in der Hashtabelle.

Beispiel: Sei $m = 5$ und $S = \{8, 22, 16, 29\}$ und $h(x) = x \bmod 5$. Die Hashtafel sieht dann wie folgt aus:

0	
1	16
2	22
3	8
4	29

Falls wir hier einen weiteren Schlüssel 21 verwenden gilt $h(21) = 1$ und wir sprechen dann von einer *Kollision*. Für eine Menge S von n Schlüsseln und für eine Menge von m Behältern $T[0], \dots, T[m-1]$ mit $n \leq m$ heißt die Hashfunktion h *perfekt für S* falls es keine Kollisionen bei der Zuordnung gibt. Die Hashfunktion soll effizient ausgewertet werden können.

7.8.1 Kollisionsbehandlung mit verketteten Listen

Eine sehr einfache Lösung in der Kollisionsbehandlung ist die Verwendung verketteter Listen in den Einträgen $T[i]$ der Hashtabelle. Doppelte Einträge werden dann an den Anfang der

Liste gesetzt und somit kann nach Auswertung der Hashfunktion in konstanter Zeit ein neuer Schlüssel eingefügt werden.

Beispiel: Für $m = 3$ betrachten wir die Schlüsselmenge $S = \{10, 22, 16, 29, 17, 3\}$ und die Hashfunktion $h(x) = x \bmod 3$. Die Abbildung beschreibt die Hashtabelle nach dem sukzessiven Einfügen der Schlüssel mit Kollisionsvermeidung mittels verketteter Listen.

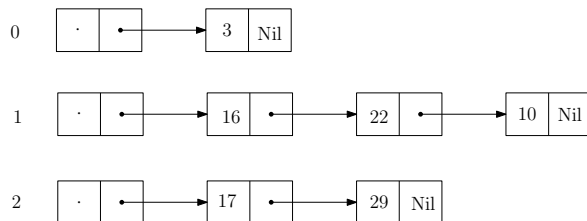


Abbildung 7.28: Die Kollisionsbehandlung durch die Verwendung von verketteten Listen $S = \{10, 22, 16, 29, 17, 3\}$ und $h(x) = x \bmod 3$.

Wir wollen nun die Verwendung von verketteten Listen für die Kollisionsbehandlung analysieren und dabei die mittleren Kosten berücksichtigen. Die Kosten einer Operation (Exakte) Suche und Entfernen ergeben sich jeweils durch das Berechnen des Hashwertes und den anschließenden linearen Suchdurchlauf durch die Liste des entsprechenden Eintrags.

Die Funktion $\delta_h(x, S)$ bezeichnet die Anzahl der Elemente $y \in S$ für die Funktion $h(x) = h(y)$ gilt. Dann kann eine Operation (Exakte) Suche und Entfernen für einen Schlüssel x in Zeit $O(1 + \delta_h(x, S))$ ausgeführt werden.

Wir betrachten dabei auch eine Folge von n möglichen zufälligen Operationen (Exakte) Suche, Einfügen und Entfernen und wollen die gesamte Laufzeit dafür nach oben abschätzen.

Wir machen dafür die folgende Annahmen:

1. Die Hashfunktion $h : U \rightarrow [0..m]$ streut das Universum gleichmäßig über das Intervall $[0..m]$. Das ist zum Beispiel für $h(x) = x \bmod m$ gegeben.
2. Sämtliche Elemente des Universums sind mit gleicher Wahrscheinlichkeit Argument der nächsten Operation.

Falls nun x_k der Schlüssel ist, der in der k -ten Aktion verwendet wird, dann folgt nun, dass die Wahrscheinlichkeit, dass $h(x_k) = i$ ist für $i \in [0..m]$ durch $\frac{1}{m}$ gegeben ist, kurz $\text{prob}(h(x_k) = i) = \frac{1}{m}$.

Theorem 19 Für die einzelnen Operationen (Exakte) Suche, Einfügen und Entfernen bezüglich einer Hashtabelle für die Schlüsselmenge S wird im worst-case $\Omega(|S|)$ Zeit benötigt.

Unter den obigen Annahmen ist der Erwartungswert für die von einer Folge von n Operationen (Exakte) Suche, Einfügen und Entfernen verursachten Kosten kleiner gleich $(1 + \frac{\alpha}{2})n$, wobei $\alpha = \frac{n}{m}$ der maximale Belegungsfaktor der Hashtafel ist.

Beweis. Der erste Teil des Theorems ist leicht zu sehen. Im schlechtesten Fall gilt für alle Schlüssel $x \in S$, dass $h(x)$ auf den gleichen Wert i_0 abbildet, dann muss im schlimmsten Fall bei jeder Operation eine Liste der Größe $|S|$ vollständig durchlaufen werden.

Für den zweiten Teil der Behauptung verwenden wir für die ersten k Operationen und für $i = 1, \dots, m$ einen Zähler $l_k(i)$, so dass $l_k(i)$ nach k Operationen in jedem Fall einen Wert größer gleich der dann aktuellen Listenlänge der Liste von $T_k[i]$ hat. Falls nun eine weitere

Operation die i -te Liste betrifft, wird $l_{k+1}(i) := l_k(i) + 1$ gesetzt und $l_k(i)$ ist dann stets eine obere Schranke der aktuellen Listenlänge nach k Operationen.

Anfangs sind alle Zähler auf 0 gesetzt. Immer wenn eine Operation die i -te Liste betrifft, erhöhen wir den Wert von $l_k(i)$ um Eins, unabhängig von der Art der Operation. Sei EK_{k+1} der Erwartungswert des Aufwandes der $(k+1)$ -ten Operation bei insgesamt n Operationen. Nehmen wir an, dass x_{k+1} der Schlüssel der $(k+1)$ -ten Operation ist, und es gelte $h(x_{k+1}) = i$. Sei nun $\text{prob}(l_k(i) = j)$ die Wahrscheinlichkeit, dass der Zähler $l_k(i)$ nach der k -ten Operation den Wert j hat. Dann ist

$$EK_{k+1} = 1 + \sum_{j=0}^k \text{prob}(l_k(i) = j) \cdot j,$$

da wir den Schlüssel auswerten müssen und relativ zur Wahrscheinlichkeit der Größe von $l_k(i) = j$ die entsprechende Länge j durchlaufen müssen.

Für $j \geq 1$ gilt nun aber

$$\text{prob}(l_k(i) = j) = \binom{k}{j} \left(\frac{1}{m}\right)^j \left(1 - \frac{1}{m}\right)^{k-j}.$$

Diese Aussage läßt sich leicht kombinatorisch erklären. Es gibt insgesamt $\binom{k}{j}$ j -elementige Sequenzen aus den ersten k -Operationen die zu beachten sind. Für jede dieser Sequenzen ist die Wahrscheinlichkeit, dass genau j -Operationen die i -te Liste betreffen $\left(\frac{1}{m}\right)^j \left(1 - \frac{1}{m}\right)^{k-j}$. Genauer, das Produkt der jeweiligen Wahrscheinlichkeit $\frac{1}{m}$, dass eine einzelne der j -Operation die i -te Liste betrifft mal dem Produkt der jeweiligen Wahrscheinlichkeit $\left(1 - \frac{1}{m}\right)$, dass eine einzelne der restlichen $k - j$ -Operationen *nicht* die i -te Liste betrifft.

Zunächst berechnen wir nun EK_{k+1} :

$$\begin{aligned} EK_{k+1} &= 1 + \sum_{j=0}^k \binom{k}{j} \left(\frac{1}{m}\right)^j \left(1 - \frac{1}{m}\right)^{k-j} \cdot j \\ &= 1 + \frac{k}{m} \sum_{j=1}^k \binom{k-1}{j-1} \left(\frac{1}{m}\right)^{j-1} \left(1 - \frac{1}{m}\right)^{k-j} \\ &\quad \left(\frac{k}{j} \text{ und } \frac{1}{m} \text{ Ausklammern}\right) \\ &= 1 + \frac{k}{m} \left(1 - \frac{1}{m}\right)^{-1} \sum_{j=1}^k \binom{k-1}{j-1} \left(\frac{1}{m}\right)^{j-1} \left(1 - \frac{1}{m}\right)^{k-(j-1)} \\ &\quad \left(\left(1 - \frac{1}{m}\right)^{-1} \text{ Ausklammern}\right) \\ &= 1 + \frac{k}{m} \frac{m}{m-1} \left(\frac{1}{m} + \left(1 - \frac{1}{m}\right)\right)^{k-1} \\ &\quad \left(\text{Binominalkoeffizienten}\right) \\ &= 1 + \frac{k}{m-1} \end{aligned}$$

Nun kann der Erwartungswert aller n Operation, $EK(n)$, folgendermaßen abgeschätzt werden:

$$\begin{aligned}
\text{EK}(n) &= \sum_{k=1}^n \text{EK}_k = \sum_{k=0}^{n-1} \text{EK}_{k+1} \\
&= \sum_{k=0}^{n-1} \left(1 + \frac{k}{m-1}\right) \\
&= n + \frac{1}{m-1} \sum_{k=0}^{n-1} k = n + \frac{1}{m-1} \frac{n(n-1)}{2} \\
&= n \cdot \left(1 + \frac{n-1}{2(m-1)}\right) < \left(1 + \frac{n}{2m}\right) \cdot n
\end{aligned}$$

Wobei wir für den letzten Schritt $n < m$ annehmen. Somit gilt für $\alpha = \frac{n}{m}$, dass $(1 + \frac{\alpha}{2})n$ eine obere Schranke für den erwarteten Aufwand der n Operationen darstellt. \square

Falls $\alpha = \frac{n}{m}$ konstant bleibt, sagt das obige Theorem, dass im Mittel nur konstant viel Aufwand zu erwarten ist, also eine deutliche Verbesserung im Vergleich zur AVL-Baumlösung.

Abschließend erwähnen wir noch, wie sich die Bedingung, dass $\alpha = \frac{n}{m}$ konstant bleibt gewährleisten können, ohne an Performanz zu verlieren.

Wir passen das m sukzessive an. Dazu verwenden wir eine Folge von Hashtafeln $T_0, T_1, T_2, \dots, T_i, \dots$ der Größen $m, 2m, \dots, 2^i m, \dots$ mit Hashfunktionen $h_i : U \rightarrow [0..2^i m - 1]$. Die Tafeln werden so verwendet, dass der aktuelle Belegungsfaktor α stets strikt zwischen $\frac{1}{4}$ und 1 liegt.

Die Tabelle T_i muss dann umkopiert werden, wenn der Belegungsfaktor aus dem Intervall herausfällt:

α wird zu 1: Wir speichern alle $2^i m$ Elemente von T_i nach T_{i+1} in $\Theta(2^i m)$ Zeit. Der Belegungsfaktor von T_{i+1} ist dann $\frac{1}{2}$. Bis zur nächsten Umspeicherung müssen mindestens $\frac{1}{4}2^{i+1}m$ Operationen durchgeführt werden bevor die nächste Umspeicherung notwendig wird.

α wird zu $\frac{1}{4}$: Wir speichern alle $\frac{1}{4}2^i m$ Elemente von T_i nach T_{i-1} in $\Theta(\frac{1}{4}2^i m)$ Zeit. Der Belegungsfaktor von T_{i-1} ist dann $\frac{1}{2}$ und mindestens $\frac{1}{4}2^{i-1}m$ Operationen können durchgeführt werden bevor die nächste Umspeicherung notwendig wird.

In beiden Fällen sind die Kosten für die Umspeicherung maximal zweimal so groß, wie die Anzahl der Operationen. Die Kosten für das Umspeichern verteilen wir deshalb auf die einzelnen Operationen. Der Erwartungswert $O(1)$ für die mittleren Kosten bleibt dann erhalten. Der Aufwand amortisiert sich hier über die Zeit. Auch wenn gelegentlich viele Elemente umkopiert werden müssen, so ist der Aufwand im Mittel gering.

7.9 Union-Find Datenstruktur

Neben der Speicherung einer Menge von Objekten ist es manchmal nützlich, ein System von Mengen zu speichern. Die Union-Find Datenstruktur dient zur Verwaltung einer Menge von disjunkten Mengen S_1, \dots, S_k . Ein Universum $U = \{x_1, x_2, \dots, x_n\}$ von Elementen ist gegeben und wir haben anfangs n disjunkte Mengen $S_i = \{x_i\}$. Über ein Array $U[1..n]$ können wir stets auf die Elemente durch $U[i] = x_i$ zugreifen. Jede Menge besitzt einen eindeutigen Repräsentanten $s_i \in S_i$. Im Folgenden beschreiben wir die Mengen der Einfachheit halber über die Indizes des Arrays U , das heißt $S_i = \{i\}$ und auch $s_i = i$.

Die Datenstruktur soll die folgenden Operationen effizient unterstützen:

- $\text{UNION}(x, y)$: Falls zwei Mengen S_x und S_y mit $x \in S_x$ und $y \in S_y$ existieren, so werden diese entfernt und durch die Menge $S_x \cup S_y$ ersetzt. Der neue Repräsentant von $S_x \cup S_y$ kann ein beliebiges Element dieser vereinigten Menge sein.
- $\text{FIND}(x)$: Liefert den Repräsentanten der Menge S mit $x \in S$ zurück.

Beispiel: Im folgenden Beispiel soll $x : A$ bedeuten, dass wir eine Menge A mit Repräsentant x gespeichert haben. Welches Element nach einer UNION-Operation Repräsentant der neuen Menge wird, ist nicht eindeutig bestimmt. Es wurde jeweils ein beliebiges Element dafür ausgewählt.

Operation	Zustand der Datenstruktur
Initialisierung	1 : {1}, 2 : {2}, 3 : {3}, 4 : {4}, 5 : {5}
FIND(3)	Ausgabe: 3, keine Zustandsänderung
UNION(1,2)	2 : {1, 2}, 3 : {3}, 4 : {4}, 5 : {5}
UNION(3,4)	2 : {1, 2}, 3 : {3, 4}, 5 : {5}
FIND(1)	Ausgabe: 2, keine Zustandsänderung
UNION(1,5)	2 : {1, 2, 5}, 3 : {3, 4}
FIND(3)	Ausgabe: 3, keine Zustandsänderung
UNION(5,4)	3 : {1, 2, 3, 4, 5}

Intern kann eine solche Struktur effizient wie folgt realisiert werden. Eine einzelne Menge S_i wird durch eine verkettete Liste dargestellt, wobei aber jedes einzelne Listenelement auch einen Zeiger auf den Kopf S_i der Liste speichert. Der Kopf der Liste enthält den Namen S_i , den Repräsentanten i und die Größe der Liste. Zusätzlich verwenden wir ein Feld $U[1..n]$ das für jedes Element j durch $U[j]$ einen Zeiger auf das Listenelement j in einer Liste S_i bereithält. Abbildung 7.29 illustriert die Situation für eine Menge $S_5 = \{1, 5, 4, 8\}$. Die Operation $\text{FIND}(x)$

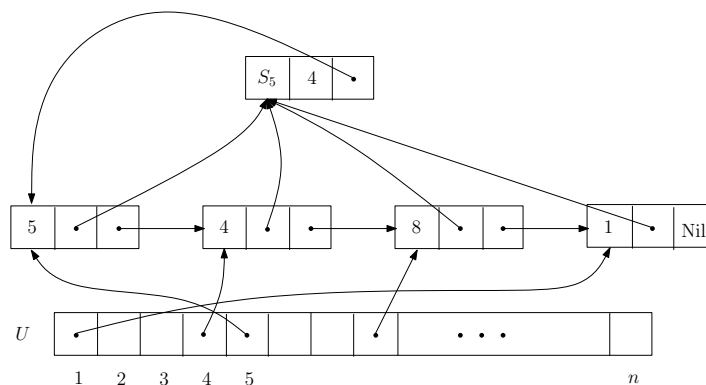


Abbildung 7.29: Die Zeigerstruktur einer Union-Find Realisierung der Menge $S_5 = \{1, 5, 4, 8\}$.

kann nun in $O(1)$ auf das Element x über $U[x]$ zugreifen und findet den Repräsentanten i von S_i durch einen Verweis.

Die Operation $\text{UNION}(x, y)$ kann wie folgt realisiert werden. Sei $|S_x| \geq |S_y|$.

- Jeder Namenszeiger von S_y wird auf den Namen von S_x gesetzt.
- Die Verkettete Liste von S_x wird an S_y angehängt. Der Mengenanzugszeiger von S_x wird auf den Anfang von S_y gesetzt. Die Größen werden addiert und im Kopfelement S_x eingetragen.

- Das alte Kopfelement S_y wird entfernt.

Es ergibt sich ein Aufwand von $O(|S_y|)$.

Theorem 20 *Eine Folge von m FIND-Operationen und $k < n$ UNION-Operationen lässt sich in Zeit $O(m + n \log n)$ ausführen.*

Beweis. Die Laufzeit $O(m)$ für die FIND-Operationen ist klar. Die abschließende Frage lautet, wie oft ein Namenszeiger eines Elementes x insgesamt umgehängt wird. Das Element gehörte dann immer zur kleineren Menge. Nach dem ersten Umhängen gehört x zu einer Menge mit größer gleich 2 Elementen. Nach dem zweiten Umhängen gehört x zu einer Menge mit größer gleich 4 Elementen. Nach dem i -ten Umhängen gehört x also zu einer Menge mit größer gleich 2^i Elementen. Eine solche Menge S_j kann aber nur maximal n Elemente enthalten. Also kann x nur maximal $\log n$ mal umgehängt worden sein. Das gilt für jedes $x \in U$ und somit ergibt sich die Laufzeit. \square

Gelegentlich wird in einigen Beschreibungen zur Union-Find Datenstruktur auch eine Operation MAKE-SET(x) angegeben, die eine Menge $S_x = \{x\}$ initialisiert. Wir haben hier diese Initialisierung vorweggenommen. An den Laufzeiten ändert sich nichts.

7.10 Datenstrukturen für Graphen

Im nächsten Kapitel wollen wir einige klassische Algorithmen betrachten, die auf einem Netzwerk von Knoten und Kanten angewendet werden und viele praktische Anwendungen haben.

Ein sogenannter *Graph* ist eine natürliche Erweiterung der Bäume. Formal beschreiben wir einen (ungerichteten) Graphen durch eine Menge $V = \{v_1, v_2, \dots, v_n\}$ von Knoten und eine Menge $E \subseteq V \times V$ von Kanten, kurz $G = (V, E)$. Jeder formale Graph G kann grafisch (oder geometrisch) in der Ebene analog zu den Bäumen realisiert werden. Für die Knoten $v_i \in V$ legen wir disjunkte Punkte in der Ebene fest und für je zwei Knoten v_1, v_2 in V mit $(v_1, v_2) \in E$ zeichnen wir einen einfachen (ohne Selbstschnitte) Weg zwischen den zugehörigen Orten. Falls wir eine bestimmte geometrische Realisation meinen, sprechen wir auch von einem *geometrischen Graphen*.

Nicht jeder Graph lässt sich so geometrisch realisieren, dass die Kantenwege paarweise schnittfrei bleiben. Graphen, die eine schnittfreie geometrische Realisation erlauben, heißen auch *planare Graphen*. Die Realisation selbst, also der geometrische Graph heißt dann auch *kreuzungsfrei*. Bei einem kreuzungsfreien geometrischen Graphen kann man in der Ebene eindeutige Flächen charakterisieren, die von Kanten eingeschlossen sind.

Wir erlauben auch Kanten der Form (v, v) , der *Grad* eines Knoten v ist die Anzahl der Kanten, die diesen Knoten als Endpunkt haben. Für *ungerichtete* Graphen ist die Kante (v_i, v_j) gleichbedeutend mit (v_j, v_i) . Bei *gerichteten* Graphen wird durch (v_i, v_j) eine Kante beschrieben, die von v_i zu v_j verläuft. Im allgemeinen kann E auch eine Multimenge sein, dass heißt, es können mehrere Kanten (v_i, v_j) in E enthalten sein. In diesem Fall wird der Graph auch als *Multigraph* bezeichnet.

Für $G = (\{v_1, v_2, v_3, v_4, v_5\}, \{(v_1, v_2), (v_1, v_3), (v_1, v_1), (v_2, v_3), (v_4, v_5), (v_3, v_4), (v_1, v_5)\})$ zeigt Abbildung 7.30 eine kreuzungsfreie geometrische Realisation von G . Wir interessieren uns hier zunächst darum, wie Graphen im allgemeinen so abgespeichert werden, dass wir effizient auf die notwendigen Informationen zugreifen können. Ein *Weg* in einem Graphen ist eine Folge von Knoten, die über Kanten miteinander verbunden sind. Bei gerichteten Graphen muss die Kantenfolge entsprechend *gerichtet* sein.

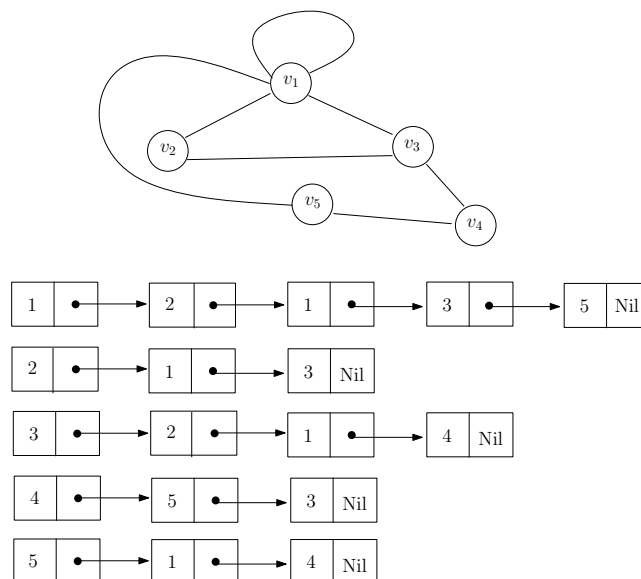


Abbildung 7.30: Eine geometrische Realisation des Graphen $G = (\{v_1, v_2, v_3, v_4, v_5\}, \{(v_1, v_2), (v_1, v_3), (v_1, v_1), (v_2, v_3), (v_4, v_5), (v_3, v_4), (v_1, v_5)\})$. Die Adjazenzliste speichert für jeden Knoten die adjazenten Knoten in einer Liste ab. Es gibt Crossreferenzen zu den jeweiligen Listen.

Beispielsweise möchten wir nun wissen, ob tatsächlich alle Knoten von einem Startknoten aus sukzessive über Kanten erreichbar sind. Wie stellen also die Frage, ob die Knoten untereinander *wegzusammenhängend* sind oder kurz ob der Graph insgesamt *zusammenhängend* ist. Zusammenhangskomponenten (maximale wegzusammenhängende Teilmengen) des Graphen werden über die Knoten definiert. Außerdem kann es sinnvoll sein, alle Knoten und Kanten (bzw. weitere Daten dazu) des Graphen sukzessive auszugeben.

7.10.1 Adjazenzlisten

Eine klassische Speichermethode ist die sogenannte *Adjazenzliste*. Falls für einen Graphen $G = (V, E)$ eine Kante (v_i, v_j) existiert, sind v_i und v_j *Nachbarknoten* und werden als *adjacent* bezeichnet.

Für jeden Knoten v gibt es in der Adjazenzliste eine verkettete Liste aller Knoten, die zu v adjacent sind, also eine Kante bilden, in Abbildung 7.30 ist ein Beispiel angegeben. Genauer handelt es sich um eine verkettete Liste von den zu v adjazenten Knoten.

Der Einfachheit halber können wir Knoten v_i auch durch Indizes $i = 1, \dots, n$ repräsentieren. Dann können wir von einem Knoten in einer Adjazenzliste leicht zu seiner Adjazenzliste gelangen und auch für jeden Knoten global eine Information abspeichern. Die Listenköpfe lassen sich dann beispielsweise durch ein Array $A[1..n]$ erreichen und die Information für die Knoten wird in $V[1..n]$ gespeichert. Das läßt sich aber genauso durch Zeiger und Namen v_i realisieren.

Offensichtlich benötigen wir für die Adjazenzliste einen Speicherplatz von $O(|V| + |E|)$. Das ist in dem Sinne optimal, dass unser Graph auch $|V| + |E|$ Objekte enthält. Wir stellen nun klassische Traversierungsmöglichkeiten für Graphen vor. Wir wollen alle Knoten des Graphen von einem Startknoten aus berichten.

Tiefensuche (DFS): Geht von einem Startknoten aus rekursiv in die *Tiefe* und besucht rekursiv den nächsten noch nicht besuchten Nachbarknoten. Falls kein solcher Knoten

mehr existiert wird zum letzten Ausgangsknoten der rekursiven Suche zurückgegangen (Backtracking). Dort wird rekursiv der nächste noch nicht besuchte Knoten in die Tiefe *exploriert*.

Breitensuche (BFS): Geht von einem Startknoten aus rekursiv in die *Breite* und besucht zunächst sukzessive durch Vorwärts- und Rückwärtsbewegung alle unbesuchten Nachbarknoten mit Abstand 1 zum Startknoten. Wenn alle Nachbarknoten besucht wurden, geht die Suche rekursiv beim ersten Nachbarknoten weiter.

Wir beschreiben das Verfahren DFS zunächst als Pseudocode und geben ein Beispiel an. Hierbei markieren wir die Knoten in der DFS Reihenfolge. Wir verwenden einen Stack S .

Prozedur 22 DFS(A, s)

Graph $G = (V, E)$ repräsentiert durch (kopierte) Adjazenzliste A . $A[v]$ ist Zeiger auf Nachbarschaftsliste von v . Vom Startknoten s aus, werden alle Knoten markiert.

```

1: Push( $s$ );
2: Markiere  $s$ 
3: Top( $v$ );
4: while  $v \neq \text{Nil}$  do
5:   if  $A[v] \neq \text{Nil}$  then
6:      $v' := \text{first}(A[v])$ ;
7:     Lösche  $v'$  aus  $A[v]$ 
8:     if  $v'$  nicht markiert then
9:       Push( $v'$ );
10:      Markiere  $v'$ 
11:    end if
12:  else
13:    Pop( $v$ );
14:  end if
15:  Top( $v$ );
16: end while

```

Wenn wir beispielsweise den Graphen aus Abbildung 7.30 betrachten, und den Knoten v_5 als Startknoten verwenden, dann können wir die Markierungen leicht in einem Array $M[1..5]$ festhalten. Die Abbildung 7.31 zeigt einen Zwischenstand bei der Durchführung des Algorithmus für Startknoten v_5 . Jeder Knoten wird nur einmal markiert und somit gibt es $|V|$ Push-Operationen. Da in jedem While-Durchlauf entweder ein Element aus einer nichtleeren Liste gestrichen wird, oder ein Element vom Stack gepoppt wird, liegt die Laufzeit des Algorithmus in $O(|V| + |E|)$.

7.10.2 Nachbarschaftsmatrix

Eine weitere simple Methode, um einen Graphen effizient abzuspeichern, ist die Verwendung einer Nachbarschaftsmatrix. Wir identifizieren die n Knoten des Graphen wiederum durch die ganzzahligen Indizes und für $V = \{1, 2, \dots, n\}$ verwenden wir eine $n \times n$ Matrix N mit Einträgen

$$N_{ij} = \begin{cases} 1 & \text{falls } (i, j) \in E \\ 0 & \text{sonst} \end{cases}$$

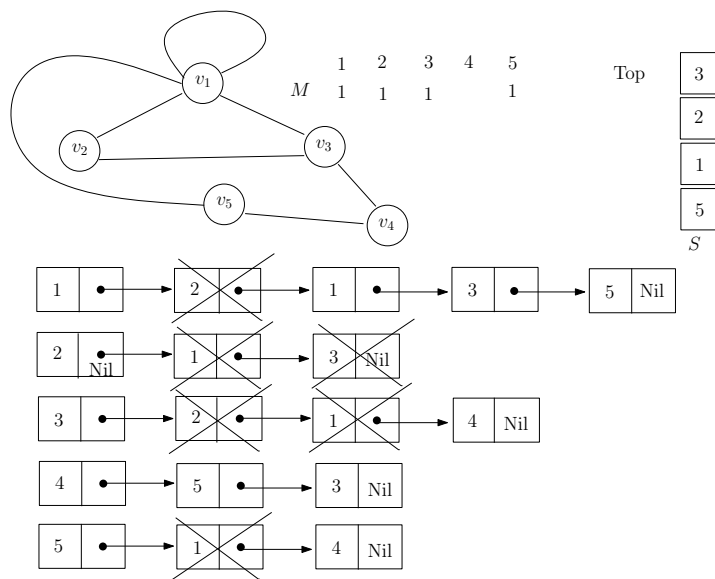


Abbildung 7.31: Ein Zwischenstand des DFS-Aufrufs für den Startknoten \$v_5\$. Im nächsten Durchlauf der While-Schleife wird das erste Element der Liste von \$v_3\$ betrachtet.

Die Matrix im Beispiel aus Abbildung 7.30 sieht dann wie folgt aus:

$$\begin{matrix}
 & 1 & 2 & 3 & 4 & 5 \\
 \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{pmatrix}
 \end{matrix}$$

Der Vorteil der Speicherung in einer Nachbarschaftsmatrix liegt darin, dass wir in konstanter Zeit testen können, ob eine Kante zwischen zwei Knoten existiert. Allerdings benötigen wir $\Omega(|V|^2)$ viel Speicherplatz.

Literatur

Eine andere Darstellung der hier betrachteten Datenstrukturen findet sich beispielsweise in [2].

Kapitel 8

Elementare Graphalgorithmen

Dieses Kapitel entstammt in wesentlichen Teilen aus dem Vorlesungsskript von Prof. Dr. Heiko Röglin, der die Vorlesung im Wintersemester 2010/2011 an der Universität Bonn gelesen hat und freundlicherweise den Quellcode zur Verfügung gestellt hat.

