

Willkommen!

Vorlesungshomepage:

http://lec.inf.ethz.ch/DA/2021

Das Team:

Assistenten Sebastian Balzer

Christoph Grötzbach

Ivana Klasovita

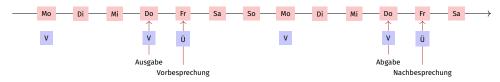
Stanislav Piasecki

Back-Office Aritra Dhar

Julia Chatain

Dozent Felix Friedrich

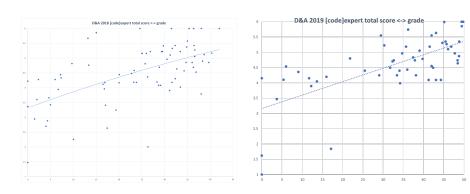
Übungsbetrieb



- Übungsblattausgabe zur Vorlesung (online).
- Vorbesprechung in der folgenden Übung.
- Bearbeitung der Übung bis spätestens am Tag vor der nächsten Übungsstunde (23:59h).
- Nachbesprechung der Übung in der nächsten Übungsstunde. Feedback zu den Abgaben innerhalb einer Woche nach Nachbesprechung.

Zu den Übungen

■ Bearbeitung der wöchentlichen Uebungsserien ist freiwillig, wird aber dringend empfohlen!



Es ist so einfach!

Für die Übungen verwenden wir eine Online-Entwicklungsumgebung, benötigt lediglich einen Browser, Internetverbindung und Ihr ETH Login.

Falls Sie keinen Zugang zu einem Computer haben: in der ETH stehen an vielen Orten öffentlich Computer bereit.

Relevantes für die Prüfung

Prüfungsstoff für die Endprüfung schliesst ein

- Vorlesungsinhalt (Vorlesung, Handout) und
- Übungsinhalte (Übungsstunden, Übungsblätter).

Prüfung (150 min) ist schriftlich. Hilfsmittel: vier A4-Blätter ohne inhaltliche und formale Anforderungen (Text, Bilder, ein-/doppelseitig, Ränder, Schriftgrössen) Die Prüfung findet voraussichtlich in hybrider Form (auf Papier und am Computer statt).

Literatur

Algorithmen und Datenstrukturen, *T. Ottmann, P. Widmayer*, Spektrum-Verlag, 5. Auflage, 2011

Algorithmen - Eine Einführung, T. Cormen, C. Leiserson, R. Rivest, C. Stein, Oldenbourg, 2010

Introduction to Algorithms, *T. Cormen, C. Leiserson, R. Rivest, C. Stein*, 3rd ed., MIT Press, 2009

The C++ Programming Language, *B. Stroustrup*, 4th ed., Addison-Wesley, 2013.

The Art of Multiprocessor Programming, *M. Herlihy*, *N. Shavit*, Elsevier, 2012.

Unser Angebot

- Bearbeitung der wöchentlichen Übungsserien \rightarrow Bonus von maximal 0.25 Notenpunkten für die Prüfung.
- Bonus proportional zur erreichten Punktzahl von **speziell markierten** Bonus-Aufgaben. Volle Punktzahl = 0.25.
- **Zulassung** zu speziell markierten Bonusaufgaben kann von der erfolgreichen Absolvierung anderer Übungsaufgaben abhängen.

Unser Angebot (Konkret)

- Insgesamt 3 Bonusaufgaben; 2/3 der Punkte reichen für 0.25 Bonuspunkte für die Prüfung
- Sie können also z.B. 2 Bonusaufgaben zu 100% lösen, oder 3 Bonusaufgaben zu je 66%, oder ...
- Bonusaufgaben müssen durch erfolgreich gelöste Übungsserien freigeschaltet (→ Experience Points) werden
- Es müssen wiederum nicht alle Übungsserien vollständig gelöst werden, um eine Bonusaufgabe freizuschalten
- Details: Übungsstunden, Online-Übungssystem (Code Expert)

Akademische Lauterkeit

Regel: Sie geben nur eigene Lösungen ab, welche Sie selbst verfasst und verstanden haben.

Wir prüfen das (zum Teil automatisiert) nach und behalten uns disziplinarische Massnahmen vor.

8

Wenn es Probleme gibt ...

- mit dem Kursinhalt
 - unbedingt alle Übungen besuchen
 - dort Fragen stellen
 - und/oder Übungsleiter kontaktieren
- alle weiteren Probleme
 - Email an Chefassistenten (Aritra Dhar) oder Dozenten (Felix Friedrich)
- Wir helfen gerne!

1. Einführung

Überblick, Algorithmen und Datenstrukturen, Korrektheit, erstes Beispiel

Ziele des Kurses

- Verständnis des Entwurfs und der Analyse grundlegender Algorithmen und Datenstrukturen.
- Vertiefter Einblick in ein modernes Programmiermodell (mit C++).
- Wissen um Chancen, Probleme und Grenzen des parallelen und nebenläufigen Programmierens.

12

Inhalte der Vorlesung

Datenstrukturen / Algorithmen

Begriff der Invariante, Kostenmodell, Landau Symbole Minimale Spannbäume, Fibonacci Heaps
Algorithmenentwurf, Induktion Kürzeste Wege, Minimaler Spannbaum Maximaler Fluss
Suchen und Auswahl, Sortieren Fundamentale Algorithmen auf Graphen
Amortisierte Analyse Wörterbücher: Hashing und Suchbäume, AVL
Dynamic Programming

Programmieren mit C++

RAII, Move Konstruktion, Smart Pointers, Promises and Futures
Templates und Generische Programmierung Threads, Mutexs and Monitors
Exceptions Funktoren und Lambdas

Parallel Programming

Parallelität vs. Concurrency, Speedup (Amdahl/Gustavson), Races, Memory Reordering, Atomic Registers, RMW (CAS,TAS), Deadlock/Starvation

1.2 Algorithmen

[Cormen et al, Kap. 1; Ottman/Widmayer, Kap. 1.1]

Algorithmus

Algorithmus

Wohldefinierte Berechnungsvorschrift, welche aus Eingabedaten (input) Ausgabedaten (output) berechnet.

Beispielproblem: Sortieren

Input: Eine Folge von n Zahlen (vergleichbaren Objekten) (a_1, a_2, \ldots, a_n) **Output**: Eine Permutation $(a'_1, a'_2, \ldots, a'_n)$ der Folge $(a_i)_{1 \le i \le n}$, so dass $a'_1 \le a'_2 \le \cdots \le a'_n$

Mögliche Eingaben

```
(1,7,3), (15,13,12,-0.5), (999,998,997,996,\ldots,2,1), (1), (1), (1)...
```

Jedes Beispiel erzeugt eine Probleminstanz.

Die Performanz (Geschwindigkeit) des Algorithmus hängt üblicherweise ab von der Probleminstanz. Es gibt oft "gute" und "schlechte" Instanzen.

Daher betrachten wir Algorithmen manchmal "im Durchschnitt" und meist "im schlechtesten Fall".

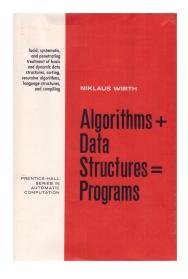
Mögliche Lösung

Wie oft werden die Zeilen jeweils ausgeführt?

```
void sort(std::vector<int> a){
   std::size n = a.size()
   for (std::size i = 0; i<n; ++i)
      for (std::size j = i+1; j<n; ++j)
        if (a[j] < a[i])
            std::swap(a[i],a[j])
}</pre>
```

Datenstrukturen

- Eine Datenstruktur organisiert Daten so in einem Computer, dass man sie (in den darauf operierenden Algorithmen) effizient nutzen kann.
- Programme = Algorithmen + Datenstrukturen.



Beispiele für Probleme in der Algorithmik

- Tabellen und Statistiken: Suchen, Auswählen und Sortieren
- Routenplanung: Kürzeste Wege Algorithmus, Heap Datenstruktur
- DNA Matching: Dynamic Programming
- Auswertungsreihenfolge: Topologische Sortierung
- Autovervollständigung: Wörterbücher/Bäume
- Schnelles Nachschlagen: Hash-Tabellen
- Der Handlungsreisende: Dynamische Programmierung, Minimal aufspannender Baum, Simulated Annealing,

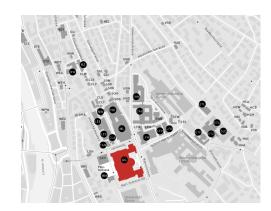
18

Charakteristik

- Extrem grosse Anzahl potentieller Lösungen
- Praktische Anwendung

Typische Design-Schritte am Beispiel

Routenplanung



21

Typische Design-Schritte

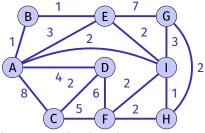
- 1. Spezifikation des Problems: Finde beste Route (kürzeste Zeit) von A nach B
- 2. Abstraktion: Graph aus Knoten, Kanten und Kantengewichten
- 3. Idee (Heureka!): Dijkstra
- 4. Datenstrukturen und Algorithmen: z.B. Min-Heap, Adjazenzmatrix / Adjazenzliste, Hash-Tabelle ...
- 5. Laufzeitanalyse: $\mathcal{O}((n+m) \cdot \log n)$
- 6. Implementation: Wahl der Repräsentation (z.B. Adjazenzmatrix/ Adjazenzliste/ Objekte



Schwieriges Problem: Der Handlungsreisende

Gegeben: Graph (Landkarte) aus Knoten (Städte) und gewichteten Kanten (Wege mit Weglänge)

Gesucht: Rundweg durch alle Städte, wobei jede Stadt genau einmal besucht wird (Hamilton-Kreis) mit kürzester Weglänge.



Der beste bekannte Algorithmus hat eine Laufzeit die mit der Anzahl der Knoten (Städte) exponentiell ansteigt.

Schon das Finden eines Hamilton-Kreises ist im allgemeinen Fall ein schwieriges Problem. Das Problem, einen Eulerkreis zu finden, nämlich einen Kreis, der jede *Kante* einmal besucht, ist hingegen ein Problem, welches in polynomieller Laufzeit gelöst werden kann.

Schwierige Probleme

- NP-vollständige Probleme: Keine bekannte effiziente Lösung (Existenz einer effizienten Lösung ist zwar sehr unwahrscheinlich es ist aber unbewiesen, dass es keine gibt!)
- Beispiel: Travelling Salesman Problem

In diesem Kurs beschäftigen wir uns hauptsächlich mit Problemen, die effizient (in Polynomialzeit) lösbar sind.

2. Effizienz von Algorithmen

Effizienz von Algorithmen, Random Access Machine Modell, Funktionenwachstum, Asymptotik [Cormen et al, Kap. 2.2,3,4.2-4.4 | Ottman/Widmayer, Kap. 1.1]

Effizienz

- Wären Rechner unendlich schnell und hätten unendlich viel Speicher ...
- ... dann bräuchten wir die Theorie der Algorithmen (nur) für Aussagen über Korrektheit (incl. Terminierung).

Realität: Ressourcen sind beschränkt und nicht umsonst:

- Rechenzeit → Effizienz
- Speicherplatz → Effizienz

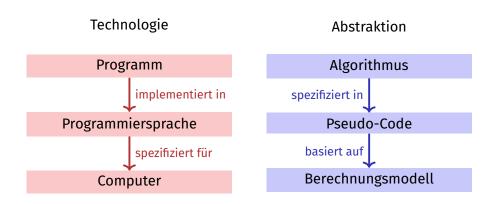
Eigentlich geht es in diesem Kurs nur um Effizienz.

Effizienz von Algorithmen

Ziele

- Laufzeitverhalten eines Algorithmus maschinenunabhängig quantifizieren.
- Effizienz von Algorithmen vergleichen.
- Abhängigkeit von der Eingabegrösse verstehen.

Programme und Algorithmen



Technologiemodell

Random Access Machine (RAM) Model

- Ausführungsmodell: Instruktionen werden der Reihe nach (auf einem Prozessorkern) ausgeführt.
- Speichermodell: Konstante Zugriffszeit (grosses Array)
- Elementare Operationen: Rechenoperation (+,-,·,...), Vergleichsoperationen, Zuweisung / Kopieroperation auf Maschinenworten (Registern), Flusskontrolle (Sprünge)
- Einheitskostenmodell: elementare Operation hat Kosten 1.
- Datentypen: Fundamentaltypen wie grössenbeschränkte Ganzzahl oder Fliesskommazahl.

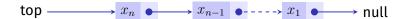
Grösse der Eingabedaten

- Typisch: Anzahl Eingabeobjekte (von fundamentalem Typ).
- Oftmals: Anzahl Bits für eine vernünftige / kostengünstige Repräsentation der Daten.
- Annahme: fundamentale Typen passen in Machinenwort (word) mit Grösse : $w \ge \log(\text{sizeof(mem)})$ Bits.

Für dynamische Datenstrukturen

Pointer Machine Modell

- Objekte beschränkter Grösse können dynamisch erzeugt werden in konstanter Zeit 1.
- Auf Felder (mit Wortgrösse) der Objekte kann in konstanter Zeit 1 zugegriffen werden.



Asymptotisches Verhalten

Genaue Laufzeit eines Algorithmus lässt sich selbst für kleine Eingabedaten kaum voraussagen.

- Betrachten das asymptotische Verhalten eines Algorithmus.
- Ignorieren alle konstanten Faktoren.

Eine Operation mit Kosten 20 ist genauso gut wie eine mit Kosten 1. Lineares Wachstum mit Steigung 5 ist genauso gut wie lineares Wachstum mit Steigung 1.

Algorithmen, Programme und Laufzeit

Programm: Konkrete Implementation eines Algorithmus.

Laufzeit des Programmes: messbarer Wert auf einer konkreten Maschine. Kann sowohl nach oben, wie auch nach unten abgeschätzt werden.

Example 1

Rechner mit 3 GHz. Maximale Anzahl Operationen pro Taktzyklus (z.B. 8). \Rightarrow untere Schranke.

Einzelne Operation dauert mit Sicherheit nie länger als ein Tag \Rightarrow obere Schranke.

Hinsichtlich des *asymptotischen Verhaltens* des Programmes spielen die Schranken keine Rolle.

33

Oberflächlich

2.2 Funktionenwachstum

 \mathcal{O} , Θ , Ω [Cormen et al, Kap. 3; Ottman/Widmayer, Kap. 1.1]

Verwende die asymptotische Notation zur Kennzeichnung der Laufzeit von Algorithmen

Wir schreiben $\Theta(n^2)$ und meinen, dass der Algorithmus sich für grosse n wie n^2 verhält: verdoppelt sich die Problemgrösse, so vervierfacht sich die Laufzeit.

Genauer: Asymptotische obere Schranke

Gegeben: Funktion $g: \mathbb{N} \to \mathbb{R}$. Definition:¹

$$\mathcal{O}(g) = \{ f : \mathbb{N} \to \mathbb{R} |$$

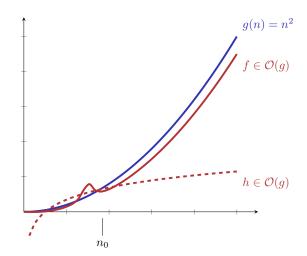
$$\exists c > 0, \exists n_0 \in \mathbb{N} :$$

$$\forall n \ge n_0 : 0 \le f(n) \le c \cdot g(n) \}$$

Schreibweise:

$$\mathcal{O}(g(n)) := \mathcal{O}(g(\cdot)) = \mathcal{O}(g).$$

Anschauung



Umkehrung: Asymptotische untere Schranke

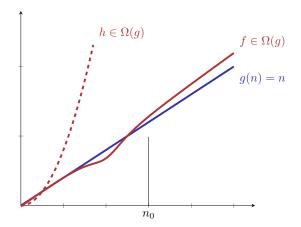
Gegeben: Funktion $g: \mathbb{N} \to \mathbb{R}$. Definition:

$$\Omega(g) = \{ f : \mathbb{N} \to \mathbb{R} |$$

$$\exists c > 0, \exists n_0 \in \mathbb{N} :$$

$$\forall n \ge n_0 : 0 \le c \cdot g(n) \le f(n) \}$$

Beispiel



¹Ausgesprochen: Menge aller reellwertiger Funktionen $f: \mathbb{N} \to \mathbb{R}$ für die gilt: es gibt ein (reellwertiges) c > 0 und ein $n_0 \in \mathbb{N}$ so dass $0 \le f(n) \le n \cdot g(n)$ für alle $n \ge n_0$.

Asymptotisch scharfe Schranke

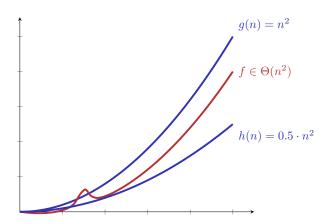
Gegeben Funktion $g:\mathbb{N}
ightarrow \mathbb{R}$.

Definition:

$$\Theta(g) := \Omega(g) \cap \mathcal{O}(g).$$

Einfache, geschlossene Form: Übung.

Beispiel

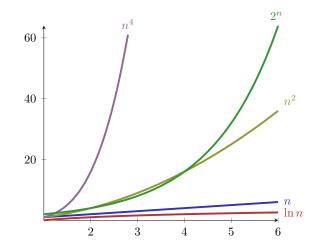


41

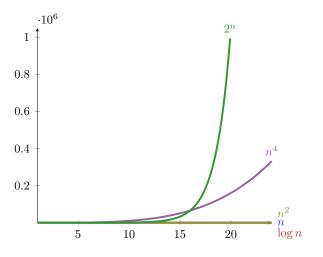
Wachstumsbezeichnungen

$\mathcal{O}(1)$	beschränkt	Array-Zugriff
$\mathcal{O}(\log \log n)$	doppelt logarithmisch	Binäre sortierte Suche interpoliert
$\mathcal{O}(\log n)$	logarithmisch	Binäre sortierte Suche
$\mathcal{O}(\sqrt{n})$	wie die Wurzelfunktion	Primzahltest (naiv)
$\mathcal{O}(n)$	linear	Unsortierte naive Suche
$\mathcal{O}(n\log n)$	superlinear / loglinear	Gute Sortieralgorithmen
$\mathcal{O}(n^2)$	quadratisch	Einfache Sortieralgorithmen
$\mathcal{O}(n^c)$	polynomial	Matrixmultiplikation
$\mathcal{O}(c^n)$	exponentiell	Travelling Salesman Dynamic Programming
$\mathcal{O}(n!)$	faktoriell	Travelling Salesman naiv

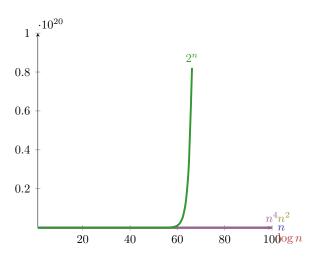
Kleine n



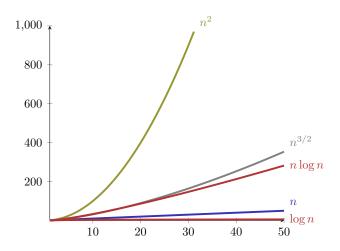
Grössere n



$\hbox{``Grosse''}\ n$



Logarithmen!



Zeitbedarf

45

Annahme: 1 Operation = $1\mu s$.

Problemgrösse	1	100	10000	10^{6}	10^{9}
$\log_2 n$	$1\mu s$	$7\mu s$	$13\mu s$	$20\mu s$	$30\mu s$
n	$1\mu s$	$100 \mu s$	1/100s	1s	17 Minuten
$n\log_2 n$	$1\mu s$	$700 \mu s$	$13/100 \mu s$	20s	8.5 Stunden
n^2	$1\mu s$	1/100s	1.7 Minuten	11.5 Tage	317 Jahrhund.
2^n	$1\mu s$	10^{14} Jahrh.	$pprox \infty$	$pprox \infty$	$pprox \infty$

Nützliches

Theorem 2 Seien $f, g: \mathbb{N} \to \mathbb{R}^+$ zwei Funktionen. Dann gilt: 1. $\lim_{n \to \infty} \frac{f(n)}{g(n)} = 0 \Rightarrow f \in \mathcal{O}(g), \, \mathcal{O}(f) \subsetneq \mathcal{O}(g).$ 2. $\lim_{n \to \infty} \frac{f(n)}{g(n)} = C > 0 \, (C \, konstant) \Rightarrow f \in \Theta(g).$ 3. $\frac{f(n)}{g(n)} \underset{n \to \infty}{\to} \infty \Rightarrow g \in \mathcal{O}(f), \, \mathcal{O}(g) \subsetneq \mathcal{O}(f).$

Erinnerung: Effizienz: Arrays vs. Verkettete Listen

- Speicher: Unser avec belegt ungefähr n ints (Vektorgrösse n), unser 11vec ungefähr 3n ints (ein Zeiger belegt i.d.R. 8 Byte)
- Laufzeit (mit avec = std::vector, llvec = std::list):

```
| Prepending (insert at front) [100,000x]:
| Avec: 675 ms | Avec: 10 ms | Avec: 13 ms | Avec: 14 ms | Avec: 525 ms | Av
```

Zur Notation

Übliche informelle Schreibweise

$$f = \mathcal{O}(g)$$

ist zu verstehen als $f \in \mathcal{O}(g)$. Es gilt nämlich

$$f_1 = \mathcal{O}(g), f_2 = \mathcal{O}(g) \not\Rightarrow f_1 = f_2!$$

$$n = \mathcal{O}(n^2), n^2 = \mathcal{O}(n^2)$$
 aber natürlich $n \neq n^2$.

Wir vermeiden die informelle "=" Schreibweise, wo sie zu Mehrdeutigkeiten führen könnte.

49

Asymptotische Laufzeiten

Mit unserer neuen Sprache $(\Omega, \mathcal{O}, \Theta)$ können wir das Verhalten der Datenstrukturen und ihrer Algorithmen präzisieren.

Typische Asymptotische Laufzeiten (Vorgriff!)

Datenstruktur	Wahlfreier Zugriff	Einfügen	Nächstes	Einfügen nach Element	Suchen
std::vector	$\Theta(1)$	$\Theta(1) A$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$
std::list	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$
std::set	-	$\Theta(\log n)$	$\Theta(\log n)$	-	$\Theta(\log n)$
std::unordered_set	_	$\Theta(1) P$	-	-	$\Theta(1) P$

A= amortisiert, P= erwartet, sonst schlechtester Fall ("worst case")

Komplexität

Komplexität eines Problems ${\cal P}$

Minimale (asymptotische) Kosten über alle Algorithmen A, die P lösen.

Komplexität der Elementarmultiplikation zweier Zahlen der Länge n ist $\Omega(n)$ und $\mathcal{O}(n^{\log_3 2})$ (Karatsuba Ofman).

Komplexität

Problem	Komplexität	$\mathcal{O}(n)$	$\mathcal{O}(n)$	$\mathcal{O}(n^2)$	$\Omega(n \log n)$
Algorithmus	Kosten ²	\uparrow $3n-4$. /	· . /	$\bigcup_{\Omega(n\log n)}$
Programm	Laufzeit	$\Theta(n)$	$\mathcal{O}(n)$	$\Theta(n^2)$	$\bigcup_{\Omega(n\log n)}$

'Anzani Elementaroperation

3. Beispiele

Korrektheit eines Algorithmus oder seiner Implementation zeigen, Rekursion und Rekurrenzen [Literaturangaben bei den Beispielen]

3.1 Altägyptische Multiplikation

Altägyptische Multiplikation – Ein Beispiel, wie man Korrektheit von Algorithmen zeigen kann.

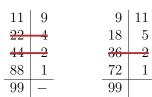
55

²Anzahl Elementaroperationen

Altägyptische Multiplikation

3

Berechnung von $11 \cdot 9$



- 1. Links verdoppeln, rechts ganzzahlig halbieren.
- 2. Gerade Zahl rechts \Rightarrow Zeile streichen.
- 3. Übrige Zeilen links addieren.

Vorteile

- Kurze Beschreibung, einfach zu verstehen.
- Effizient für Computer im Dualsystem: Verdoppeln = Left Shift, Halbieren = Right Shift

left shift
$$9 = 01001_2 \rightarrow 10010_2 = 18$$

right shift $9 = 01001_2 \rightarrow 00100_2 = 4$

Fragen

- Für welche Eingaben liefert der Algorithums das richtige Resultat (in endlicher Zeit)?
- Wie beweist man seine Korrektheit?
- Was ist ein gutes Mass für seine Effizienz?

Die Essenz

Wenn b > 1, $a \in \mathbb{Z}$, dann:

$$a \cdot b = egin{cases} 2a \cdot rac{b}{2} & \text{falls } b \text{ gerade,} \ a + 2a \cdot rac{b-1}{2} & \text{falls } b \text{ ungerade.} \end{cases}$$

58

³Auch bekannt als Russiche Bauernmulltiplikation

Terminierung

$a\cdot b= egin{cases} a & \text{falls } b=1 \text{,} \\ 2a\cdot rac{b}{2} & \text{falls } b ext{ gerade,} \\ a+2a\cdot rac{b-1}{2} & \text{falls } b ext{ ungerade.} \end{cases}$

Rekursiv funktional notiert

$$f(a,b) = \begin{cases} a & \text{falls } b = 1\text{,} \\ f(2a, \frac{b}{2}) & \text{falls } b \text{ gerade,} \\ a + f(2a, \frac{b-1}{2}) & \text{falls } b \text{ ungerade.} \end{cases}$$

61

Als Funktion programmiert

```
// pre: b>0
// post: return a*b
int f(int a, int b){
  if(b==1)
    return a;
  else if (b%2 == 0)
    return f(2*a, b/2);
  else
    return a + f(2*a, (b-1)/2);
}
```

Korrektheit: Mathematischer Beweis

$$f(a,b) = \begin{cases} a & \text{falls } b = 1\text{,} \\ f(2a,\frac{b}{2}) & \text{falls } b \text{ gerade,} \\ a + f(2a \cdot \frac{b-1}{2}) & \text{falls } b \text{ ungerade.} \end{cases}$$

Zu zeigen: $f(a,b) = a \cdot b$ für $a \in \mathbb{Z}$, $b \in \mathbb{N}^+$.

Korrektheit: Mathematischer Beweis per Induktion

Sei $a \in \mathbb{Z}$, zu zeigen $f(a,b) = a \cdot b \quad \forall b \in \mathbb{N}^+$. Anfang: $f(a,1) = a = a \cdot 1$ Hypothese: $f(a,b') = a \cdot b' \quad \forall 0 < b' \le b$ Schritt: $f(a,b') = a \cdot b' \quad \forall 0 < b' \le b \stackrel{!}{\Rightarrow} f(a,b+1) = a \cdot (b+1)$ $f(a,b+1) = \begin{cases} f(2a,\underbrace{b+1}_{0 < \cdot \le b}) \stackrel{i.H.}{=} a \cdot (b+1) & \text{falls } b > 0 \text{ ungerade,} \\ a + f(2a,\underbrace{b}_{0 < \cdot \le b}) \stackrel{i.H.}{=} a + a \cdot b & \text{falls } b > 0 \text{ gerade.} \end{cases}$

[Code-Umformung: Endrekursion]

Die Rekursion lässt sich endrekursiv schreiben

```
// pre: b>0
// pre: b>0
                                        // post: return a*b
                                        int f(int a, int b){
// post: return a*b
int f(int a, int b){
                                          if(b==1)
 if(b==1)
                                           return a;
                                          int z=0;
   return a:
 else if (b\%2 == 0)
                                          if (b%2 != 0){
   return f(2*a, b/2);
                                            --b;
                                            z=a:
   return a + f(2*a, (b-1)/2);
                                          return z + f(2*a, b/2);
```

[Code-Umformung: Endrekursion ⇒ Iteration]

```
int f(int a, int b) {
                                        int res = 0;
// pre: b>0
                                        while (b != 1) {
// post: return a*b
                                          int z = 0;
int f(int a, int b){
                                          if (b % 2 != 0){
 if(b==1)
                                            --b:
   return a;
                                            z = a;
 int z=0;
 if (b%2 != 0){
                                          res += z;
   --b:
                                          a *= 2; // neues a
   z=a;
                                          b /= 2; // neues b
 return z + f(2*a, b/2);
                                        res += a; // Basisfall b=1
                                        return res;
                                      }
```

[Code-Umformung: Vereinfachen]

```
int f(int a, int b) {
 int res = 0;
                                        // pre: b>0
 while (b != 1) {
                                        // post: return a*b
   int z = 0;
                                        int f(int a, int b) {
   if (b % 2 != 0){
                                          int res = 0;
     --b; --> Teil der Division
                                          while (b > 0) {
     z = a; Direkt in res
                                            if (b % 2 != 0)
                                             res += a;
   res += z;
                                            a *= 2;
   a *= 2;
                                            b /= 2;
   b /= 2;
                                          return res;
 res += a; ---- in den Loop
 return res;
```

66

Korrektheit: Argumentation mit Invarianten!

```
// pre: b>0
// post: return a*b
int f(int a, int b) {
                                          Sei x := a \cdot b.
  int res = 0;
                                         Hier gilt x = a \cdot b + res
  while (b > 0) {
    if (b % 2 != 0){
                                          Wenn hier x = a \cdot b + res \dots
      res += a;
      --b;
                                          ... dann auch hier x = a \cdot b + res
                                          b gerade
    a *= 2;
    b /= 2;
                                          Hier gilt x = a \cdot b + res
                                          Hier gilt x = a \cdot b + res und b = 0
  return res;
                                          Also res = x.
```

Zusammenfassung

Der Ausdruck $a \cdot b + res$ ist eine Invariante.

- Werte von *a, b, res* ändern sich, aber die Invariante bleibt "im Wesentlichen" unverändert: Invariante vorübergehend durch eine Anweisung zerstört, aber dann darauf wieder hergestellt. Betrachtet man solche Aktionsfolgen als atomar, bleibt der Wert tatsächlich invariant
- Insbesondere erhält die Schleife die Invariante (Schleifeninvariante), sie wirkt dort wie der Induktionsschritt bei der vollständigen Induktion
- Invarianten sind offenbar mächtige Beweishilfsmittel!

[Weiteres Kürzen]

```
// pre: b>0
// post: return a*b
                                       // pre: b>0
int f(int a, int b) {
                                       // post: return a*b
 int res = 0;
                                       int f(int a, int b) {
 while (b > 0) {
                                         int res = 0;
   if (b % 2 != 0){
                                         while (b > 0) {
     res += a;
                                          res += a * (b\%2);
     --b;
                                          a *= 2;
                                          b /= 2;
   a *= 2;
   b /= 2;
                                         return res;
 return res;
```

[Analyse]

69

```
// pre: b>0
// post: return a*b
int f(int a, int b) {
  int res = 0;
  while (b > 0) {
    res += a * (b%2);
    a *= 2;
    b /= 2;
  }
  return res;
}
```

Altägyptische Multiplikation entspricht der Schulmethode zur Basis 2.

1	0	0	1	×	1	0	1	1	
					1	0	0	1	(9)
				1	0	0	1		(18)
				1	1	0	1	1	
		1	0	0	1				(72)
		1	1	0	0	0	1	1	(99)

Effizienz

Frage: Wie lange dauert eine Multiplikation von a und b?

- Mass für die Effizienz
 - Gesamtzahl der elementaren Operationen: Verdoppeln, Halbieren, Test auf "gerade", Addition
 - Im rekursiven wie im iterativen Code: maximal 6 Operationen pro Aufruf bzw. Durchlauf
- Wesentliches Kriterium:
 - Anzahl rekursiver Aufrufe oder
 - Anzahl Schleifendurchläufe(im iterativen Fall)
- $\frac{b}{2^n} \le 1$ gilt für $n \ge \log_2 b$. Also nicht mehr als $6\lceil \log_2 b \rceil$ elementare Operationen.

3.2 Schnelle Multiplikation von Zahlen

[Ottman/Widmayer, Kap. 1.2.3]

Beispiel 2: Multiplikation grosser Zahlen

Primarschule:

 $2\cdot 2=4$ einstellige Multiplikationen. \Rightarrow Multiplikation zweier n-stelliger Zahlen: n^2 einstellige Multiplikationen

Beobachtung

$$ab \cdot cd = (10 \cdot a + b) \cdot (10 \cdot c + d)$$
$$= 100 \cdot a \cdot c + 10 \cdot a \cdot c$$
$$+ 10 \cdot b \cdot d + b \cdot d$$
$$+ 10 \cdot (a - b) \cdot (d - c)$$

Verbesserung?

а	,	b		c	d	
6	,	2		3	7	
				1	4	$d \cdot b$
			1	4		$d \cdot b$
			1	6		$(a-b)\cdot (d-c)$
			1	8		$c \cdot a$
		1	8			$c \cdot a$
=		2	2	9	4	

 $\rightarrow 3$ einstellige Multiplikationen.

Verallgemeinerung

Annahme: zwei n-stellige Zahlen, $n = 2^k$ für ein k.

$$(10^{n/2}a + b) \cdot (10^{n/2}c + d) = 10^n \cdot a \cdot c + 10^{n/2} \cdot a \cdot c + 10^{n/2} \cdot b \cdot d + b \cdot d + 10^{n/2} \cdot (a - b) \cdot (d - c)$$

Rekursive Anwendung dieser Formel: Algorithmus von Karatsuba und Ofman (1962).

Grosse Zahlen

$$6237 \cdot 5898 = \underbrace{62}_{a'} \underbrace{37}_{b'} \cdot \underbrace{58}_{c'} \underbrace{98}_{d'}$$

Rekursive / induktive Anwendung: $a' \cdot c'$, $a' \cdot d'$, $b' \cdot c'$ und $c' \cdot d'$ wie oben berechnen.

 $\rightarrow 3 \cdot 3 = 9$ statt 16 einstellige Multiplikationen.

3

Algorithmus Karatsuba Ofman

```
\begin{array}{ll} \text{Input:} & \text{Zwei } n\text{-stellige } (n>0) \text{ ganze positive Zahlen } x \text{ und } y \text{ mit dezimalen} \\ & \text{Ziffern } (x_i)_{1 \leq i \leq n} \text{ und } (y_i)_{1 \leq i \leq n} \\ \text{Output:} & \text{Produkt } x \cdot y \\ \text{if } n=1 \text{ then} \\ & \text{return } x_1 \cdot y_1 \\ \text{else} \\ & \text{Sei } m:= \left \lfloor \frac{n}{2} \right \rfloor \\ & \text{Unterteile } a:=(x_1,\ldots,x_m), \ b:=(x_{m+1},\ldots,x_n), \ c:=(y_1,\ldots,y_m), \\ & d:=(y_{m+1},\ldots,y_n) \\ & \text{Berechne rekursiv } A:=a \cdot c, \ B:=b \cdot d, \ C:=(a-b) \cdot (d-c) \\ & \text{Berechne } R:=10^n \cdot A+10^m \cdot A+10^m \cdot B+B+10^m \cdot C \\ & \text{return } R \\ \end{array}
```

Analyse

M(n): Anzahl einstelliger Multiplikationen.

Rekursive Anwendung des obigen Algorithmus ⇒ Rekursionsgleichung:

$$M(2^k) = \begin{cases} 1 & \text{falls } k = 0, \\ 3 \cdot M(2^{k-1}) & \text{falls } k > 0. \end{cases}$$
 (R)

Teleskopieren

Iteratives Einsetzen der Rekursionsformel zum Lösen der Rekursionsgleichung.

$$\begin{split} M(2^k) &= 3 \cdot M(2^{k-1}) = 3 \cdot 3 \cdot M(2^{k-2}) = 3^2 \cdot M(2^{k-2}) \\ &= \dots \\ &\stackrel{!}{=} 3^k \cdot M(2^0) = 3^k. \end{split}$$

82

Beweis: Vollständige Induktion

Hypothese H(k):

$$M(2^k) = F(k) := 3^k.$$
 (H)

Behauptung:

H(k) gilt für alle $k \in \mathbb{N}_0$.

Induktionsanfang k = 0:

$$M(2^0) \stackrel{R}{=} 1 = F(0)$$
.