8.1 Minimale Spann bäume

Sei $G = (V, E)$ ein ungerichteter zusammenhängender Graph. Sei $w : E \rightarrow \mathbb{R}_{>0}$ eine *Gewichtung der Kanten*, d.h. eine Funktion, die jeder Kante $e \in E$ ein Gewicht $w(e)$ zuordnet. Ein Graph heißt *kreisfrei*, falls keine Folge von Kanten $(v_1, v_2)(v_2, v_3), \dots, (v_{n-1}, v_n)$ mit $v_1 = v_n$ existiert.

Definition 21 Eine Kantenmenge $T \subseteq E$ heißt *Spannbaum von G* , wenn der Graph $G = (V, T)$ zusammenhängend und kreisfrei ist. Wir erweitern die Funktion w auf Teilmengen von E , indem wir definieren

$$w(T) = \sum_{e \in T} w(e),$$

und wir nennen $w(T)$ das Gewicht von T . Ein Spannbaum $T \subseteq E$ heißt *minimaler Spannbaum von G* (kurz auch *MST für Minimum Spanning Tree*) wenn es keinen Spannbaum von G gibt, der ein kleineres Gewicht als T hat.

Die Berechnung eines minimalen Spannbaumes ist ein Problem, das in vielen Anwendungskontexten auftritt, da der MST den Graphen bezüglich des Zusammenhanges gut approximiert und eine geringere Komplexität hat. Für unser Anwendungsbeispiel der optimalen Platzierung von Agenten aus Kapitel 1 und Abschnitt 1.1.1 läßt sich beispielsweise zeigen, dass die beste Lösung für den MST des Graphen stets eine 2-Approximation der optimalen Lösung für den gesamten Graphen ergibt.

Wir entwerfen in diesem Abschnitt einen Algorithmus, der effizient für einen gegebenen Graphen einen minimalen Spannbaum berechnet.

8.1.1 Algorithmus von Kruskal

Mit Hilfe einer Union-Find-Datenstruktur können wir den *Algorithmus von Kruskal* zur Berechnung eines minimalen Spannbaumes angeben. Er ist nach Joseph Kruskal benannt, der diesen Algorithmus 1956 veröffentlichte.

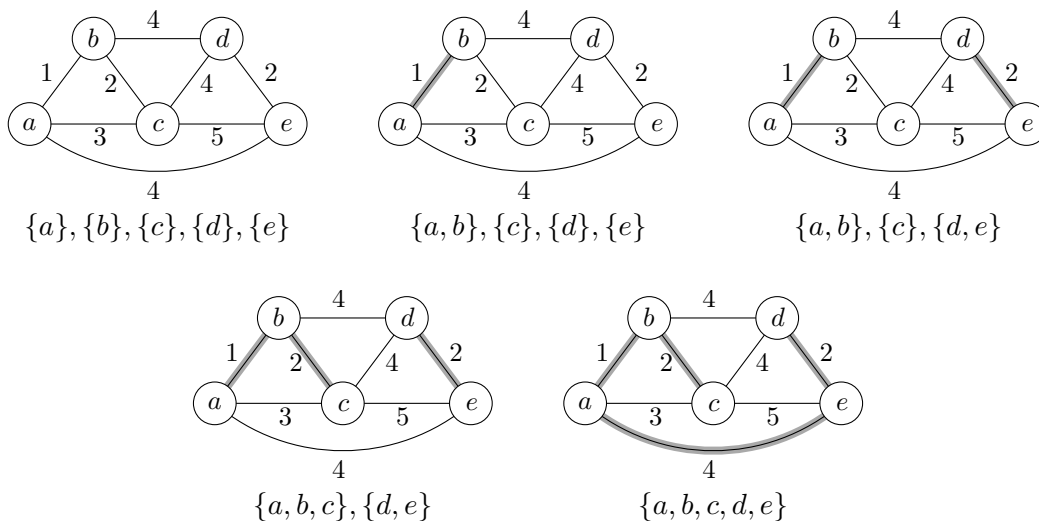
```

KRUSKAL( $G, w$ )
1  Teste mittels DFS, ob  $G$  zusammenhängend ist. Falls nicht Abbruch.
2  for all  $v \in V$  { MAKE-SET( $v$ ) }
3   $T = \emptyset$ 
4  Sortiere die Kanten in  $E$  gemäß ihrem Gewicht. Danach gelte  $E = \{e_1, \dots, e_m\}$  mit
    $w(e_1) \leq w(e_2) \leq \dots \leq w(e_m)$ . Außerdem sei  $e_i = (u_i, v_i)$ .
5  for  $i = 1$  to  $m$  {
6      if FIND( $u_i$ )  $\neq$  FIND( $v_i$ ) {
7           $T = T \cup \{e_i\}$ 
8          UNION( $u_i, v_i$ )
9      }
10 }
11 return  $T$ 

```

Die Mengen, die in der Union-Find-Datenstruktur im Algorithmus verwendet werden, entsprechen den Zusammenhangskomponenten des Graphen $G = (V, T)$. Wir gehen die Kanten in der Reihenfolge ihres Gewichtes durch und fügen eine Kante zur Menge T hinzu, wenn sie nicht innerhalb einer bereits zusammenhängenden Menge verläuft.

Bevor wir die Korrektheit beweisen und die Laufzeit analysieren, demonstrieren wir das Verhalten an einem kleinen Beispiel.



Korrektheit

Die Korrektheit des Algorithmus folgt aus dem folgenden Lemma.

Lemma 22 Sei V_1, \dots, V_k eine disjunkte Zerlegung von V . Sei $E_i \subseteq E$ die Menge der Kanten $(u, v) \in E$ mit $u, v \in V_i$ und sei $T_i \subseteq E_i$ ein Spannbaum des Graphen (V_i, E_i) .

Ferner sei $E' = E \setminus (E_1 \cup \dots \cup E_k)$, also

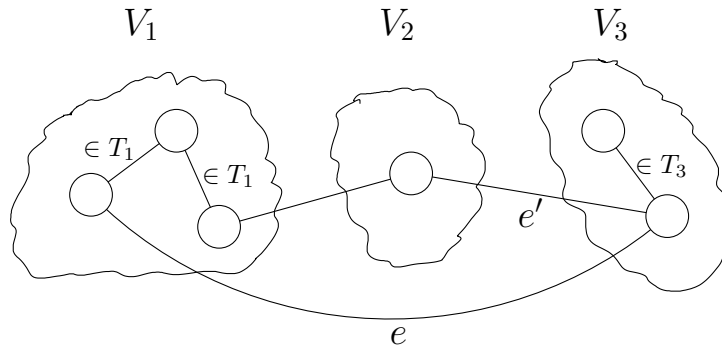
$$E' = \{(u, v) \in E \mid u \in V_i, v \in V_j, i \neq j\},$$

und sei e eine Kante mit dem kleinsten Gewicht aus E' .

Unter allen Spannbaum von $G = (V, E)$, die die Kanten aus $T_1 \cup \dots \cup T_k$ enthalten, gibt es einen mit kleinstem Gewicht, der zusätzlich auch die Kante e enthält.

Beweis. Sei T ein Spannbaum von $G = (V, E)$, der $T_1 \cup \dots \cup T_k$ aber nicht e enthält. Wir konstruieren einen Spannbaum T' , der $T_1 \cup \dots \cup T_k$ und e enthält und für den gilt $w(T') \leq w(T)$. Daraus folgt direkt das Lemma, denn angenommen wir haben einen Spannbaum T mit kleinstmöglichem Gewicht, der e nicht enthält, so können wir einen anderen Spannbaum T' konstruieren, der e enthält und ebenfalls das kleinstmögliche Gewicht hat.

Sei $T'' = T \cup \{e\}$. Dann enthält der Graph (V, T'') einen Kreis, auf dem e liegt. Da e zwei verschiedene Komponenten V_i und V_j mit $i \neq j$ verbindet, muss auf diesem Kreis noch eine weitere Kante $e' \in E'$ liegen, die zwei verschiedene Komponenten verbindet. Diese Situation ist in dem folgenden Bild dargestellt:



Aufgrund der Voraussetzung, dass e die billigste Kante aus E' ist, muss $w(e) \leq w(e')$ gelten. Wir betrachten nun die Menge $T' = T'' \setminus \{e'\} = (T \cup \{e\}) \setminus \{e'\}$.

Zunächst können wir festhalten, dass (V, T') zusammenhängend ist. Der Graph (V, T) ist nach Voraussetzung zusammenhängend. Gibt es in (V, T) zwischen zwei Knoten $u, v \in V$ einen Weg, der Kante e' nicht benutzt, so gibt es diesen Weg auch in (V, T') . Gibt es in (V, T) einen Weg P zwischen u und v , der Kante $e' = (w_1, w_2)$ benutzt, so erhält man auch in (V, T') einen Weg zwischen u und v , indem man P nimmt und die Kante e' durch den Weg zwischen w_1 und w_2 über die Kante e ersetzt.

Ebenso können wir argumentieren, dass (V, T') kreisfrei ist: Angenommen es gibt in (V, T') einen Kreis C . Wenn C die Kante e nicht enthält, so gibt es den Kreis C auch in (V, T) , im Widerspruch zur Definition von T . Enthält der Kreis die Kante $e = (u_1, u_2)$, so erhalten wir einen Kreis in (V, T) , indem wir in C die Kante e durch den Weg zwischen u_1 und u_2 über e' ersetzen.

Somit können wir festhalten, dass (V, T') ein Spannbaum ist. Außerdem gilt

$$w(T') = w(T) + w(e) - w(e') \leq w(T).$$

Damit ist das Lemma bewiesen. □

Wir können nun das folgende Theorem beweisen.

Theorem 23 *Der Algorithmus von Kruskal berechnet einen minimalen Spannbaum.*

Beweis. Wir benutzen die folgende Invariante, die aus drei Teilen besteht:

1. Die Mengen, die in der Union-Find-Datenstruktur gespeichert werden, entsprechen am Ende des Schleifenrumpfes der for-Schleife stets den Zusammenhangskomponenten des Graphen (V, T) , wobei T die in Zeilen 3 und 7 betrachtete Kantenmenge ist.
2. Die Zusammenhangskomponenten von (V, T) stimmen am Ende des Schleifenrumpfes der for-Schleife stets mit den Zusammenhangskomponenten von $(V, \{e_1, \dots, e_i\})$ überein.

- Die Kantenmenge $T \subseteq E$ kann am Ende des Schleifenrumpfes der for-Schleife stets zu einer Kantenmenge $T' \supseteq T$ erweitert werden, sodass T' ein minimaler Spannbaum von G ist. Insbesondere ist (V, T) azyklisch.

Ist diese Invariante korrekt, so ist auch der Algorithmus korrekt. Das sehen wir wie folgt: Am Ende des letzten Schleifendurchlaufes garantiert uns der zweite Teil der Invariante, dass die Zusammenhangskomponenten von (V, T) mit denen von $G = (V, E)$ übereinstimmen. Da G laut Voraussetzung zusammenhängend ist, ist also auch (V, T) zusammenhängend. Ist G keine korrekte Eingabe, da er nicht zusammenhängend ist, wird dies im ersten Schritt des Algorithmus erkannt.

Der dritte Teil der Invariante garantiert uns, dass T zu einer Kantenmenge $T' \supseteq T$ erweitert werden kann, sodass T' ein minimaler Spannbaum von G ist. Da (V, T) bereits zusammenhängend ist, würde die Hinzunahme jeder Kante einen Kreis schließen. Somit muss $T' = T$ gelten, woraus direkt folgt, dass $T = T'$ ein minimaler Spannbaum von G ist.

Nun müssen wir nur noch die Korrektheit der Invariante zeigen. Dazu überlegen wir uns zunächst, dass die Invariante nach dem ersten Schleifendurchlauf gilt. In diesem ersten Durchlauf fügen wir die billigste Kante e_1 in die bisher leere Menge T ein, da ihre beiden Endknoten u_1 und v_1 zu Beginn in verschiedenen einelementigen Mengen liegen. Nach dem Einfügen von e_1 vereinigen wir diese beiden einelementigen Mengen, so dass nach dem ersten Schleifendurchlauf die Knoten u_1 und v_1 in derselben und alle anderen Knoten in getrennten einelementigen Mengen liegen. Dies sind genau die Zusammenhangskomponenten von $(V, T) = (V, \{e_1\})$. Damit sind der erste und zweite Teil der Invariante gezeigt. Für den dritten Teil wenden wir Lemma 22 an. Wählen wir in diesem Lemma $V_1 = \{v_1\}, \dots, V_m = \{v_m\}$, wobei wir $V = \{v_1, \dots, v_m\}$ annehmen, so impliziert es direkt, dass wir die Menge $T = \{e_1\}$ zu einem minimalen Spannbaum von G erweitern können.

Schauen wir uns nun einen Schleifendurchlauf an und nehmen an, dass die Invariante am Ende des vorherigen Schleifendurchlaufs erfüllt war.

- Falls wir Kante e_i nicht in T einfügen, so ändern sich weder die Mengen in der Union-Find-Datenstruktur noch die Zusammenhangskomponenten von (V, T) . Teil eins der Invariante bleibt dann also erhalten. Fügen wir e_i in T ein, so fallen die Zusammenhangskomponenten von u_i und v_i in (V, T) zu einer gemeinsamen Komponente zusammen. Diese Veränderung bilden wir in der Union-Find-Datenstruktur korrekt ab. Also bleibt auch dann der erste Teil erhalten.
- Bezeichne V_1, \dots, V_k die Zusammenhangskomponenten von (V, T) am Ende des vorherigen Schleifendurchlaufs und bezeichne T_1, \dots, T_k die Kanten innerhalb dieser Komponenten. Falls wir e_i nicht zu T hinzufügen, so verläuft es innerhalb einer der gerade definierten Komponenten V_1, \dots, V_k . Das heißt die Komponenten von $(V, \{e_1, \dots, e_{i-1}\})$ und $(V, \{e_1, \dots, e_i\})$ sind identisch und die Invariante bleibt erhalten. Fügen wir e_i zu T hinzu, so verschmelzen in $(V, \{e_1, \dots, e_i\})$ die Komponenten, die u_i und v_i enthalten. Diese Verschmelzung passiert aber genauso in (V, T) , da wir e_i zu T hinzufügen. Also stimmen die Zusammenhangskomponenten von (V, T) wieder mit den Zusammenhangskomponenten von $(V, \{e_1, \dots, e_i\})$ überein.
- Aufgrund der Invariante können wir $T = T_1 \cup \dots \cup T_k$ zu einem minimalen Spannbaum von $G = (V, E)$ erweitern. Fügen wir in dieser Iteration e_i zu T hinzu, so garantiert Lemma 22, dass wir auch $T \cup \{e_i\}$ zu einem minimalen Spannbaum erweitern können: Da alle Kanten e_1, \dots, e_{i-1} aufgrund des zweiten Teils der Invariante innerhalb der Komponenten V_1, \dots, V_k verlaufen, ist e_i die billigste Kante, die zwei verschiedene Komponenten miteinander verbindet. Somit sind die Voraussetzungen von Lemma 22 erfüllt.

□

Laufzeit

Nachdem wir die Korrektheit des Algorithmus von Kruskal bewiesen haben, analysieren wir nun noch die Laufzeit. Wir haben bereits im letzten Kapitel gesehen, wie effizient wir die Operationen der Union-Find-Datenstruktur implementieren können, siehe Theorem 20.

Nun können wir die Laufzeit des Algorithmus von Kruskal abschließend bestimmen.

Theorem 24 *Für zusammenhängende ungerichtete Graphen mit m Kanten benötigt der Algorithmus von Kruskal Laufzeit $O(m \log m) = O(|E| \cdot \log |V|)$.*

Beweis. Die Tiefensuche im ersten Schritt benötigt Zeit $O(|V| + m)$. Da der Graph laut Voraussetzung zusammenhängend ist, gilt $m \geq |V| - 1$, also $O(|V| + m) = O(m)$. Das Sortieren der Kanten benötigt Zeit $O(m \log m)$. Die for-Schleife wird m mal durchlaufen und abgesehen von UNION und FIND werden in jedem Durchlauf nur konstant viele Operationen durchgeführt, so dass wir eine Laufzeit von $O(m)$ für alle Operationen in der Schleife abgesehen von UNION und FIND ansetzen können.

Wir führen in der Union-Find-Datenstruktur $2m$ FIND-Operationen und $m - 1$ UNION-Operationen durch. Mit Theorem 20 ergibt sich also eine Laufzeit von $O(m \log m + 2m) = O(m \log m)$. Wegen $|E| \leq |V|^2$ folgt $O(m \log m) = O(|E| \cdot \log |V|^2) = O(|E| \cdot \log |V|)$. Damit ist das Theorem bewiesen. □

Literatur

Der Algorithmus von Kruskal kann in Kapitel 23 von [1] nachgelesen werden. Kapitel 21 enthält zudem eine Beschreibung von Union-Find-Datenstrukturen.

8.2 Kürzeste Wege

Sei $G = (V, E)$ ein gerichteter Graph und sei $w : E \rightarrow \mathbb{R}$ eine Gewichtung der Kanten. Da wir uns in diesem Abschnitt mit kürzesten Wegen beschäftigen, werden wir $w(e)$ auch als die *Länge der Kante e* bezeichnen.

Definition 25 *Für zwei gegebene Knoten $s, t \in V$ ist $P = (v_0, v_1, \dots, v_\ell)$ ein Weg von s nach t , wenn $v_0 = s$, $v_\ell = t$ und wenn $(v_i, v_{i+1}) \in E$ für alle $i \in \{0, 1, \dots, \ell - 1\}$ gilt. Wir definieren die Länge von P als $w(P) = \sum_{i=0}^{\ell-1} w(v_i, v_{i+1})$. Wir sagen, dass P ein kürzester Weg von s nach t ist, falls es keinen anderen Weg P' von s nach t mit $w(P') < w(P)$ gibt. Wir nennen die Länge $w(P)$ des kürzesten Weges P die Entfernung von s nach t und bezeichnen diese mit $\delta(s, t)$. Existiert kein Weg von s nach t , so gelte $\delta(s, t) = \infty$.*

Zur besseren Intuition können wir uns im Folgenden vorstellen, dass die Knoten des Graphen Kreuzungen entsprechen und dass für jede Straße von Kreuzung $u \in V$ zu Kreuzung $v \in V$ eine Kante $e = (u, v)$ in E enthalten ist, wobei $w(e) = w((u, v))$ die Länge der Straße angibt. Es gibt jedoch auch zahlreiche andere Anwendungskontexte, in denen kürzeste Wege gefunden werden müssen. Insbesondere auch solche, in denen negative Kantengewichte auftreten können.

Wir sind an der Lösung der folgenden zwei Probleme interessiert:

1. Im *Single-Source Shortest Path Problem (SSSP)* ist zusätzlich zu G und w ein Knoten $s \in V$ gegeben und wir möchten für jeden Knoten $v \in V$ einen kürzesten Weg von s nach v und die Entfernung von s nach v berechnen.
2. Im *All-Pairs Shortest Path Problem (APSP)* sind nur G und w gegeben und wir möchten für jedes Paar $u, v \in V$ einen kürzesten Weg von u nach v und die Entfernung von u nach v berechnen.

Für den Spezialfall, dass $w(e) = 1$ für alle Kanten $e \in E$ gilt, haben wir mit Breitensuche bereits einen Algorithmus kennengelernt, der das SSSP löst. Breitensuche lässt sich allerdings nicht direkt auf den Fall verallgemeinern, dass wir beliebige Kantengewichte haben.

Bevor wir Algorithmen für das SSSP und das APSP beschreiben, schauen wir uns zunächst einige strukturelle Eigenschaften von kürzesten Wegen an.

Optimale Substruktur

Für die Algorithmen, die wir vorstellen werden, ist es wichtig, dass ein kürzester Weg von $v_0 \in V$ nach $v_\ell \in V$ aus kürzesten Wegen zwischen anderen Knoten zusammengesetzt ist.

Lemma 26 *Sei $P = (v_0, \dots, v_\ell)$ ein kürzester Weg von $v_0 \in V$ nach $v_\ell \in V$. Für jedes Paar i, j mit $0 \leq i < j \leq \ell$ ist $P_{ij} = (v_i, v_{i+1}, \dots, v_j)$ ein kürzester Weg von v_i nach v_j .*

Beweis. Wir zerlegen den Weg $P = P_{0\ell}$ in die drei Teilwege P_{0i} , P_{ij} und $P_{j\ell}$. Dies schreiben wir als

$$P = v_0 \xrightarrow{P_{0i}} v_i \xrightarrow{P_{ij}} v_j \xrightarrow{P_{j\ell}} v_\ell.$$

Insbesondere gilt $w(P) = w(P_{0i}) + w(P_{ij}) + w(P_{j\ell})$.

Nehmen wir nun an, dass P_{ij} kein kürzester Weg von v_i nach v_j ist. Dann gibt es einen Weg P'_{ij} von i nach j mit $w(P'_{ij}) < w(P_{ij})$. Dann erhalten wir ebenfalls einen anderen Weg P' von v_0 nach v_ℓ :

$$P' = v_0 \xrightarrow{P_{0i}} v_i \xrightarrow{P'_{ij}} v_j \xrightarrow{P_{j\ell}} v_\ell.$$

Für diesen Weg gilt $w(P') = w(P) - w(P_{ij}) + w(P'_{ij}) < w(P)$ im Widerspruch zur Voraussetzung, dass P ein kürzester Weg von v_0 nach v_ℓ ist. Also kann es keinen kürzeren Weg als P_{ij} von v_i nach v_j geben. \square

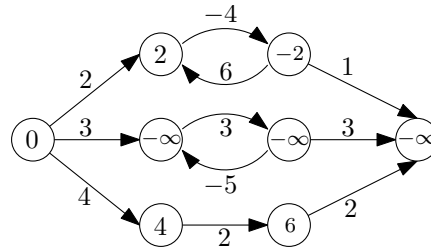
Daraus erhalten wir direkt das folgende Korollar.

Corollary 27 *Sei $P = (v_0, \dots, v_\ell)$ ein kürzester Weg von $v_0 \in V$ nach $v_\ell \in V$. Dann gilt*

$$\delta(v_0, v_\ell) = \delta(v_0, v_{\ell-1}) + w(v_{\ell-1}, v_\ell).$$

Negative Kantengewichte

Wir haben in der obigen Definition auch erlaubt, dass Kantengewichte negativ sein können. Dies kann in einigen Anwendungskontexten tatsächlich auftreten, es führt jedoch zu zusätzlichen Schwierigkeiten, insbesondere dann, wenn der Graph einen Kreis mit negativem Gesamtgewicht enthält. Dann ist nämlich für einige Paare von Knoten kein kürzester Weg mehr definiert, da man beliebig oft den Kreis mit negativem Gesamtgewicht entlang laufen könnte. Die folgende Abbildung zeigt einen Graphen mit einem negativen Kreis. Die Zahlen in den Knoten entsprechen dabei den Entfernungen der Knoten vom linken Knoten.



Der *Algorithmus von Dijkstra*, den wir für das SSSP kennenlernen werden, funktioniert nur für Graphen, die keine negativen Kantengewichte haben. Der *Floyd-Warshall-Algorithmus* für APSP funktioniert auch für negative Kantengewichte solange es keine negativen Kreise gibt. Gibt es einen negativen Kreis, so kann der Floyd-Warshall-Algorithmus zumindest benutzt werden, die Existenz eines solchen zu verifizieren.

8.2.1 Single-Source Shortest Path Problem

Es sei ein Graph $G = (V, E)$, eine Gewichtung $w : E \rightarrow \mathbb{R}_{\geq 0}$ und ein Startknoten $s \in V$ gegeben. In diesem Abschnitt gehen wir davon aus, dass alle Kantengewichte nicht-negativ sind.

Darstellung kürzester Wege

Zunächst halten wir fest, dass ein kürzester Weg in einem Graphen ohne negativen Kreis keinen Knoten mehrfach besuchen muss. Wäre $P = (v_0, \dots, v_\ell)$ nämlich ein kürzester Weg von $v_0 \in V$ nach $v_\ell \in V$ und gäbe es v_i und v_j mit $i < j$ und $v_i = v_j$, so wäre der Teilweg $P_{ij} = (v_i, \dots, v_j)$ ein Kreis, der laut Voraussetzung keine negativen Gesamtkosten hat. Das heißt, er hat entweder Kosten 0 oder echt positive Kosten. In jedem Falle können wir einen neuen Weg P' von v_0 nach v_ℓ konstruieren, indem wir den Kreis überspringen: $P' = (v_0, \dots, v_i, v_{j+1}, \dots, v_\ell)$. Für diesen neuen Weg P' gilt $w(P') \leq w(P)$. Somit können wir davon ausgehen, dass die kürzesten Wege zwischen zwei Knoten, die wir berechnen werden, aus höchstens $n - 1$ Kanten bestehen.

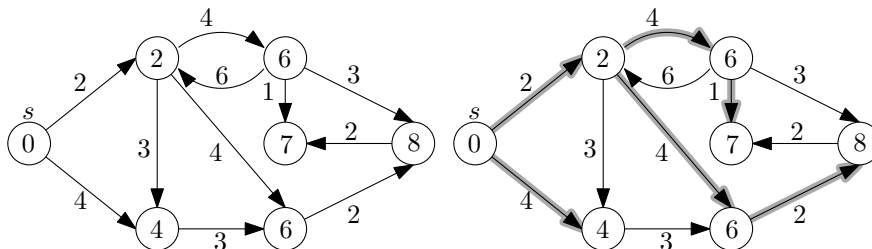
Beim SSSP berechnen wir n kürzeste Wege, nämlich einen vom Startknoten s zu jedem anderen Knoten. Beim APSP berechnen wir sogar $\binom{n}{2}$ kürzeste Wege. Da wir als obere Abschätzung für die Anzahl Kanten auf einem kürzesten Weg nur $n - 1$ haben, müssten wir damit rechnen, dass wir Platz $\Theta(n^2)$ bzw. $\Theta(n^3)$ brauchen, um alleine die Ausgabe aufzuschreiben. Wir wollen uns nun überlegen, wie man die Ausgabe effizienter kodieren kann. Dabei machen wir uns die Eigenschaft der optimalen Substruktur aus Lemma 26 zu Nutze.

Wir betrachten in diesem Abschnitt nur das SSSP. Statt alle n Wege separat zu speichern, wird der Algorithmus von Dijkstra einen Kürzeste-Wege-Baum mit Wurzel s gemäß der folgenden Definition berechnen.

Definition 28 Wir nennen (V', E') mit $V' \subseteq V$ und $E' \subseteq E$ einen Kürzeste-Wege-Baum mit Wurzel s , wenn die folgenden Eigenschaften erfüllt sind.

1. V' ist die Menge der Knoten, die von s aus in G erreichbar sind.
2. G' ist ein gewurzelter Baum mit Wurzel s . (Das bedeutet, dass G' azyklisch ist, selbst wenn wir die Richtung der Kanten ignorieren, und dass alle Kanten von s weg zeigen.)
3. Für alle $v \in V'$ ist der eindeutige Weg von s nach v in G' ein kürzester Weg von s nach v in G .

Kürzeste-Wege-Bäume benötigen Platz $O(n)$, da wir für jeden Knoten nur seinen Vorgänger speichern müssen. Sie sind also eine platzsparende Möglichkeit, um alle kürzesten Wege von s zu den anderen Knoten des Graphen zu kodieren. Die folgende Abbildung zeigt rechts einen Kürzeste-Wege-Baum für den Graphen auf der linken Seite.



Relaxation von Kanten

Der Algorithmus von Dijkstra verwaltet für jeden Knoten $v \in V$ zwei Attribute: $\pi(v) \in V \cup \{\text{nil}\}$ und $d(v) \in \mathbb{R}$. Der Knoten $\pi(v)$ wird am Ende des Algorithmus der Vorgänger von v in einem Kürzeste-Wege-Baum mit Wurzel s sein (sofern $\pi(v) \neq \text{nil}$) und die Zahl $d(v) \in \mathbb{R}$ wird am Ende die Entfernung von s nach v angeben. Wir initialisieren $\pi(v)$ und $d(v)$ wie folgt:

```

INITIALIZE-DIJKSTRA( $G, s$ )
1  for each  $v \in V$  {
2       $d(v) = \infty$ 
3       $\pi(v) = \text{nil}$ 
4  }
5   $d(s) = 0$ 

```

Im Algorithmus wird $d(v) \geq \delta(s, v)$ zu jedem Zeitpunkt und für jeden Knoten $v \in V$ gelten. Das heißt, $d(v)$ ist stets eine obere Schranke für die Länge des kürzesten Weges von s nach v . Direkt nach der Initialisierung ist dies sicherlich der Fall. Die Idee ist nun, die oberen Schranken Schritt für Schritt zu verbessern, bis irgendwann $d(v) = \delta(s, v)$ für alle Knoten v gilt. Dazu benutzen wir die folgende Beobachtung.

Lemma 29 Für $x, y, z \in V$ gilt stets

$$\delta(x, y) \leq \delta(x, z) + w(z, y).$$

Beweis. Wir konstruieren einen Weg von x nach y , indem wir von x auf einem kürzesten Weg nach z laufen und dann der direkten Kanten von z zu y folgen. Dieser Pfad hat Länge $\delta(x, z) + w(z, y)$. Da er offensichtlich nicht kürzer sein kann als der kürzeste Weg von x nach y , muss $\delta(x, z) + w(z, y) \geq \delta(x, y)$ gelten. \square

Aus dieser Beobachtung folgt direkt eine Möglichkeit, neue obere Schranken zu erhalten: Seien $u, v \in V$ und gelte $d(u) \geq \delta(s, u)$. Dann gilt

$$\delta(s, v) \leq d(u) + w(u, v).$$

Das nutzen wir, um die oberen Schranken mittels der folgenden Methode zu verbessern.

```
RELAX( $u, v$ )
```

```

1  if ( $d(v) > d(u) + w(u, v)$ ) {
2       $d(v) = d(u) + w(u, v)$ 
3       $\pi(v) = u$ 
4  }
```

Aus den bisherigen Überlegungen folgt direkt das folgende Lemma.