Induktionsschritt $H(k) \Rightarrow H(k+1)$:

$$M(2^{k+1}) \stackrel{R}{=} 3 \cdot M(2^k) \stackrel{H(k)}{=} 3 \cdot F(k) = 3^{k+1} = F(k+1).$$
 \checkmark

Vergleich

Primarschulmethode: n^2 einstellige Multiplikationen. Karatsuba/Ofman:

$$M(n) = 3^{\log_2 n} = (2^{\log_2 3})^{\log_2 n} = 2^{\log_2 3 \log_2 n} = n^{\log_2 3} \approx n^{1.58}.$$

Beispiel: 1000-stellige Zahl: $1000^2/1000^{1.58} \approx 18$.

Bestmöglicher Algorithums?

Wir kennen nun eine obere Schranke $n^{\log_2 3}$.

Es gibt praktisch (für grosses n) relevante, schnellere Algorithmen. Beispiel: Schönhage-Strassen-Algorithmus (1971) basierend auf schneller Fouriertransformation mit Laufzeit $\mathcal{O}(n\log n \cdot \log\log n)$. Die beste obere Schranke ist nicht bekannt. ⁴

Untere Schranke: n. Jede Ziffer muss zumindest einmal angeschaut werden.

Anhang: Asymptotik mit Additionen und Shifts

Bei jeder Multiplikation zweier n-stelliger Zahlen kommt auch noch eine konstante Anzahl Additionen, Subtraktionen und Shifts dazu Additionen, Subtraktionen und Shifts von n stelligen Zahlen kosten $\mathcal{O}(n)$ Daher ist die asymptotische Laufzeit eigentlich (mit geeignetem c>1, d>0) bestimmt durch die folgende Rekurrenz

$$T(n) = \begin{cases} 3 \cdot T\left(\frac{1}{2}n\right) + c \cdot n & \text{falls } n > 1 \\ d & \text{sonst} \end{cases}$$

87

Anhang: Asymptotik mit Additionen und Shifts

Annahme: $n=2^k$, k>0

$$\begin{split} T(2^k) &= 3 \cdot T\left(2^{k-1}\right) + c \cdot 2^k \\ &= 3 \cdot \left(3 \cdot T(2^{k-2}) + c \cdot 2^{k-1}\right) + c \cdot 2^k \\ &= 3 \cdot \left(3 \cdot \left(3 \cdot T(2^{k-3}) + c \cdot 2^{k-2}\right) + c \cdot 2^{k-1}\right) + c \cdot 2^k \\ &= 3 \cdot \left(3 \cdot \left(\dots \left(3 \cdot T(2^{k-k}) + c \cdot 2^1\right)\dots\right) + c \cdot 2^{k-1}\right) + c \cdot 2^k \\ &= 3^k \cdot d + c \cdot 2^k \sum_{i=0}^{k-1} \frac{3^i}{2^i} = 3^k \cdot d + c \cdot 2^k \frac{\frac{3^k}{2^k} - 1}{\frac{3}{2} - 1} \\ &= 3^k (d + 2c) - 2c \cdot 2^k \end{split}$$

Somit $T(2^k) = 3^k(d+2c) - 2c \cdot 2^k \in \Theta(3^k) = \Theta(3^{\log_2 n}) = \Theta(n^{\log_2 3}).$

3.3 Maximum Subarray Problem

Algorithmenentwurf – Maximum Subarray Problem [Ottman/Widmayer, Kap. 1.3] Divide and Conquer [Ottman/Widmayer, Kap. 1.2.2. S.9; Cormen et al, Kap. 4-4.1]

 $^{^4}$ Im März 2019 haben David Harvey and Joris van der Hoeven einen, praktisch noch irrelevanten, $\mathcal{O}(n\log n)$ Algorithmus vorgestellt. Man vermutet, dass $n\log n$ die beste obere Grenze ist, hat es aber noch nicht bewiesen.

Algorithmenentwurf

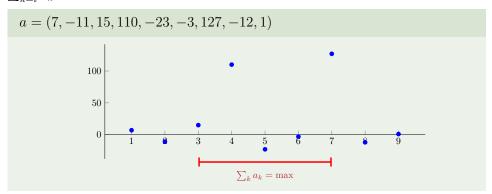
Induktive Entwicklung eines Algorithmus: Zerlegung in Teilprobleme, Verwendung der Lösungen der Teilproblem zum Finden der endgültigen Lösung.

Ziel: Entwicklung des asymptotisch effizientesten (korrekten) Algorithmus. Effizienz hinsichtlich der Laufzeitkosten (# Elementaroperationen) oder / und Speicherbedarf.

Maximum Subarray Problem

Gegeben: ein Array von n reellen Zahlen (a_1,\ldots,a_n) .

Gesucht: Teilstück [i,j], $1 \le i \le j \le n$ mit maximaler positiver Summe $\sum_{k=i}^{j} a_k$.



Naiver Maximum Subarray Algorithmus

 $\begin{array}{ll} \textbf{Input:} & \text{Eine Folge von } n \text{ Zahlen } (a_1,a_2,\ldots,a_n) \\ \textbf{Output:} & I, \ J \text{ mit } \sum_{k=I}^J a_k \text{ maximal.} \\ M \leftarrow 0; \ I \leftarrow 1; \ J \leftarrow 0 \\ \textbf{for } i \in \{1,\ldots,n\} \textbf{ do} \\ & \quad \text{for } j \in \{i,\ldots,n\} \textbf{ do} \\ & \quad \quad \left[\begin{array}{c} m = \sum_{k=i}^j a_k \\ \textbf{if } m > M \textbf{ then} \\ & \quad \quad \ \ \, \bot \\ M \leftarrow m; \ I \leftarrow i; \ J \leftarrow j \end{array} \right] \\ \textbf{return } I,J \end{array}$

Analyse

Theorem 3

Der naive Algorithmus für das Maximum Subarray Problem führt $\Theta(n^3)$ Additionen durch.

Beweis:

$$\sum_{i=1}^{n} \sum_{j=i}^{n} (j-i+1) = \sum_{i=1}^{n} \sum_{j=0}^{n-i} (j+1) = \sum_{i=1}^{n} \sum_{j=1}^{n-i+1} j = \sum_{i=1}^{n} \frac{(n-i+1)(n-i+2)}{2}$$

$$= \sum_{i=0}^{n} \frac{i \cdot (i+1)}{2} = \frac{1}{2} \left(\sum_{i=1}^{n} i^2 + \sum_{i=1}^{n} i \right)$$

$$= \frac{1}{2} \left(\frac{n(2n+1)(n+1)}{6} + \frac{n(n+1)}{2} \right) = \frac{n^3 + 3n^2 + 2n}{6} = \Theta(n^3).$$

Beobachtung

$$\sum_{k=i}^{j} a_k = \underbrace{\left(\sum_{k=1}^{j} a_k\right)}_{S_j} - \underbrace{\left(\sum_{k=1}^{i-1} a_k\right)}_{S_{i-1}}$$

Präfixsummen

$$S_i := \sum_{k=1}^i a_k.$$

Analyse

Theorem 4

Der Präfixsummen Algorithmus für das Maximum Subarray Problem führt $\Theta(n^2)$ Additionen und Subtraktionen durch.

Beweis:

$$\sum_{i=1}^{n} 1 + \sum_{i=1}^{n} \sum_{j=i}^{n} 1 = n + \sum_{i=1}^{n} (n-i+1) = n + \sum_{i=1}^{n} i = \Theta(n^2)$$

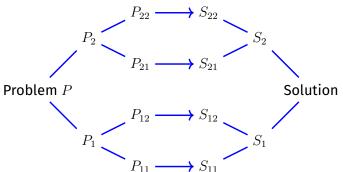
Maximum Subarray Algorithmus mit Präfixsummen

$$\begin{array}{lll} \textbf{Input:} & \text{Eine Folge von } n \text{ Zahlen } (a_1,a_2,\ldots,a_n) \\ \textbf{Output:} & I, \ J \text{ mit } \sum_{k=J}^J a_k \text{ maximal.} \\ \mathcal{S}_0 \leftarrow 0 \\ \textbf{for } i \in \{1,\ldots,n\} \text{ do } // \text{ Präfixsumme} \\ & \ \ \, \bigcup \text{ } \mathcal{S}_i \leftarrow \mathcal{S}_{i-1} + a_i \\ M \leftarrow 0; \ I \leftarrow 1; \ J \leftarrow 0 \\ \textbf{for } i \in \{1,\ldots,n\} \text{ do } \\ & \ \ \, \big| \begin{array}{ll} m = \mathcal{S}_j - \mathcal{S}_{i-1} \\ \textbf{if } m > M \text{ then} \\ & \ \ \, \bigcup \text{ } M \leftarrow m; \ I \leftarrow i; \ J \leftarrow j \end{array}$$

divide et impera

Teile und (be)herrsche (engl. divide and conquer)

Zerlege das Problem in Teilprobleme, deren Lösung zur vereinfachten Lösung des Gesamtproblems beitragen.



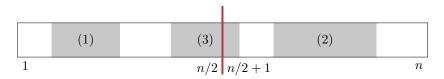
Maximum Subarray - Divide

- Divide: Teile das Problem in zwei (annähernd) gleiche Hälften auf: $(a_1, \ldots, a_n) = (a_1, \ldots, a_{\lfloor n/2 \rfloor}, a_{\lfloor n/2 \rfloor+1}, \ldots, a_1)$
- Vereinfachende Annahme: $n = 2^k$ für ein $k \in \mathbb{N}$.

Maximum Subarray - Conquer

Sind i, j die Indizes einer Lösung \Rightarrow Fallunterscheidung:

- 1. Lösung in linker Hälfte $1 \le i \le j \le n/2 \Rightarrow$ Rekursion (linke Hälfte)
- 2. Lösung in rechter Hälfte $n/2 < i \le j \le n \Rightarrow$ Rekursion (rechte Hälfte)
- 3. Lösung in der Mitte $1 \le i \le n/2 < j \le n \Rightarrow$ Nachfolgende Beobachtung



8

Maximum Subarray - Beobachtung

Annahme: Lösung in der Mitte $1 \le i \le n/2 < j \le n$

$$\begin{split} S_{\text{max}} &= \max_{\substack{1 \leq i \leq n/2 \\ n/2 < j \leq n}} \sum_{k=i}^{j} a_k = \max_{\substack{1 \leq i \leq n/2 \\ n/2 < j \leq n}} \left(\sum_{k=i}^{n/2} a_k + \sum_{k=n/2+1}^{j} a_k \right) \\ &= \max_{\substack{1 \leq i \leq n/2 \\ 1 \leq i \leq n/2}} \sum_{k=i}^{n/2} a_k + \max_{\substack{n/2 < j \leq n \\ k=n/2+1}} \sum_{k=n/2+1}^{j} a_k \\ &= \max_{\substack{1 \leq i \leq n/2 \\ 1 \leq i \leq n/2}} \underbrace{S_{n/2} - S_{i-1}}_{\text{Suffixsumme}} + \max_{\substack{n/2 < j \leq n \\ Pr\"{affixsumme}}} \underbrace{S_j - S_{n/2}}_{\text{Pr\"{affixsumme}}} \end{split}$$

Maximum Subarray Divide and Conquer Algorithmus

```
\begin{array}{ll} \textbf{Input:} & \text{Eine Folge von } n \text{ Zahlen } (a_1,a_2,\ldots,a_n) \\ \textbf{Output:} & \text{Maximales } \sum_{k=i'}^{j'} a_k. \\ \textbf{if } n=1 \textbf{ then} \\ | \textbf{ return } \max\{a_1,0\} \\ \textbf{else} \\ | & \text{Unterteile } a=(a_1,\ldots,a_n) \text{ in } A_1=(a_1,\ldots,a_{n/2}) \text{ und } A_2=(a_{n/2+1},\ldots,a_n) \\ & \text{Berechne rekursiv beste Lösung } W_1 \text{ in } A_1 \\ & \text{Berechne grösste Suffixsumme } S \text{ in } A_2 \\ & \text{Berechne grösste Präfixsumme } P \text{ in } A_2 \\ & \text{Setze } W_3 \leftarrow S+P \\ & \textbf{ return } \max\{W_1,W_2,W_3\} \end{array}
```

Analyse

Theorem 5

Der Divide and Conquer Algorithmus für das Maximum Subarray Sum Problem führt $\Theta(n\log n)$ viele Additionen und Vergleiche durch.

Analyse

Rekursionsgleichung

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

Analyse

Input: Eine Folge von n Zahlen (a_1,a_2,\ldots,a_n) Output: Maximales $\sum_{k=i'}^{j'} a_k$.

if n=1 then $\Theta(1)$ return $\max\{a_1,0\}$ else $\Theta(1)$ Unterteile $a=(a_1,\ldots,a_n)$ in $A_1=(a_1,\ldots,a_{n/2})$ und $A_2=(a_{n/2+1},\ldots,a_n)$ T(n/2) Berechne rekursiv beste Lösung W_1 in A_1 T(n/2) Berechne rekursiv beste Lösung W_2 in A_2 $\Theta(n)$ Berechne grösste Suffixsumme S in A_1

Analyse

102

 $\Theta(1)$ Setze $W_3 \leftarrow S + P$ $\Theta(1)$ return $\max\{W_1, W_2, W_3\}$

 $\Theta(n)$ Berechne grösste Präfixsumme P in A_2

Mit $n=2^k$:

$$\overline{T}(k) := T(2^k) = \begin{cases} c & \text{falls } k = 0 \\ 2\overline{T}(k-1) + a \cdot 2^k & \text{falls } k > 0 \end{cases}$$

Lösung:

$$\overline{T}(k) = 2^k \cdot c + \sum_{i=0}^{k-1} 2^i \cdot a \cdot 2^{k-i} = c \cdot 2^k + a \cdot k \cdot 2^k = \Theta(k \cdot 2^k)$$

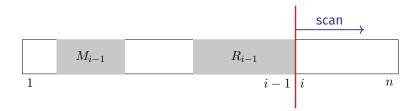
also

$$T(n) = \Theta(n \log n)$$

104

Maximum Subarray Sum Problem – Induktiv

Annahme: Maximaler Wert M_{i-1} der Subarraysumme für (a_1,\ldots,a_{i-1}) $(1< i\leq n)$ bekannt.



 a_i : erzeugt höchstens Intervall am Rand (Präfixsumme).

$$R_{i-1} \Rightarrow R_i = \max\{R_{i-1} + a_i, 0\}$$

Analyse

Theorem 6

Der induktive Algorithmus für das Maximum Subarray Sum Problem führt $\Theta(n)$ viele Additionen und Vergleiche durch.

Induktiver Maximum Subarray Algorithmus

```
\begin{array}{ll} \textbf{Input:} & \text{Eine Folge von } n \text{ Zahlen } (a_1,a_2,\ldots,a_n). \\ \textbf{Output:} & \max\{0,\max_{i,j}\sum_{k=i}^{j}a_k\}. \\ M \leftarrow 0 \\ R \leftarrow 0 \\ \textbf{for } i=1\ldots n \text{ do} \\ & R \leftarrow R+a_i \\ & \textbf{if } R < 0 \text{ then} \\ & \bot R \leftarrow 0 \\ & \textbf{if } R > M \text{ then} \\ & \bot M \leftarrow R \\ \\ \textbf{return } M; \end{array}
```

106

Komplexität des Problems?

Geht es besser als $\Theta(n)$?

Jeder korrekte Algorithmus für das Maximum Subarray Sum Problem muss jedes Element im Algorithmus betrachten.

Annahme: der Algorithmus betrachtet nicht a_i .

- 1. Lösung des Algorithmus enthält a_i . Wiederholen den Algorithmus mit genügend kleinem a_i , so dass die Lösung den Punkt nicht enthalten hätte dürfen.
- 2. Lösung des Algorithmus enthält a_i nicht. Wiederholen den Algorithmus mit genügend grossem a_i , so dass die Lösung a_i hätten enthalten müssen.

Komplexität des Maximum Subarray Sum Problems

Theorem 7

Das Maximum Subarray Sum Problem hat Komplexität $\Theta(n)$.

Beweis: Induktiver Algorithmus mit asymptotischer Laufzeit $\mathcal{O}(n)$. Jeder Algorithmus hat Laufzeit $\Omega(n)$. Somit ist die Komplexität $\Omega(n) \cap \mathcal{O}(n) = \Theta(n)$.

3.4 Anhang

Herleitung und Wiederholung einiger mathematischen Formeln

110

111

Logarithmen

$$\log_a y = x \Leftrightarrow a^x = y \quad (a > 0, y > 0)$$

$$\log_a (x \cdot y) = \log_a x + \log_a y \qquad a^x \cdot a^y = a^{x+y}$$

$$\log_a \frac{x}{y} = \log_a x - \log_a y \qquad \frac{a^x}{a^y} = a^{x-y}$$

$$\log_a x^y = y \log_a x \qquad a^{x \cdot y} = (a^x)^y$$

$$\log_a n! = \sum_{i=1}^n \log i$$

$$\log_b x = \log_b a \cdot \log_a x \qquad a^{\log_b x} = x^{\log_b a}$$

Letzteres sieht man durch Einsetzen von $x \to a^{\log_a x}$

Summen

$$\sum_{i=0}^{n} i = \frac{n \cdot (n+1)}{2} \in \Theta(n^2)$$

Trick

$$\sum_{i=0}^{n} i = \frac{1}{2} \left(\sum_{i=0}^{n} i + \sum_{i=0}^{n} n - i \right) = \frac{1}{2} \sum_{i=0}^{n} i + n - i$$
$$= \frac{1}{2} \sum_{i=0}^{n} n = \frac{1}{2} (n+1) \cdot n$$

Summen

$$\sum_{i=0}^{n} i^2 = \frac{n \cdot (n+1) \cdot (2n+1)}{6}$$

Trick:

$$\sum_{i=1}^{n} i^3 - (i-1)^3 = \sum_{i=0}^{n} i^3 - \sum_{i=0}^{n-1} i^3 = n^3$$

$$\sum_{i=1}^{n} i^3 - (i-1)^3 = \sum_{i=1}^{n} i^3 - i^3 + 3i^2 - 3i + 1 = n - \frac{3}{2}n \cdot (n+1) + 3\sum_{i=0}^{n} i^2$$

$$\Rightarrow \sum_{i=0}^{n} i^2 = \frac{1}{6}(2n^3 + 3n^2 + n) \in \Theta(n^3)$$

Kann einfach verallgemeinert werden: $\sum_{i=1}^{n} i^k \in \Theta(n^{k+1})$.