Lemma 30 *Nach einem Aufruf von INITIALIZE-DIJKSTRA gefolgt von einer beliebigen Sequenz von RELAX-Aufrufen gilt $d(v) \geq \delta(s, v)$ für alle Knoten v . Das heißt, der Wert $d(v)$ ist zu jedem Zeitpunkt eine obere Schranke für die Entfernung $\delta(s, v)$ von s nach v .*

Algorithmus von Dijkstra

Der Algorithmus von Dijkstra wird eine Folge von RELAX-Aufrufen durchführen und sicherstellen, dass am Ende tatsächlich $d(v) = \delta(s, v)$ für alle Knoten v gilt. In der folgenden Implementierung wird $S \subseteq V$ eine Menge von Knoten bezeichnen, sodass für jeden Knoten $v \in S$ bereits $d(v) = \delta(s, v)$ gilt. Es sei außerdem $Q = V \setminus S$. Wir gehen davon aus, dass der Graph G in Adjazenzlistendarstellung gegeben ist. Die Methode EXTRACT-MIN(Q) entfernt einen Knoten $u \in Q$ aus Q , der unter all diesen Knoten den kleinsten Wert $d(u)$ hat, und gibt diesen Knoten zurück.

```

DIJKSTRA( $G, w, s$ )
1  INITIALIZE-DIJKSTRA( $G, s$ )
2   $S = \emptyset; Q = V;$ 
3  while  $Q \neq \emptyset$  {
4       $u = \text{EXTRACT-MIN}(Q)$ 
5       $S = S \cup \{u\}$ 
6      for each  $v \in \text{Adj}[u]$  {
7          RELAX( $u, v$ )
8      }
9  }
```

Ein Beispiel für die Ausführung des Algorithmus von Dijkstra findet sich in Abbildung 8.1.

Theorem 31 *Der Algorithmus von Dijkstra terminiert auf gerichteten Graphen $G = (V, E)$ mit nicht-negativen Kantengewichten $w : E \rightarrow \mathbb{R}_{\geq 0}$ und Startknoten $s \in V$ in einem Zustand, in dem $d(v) = \delta(s, v)$ für alle $v \in V$ gilt.*

Beweis. Wir beweisen das Theorem mit Hilfe der folgenden Invariante, bei der wir die Konvention $w((u, v)) = \infty$ für $(u, v) \notin E$ anwenden.

Am Ende jeder Ausführung der **while**-Schleife in Zeilen 4–8 gilt:

1. $\forall v \in S : d(v) = \delta(s, v),$
2. $\forall v \in Q = V \setminus S : d(v) = \min\{\delta(s, x) + w(x, v) \mid x \in S\}.$

Gilt diese Invariante, so folgt aus ihr direkt das Theorem: In jedem Schritt der **while**-Schleife fügen wir einen weiteren Knoten aus Q zu der Menge S hinzu. Da zu jedem Zeitpunkt $S \cup Q = V$ gilt, muss nach n Durchläufen der **while**-Schleife $S = V$ gelten. Dann ist $Q = \emptyset$ und

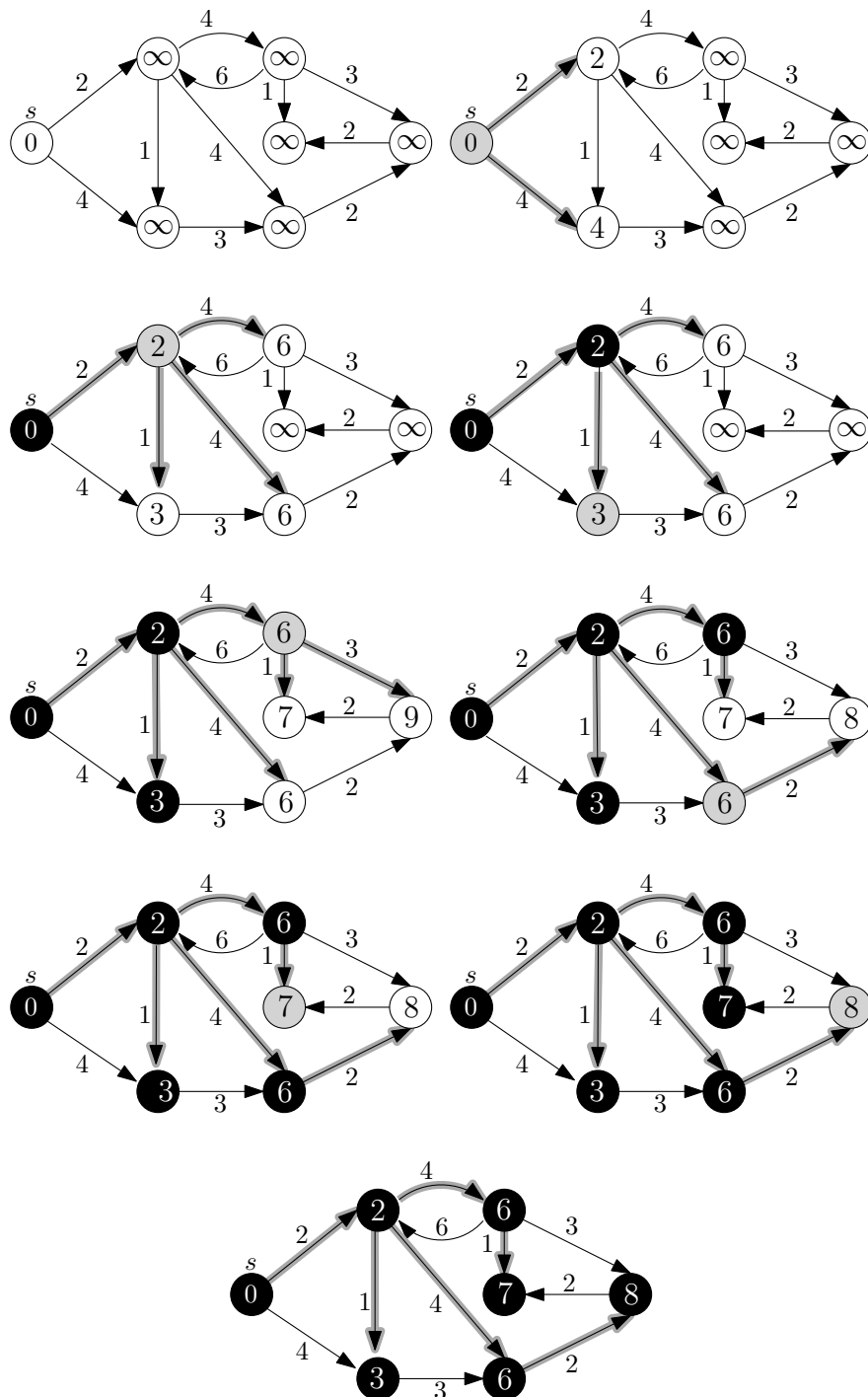


Abbildung 8.1: Beispiel für den Ablauf des Dijkstra-Algorithmus. Schwarze Knoten gehören zu S und der graue Knoten ist jeweils der Knoten u , der aus Q extrahiert wird. Eine graue Kanten (u, v) bedeutet, dass $\pi(v) = u$ gilt.

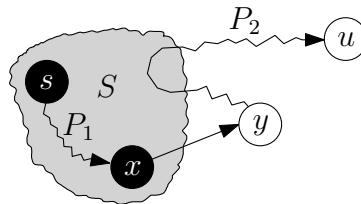
der Algorithmus terminiert. Die Invariante besagt dann, dass $d(v) = \delta(s, v)$ für alle Knoten $v \in S = V$ gilt.

Am Ende des ersten Schleifendurchlaufs gilt $S = \{s\}$ und $d(s) = \delta(s, s) = 0$. Somit ist der erste Teil der Invariante erfüllt. Den zweiten Teil der Invariante können wir in diesem Falle umschreiben zu: $\forall v \in V \setminus \{s\} : d(v) = \delta(s, s) + w(s, v) = w(s, v)$. Die Gültigkeit dieser Aussage folgt daraus, dass wir für jeden Knoten $v \in \text{Adj}[s]$ die Methode $\text{RELAX}(s, v)$ aufgerufen haben, die laut ihrer Definition $d(v) = \min\{d(s) + w(s, v), \infty\} = w(s, v)$ gesetzt hat.

Wir müssen uns nun nur noch überlegen, dass die Invariante erhalten bleibt. Gehen wir also davon aus, wir befinden uns in einem Durchlauf der **while**-Schleife und am Ende des vorherigen Durchlaufs galt die Invariante. Während dieser Iteration fügen wir den Knoten $u \in Q$, der momentan unter allen Knoten aus Q den kleinsten Wert $d(u)$ hat, zur Menge S hinzu.

Wir zeigen zunächst den ersten Teil der Invariante, indem wir zeigen, dass für diesen Knoten $d(u) = \delta(s, u)$ gilt. Sei dazu P ein kürzester Weg von s nach u und sei y der erste Knoten auf dem Pfad P , der nicht zur Menge S gehört. Da u selbst noch nicht zu S gehört, muss es einen solchen Knoten y geben. Weiterhin muss y einen Vorgänger haben, da der erste Knoten auf dem Weg P der Knoten $s \in S$ ist. Wir nennen diesen Vorgänger x . Den Teilweg von P von s nach x nennen wir P_1 und den Teilweg von y nach u nennen wir P_2 . Diese Teilwege können auch leer sein.

Die Notation ist in der folgenden Abbildung veranschaulicht.



Wegen Lemma 26 muss P_1 ein kürzester Weg von s nach x sein. Somit gilt

$$\delta(s, u) = w(P_1) + w(x, y) + w(P_2) = \delta(s, x) + w(x, y) + w(P_2) \geq d(y) + w(P_2) \geq d(y),$$

wobei wir bei der vorletzten Ungleichung den zweiten Teil der Invariante für $y \in Q$ ausgenutzt haben. Wegen Lemma 30 gilt $d(u) \geq \delta(s, u)$ und somit folgt insgesamt $d(u) \geq d(y)$. Da u ein Knoten aus Q mit kleinstem d -Wert ist, muss auch gelten $d(u) \leq d(y)$. Alles in allem erhalten wir die Ungleichungskette

$$d(u) \geq \delta(s, u) \geq d(y) \geq d(u),$$

deren einzige Lösung $d(u) = \delta(s, u)$ ist. Dies ist genau die Aussage, die wir für den ersten Teil der Invariante zeigen mussten.

Der zweite Teil der Invariante folgt nun einfach daraus, dass vorher $d(v) = \min\{\delta(s, x) + w(x, v) \mid x \in S\}$ für all $v \in Q$ galt und dass nun für jeden Knoten aus $v \in \text{Adj}[u]$ die Operation $\text{RELAX}(u, v)$ aufgerufen wird, die dazu führt, dass gilt:

$$\begin{aligned} d(v) &= \min\{\min\{\delta(s, x) + w(x, v) \mid x \in S\}, \delta(s, u) + w(u, v)\} \\ &= \min\{\delta(s, x) + w(x, v) \mid x \in S \cup \{u\}\}. \end{aligned}$$

□

Nun fehlt nur noch der Beweis, dass die Vorgänger-Relation, die in π berechnet wird, wirklich einen Kürzesten-Wege-Baum mit Wurzel s ergibt. Dies folgt unmittelbar aus dem folgenden Lemma, welches wir in der Vorlesung jedoch nicht beweisen werden.

Lemma 32 *Es gelte nach einem Aufruf von INITIALIZE-DIJKSTRA gefolgt von einer Sequenz von RELAX-Aufrufen $d(v) = \delta(s, v)$ für jeden Knoten v . Sei*

$$V_\pi = \{v \in V \mid \pi(v) \neq \text{nil}\} \cup \{s\} \quad \text{und} \quad E_\pi = \{(\pi(v), v) \in E \mid v \in V_\pi \setminus \{s\}\}.$$

Dann ist der Graph (V_π, E_π) ein Kürzeste-Wege-Baum mit Wurzel s .

Implementierung und Laufzeit Wir werden nun die Laufzeit des Algorithmus von Dijkstra analysieren. Dazu müssen wir uns zunächst überlegen, in welcher Datenstruktur wir die Menge Q verwalten. Diese Datenstruktur muss eine Menge verwalten können und sollte die folgenden drei Operationen möglichst effizient unterstützen:

- INSERT(Q, x, d): Füge ein neues Element x mit Schlüssel $d \in \mathbb{R}$ in die Menge Q ein.
- EXTRACT-MIN(Q): Entferne aus Q ein Element mit dem kleinsten Schlüssel und gib dieses Element zurück.
- CHANGE-KEY(Q, x, d_1, d_2): Ändere den Schlüssel des Objektes $x \in Q$ von d_1 auf d_2 .

Eine solche Datenstruktur heißt *Min-Priority Queue* und sie lässt sich leicht durch z. B. eine verkettete Liste realisieren, wenn wir auf Effizienz keinen Wert legen. Sei n die Anzahl der Einträge der Priority Queue, dann würde eine Realisierung als verkettete Liste Zeit $\Theta(n)$ für EXTRACT-MIN und DECREASE-KEY benötigen. Dies wollen wir nun verbessern.

Wir haben mit AVL-Bäumen schon eine Möglichkeit in der Vorlesung kennengelernt, um Priority Queues effizienter zu realisieren. Wir legen dafür einen AVL-Baum an und fügen bei INSERT(Q, x, d) einfach das Objekt x mit Schlüssel d in den AVL-Baum ein. Wir müssen auf den Fall aufpassen, dass wir zwei Objekte x und y mit dem gleichen Schlüssel einfügen wollen. Wir nehmen für diesen Fall an, dass es eine inhärente Ordnung auf den Objekten gibt und dass bei eigentlich gleichem Schlüssel d diese inhärente Ordnung entscheidet, welcher Schlüssel als kleiner betrachtet wird. In unserem Beispiel sind die Objekte in Q die Knoten aus V . Wir können davon ausgehen, dass diese mit v_1, \dots, v_n durchnummeriert sind und falls wir v_i und v_j mit gleichem Schlüssel d in Q einfügen, so wird v_i als kleiner als v_j betrachtet, wenn $i < j$ gilt. Auf diese Weise ist stets eine eindeutige Ordnung auf den Schlüsseln gegeben. Da die Höhe eines AVL-Baumes mit n Knoten durch $O(\log(n))$ beschränkt ist, können wir INSERT(Q, x, d) in Zeit $O(\log(n))$ ausführen.

EXTRACT-MIN(Q) können wir in einem AVL-Baum realisieren, indem wir ausgehend von der Wurzel immer dem Zeiger zum linken Kind folgen. Auf diese Weise finden wir das Objekt mit dem kleinsten Schlüssel. Auch dies geht wegen der beschränkten Höhe in Zeit $O(\log(n))$. Haben wir es einmal gefunden, so können wir es in Zeit $O(\log n)$ löschen und zurückgeben. Ebenso können wir DECREASE-KEY in Zeit $O(\log(n))$ realisieren, indem wir zunächst nach x suchen (dafür ist wichtig, dass wir seinen momentanen Schlüssel d_1 übergeben bekommen), es dann löschen und mit dem neuen Schlüssel d_2 wieder einfügen.

Theorem 33 *Die Laufzeit des Algorithmus von Dijkstra beträgt $O((n + m) \log n)$.*

Beweis. Die Initialisierung erfolgt in Zeit $O(n)$. Da die Adjazenzliste eines Knoten $v \in V$ dann das einzige Mal durchlaufen wird, wenn der Knoten v der Menge S hinzugefügt wird, wird für jede Kante $e = (u, v) \in E$ genau einmal die Operation RELAX(u, v) aufgerufen. Innerhalb des Aufrufs von RELAX(u, v) wird gegebenenfalls der Wert $d(v)$ reduziert. Dies entspricht einem Aufruf von CHANGE-KEY, welcher Laufzeit $O(\log n)$ benötigt. Alle Aufrufe der RELAX-Methode zusammen haben somit eine Laufzeit von $O(m \log n)$. Wir benötigen außerdem genau n Aufrufe von EXTRACT-MIN. Diese haben eine Gesamtlaufzeit von $O(n \log n)$. Damit folgt das Theorem. \square

8.2.2 All-Pairs Shortest Path Problem

In diesem Abschnitt beschäftigen wir uns mit dem APSP. Dazu sei ein gerichteter Graph $G = (V, E)$ mit Kantengewichten $w : E \rightarrow \mathbb{R}$ gegeben. Im Gegensatz zum Dijkstra-Algorithmus funktioniert der Floyd-Warshall-Algorithmus, den wir gleich beschreiben werden, auch, wenn es negative Kantengewichte gibt. Enthält der Graph einen negativen Kreis, so stellt der Algorithmus dies fest und bricht ab.

Für den Floyd-Warshall-Algorithmus gehen wir davon aus, dass $V = \{1, \dots, n\}$ gilt, und wir untersuchen Wege, die nur bestimmte Knoten benutzen dürfen. Sei $P = (v_0, \dots, v_\ell)$ ein kürzester Weg von v_0 nach v_ℓ . Wir haben bereits im SSSP argumentiert, dass es für jedes Paar von Knoten u und v in einem Graphen ohne negativen Kreis stets einen kürzesten Weg von u nach v gibt, der keinen Knoten mehrfach besucht. Wir nehmen an, dass dies auf P zutrifft, das heißt, die Knoten v_0, \dots, v_ℓ seien paarweise verschieden. Wir nennen einen solchen Weg auch *einfachen Weg*. Dann nennen wir $v_1, \dots, v_{\ell-1}$ die *Zwischenknoten von P* . Wir schränken nun die Menge der erlaubten Zwischenknoten ein.

Wir betrachten für jedes Paar $i, j \in V$ alle einfachen Wege von i nach j , die nur Zwischenknoten aus der Menge $\{1, \dots, k\}$ für ein bestimmtes k benutzen dürfen. Sei P_{ij}^k der kürzeste solche Weg von i nach j und sei $\delta_{ij}^{(k)}$ seine Länge. Wir überlegen uns, wie wir den Weg P_{ij}^k und $\delta_{ij}^{(k)}$ bestimmen können, wenn wir für jedes Paar $i, j \in V$ bereits den Weg P_{ij}^{k-1} kennen, also den kürzesten Weg von i nach j , der als Zwischenknoten nur Knoten aus $\{1, \dots, k-1\}$ benutzen darf. Hierzu unterscheiden wir zwei Fälle:

- Enthält der Weg P_{ij}^k den Knoten k nicht als Zwischenknoten, so enthält er nur die Knoten aus $\{1, \dots, k-1\}$ als Zwischenknoten und somit ist P_{ij}^k auch dann ein kürzester Weg, wenn wir nur diese Knoten als Zwischenknoten erlauben. Es gilt dann also $\delta_{ij}^{(k)} = \delta_{ij}^{(k-1)}$ und $P_{ij}^k = P_{ij}^{k-1}$.
- Enthält der Weg P_{ij}^k den Knoten k als Zwischenknoten, so zerlegen wir ihn wie folgt in zwei Teile:

$$P_{ij}^k = i \xrightarrow{P} k \xrightarrow{P'} j.$$

Da P_{ij}^k ein einfacher Weg ist, kommt k nur einmal vor und ist somit nicht in den Teilwegen P und P' enthalten. Diese müssen laut Lemma 26 kürzeste Wege von i nach k bzw. von k nach j sein, die nur Knoten aus $\{1, \dots, k-1\}$ als Zwischenknoten benutzen. Somit gilt in diesem Fall

$$\delta_{ij}^{(k)} = \delta_{ik}^{(k-1)} + \delta_{kj}^{(k-1)}.$$

und

$$P_{ij}^k = i \xrightarrow{P^{k-1}} k \xrightarrow{P_{kj}^{k-1}} j.$$

Aus dieser Beobachtung können wir direkt rekursive Formeln zur Berechnung der Entfernungen $\delta^{(k)}(i, j)$ und der Wege P_{ij}^k herleiten. Wir betrachten zunächst nur die Entfernungen. Für alle $i, j \in V$ gilt

$$\delta_{ij}^{(k)} = \begin{cases} w(i, j) & \text{für } k = 0 \\ \min\{\delta_{ij}^{(k-1)}, \delta_{ik}^{(k-1)} + \delta_{kj}^{(k-1)}\} & \text{für } k > 0. \end{cases}$$

Für $k = n$ erlauben wir alle Knoten $\{1, \dots, n\}$ als Zwischenknoten. Demzufolge besteht für $n = k$ keine Einschränkung mehr und es gilt $\delta(i, j) = \delta_{ij}^{(n)}$ für alle $i, j \in V$.

Wir weisen an dieser Stelle nochmals explizit darauf hin, dass wir für diese Rekursionsformel essentiell ausgenutzt haben, dass es im Graphen G keinen Kreis mit negativem Gesamtgewicht

gibt. Gibt es nämlich einen solchen Kreis, der von einem Knoten $i \in V$ erreichbar ist und von dem aus der Knoten $j \in V$ erreichbar ist, so ist der kürzeste Weg von i nach j nicht mehr einfach. Er würde stattdessen den Kreis unendlich oft durchlaufen und demzufolge die Knoten auf dem Kreis mehrfach besuchen.

Floyd-Warshall-Algorithmus

Der Floyd-Warshall-Algorithmus nutzt die rekursive Darstellung der $\delta_{ij}^{(k)}$, die wir soeben hergeleitet haben, um das APSP zu lösen. Zur Vereinfachung der Notation gehen wir davon aus, dass $V = \{1, \dots, n\}$ gilt und dass der Graph als $(n \times n)$ -Adjazenzmatrix $W = (w_{ij})$ gegeben ist. Die Adjazenzmatrix W enthält bei gewichteten Graphen an der Stelle w_{ij} das Gewicht der Kante (i, j) . Existiert diese Kante nicht, so gilt $w_{ij} = \infty$. Des Weiteren gilt $w_{ii} = 0$ für alle $i \in V$.

```

FLOYD-WARSHALL( $W$ )
1   $D^{(0)} = W$ 
2  for  $k = 1$  to  $n$  {
3      Erzeuge  $(n \times n)$ -Nullmatrix  $D^{(k)} = (\delta_{ij}^{(k)})$ 
4      for  $i = 1$  to  $n$  {
5          for  $j = 1$  to  $n$  {
6               $\delta_{ij}^{(k)} = \min\{\delta_{ij}^{(k-1)}, \delta_{ik}^{(k-1)} + \delta_{kj}^{(k-1)}\}$ 
7          }
8      }
9  }
10 return  $D^{(n)}$ 

```

Theorem 34 *Der Floyd-Warshall-Algorithmus löst das APSP für Graphen ohne Kreise mit negativem Gesamtgewicht, die als Adjazenzmatrix dargestellt sind und n Knoten enthalten, in Zeit $\Theta(n^3)$.*

Beweis. Die Korrektheit des Algorithmus folgt direkt aus der rekursiven Formel für $\delta_{ij}^{(k)}$, die wir oben hergeleitet haben, da der Algorithmus exakt diese rekursive Formel berechnet. Die Laufzeit von $\Theta(n^3)$ ist leicht ersichtlich, da wir drei ineinander geschachtelte **for**-Schleifen haben, die jeweils über n verschiedene Werte zählen, und innerhalb der innersten Schleifen nur Operationen ausgeführt werden, die konstante Zeit benötigen. Das Erstellen der Nullmatrizen benötigt insgesamt ebenfalls Zeit $\Theta(n^3)$. \square

Wenn wir zusätzlich zu den Entfernungen auch die kürzesten Wege ermitteln wollen, so können wir dies auch wieder mittels einer rekursiv definierten Vorgängerfunktion π machen. Sei $\pi_{ij}^{(k)}$ der Vorgänger von j auf dem Pfad P_{ij}^k , das heißt auf dem kürzesten Pfad von i nach j , der als Zwischenknoten nur Knoten aus $\{1, \dots, k\}$ benutzen darf. Da für $k = 0$ überhaupt keine Zwischenknoten erlaubt sind, gilt für alle Knoten $i, j \in V$

$$\pi_{ij}^{(0)} = \begin{cases} \text{nil} & \text{falls } i = j \text{ oder } w_{ij} = \infty, \\ i & \text{falls } i \neq j \text{ und } w_{ij} < \infty. \end{cases}$$

Das heißt, $\pi_{ij}^{(0)}$ ist genau dann gleich i , wenn es eine direkte Kante von i nach j gibt. Für $k > 0$ orientieren wir uns an der rekursiven Formel für $\delta_{ij}^{(k)}$, die wir oben hergeleitet haben.

Ist $\delta_{ij}^{(k)} = \delta_{ij}^{(k-1)}$, so ist auch $P_{ij}^k = P_{ij}^{k-1}$ und somit ist der Vorgänger von j auf dem Pfad P_{ij}^k derselbe wie der auf dem Pfad P_{ij}^{k-1} . Ist $\delta_{ij}^{(k)} < \delta_{ij}^{(k-1)}$, so gilt

$$P_{ij}^k = i \xrightarrow{P_{ik}^{k-1}} k \xrightarrow{P_{kj}^{k-1}} j,$$

und damit ist insbesondere der Vorgänger von j auf dem Pfad P_{ij}^k derselbe wie auf dem Pfad P_{kj}^{k-1} . Zusammenfassend halten wir für $k > 0$ fest

$$\pi_{ij}^{(k)} = \begin{cases} \pi_{ij}^{(k-1)} & \text{falls } \delta_{ij}^{(k-1)} \leq \delta_{ik}^{(k-1)} + \delta_{kj}^{(k-1)}, \\ \pi_{kj}^{(k-1)} & \text{falls } \delta_{ij}^{(k-1)} > \delta_{ik}^{(k-1)} + \delta_{kj}^{(k-1)}. \end{cases}$$

Die $\pi_{ij}^{(k)}$ können zusammen mit den $\delta_{ij}^{(k)}$ im Floyd-Warshall-Algorithmus berechnet werden, ohne dass sich die Laufzeit signifikant ändert.

Negative Kreise Wir wollen nun noch kurz diskutieren, wie man den Floyd-Warshall-Algorithmus benutzen kann, um die Existenz eines Kreises mit negativem Gesamtgewicht zu detektieren. Dazu müssen wir uns nur überlegen, was die Einträge $\delta_{ii}^{(k)}$ auf der Diagonalen der Matrix $D^{(k)}$ bedeuten. Für $k = 0$ werden diese Einträge alle mit 0 initialisiert. Gibt es nun für $k > 0$ einen Kreis mit negativem Gesamtgewicht, der den Knoten i enthält und sonst nur Knoten aus der Menge $\{1, \dots, k\}$, so gilt $\delta_{ii}^{(k)} < 0$ und somit auch $\delta_{ii}^{(n)} < 0$, da es einen Pfad negativer Länge von i nach i gibt. Also genügt es, am Ende des Algorithmus, die Einträge auf der Hauptdiagonalen zu betrachten. Ist mindestens einer von diesen negativ, so wissen wir, dass ein Kreis mit negativem Gesamtgewicht existiert. In diesem Fall, gibt die vom Algorithmus berechnete Matrix $D^{(n)}$ nicht die korrekten Abstände an.

Vergleich mit Algorithmus von Dijkstra Wir haben weiter oben bereits angesprochen, dass das APSP auch dadurch gelöst werden kann, dass wir einen Algorithmus für das SSSP für jeden möglichen Startknoten einmal aufrufen. Hat der Algorithmus für das SSSP eine Laufzeit von $O(f)$, so ergibt sich zur Lösung des APSP eine Laufzeit von $O(nf)$. Für den Dijkstra-Algorithmus ergibt sich also eine Laufzeit von $O(nm \log n)$.

Für Graphen mit nicht-negativen Kantengewichten stellt sich somit die Frage, ob es effizienter ist, den Floyd-Warshall-Algorithmus mit Laufzeit $O(n^3)$ oder den Dijkstra-Algorithmus mit Lösung $O(nm \log n)$ zur Lösung des APSP zu nutzen. Diese Frage lässt sich nicht allgemein beantworten. Für dichte Graphen mit $m = \Theta(n^2)$ vielen Kanten ist zum Beispiel der Floyd-Warshall-Algorithmus schneller. Für dünne Graphen mit $m = O(n)$ Kanten ist hingegen die n -malige Ausführung des Dijkstra-Algorithmus schneller.

Literatur

Das SSSP und APSP werden in den Kapiteln 24 und 25 von [1] beschrieben. Außerdem sind sie in den Kapiteln 4.3.1 und 4.3.3 von [2] erklärt.

8.3 Flussprobleme

In diesem Abschnitt lernen wir mit Flussproblemen eine weitere wichtige Klasse von algorithmischen Graphproblemen kennen, die in vielen Bereichen Anwendung findet. Ein *Flussnetzwerk* ist ein gerichteter Graph $G = (V, E)$ mit zwei ausgezeichneten Knoten $s, t \in V$ und einer

Kapazitätsfunktion $c : V \times V \rightarrow \mathbb{N}_0$, die jeder Kante $(u, v) \in E$ eine ganzzahlige nicht-negative Kapazität $c(u, v)$ zuweist. Wir definieren $c(u, v) = 0$ für alle $(u, v) \notin E$. Außerdem nennen wir den Knoten s die *Quelle* und den Knoten t die *Senke*. Wir definieren $n = |V|$ und $m = |E|$. Außerdem nehmen wir an, dass jeder Knoten von der Quelle s zu erreichen ist. Das bedeutet insbesondere, dass $m \geq n - 1$ gilt.

Definition 35 Ein Fluss in einem Flussnetzwerk ist eine Funktion $f : V \times V \rightarrow \mathbb{R}$, die die folgenden beiden Eigenschaften erfüllt.

1. Flusserhaltung: Für alle Knoten $u \in V \setminus \{s, t\}$ gilt

$$\sum_{v \in V} f(v, u) = \sum_{v \in V} f(u, v).$$

2. Kapazitätsbeschränkung: Für alle Knoten $u, v \in V$ gilt

$$0 \leq f(u, v) \leq c(u, v).$$

Wir definieren den Wert eines Flusses $f : V \times V \rightarrow \mathbb{R}$ als

$$|f| := \sum_{v \in V} f(s, v) - \sum_{v \in V} f(v, s)$$

und weisen explizit darauf hin, dass die Notation $|\cdot|$ hier keinen Betrag bezeichnet, sondern den Wert des Flusses gemäß obiger Definition.

Das Flussproblem, das wir uns in diesem Abschnitt anschauen werden, lässt sich nun wie folgt formulieren:

Gegeben sei ein Flussnetzwerk G . Berechne einen maximalen Fluss auf G , d. h. einen Fluss f mit größtmöglichem Wert $|f|$.

Anschaulich ist die Frage also, wie viele Einheiten Fluss man in dem gegebenen Netzwerk maximal von der Quelle zur Senke transportieren kann.

Als erste einfache Beobachtung halten wir fest, dass der Fluss, der in der Quelle entsteht, komplett in der Senke ankommt und nicht zwischendurch versickert.

Lemma 36 Sei f ein Fluss in einem Flussnetzwerk G . Dann gilt

$$|f| = \sum_{v \in V} f(v, t) - \sum_{v \in V} f(t, v).$$

Beweis. Als erstes beobachten wir, dass die folgenden beiden Summen gleich sind:

$$\sum_{u \in V} \sum_{v \in V} f(v, u) = \sum_{u \in V} \sum_{v \in V} f(u, v).$$

Dies folgt einfach daraus, dass wir sowohl links als auch rechts über jedes Knotenpaar genau einmal summieren. Nun nutzen wir die Flusserhaltung für alle Knoten $u \in V \setminus \{s, t\}$ und erhalten aus obiger Gleichung:

$$\begin{aligned} \sum_{u \in \{s, t\}} \sum_{v \in V} f(v, u) &= \sum_{u \in \{s, t\}} \sum_{v \in V} f(u, v) \\ \iff \sum_{v \in V} f(v, s) + \sum_{v \in V} f(v, t) &= \sum_{v \in V} f(s, v) + \sum_{v \in V} f(t, v) \\ \iff \sum_{v \in V} f(s, v) - \sum_{v \in V} f(v, s) &= \sum_{v \in V} f(v, t) - \sum_{v \in V} f(t, v), \end{aligned}$$

woraus mit der Definition von $|f|$ direkt das Lemma folgt. \square

Anwendungsbeispiele

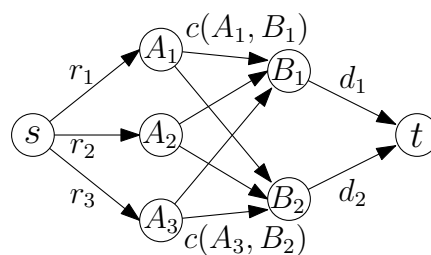
Bevor wir Algorithmen zur Lösung des Flussproblems vorstellen, präsentieren wir ein Anwendungsbeispiel, das wir aus Kapitel 4.5 von [2] übernommen haben. Nehmen wir an, in den Seehäfen A_1, \dots, A_p liegen Bananen zur Verschiffung bereit. Es gibt Zielhäfen B_1, \dots, B_q , zu denen die Bananen transportiert werden sollen. Für $i = 1, \dots, p$ bezeichne r_i , wie viele Tonnen Bananen im Hafen A_i bereit stehen, und für $j = 1, \dots, q$ bezeichne d_j wie viele Tonnen Bananen am Zielfhafen B_j angefordert wurden. Für $i = 1, \dots, p$ und $j = 1, \dots, q$ gibt es zwischen den Häfen A_i und B_j eine Schifffahrtslinie, die maximal $c(A_i, B_j)$ Tonnen Bananen transportieren kann. Es kann auch sein, dass es zwischen zwei Häfen A_i und B_j keine Schifffahrtslinie gibt; in diesem Falle definieren wir $c(A_i, B_j)$ als 0. Die folgenden Fragen stellen sich dem Handelsunternehmen:

1. Ist es möglich, alle Anforderungen zu befriedigen?
2. Falls nein, wie viele Tonnen Bananen können maximal zu den Zielhäfen transportiert werden?
3. Wie sollen die Bananen verschifft werden?

Alle diese Fragen können als Flussproblem modelliert werden. Dazu konstruieren wir einen Graphen $G = (V, E)$ mit den Knoten $V = \{s, t, A_1, \dots, A_p, B_1, \dots, B_q\}$. Es gibt drei Typen von Kanten:

- Für $i = 1, \dots, p$ gibt es eine Kante (s, A_i) mit Kapazität r_i .
- Für $j = 1, \dots, q$ gibt es eine Kante (B_j, t) mit Kapazität d_j .
- Für jedes $i = 1, \dots, p$ und $j = 1, \dots, q$ gibt es eine Kante (A_i, B_j) mit Kapazität $c(A_i, B_j)$.

Die folgende Abbildung zeigt ein Beispiel für $p = 3$ und $q = 2$.



Sei $f : E \rightarrow \mathbb{R}$ ein maximaler Fluss im Graphen G . Man kann sich anschaulich schnell davon überzeugen, dass man mit Hilfe dieses Flusses alle drei Fragen, die oben gestellt wurden, beantworten kann:

1. Wir können alle Anforderungen befriedigen, wenn der Wert $|f|$ des Flusses gleich $\sum_{j=1}^q d_j$ ist.
2. Wir können maximal $|f|$ Tonnen Bananen zu den Zielhäfen transportieren.
3. Der Fluss $f(e)$ auf Kanten $e = (A_i, B_j)$ gibt an, wie viele Tonnen Bananen entlang der Schifffahrtslinie von A_i nach B_j transportiert werden sollen.

Wir verzichten an dieser Stelle auf einen formalen Beweis, dass dies die korrekten Antworten auf obige Fragen sind, möchten aber den Leser dazu anregen, selbst darüber nachzudenken.

Flussprobleme können auch dazu genutzt werden, andere Probleme zu modellieren, die auf den ersten Blick ganz anders geartet sind. Ein Beispiel (siehe Übungen) ist es, zu entscheiden, ob eine Fußballmannschaft bei einer gegebenen Tabellensituation und einem gegebenen restlichen Spielplan gemäß der alten Zwei-Punkte-Regel noch Meister werden kann.

8.3.1 Algorithmus von Ford und Fulkerson

Der folgende Algorithmus von Ford und Fulkerson berechnet für ein gegebenes Flussnetzwerk $G = (V, E)$ mit Kapazitäten $c : V \times V \rightarrow \mathbb{N}_0$ einen maximalen Fluss. Wir gehen in diesem Kapitel davon aus, dass das Flussnetzwerk dergestalt ist, dass für kein Paar von Knoten $u, v \in V$ die beiden Kanten (u, v) und (v, u) in E enthalten sind. Dies ist keine wesentliche Einschränkung, da jedes Netzwerk, das diese Eigenschaft verletzt, einfach in ein äquivalentes umgebaut werden kann, das keine entgegengesetzten Kanten enthält. Es ist eine gute Übung für den Leser, sich zu überlegen, wie das funktioniert.

```

FORD-FULKERSON( $G, c, s \in V, t \in V$ )
1  Setze  $f(e) = 0$  für alle  $e \in E$ . //  $f$  ist gültiger Fluss mit Wert 0
2  while  $\exists$  flussvergrößernder Weg  $P$  {
3      Erhöhe den Fluss  $f$  entlang  $P$ .
4  }
5  return  $f$ ;

```

Wir müssen noch beschreiben, was ein *flussvergrößernder Weg* ist und was es bedeutet, den Fluss entlang eines solchen Weges zu erhöhen. Dazu führen wir den Begriff des *Restnetzwerkes* ein.

Definition 37 Sei $G = (V, E)$ ein Flussnetzwerk mit Kapazitäten $c : V \times V \rightarrow \mathbb{N}_0$ und sei f ein Fluss in G . Das dazugehörige Restnetzwerk $G_f = (V, E_f)$ ist auf der gleichen Menge von Knoten V definiert wie das Netzwerk G . Wir definieren eine Funktion $rest_f : V \times V \rightarrow \mathbb{R}$ mit

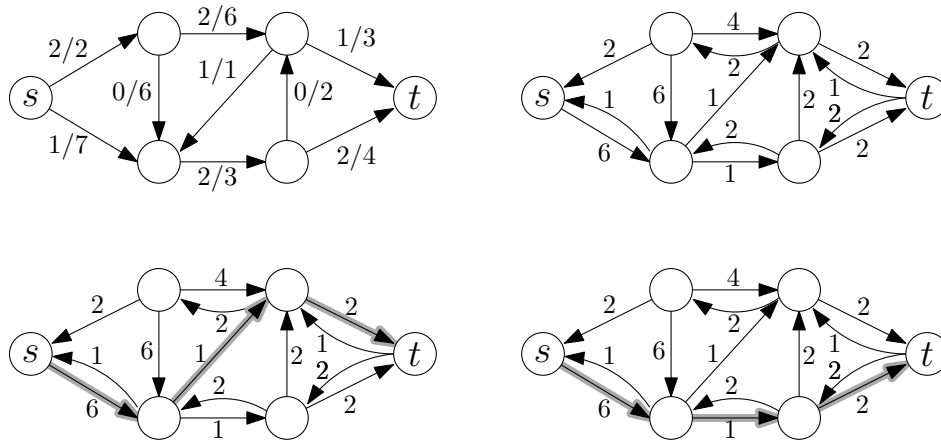
$$rest_f(u, v) = \begin{cases} c(u, v) - f(u, v) & \text{falls } (u, v) \in E, \\ f(v, u) & \text{falls } (v, u) \in E, \\ 0 & \text{sonst.} \end{cases}$$

Die Kantenmenge E_f ist definiert als

$$E_f = \{(u, v) \in V \times V \mid rest_f(u, v) > 0\}.$$

Ein flussvergrößernder Weg ist ein einfacher Weg von s nach t im Restnetzwerk G_f .

Die folgende Abbildung zeigt oben links ein Beispiel für ein Flussnetzwerk mit einem Fluss f . Dabei bedeutet die Beschriftung a/b an einer Kante e , dass $f(e) = a$ und $c(e) = b$ gilt. Rechts oben ist das zugehörige Restnetzwerk dargestellt, und in der unteren Reihe sind zwei flussvergrößernde Wege zu sehen.



Nun müssen wir noch klären, was es bedeutet, den Fluss entlang eines Weges zu erhöhen.

Definition 38 Sei G ein Flussnetzwerk mit Kapazitäten $c : V \times V \rightarrow \mathbb{N}$, sei $f : V \times V \rightarrow \mathbb{R}$ ein Fluss und sei P ein Weg im Restnetzwerk G_f von s nach t . Wir bezeichnen mit $f \uparrow P : V \times V \rightarrow \mathbb{R}$ den Fluss, der entsteht, wenn wir f entlang P erhöhen. Dieser Fluss ist definiert durch

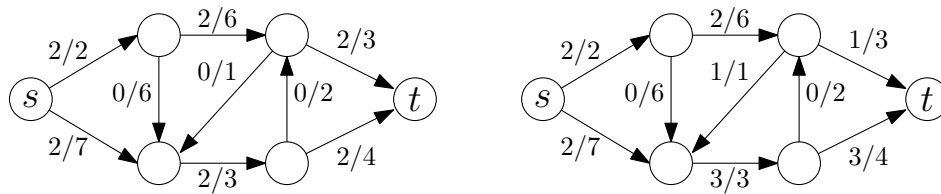
$$(f \uparrow P)(u, v) = \begin{cases} f(u, v) + \delta & \text{falls } (u, v) \in E \text{ und } (u, v) \in P, \\ f(u, v) - \delta & \text{falls } (u, v) \in E \text{ und } (v, u) \in P, \\ f(u, v) & \text{sonst,} \end{cases}$$

wobei

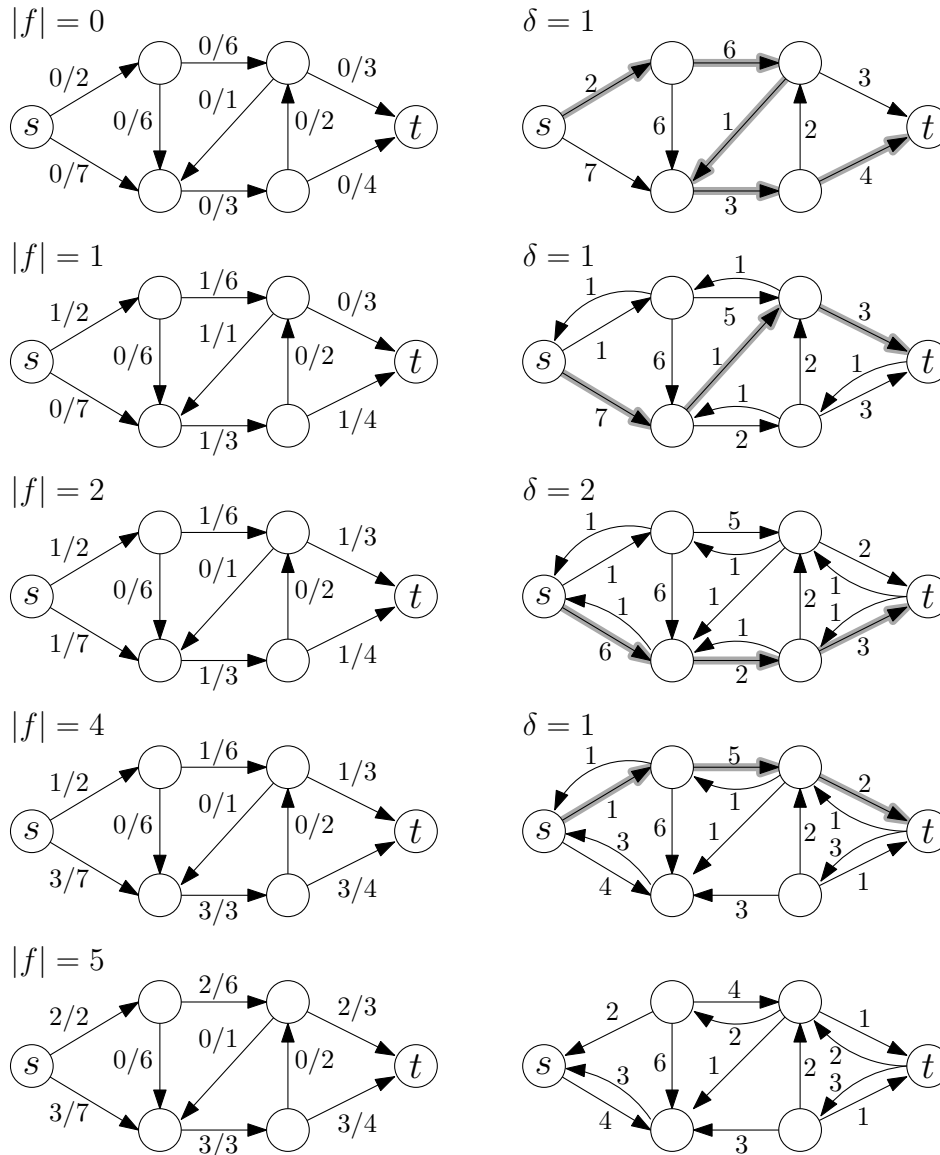
$$\delta = \min_{e \in P} (\text{rest}_f(e)).$$

Aus der Definition von E_f folgt, dass $\delta > 0$ gilt.

In der folgenden Abbildung sind die Flüsse dargestellt, die entstehen, wenn wir den Fluss in obiger Abbildung entlang der beiden oben dargestellten Wege verbessern. In beiden Beispielen ist $\delta = 1$.



Die folgende Abbildung zeigt ein Beispiel für die Ausführung des Algorithmus von Ford und Fulkerson.



Korrektheit

Um die Korrektheit des Algorithmus von Ford und Fulkerson zu zeigen, müssen wir zunächst zeigen, dass $f \uparrow P$ wieder ein Fluss ist.

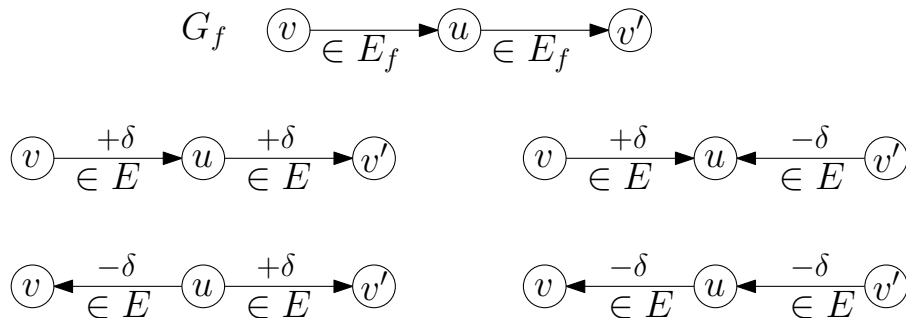
Lemma 39 Sei f ein Fluss in einem Netzwerk G und sei P ein Weg von s nach t im Restnetzwerk G_f . Die Funktion $f \uparrow P : V \times V \rightarrow \mathbb{R}$ ist wieder ein Fluss. Für diesen Fluss gilt

$$|f \uparrow P| = |f| + \delta.$$

Beweis. Wir müssen zeigen, dass die beiden Eigenschaften aus Definition 35 erfüllt sind.

- *Flusserhaltung:* Sei $u \in V \setminus \{s, t\}$. Falls u auf dem Weg P nicht vorkommt, so gilt $f(u, v) = (f \uparrow P)(u, v)$ und $f(v, u) = (f \uparrow P)(v, u)$ für alle Knoten $v \in V$ und somit folgt die Flusserhaltung für u in $f \uparrow P$ aus der Flusserhaltung für u im Fluss f .

Kommt u auf dem Weg P vor, so gibt es genau eine Kante $e \in P$ der Form (v, u) für ein $v \in V$ und es gibt genau eine Kante e' der Form (u, v') für ein $v' \in V$. Die Eindeutigkeit folgt daraus, dass P per Definition einfach ist und somit den Knoten u nur einmal besucht. Die Kanten e und e' sind die einzigen zu u inzidenten Kanten, auf denen sich der Fluss ändert. Wir brauchen also nur diese beiden Kanten näher zu betrachten. Es gilt entweder $(v, u) \in E$ oder $(u, v) \in E$. Ebenso gilt entweder $(v', u) \in E$ oder $(u, v') \in E$. Wir schauen uns alle vier Fälle einzeln an und überprüfen, dass die Flusserhaltung erhalten bleibt. In der folgenden Abbildung ist oben der relevante Ausschnitt des Restnetzwerkes dargestellt und darunter die vier Fälle. In jedem der Fälle ist leicht zu sehen, dass die Flusserhaltung an u erhalten bleibt.



- *Kapazitätsbeschränkung:* Seien $u, v \in V$ beliebig so gewählt, dass $e = (u, v) \in E_f$ eine Kante auf dem Weg P in G_f ist. Wir zeigen, dass $0 \leq (f \uparrow P)(e) \leq c(e)$ gilt. Dafür unterscheiden wir zwei Fälle.

– Ist $e = (u, v) \in E$, so erhöhen wir den Fluss auf e um δ . Auf Grund der Definition von δ und wegen $f(e) \geq 0$ und $\delta \geq 0$ gilt

$$0 \leq f(e) + \delta \leq f(e) + \text{rest}_f(e) = f(e) + (c(e) - f(e)) = c(e).$$

– Ist $e' = (v, u) \in E$, so verringern wir den Fluss auf e' um δ . Auf Grund der Definition von δ und wegen $f(e') \leq c(e')$ und $\delta \geq 0$ gilt

$$c(e') \geq f(e') - \delta \geq f(e') - \text{rest}_f(e) \geq f(e') - f(e') = 0.$$

□

Aus Lemma 39 folgt direkt, dass die Funktion f , die der Algorithmus von Ford und Fulkerson verwaltet, nach jeder Iteration ein Fluss ist. Wir zeigen nun, dass der Algorithmus, wenn er terminiert, einen maximalen Fluss berechnet hat. Dazu führen wir zunächst die folgende Definition ein.

Definition 40 Sei $G = (V, E)$ ein Flussnetzwerk mit Kapazitäten $c : V \times V \rightarrow \mathbf{N}$. Sei $s \in V$ die Quelle und $t \in V$ die Senke. Ein Schnitt von G ist eine Partitionierung der Knotenmenge V in zwei Teile $S \subseteq V$ und $T \subseteq V$ mit $s \in S$, $t \in T$ und $T = V \setminus S$. Wir nennen

$$c(S, T) = \sum_{u \in S} \sum_{v \in T} c(u, v)$$

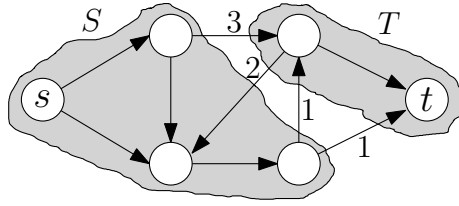
die Kapazität des Schnittes (S, T) . Ein Schnitt (S, T) heißt minimal, wenn es keinen Schnitt (S', T') mit $c(S', T') < c(S, T)$ gibt.

Für einen Fluss f und einen Schnitt (S, T) sei

$$f(S, T) = \sum_{u \in S} \sum_{v \in T} f(u, v) - \sum_{u \in S} \sum_{v \in T} f(v, u)$$

der Fluss über den Schnitt.

Die folgende Abbildung zeigt einen Schnitt mit Kapazität 5.



Anschaulich ist leicht einzusehen, dass die Kapazität $c(S, T)$ eine obere Schranke dafür ist, wie viel Fluss von der Quelle zur Senke geschickt werden kann, und dass für einen Fluss f der Fluss über den Schnitt (S, T) genau $|f|$ sein muss. Formal zeigen wir das in folgendem Lemma.

Lemma 41 Sei (S, T) ein Schnitt eines Flussnetzwerkes G und sei f ein Fluss in G . Dann gilt

$$|f| = f(S, T) \leq c(S, T).$$

Beweis. Ähnlich zu dem Beweis von Lemma 36 summieren wir zunächst über alle Kanten, die innerhalb der Menge S verlaufen. Wir erhalten die folgende Gleichung:

$$\begin{aligned} \sum_{v \in S} f(s, v) + \sum_{u \in S \setminus \{s\}} \sum_{v \in S} f(u, v) &= \sum_{v \in S} f(v, s) + \sum_{u \in S \setminus \{s\}} \sum_{v \in S} f(v, u) \\ &\iff \\ \sum_{v \in S} f(s, v) - \sum_{v \in S} f(v, s) &= \sum_{u \in S \setminus \{s\}} \sum_{v \in S} f(v, u) - \sum_{u \in S \setminus \{s\}} \sum_{v \in S} f(u, v). \end{aligned} \quad (8.1)$$

Für alle Knoten $u \in S \setminus \{s\}$ gilt die Flusserhaltung. Dafür müssen wir aber auch die Kanten betrachten, die von u in die Menge T führen, und die Kanten, die von der Menge T zu u führen. Wir erhalten

$$\sum_{v \in S} f(u, v) + \sum_{v \in T} f(u, v) = \sum_{v \in S} f(v, u) + \sum_{v \in T} f(v, u).$$

Damit können wir Gleichung 8.1 schreiben als

$$\sum_{v \in S} f(s, v) - \sum_{v \in S} f(v, s) = \sum_{u \in S \setminus \{s\}} \sum_{v \in T} f(u, v) - \sum_{u \in S \setminus \{s\}} \sum_{v \in T} f(v, u). \quad (8.2)$$

Für die Quelle s gilt entsprechend

$$\begin{aligned} \sum_{v \in S} f(s, v) + \sum_{v \in T} f(s, v) &= |f| + \sum_{v \in S} f(v, s) + \sum_{v \in T} f(v, s) \\ &\iff \\ \sum_{v \in S} f(s, v) - \sum_{v \in S} f(v, s) &= |f| + \sum_{v \in S} f(v, s) - \sum_{v \in T} f(s, v). \end{aligned} \quad (8.3)$$

Wir setzen Gleichung 8.3 in Gleichung 8.2 ein und erhalten

$$\begin{aligned} |f| + \sum_{v \in T} f(v, s) - \sum_{v \in T} f(s, v) &= \sum_{u \in S \setminus \{s\}} \sum_{v \in T} f(u, v) - \sum_{u \in S \setminus \{s\}} \sum_{v \in T} f(v, u) \\ &\iff \\ |f| &= \sum_{u \in S} \sum_{v \in T} f(u, v) - \sum_{u \in S} \sum_{v \in T} f(v, u) = f(S, T). \end{aligned}$$

Die Ungleichung folgt einfach aus obiger Gleichung, da $f(u, v) \leq c(u, v)$ und $f(v, u) \geq 0$ für alle $u, v \in V$:

$$f(S, T) \leq \sum_{u \in S} \sum_{v \in T} f(u, v) \leq \sum_{u \in S} \sum_{v \in T} c(u, v) \leq c(S, T).$$

□

Mit Hilfe von Lemma 41 zeigen wir das folgende Theorem, aus dem direkt folgt, dass der Algorithmus von Ford und Fulkerson einen maximalen Fluss gefunden hat, wenn er terminiert.

Theorem 42 *Sei f ein Fluss in einem Netzwerk G . Dann sind die folgenden drei Aussagen äquivalent.*

- a) f ist ein maximaler Fluss.
- b) Das Restnetzwerk G_f enthält keinen flussvergrößernden Weg.
- c) Es gibt einen Schnitt (S, T) mit $|f| = c(S, T)$.

Dieses Theorem besagt insbesondere, dass es einen Schnitt (S, T) geben muss, dessen Kapazität so groß ist wie der Wert $|f|$ des maximalen Flusses f . Zusammen mit Lemma 41, das besagt, dass es keinen Schnitt (S', T') mit kleinerer Kapazität als $|f|$ geben kann, folgt die wichtige Aussage, dass die Kapazität des minimalen Schnittes gleich dem Wert des maximalen Flusses ist. Die Korrektheit des Algorithmus von Ford und Fulkerson bei Terminierung folgt einfach daraus, dass der Algorithmus erst dann terminiert, wenn es keinen flussvergrößernden Weg mehr gibt. Nach dem Theorem ist das aber genau dann der Fall, wenn der Fluss maximal ist.

Beweis.[Beweis von Theorem 42] Wir werden zeigen, dass b) aus a) folgt, dass c) aus b) folgt und dass a) aus c) folgt. Diese drei Implikationen bilden einen Ringschluss und damit ist gezeigt, dass die drei Aussagen äquivalent sind.

- **a) \Rightarrow b)** Wir führen einen Widerspruchsbeweis und nehmen an, dass es einen flussvergrößernden Weg P in G_f gibt. Dann können wir den Fluss f entlang P erhöhen und erhalten den Fluss $f \uparrow P$, welcher nach Lemma 39 einen um $\delta > 0$ höheren Wert als f hat. Damit kann f kein maximaler Fluss sein.
- **c) \Rightarrow a)** Sei f ein Fluss mit $|f| = c(S, T)$ für einen Schnitt (S, T) und sei f' ein maximaler Fluss. Dann gilt $|f| \leq |f'|$. Wenn wir Lemma 41 für den Fluss f' und den Schnitt (S, T) anwenden, dann erhalten wir $|f'| \leq c(S, T)$. Zusammen ergibt das

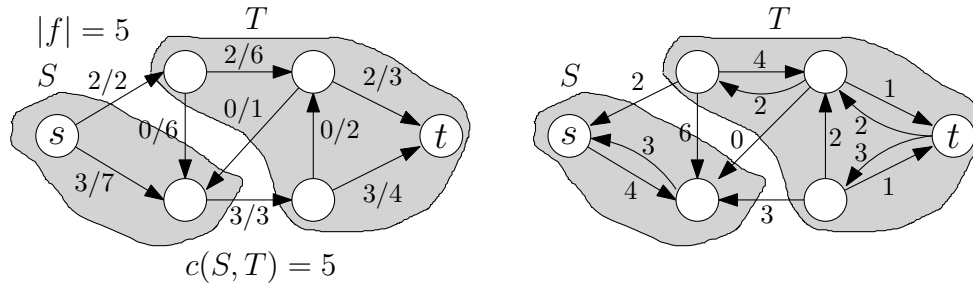
$$c(S, T) = |f| \leq |f'| \leq c(S, T),$$

woraus $c(S, T) = |f| = |f'|$ folgt. Damit ist auch f ein maximaler Fluss.

- **b) \Rightarrow c)** Laut Voraussetzung gibt es keinen flussvergrößernden Weg, das heißt, es gibt keinen Weg von s nach t im Restnetzwerk G_f . Wir teilen die Knoten V des Restnetzwerkes in zwei Klassen S und T ein:

$$S = \{v \in V \mid \text{es gibt Weg in } G_f \text{ von } s \text{ nach } v\} \quad \text{und} \quad T = V \setminus S.$$

Da der Knoten t in G_f nicht von s aus erreichbar ist, gilt $s \in S$ und $t \in T$. Damit ist (S, T) ein Schnitt des Graphen. Die Situation ist in folgender Abbildung dargestellt.



Wir argumentieren nun, dass $|f| = c(S, T)$ für diesen Schnitt (S, T) gilt.

- Sei $(u, v) \in E$ eine Kante mit $u \in S$ und $v \in T$. Da u in G_f von s aus erreichbar ist, v aber nicht, gilt $(u, v) \notin E_f$. Mit der Definition des Restnetzwerkes folgt $\text{rest}_f(u, v) = 0$, also $f(u, v) = c(u, v)$.
- Sei $(v, u) \in E$ eine Kante mit $u \in S$ und $v \in T$. Da u in G_f von s aus erreichbar ist, v aber nicht, gilt $(u, v) \notin E_f$. Mit der Definition des Restnetzwerkes folgt $\text{rest}_f(u, v) = 0$, also $f(v, u) = 0$.

Somit gilt mit Lemma 41

$$\begin{aligned} c(S, T) &= \sum_{u \in S} \sum_{v \in T} c(u, v) = \sum_{u \in S} \sum_{v \in T} f(u, v) - \sum_{u \in S} \sum_{v \in T} f(v, u) \\ &= f(S, T) = |f| \leq c(S, T). \end{aligned}$$

Damit folgt $|f| = f(S, T) = c(S, T)$.