4. Suchen

Lineare Suche, Binäre Suche, (Interpolationssuche,) Untere Schranken [Ottman/Widmayer, Kap. 3.2, Cormen et al, Kap. 2: Problems 2.1-3,2.2-3,2.3-5]

Geometrische Reihe

$$\sum_{i=0}^{n} \rho^{i} \stackrel{!}{=} \frac{1 - \rho^{n+1}}{1 - \rho}$$

$$\sum_{i=0}^{n} \rho^{i} \cdot (1 - \rho) = \sum_{i=0}^{n} \rho^{i} - \sum_{i=0}^{n} \rho^{i+1} = \sum_{i=0}^{n} \rho^{i} - \sum_{i=1}^{n+1} \rho^{i}$$
$$= \rho^{0} - \rho^{n+1} = 1 - \rho^{n+1}.$$

Für $0 \le \rho < 1$:

$$\sum_{i=0}^{\infty} \rho^i = \frac{1}{1-\rho}$$

Das Suchproblem

Gegeben

■ Menge von Datensätzen.

Telefonverzeichnis, Wörterbuch, Symboltabelle

- Jeder Datensatz hat einen Schlüssel k.
- Schlüssel sind vergleichbar: eindeutige Antwort auf Frage $k_1 \le k_2$ für Schlüssel k_1 , k_2 .

Aufgabe: finde Datensatz nach Schlüssel k.

Suche in Array

Gegeben

- Array A mit n Elementen (A[1], ..., A[n]).
- Schlüssel b

Gesucht: Index k, $1 \le k \le n$ mit A[k] = b oder "nicht gefunden".



Lineare Suche

Durchlaufen des Arrays von A[1] bis A[n].

- Bestenfalls 1 Vergleich.
- lacksquare Schlimmstenfalls n Vergleiche.
- Annahme: Jede Anordnung der n Schlüssel ist gleichwahrscheinlich. Erwartete Anzahl Vergleiche für die erfolgreiche Suche:

$$\frac{1}{n} \sum_{i=1}^{n} i = \frac{n+1}{2}.$$

118

Suche im sortierten Array

Gegeben

- Sortiertes Array A mit n Elementen $(A[1], \ldots, A[n])$ mit $A[1] \leq A[2] \leq \cdots \leq A[n]$.
- Schlüssel b

Gesucht: Index k, $1 \le k \le n$ mit A[k] = b oder "nicht gefunden".

10	20	22	24	28	32	35	38	41	42
1	2	3	4	5	6	7	8	9	10

Divide and Conquer!

Suche b=23.

b < 28	42	41	38	35	32	28	24	22	20	10
	10	9	8	7	6	5	4	3	2	1
b > 20	42	41	38	35	32	28	24	22	20	10
•	10	9	8	7	6	5	4	3	2	1
b > 22	42	41	38	35	32	28	24	22	20	10
	10	9	8	7	6	5	4	3	2	1
b < 24	42	41	38	35	32	28	24	22	20	10
	10	9	8	7	6	5	4	3	2	1
erfolglos	42	41	38	35	32	28	24	22	20	10
	10	9	8	7	6	5	4	3	2	1

Binärer Suchalgorithmus BSearch(A, l, r, b)

```
Input: Sortiertes Array A von n Schlüsseln. Schlüssel b. Bereichsgrenzen 1 \leq l, r \leq n mit l \leq r oder l = r + 1.

Output: Index m \in [l, \ldots, r + 1], so dass A[i] \leq b für alle l \leq i < m und A[i] \geq b für alle m < i \leq r. m \leftarrow \lfloor (l+r)/2 \rfloor if l > r then l = l erfolglose Suche return l = l else if l = l then l = l gefunden return l else if l = l then l Element liegt links return BSearch l return BSearch l else l so l search l else l so l search l search l return BSearch l search l return BSearch l return l return
```

Analyse (schlechtester Fall)

$$T(n) = egin{cases} d & \text{falls } n=1 \text{,} \\ T(n/2) + c & \text{falls } n>1 \text{.} \end{cases}$$

Vermutung : $T(n) = d + c \cdot \log_2 n$

Beweis durch Induktion:

- Induktionsanfang: T(1) = d.
- Hypothese: $T(n/2) = d + c \cdot \log_2 n/2$
- \blacksquare Schritt $(n/2 \rightarrow n)$

$$T(n) = T(n/2) + c = d + c \cdot (\log_2 n - 1) + c = d + c \log_2 n.$$

Analyse (schlechtester Fall)

Rekurrenz ($n=2^k$)

$$T(n) = egin{cases} d & \text{falls } n = 1, \\ T(n/2) + c & \text{falls } n > 1. \end{cases}$$

Teleskopieren:

$$\begin{split} T(n) &= T\left(\frac{n}{2}\right) + c = T\left(\frac{n}{4}\right) + 2c = \dots \\ &= T\left(\frac{n}{2^i}\right) + i \cdot c \\ &= T\left(\frac{n}{n}\right) + \log_2 n \cdot c = d + c \cdot \log_2 n \in \Theta(\log n) \end{split}$$

123

122

Resultat

Theorem 8

Der Algorithmus zur binären sortierten Suche benötigt im schlechtesten Fall $\Theta(\log n)$ Elementarschritte.

Iterativer binärer Suchalgorithmus

Korrektheit

Algorithmus bricht nur ab, falls A[l..r] leer oder b gefunden.

Invariante: Falls b in A, dann im Bereich A[l..r]

Beweis durch Induktion

- Induktionsanfang: $b \in A[1..n]$ (oder nicht)
- Hypothese: Invariante gilt nach *i* Schritten
- Schritt:

$$b < A[m] \Rightarrow b \in A[l..m-1]$$

 $b > A[m] \Rightarrow b \in A[m+1..r]$

126

[Geht es noch besser?]

Annahme: Gleichverteilung der Werte im Array.

Beispiel

Name "Becker" würde man im Telefonbuch vorne suchen. "Wawrinka" wohl ziemlich weit hinten.

Binäre Suche vergleicht immer zuerst mit der Mitte.

Binäre Suche setzt immer $m = \left| l + \frac{r-l}{2} \right|$.

[Interpolations suche]

Erwartete relative Position von b im Suchintervall [l, r]

$$\rho = \frac{b - A[l]}{A[r] - A[l]} \in [0, 1].$$

Neue "Mitte": $l + \rho \cdot (r - l)$

Anzahl Vergleiche im Mittel $\mathcal{O}(\log \log n)$ (ohne Beweis).

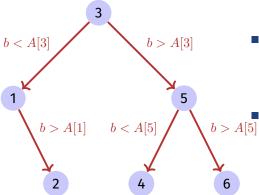
Ist Interpolationssuche also immer zu bevorzugen?

Nein: Anzahl Vergleiche im schlimmsten Fall $\Omega(n)$.

Untere Schranke

Binäre Suche (im schlechtesten Fall): $\Theta(\log n)$ viele Vergleiche. Gilt für *jeden* Suchalgorithms in sortiertem Array (im schlechtesten Fall): Anzahl Vergleiche = $\Omega(\log n)$?

Entscheidungsbaum



- Für jede Eingabe b = A[i]muss Algorithmus erfolgreich sein \Rightarrow Baum enthält mindestens n Knoten.
- Anzahl Vergleiche im
 schlechtesten Fall = Höhe des Baumes = maximale Anzahl Knoten von Wurzel zu Blatt.

Entscheidungsbaum

Binärer Baum der Höhe h hat höchstens $2^0+2^1+\cdots+2^{h-1}=2^h-1<2^h$ Knoten.

$$2^h > n \Rightarrow h > \log_2 n$$

Entscheidungsbaum mit n Knoten hat mindestens Höhe $\log_2 n$. Anzahl Entscheidungen = $\Omega(\log n)$.

Theorem 9

Jeder Algorithmus zur vergleichsbasierten Suche in sortierten Daten der Länge n benötigt im schlechtesten Fall $\Omega(\log n)$ Vergleichsschritte.

Untere Schranke für Suchen in unsortiertem Array

Theorem 10

Jeder vergleichsbasierte Algorithmus zur Suche in unsortierten Daten der Länge n benötigt im schlechtesten Fall $\Omega(n)$ Vergleichsschritte.

Versuch

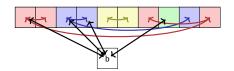
Korrekt?

"Beweis": Um b in A zu finden, muss b mit jedem Element A[i] ($1 \le i \le n$) verglichen werden.

Falsch! Vergleiche zwischen Elementen von A möglich!

134

Besseres Argument



- Unterteilung der Vergleiche: Anzahl Vergleiche mit b: e Anzahl Vergleiche untereinander ohne b: i
- Vergleiche erzeugen q Gruppen. Initial: q = n.
- Vereinigen zweier Gruppen benötigt mindestens einen (internen Vergleich: n - q < i.
- \blacksquare Mindestens ein Element pro Gruppe muss mit b verglichen werden: e > q.
- Anzahl Vergleiche $i + e \ge n g + g = n$.

Das Auswahlproblem

Eingabe

- Unsortiertes Array $A = (A_1, \ldots, A_n)$ paarweise verschiedener Werte
- \blacksquare Zahl $1 \le k \le n$.

Ausgabe: A[i] mit $|\{j : A[j] < A[i]\}| = k - 1$

Spezialfälle

k=1: Minimum: Algorithmus mit n Vergleichsoperationen trivial. k=n: Maximum: Algorithmus mit n Vergleichsoperationen trivial. $k = \lfloor n/2 \rfloor$: Median.

5. Auswählen

Das Auswahlproblem, Randomisierte Berechnung des Medians, Lineare Worst-Case Auswahl [Ottman/Widmayer, Kap. 3.1, Cormen et al, Kap. 9]

Naiver Algorithmus

```
Wiederholt das Minimum entfernen / auslesen: \Theta(k \cdot n). \to Median in \Theta(n^2)
```

Bessere Ansätze

- Sortieren (kommt bald): $\Theta(n \log n)$
- Pivotieren: $\Theta(n)$!

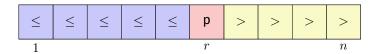
Min und Max

 $oldsymbol{?}$ Separates Finden von Minimum und Maximum in $(A[1],\ldots,A[n])$ benötigt insgesamt 2n Vergleiche. (Wie) geht es mit weniger als 2n Vergleichen für beide gemeinsam?

① Es geht mit $\frac{3}{2}n$ Vergleichen: Vergleiche jeweils 2 Elemente und deren kleineres mit Min und grösseres mit Max.⁵ Possible with $\frac{3}{2}n$ comparisons: compare 2 elements each and then the smaller one with min and the greater one with max.⁶

Pivotieren

- 1. Wähle ein (beliebiges) Element p als Pivotelement
- 2. Teile A in zwei Teile auf, bestimme dabei den Rang von p, indem die Anzahl der Indizes i mit $A[i] \le p$ gezählt werden.
- 3. Rekursion auf dem relevanten Teil. Falls k=r, dann gefunden.



⁵Das liefert einen Hinweis darauf, dass der naive Algorithmus verbessert werden kann.

⁶An indication that the naive algorithm can be improved.

Algorithmus Partition(A, l, r, p)

return |-1

return |-1

Korrektheit: Fortschritt

Korrektheit: Invariante

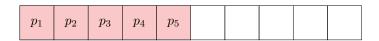
 $\begin{aligned} & \text{Invariante } I \text{: } A_i \leq p \ \forall i \in [0, l) \text{, } A_i \geq p \ \forall i \in (r, n] \text{, } \exists k \in [l, r] \text{ : } A_k = p \text{.} \\ & \text{while } l \leq r \ \text{do} & & & I \\ & & \text{while } A[l]$

142

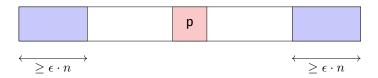
Wahl des Pivots

return |-1

Das Minimum ist ein schlechter Pivot: worst Case $\Theta(n^2)$



Ein guter Pivot hat linear viele Elemente auf beiden Seiten.



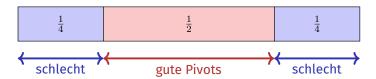
Analyse

Unterteilung mit Faktor q (0 < q < 1): zwei Gruppen mit $q \cdot n$ und $(1 - q) \cdot n$ Elementen (ohne Einschränkung $q \ge 1 - q$).

$$\begin{split} T(n) &\leq T(q \cdot n) + c \cdot n \\ &\leq c \cdot n + q \cdot c \cdot n + T(q^2 \cdot n) \leq \ldots = c \cdot n \sum_{i=0}^{\log_q(n)-1} q^i + T(1) \\ &\leq c \cdot n \sum_{i=0}^{\infty} q^i \quad + d = c \cdot n \cdot \frac{1}{1-q} + d = \mathcal{O}(n) \end{split}$$
 geom. Reihe

Wie bekommen wir das hin?

Der Zufall hilft uns (Tony Hoare, 1961). Wähle in jedem Schritt einen zufälligen Pivot.



Wahrscheinlichkeit für guten Pivot nach einem Versuch: $\frac{1}{2} =: \rho$. Wahrscheinlichkeit für guten Pivot nach k Versuchen: $(1-\rho)^{k-1} \cdot \rho$. Erwartete Anzahl Versuche: $1/\rho = 2$ (Erwartungswert der geometrischen Verteilung:)

147

Algorithmus Quickselect (A, l, r, k)

```
Input: Array A der Länge n. Indizes 1 \le l \le k \le r \le n, so dass für alle
        x \in A[l..r] : |\{j|A[j] \le x\}| \ge l \text{ und } |\{j|A[j] \le x\}| \le r.
Output: Wert x \in A[l..r] mit |\{j|A[j] \le x\}| \ge k und |\{j|x \le A[j]\}| \ge n - k + 1
if l=r then
   return A[l];
x \leftarrow \mathtt{RandomPivot}(A, l, r)
m \leftarrow \mathtt{Partition}(A, l, r, x)
if k < m then
    return QuickSelect(A, l, m-1, k)
else if k > m then
    return QuickSelect(A, m+1, r, k)
else
    return A[k]
```

Algorithmus RandomPivot (A, l, r)

```
Input: Array A der Länge n. Indizes 1 \le l \le r \le n
Output: Zufälliger "guter" Pivot x \in A[l, ..., r]
repeat
    wähle zufälligen Pivot x \in A[l..r]
    p \leftarrow l
     for i = l to r do
        if A[j] \le x then p \leftarrow p + 1
until \left| \frac{3l+r}{4} \right| \le p \le \left\lceil \frac{l+3r}{4} \right\rceil
```

Dieser Algorithmus ist nur von theoretischem Interesse und liefert im Erwartungswert nach 2 Durchläufen einen guten Pivot. Praktisch kann man im Algorithmus Quickselect direkt einen zufälligen Pivot uniformverteilt ziehen oder einen deterministischen Pivot wählen, z.B. den Median von drei Elementen.

Median der Mediane

Ziel: Finde einen Algorithmus, welcher im schlechtesten Fall nur linear viele Schritte benötigt.

Algorithmus Select (k-smallest)

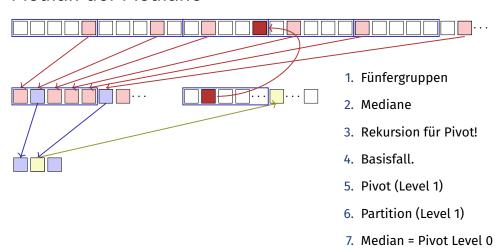
- Fünfergruppen bilden.
- Median jeder Gruppe bilden (naiv).
- Select rekursiv auf den Gruppenmedianen.
- Partitioniere das Array um den gefundenen Median der Mediane. Resultat: i
- Wenn i = k, Resultat. Sonst: Select rekursiv auf der richtigen Seite.

Algorithmus MMSelect(A, l, r, k)

Input: Array A der Länge n mit paarweise verschiedenen Einträgen. $1 \leq l \leq k \leq r \leq n, \ A[i] < A[k] \ \forall \ 1 \leq i < l, \ A[i] > A[k] \ \forall \ r < i \leq n$ Output: Wert $x \in A$ mit $|\{j|A[j] \leq x\}| = k$ $m \leftarrow \operatorname{MMChoose}(A,l,r)$ $i \leftarrow \operatorname{Partition}(A,l,r,m)$ if k < i then $| \operatorname{return \ MMSelect}(A,l,i-1,k)$ else if k > i then $| \operatorname{return \ MMSelect}(A,i+1,r,k)$ else $| \operatorname{return \ MMSelect}(A,i+1,r,k)$ else $| \operatorname{return \ A}[i]$

Median der Mediane

150

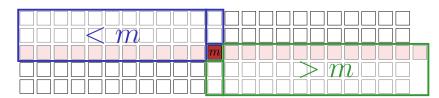


8. 2. Rekursion startet

Algorithmus $\mathtt{MMChoose}(A, l, r)$

```
\begin{array}{c} \textbf{Input:} \  \, \mathsf{Array} \ A \ \mathsf{der} \ \mathsf{L\"{a}nge} \ n \ \mathsf{mit} \ \mathsf{paarweise} \ \mathsf{verschiedenen} \ \mathsf{Eintr\"{a}gen}. \\ 1 \leq l \leq r \leq n. \\ \mathbf{Output:} \ \mathsf{Median} \ m \ \mathsf{der} \ \mathsf{Mediane} \\ \mathbf{if} \ r - l \leq 5 \ \mathbf{then} \\ \mid \ \mathsf{return} \ \mathsf{MedianOf5}(A[l,\ldots,r]) \\ \mathbf{else} \\ \mid \ A' \leftarrow \mathsf{MedianOf5Array}(A[l,\ldots,r]) \\ \mathbf{return} \ \mathsf{MMSelect}(A',1,|A'|,\left\lfloor \frac{|A'|}{2}\right\rfloor) \end{array}
```

Was bringt das?