□

Laufzeit

Wir haben den Korrektheitsbeweis noch nicht abgeschlossen. Wir haben lediglich gezeigt, dass der Algorithmus *partiell korrekt* ist. Das bedeutet, wir haben gezeigt, dass der Algorithmus das richtige Ergebnis liefert, falls er terminiert. Wir haben aber noch nicht gezeigt, dass er wirklich auf jeder Eingabe terminiert. Das erledigen wir nun implizit bei der Laufzeitabschätzung.

Theorem 43 Für ganzzahlige Kapazitäten $c : V \times V \rightarrow \mathbb{N}_0$ ist die Anzahl an Iterationen der **while**-Schleife im Algorithmus von Ford und Fulkerson durch $C = \sum_{e \in E} c(e)$ nach oben beschränkt. Die Laufzeit des Algorithmus ist $O(mC)$.

Beweis. Man kann sehr einfach (z. B. mit Hilfe einer Invariante) zeigen, dass der aktuelle Fluss f , den der Algorithmus verwaltet stets ganzzahlig ist, dass also stets $f : V \times V \rightarrow \mathbb{N}_0$ gilt. Dies folgt einfach daraus, dass für einen ganzzahligen Fluss f und einen flussvergrößernden Weg P stets $\delta \in \mathbb{N}$ gilt. Demzufolge steigt der Wert des Flusses mit jeder Iteration der **while**-Schleife um mindestens eins an. Da der Wert eines jeden Flusses durch $\sum_{e \in E} c(e)$ beschränkt ist, garantiert dies die Terminierung nach höchstens C vielen Schritten.

Für einen gegebenen Fluss f kann das Restnetzwerk G_f in Zeit $O(m)$ berechnet werden. Ein Weg P von s nach t in G_f kann mittels Tiefensuche in Zeit $O(m)$ gefunden werden. Ist ein solcher Weg P gefunden, so kann in Zeit $O(m)$ der neue Fluss $f \uparrow P$ berechnet werden. Insgesamt dauert also jede Iteration Zeit $O(m)$, woraus das Theorem folgt. □

Aus diesem Theorem und seinem Beweis folgt direkt, dass es bei ganzzahligen Kapazitäten auch stets einen ganzzahligen maximalen Fluss gibt. Dies ist eine sehr wichtige Eigenschaft für viele Anwendungen und deswegen halten wir sie explizit als Korollar fest.

Corollary 44 *Sind alle Kapazitäten ganzzahlig, so gibt es stets einen ganzzahligen maximalen Fluss.*

8.3.2 Algorithmus von Edmonds und Karp

Die Laufzeit $O(mC)$, die wir in Theorem 43 bewiesen haben, ist nicht besonders gut, weil sie mit den Kapazitäten der Kanten wächst. Dies steht im Gegensatz zu den Laufzeiten von zum Beispiel den Algorithmen von Kruskal und Dijkstra. Für diese Algorithmen haben wir Laufzeitschranken bewiesen, die nur von der Größe des Graphen abhängen, nicht jedoch von den Gewichten oder Längen der Kanten. Eine Laufzeit wie $O(mC)$, die von den Kapazitäten abhängt, heißt *pseudopolynomielle Laufzeit*. Wir werden diesen Begriff im nächsten Semester weiter vertiefen. Hier sei nur angemerkt, dass man, wenn möglich, Laufzeiten anstrebt, die nur von der Größe des Graphen, nicht jedoch von den Kapazitäten abhängen, da diese sehr groß werden können. Wir werden nun den *Algorithmus von Edmonds und Karp* kennenlernen. Dabei handelt es sich um eine leichte Modifikation des Algorithmus von Ford und Fulkerson, deren Laufzeit nur von der Größe des Graphen abhängt.

Wir haben in der Beschreibung des Algorithmus von Ford und Fulkerson offen gelassen, wie die flussvergrößernden Wege gewählt werden, wenn es mehrere zur Auswahl gibt. Unsere Analyse von Korrektheit und Laufzeit war unabhängig von dieser Wahl. Wir zeigen nun, dass für eine geeignete Wahl die Laufzeit reduziert werden kann.

```

EDMONDS-KARP( $G, c, s \in V, t \in V$ )
1  Setze  $f(e) = 0$  für alle  $e \in E$ . //  $f$  ist gültiger Fluss mit Wert 0
2  while  $\exists$  flussvergrößernder Weg  $P$  {
3      Wähle einen flussvergrößernden Weg  $P$  mit so wenig Kanten wie möglich.
4      Erhöhe den Fluss  $f$  entlang  $P$ .
5  }
6  return  $f$ ;

```

Der einzige Unterschied zum Algorithmus von Ford und Fulkerson ist also, dass immer ein kürzester flussvergrößernder Weg gewählt wird, wobei die Länge eines Weges als die Anzahl seiner Kanten definiert ist. Wir werden in diesem Zusammenhang auch von der *Distanz von s zu $v \in V$ im Restnetzwerk G_f* sprechen. Damit meinen wir, wie viele Kanten der kürzeste Weg (d. h. der Weg mit den wenigsten Kanten) von s nach v in G_f hat.

Theorem 45 *Die Laufzeit des Algorithmus von Edmonds und Karp ist $O(m^2n) = O(n^5)$.*

Beweis. Eine Iteration der **while**-Schleife kann immer noch in Zeit $O(m)$ durchgeführt werden. Einen kürzesten flussvergrößernden Weg findet man nämlich zum Beispiel, indem man von s startend eine Breitensuche im Restnetzwerk durchführt. Diese benötigt Zeit $O(m)$. Zum Beweis des Theorems genügt es also zu zeigen, dass der Algorithmus nach $O(mn)$ Iterationen terminiert. Dazu zeigen wir zunächst die folgende strukturelle Eigenschaft über die Distanzen im Restnetzwerk.

Lemma 46 *Sei $x \in V$ beliebig. Die Distanz von s zu x in G_f wird im Laufe des Algorithmus nicht kleiner.*

Beweis. Wir betrachten eine Iteration der **while**-Schleife, in der der Fluss f entlang des Weges P zu dem Fluss $f \uparrow P$ erhöht wird. Die Kantenmenge $E_{f \uparrow P}$ unterscheidet sich von der Kantenmenge E_f wie folgt:

- Für jede Kante $(u, v) \in P \subseteq E_f$ verringert sich $\text{rest}_f(u, v)$ um δ , das heißt $\text{rest}_{f \uparrow P}(u, v) = \text{rest}_f(u, v) - \delta$. Eine Kante mit $\text{rest}_f(u, v) = \delta$ heißt *Flaschenhalskante* und sie ist in $E_{f \uparrow P}$ nicht mehr enthalten.
- Für jede Kante $(u, v) \in P \subseteq E_f$ erhöht sich die Restkapazität $\text{rest}_f(v, u)$ der entgegengesetzten Kante um δ , das heißt $\text{rest}_{f \uparrow P}(v, u) = \text{rest}_f(v, u) + \delta$. War $\text{rest}_f(v, u) = 0$, so war $(v, u) \notin E_f$, nun gilt aber $(v, u) \in E_{f \uparrow P}$.

Wir spalten den Übergang von E_f zu $E_{f \uparrow P}$ in mehrere Schritte auf. Zunächst fügen wir alle neuen Kanten nach und nach ein, die gemäß dem zweiten Aufzählungspunkt entstehen. Anschließend entfernen wir die Kanten nach und nach gemäß dem ersten Aufzählungspunkt, um die Kantenmenge $E_{f \uparrow P}$ zu erhalten. Nach jedem Einfügen und Löschen einer Kante untersuchen wir, wie sich die Distanz von s zu x geändert hat.

Als erstes betrachten wir den Fall, dass wir eine Kante (v, u) einfügen. Gemäß dem zweiten Aufzählungspunkt bedeutet das, dass die Kante (u, v) auf dem Weg P liegen muss. In dem Algorithmus wurde der Weg P so gewählt, dass er ein kürzester Weg von s nach t ist. Wegen der optimalen Substruktur, die wir uns für kürzeste Wege überlegt haben, bedeutet das insbesondere, dass die Teilwege von P , die von s nach u und v führen, kürzeste Wege von s nach u bzw. v sind. Die Länge des Teilweges nach u bezeichnen wir mit ℓ . Dann ist die Länge des Teilweges nach v genau $\ell + 1$. Somit verbindet die neue Kante einen Knoten mit Distanz $\ell + 1$ von s mit einem Knoten mit Distanz ℓ von s . Dadurch kann sich jedoch die Distanz von s zu keinem Knoten des Graphen verringern. Um die Distanz zu verringern, müsste die neue Kante einen Knoten mit Distanz k von s für ein $k \geq 0$ mit einem Knoten mit Distanz $k + i$ von s für ein $i \geq 2$ verbinden. Damit folgt, dass selbst nach dem Einfügen aller neuen Kanten der Knoten x immer noch die gleiche Distanz von s hat.

Es ist leicht einzusehen, dass das Löschen von Kanten keine Distanzen verringern kann. Deshalb hat der Knoten x nach dem Löschen mindestens die gleiche Distanz von s wie in G_f . \square

Mit Hilfe von Lemma 46 können wir nun das Theorem beweisen. Bei jeder Iteration der **while**-Schleife gibt es mindestens eine Flaschenhalskante (u, v) . Diese wird aus dem Restnetzwerk entfernt. Sei ℓ die Distanz von s zu u , bevor die Kante entfernt wird. Die Distanz zu v ist dann $\ell + 1$. Es ist möglich, dass diese Kante zu einem späteren Zeitpunkt wieder in das Restnetzwerk eingefügt wird, aber nur dann, wenn die Kante (v, u) in einer späteren Iteration auf dem gewählten flussvergrößernden P' liegt. Da sich die Distanz von s zu v nicht verringern kann, muss zu diesem Zeitpunkt die Distanz von s nach v immer noch mindestens $\ell + 1$ sein. Da P' ein kürzester Weg ist, können wir wegen der optimalen Substruktur wieder folgern, dass die Distanz von s zu u zu diesem Zeitpunkt mindestens $\ell + 2$ sein muss.

Zusammenfassend können wir also festhalten, dass das Entfernen und Wiedereinfügen einer Kante (u, v) im Restnetzwerk die Distanz von s zu u um mindestens 2 erhöht. Da die maximale Distanz $n - 1$ ist, kann eine Kante demnach nicht öfter als $n/2$ mal entfernt werden. In jeder Iteration wird mindestens eine Kante entfernt und es gibt $2m$ potentielle Kanten im Restnetzwerk. Damit ist die Anzahl an Iterationen der **while**-Schleife nach oben beschränkt durch $\frac{n}{2} \cdot 2m = nm$. \square

Literatur

Flussprobleme werden in Kapitel 26 von [1] beschrieben. Außerdem sind sie in Kapitel 4.5 von [2] erklärt.

8.3.3 Der Heiratssatz

Flussmaximierung hat zahlreiche überraschende Anwendungen in anderen Gebieten. Als Beispiel betrachten wir den Satz von Hall, der auch als *Heiratssatz* bekannt ist. Gegeben sind eine Menge V_1 von Damen und eine Menge V_2 von Herren. Jede Dame $a \in V_1$ hat eine nicht leere Menge $F(a) \subseteq V_2$ von Freunden in V_2 .

Die Mengen $F(a)$ brauchen nicht disjunkt zu sein. Deshalb stellt sich die Frage, ob jede Dame a einen ihrer Freunde in $F(a)$ heiraten kann, ohne daß es zu Polygamie kommt; gesucht wird also eine *injektive* Abbildung

$$h : V_1 \longrightarrow V_2$$

mit $h(a) \in F(a)$ für alle $a \in V_1$.

Eine (scheinbar) schwächere Forderung ist durch die sogenannte *Party-Bedingung* gegeben: Wenn eine Teilmenge von Damen alle ihre Freunde einlädt, gibt es genügend viele Tanzpartner (die aber keine Freunde zu sein brauchen). Das heißt, für jede Teilmenge $W \subseteq V_1$ gilt

$$|W| \leq \left| \bigcup_{a \in W} F(a) \right|$$

Überraschenderweise gilt nun folgendes

Theorem 47 *Eine Heirat h existiert genau dann, wenn die Party-Bedingung erfüllt ist.*

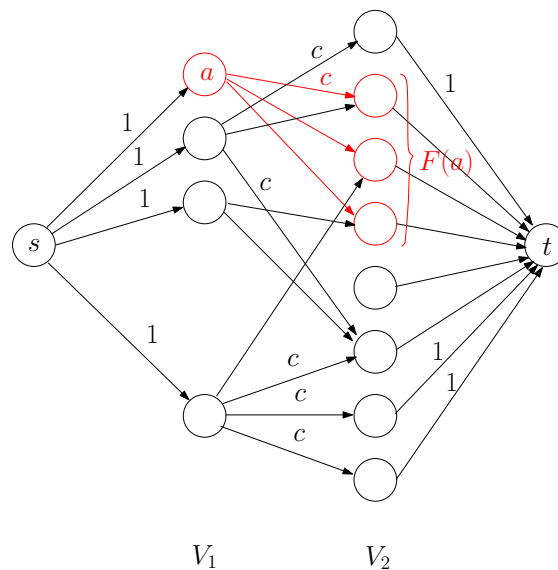


Abbildung 8.2: Die Kapazität c hat den Wert $|V_1| + 1$.

Beweis. Wenn eine Heirat h existiert, kann jede Dame a auf jeder Party mit ihrem Ehemann $h(a) \in F(a)$ tanzen; deshalb ist die Party-Bedingung erfüllt. Umgekehrt nehmen wir nun an, daß die Party-Bedingung erfüllt ist, und betrachten das in Abbildung 8.2 gezeigte Flussnetzwerk. Links stehen die Damen, rechts die Herren, und von jeder Dame gehen Kanten mit Kapazität $c := |V_1| + 1$ zu ihren Freunden aus. Quelle s ist mit allen Damen verbunden, und alle Herren sind mit der Senke verbunden; diese Kanten haben alle die Kapazität 1.

Behauptung: Die von s ausgehenden Kanten bilden einen minimalen Schnitt mit Kapazität $|V_1|$.

Zum Beweis nehmen wir an, es gäbe einen Schnitt mit Kapazität $< |V_1|$. Kanten von V_1 nach V_2 können darin nicht vorkommen, denn deren Kapazität c ist zu groß. Also kann der Schnitt nur Kanten von s zu einer Teilmenge V_1^0 von V_1 enthalten sowie Kanten von einer Teilmenge $V_2^0 \subseteq V_2$ nach t , und es gilt $|V_1^0| + |V_2^0| < |V_1|$.

Sei nun $a \in V_1 \setminus V_1^0$ und $w \in F(a)$. Dann muß w zu V_2^0 gehören, denn sonst gäbe es ja den Pfad $s \rightarrow a \rightarrow w \rightarrow t$ im Widerspruch zur Eigenschaft eines Schnitts. Das bedeutet

$$\bigcup_{a \in V_1 \setminus V_1^0} F(a) \subseteq V_2^0,$$

also

$$\left| \bigcup_{a \in V_1 \setminus V_1^0} F(a) \right| \leq |V_2^0| < |V_1| - |V_1^0| = |V_1 \setminus V_1^0|$$

im Widerspruch zur Party-Bedingung. Also ist die Behauptung richtig, und aus Theorem 42 folgt, daß es einen Fluss mit Wert $|V_1|$ geben muß. Dieser Fluss muß über $|V_1|$ *verschiedene* Knoten von V_2 zur Senke laufen und markiert deshalb eine Heirat. \square

Kapitel 9

Lineare Programmierung

Wir haben bereits einige allgemeine Techniken zum Entwurf von Algorithmen wie zum Beispiel dynamische Programmierung und Greedy-Algorithmen kennengelernt. Dennoch mussten wir bei jedem neuen Problem separat prüfen, ob und wie genau sich diese Methoden anwenden lassen. Wir werden nun *lineare Programmierung* kennenlernen. Dabei handelt es sich um eine sehr allgemeine Technik zur Lösung von Optimierungsproblemen, die noch standardisierter ist als die Techniken, die wir bisher gesehen haben.

Ein *lineares Programm (LP)* ist ein Optimierungsproblem, bei dem es darum geht, die optimalen Werte für d reelle *Variablen* $x_1, \dots, x_d \in \mathbb{R}$ zu bestimmen. Dabei gilt es, eine lineare *Zielfunktion*

$$c_1x_1 + \dots + c_dx_d \tag{9.1}$$

für gegebene Koeffizienten $c_1, \dots, c_d \in \mathbb{R}$ zu minimieren oder zu maximieren. Es müssen dabei m lineare *Nebenbedingungen* eingehalten werden. Für jedes $i \in \{1, \dots, m\}$ sind Koeffizienten $a_{i1}, \dots, a_{id} \in \mathbb{R}$ und $b_i \in \mathbb{R}$ gegeben, und eine Belegung der Variablen ist nur dann gültig, wenn sie die Nebenbedingungen

$$\begin{aligned} a_{11}x_1 + \dots + a_{1d}x_d &\leq b_1 \\ &\vdots \\ a_{m1}x_1 + \dots + a_{md}x_d &\leq b_m \end{aligned} \tag{9.2}$$

einhält. Statt dem Relationszeichen \leq erlauben wir prinzipiell auch Nebenbedingungen mit dem Relationszeichen \geq . Da man das Relationszeichen einer Nebenbedingung aber einfach dadurch umdrehen kann, dass man beide Seiten mit -1 multipliziert, gehen wir im Folgenden ohne Beschränkung der Allgemeinheit davon aus, dass alle Nebenbedingungen das Relationszeichen \leq enthalten.

Um die Notation zu vereinfachen, definieren wir $x^T = (x_1, \dots, x_d)$ und $c^T = (c_1, \dots, c_d)$. Damit sind $x \in \mathbb{R}^d$ und $c \in \mathbb{R}^d$ Spaltenvektoren und wir können die Zielfunktion (9.1) als Skalarprodukt $c \cdot x$ schreiben. Außerdem definieren wir $A \in \mathbb{R}^{m \times d}$ als die Matrix mit den Einträgen a_{ij} und wir definieren $b^T = (b_1, \dots, b_m) \in \mathbb{R}^m$. Dann entspricht jede Zeile der Matrix einer Nebenbedingung und wir können die Nebenbedingungen (9.2) als $Ax \leq b$ schreiben. Zusammengefasst können wir ein lineares Programm also wie folgt beschreiben.

Lineares Programm

Eingabe: $c \in \mathbb{R}^d, b \in \mathbb{R}^m, A \in \mathbb{R}^{m \times d}$

Lösungen: alle $x \in \mathbb{R}^d$ mit $Ax \leq b$

Zielfunktion: minimiere/maximiere $c \cdot x$

Wir präsentieren nun einige Beispiele, die den Leser davon überzeugen sollen, dass man viele Probleme in Form eines linearen Programms ausdrücken kann.

Hobbygärtner

Ein Hobbygärtner besitzt 100 m^2 Land, auf dem er Blumen und Gemüse anbauen möchte.

- Nur 60 m^2 des Landes sind für den Anbau von Blumen und Gemüse geeignet. Auf den restlichen 40 m^2 kann nur Gemüse angebaut werden.
- Er hat ein Budget in Höhe von 720 € , das er investieren kann. Das Pflanzen und Düngen der Blumen kostet ihn 9 € pro Quadratmeter; bei Gemüse sind es nur 6 € pro Quadratmeter.
- Der Erlös, den er pro Quadratmeter für Blumen erzielt, beträgt 20 € ; bei Gemüse sind es nur 10 € pro Quadratmeter.

Auf wie vielen Quadratmetern sollte der Hobbygärtner Blumen und auf wie vielen sollte er Gemüse anbauen, wenn er seinen Gewinn maximieren möchte?

Um dieses Problem als LP zu formulieren, führen wir zwei Variablen $x_B \in \mathbb{R}$ und $x_G \in \mathbb{R}$ ein, die angeben, auf wie vielen Quadratmetern Blumen bzw. Gemüse angebaut werden. Dann können wir die Nebenbedingungen wie folgt formulieren:

$$\begin{aligned} x_B \geq 0, x_G \geq 0 & \quad (\text{keine negative Anbaufläche}) \\ x_B + x_G \leq 100 & \quad (\text{maximal } 100 \text{ m}^2) \\ x_B \leq 60 & \quad (\text{maximal } 60 \text{ m}^2 \text{ Blumen}) \\ 9x_B + 6x_G \leq 720 & \quad (\text{Budget von } 720 \text{ €}) \end{aligned}$$

Der Gewinn ist die Differenz von Erlös und Ausgaben. Somit soll die folgende Zielfunktion maximiert werden:

$$(20 - 9)x_B + (10 - 6)x_G = 11x_B + 4x_G.$$

Wir können auch einige Probleme, die wir bereits kennen, als lineares Programm formulieren.

Maximaler Fluss

Gegeben sei ein Flussnetzwerk $G = (V, E)$ mit Quelle $s \in V$ und Senke $t \in V$ und einer *Kapazitätsfunktion* $c : E \rightarrow \mathbb{N}_0$. Das Problem, einen maximalen Fluss von s nach t zu berechnen, können wir als LP darstellen, indem wir für jede Kante $e \in E$ eine Variable $x_e \in \mathbb{R}$ einführen, die dem Fluss auf Kante e entspricht. Wir wollen dann die Zielfunktion

$$\sum_{e=(s,v)} x_e - \sum_{e=(v,s)} x_e$$

unter den folgenden Bedingungen maximieren.

$$\begin{aligned} \forall e \in E : x_e \geq 0 & \quad (\text{Fluss nicht negativ}) \\ \forall e \in E : x_e \leq c(e) & \quad (\text{Fluss nicht größer als Kapazität}) \\ \forall v \in V \setminus \{s, t\} : \sum_{e=(u,v)} x_e - \sum_{e=(v,u)} x_e = 0 & \quad (\text{Flusserhaltung}) \end{aligned}$$

Die letzten Nebenbedingungen sind als Gleichungen und nicht als Ungleichungen formuliert. Dies ist ohne Erweiterung des Modells möglich, da wir jede Gleichung durch zwei Ungleichungen mit den Relationszeichen \leq und \geq ersetzen können.

9.1 Grundlagen

9.1.1 Kanonische Form und Gleichungsform

Um Algorithmen zur Lösung von linearen Programmen zu beschreiben, ist es hilfreich, die folgenden beiden Formen zu betrachten. Dabei sei $c \in \mathbb{R}^d$, $b \in \mathbb{R}^m$ und $A \in \mathbb{R}^{m \times d}$, und $x \in \mathbb{R}^d$ sei ein Vektor aus d Variablen.

<i>kanonische Form</i>	<i>Gleichungsform</i>
$\min c \cdot x$	$\min c \cdot x$
$Ax \leq b$	$Ax = b$
$x \geq 0$	$x \geq 0$

Wir können uns ohne Beschränkung der Allgemeinheit auf eine dieser beiden Formen konzentrieren, denn jedes lineare Programm kann in ein äquivalentes lineares Programm in kanonischer Form und in ein äquivalentes lineares Programm in Gleichungsform transformiert werden. Dazu müssen nur die folgenden Operationen in der richtigen Reihenfolge angewendet werden. Wir bezeichnen bei der Beschreibung dieser Operationen mit $a_i = (a_{i1}, \dots, a_{id})^T \in \mathbb{R}^d$ den Vektor, der der i -ten Zeile der Matrix A entspricht. Dann können wir die i -te Nebenbedingung in kanonischer Form beispielsweise kurz als $a_i \cdot x \leq b_i$ schreiben.

- Haben wir ein lineares Programm, in dem eine Zielfunktion $c \cdot x$ maximiert werden soll, so können wir dies äquivalent dadurch ausdrücken, dass die Zielfunktion $-c \cdot x$, minimiert werden soll.
- Für jede Variable x_i , für die es keine Nebenbedingung $x_i \geq 0$ gibt, können wir zwei Variablen x'_i und x''_i mit Nebenbedingungen $x'_i \geq 0$ und $x''_i \geq 0$ einführen und jedes Vorkommen von x_i in Zielfunktion und Nebenbedingungen durch $x'_i - x''_i$ substituieren.
- Bei Nebenbedingungen der Form $a_i \cdot x \geq b_i$ können wir das Relationszeichen umdrehen, indem wir beide Seiten der Ungleichung mit -1 multiplizieren.
- Eine Gleichung $a_i \cdot x = b_i$ kann durch die beiden Ungleichungen $a_i \cdot x \leq b_i$ und $a_i \cdot x \geq b_i$ ersetzt werden.
- Eine Ungleichung $a_i \cdot x \leq b_i$ können wir in eine Gleichung umwandeln, indem wir eine *Schlupfvariable* s_i mit Nebenbedingung $s_i \geq 0$ einführen und die Ungleichung durch die Gleichung $s_i + a_i \cdot x = b_i$ ersetzen.

9.1.2 Geometrische Interpretation

Bevor wir einen Algorithmus beschreiben, der für ein gegebenes lineares Programm eine optimale Belegung der Variablen berechnet, ist es zunächst hilfreich, eine geometrische Interpretation von linearen Programmen zu erarbeiten. Dazu gehen wir davon aus, dass ein lineares Programm in kanonischer Form gegeben ist.

Nebenbedingungen

Jede Variablenbelegung $x \in \mathbb{R}^d$ kann als ein Punkt in dem d -dimensionalen Raum \mathbb{R}^d aufgefasst werden. Die Nebenbedingungen bestimmen dann, welche Punkte zulässige Lösungen des linearen Programms sind. Eine Gleichung der Form $a_i \cdot x = b_i$ definiert eine *affine Hyperebene* $\{x \in \mathbb{R}^d \mid a_i \cdot x = b_i\}$. Dabei handelt es sich um die Verallgemeinerung einer normalen

Ebene auf beliebige Dimensionen $d \in \mathbb{N}$. Für $d = 3$ entsprechen Hyperebenen normalen Ebenen; für $d = 2$ sind es beispielsweise Geraden. Jede solche affine Hyperebene definiert den *abgeschlossenen Halbraum*

$$\mathcal{H}_i = \{x \in \mathbb{R}^d \mid a_i \cdot x \leq b_i\}.$$

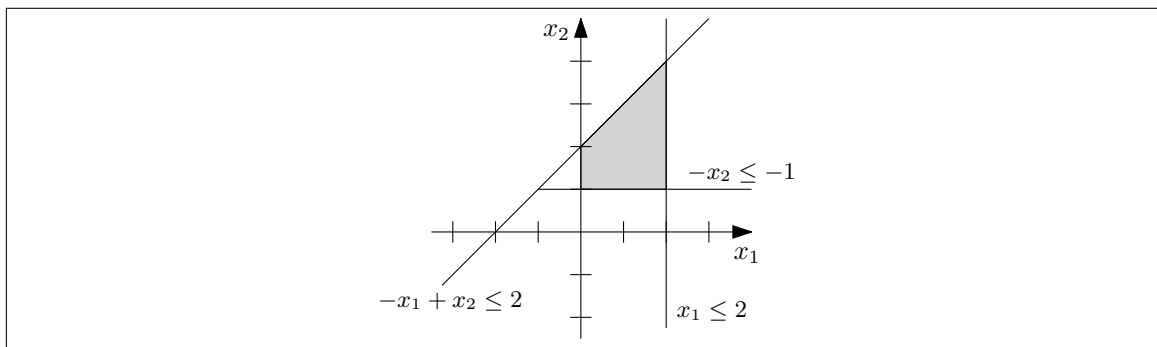
Eine Variablenbelegung $x \in \mathbb{R}^d$ erfüllt genau dann Nebenbedingung i , wenn $x \in \mathcal{H}_i$ gilt. Insgesamt ist eine Variablenbelegung $x \in \mathbb{R}^d$ in einem linearen Programm in kanonischer Form genau dann gültig, wenn $x \in \mathcal{P} := \mathcal{H}_1 \cap \dots \cap \mathcal{H}_m \cap \mathbb{R}_{\geq 0}^d$ gilt, wobei durch den Schnitt mit $\mathbb{R}_{\geq 0}^d$ der Bedingung $x \geq 0$ Rechnung getragen wird.

Beispiel

In der folgenden Abbildung ist ein lineares Programm mit den Nebenbedingungen

$$\begin{aligned} x_1 &\geq 0, & x_2 &\geq 0, \\ x_1 &\leq 2, & -x_2 &\leq -1, \\ & & -x_1 + x_2 &\leq 2 \end{aligned}$$

dargestellt. Der gültige Bereich \mathcal{P} ist in der Abbildung grau hervorgehoben.



Wir können die Menge $\mathbb{R}_{\geq 0}^d$ selbst als einen Schnitt von d Halbräumen darstellen:

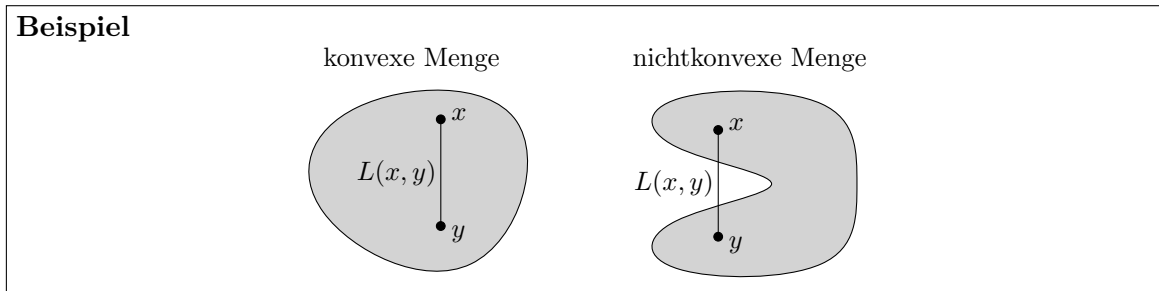
$$\mathbb{R}_{\geq 0}^d = \{x \in \mathbb{R}^d \mid -x_1 \leq 0\} \cap \dots \cap \{x \in \mathbb{R}^d \mid -x_d \leq 0\}.$$

Somit können wir insgesamt festhalten, dass \mathcal{P} der Schnitt von endlich vielen Halbräumen ist. Einen solchen Schnitt nennt man *konvexes Polyeder*, und wir werden das zu einem linearen Programm gehörende konvexe Polyeder auch als sein *Lösungspolyeder* bezeichnen. Wir sagen, dass ein lineares Programm *zulässig* ist, wenn sein Lösungspolyeder nicht leer ist.