- Anzahl Fünfergruppen: $\lceil \frac{n}{5} \rceil$, ohne Mediangruppe: $\lceil \frac{n}{5} \rceil 1$
- Minimale Anzahl Gruppen links / rechts von Mediangruppe $\left|\frac{1}{2}\left(\left[\frac{n}{5}\right]-1\right)\right|$
- lacksquare Minimale Anzahl Punkte kleiner / grösser als m

$$3\left\lfloor \frac{1}{2} \left(\left\lceil \frac{n}{5} \right\rceil - 1 \right) \right\rfloor \ge 3\left\lfloor \frac{1}{2} \left(\frac{n}{5} - 1 \right) \right\rfloor \ge 3\left(\frac{n}{10} - \frac{1}{2} - 1 \right) > \frac{3n}{10} - 6$$

(Fülle Restgruppe konzeptuell mit Punkten aus Mediangruppe auf)

 \Rightarrow Rekursiver Aufruf mit maximal $\lceil \frac{7n}{10} + 6 \rceil$ Elementen.

Analyse

Rekursionsungleichung:

$$T(n) \le T\left(\left\lceil \frac{n}{5}\right\rceil\right) + T\left(\left\lceil \frac{7n}{10} + 6\right\rceil\right) + d \cdot n.$$

mit einer Konstanten d. Behauptung:

$$T(n) = \mathcal{O}(n).$$

Beweis

Induktionsanfang: 7 Wähle c so gross, dass

$$T(n) \leq c \cdot n$$
 für alle $n \leq n_0$.

Induktionsannahme: H(n)

$$T(i) < c \cdot i$$
 für alle $i < n$.

Induktionsschritt: $H(k)_{k < n} \to H(n)$

$$\begin{split} T(n) &\leq T\bigg(\bigg\lceil\frac{n}{5}\bigg\rceil\bigg) + T\bigg(\bigg\lceil\frac{7n}{10} + 6\bigg\rceil\bigg) + d \cdot n \\ &\leq c \cdot \bigg\lceil\frac{n}{5}\bigg\rceil + c \cdot \bigg\lceil\frac{7n}{10} + 6\bigg\rceil + d \cdot n \qquad \text{(für } n > 20\text{)}. \end{split}$$

Beweis

Induktionsschritt:

$$T(n) \stackrel{n>20}{\leq} c \cdot \left\lceil \frac{n}{5} \right\rceil + c \cdot \left\lceil \frac{7n}{10} + 6 \right\rceil + d \cdot n$$

$$\leq c \cdot \frac{n}{5} + c + c \cdot \frac{7n}{10} + 6c + c + d \cdot n = \frac{9}{10} \cdot c \cdot n + 8c + d \cdot n.$$

Zu zeigen

$$\exists n_0, \exists c \mid \frac{9}{10} \cdot c \cdot n + 8c + d \cdot n \le cn \quad \forall n \ge n_0$$

Also

$$8c + d \cdot n \le \frac{1}{10}cn \quad \Leftrightarrow \quad n \ge \frac{80c}{c - 10d}$$

Setze z.B.
$$c = 90d, n_0 = 91$$
 $\Rightarrow T(n) \le cn \ \forall \ n \ge n_0$

 $^{^{7}}$ Es wird sich im Induktionsschritt herausstellen, dass der Basisfall für alle $n \leq n_0$ und ein bestimmtes (aber festes) $n_0 > 0$ betrachtet werden muss. Da c beliebig gross gewählt werden darf und eine beschränkte Anzahl von Termen eingeht, ist das eine einfache Erweiterung des Basisfalles n=1

Resultat

Theorem 11

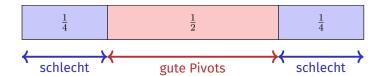
Das i-te Element einer Folge von n Elementen kann im schlechtesten Fall in $\Theta(n)$ Schritten gefunden werden.

5.1 Anhang

Herleitung einiger mathematischen Formeln

Überblick

1.	Wiederholt Minimum finden	$\mathcal{O}(n^2)$
2.	Sortieren und $A[i]$ ausgeben	$\mathcal{O}(n \log n)$
3.	Quickselect mit zufälligem Pivot	$\mathcal{O}(n)$ im Mittel
4.	Median of Medians (Blum)	$\mathcal{O}(n)$ im schlimmsten Fall



158

[Erwartungswert der geometrischen Verteilung]

Zufallsvariable $X \in \mathbb{N}^+$ mit $\mathbb{P}(X=k) = (1-p)^{k-1} \cdot p$. Erwartungswert

$$\mathbb{E}(X) = \sum_{k=1}^{\infty} k \cdot (1-p)^{k-1} \cdot p = \sum_{k=1}^{\infty} k \cdot q^{k-1} \cdot (1-q)$$

$$= \sum_{k=1}^{\infty} k \cdot q^{k-1} - k \cdot q^k = \sum_{k=0}^{\infty} (k+1) \cdot q^k - k \cdot q^k$$

$$= \sum_{k=0}^{\infty} q^k = \frac{1}{1-q} = \frac{1}{p}.$$

6. C++ vertieft (I)

Kurzwiederholung: Vektoren, Zeiger und Iteratoren Bereichsbasiertes for, Schlüsselwort auto, eine Klasse für Vektoren, Indexoperator, Move-Konstruktion, Iterator.

Was lernen wir heute?

- Schlüsselwort auto
- Bereichsbasiertes for
- Kurzwiederholung der Dreierregel
- Indexoperator
- Move Semantik, X-Werte und Fünferregel
- Eigene Iteratoren

2

Wir erinnern uns...

6.1 Nützliche Tools

Auf dem Weg zu elegantem, weniger komplizierten Code

auto

Das Schlüsselwort auto (ab C++11): Der Typ einer Variablen wird inferiert vom Initialisierer.

```
int x = 10;
auto y = x; // int
auto z = 3; // int
std::vector<double> v(5);
auto i = v[3]; // double
```

Schon etwas besser...

```
#include <iostream>
#include <vector>

int main(){
   std::vector<int> v(10); // Vector of length 10

for (int i = 0; i < v.size(); ++i)
   std::cin >> v[i];

for (auto it = v.begin(); it != v.end(); ++it) {
   std::cout << *it << " ";
}
}</pre>
```

6

Bereichsbasiertes for (C++11)

```
for (range-declaration : range-expression)
    statement;
```

- range-declaration: benannte Variable vom Elementtyp der durch range-expression spezifizierten Folge.
- range-expression: Ausdruck, der eine Folge von Elementen repräsentiert via Iterator-Paar begin(), end(), oder in Form einer Initialisierungsliste.

```
std::vector<double> v(5);
for (double x: v) std::cout << x; // 00000
for (int x: {1,2,5}) std::cout << x; // 125
for (double& x: v) x=5;</pre>
```

Cool!

```
#include <iostream>
#include <vector>
int main(){
  std::vector<int> v(10); // Vector of length 10

for (auto& x: v)
  std::cin >> x;

for (const auto x: v)
  std::cout << x << " ";
}</pre>
```

6.2 Speicherallokation

Bau einer Vektorklasse

Eine Klasse für (double) Vektoren

```
class Vector{
public:
    // constructors
    Vector(): sz{0}, elem{nullptr} {};
    Vector(std::size_t s): sz{s}, elem{new double[s]} {}
    // destructor
    ~Vector(){
        delete[] elem;
    }
    // (something is missing here)
private:
    std::size_t sz;
    double* elem;
}
```

Für unser genaues Verständis

Wir bauen selbst eine Vektorklasse, die so etwas kann!

Auf dem Weg lernen wir etwas über

- RAII (Resource Acquisition is Initialization) und Move-Konstruktion
- Index-Operatoren und andere Nützlichkeiten
- Templates
- Exception Handling
- Funktoren und Lambda-Ausdrücke

1

Elementzugriffe

172

```
class Vector{
    ...
    // getter. pre: 0 <= i < sz;
    double get(std::size_t i) const{
        return elem[i];
    }
    // setter. pre: 0 <= i < sz;
    void set(std::size_t i, double d){
        elem[i] = d;
    }
    // size property
    std::size_t size() const {
        return sz;
    }
}</pre>
```

```
class Vector{
public:
    Vector();
    Vector(std::size_t s);
    -Vector();
    double get(std::size_t i) const;
    void set(std::size_t i, double d);
    std::size_t size() const;
}
```

(Vector Schnittstelle)

}

Was läuft schief?

```
int main(){
                                                 class Vector{
                                                 public:
  Vector v(32):
                                                  Vector();
  for (std::size_t i = 0; i!=v.size(); ++i)
                                                  Vector(std::size_t s);
                                                  ~Vector();
   v.set(i, i);
                                                  double get(std::size_t i) const;
  Vector w = v;
                                                  void set(std::size t i, double d);
  for (std::size_t i = 0; i!=w.size(); ++i)
                                                  std::size_t size() const;
   w.set(i, i*i);
 return 0;
                                                       (Vector Schnittstelle)
*** Error in 'vector1': double free or corruption
(!prev): 0x000000000d23c20 ***
====== Backtrace: ======
/lib/x86 64-linux-gnu/libc.so.6(+0x777e5)[0x7fe5a5ac97e5]
```

Rule of Three!

```
class Vector{
...
  public:
  // copy constructor
  Vector(const Vector &v)
    : sz{v.sz}, elem{new double[v.sz]} {
    std::copy(v.elem, v.elem + v.sz, elem);
  }
}
```

```
class Vector{
public:
    Vector();
    Vector(std::size_t s);
    ~Vector();
    Vector(const Vector &v);
    double get(std::size_t i) const;
    void set(std::size_t i, double d);
    std::size_t size() const;
}
```

(Vector Schnittstelle)

Rule of Three!

```
class Vector{
                                                     class Vector(
                                                     public:
                                                      Vector();
  // assignment operator
                                                      Vector(std::size_t s);
  Vector& operator=(const Vector& v){
                                                      ~Vector():
                                                      Vector(const Vector &v);
    if (v.elem == elem) return *this;
                                                      Vector operator=(const Vector&v);
    if (elem != nullptr) delete[] elem;
                                                      double get(std::size_t i) const;
                                                      void set(std::size_t i, double d);
    sz = v.sz:
                                                      std::size_t size() const;
    elem = new double[sz];
    std::copy(v.elem, v.elem+v.sz, elem);
                                                           (Vector Schnittstelle)
   return *this;
}
```

letzt ist es zumindest korrekt. Aber umständlich.

Weiterleitung des Konstruktors

```
public:
// copy constructor
// (with constructor delegation)
Vector(const Vector &v): Vector(v.sz)
{
   std::copy(v.elem, v.elem + v.sz, elem);
}
```

176

Copy-&-Swap Idiom

```
class Vector{
 // Assignment operator
 Vector& operator= (const Vector&v){
   Vector cpy(v);
   swap(cpy);
                                copy-and-swap idiom: alle Felder von *this
   return *this:
                                tauschen mit den Daten von cpy. Beim Verlassen
                                von operator= wird cpy aufgeräumt (dekonstru-
private:
                                iert), während die Kopie der Daten von v in *this
 // helper function
                                verbleiben.
 void swap(Vector& v){
   std::swap(sz, v.sz);
   std::swap(elem, v.elem);
}
```

Arbeit an der Fassade.

Getter und Setter unschön. Wir wollen einen Indexoperator. Überladen! So?

```
class Vector{
 double operator[] (std::size_t pos) const{
   return elem[pos];
 void operator[] (std::size t pos, double value){
   elem[pos] = value;
                                                               Nein!
```

Referenztypen!

```
class Vector{
 // for non-const objects
 double& operator[] (std::size_t pos){
   return elem[pos]; // return by reference!
 // for const objects
 const double& operator[] (std::size t pos) const{
   return elem[pos];
}
```

Soweit, so gut.

```
int main(){
 Vector v(32); // constructor
 for (int i = 0; i<v.size(); ++i)</pre>
   v[i] = i; // subscript operator
 Vector w = v; // copy constructor
 for (int i = 0; i<w.size(); ++i)</pre>
   w[i] = i*i:
 const auto u = w;
 for (int i = 0; i<u.size(); ++i)</pre>
   std::cout << v[i] << ":" << u[i] << " "; // 0:0 1:1 2:4 ...
 return 0;
}
```

Bereichsbasiertes for

6.3 Iteratoren

Wie man bereichsbasiertes for unterstützt.

Wir wollten doch das:

182

```
Vector v = ...;
for (auto x: v)
  std::cout << x << " ";</pre>
```

Dafür müssen wir einen Iterator über begin und end bereitstellen.

Iterator für den Vektor

```
class Vector{
...
    // Iterator
    double* begin(){
        return elem;
    }
    double* end(){
        return elem+sz;
    }
}
```

Const Iterator für den Vektor

```
class Vector{
...
   // Const-Iterator
   const double* begin() const{
      return elem;
   }
   const double* end() const{
      return elem+sz;
   }
}
```

Zwischenstand

6.4 Effizientes Speicher-Management*

Wie man Kopien vermeidet

Vector Schnittstelle

```
class Vector{
public:
    Vector(); // Default Constructor
    Vector(std::size_t s); // Constructor
    ~Vector(); // Destructor
    Vector(const Vector &v); // Copy Constructor
    Vector& operator=(const Vector&v); // Assignment Operator
    double& operator[] (std::size_t pos); // Subscript operator (read/write)
    const double& operator[] (std::size_t pos) const; // Subscript operator
    std::size_t size() const;
    double* begin(); // iterator begin
    double* end(); // iterator end
    const double* begin() const; // const iterator begin
    const double* end() const; // const iterator end
}
```

Anzahl Kopien

Wie oft wird v kopiert?

Move-Konstruktor und Move-Zuweisung

```
class Vector{
...
    // move constructor
    Vector (Vector&& v): Vector() {
        swap(v);
    };
    // move assignment
    Vector& operator=(Vector&& v){
        swap(v);
        return *this;
    };
}
```

Erklärung

Wenn das Quellobjekt einer Zuweisung direkt nach der Zuweisung nicht weiter existiert, dann kann der Compiler den Move-Zuweisungsoperator anstelle des Zuweisungsoperators einsetzen.⁸ Damit wird eine potentiell teure Kopie vermieden.

Anzahl der Kopien im vorigen Beispiel reduziert sich zu 1.

Vector Schnittstelle

```
class Vector{
public:
    Vector();
    Vector(std::size_t s);
    ~Vector();
    Vector(const Vector &v);
    Vector& operator=(const Vector&v);
    Vector (Vector&& v);
    Vector& operator=(Vector&& v);
    const double& operator[] (std::size_t pos) const;
    double& operator[] (std::size_t pos);
    std::size_t size() const;
}
```

190

Illustration zur Move-Semantik

```
// nonsense implementation of a "vector" for demonstration purposes
class Vec{
public:
    Vec () {
        std::cout << "default constructor\n";}
    Vec (const Vec&) {
        std::cout << "copy constructor\n";}
    Vec& operator = (const Vec&) {
        std::cout << "copy assignment\n"; return *this;}
    ~Vec() {}
};</pre>
```

⁸Analoges gilt für den Kopier-Konstruktor und den Move-Konstruktor.

Wie viele Kopien?

```
Vec operator + (const Vec& a, const Vec& b){
                                                Ausgabe
   Vec tmp = a;
                                                default constructor
   // add b to tmp
                                                copy constructor
   return tmp;
                                                copy constructor
}
                                                copy constructor
                                                copy assignment
int main (){
   Vec f;
                                                4 Kopien des Vektors
   f = f + f + f + f;
}
```

Wie viele Kopien?

```
Vec operator + (const Vec& a, const Vec& b){
                                                Ausgabe
   Vec tmp = a;
                                                default constructor
   // add b to tmp
                                                copy constructor
   return tmp;
                                                copy constructor
}
                                                copy constructor
                                                move assignment
int main (){
   Vec f;
                                               3 Kopien des Vektors
   f = f + f + f + f:
}
```

Illustration der Move-Semantik

```
// nonsense implementation of a "vector" for demonstration purposes
class Vec{
public:
    Vec () { std::cout << "default constructor\n";}
    Vec (const Vec&) { std::cout << "copy constructor\n";}
    Vec& operator = (const Vec&) {
        std::cout << "copy assignment\n"; return *this;}
    ~Vec() {}
    // new: move constructor and assignment
    Vec (Vec&&) {
        std::cout << "move constructor\n";}
    Vec& operator = (Vec&&) {
        std::cout << "move assignment\n"; return *this;}
};</pre>
```

194

Wie viele Kopien?

```
Vec operator + (Vec a, const Vec& b){
                                        Ausgabe
   // add b to a
                                        default constructor
   return a:
                                        copy constructor
}
                                        move constructor
                                        move constructor
int main (){
                                        move constructor
   Vec f;
                                        move assignment
   f = f + f + f + f;
}
                                        1 Kopie des Vektors
```

Erklärung: Move-Semantik kommt zum Einsatz, wenn ein x-wert (expired) zugewiesen wird. R-Wert-Rückgaben von Funktionen sind x-Werte.

http://en.cppreference.com/w/cpp/language/value_category

Wie viele Kopien

```
void swap(Vec& a, Vec& b){
    Vec tmp = a;
    a=b;
    b=tmp;
}
int main (){
    Vec f;
    Vec g;
    swap(f,g);
}
```

Ausgabe default constructor default constructor copy constructor copy assignment copy assignment

3 Kopien des Vektors

std::swap & std::move

std::swap ist (mit Templates) genau wie oben gesehen implementiert std::move kann verwendet werden, um die Elemente eines Containers in einen anderen zu verschieben

```
std::move(va.begin(),va.end(),vb.begin())
```

X-Werte erzwingen

```
void swap(Vec& a, Vec& b){
   Vec tmp = std::move(a);
   a=std::move(b);
   b=std::move(tmp);
}
int main (){
   Vec f;
   Vec g;
   swap(f,g);
}
```

Ausgabe default constructor default constructor move constructor move assignment move assignment

0 Kopien des Vektors

Erklärung: Mit std::move kann man einen L-Wert Ausdruck zu einem X-Wert machen. Dann kommt wieder Move-Semantik zum Einsatz.

http://en.cppreference.com/w/cpp/utility/move

198

Heutige Zusammenfassung

- Benutze auto um Typen vom Initialisierer zu inferieren.
- X-Werte sind solche, bei denen der Compiler weiss, dass Sie ihre Gültigkeit verlieren.
- Benutze Move-Konstruktion, um X-Werte zu verschieben statt zu kopieren.
- Wenn man genau weiss, was man tut, kann man X-Werte auch erzwingen.
- Indexoperatoren können überladen werden. Zum Schreiben benutzt man Referenzen.
- Hinter bereichsbasiertem for wirkt ein Iterator.
- Iteration wird unterstützt, indem man einen Iterator nach Konvention der Standardbibliothek implementiert.

7. Sortieren I

Einfache Sortierverfahren

7.1 Einfaches Sortieren

Sortieren durch Auswahl, Sortieren durch Einfügen, Bubblesort [Ottman/Widmayer, Kap. 2.1, Cormen et al, Kap. 2.1, 2.2, Exercise 2.2-2, Problem 2-2

202

Problemstellung

Eingabe: Ein Array A=(A[1],...,A[n]) der Länge n. **Ausgabe:** Eine Permutation A' von A, die sortiert ist: $A'[i] \leq A'[j]$ für alle $1 \leq i \leq j \leq n$.

Algorithmus: IsSorted(A)

Beobachtung

IsSorted(A): "nicht sortiert", wenn A[i] > A[i+1] für ein i.

\Rightarrow Idee:

```
\begin{array}{c|c} \text{for } j \leftarrow 1 \text{ to } n-1 \text{ do} \\ & \text{if } A[j] > A[j+1] \text{ then} \\ & \quad \lfloor \text{ swap}(A[j], A[j+1]); \end{array}
```

Ausprobieren

- $5 \leftrightarrow 6$ 2 8 4 1 (j = 1)
- $\boxed{5} \quad \boxed{6} \longleftrightarrow \boxed{2} \quad \boxed{8} \quad \boxed{4} \quad \boxed{1} \quad (j=2)$
- [5] [2] $[6 \leftrightarrow 8]$ [4] [1] (j=3)
- [5] [2] [6] [8] \longleftrightarrow [4] (j=4)
- [5] [2] [6] [4] [8] \longleftrightarrow [1] (j=5)
- 5 2 6 4 1 8

- Nicht sortiert! (੨).
- Aber das grösste Element wandert ganz nach rechts.
 - ⇒ Neue Idee! ©

206

Ausprobieren

```
8
8
8
4
4
4
6
1
1
5
5
                                               (j = 1, i = 1)
5
5
5
5
2
2
2
2
2
2
2
2
2
       6
6
2
2
2
5
5
5
5
4
4
4
               2
6
6
6
                              4
4
8
1
1
1
6
6
                                      1
1
1
1
8
                                               (j=2)
                                               (j = 3)
                                               (j=4)
                                                               ■ Wende das Verfahren
                                               (j = 5)
                                               (j = 1, i = 2)
                                                                  iterativ an.
                                      8
               6
4
4
4
5
                                              (j=2)
                                              (j=3)
                                                               \blacksquare Für A[1,\ldots,n],
                                              (j = 4)
                                                                  dann A[1,\ldots,n-1],
                                      8 8
                                              (j = 1, i = 3)
                                                                  dann A[1,\ldots,n-2],
                                              (j=2)
                                               (j = 3)
                                                                   etc.
                              6
                                      8
                                               (j = 1, i = 4)
               1
                              6
                                      8
                                              (j = 2)
       1
                                      8
                                              (i = 1, j = 5)
```

Algorithmus: Bubblesort

Analyse

Anzahl Schlüsselvergleiche $\sum_{i=1}^{n-1}(n-i)=\frac{n(n-1)}{2}=\Theta(n^2)$. Anzahl Vertauschungen im schlechtesten Fall: $\Theta(n^2)$

Was ist der schlechteste Fall?

Wenn A absteigend sortiert ist.

Sortieren durch Auswahl

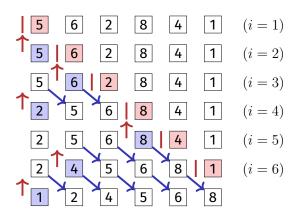
- 1 (i = 1)(i = 2)(i = 3)(i = 4)(i = 5)(i = 6)
- Auswahl des kleinsten Elementes durch Suche im unsortierten Teil A[i..n] des Arrays.
- Tausche kleinstes
 Element an das erste
 Element des unsortierten
 Teiles.
- Unsortierter Teil wird ein Element kleiner $(i \rightarrow i+1)$. Wiederhole bis alles sortiert. (i=n)

Algorithmus: Sortieren durch Auswahl

Analyse

Anzahl Vergleiche im schlechtesten Fall: $\Theta(n^2)$. Anzahl Vertauschungen im schlechtesten Fall: $n-1=\Theta(n)$

Sortieren durch Einfügen



Iteratives Vorgehen:

$$i = 1...n$$

- Einfügeposition für Element *i* bestimmen.
- Element i einfügen, ggfs. Verschiebung nötig.

Sortieren durch Einfügen

Welchen Nachteil hat der Algorithmus im Vergleich zum Sortieren durch Auswahl?

Im schlechtesten Fall viele Elementverschiebungen.

Welchen Vorteil hat der Algorithmus im Vergleich zum Sortieren durch Auswahl?

Der Suchbereich (Einfügebereich) ist bereits sortiert. Konsequenz: binäre Suche möglich.

214

Algorithmus: Sortieren durch Einfügen

Analyse

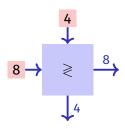
Anzahl Vergleiche im schlechtesten Fall:

 $\sum_{k=1}^{n-1} a \cdot \log k = a \log((n-1)!) \in \Theta(n \log n).$

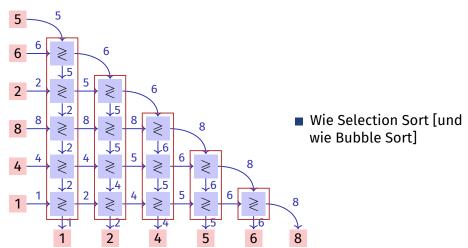
Anzahl Vertauschungen im schlechtesten Fall: $\sum_{k=2}^n (k-1) \in \Theta(n^2)$

Anderer Blickwinkel

Sortierknoten:

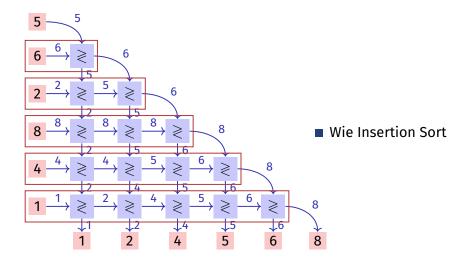


Anderer Blickwinkel



218

Anderer Blickwinkel



Schlussfolgerung

Selection Sort, Bubble Sort und Insertion Sort sind in gewissem Sinne dieselben Sortieralgorithmen. Wird später präzisiert. ⁹

221

⁹Im Teil über parallele Sortiernetzwerke. Für sequentiellen Code gelten natürlich weiterhin die zuvor gemachten Feststellungen.

Shellsort (Donald Shell 1959)

Intution: Verschieben weit entfernter Elemente dauert lange bei obigen naiven Verfahren

Insertion Sort auf Teilfolgen der Form $(A_{k\cdot i})$ ($i\in\mathbb{N}$) mit absteigenden Abständen k. Letzte Länge ist zwingend k=1.

Worst-case Performance hängt kritisch von den gewählten Teilfolgen ab. Beispiele:

- Ursprünglich mit Folge $1, 2, 4, 8, ..., 2^k$ konzipiert. Laufzeit: $\mathcal{O}(n^2)$
- Folge $1, 3, 7, 15, ..., 2^{k-1}$ (Hibbard 1963). $\mathcal{O}(n^{3/2})$
- Folge $1, 2, 3, 4, 6, 8, ..., 2^p 3^q$ (Pratt 1971). $\mathcal{O}(n \log^2 n)$

Shellsort

	0	1	2	3	4	5	6	7	8	9
insertion sort, $\boldsymbol{k}=7$	0	1	9	3	4	5	6	7	8	2
	0	8	9	3	4	5	6	7	1	2
	7	8	9	3	4	5	6	0	1	2
insertion sort, $k=3$	7	8	9	6	4	5	3	0	1	2
	7	8	9	6	4	5	3	0	1	2
	7	8	9	6	4	5	3	0	1	2
insertion sort, $k=1$	9	8	7	6	5	4	3	2	1	0

8. Sortieren II

Mergesort, Quicksort

8.1 Mergesort

222

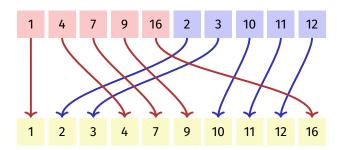
[Ottman/Widmayer, Kap. 2.4, Cormen et al, Kap. 2.3],

Mergesort (Sortieren durch Verschmelzen)

Divide and Conquer!

- Annahme: Zwei Hälften eines Arrays A bereits sortiert.
- Folgerung: Minimum von A kann mit 2 Vergleichen ermittelt werden.
- Iterativ: Füge die beiden vorsortierten Hälften von A zusammen in $\mathcal{O}(n)$.

Merge



226

Algorithmus Merge(A, l, m, r)

```
 \begin{array}{lll} \textbf{Input:} & \text{Array } A \text{ der L\"ange } n, \text{ Indizes } 1 \leq l \leq m \leq r \leq n. \\ & A[l,\ldots,m], \ A[m+1,\ldots,r] \text{ sortiert} \\ \textbf{Output:} & A[l,\ldots,r] \text{ sortiert} \\ 1 & B \leftarrow \text{new Array}(r-l+1) \\ 2 & i \leftarrow l; \ j \leftarrow m+1; \ k \leftarrow 1 \\ 3 & \textbf{while } i \leq m \text{ and } j \leq r \text{ do} \\ 4 & & \textbf{if } A[i] \leq A[j] \text{ then } B[k] \leftarrow A[i]; \ i \leftarrow i+1 \\ 5 & & \textbf{else } B[k] \leftarrow A[j]; \ j \leftarrow j+1 \\ 6 & & k \leftarrow k+1; \\ 7 & \textbf{while } i \leq m \text{ do } B[k] \leftarrow A[i]; \ i \leftarrow i+1; \ k \leftarrow k+1 \\ 8 & \textbf{while } j \leq r \text{ do } B[k] \leftarrow A[j]; \ j \leftarrow j+1; \ k \leftarrow k+1 \\ 9 & \textbf{for } k \leftarrow l \text{ to } r \text{ do } A[k] \leftarrow B[k-l+1] \\ \end{array}
```

Korrektheit

Hypothese: Nach k Durchläufen der Schleife von Zeile 3 ist $B[1,\ldots,k]$ sortiert und $B[k] \leq A[i]$, falls $i \leq m$ und $B[k] \leq A[j]$ falls $j \leq r$. Beweis per Induktion:

Induktions an fang: Das leere Array $B[1, \ldots, 0]$ ist trivialerweise sortiert. Induktions schluss $(k \to k+1)$:

- lacksquare oBdA $A[i] \leq A[j]$, $i \leq m, j \leq r$.
- B[1,...,k] ist nach Hypothese sortiert und $B[k] \leq A[i]$.
- Nach $B[k+1] \leftarrow A[i]$ ist $B[1, \ldots, k+1]$ sortiert.
- $B[k+1] = A[i] \le A[i+1]$ (falls $i+1 \le m$) und $B[k+1] \le A[j]$ falls $j \le r$.
- $k \leftarrow k+1, i \leftarrow i+1$: Aussage gilt erneut.

Analyse (Merge)

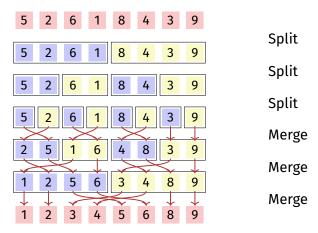
Lemma 12

Wenn: Array A der Länge n, Indizes $1 \le l < r \le n$. $m = \lfloor (l+r)/2 \rfloor$ und $A[l, \ldots, m]$, $A[m+1, \ldots, r]$ sortiert.

Dann: im Aufruf Merge(A, l, m, r) werden $\Theta(r - l)$ viele Schlüsselbewegungen und Vergleiche durchgeführt.

Beweis: (Inspektion des Algorithmus und Zählen der Operationen).

Mergesort



231

230

Algorithmus (Rekursives 2-Wege) Mergesort(A, l, r)

$\begin{array}{lll} \textbf{Input:} & \mathsf{Array} \ A \ \mathsf{der} \ \mathsf{L\"{a}nge} \ n. \ 1 \leq l \leq r \leq n \\ \textbf{Output:} & A[l, \dots, r] \ \mathsf{sortiert.} \\ & \textbf{if} \ l < r \ \textbf{then} \\ & & m \leftarrow \lfloor (l+r)/2 \rfloor \qquad // \ \mathsf{Mittlere} \ \mathsf{Position} \\ & & \mathsf{Mergesort}(A, l, m) \qquad // \ \mathsf{Sortiere} \ \mathsf{vordere} \ \mathsf{H\"{a}lfte} \\ & & \mathsf{Mergesort}(A, m+1, r) \ // \ \mathsf{Sortiere} \ \mathsf{hintere} \ \mathsf{H\"{a}lfte} \\ & & \mathsf{Merge}(A, l, m, r) \qquad // \ \mathsf{Verschmelzen} \ \mathsf{der} \ \mathsf{Teilfolgen} \\ \end{array}$

Analyse

Rekursionsgleichung für die Anzahl Vergleiche und Schlüsselbewegungen:

$$T(n) = T(\left\lceil \frac{n}{2} \right\rceil) + T(\left\lfloor \frac{n}{2} \right\rfloor) + \Theta(n) \in \Theta(n \log n)$$

Algorithmus StraightMergesort(A)

Rekursion vermeiden: Verschmelze Folgen der Länge 1, 2, 4... direkt

```
\begin{array}{lll} \textbf{Input:} & \mathsf{Array} \ A \ \mathsf{der} \ \mathsf{L\"{a}nge} \ n \\ \textbf{Output:} & \mathsf{Array} \ A \ \mathsf{sortiert} \\ \textit{length} \leftarrow 1 \\ \textbf{while} \ \textit{length} < n \ \textbf{do} \\ & r \leftarrow 0 \\ \textbf{while} \ r + \textit{length} < n \ \textbf{do} \\ & | l \leftarrow r + 1 \\ & m \leftarrow l + \textit{length} - 1 \\ & r \leftarrow \min(m + \textit{length}, n) \\ & | \mathsf{Merge}(A, l, m, r) \\ & | \textit{length} \leftarrow \textit{length} \cdot 2 \\ \end{array}
```

Analyse

Wie rekursives Mergesort führt reines 2-Wege-Mergesort immer $\Theta(n \log n)$ viele Schlüsselvergleiche und -bewegungen aus.

234

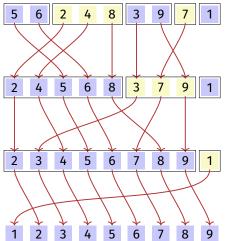
Natürliches 2-Wege Mergesort

Beobachtung: Obige Varianten nutzen nicht aus, wenn vorsortiert ist und führen immer $\Theta(n \log n)$ viele Bewegungen aus.

Wie kann man teilweise vorsortierte Folgen besser sortieren?

 \bigcirc Rekursives Verschmelzen von bereits vorsortierten Teilen (Runs) von A.

Natürliches 2-Wege Mergesort



Algorithmus NaturalMergesort(A)

```
\begin{array}{l} \textbf{Input:} & \text{Array } A \text{ der L\"ange } n > 0 \\ \textbf{Output:} & \text{Array } A \text{ sortiert} \\ \textbf{repeat} \\ \hline \\ r \leftarrow 0 \\ \textbf{while } r < n \text{ do} \\ \hline \\ l \leftarrow r+1 \\ m \leftarrow l; \textbf{ while } m < n \text{ and } A[m+1] \geq A[m] \text{ do } m \leftarrow m+1 \\ \textbf{if } m < n \text{ then} \\ \hline \\ r \leftarrow m+1; \textbf{ while } r < n \text{ and } A[r+1] \geq A[r] \text{ do } r \leftarrow r+1 \\ \hline \\ \textbf{Merge}(A,l,m,r); \\ \textbf{else} \\ \hline \\ r \leftarrow n \\ \hline \\ \textbf{until } l=1 \\ \hline \end{array}
```

8.2 Quicksort

[Ottman/Widmayer, Kap. 2.2, Cormen et al, Kap. 7]

Analyse

Ist es auch im Mittel asymptotisch besser als StraightMergesort?

f ONein. Unter Annahme der Gleichverteilung der paarweise unterschiedlichen Schlüssel haben wir im Mittel n/2 Stellen i mit $k_i > k_{i+1}$, also n/2 Runs und sparen uns lediglich einen Durchlauf, also n Vergleiche.

Natürliches Mergesort führt im schlechtesten und durchschnittlichen Fall $\Theta(n \log n)$ viele Vergleiche und Bewegungen aus.

Quicksort

238

Was ist der Nachteil von Mergesort?

Benötigt zusätzlich $\Theta(n)$ Speicherplatz für das Verschmelzen.

Wie könnte man das Verschmelzen einsparen?

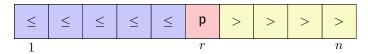
Sorge dafür, dass jedes Element im linken Teil kleiner ist als im rechten Teil.

Wie?

Pivotieren und Aufteilen!