Wie der Name schon sagt, ist jedes konvexe Polyeder \mathcal{P} *konvex*. Das bedeutet, dass für zwei beliebige Punkte $x, y \in \mathcal{P}$ auch alle Punkte auf der direkten Verbindungslinie zwischen x und y zum konvexen Polyeder \mathcal{P} gehören. Formal ausgedrückt ist eine Menge $X \subseteq \mathbb{R}^d$ konvex, wenn für jedes Paar $x, y \in X$ von Punkten

$$L(x, y) := \{\lambda x + (1 - \lambda)y \mid \lambda \in [0, 1]\} \subseteq X.$$

gilt. Hierbei ist $L(x, y)$ die Verbindungslinie zwischen x und y .



Lemma 48 *Jedes konvexe Polyeder \mathcal{P} ist konvex.*

Beweis. Sei \mathcal{P} der Durchschnitt der abgeschlossenen Halbräume $\mathcal{H}_1, \dots, \mathcal{H}_n$. Zunächst beweisen wir, dass jeder abgeschlossene Halbraum konvex ist. Seien $u \in \mathbb{R}^d$ und $w \in \mathbb{R}$ beliebig und seien x und y beliebige Punkte aus dem abgeschlossenen Halbraum $\mathcal{H} = \{z \mid u \cdot z \leq w\}$. Für jedes $\lambda \in [0, 1]$ gehört dann auch der Punkt $\lambda x + (1 - \lambda)y$ zu \mathcal{H} , denn

$$u \cdot (\lambda x + (1 - \lambda)y) = \lambda(u \cdot x) + (1 - \lambda)(u \cdot y) \leq \lambda w + (1 - \lambda)w = w.$$

Somit gehören alle Punkte aus $L(x, y)$ zum abgeschlossenen Halbraum \mathcal{H} , und damit ist bewiesen dass jeder beliebige abgeschlossene Halbraum konvex ist.

Nun müssen wir nur noch zeigen, dass der Schnitt zweier konvexer Mengen wieder konvex ist. Durch mehrfache Anwendung dieser Beobachtung folgt dann direkt das Lemma. Seien also $X \subseteq \mathbb{R}^d$ und $Y \subseteq \mathbb{R}^d$ konvexe Mengen, und seien $x, y \in X \cap Y$ beliebig. Da die Mengen X und Y konvex sind, gilt $L(x, y) \subseteq X$ und $L(x, y) \subseteq Y$. Dementsprechend gilt auch $L(x, y) \subseteq X \cap Y$ und somit ist $X \cap Y$ konvex. \square

Sei ein lineares Programm in kanonischer Form mit Lösungspolyeder \mathcal{P} gegeben. Die Tatsache, dass bei einem linearen Programm die erlaubten Variablenbelegungen eine konvexe Menge bilden, hat eine wichtige Konsequenz für die optimale Variablenbelegung. Wir sagen, eine Variablenbelegung $x \in \mathcal{P}$ ist *lokal optimal*, wenn es ein $\epsilon > 0$ gibt, für das es kein $y \in \mathcal{P}$ mit $\|x - y\| \leq \epsilon$ und $c \cdot y < c \cdot x$ gibt. Hierbei bezeichnet $\|x - y\|$ den euklidischen Abstand zwischen x und y , also $\|x - y\| = \sqrt{(x_1 - y_1)^2 + \dots + (x_d - y_d)^2}$. Eine Variablenbelegung ist also lokal optimal, wenn es eine Umgebung um sie herum gibt, in der es keine echt bessere zulässige Variablenbelegung gibt.

Theorem 49 *Sei ein lineares Programm in kanonischer Form mit Lösungspolyeder \mathcal{P} gegeben und sei $x \in \mathcal{P}$ eine lokal optimale Variablenbelegung. Dann ist x auch global optimal, d. h. es gibt kein $y \in \mathcal{P}$ mit $c \cdot y < c \cdot x$.*

Beweis. Angenommen, es gäbe ein $y \in \mathcal{P}$ mit $c \cdot y < c \cdot x$. Da das Lösungspolyeder \mathcal{P} konvex ist, liegt auch jeder Punkt auf der Verbindungslinie $L(x, y)$ in \mathcal{P} . Sei $z = \lambda x + (1 - \lambda)y$ mit $\lambda \in [0, 1)$ ein solcher Punkt. Dann können wir den Wert der Zielfunktion an der Stelle z schreiben als

$$c \cdot z = c \cdot (\lambda x + (1 - \lambda)y) = \lambda(c \cdot x) + (1 - \lambda)(c \cdot y) < c \cdot x,$$

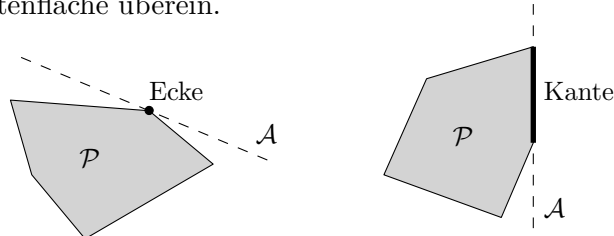
wobei wir bei der letzten Ungleichung $\lambda < 1$ und $c \cdot y < c \cdot x$ ausgenutzt haben. Das bedeutet, jeder Punkt $z \in L(x, y)$ mit $z \neq x$ ist eine echt bessere Variablenbelegung als x . Da jeder solche Punkt zu \mathcal{P} gehört und damit gültig ist, gibt es für jedes $\epsilon > 0$ einen Punkt $z \in \mathcal{P}$ mit $c \cdot z < c \cdot x$ und $\|x - z\| \leq \epsilon$. Damit ist x im Widerspruch zur Annahme nicht lokal optimal. \square

Wir betrachten nun den Rand von konvexen Polyedern genauer. Sei $\mathcal{P} \subseteq \mathbb{R}^d$ ein beliebiges konvexes Polyeder und sei eine Hyperebene $\mathcal{A} = \{x \in \mathbb{R}^d \mid u \cdot x = w\}$ mit beliebigem $u \in \mathbb{R}^d$

und $w \in \mathbb{R}$ gegeben. Ist $\mathcal{P} \cap \mathcal{A} \neq \emptyset$ und gilt $\mathcal{P} \subseteq \mathcal{H}$, wobei $\mathcal{H} = \{x \in \mathbb{R}^d \mid u \cdot x \leq w\}$ den zu \mathcal{A} gehörenden Halbraum bezeichnet, so sagen wir, dass \mathcal{A} eine *stützende Hyperebene* von \mathcal{P} ist. Für jede stützende Hyperebene \mathcal{A} nennen wir die Menge $\mathcal{P} \cap \mathcal{A}$ *Fläche von \mathcal{P}* . Eine nulldimensionale Fläche heißt auch *Ecke von \mathcal{P}* , eine eindimensionale Fläche heißt auch *Kante von \mathcal{P}* und eine Fläche der Dimension $d - 1$ heißt auch *Seitenfläche von \mathcal{P}* .

Beispiel

Die folgende Abbildung zeigt ein zweidimensionales Beispiel. Für $d = 2$ stimmen die Begriffe Kante und Seitenfläche überein.



Zielfunktion

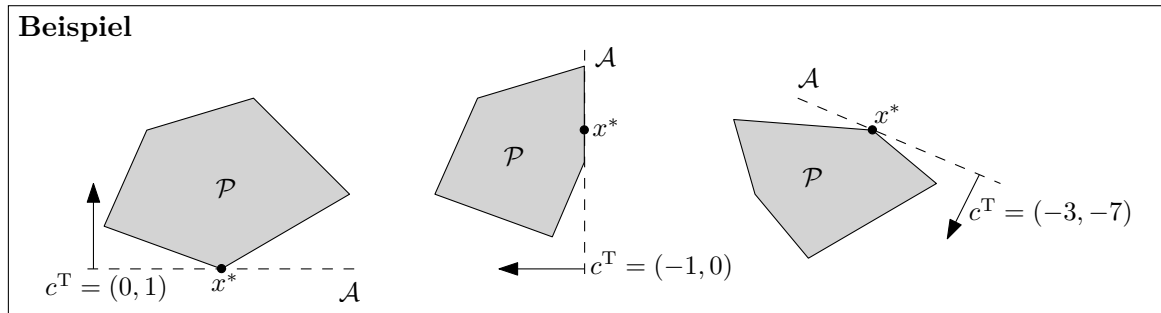
Wir bezeichnen ein lineares Programm als *unbeschränkt*, wenn der zu minimierende Zielfunktionswert innerhalb des Lösungspolyeders \mathcal{P} beliebig klein werden kann, d. h. wenn für alle $z \in \mathbb{R}$ eine Variablenbelegung $x \in \mathcal{P}$ mit $c \cdot x \leq z$ existiert. Ansonsten heißt das lineare Programm *beschränkt*.

Bisher haben wir lediglich eine geometrische Interpretation der Nebenbedingungen erarbeitet. Wir können auch die Zielfunktion geometrisch interpretieren. Dazu überlegen wir uns zunächst, welche Teilmengen von \mathbb{R}^d den gleichen Zielfunktionswert bezüglich einer linearen Zielfunktion haben. Sei $c \cdot x$ eine beliebige lineare Zielfunktion und sei $w \in \mathbb{R}$ beliebig. Die Menge

$$\{x \in \mathbb{R}^d \mid c \cdot x = w\}$$

bildet eine affine Hyperebene mit Normalenvektor c . Der Vektor c steht also senkrecht auf der Hyperebene, die die Punkte mit gleichem Zielfunktionswert enthält. Wir können nun wie folgt grafisch eine optimale Lösung eines linearen Programms bestimmen.

1. Finde ein $w \in \mathbb{R}$, für das der Schnitt der Hyperebene $\{x \in \mathbb{R}^d \mid c \cdot x = w\}$ mit dem Lösungspolyeder nicht leer ist.
2. Verschiebe die Hyperebene $\{x \in \mathbb{R}^d \mid c \cdot x = w\}$ solange parallel in Richtung $-c$ wie sie einen nichtleeren Schnitt mit dem Lösungspolyeder \mathcal{P} besitzt. Das ist äquivalent dazu, den Wert w solange zu verringern wie $\{x \in \mathbb{R}^d \mid c \cdot x = w\} \cap \mathcal{P} \neq \emptyset$ gilt.
3. Terminiert der zweite Schritt nicht, da jede betrachtete Hyperebene einen nichtleeren Schnitt mit \mathcal{P} besitzt, so ist das lineare Programm unbeschränkt. Ansonsten sei $\mathcal{A} = \{x \in \mathbb{R}^d \mid c \cdot x = w\}$ die letzte Hyperebene, für die $\mathcal{A} \cap \mathcal{P} \neq \emptyset$ gilt. Dann ist \mathcal{A} eine stützende Hyperebene von \mathcal{P} . Jeder Punkt $x^* \in \mathcal{A} \cap \mathcal{P}$ stellt eine optimale Variablenbelegung des linearen Programms dar, da $\mathcal{P} \subseteq \{x \in \mathbb{R}^d \mid c \cdot x \geq w\}$.



An dieser Stelle weisen wir darauf hin, dass man mit dem oben beschriebenen grafischen Verfahren zwar jedes lineare Programm optimal lösen kann, dass es sich jedoch um keinen Algorithmus handelt, der so ohne Weiteres implementiert werden kann. In höheren Dimensionen ist es beispielsweise bereits nicht trivial, im ersten Schritt ein geeignetes $w \in \mathbb{R}$ zu finden, für das der Schnitt von $\{x \in \mathbb{R}^d \mid c \cdot x = w\}$ und \mathcal{P} nicht leer ist.

9.2 Anwendung des Simplex Algorithmus

Wir beschreiben nun zunächst an einem Beispiel, wie das Simplex Verfahren von Dantzig im wesentlichen funktioniert, ein weiteres Kapitel beschreibt dieses Verfahren in einer Matrixschreibweise.

Betrachte das folgende LP System.

$$\begin{array}{ll}
 \text{Max} & 5x_1 + 4x_2 + 3x_3 \\
 \text{subject to} & 2x_1 + 3x_2 + x_3 \leq 5 \\
 & 4x_1 + x_2 + 2x_3 \leq 11 \\
 & 3x_1 + 4x_2 + 2x_3 \leq 8 \\
 & x_1, x_2, x_3 \geq 0
 \end{array} \tag{9.3}$$

Als erstes führen wir die Slackvariablen x_4 , x_5 und x_6 ein, um aus (9.3) ein Gleichungssystem zu erstellen. Betrachten wir dazu die Ungleichung $2x_1 + 3x_2 + x_3 \leq 5$ so sehen wir, dass die linke Seite maximal so groß ist, wie die rechte Seite und diesen Unterschied können wir als *Slack* bezeichnen und mit einer Variable $x_4 \geq 0$ versehen. Genauso gehen wir mit den anderen Ungleichungen um. Den Zielfunktionswert können wir bequem durch eine Variable z beschreiben und erhalten das System:

$$\begin{array}{ll}
 \text{Max} & z = 5x_1 + 4x_2 + 3x_3 \\
 \text{subject to} & x_4 = 5 - 2x_1 - 3x_2 - x_3 \\
 & x_5 = 11 - 4x_1 - x_2 - 2x_3 \\
 & x_6 = 8 - 3x_1 - 4x_2 - 2x_3 \\
 & x_1, x_2, x_3, x_4, x_5, x_6 \geq 0
 \end{array} \tag{9.4}$$

Jetzt ist klar, dass jede zulässige Lösung (x_1, x_2, x_3) von (9.3) zu einer zulässigen Lösung $(x_1, x_2, x_3, x_4, x_5, x_6)$ von (9.4) erweitert werden kann. Umgekehrt kann auch jede zulässige Lösung $(x_1, x_2, x_3, x_4, x_5, x_6)$ von (9.4) zu einer zulässigen Lösung (x_1, x_2, x_3) von (9.3) eingeschränkt werden. Desweiteren übertragen sich die Zielfunktionswerte und optimale Lösungen übertragen sich somit ebenfalls.

Wir notieren (9.4) in der folgenden Form:

$$\begin{array}{rcl}
 x_4 & = & 5 - 2x_1 - 3x_2 - x_3 \\
 x_5 & = & 11 - 4x_1 - x_2 - 2x_3 \\
 x_6 & = & 8 - 3x_1 - 4x_2 - 2x_3 \\
 \hline
 z & = & 5x_1 + 4x_2 + 3x_3
 \end{array} \tag{9.5}$$

Die Simplex Methode versucht nun eine zulässige Lösung sukzessive zu verbessern. D.h., wir wollen von einer zulässigen Lösung $(x_1, x_2, x_3, x_4, x_5, x_6)$ zu einer weiteren zulässigen Lösung $(x'_1, x'_2, x'_3, x'_4, x'_5, x'_6)$ übergehen, so dass

$$5x'_1 + 4x'_2 + 3x'_3 > 5x_1 + 4x_2 + 3x_3$$

gilt. Wenn wir das sukzessive wiederholen, soll die optimale Lösung erreicht werden. Zunächst müssen wir eine zulässige Lösung finden. Unser Beispiel ist unproblematisch, wir setzen die *Entscheidungsvariablen* x_1, x_2, x_3 auf Null und erhalten die zulässige Lösung

$$x_1 = 0, x_2 = 0, x_3 = 0, x_4 = 5, x_5 = 11, x_6 = 8 \tag{9.6}$$

mit Wert $z = 0$. Wie wir später sehen werden gehören hierbei $x_4 = 5, x_5 = 11, x_6 = 8$ zur zulässigen *Basislösung*. Wenn wir nun den Zielfunktionswert z verbessern wollen, so können wir beispielsweise x_1 erhöhen und $x_2 = x_3 = 0$ beibehalten und bekommen $z = 5x_1 > 0$. Wie groß x_1 werden darf, hängt von den Nebenbedingungen ab. Falls wir beispielsweise $x_1 = 2$ setzen, erhalten wir $z = 10$ und mit $x_4 = 0, x_5 = 3$ und $x_6 = 2$ eine zulässige Lösung. Für $x_1 = 3$ wäre $z = 15$ zwar besser, aber $x_4 = -1$ ist nicht zulässig. Wie groß dürfen wir also x_1 machen? Wir sagen auch, dass x_1 nun das *Pivotelement* ist, dass wir ausgewählt haben, um die Lösung zu verbessern.

Die erste Gleichung sagt, dass $x_4 = 5 - 2x_1 - 3x_2 - x_3 \geq 0$ gelten soll, daraus folgt dann $x_1 \leq \frac{5}{2}$. Aus der zweiten Gleichung folgt analog $x_1 \leq \frac{11}{4}$ und aus der dritten Gleichung $x_1 \leq \frac{8}{3}$. Dann ist $x_1 \leq \frac{5}{2}$ die stringenteste Bedingung und wir haben die nächste zulässige Lösung:

$$x_1 = \frac{5}{2}, x_2 = 0, x_3 = 0, x_4 = 0, x_5 = 1, x_6 = \frac{1}{2} \tag{9.7}$$

mit der Verbesserung $z = \frac{25}{2}$.

Wie machen wir nun weiter? Damit wir ähnlich verfahren können, formen wir (9.5) so um, dass x_1, x_5 und x_6 die neue *Basislösung* bilden. Aus der ersten Gleichung erhalten wir:

$$x_1 = \frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 \tag{9.8}$$

Diese können wir nun in die weiteren Gleichungen von (9.5) einsetzen und bekommen so:

$$\begin{aligned}
 x_5 &= 11 - 4 \left(\frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 \right) - x_2 - 2x_3 \\
 &= 1 + 5x_2 + x_4 \\
 x_6 &= 8 - 3 \left(\frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 \right) - 4x_2 - 2x_1 \\
 &= \frac{1}{2} + \frac{1}{2}x_2 - \frac{1}{2}x_3 + \frac{3}{2}x_4 \\
 z &= 5 \left(\frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 \right) + 4x_2 + 3x_3 \\
 &= \frac{25}{2} - \frac{7}{2}x_2 + \frac{1}{2}x_3 - \frac{5}{2}x_4
 \end{aligned} \tag{9.9}$$

Daraus erhalten wir insgesamt das neue System

$$\begin{array}{rcl}
 x_1 & = & \frac{5}{2} - \frac{3}{2}x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4 \\
 x_5 & = & 1 + 5x_2 + 2x_4 \\
 x_6 & = & \frac{1}{2} + \frac{1}{2}x_2 - \frac{1}{2}x_3 + \frac{3}{2}x_4 \\
 \hline
 z & = & \frac{25}{2} - \frac{7}{2}x_2 + \frac{1}{2}x_3 - \frac{5}{2}x_4
 \end{array} \tag{9.10}$$

Wenn wir nach der bisherigen Methode die Nicht-Basisvariablen vergrößern wollen, um den Wert $z = \frac{25}{2} - \frac{7}{2}x_2 + \frac{1}{2}x_3 - \frac{5}{2}x_4$ zu erhöhen, so bleibt uns als Auswahl nur x_3 übrig, x_2 und x_4 haben negative Koeffizienten. Mit $x_2 = x_4 = 0$ folgt aus $x_1 \geq 0$ die Beschränkung $x_3 \leq 5$. $x_5 \geq 0$ gibt keine Beschränkung und aus $x_6 \geq 0$ folgt $x_3 \leq 1$, die stärkste Einschränkung. Damit ergibt sich die Lösung:

$$x_1 = 2, x_2 = 0, x_3 = 1, x_4 = 0, x_5 = 1, x_6 = 0 \tag{9.11}$$

mit der Verbesserung $z = 13$. Nun erstellen wir wieder ein System, bei dem die Basisvariablen x_1, x_3 und x_5 auf der linken Seite stehen und die Variablen x_2, x_4 und x_6 mit Wert 0 rechts. Dazu betrachten wir wieder die neue Variable $x_3 = 1 + x_2 + 3x_4 - 2x_6$, die wir aus der Gleichung für x_6 auf (9.10) gewinnen. Einsetzen in die verbleibenden Gleichungen ergibt dann:

$$\begin{array}{rcl}
 x_3 & = & 1 + x_2 + 3x_4 - 2x_6 \\
 x_1 & = & 2 - 2x_2 - 2x_4 + x_6 \\
 x_5 & = & 1 + 5x_2 + 2x_4 \\
 \hline
 z & = & 13 - 3x_2 - x_4 - x_6
 \end{array} \tag{9.12}$$

Wir erhalten die Lösung:

$$x_1 = 2, x_2 = 0, x_3 = 1, x_4 = 0, x_5 = 1, x_6 = 0. \tag{9.13}$$

Jetzt kann das Verfahren nicht weitergeführt werden, da eine Erhöhung der Werte von x_2, x_4 oder x_6 den Zielfunktionswert z nur verringert. Jede zulässige Lösung von (9.12) ist eine zulässige Lösung von (9.5) und umgekehrt, da wir nur Äquivalenzumformungen vorgenommen haben. Jede zulässige Lösung $(x_1, x_2, x_3, x_4, x_5, x_6)$ mit $x_2 \geq 0, x_4 \geq 0$ und $x_6 \geq 0$ ergibt somit $z \leq 13$ und $z = 13$ muss der optimale Wert sein! Das Verfahren endet also via Konstruktion mit der optimalen Lösung.

Wir benutzen die folgende Terminologie: Die Variablen auf der linken Seite sind stets die Basisvariablen. Die restlichen Variablen die Nicht-Basisvariablen. In jedem Schritt wird eine Nicht-Basisvariable erhöht und gelangt in die Basis. Dafür muss eine Variable aus der Basis weichen. Welche Variable jeweils die Basis betritt, nennen wir die *Pivotregel*. Verlassen muss die Basis die Variable mit der stärksten Beschränkung. Der Wert dieser Variable wird in der nächsten Lösung auf 0 gesetzt.

9.2.1 Allgemeines Verfahren: Dictionary

Das oben beschriebene Verfahren können wir allgemein wie folgt beschreiben. Aus einem System in kanonischer Form

$$\begin{aligned}
\text{Max} \quad & \sum_{j=1}^n c_j x_j \\
\text{subject to} \quad & \sum_{j=1}^n a_{ij} x_j \leq b_i \quad (i = 1, 2, \dots, m) \\
& x_j \geq 0 \quad (i = 1, 2, \dots, n)
\end{aligned} \tag{9.14}$$

Wenn wir nun die Slackvariablen $x_{n+1}, x_{n+2}, \dots, x_{n+m}$ einführen und den Zielfunktionswert mit z beschreiben, erhalten wir das äquivalente *Dictionary*:

Im Rahmen des Simplex Algorithmus wird also folgendes System mit insgesamt m Basisvariablen $x_{n+1}, x_{n+2}, \dots, x_{n+m}$ aufgebaut:

$$\begin{aligned}
x_{n+i} &= b_i - \sum_{j=1}^n a_{ij} x_j \quad (i = 1, 2, \dots, m) \\
z &= \sum_{j=1}^n c_j x_j
\end{aligned} \tag{9.15}$$

Im Rahmen des Simplex-Algorithmus kann jede zulässige Lösung x_1, x_2, \dots, x_n von (9.19) durch eine zulässige Lösung x_1, x_2, \dots, x_{n+m} repräsentiert werden, wobei $x_{n+1}, x_2, \dots, x_{n+m}$ gemäß (9.15) definiert sind. In jeder Iteration wird eine zulässige Lösung x_1, x_2, \dots, x_{n+m} zu einer zulässigen Lösung $x'_1, x'_2, \dots, x'_{n+m}$ verändert, so dass

$$\sum_{j=1}^n c_j x'_j > \sum_{j=1}^n c_j x_j$$

gilt. Dabei formen wir das Dictionary (9.15) entsprechend um.

Nach Übertragung von (9.19) nach (9.15) sollen alle weiteren Umformungen in Dictionaries folgendermaßen korrekt sein:

- Jede Lösung eines weiteren Dictionaries muss eine Lösung von (9.15) sein und umgekehrt.
- In den Gleichungen des Dictionaries müssen m Variablen von x_1, x_2, \dots, x_{n+m} und z durch die verbleibenden n Variablen dargestellt werden.

Diese Bedingungen haben wir durch Äquivalenzumformungen im vorherigen Kapitel erreicht. Zum Beispiel war die Basislösung (9.6) des Dictionaries (9.5) immer noch eine Lösung des letzten Dictionaries (9.12). Umgekehrt war die optimale Lösung

$$x_1 = 2, x_2 = 0, x_3 = 1, x_4 = 0, x_5 = 1, x_6 = 0 \tag{9.16}$$

von (9.12) auch eine Lösung von (9.5)

Wenn wir ein Dictionary so wie oben festlegen, dann erhalten wir (möglicherweise) durch das Festlegen der n Variablen der linken Seite auf 0 eine *zulässige* Lösung. Wenn das der Fall ist, dann sprechen wir von einem *zulässigen* Dictionary. Dann beschreibt jedes zulässige Dictionary eine zulässige Lösung. Es ist möglich, dass zulässige Lösungen gar nicht durch ein zulässiges Dictionary beschrieben werden können. Beispielsweise kann die zulässige Lösung

$$x_1 = 1, x_2 = 0, x_3 = 1, x_4 = 2, x_5 = 5, x_6 = 3 \tag{9.17}$$

wird, haben wir danach bereits ein zulässiges Dictionary und können den Simplex-Algorithmus starten mit x_0 als Basisvariable.

Wenn wir es nun im Laufe der Iteration wieder schaffen, x_0 aus der Basis zu entfernen, dann haben wir immer noch eine zulässige Lösung für das Hilfs-Dictionary. Da x_0 nun eine Nicht-Basisvariable ist, muss der Wert der Lösung $w = -x_0 = 0$ sein, weil wir die Lösungen nur verbessert haben und $w = -x_0$ stets auch die gültige Lösung ist. Diese Lösung ist also optimal mit $x_0 = 0$ und wir können dann x_0 auch wieder entfernen und erhalten dann ein zulässiges Dictionary für das Ausgangssystem.

Dazu schauen wir uns das folgende Beispiel an! Betrachte:

$$\begin{array}{rcl}
 \text{Max} & x_1 & - 4x_2 + x_3 \\
 \text{subject to} & 2x_1 & - x_2 + 2x_3 \leq 4 \\
 & 2x_1 & - 3x_2 + x_3 \leq -5 \\
 & -x_1 & + x_2 - 2x_3 \leq -1 \\
 & & x_1, x_2, x_3 \geq 0
 \end{array} \tag{9.21}$$

Daraus bauen wir in *Phase I* das unzulässige Dictionary:

$$\begin{array}{rcl}
 x_4 & = & 4 - 2x_1 + x_2 - 2x_3 + x_0 \\
 x_5 & = & -5 - 2x_1 + 3x_2 - x_3 + x_0 \\
 x_6 & = & -1 + x_1 - x_2 + 2x_3 + x_0 \\
 \hline
 w & = & -x_0
 \end{array} \tag{9.22}$$

Wenn wir aber x_0 in die Basis bringen und die am wenigsten zulässige Variable x_5 ersetzen, erhalten wir das zulässige Dictionary:

$$\begin{array}{rcl}
 x_0 & = & 5 + 2x_1 - 3x_2 + x_3 + x_5 \\
 x_4 & = & 9 - 2x_2 - x_3 + x_5 \\
 x_6 & = & 4 + 3x_1 - 4x_2 + 3x_3 + x_5 \\
 \hline
 w & = & -5 - 2x_1 + 3x_2 - x_3 - x_5
 \end{array} \tag{9.23}$$

Pivotieren mit x_2 und entfernen von x_6 (Beschränkung $x_2 \leq 1$) ergibt:

$$\begin{array}{rcl}
 x_2 & = & 1 + 0.75x_1 + 0.75x_3 + 0.25x_5 - 0.25x_6 \\
 x_0 & = & 2 - 0.25x_1 - 1.25x_3 + 0.25x_5 + 0.75x_6 \\
 x_4 & = & 7 - 1.5x_1 - 2.5x_3 + 0.5x_5 + 0.5x_6 \\
 \hline
 w & = & -2 + 0.25x_1 + 1.25x_3 - 0.25x_5 - 0.75x_6
 \end{array} \tag{9.24}$$

Nun haben wir die Wahl zwischen x_1 und x_3 als Pivotelement. x_1 müsste mit x_4 tauschen, da es die größte Einschränkung ergibt, sonst gäbe es danach wieder eine unzulässige Lösung. Wählen wir aber x_3 als Pivotelement, so können wir x_0 aus der Basis entfernen und erhalten:

$$\begin{array}{rcl}
 x_3 & = & 1.6 - 0.2x_1 + 0.2x_5 + 0.6x_6 - 0.8x_0 \\
 x_2 & = & 2.2 + 0.6x_1 + 0.4x_5 + 0.2x_6 + 0.6x_0 \\
 x_4 & = & 3 - x_1 - x_6 + 2x_0 \\
 \hline
 w & = & -x_0
 \end{array} \tag{9.25}$$

Der optimale Wert des Hilfsystems ist offensichtlich 0 (und das muss auch so sein, wenn x_0 die Basis verlässt) und wir bekommen eine zulässige Lösung $x_1 = 0$, $x_2 = 2.2$ und $x_3 = 3$ des Originalproblems.