Pivotieren

- 1. Wähle ein (beliebiges) Element p als Pivotelement
- 2. Teile A in zwei Teile auf: einen Teil L der Elemente mit $A[i] \leq p$ und einen Teil R der Elemente mit A[i] > p.
- 3. Quicksort: Rekursion auf Teilen L und R



Algorithmus Partition(A, l, r, p)

Input: Array A, welches den Pivot p in $A[l,\ldots,r]$ mindestens einmal enthält. **Output:** Array A partitioniert in $A[l,\ldots,r]$ um p. Rückgabe der Position von p. while l < r do

```
\begin{aligned} & \textbf{while} \ A[l]  <math display="block">& \textbf{while} \ A[r] > p \ \textbf{do} \\ & \  \  \, \bigsqcup \ r \leftarrow r-1  & \textbf{swap}(A[l], \ A[r]) \\ & \textbf{if} \ A[l] = A[r] \ \textbf{then} \\ & \  \  \, \bigsqcup \ l \leftarrow l+1 \end{aligned}
```

return |-1

242

Algorithmus Quicksort(A, l, r)

```
\begin{array}{ll} \textbf{Input:} & \text{Array } A \text{ der L\"ange } n. \ 1 \leq l \leq r \leq n. \\ \textbf{Output:} & \text{Array } A, \text{ sortiert in } A[l, \ldots, r]. \\ \textbf{if } l < r \text{ then} \\ & \text{W\"ahle Pivot } p \in A[l, \ldots, r] \\ & k \leftarrow \texttt{Partition}(A, l, r, p) \\ & \text{Quicksort}(A, l, k - 1) \\ & \text{Quicksort}(A, k + 1, r) \end{array}
```

Quicksort (willkürlicher Pivot)

2 4 5 6 8 3 7 9 1

2 1 3 6 8 5 7 9 4

1 2 3 4 5 8 7 9 6

1 2 3 4 5 6 7 9 8

1 2 3 4 5 6 7 8 9

1 2 3 4 5 6 7 8 9

Analyse: Anzahl Vergleiche

Schlechtester Fall. Pivotelement = Minimum oder Maximum; Anzahl Vergleiche:

$$T(n) = T(n-1) + c \cdot n, \ T(1) = d \implies T(n) \in \Theta(n^2)$$

Analyse: Anzahl Vertauschungen

Resultat eines Aufrufes an Partition (Pivot 3):

2 1 3 6 8 5 7 9 4

- ? Wie viele Vertauschungen haben hier maximal stattgefunden?
- ① 2. Die maximale Anzahl an Vertauschungen ist gegeben durch die Anzahl Schlüssel im kleineren Bereich.

16

Analyse: Anzahl Vertauschungen

Gedankenspiel

- Jeder Schlüssel aus dem kleineren Bereich zahlt bei einer Vertauschung eine Münze.
- Wenn ein Schlüssel eine Münze gezahlt hat, ist der Bereich, in dem er sich befindet maximal halb so gross wie zuvor.
- Jeder Schlüssel muss also maximal $\log n$ Münzen zahlen. Es gibt aber nur n Schlüssel.

Folgerung: Es ergeben sich $\mathcal{O}(n\log n)$ viele Schlüsselvertauschungen im schlechtesten Fall!

Randomisiertes Quicksort

Quicksort wird trotz $\Theta(n^2)$ Laufzeit im schlechtesten Fall oft eingesetzt. Grund: Quadratische Laufzeit unwahrscheinlich, sofern die Wahl des Pivots und die Vorsortierung nicht eine ungünstige Konstellation aufweisen. Vermeidung: Zufälliges Ziehen eines Pivots. Mit gleicher Wahrscheinlichkeit aus [l,r].

 2^{i}

Analyse (Randomisiertes Quicksort)

Erwartete Anzahl verglichener Schlüssel bei Eingabe der Länge n:

$$T(n) = (n-1) + \frac{1}{n} \sum_{k=1}^{n} (T(k-1) + T(n-k)), \ T(0) = T(1) = 0$$

Behauptung $T(n) \leq 4n \log n$.

Beweis per Induktion:

Induktionsanfang: klar für n=0 (mit $0 \log 0 := 0$) und für n=1.

Hypothese: $T(n) \le 4n \log n$ für ein n.

Induktionsschritt: $(n-1 \rightarrow n)$

Analyse (Randomisiertes Quicksort)

$$\begin{split} T(n) &= n - 1 + \frac{2}{n} \sum_{k=0}^{n-1} T(k) \overset{\mathsf{H}}{\leq} n - 1 + \frac{2}{n} \sum_{k=0}^{n-1} 4k \log k \\ &= n - 1 + \sum_{k=1}^{n/2} 4k \underbrace{\log k}_{\leq \log n - 1} + \sum_{k=n/2+1}^{n-1} 4k \underbrace{\log k}_{\leq \log n} \\ &\leq n - 1 + \frac{8}{n} \left((\log n - 1) \sum_{k=1}^{n/2} k + \log n \sum_{k=n/2+1}^{n-1} k \right) \\ &= n - 1 + \frac{8}{n} \left((\log n) \cdot \frac{n(n-1)}{2} - \frac{n}{4} \left(\frac{n}{2} + 1 \right) \right) \\ &= 4n \log n - 4 \log n - 3 \leq 4n \log n \end{split}$$

251

250

Analyse (Randomisiertes Quicksort)

Theorem 13

Im Mittel benötigt randomisiertes Quicksort $\mathcal{O}(n \log n)$ Vergleiche.

Praktische Anmerkungen

Rekursionstiefe im schlechtesten Fall: $n-1^{10}$. Dann auch Speicherplatzbedarf $\mathcal{O}(n)$.

Kann vermieden werden: Rekursion nur auf dem kleineren Teil. Dann garantiert $\mathcal{O}(\log n)$ Rekursionstiefe und Speicherplatzbedarf.

¹⁰Stack-Overflow möglich!

Quicksort mit logarithmischem Speicherplatz

```
\begin{array}{ll} \textbf{Input:} & \text{Array } A \text{ der L\"ange } n. \ 1 \leq l \leq r \leq n. \\ \textbf{Output:} & \text{Array } A \text{, sortiert zwischen } l \text{ und } r. \\ \textbf{while } l < r \text{ do} \\ & \text{W\"ahle Pivot } p \in A[l, \ldots, r] \\ & k \leftarrow \text{Partition}(A, l, r, p) \\ & \textbf{if } k - l < r - k \textbf{ then} \\ & \text{Quicksort}(A[l, \ldots, k-1]) \\ & l \leftarrow k + 1 \\ & \textbf{else} \\ & \text{Quicksort}(A[k+1, \ldots, r]) \\ & r \leftarrow k - 1 \end{array}
```

Der im ursprünglichen Algorithmus verbleibende Aufruf an Quicksort $(A[l,\ldots,r])$ geschieht iterativ (Tail Recursion ausgenutzt!): die If-Anweisung wurde zur While Anweisung.

8.3 Anhang

Herleitung einiger mathematischen Formeln

Praktische Anmerkungen

- Für den Pivot wird in der Praxis oft der Median von drei Elementen genommen. Beispiel: Median3(A[l], A[r], A[|l+r/2|]).
- Es existiert eine Variante von Quicksort mit konstanten Speicherplatzbedarf. Idee: Zwischenspeichern des alten Pivots am Ort des neuen Pivots.
- Komplizierte Divide-And-Conquer-Algorithmen verwenden oft als Basisfall einen trivialen ($\Theta(n^2)$) Algorithmus für kleine Problemgrössen.

 $\log n! \in \Theta(n \log n)$

254

$$\log n! = \sum_{i=1}^{n} \log i \le \sum_{i=1}^{n} \log n = n \log n$$

$$\sum_{i=1}^{n} \log i = \sum_{i=1}^{\lfloor n/2 \rfloor} \log i + \sum_{\lfloor n/2 \rfloor + 1}^{n} \log i$$

$$\ge \sum_{i=2}^{\lfloor n/2 \rfloor} \log 2 + \sum_{\lfloor n/2 \rfloor + 1}^{n} \log \frac{n}{2}$$

$$= (\underbrace{\lfloor n/2 \rfloor}_{>n/2 - 1} - 2 + 1) + (\underbrace{n - \lfloor n/2 \rfloor}_{\ge n/2})(\log n - 1)$$

$$> \frac{n}{2} \log n - 2.$$

255

$[n! \in o(n^n)]$

$$n\log n \ge \sum_{i=1}^{\lfloor n/2\rfloor} \log 2i + \sum_{i=\lfloor n/2\rfloor+1}^{n} \log i$$

$$= \sum_{i=1}^{n} \log i + \left\lfloor \frac{n}{2} \right\rfloor \log 2$$

$$> \sum_{i=1}^{n} \log i + n/2 - 1 = \log n! + n/2 - 1$$

$$\begin{split} n^n &= 2^{n \log_2 n} \geq 2^{\log_2 n!} \cdot 2^{n/2} \cdot 2^{-1} = n! \cdot 2^{n/2 - 1} \\ &\Rightarrow \frac{n!}{n^n} \leq 2^{-n/2 + 1} \overset{n \to \infty}{\longrightarrow} 0 \Rightarrow n! \in o(n^n) = \mathcal{O}(n^n) \backslash \Omega(n^n) \end{split}$$

[Sogar $n! \in o((n/c)^n) \forall 0 < c < e$]

Konvergenz oder Divergenz von $f_n = \frac{n!}{(n/c)^n}$. Quotientenkriterium

$$\frac{f_{n+1}}{f_n} = \frac{(n+1)!}{\left(\frac{n+1}{c}\right)^{n+1}} \cdot \frac{\left(\frac{n}{c}\right)^n}{n!} = c \cdot \left(\frac{n}{n+1}\right)^n \longrightarrow c \cdot \frac{1}{e} \lessgtr 1 \text{ wenn } c \lessgtr e$$

 $\begin{array}{l} \operatorname{denn}\left(1+\frac{1}{n}\right)^n \to e. \text{ Sogar die Reihe } \sum_{i=1}^n f_n \text{ konvergiert / divergiert für } c \leqslant e. \\ f_n \text{ divergiert für } c = e \text{, denn (Stirling): } n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n. \end{array}$

258

[Quotientenkriterium]

Quotientenkriterium für eine Folge $(f_n)_{n\in\mathbb{N}}$: Wenn $\xrightarrow{f_{n+1}} \xrightarrow[n\to\infty]{} \lambda$, dann sind die Folge f_n und auch die Reihe $\sum_{i=1}^n f_i$

- **•** konvergent, falls $\lambda < 1$ und
- divergent, falls $\lambda > 1$.

[Quotientenkriterium Herleitung]

Quotientenkriterium ergibt sich aus: Geometrische Reihe

$$S_n(r) := \sum_{i=0}^n r^i = \frac{1 - r^{n+1}}{1 - r}.$$

konvergiert für $n \to \infty$ genau dann wenn -1 < r < 1. Sei nämlich $0 \le \lambda < 1$:

$$\forall \varepsilon > 0 \,\exists n_0 : f_{n+1}/f_n < \lambda + \varepsilon \,\forall n \ge n_0$$

$$\Rightarrow \exists \varepsilon > 0, \exists n_0 : f_{n+1}/f_n \le \mu < 1 \,\forall n \ge n_0$$

Somit

260

$$\sum_{n=n_0}^{\infty} f_n \leq f_{n_0} \cdot \sum_{n=n_0}^{\infty} \cdot \mu^{n-n_0} \quad \text{konvergiert}.$$

(Analog für Divergenz)

,3

9. Sortieren III

Untere Schranken für das vergleichsbasierte Sortieren, Radix- und Bucketsort

9.1 Untere Grenzen für Vergleichbasiertes Sortieren

[Ottman/Widmayer, Kap. 2.8, Cormen et al, Kap. 8.1]

262

Untere Schranke für das Sortieren

Bis hierher: Sortieren im schlechtesten Fall benötigt $\Omega(n \log n)$ Schritte. Geht es besser? Nein:

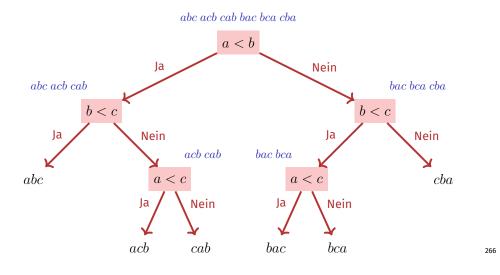
Theorem 14

Vergleichsbasierte Sortierverfahren benötigen im schlechtesten Fall und im Mittel mindestens $\Omega(n\log n)$ Schlüsselvergleiche.

Vergleichsbasiertes Sortieren

- Algorithmus muss unter n! vielen Anordnungsmöglichkeiten einer Folge $(A_i)_{i=1,\dots,n}$ die richtige identifizieren.
- Zu Beginn weiss der Algorithmus nichts.
- Betrachten den "Wissensgewinn" des Algorithmus als Entscheidungsbaum:
 - Knoten enthalten verbleibende Möglichkeiten
 - Kanten enthalten Entscheidungen

Entscheidungsbaum

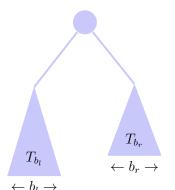


Entscheidungsbaum

Ein binärer Baum mit L Blättern hat K=L-1 innere Knoten. Die Höhe eines binären Baumes mit L Blättern ist mindestens $\log_2 L. \Rightarrow$ Höhe des Entscheidungsbaumes $h \geq \log n! \in \Omega(n \log n)$. Somit auch die Länge des längsten Pfades im Entscheidungsbaum $\in \Omega(n \log n)$.

Bleibt zu zeigen: mittlere Länge M(n) eines Pfades $M(n) \in \Omega(n \log n)$.

Untere Schranke im Mittel



- Entscheidungsbaum T_n mit n Blättern, mittlere Tiefe eines Blatts $m(T_n)$
- Annahme: $m(T_n) \ge \log n$ nicht für alle n.
- Wähle kleinstes b mit $m(T_b) < \log b \Rightarrow b \geq 2$
- $b_l + b_r = b \text{ mit } b_l > 0 \text{ und } b_r > 0 \Rightarrow b_l < b, b_r < b$ ⇒ $m(T_{b_l}) > \log b_l \text{ und } m(T_{b_n}) > \log b_r$

Untere Schranke im Mittel

Mittlere Tiefe eines Blatts:

$$m(T_b) = \frac{b_l}{b}(m(T_{b_l}) + 1) + \frac{b_r}{b}(m(T_{b_r}) + 1)$$

$$\geq \frac{1}{b}(b_l(\log b_l + 1) + b_r(\log b_r + 1)) = \frac{1}{b}(b_l \log 2b_l + b_r \log 2b_r)$$

$$\geq \frac{1}{b}(b \log b) = \log b.$$

Widerspruch.

Die letzte Ungleichung gilt, da $f(x) = x \log x$ konvex ist (f''(x) = 1/x > 0) und für eine konvexe Funktion gilt $f((x+y)/2) \le 1/2f(x) + 1/2f(y)$ ($x = 2b_l$, $y = 2b_r$ einsetzen). Te Einsetzen von $x = 2b_l$, $y = 2b_r$, und $b_l + b_r = b$.

 $[\]overline{\ \ }^{11}$ Beweis: starte mit leerem Baumm, K=0, L=1. Jeder hinzugefügte Knoten ersetzt ein Blatt durch 2 Blätter. Also.

¹²allgemein $f(\lambda x + (1 - \lambda)y) \le \lambda f(x) + (1 - \lambda)f(y)$ für $0 \le \lambda \le 1$.

Radix Sort

9.2 Radixsort und Bucketsort

Radixsort, Bucketsort [Ottman/Widmayer, Kap. 2.5, Cormen et al, Kap. 8.3]

Vergleichsbasierte Sortierverfahren: Schlüssel vergleichbar (< oder >, =). Ansonsten keine Voraussetzung.

Andere Idee: nutze mehr Information über die Zusammensetzung der Schlüssel.

270

Annahmen

Annahme: Schlüssel darstellbar als Wörter aus einem Alphabet mit m Elementen.

Beispiele

•		
m = 2 $m = 16$	Dezimalzahlen Dualzahlen Hexadezimalzahlen Wörter	$183 = 183_{10}$ 101_2 $A0_{16}$ "INFORMATIK"

m heisst die Wurzel (lateinisch $\it Radix$) der Darstellung.

Annahmen

- \blacksquare Schlüssel =m-adische Zahlen mit gleicher Länge.
- \blacksquare Verfahren z zur Extraktion der k-ten Ziffer eines Schlüssels in $\mathcal{O}(1)$ Schritten.

Beispiel

```
z_{10}(0, 85) = 5

z_{10}(1, 85) = 8

z_{10}(2, 85) = 0
```

Radix-Exchange-Sort

Schlüssel mit Radix 2. Beobachtung: Wenn für ein k > 0:

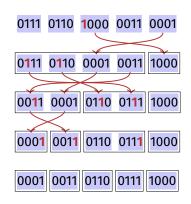
$$z_2(i,x) = z_2(i,y)$$
 für alle $i > k$

und

$$z_2(k,x) < z_2(k,y),$$

dann ist x < y.

Radix-Exchange-Sort



Radix-Exchange-Sort

Idee:

- Starte mit maximalem k.
- Binäres Aufteilen der Datensätze mit $z_2(k,\cdot)=0$ vs. $z_2(k,\cdot)=1$ wie bei Quicksort.
- $k \leftarrow k-1$.

274

Algorithmus RadixExchangeSort(A, l, r, b)

```
Input: Array A der Länge n, linke und rechte Grenze 1 \leq l \leq r \leq n, Bitposition b

Output: Array A, im Bereich [l,r] nach Bits [0,\ldots,b] sortiert. if l < r and b \geq 0 then  \begin{vmatrix} i \leftarrow l - 1 \\ j \leftarrow r + 1 \end{vmatrix}  repeat  \begin{vmatrix} \text{repeat } i \leftarrow i + 1 \text{ until } z_2(b,A[i]) = 1 \text{ or } i \geq j \\ \text{repeat } j \leftarrow j - 1 \text{ until } z_2(b,A[j]) = 0 \text{ or } i \geq j \\ \text{if } i < j \text{ then swap}(A[i],A[j]) \end{vmatrix}  until i \geq j
```

76 277

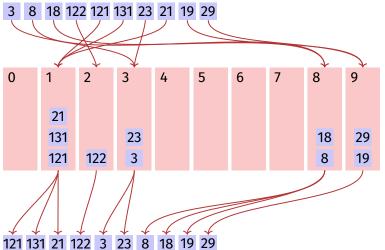
RadixExchangeSort(A, l, i - 1, b - 1)RadixExchangeSort(A, i, r, b - 1)

Analyse

RadixExchangeSort ist rekursiv mit maximaler Rekursionstiefe = maximaler Anzahl Ziffern p.