Jetzt können wir einfach x_0 streichen und die ersten drei Gleichungen verwenden. Lediglich die Zielfunktion muss durch

$$\begin{aligned} z &= x_1 - x_2 + x_3 & (9.26) \\ &= x_1 - (2.2 + 0.6x_1 + 0.4x_5 + 0.2x_6) + (1.6 - 0.2x_1 + 0.2x_5 + 0.6x_6) \\ &= -0.6 + 0.2x_1 - 0.2x_5 + 0.4x_6 \end{aligned}$$

angepasst werden und wir erhalten das zulässige Dictionary

$$\begin{array}{rcll} x_3 & = & 1.6 & - & 0.2x_1 & + & 0.2x_5 & + & 0.6x_6 \\ x_2 & = & 2.2 & + & 0.6x_1 & + & 0.4x_5 & + & 0.2x_6 \\ x_4 & = & 3 & - & x_1 & & & - & x_6 \\ \hline w & = & -0.6 & + & 0.2x_1 & - & 0.2x_5 & + & 0.4x_6 \end{array} \quad (9.27)$$

des Ausgangs-LP.

Die Policy ist, dass wir immer sofort x_0 wieder aus der Basis entfernen, wenn es aufgrund der Beschränkungen geht! Es muss noch argumentiert werden, was passiert, wenn das Dictionary dieser ersten Phase in eine optimale Lösung läuft, in der x_0 in der Basis verbleibt.

- Falls x_0 eine Basisvariable bleibt, das Verfahren terminiert und der Wert $w = -x_0$ am Ende ungleich 0 ist, so ist das Ausgangsproblem gar nicht zulässig gewesen.
- Falls x_0 eine Basisvariable bleibt, das Verfahren terminiert und der Wert w gleich 0 ist, so hatten wir einen letzten Schritt, in dem der Wert $w = -x_0$ zu 0 wurde. Dann muss aber x_0 im letzten Schritt zu 0 geworden sein und war somit ein Kandidat für das Verlassen der Basis, ein Widerspruch, dieser Fall kann gar nicht eintreten!

Iteration

Die Iteration besteht im wesentlichen daraus, die Variable auszuwählen, die in die Basis eintritt und die Variable zu wählen, die die Basis verlässt.

Wahl der eintretenden Variable: Falls in der Beschreibung von $z = z^* + \sum_{j \in N} \bar{c}_j x_j$ für die Menge der Nicht-Basisvariablen N gilt, dass alle $\bar{c}_j \leq 0$ sind, dann haben wir bereits argumentiert, dass die Lösung z^* (für ein Maximierungsproblem) optimal sein muss, denn für jede andere zulässige Lösung gilt diese Gleichung ebenfalls. Ansonsten können mehrere Variablen zur Auswahl stehen. Die Wahl der eintretenden Variable hat einen wesentlichen Anteil an der Terminierung und Laufzeit der Simplex Methode, das werden wir am Ende des Kapitels noch diskutieren.

Wahl der austretenden Variable: Im Prinzip ist die Wahl der austretenden Variable dadurch fixiert, dass wir den Wert der eintretenden Variable so maximieren wollen, dass die austretende Variable in der Lösung den Wert 0 erhält und keine andere Variable negativ wird. D.h., die Variable mit der stärksten Einschränkung verlässt automatisch die Basis. Deshalb ist die Iteration im wesentlichen durch die eintretende Variable bestimmt. Besonderheiten sind die folgenden.

Im folgenden Dictionary muss im nächsten Schritt x_3 als eintretende Variable gewählt werden. Diese Variable hat aber in dem Dictionary gar keine Beschränkung. Das macht aber nichts, denn der Simplex Algorithmus bemerkt an dieser Stelle, dass das Problem unbeschränkt ist, wir können $x_1 = d$ beliebig erhöhen, und erhalten stets eine zulässige Lösung.

$$\begin{array}{rcll} x_2 & = & 5 & + & x_3 & + & 3x_1 \\ x_4 & = & 7 & - & 2x_3 & & \\ \hline z & = & 15 & - & 3x_3 & + & x_1 \end{array} \quad (9.28)$$

Das gilt natürlich im allgemeinen. Falls wir keine austretende Variable finden ist das Problem sogleich *unbeschränkt*.

Degeneriertheit: Es kann passieren, dass es mehrere Kandidaten für die austretende Variable gibt. Das bedeutet, dass dann auch eine Basisvariable den Wert 0 annehmen kann.

$$\begin{array}{rccccrcr}
 x_4 & = & 5 & & & - & 2x_3 \\
 x_5 & = & 3 & - & 2x_1 & + & 4x_2 & - & 6x_3 \\
 x_6 & = & 2 & + & x_1 & - & 3x_2 & - & 4x_3 \\
 \hline
 z & = & & & 2x_1 & - & x_2 & + & 8x_3
 \end{array} \tag{9.29}$$

Wählen wir hier x_3 als eintretende Variable, so bieten alle drei Gleichungen die gleiche Einschränkung $x_3 \leq \frac{1}{2}$ und wir erhalten das System

$$\begin{array}{rccccrcr}
 x_3 & = & 0.5 & & & - & 0.5x_4 \\
 x_5 & = & & - & 2x_1 & + & 4x_2 & + & 3x_4 \\
 x_6 & = & & & x_1 & - & 3x_2 & + & 2x_4 \\
 \hline
 z & = & 4 & & 2x_1 & - & x_2 & - & 4x_3
 \end{array} \tag{9.30}$$

und in der nächsten Iteration können wir den Zielfunktionswert nicht erhöhen, da die eintretende Variable x_1 den Wert 0 annimmt. Wir erhalten:

$$\begin{array}{rccccrcr}
 x_1 & = & & 2x_2 & + & 1.5x_4 & - & 0.5x_5 \\
 x_3 & = & 0.5 & & & - & 0.5x_4 \\
 x_6 & = & & x_2 & + & 3.5x_4 & - & 0.5x_5 \\
 \hline
 z & = & 4 & + & 3x_2 & - & x_4 & - & x_5
 \end{array} \tag{9.31}$$

Wir nennen Basislösungen, in denen eine oder mehrere Basisvariablen den Wert 0 haben *degenerierte* Basislösungen. Eine Iteration des Algorithmus nennen wir *degeneriert*, falls der Wert der Lösung nicht steigt. Das kann nur passieren, wenn vorher eine degenerierte Basislösung entstanden ist. Es kann passieren, dass der Simplex Algorithmus eine Weile degeneriert und der Zielfunktionswert gleich bleibt, bis schließlich eine weitere Erhöhung möglich ist. Der Leser darf gerne prüfen, dass im obigen Dictionary noch eine weitere degenerierte Iteration stattfindet, dann aber die nächste Iteration zur optimalen Lösung führt.

Terminierung

Das obige Beispiel führt uns dazu, zu untersuchen, ob der Algorithmus immer terminiert. In der Tat ist es möglich, dass Zykel entstehen, die immer wieder und ohne Verbesserung durchlaufen werden. Es ist klar, dass es sich dabei um degenerierte Iterationen handeln muss, sonst würde der Wert ja erhöht. Wir geben im folgenden spezielle Regeln vor, um einen Zykel zu erzwingen.

1. Die eintretende Variable ist immer diejenige Nicht-Basisvariable mit dem größten Koeffizienten.
2. Bei mehreren Kandidaten für die austretende Variable, wähle diejenige mit kleinsten Index.

Nach sechs Iterationen gelangen wir zum Ausgangs-Dictionary zurück

$$\begin{array}{rccccrcr}
 x_5 & = & & - & 0.5x_1 & + & 5.5x_2 & + & 2.5x_3 & - & 9x_4 \\
 x_6 & = & & - & 0.5x_1 & + & 1.5x_2 & + & 0.5x_3 & - & x_4 \\
 x_7 & = & 1 & - & x_1 & & & & & & \\
 \hline
 z & = & & & 10x_1 & - & 57x_2 & - & 9x_3 & - & 24x_4
 \end{array} \tag{9.32}$$

Iteration I ergibt:

$$\begin{array}{rcl}
 x_1 & = & 11x_2 + 5x_3 - 18x_4 - 2x_5 \\
 x_6 & = & -4x_2 - 2x_3 + 8x_4 + x_5 \\
 x_7 & = & 1 - 11x_2 - 5x_3 + 18x_4 + 2x_5 \\
 \hline
 z & = & 53x_2 + 41x_3 - 204x_4 - 20x_5
 \end{array} \tag{9.33}$$

Iteration II ergibt:

$$\begin{array}{rcl}
 x_2 & = & -0.5x_3 + 2x_4 + 0.25x_5 - 0.25x_6 \\
 x_1 & = & -0.5x_3 + 4x_4 + 0.75x_5 - 2.75x_6 \\
 x_7 & = & 1 + 0.5x_3 - 4x_4 - 0.75x_5 - 13.25x_6 \\
 \hline
 z & = & 14.5x_3 - 98x_4 - 6.75x_5 - 13.25x_6
 \end{array} \tag{9.34}$$

Iteration III ergibt:

$$\begin{array}{rcl}
 x_3 & = & 8x_4 + 1.5x_5 - 5.5x_6 - 2x_1 \\
 x_2 & = & -2x_4 - 0.5x_5 - 2.5x_6 + x_1 \\
 x_7 & = & 1 - x_1 \\
 \hline
 z & = & 18x_4 + 15x_5 - 93x_6 - 29x_1
 \end{array} \tag{9.35}$$

Iteration IV ergibt:

$$\begin{array}{rcl}
 x_4 & = & -0.25x_5 + 1.25x_6 + 0.5x_1 - 0.5x_2 \\
 x_3 & = & -0.5x_5 + 4.5x_6 - 2x_1 - 4x_2 \\
 x_7 & = & 1 - x_1 \\
 \hline
 z & = & 10.5x_5 - 70.5x_6 - 20x_1 - 9x_2
 \end{array} \tag{9.36}$$

Iteration V ergibt:

$$\begin{array}{rcl}
 x_5 & = & 9x_6 + 4x_1 - 8x_2 - 2x_3 \\
 x_4 & = & -x_6 - 0.5x_1 + 1.5x_2 + 0.5x_3 \\
 x_7 & = & 1 - x_1 \\
 \hline
 z & = & 24x_6 + 22x_1 - 93x_2 - 21x_3
 \end{array} \tag{9.37}$$

Iteration VI ergibt:

$$\begin{array}{rcl}
 x_5 & = & -0.5x_1 + 5.5x_2 + 2.5x_3 - 9x_4 \\
 x_6 & = & -0.5x_1 + 1.5x_2 + 0.5x_3 - x_4 \\
 x_7 & = & 1 - x_1 \\
 \hline
 z & = & 10x_1 - 57x_2 - 9x_3 - 24x_4
 \end{array} \tag{9.38}$$

Es stellt sich sofort die Frage, ob es noch andere Möglichkeiten der Nicht-Terminierung gibt. Glücklicherweise ist das nicht der Fall und der Algorithmus kann somit darauf reagieren. Zudem kommen Zykel in der Praxis kaum vor.

Es gilt das folgende Theorem:

Theorem 50 Falls das Simplex Verfahren nicht terminiert, muss es einen Zykel geben.

Beweis. Es gibt stets nur endlich viele Möglichkeiten m Basisvariablen aus den $n + m$ Variablen auszuwählen. Wenn die Simplex Methode also nicht terminiert, dann muss die gleiche Basis in zwei verschiedenen Iterationen verwendet werden.

Wir wollen nun zeigen, dass dann die Lösungen identisch sind und sich somit der Wert die ganze Zeit nicht geändert hat. (Wenn wir die Simplex Methode wie später angegeben direkt in Matrix-Schreibweise auffassen sind wir hier bereits fertig, wir können es aber auch direkt beweisen.)

Wenn wir für die Menge der Basisvariablen B die zwei Dictionaries

$$\begin{array}{r} x_i = b_i - \sum_{j \notin B}^n a_{ij} x_j \quad (i \in B) \\ \hline z = \sum_{j \notin B}^n c_j x_j \end{array} \quad (9.39)$$

$$\begin{array}{r} x_i = b_i^* - \sum_{j \notin B}^n a_{ij}^* x_j \quad (i \in B) \\ \hline z = \sum_{j \notin B}^n c_j^* x_j \end{array} \quad (9.40)$$

betrachten, so wissen wir, dass jede Lösung x_1, x_2, \dots, x_{n+m} von (9.39) eine Lösung von (9.40) und umgekehrt.

Dann sind für $x_k \notin B$ und ein beliebiges t die Werte

$$x_k = t, \quad x_j = 0 \text{ für } j \neq k, j \notin B, \quad x_i = b_i - a_{ij}t \text{ für } i \in B, \quad z = v + c_k t \quad (9.41)$$

eine Lösung von (9.39). Da es auch eine Lösung von (9.39) ist, gilt:

$$b_i - a_{ij}t = b_i^* - a_{ij}^*t \text{ für alle } i \in B \text{ und } v + c_k t = v^* + c_k^* t \quad (9.42)$$

für beliebige Werte von t . Mit $t = 0$ folgt die Behauptung! \square

9.3 Beispiel Anwendung: Integer LP

Linearer Programmierung hat viele Anwendungen, wir zeigen hier ein Beispiel aus unseren aktuellen Forschungsgebieten. Es geht darum ein ausbrechendes Feuer zu stoppen und wir werden die Zeitkomponente geschickt modellieren.

In einem diskreten, einfachen Zellenmodell (Gitterzellen) kann man sich vorstellen, dass das Feuer sich von einer Zelle aus in jedem Zeitschritt zu den 4 direkten Nachbarzellen ausbreitet. Die Feuerwehr kann in jedem Schritt k Zellen außerhalb blocken. Das Feuer startet bei s und die Feuerbekämpfer dürfen in jedem Schritt erst $k = 2$ Zellen blocken, danach breitet sich das Feuer aus, siehe Abbildung 9.1.

Die Frage ist nun, wie schnell kann man das Feuer eindämmen bzw. bei welcher Strategie verbrennen die wenigsten Zellen. Interessanterweise kann man zur Berechnung optimaler Lösungen auf Lineare Programmierung zurückgreifen. Entscheidend ist, dass wir eine Idee für die Modellierung entwickeln.

Für $k = 2$ kann man tatsächlich eine optimale Strategie für die Zeit- und auch die Flächenoptimierung finden. Abbildung 9.2 zeigt die optimale Strategie!

Wie kann man nun zeigen, dass diese Strategie optimal ist? Wir verwenden für den Beweis einfach ein Lineares Programm mit ganzzahliger Lösung ein sogenanntes Integer LP.

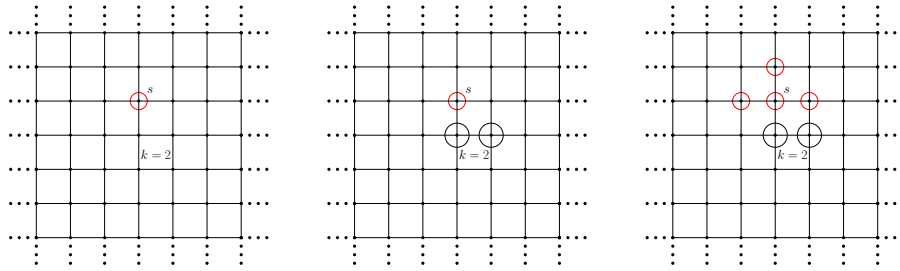


Abbildung 9.1: Das Feuer startet bei s . In jedem Schritt dürfen erst $k = 2$ Zellen außerhalb des Feuers blockiert werden, dann breitet sich das Feuer von jeder Position auf die noch möglichen Nachbarzellen aus.

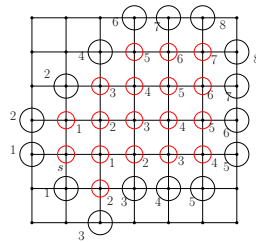


Abbildung 9.2: Für $k = 2$ kann das Feuer in 8 Schritten gestoppt werden, es verbrennen nur 18 Zellen.

Da wir bereits wissen, dass wir eine Lösung mit 8 Schritten haben, beschränken wir die Zeit auf ein festes $T > 8$ und betrachten ein festes $l \times l$ für eine Konstante l .

Sei $L = \{(x, y) \mid |x| \leq l \text{ and } |y| \leq l\}$ und $0 \leq t \leq T$. Für $v \in L$ sei $N(v)$ die 4-Nachbarschaft in L . Für das Integer LP benutzen wir folgende Menge von Variablen für *alle* $v \in L$ und *alle* Zeitschritte t !

$$b_{v,t} = \begin{cases} 1 & : \text{Knoten } v \in L \text{ brennt vor oder zum Zeitpunkt } t \\ 0 & : \text{sonst} \end{cases}$$

$$d_{v,t} = \begin{cases} 1 & : \text{Knoten } v \in L \text{ wird vor oder zum Zeitpunkt } t \text{ blockiert} \\ 0 & : \text{sonst} \end{cases}$$

Wir wollen beispielsweise die Anzahl der brennenden Zellen minimieren, dazu definieren wir folgendes LP:

$$\begin{aligned} &\text{Minimize} && \sum_{v \in L} b_{v,T} \\ \text{s. t.} &&& b_{v,t} + d_{v,t} - b_{w,t-1} \geq 0 && : \forall v \in L, v \in N(w), 1 \leq t \leq T \\ &&& b_{v,t} + d_{v,t} \leq 1 && : \forall v \in L, 1 \leq t \leq T \\ &&& b_{v,t} - b_{v,t-1} \geq 0 && : \forall v \in L, 1 \leq t \leq T \\ &&& d_{v,t} - d_{v,t-1} \geq 0 && : \forall v \in L, 1 \leq t \leq T \\ &&& \sum_{v \in L} (d_{v,t} - d_{v,t-1}) \geq 2 && : \forall 1 \leq t \leq T \\ &&& b_{v,0} = 1 && : v \in L \text{ ist der Start } (0, 0) \\ &&& b_{v,0} = 0 && : v \in L \text{ ist nicht der Start } (0, 0) \\ &&& d_{v,0} = 0 && : \forall v \in L \\ &&& d_{v,t}, b_{v,t} \in \{0, 1\} && : \forall v \in L, 1 \leq t \leq T \end{aligned}$$

- Die erste Ungleichung steuert die Ausbreitung des Feuers. Der Nachbar einer zum Zeitpunkt $t - 1$ brennenden Zelle ist zum nächsten Zeitpunkt t entweder blockiert oder brennt im nächsten Zeitschritt. Durch die Ungleichung kann auch eine Zelle zufällig Feuer fangen, das passiert aber nicht in der optimalen Lösung.
- Die nächste Ungleichung verhindert, dass eine beschützte Zelle brennt und dass eine brennende Zelle geblockt wird, es kann nur eine Bedingung oder keine gelten.
- Die nächsten beiden Ungleichungen sorgen dafür, dass der Status der Zellen erhalten bleibt.
- Die Summe kontrolliert, dass wir zu jedem Zeitpunkt nur zwei weitere Zellen blockieren dürfen. In einer optimalen Lösung werden die Blocker sofort eingesetzt.
- Zu Beginn brennt die Ursprungszelle und keine andere Zelle ist beschützt worden.
- Die Werte der Variablen sind entweder 1 oder 0.

Wird dieses Programm mit $l = 6$ und $T = 9$ und Wurzel $(0, 0)$ zum Beispiel mit dem in der Vorlesung vorgestellten ILP-Solver SCIP gestartet, so erhält man die Lösung aus Abbildung 9.2. Dieses Beispiel zeigt wiederum, dass man auch in der Theorie computerunterstützt Beweise führen kann.

9.4 Algebraische Interpretation mit Lösungspolyeder

In der geometrischen Erklärung der Linearen Programme haben wir uns die Lösungspolyeder und die optimalen Lösungen als Rand des Lösungspolyeders geometrisch erklärt. Im Simplex-Verfahren haben wir Gleichungssysteme (Dictionaries) solange umgeformt, bis die Lösung nicht mehr verbessert werden konnte. Diese Dinge sind natürlich identisch und die folgende algebraische Interpretation bringt die Sachverhalte zusammen.

Die kanonische Form eignet sich gut zur anschaulichen Darstellung von linearen Programmen. Bei dem Entwurf von Algorithmen ist es jedoch einfacher, lineare Programme in Gleichungsform zu betrachten. In diesen Abschnitt gehen wir deshalb davon aus, dass ein solches lineares Programm gegeben ist, in dem die Zielfunktion $c \cdot x$ unter den Nebenbedingungen $Ax = b$ und $x \geq 0$ minimiert werden soll. Dabei seien $c \in \mathbb{R}^d$, $b \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times d}$ und $x \in \mathbb{R}^d$ sei ein Vektor, der aus d Variablen besteht. Es bezeichne \mathcal{P} das zugehörige Lösungspolyeder.

Alternativ können wir hier auch sehr leicht die Gleichungsvariante (Dictionary) des Simplex-Verfahrens betrachten indem wir die Variablen alle auf eine Seite setzen. Beispielsweise erhalten wir für das Start-Dictionary (9.15) das Gleichungssystem

$$\begin{aligned} x_{n+i} - \sum_{j=1}^n a_{ij}x_j &= b_i & (i = 1, 2, \dots, m) \\ z &= \sum_{j=1}^n c_jx_j \end{aligned} \tag{9.43}$$

und die kombinierte Matrix $A' := E^{m,m}(-A) \in \mathbb{R}^{m,n+m}$ nimmt dabei den Platz der Matrix A ein, wobei $E^{m,m}$ die Einheitsmatrix mit m Spalten und Zeilen bezeichnet.

Wir haben bereits gesehen, dass eine Nebenbedingung der Form $a_i \cdot x = b_i$ auch durch die beiden Ungleichungen $a_i \cdot x \leq b_i$ und $a_i \cdot x \geq b_i$ ausgedrückt werden kann. Das bedeutet, dass auch die Nebenbedingungen in einem linearen Programm in Gleichungsform ein Lösungspolyeder bestimmen. Alle Begriffe, die wir in den vorangegangenen Abschnitten für Lösungspolyeder

definiert haben (stützende Hyperebene, Fläche, Ecke usw.), können wir also weiterhin benutzen.

Die Schwierigkeit beim Lösen von linearen Programmen besteht darin, dass der Lösungspolyeder im Allgemeinen eine unendliche Menge ist. Wir können also nicht einfach alle zulässigen Variablenbelegungen testen und die beste wählen. Auch das Simplex-Verfahren trifft eine Auswahl der Lösungen durch eine Iteration der Basislösungen. Nach der Diskussion des grafischen Lösungsverfahrens mag der Leser aber bereits die Intuition gewonnen haben, dass es stets eine Ecke des Lösungspolyeders gibt, die eine optimale Variablenbelegung beschreibt, sofern das lineare Programm beschränkt ist. Wir zeigen als erstes, dass dies tatsächlich so ist und dass es nur endlich viele Ecken gibt. Daraus ergibt sich sogar ein (nicht effizienter) Algorithmus zur Lösung von beschränkten linearen Programmen.

Wir benötigen zunächst zwei Definitionen.

- Ein Punkt $x \in \mathcal{P}$ heißt *Extrempunkt*, wenn es kein $y \in \mathbb{R}^d$ mit $y \neq 0$ gibt, für das $x + y \in \mathcal{P}$ und $x - y \in \mathcal{P}$ gilt.
- Für $B \subseteq \{1, \dots, d\}$ bezeichnen wir mit A_B die Teilmatrix von A , die aus genau den Spalten von A besteht, deren Index in B enthalten ist. Analog bezeichnen wir für $x \in \mathbb{R}^d$ mit x_B den $|B|$ -dimensionalen Subvektor von x , der nur aus den Komponenten aus B besteht. Sei nun $x \in \mathbb{R}^d$ und $B(x) = \{i \in \{1, \dots, d\} \mid x_i \neq 0\}$. Anstatt $A_{B(x)}$ schreiben wir im Folgenden kurz A_x . Sind die Spalten von A_x linear unabhängig und gilt $Ax = b$, so heißt x *Basislösung*. Gilt zusätzlich $x \in \mathcal{P}$ (d. h. $x \geq 0$), so nennen wir x *zulässige Basislösung*. Für jede Basislösung x gilt $|B(x)| \leq m$, da der Rang der Matrix A durch m nach oben beschränkt ist.

Die Definition der Basislösung stimmt hier mit der Dictionary-Beschreibung überein. Aus der oben beschriebenen Matrix A' werden jeweils die m -Spalten der Basisvariablen betrachtet (zu Beginn die Einheitsmatrix $E^{m,m}$). Für diese Matrix wird dann die Basislösung berechnet.

Lemma 51 *Sei \mathcal{P} das Lösungspolyeder eines linearen Programms in Gleichungsform und sei $x \in \mathcal{P}$. Dann sind die folgenden drei Aussagen äquivalent.*

- a) x ist eine Ecke von \mathcal{P} .
- b) x ist ein Extrempunkt von \mathcal{P} .
- c) x ist eine zulässige Basislösung.

Beweis. Wir beweisen das Lemma mit einem Ringschluss.

- a) \Rightarrow b)** Sei x eine Ecke und sei $\mathcal{A} = \{z \in \mathbb{R}^d \mid u \cdot z = w\}$ für $u \in \mathbb{R}^d$ und $w \in \mathbb{R}$ eine stützende Hyperebene mit $\mathcal{A} \cap \mathcal{P} = \{x\}$ und $\mathcal{P} \subseteq \{z \in \mathbb{R}^d \mid u \cdot z \leq w\}$. Angenommen es gäbe ein $y \neq 0$ mit $x + y \in \mathcal{P}$ und $x - y \in \mathcal{P}$. Wegen $x + y \notin \mathcal{A} \cap \mathcal{P}$ muss $u \cdot (x + y) < w$ und damit $u \cdot y < 0$ gelten. Dann gilt aber $u \cdot (x - y) > w$ im Widerspruch zur Annahme $x - y \in \mathcal{P}$.
- b) \Rightarrow c)** Angenommen die Spalten in A_x sind nicht linear unabhängig. Wir zeigen, dass x dann kein Extrempunkt sein kann. Da die Spalten von A_x nicht linear unabhängig sind, besitzt das Gleichungssystem $A_x y' = 0$ mit $y' \in \mathbb{R}^{|B(x)|}$ eine Lösung $y' \neq 0$. Wir erweitern diese zu einer Lösung $y \in \mathbb{R}^d$ des Gleichungssystem $Ay = 0$, indem wir $y_i = 0$ für alle $i \notin B(x)$ setzen. Für jedes $\lambda \in \mathbb{R}$ gilt dann $A(x + \lambda y) = b$ und $A(x - \lambda y) = b$. Außerdem gilt für hinreichend kleine $\lambda > 0$ auch $x + \lambda y \geq 0$ und $x - \lambda y \geq 0$, da $y_i = 0$ für alle $i \notin B(x)$ und $x_i > 0$ für alle $i \in B(x)$ gilt. Damit ist gezeigt, dass x kein Extrempunkt ist.

c) \Rightarrow a) Sei $x \in \mathcal{P}$ eine Basislösung und seien $B = B(x)$ und $N = \{1, \dots, d\} \setminus B$. Dann zeigt der folgende Gedankengang, dass $\mathcal{A} = \{z \in \mathbb{R}^d \mid \sum_{i \in N} z_i = 0\}$ eine stützende Hyperebene mit $\mathcal{P} \cap \mathcal{A} = \{x\}$ ist. Zunächst folgt direkt aus der Definition, dass $x \in \mathcal{A}$ gilt, und wegen $\mathcal{P} \subseteq \mathbb{R}_{\geq 0}^d$ folgt auch $\mathcal{P} \subseteq \{z \in \mathbb{R}^d \mid \sum_{i \in N} z_i \geq 0\}$. Es bleibt also lediglich zu zeigen, dass es kein $y \in \mathcal{P}$ mit $y \neq x$ und $\sum_{i \in N} y_i = 0$ gibt. Nehmen wir an, es gäbe ein solches y . Es gilt dann $y_i = 0$ für alle $i \in N$ und somit besitzt das Gleichungssystem $A_B z = b$ zwei verschiedene Lösungen, nämlich x_B und y_B . Dies steht aber im Widerspruch dazu, dass die Spalten von A_B linear unabhängig sind. □

Mithilfe der gerade gezeigten Charakterisierung von Ecken können wir zeigen, dass der optimale Wert eines linearen Programms stets an einer Ecke angenommen wird.

Theorem 52 *Zu jedem zulässigen und beschränkten linearen Programm mit Lösungspolyeder \mathcal{P} gibt es eine Ecke von \mathcal{P} , die einer optimalen Variablenbelegung entspricht.*

Beweis. Sei $x \in \mathcal{P}$ eine optimale Variablenbelegung. Ist x eine Ecke, so ist nichts zu tun. Ansonsten gibt es laut Lemma 51 ein $y \in \mathbb{R}^d$ mit $y \neq 0$, für das $x + y \in \mathcal{P}$ und $x - y \in \mathcal{P}$ gelten. Das bedeutet, dass für dieses y insbesondere $A(x + y) = b$ und $A(x - y) = b$ gelten, woraus direkt $Ay = 0$ folgt. Damit gilt $A(x + \lambda y) = Ax + \lambda Ay = b$ für jedes $\lambda \geq 0$.

Gilt $c \cdot y < 0$ oder $c \cdot y > 0$, so ist im Widerspruch zur Optimalität von x entweder $c \cdot (x + y) < c \cdot x$ oder $c \cdot (x - y) < c \cdot x$. Somit muss $c \cdot y = 0$ gelten. Dann gilt $c \cdot (x + \lambda y) = c \cdot x$ für jedes $\lambda \geq 0$. Wir können ohne Beschränkung der Allgemeinheit davon ausgehen, dass es einen Index j mit $y_j < 0$ gibt, da wir ansonsten y durch $-y$ ersetzen könnten.

Wir werden eine Lösung $z \in \mathcal{P}$ konstruieren, für die $c \cdot z = c \cdot x$ gilt und in der mindestens eine Variable mehr den Wert Null annimmt als in der Lösung x . Wir können dann x durch z ersetzen. Dies können wir höchstens d mal wiederholen, bis wir bei einer Ecke angekommen sind.

Wir wählen $\lambda \geq 0$ größtmöglich unter der Einschränkung $z := x + \lambda y \geq 0$. Dies entspricht der Wahl

$$\lambda := \min \left\{ -\frac{x_j}{y_j} \mid j \in \{1, \dots, d\}, y_j < 0 \right\}.$$

Wir bezeichnen mit k einen Index, für den dieses Minimum angenommen wird, d. h. $y_k < 0$ und $-x_k/y_k = \lambda$. Wegen $x \geq 0$ folgt $\lambda \geq 0$. Es gilt sogar $\lambda > 0$, denn ist $y_j < 0$, so muss $x_j > 0$ sein, denn ansonsten wäre $x_j + y_j < 0$ und damit $x + y \notin \mathcal{P}$.

Wegen $z \geq 0$ und $Az = Ax + \lambda Ay = b$ folgt insgesamt, dass z eine gültige Lösung für das gegebene lineare Programm ist. Ferner gilt $c \cdot z = c \cdot x$ wegen $c \cdot y = 0$. Die Lösung z hat im Gegensatz zur Lösung x an der Stelle k eine Null. Wir haben oben bereits argumentiert, dass $x_j = 0$ impliziert, dass auch $y_j = 0$ gilt, da ansonsten $x + y \geq 0$ oder $x - y \geq 0$ nicht gilt. Somit hat z an jeder Stelle, an der x eine Null hat, ebenfalls eine. Insgesamt ist die Anzahl an Nullen in z also um mindestens eins größer als in x . □

Wir gehen im Folgenden ohne Beschränkung der Allgemeinheit davon aus, dass die Matrix $A \in \mathbb{R}^{m \times d}$ Rang m hat, d. h. dass ihre Zeilen linear unabhängig sind. Ist das nicht der Fall, so enthält das lineare Programm entweder eine redundante Nebenbedingung, die entfernt werden kann, oder es ist nicht zulässig.

Theorem 52 besagt zusammen mit Lemma 51, dass wir zur Lösung eines beschränkten linearen Programms nur alle Basislösungen testen müssen. Sei $x \in \mathbb{R}^d$ eine Basislösung, und sei $B \supseteq B(x)$ so gewählt, dass $|B| = m$ gilt und die Spalten von A_B linear unabhängig sind. Die Existenz eines solchen B folgt daraus, dass Matrix A Rang m hat. Es gilt dann $A_B x_B = b$,

und wegen der linearen Unabhängigkeit der Spalten ist x_B sogar die eindeutige Lösung dieses Gleichungssystems. Diese Interpretation ist vollkommen identisch zu der Beschreibung der Basislösungen aus den Dictionaries.

Um einfach alle Basislösungen zu direkt testen, genügt es also alle m -elementigen Teilmengen $B \subseteq \{1, \dots, d\}$ darauf zu testen, ob die Spalten von A_B linear unabhängig sind und ob das Gleichungssystem $A_B x_B = b$ eine Lösung $x_B \in \mathbb{R}_{\geq 0}^m$ besitzt.

Teste-alle-Basislösungen

1. $\min := \infty$; $x^* := \perp$;
2. **for all** $B \subseteq \{1, \dots, d\}$ mit $|B| = m$ **do**
3. **if** Spalten von A_B linear unabhängig **then**
4. Löse das Gleichungssystem $A_B x_B = b$ mit $x_B \in \mathbb{R}^m$.
5. Erweitere x_B zu $x \in \mathbb{R}^d$ durch Ergänzung von Nullen für alle $i \notin B$.
6. **if** $(x \geq 0)$ **and** $(c \cdot x < \min)$ **then**
7. $\min := c \cdot x$; $x^* := x$;
8. Ausgabe: x^*

Dieser Algorithmus ist natürlich nicht effizient, da er alle $\binom{d}{m}$ vielen m -elementigen Teilmengen von $\{1, \dots, d\}$ testet. Außerdem liefert er nur bei beschränkten linearen Programmen das richtige Ergebnis.

9.4.1 Anschauliches Beispiel der Basislösung

Zum Schluss beantworten wir noch die Frage, wie man sich eine Basislösung anschaulich in einem linearen Programm in kanonischer Form vorstellen kann und wiederholen nochmal kurz die Einführung der Slackvariablen. Sei ein lineares Programm in kanonischer Form mit d Variablen und m Nebenbedingungen gegeben, das wir in Gleichungsform transformieren, indem wir m Schlupfvariablen einfügen. Das entsprechende lineare Programm in Gleichungsform besteht dann aus $d + m$ Variablen und m Nebenbedingungen. In jeder Basislösung $x \in \mathbb{R}^{d+m}$ dieses linearen Programms sind höchstens m Variablen ungleich Null und damit mindestens $(d+m) - m = d$ Variablen gleich Null. Jede dieser Variablen ist entweder eine Variable aus dem ursprünglichen linearen Programm in kanonischer Form oder eine neu eingeführte Schlupfvariable. Ist eine Schlupfvariable gleich Null, so bedeutet das, dass die zugehörige Ungleichung mit Gleichheit erfüllt ist. Ist eine Variable aus dem ursprünglichen linearen Programm gleich Null, so bedeutet das, dass die Nichtnegativitätsbedingung mit Gleichheit erfüllt ist. Auf jeden Fall sind d Nebenbedingungen mit Gleichheit erfüllt und die Basislösung entspricht dem Schnittpunkt der zugehörigen Hyperebenen.

Beispiel

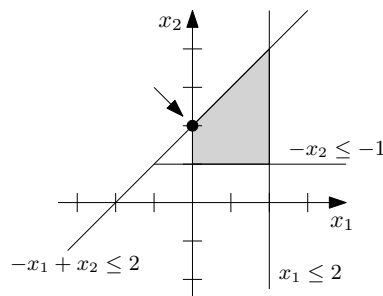
Wie betrachten die folgenden Nebenbedingungen in kanonischer Form:

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \geq 0 \quad \text{und} \quad \begin{pmatrix} 1 & 0 \\ 0 & -1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \leq \begin{pmatrix} 2 \\ -1 \\ 2 \end{pmatrix}$$

In Gleichungsform können wir diese Nebenbedingungen wie folgt schreiben:

$$\begin{pmatrix} x_1 \\ x_2 \\ s_1 \\ s_2 \\ s_3 \end{pmatrix} \geq 0 \quad \text{und} \quad \begin{pmatrix} 1 & 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 1 & 0 \\ -1 & 1 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ s_1 \\ s_2 \\ s_3 \end{pmatrix} = \begin{pmatrix} 2 \\ -1 \\ 2 \end{pmatrix}$$

Das Lösungspolyeder \mathcal{P} ist in der Abbildung grau hervorgehoben. Die markierte Ecke $(0, 2)$ entspricht dem Schnittpunkt der Hyperebenen $x_1 = 0$ und $-x_1 + x_2 = 2$. In dem linearen Programm in Gleichungsform entspricht diese Ecke der Basislösung $(0, 2, 2, 1, 0)$, wobei die letzten drei Komponenten s_1 , s_2 und s_3 entsprechen.



Für $x = (0, 2, 2, 1, 0)$ gilt $B = B(x) = \{2, 3, 4\}$ und dementsprechend ist x_B die eindeutige Lösung des folgenden Gleichungssystems:

$$\begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_2 \\ s_1 \\ s_2 \end{pmatrix} = \begin{pmatrix} 2 \\ -1 \\ 2 \end{pmatrix}$$

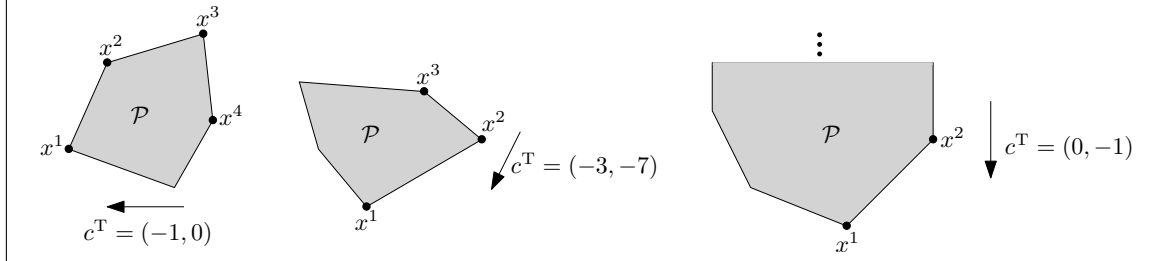
9.5 Simplex-Algorithmus in der Matrix-Variante

Der Simplex-Algorithmus, den wir in diesem Abschnitt beschreiben werden, wurde bereits 1947 von George Dantzig vorgeschlagen. Es handelt sich dabei auch heute noch um einen der erfolgreichsten Algorithmen zur Lösung von linearen Programmen in der Praxis. Wir beschreiben den Algorithmus zunächst wieder anschaulich, bevor wir eine formale Beschreibung angeben.

Sei ein lineares Programm mit Lösungspolyeder \mathcal{P} gegeben. Der Simplex-Algorithmus führt eine lokale Suche auf den Ecken des Lösungspolyeders durch. Das bedeutet, er startet mit einer beliebigen Ecke $x^1 \in \mathcal{P}$ und testet, ob es eine benachbarte Ecke mit besserem Zielfunktionswert gibt. Dabei sind zwei Ecken benachbart, wenn sie auf der gleichen Kante liegen. Gibt es eine solche benachbarte Ecke $x^2 \in \mathcal{P}$, so macht der Algorithmus mit x^2 analog weiter und testet wieder, ob es eine bessere benachbarte Ecke gibt. Den Übergang von einer Ecke zu einer besseren in der Nachbarschaft nennen wir auch *Pivotschritt*. Da der Algorithmus den Zielfunktionswert mit jedem Schritt verbessert, kann er keine Ecke mehrfach besuchen. Da es nur endliche viele Ecken gibt, muss also nach einer endlichen Anzahl an Schritten eine Ecke $x^k \in \mathcal{P}$ erreicht sein, in deren Nachbarschaft es keine bessere Ecke gibt. Es gibt dann zwei Fälle, die eintreten können. Entweder die Ecke x^k ist in einer unbeschränkten Kante enthalten, entlang derer der Zielfunktionswert beliebig gut werden kann, oder, wenn das nicht der Fall ist, dann ist die Ecke x^k optimal. Wir werden später beweisen, dass dies wirklich so ist.

Beispiel

Bei den ersten beiden Beispielen wird das Optimum an den Ecken x^4 bzw. x^3 angenommen. Im dritten Beispiel ist x^2 die letzte besuchte Ecke, von der aus es eine unendliche Kante gibt, entlang derer der Wert der Zielfunktion beliebig klein wird.



Genau wie bei dem geometrischen Algorithmus zur Lösung von linearen Programmen weisen wir auch hier darauf hin, dass sich aus der obigen anschaulichen Beschreibung des Simplex-Algorithmus noch keine direkte Implementierung ergibt. So ist zum Beispiel nicht klar, wie im ersten Schritt überhaupt eine beliebige Ecke von \mathcal{P} gefunden werden kann oder wie man die Nachbarschaft einer Ecke durchsuchen kann.

9.5.1 Formale Beschreibung

Zur formalen Beschreibung betrachten wir wieder ein lineares Programm in Gleichungsform, in dem die Zielfunktion $c \cdot x$ unter den Nebenbedingungen $Ax = b$ und $x \geq 0$ minimiert werden soll. Dabei seien $c \in \mathbb{R}^d$, $b \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times d}$ und $x \in \mathbb{R}^d$ sei ein Vektor, der aus d Variablen besteht. Es bezeichne \mathcal{P} das zugehörige Lösungspolyeder. Wie bereits im letzten Abschnitt gehen wir davon aus, dass alle Zeilen der Matrix A linear unabhängig sind.

Wir sagen, dass das lineare Programm *degeneriert* ist, wenn es eine zulässige Basislösung $x \in \mathcal{P}$ mit $|B(x)| < m$ gibt. Ist ein lineares Programm nicht degeneriert, so sind in jeder zulässigen Basislösung genau m Variablen ungleich Null. Wir werden in diesem Abschnitt davon ausgehen, dass das vorliegende lineare Programm nicht degeneriert ist. Diese Annahme erleichtert zwar die Beschreibung, sie ist aber nicht notwendig, denn an sich können mit dem Simplex-Algorithmus auch degenerierte lineare Programme gelöst werden. Wir haben gesehen, dass Basislösungen in linearen Programmen in kanonischer Form Punkten entsprechen, in denen sich mindestens d der durch die Nebenbedingungen beschriebenen Hyperebenen schneiden. Ein solches lineares Programm ist genau dann degeneriert, wenn es einen Punkt in \mathcal{P} gibt, in dem sich mehr als d Hyperebenen schneiden.

Wir gehen zunächst davon aus, dass wir eine beliebige Ecke $x = x^1 \in \mathcal{P}$ kennen. Diese Ecke entspricht gemäß Lemma 51 einer zulässigen Basislösung. Für $B = B(x)$ und $N = \{1, \dots, d\} \setminus B$ ergibt sich diese Basislösung als eindeutige Lösung des Gleichungssystems $A_B x_B = b$ und $x_N = 0$. Aufgrund der Annahme, dass das lineare Programm nicht degeneriert ist, besteht die Matrix A_B aus genau m linear unabhängigen Spalten. Somit ist A_B einer invertierbare $m \times m$ -Matrix. Es gilt also insgesamt

$$x = \begin{pmatrix} x_B \\ x_N \end{pmatrix} = \begin{pmatrix} A_B^{-1} b \\ 0 \end{pmatrix}.$$

Wir nennen die Variablen in x_B *Basisvariablen* und die Variablen in x_N *Nichtbasisvariablen*. Um zu einer benachbarten Ecke zu gelangen, wollen wir jetzt den Wert einer Nichtbasisvariablen erhöhen. Dazu ist es zunächst hilfreich, die Werte der Basisvariablen in Abhängigkeit von den Werten der Nichtbasisvariablen auszudrücken. Dazu bezeichne A_N die Teilmatrix von A ,

die aus genau den Spalten besteht, die nicht in A_B enthalten sind. Es gilt dann

$$A_B x_B + A_N x_N = b \iff x_B = A_B^{-1}(b - A_N x_N).$$

Wählen wir $x_N \in \mathbb{R}_{\geq 0}^{d-m}$ beliebig, so ergibt sich daraus eine eindeutige Wahl von x_B . Wir können deshalb die Zielfunktion in Abhängigkeit von x_N ausdrücken:

$$\begin{aligned} c \cdot x &= c_B \cdot x_B + c_N \cdot x_N = c_B \cdot (A_B^{-1}(b - A_N x_N)) + c_N \cdot x_N \\ &= c_B \cdot (A_B^{-1}b) + (c_N^T - c_B^T A_B^{-1} A_N) x_N. \end{aligned} \quad (9.44)$$

Eine Wahl $x_N \in \mathbb{R}_{\geq 0}^{d-m}$ ist nur dann zulässig, wenn sich daraus nichtnegative Werte für die Basisvariablen ergeben.

Wir nennen zwei Basislösungen *benachbart*, wenn sich die Mengen der Basisvariablen nur in einem Element unterscheiden. Wir können von einer Basislösung x zu einer benachbarten Basislösung gelangen, indem wir eine Nichtbasisvariable $j \in N$ auswählen und den Wert von x_j erhöhen und die Werte aller anderen Nichtbasisvariablen auf Null lassen. Die zu j gehörende Komponente des Vektors $\bar{c} := c_N^T - c_B^T A_B^{-1} A_N$ aus Gleichung (9.44) beschreibt dann, wie sich der Wert der Zielfunktion ändert. Wir nennen \bar{c} auch den Vektor der *reduzierten Kosten*. Nur wenn die zu j gehörende Komponente des Vektors \bar{c} negativ ist, verbessert sich die Lösung durch die Erhöhung von x_j . Das folgende Lemma besagt, dass wir anhand der reduzierten Kosten erkennen können, ob x eine optimale Lösung ist.

Lemma 53 *Falls der Vektor der reduzierten Kosten \bar{c} zu einer Basislösung x keinen negativen Eintrag enthält, so ist x eine optimale Lösung für das lineare Programm.*

Beweis. Es sei $\bar{c} := c_N^T - c_B^T A_B^{-1} A_N \geq 0$. Sei $x' \in \mathcal{P}$ eine beliebige zulässige Lösung. Dann ist $x'_N \geq 0$ und Gleichung (9.44) impliziert

$$c \cdot x' = c_B \cdot (A_B^{-1}b) + \bar{c} x'_N \geq c_B \cdot (A_B^{-1}b) = c \cdot x.$$

Also ist x' nicht besser als x , und somit muss x eine optimale Lösung sein. \square

Gilt $\bar{c} \geq 0$ für die aktuelle Basislösung x , so terminiert der Simplex-Algorithmus und gibt aus, dass x optimal ist. Ansonsten betrachten wir eine beliebige Nichtbasisvariable x_j mit $\bar{c}_j < 0$ und erhöhen diese kontinuierlich von 0. Wenn wir die Basisvariablen gemäß $x_B = A_B^{-1}(b - A_N x_N)$ entsprechend anpassen, dann verringert sich durch die Erhöhung von x_j gemäß Gleichung (9.44) der Zielfunktionswert. Es kann jedoch sein, dass die Erhöhung von x_j dazu führt, dass sich der Wert einiger Basisvariablen verringert und sogar negativ wird. Dann erhalten wir keine zulässige Lösung mehr. Da wir davon ausgehen, dass das lineare Programm nicht degeneriert ist, haben in x aber alle Basisvariablen echt positive Werte. Das bedeutet, wir können x_j zumindest von Null auf einen eventuell kleinen, aber zumindest echt positiven Wert setzen, ohne dadurch eine unzulässige Lösung zu erhalten.

Die zu x_j gehörende Spalte der Matrix $A_B^{-1} A_N$ nennen wir im Folgenden $u \in \mathbb{R}^m$. Ist keine Komponente von u echt positiv, so können wir x_j beliebig groß setzen, ohne dadurch eine unzulässige Lösung zu erhalten. Das bedeutet, dass das lineare Programm unbeschränkt ist, da wir beliebig gute zulässige Lösungen finden können. Ansonsten setzen wir x_j größtmöglich unter der Bedingung, dass alle Basisvariablen nichtnegativ bleiben. Das bedeutet

$$x_j := \min \left\{ \frac{x_\ell}{u_\ell} \mid \ell \in B, u_\ell > 0 \right\}.$$

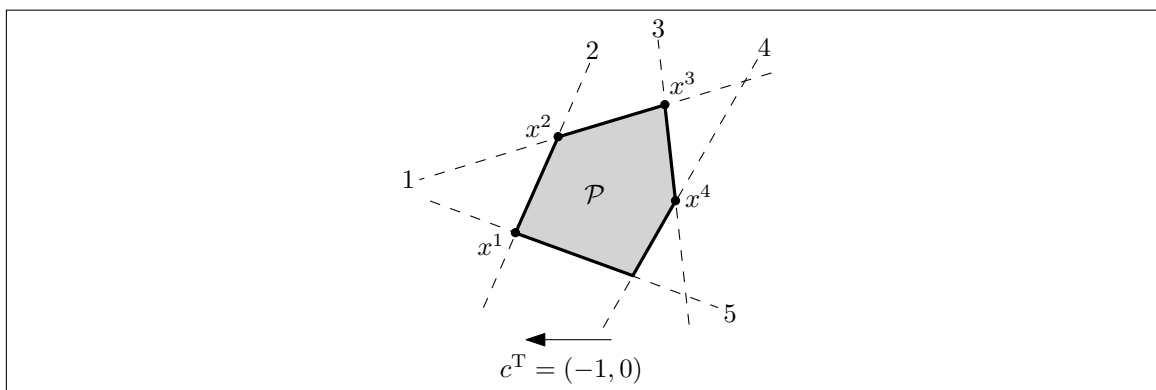
Sei x_k eine Basisvariable, für die das obige Minimum angenommen wird. Nachdem wir x_j erhöht haben, nimmt x_k den Wert Null an. Auf diese Weise erhalten wir eine neue Basislösung x'

mit $B(x') = B(x) \cup \{j\} \setminus \{k\}$. Damit ist x' eine zu x benachbarte Basislösung mit besserem Zielfunktionswert. Hier haben wir ausgenutzt, dass das lineare Programm nicht degeneriert ist und dass deshalb $|B(x')| = m$ gilt.

Wir haben uns bereits im vorigen Abschnitt überlegt, dass Basislösungen in linearen Programmen in kanonischer Form Schnittpunkten von d durch die Nebenbedingungen beschriebenen Hyperebenen entsprechen. Diese Hyperebenen sind durch die Nichtbasisvariablen bestimmt. Nun haben wir gerade definiert, dass zwei Basislösungen benachbart sind, wenn genau eine Nichtbasisvariable durch eine andere ausgetauscht wird. Für lineare Programme in kanonischer Form können wir also sagen, dass zwei Basislösungen benachbart sind, wenn sie jeweils ein Schnittpunkt von d Hyperebenen sind und wenn sie auf $d - 1$ dieser Hyperebenen gemeinsam liegen.

Beispiel

In dem folgenden Beispiel ist eine Ausführung des Simplex-Algorithmus zu sehen, in der die Folge x^1, x^2, x^3, x^4 von Basislösungen durchlaufen wird. Jede dieser Basislösungen können wir als Schnittpunkt zweier Hyperebenen darstellen. Die Basislösung x^1 entspricht zum Beispiel der Menge $\{2, 5\}$ von Hyperebenen. Insgesamt ergibt sich die Folge $\{2, 5\}, \{1, 2\}, \{1, 3\}, \{3, 4\}$. Dort ist zu erkennen, dass jeweils eine Hyperebene ausgetauscht wird.



9.5.2 Berechnung der initialen Basislösung

Wir haben noch nicht diskutiert, wie eine initiale Basislösung gefunden werden kann, mit der der Simplex-Algorithmus gestartet werden kann. Sei ein lineares Programm mit den Nebenbedingungen $Ax = b$ gegeben und sei ohne Beschränkung der Allgemeinheit $b \geq 0$. Ist das nicht der Fall, so können die entsprechenden Nebenbedingungen mit -1 multipliziert werden. Für die m Nebenbedingungen führen wir Hilfsvariablen $h_1 \geq 0, \dots, h_m \geq 0$ ein. Die Nebenbedingung $a_i \cdot x = b_i$ ersetzen wir für jedes $i \in \{1, \dots, m\}$ durch die Nebenbedingung $a_i \cdot x + h_i = b_i$. Wir ignorieren außerdem die gegebene Zielfunktion des linearen Programms und definieren als neue Zielfunktion $h_1 + \dots + h_m$.

Das so entstandene lineare Programm hat die Eigenschaft, dass wir einfach eine zulässige Basislösung angeben können. Dazu setzen wir $h_i = b_i$ für jedes $i \in \{1, \dots, m\}$ und $x_i = 0$ für alle $i \in \{1, \dots, d\}$. Wir können somit den Simplex-Algorithmus benutzen, um für dieses lineare Programm eine optimale Basislösung (x^*, h^*) zu berechnen. Dieser Schritt wird als *erste Phase* des Simplex-Algorithmus bezeichnet. Ist der Wert der Zielfunktion in dieser optimalen Lösung echt größer als Null, dann gibt es keine Lösung, die alle Gleichungen des ursprünglichen linearen Programms erfüllt, und somit ist dieses lineare Programm nicht zulässig. Ansonsten, wenn $h^* = 0$ gilt, so ist x^* eine zulässige Basislösung für das ursprüngliche lineare

Programm. Mit dieser Lösung x^* kann dann der Simplex-Algorithmus für das ursprüngliche lineare Programm initialisiert werden. Die Optimierung der eigentlichen Zielfunktion wird dann als *zweite Phase* des Simplex-Algorithmus bezeichnet.

9.5.3 Laufzeit

Die Laufzeit des Simplex-Algorithmus setzt sich aus zwei Komponenten zusammen. Zum einen muss geklärt werden, was die Komplexität eines einzelnen Pivotschrittes ist und zum anderen muss geklärt werden, wie viele Pivotschritte es gibt. Die Antwort auf die erste Frage gibt folgendes Theorem, das wir nicht beweisen werden.

Theorem 54 *Die Laufzeit eines einzelnen Pivotschrittes ist polynomiell in der Eingabelänge des linearen Programms beschränkt.*

Aus der Beschreibung des Simplex-Algorithmus geht bereits hervor, dass zur Durchführung eines Pivotschrittes nur einige elementare Rechenoperationen durchgeführt werden müssen. Trotzdem ist das obige Theorem nicht trivial, denn die Laufzeit dieser Rechenoperationen hängt von der Darstellungslänge der vorkommenden Zahlen ab. Diese muss man in Beziehung zu der Eingabelänge setzen. Man muss also zeigen, dass die Zahlen, mit denen man beim Simplex-Algorithmus rechnet, nicht zu groß werden können.

Die Antwort auf die zweite Frage ist hingegen ernüchternd. Klee und Minty haben 1972 folgendes Theorem gezeigt.

Theorem 55 *Für jedes $n \in \mathbb{N}$ gibt es ein lineares Programm in Gleichungsform mit $3n$ Variablen und $2n$ Nebenbedingungen, in dem alle Koeffizienten ganzzahlig sind und Absolutwert höchstens 4 haben, und auf dem der Simplex-Algorithmus $2^n - 1$ Pivotschritte durchführen kann.*

In diesem Theorem haben wir davon gesprochen, dass der Simplex-Algorithmus $2^n - 1$ Pivotschritte „durchführen kann“. Diese Formulierung ist so gewählt, da es für eine Basislösung oft mehrere mögliche Pivotschritte gibt und wir bisher nicht spezifiziert haben, welcher davon durchgeführt wird. Das Theorem ist so zu lesen, dass eine exponentiell lange Folge von Pivotschritten existiert, die der Simplex-Algorithmus durchführen kann. Es ist aber zunächst nicht ausgeschlossen, dass man durch eine geschickte Wahl der Pivotschritte garantieren kann, dass der Simplex-Algorithmus stets polynomiell viele Schritte benötigt. Eine Regel, die die Wahl des nächsten Pivotschrittes bestimmt, wenn mehrere zur Auswahl stehen, heißt *Pivotregel*. Zahlreiche solche Regeln wurden vorgeschlagen, für fast alle von ihnen kennt man aber mittlerweile Instanzen, auf denen der Simplex-Algorithmus exponentiell viele Pivotschritte benötigt. Ob es eine Pivotregel gibt, die garantiert, dass der Simplex-Algorithmus nach polynomiell vielen Schritten terminiert, ist bis heute unklar.

9.6 Komplexität von linearer Programmierung

Lineare Programmierung ist ein sehr gut untersuchtes Thema der Optimierung, und für lange Zeit war es eine große offene Frage, ob es Algorithmen gibt, die lineare Programme in polynomieller Zeit lösen. Diese Frage wurde 1979 von Leonid Khachiyan positiv beantwortet. Er stellte die *Ellipsoidmethode* vor, die in polynomieller Zeit testet, ob ein gegebenes Lösungspolyeder \mathcal{P} leer ist oder nicht. Die Ellipsoidmethode erzeugt zunächst ein Ellipsoid (die mehrdimensionale Entsprechung einer Ellipse), das \mathcal{P} komplett enthält. In jeder Iteration

wird dann getestet, ob das Zentrum des Ellipsoides in \mathcal{P} liegt. Falls ja, so ist \mathcal{P} nicht leer. Ansonsten wird ein neues Ellipsoid berechnet, das \mathcal{P} immer noch enthält, dessen Volumen, aber um einen gewissen Faktor kleiner ist als das Volumen des alten Ellipsoides. Man kann nachweisen, dass nach polynomiell vielen Iterationen entweder ein Punkt in \mathcal{P} gefunden wird oder das aktuelle Ellipsoid so klein geworden ist, dass $\mathcal{P} = \emptyset$ gelten muss.

Eine beliebige Variablenbelegung $x \in \mathcal{P}$ zu berechnen oder zu entscheiden, dass keine solche existiert, ist ein scheinbar einfacheres Problem, als die optimale Lösung eines linearen Programms zu bestimmen. Tatsächlich gilt aber, dass sich ein polynomieller Algorithmus \mathcal{A} für das erste Problem in einen polynomiellen Algorithmus für das zweite Problem transformieren lässt. Aus einem linearen Programm, in dem die Zielfunktion $c \cdot x$ unter den Nebenbedingungen $Ax \leq b$ minimiert werden soll, erzeugen wir ein neues lineares Programm mit den Nebenbedingungen $Ax \leq b$ und zusätzlich $c \cdot x \leq z$ für ein $z \in \mathbb{R}$. Ist das zu diesem linearen Programm gehörende Lösungspolyeder nicht leer, so besitzt das ursprüngliche lineare Programm eine Lösung mit Wert höchstens z . Wir suchen jetzt mit Hilfe von Algorithmus \mathcal{A} und einer binären Suche das kleinste z für das das oben genannte lineare Programm eine Lösung besitzt. Für dieses lineare Programm bestimmen wir dann mit Hilfe von \mathcal{A} eine zulässige Lösung. Diese Lösung entspricht dann einer optimalen Lösung des ursprünglichen linearen Programms.

Obwohl die Ellipsoidmethode ein großer theoretischer Durchbruch war, ist ihr Einfluss auf die Praxis gering, denn in praktischen Anwendungen ist der Simplex-Algorithmus trotz seiner exponentiellen Worst-Case-Laufzeit der polynomiellen Ellipsoidmethode deutlich überlegen. Im Jahre 1984 gelang es Narendra Karmarkar mit dem *Innere-Punkte-Verfahren* einen polynomiellen Algorithmus zur Lösung von linearen Programmen vorzustellen, der auch in einigen praktischen Anwendungen ähnlich gute Laufzeiten wie der Simplex-Algorithmus erzielt. Dennoch ist der Simplex-Algorithmus auch heute noch der am meisten benutzte Algorithmus zur Lösung von linearen Programmen.

Literaturverzeichnis

- [1] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, Clifford Stein. Introduction to Algorithms. MIT Press, 3. Auflage, 2009.
- [2] Norbert Blum. Algorithmen und Datenstrukturen: Eine anwendungsorientierte Einführung. Oldenbourg Verlag, 2004.