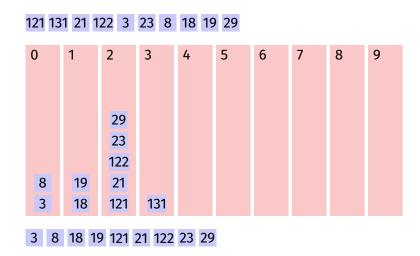
Laufzeit im schlechtesten Fall $\mathcal{O}(p \cdot n)$.

Bucket Sort (Sortieren durch Fachverteilen)

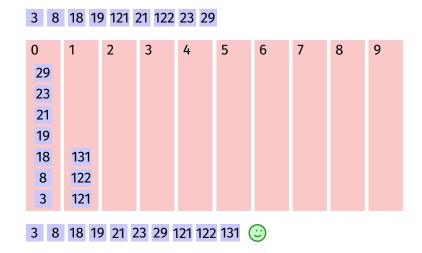


278

Bucket Sort (Sortieren durch Fachverteilen)



Bucket Sort (Sortieren durch Fachverteilen)



Implementations details

Bucketgrösse sehr unterschiedlich. Möglichkeiten

- Verkettete Liste oder dynamisches Array für jede Ziffer.
- Ein Array der Länge *n*, Offsets für jede Ziffer in erstem Durchlauf bestimmen.

Annahmen: Eingabelänge n , Anzahl Bits / Ganzzahl: k , Anzahl Buckets: 2^b

Asymptotische Laufzeit $\mathcal{O}(\frac{k}{h} \cdot (n+2^b)$.

Zum Beispiel: k = 32, $2^b = 256$: $\frac{k}{b} \cdot (n+2^b) = 4n + 1024$.

Bucket Sort - Andere Voraussetzung

```
Annahme: gleichmässig verteilte Daten, z.B. aus [0,1)
```

 $\textbf{Input} \colon \quad \text{Array A der L\"ange n, $A_i \in [0,1)$, Konstante $M \in \mathbb{N}^+$}$

Output: Sortiertes Array

 $k \leftarrow \lceil n/M \rceil$

 $B \leftarrow \text{new array of } k \text{ empty lists}$

for $i \leftarrow 1$ to n do

 $B[\lfloor A_i \cdot k \rfloor]$.append(A[i])

for $i \leftarrow 1$ to k do

sort B[i] // z.B. insertion sort, mit Laufzeit $\mathcal{O}(M^2)$

return $B[0] \circ B[1] \circ \cdots \circ B[k]$ // konkateniert

Erwartete asymptotische Laufzeit $\mathcal{O}(n)$ (Beweis in Cormen et al, Kap. 8.4)

282

Was lernen wir heute?

10. C++ vertieft (II): Templates

- Templates von Klassen
- Funktionentemplates
- Smart Pointers

Motivation

Ziel: generische Vektor-Klasse und Funktionalität.

```
Vector<double> vd(10);
Vector<int> vi(10);
Vector<char> vi(20);
auto nd = vd * vd; // norm (vector of double)
auto ni = vi * vi; // norm (vector of int)
```

Typen als Template Parameter

- Ersetze in der konkreten Implementation einer Klasse den Typ, der generisch werden soll (beim Vektor: double) durch einen Stellvertreter, z.B. T.
- 2. Stelle der Klasse das Konstrukt template<typename T> voran (ersetze T ggfs. durch den Stellvertreter)..

Das Konstrukt template<typename T> kann gelesen werden als "für alle Typen T".

86

Typen als Template Parameter

```
template <typename ElementType>
class Vector{
    std::size_t size;
    ElementType* elem;
public:
        ...
    Vector(std::size_t s):
        size{s},
        elem{new ElementType[s]}{}
    ...
    ElementType& operator[](std::size_t pos){
        return elem[pos];
    }
    ...
}
```

Template Instanzierung

Vector<typeName> erzeugt Typinstanz von Vector mit ElementType=typeName.
Bezeichnung: Instanzierung.

```
Vector<double> x;  // vector of double
Vector<int> y;  // vector of int
```

Vector<Vector<double>> x; // vector of vector of double

Type-checking

Templates sind weitgehend Ersetzungsregeln zur Instanzierungszeit und während der Kompilation. Es wird immer so wenig geprüft wie nötig und so viel wie möglich.

Beispiel

```
template <typename T>
class Pair{
    T left; T right;
public:
    Pair(T 1, T r):left{1}, right{r}{}
    T min(){
       return left < right ? left : right;
    }
};

Pair<int> a(10,20); // ok
auto m = a.min(); // ok
Pair<Pair<int>> b(a,Pair<int>(20,30)); // ok
auto n = b.min(); // no match for operator<!</pre>
```

290

291

Generische Programmierung

Generische Komponenten sollten eher als Generalisierung eines oder mehrerer Beispiele entwickelt werden als durch Ableitung von Grundprinzipien.

```
template <typename T>
class Vector{
public:
 Vector();
 Vector(std::size_t);
 ~Vector();
 Vector(const Vector&):
 Vector& operator=(const Vector&);
 Vector (Vector&&);
 Vector& operator=(Vector&&);
 const T& operator[] (std::size_t) const;
 T& operator[] (std::size_t);
 std::size t size() const;
 T* begin();
 T* end();
 const T* begin() const;
 const T* end() const;
```

Funktionentemplates

- 1. Ersetze in der konkreten Implementation einer Funktion den Typ, der generisch werden soll durch einen Namen, z.B. T,
- 2. Stelle der Funktion das Konstrukt template<typename T> voran (ersetze T ggfs. durch den gewählten Namen).

Funktionentemplates

```
template <typename T>
void swap(T& x, T&y){
   T temp = x;
   x = y;
   y = temp;
}
```

Typen der Aufrufparameter determinieren die Version der Funktion, welche (kompiliert und) verwendet wird:

```
int x=5;
int y=6;
swap(x,y); // calls swap with T=int
```

.. auch mit Operatoren

```
template <typename T>
                             Pair<int> a(10,20); // ok
class Pair{
                             std::cout << a; // ok
   T left; T right;
public:
   Pair(T 1, T r):left{1}, right{r}{}
   T min(){ return left < right? left: right; }</pre>
   std::ostream& print (std::ostream& os) const{
       return os << "("<< left << "," << right<< ")";</pre>
   }
};
template <typename T>
std::ostream& operator<< (std::ostream& os, const Pair<T>& pair){
   return pair.print(os);
}
```

Sicherheiten

```
template <typename T>
void swap(T& x, T&y){
   T temp = x;
   x = y;
   y = temp;
}
```

Eine unverträgliche Version der Funktion wird nicht erzeugt:

```
int x=5;
double y=6;
swap(x,y); // error: no matching function for ...
```

294

Praktisch!

```
// Output of an arbitrary container
template <typename T>
void output(const T& t){
   for (auto x: t)
      std::cout << x << " ";
   std::cout << "\n";
}
int main(){
   std::vector<int> v={1,2,3};
   output(v); // 1 2 3
}
```

Explizite Typangabe

```
// input of an arbitrary pair
template <typename T>
Pair<T> read(){
   T left;
   T right;
   std::cin << left << right;
   return Pair<T>(left,right);
}
...
auto p = read<double>();
```

Wenn der Typ bei der Instanzierung nicht inferiert werden kann, muss er explizit angegeben werden.

Mächtig!

```
template <typename T> // square number
T sq(T x){
    return x*x;
}
template <typename Container, typename F>
void apply(Container& c, F f){ // x <- f(x) forall x in c
    for(auto& x: c)
        x = f(x);
}
int main(){
    std::vector<int> v={1,2,3};
    apply(v,sq<int>);
    output(v); // 1 4 9
}
```

298

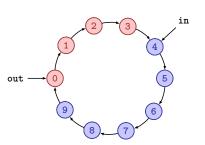
Spezialisierung

```
template <>
class Pair<bool>{
    short both;
public:
    Pair(bool 1, bool r):both{(1?1:0) + (r?2:0)} {};
    std::ostream& print (std::ostream& os) const{
        return os << "("<< both % 2 << "," << both /2 << ")";
    }
};

Pair<int> i(10,20); // ok -- generic template
    std::cout << i << std::endl; // (10,20);
    Pair<bool> b(true, false); // ok -- special bool version
    std::cout << b << std::endl; // (1,0)</pre>
```

Templateparametrisierung mit Werten

```
template <typename T, int size>
class CircularBuffer{
  T buf[size];
  int in; int out;
public:
    CircularBuffer():in{0},out{0}{};
  bool empty(){
    return in == out;
  }
  bool full(){
    return (in + 1) % size == out;
  }
  void put(T x); // declaration
  T get(); // declaration
};
```



o **}**;

Templateparametrisierung mit Werten

Nochmal Speichermanagement

Richtlinie "Dynamischer Speicher"

Zu jedem new gibt es ein passendes delete!

Vermeide:

- Speicherlecks: "alte" Objekte, die den Speicher blockieren
- Zeiger auf freigegebene Objekte: hängende Zeiger (dangling pointers)
- Mehrfache Freigabe eines Objektes mit delete.

Wie?

302

Smart Pointers

- Können sicherstellen, dass ein Objekt gelöscht wird genau dann, wenn es nicht mehr genutzt wird
- Basieren auf dem RAII (Resource Acquisition is Initialization) Paradigma.
- Können an die Stelle jedes gewöhnlichen Pointers treten: sind als Klassentemplates implementiert.
- Es gibt std::unique_ptr<>, std::shared_ptr<> (und std::weak_ptr<>)

```
std::unique_ptr<Node> nodeU(new Node()); // unique pointer
std::shared_ptr<Node> nodeS(new Node()); // shared pointer
```

Unique Pointer

- Der Dekonstruktor von std::unique_ptr<T> löscht den enthaltenen Zeiger.
- \blacksquare std::unique_ptr<T> hat exklusiv Zugriff auf den enthaltenen Zeiger auf T.
- Kopierkonstruktor und Assignment Operator sind gelöscht. Ein Unique Pointer kann nicht als Wert kopiert werden. Movekonstruktor ist vorhanden: der Zeiger kann verschoben werden.
- Kein Zusatzaufwand zur Laufzeit im Vergleich zu einem normalen Zeiger.

```
std::unique_ptr<Node> nodeU(new Node()); // unique pointer
std::unique_ptr<Node> node2 = std::move(nodeU); // ok
std::unique_ptr<Node> node3 = nodeU; // error
```

Shared Pointer

- std::shared_ptr<T> zählt die Anzahl von Besitzern eines Zeigers (Referenzzähler). Wenn der Referenzzähler auf 0 fällt, wird der Zeiger gelöscht.
- Shared Pointers können kopiert werden.
- Shared Pointers haben zusätzlichen Speicher- und Laufzeitbedarf: sie verwalten den Referenzzähler zur Laufzeit und enthalten jeweils einen Zeiger auf den Referenzzähler.

RefCount (2)

Std::shared_ptr<Node>

Shared Pointer

```
std::shared_ptr<Node> nodeS(new Node()); // shared pointer, rc = 1
std::shared_ptr<Node> node2 = std::move(nodeS); // ok, rc unchanged
std::shared_ptr<Node> node3 = node2; // ok, rc = 2
```

306

Smart Pointers

Einige Regeln

- Niemals delete auf einen Zeiger im Smart Pointer aufrufen.
- new vermeiden, stattdessen:

```
std::unique_ptr<Node> nodeU = std::make_unique<Node>()
std::shared_ptr<Node> nodeS = std::make_shared<Node>()
```

- Wo möglich, std::unique_ptr verwenden.
- Bei der Verwendung von std::shared_ptr sicherstellen, dass es keine Zyklen im Zeigergraphen gibt.

11. Elementare Datenstrukturen

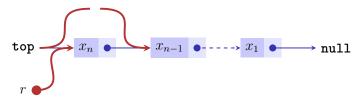
Abstrakte Datentypen Stapel, Warteschlange, Implementationsvarianten der verketteten Liste [Ottman/Widmayer, Kap. 1.5.1-1.5.2, Cormen et al, Kap. 10.1.-10.2]

Abstrakte Datentypen

Wir erinnern uns¹³ (Vorlesung Informatik I) Ein Stack ist ein abstrakter Datentyp (ADT) mit Operationen

- **push**(x, S): Legt Element x auf den Stapel S.
- \blacksquare pop(S): Entfernt und liefert oberstes Element von S, oder null.
- \blacksquare top(S): Liefert oberstes Element von S, oder null.
- \blacksquare is Empty(S): Liefert true wenn Stack leer, sonst false.
- emptyStack(): Liefert einen leeren Stack.

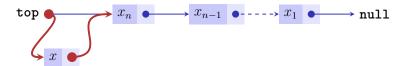
Implementation Pop



pop(S):

- 1. Ist top=null, dann gib null zurück
- 2. Andernfalls merke Zeiger p von top in r.
- 3. Setze top auf p.next und gib r zurück

Implementation Push



 $\operatorname{\mathtt{push}}(x,S)$:

- 1. Erzeuge neues Listenelement mit x und Zeiger auf den Wert von top.
- 2. Setze top auf den Knotem mit x.

10

Analyse

Jede der Operationen push, pop, top und is ${\tt Empty}$ auf dem Stack ist in $\mathcal{O}(1)$ Schritten ausführbar.

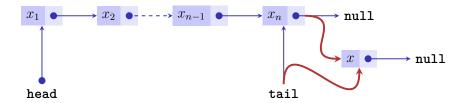
¹³hoffentlich

Queue (Schlange / Warteschlange / Fifo)

Queue ist ein ADT mit folgenden Operationen:

- \blacksquare enqueue(x,Q): fügt x am Ende der Schlange an.
- dequeue(Q): entfernt x vom Beginn der Schlange und gibt x zurück (null sonst.)
- lacktriangle head(Q): liefert das Objekt am Beginn der Schlage zurück (null sonst.)
- lacktriangle is $\operatorname{Empty}(Q)$: liefert true wenn Queue leer, sonst false.
- emptyQueue(): liefert leere Queue zurück.

Implementation Queue

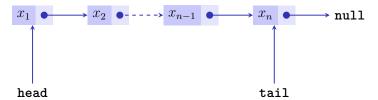


enqueue(x, S):

314

- 1. Erzeuge neues Listenelement mit x und Zeiger auf null.
- 2. Wenn tail \neq null, setze tail.next auf den Knoten mit x.
- 3. Setze tail auf den Knoten mit x.
- 4. Ist head = null, dann setze head auf tail.

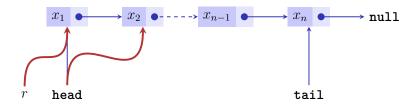
Invarianten!



Mit dieser Implementation gilt

- \blacksquare entweder head = tail = null,
- lacktriangledown oder head = tail eq null und head.next = null
- oder head \neq null und tail \neq null und head \neq tail und head.next \neq null.

Implementation Queue



dequeue(S):

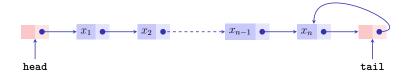
- 1. Merke Zeiger von head in r. Wenn r = null, gib r zurück.
- 2. Setze den Zeiger von head auf head.next.
- 3. Ist nun head = null, dann setze tail auf null.
- 4. Gib den Wert von r zurück.

Analyse

Jede der Operationen enqueue, dequeue, head und is
Empty auf der Queue ist in $\mathcal{O}(1)$ Schritten ausführbar.

Implementationsvarianten verketteter Listen

Liste mit Dummy-Elementen (Sentinels).



Vorteil: Weniger Spezialfälle!

Variante davon: genauso, dabei Zeiger auf ein Element immer einfach indirekt gespeichert. (Bsp: Zeiger auf x_3 zeigt auf x_2 .)

318

319

Implementationsvarianten verketteter Listen

Doppelt verkettete Liste



Übersicht

	enqueue	delete	search	concat
(A)	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
(B)	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$
(C)	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(1)$
(D)	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$	$\Theta(1)$

- (A) = Einfach verkettet
- (B) = Einfach verkettet, mit Dummyelement am Anfang und Ende
- (C) = Einfach verkettet, mit einfach indirekter Elementaddressierung
- (D) = Doppelt verkettet

12. Amortisierte Analyse

Amortisierte Analyse: Aggregat Analyse, Konto-Methode, Potentialmethode [Ottman/Widmayer, Kap. 3.3, Cormen et al, Kap. 17]

Multistack

322

Multistack unterstützt neben push und pop noch multipop(k,S): Entferne die $\min(\operatorname{size}(S),k)$ zuletzt eingefügten Objekte und liefere diese zurück.

Implementation wie beim Stack. Laufzeit von multipop ist O(k).

Akademische Frage

Führen wir auf einem Stack mit n Elementen n mal multipop(k,S) aus, kostet das dann $\mathcal{O}(n^2)$? Sicher richtig, denn jeder multipop kann Zeit $\mathcal{O}(n)$ haben. Wie bekommen wir eine schärfere Abschätzung?

Amortisierte Analyse

 Obere Schranke: Abschätzung der durchschnittlichen Laufzeit jeder betrachteten Operation im schlechtesten Fall.

$$\frac{1}{n} \sum_{i=1}^{n} \mathsf{Kosten}(\mathsf{op}_i)$$

- Nutzt aus, dass wenige teure Operationen vielen billigen Operationen gegenüberstehen.
- In der amortisierten Analyse sucht man nach einer Kostenfunktion / einem Potential, um zu zeigen, wie die billigen Operationen für die teuren Operationen "aufkommen" können.

Aggregierte Analyse

Direkte Argumentation: berechne eine Schranke für die Gesamtzahl der Elementaroperationen und teile durch die Anzahl der Operationen

Aggregierte Analyse: (Stack)

- Bei n Operationen können insgesamt maximal n Elemente auf den Stack gelegt werden. Also können auch insgesamt nur maximal n Elemente vom Stack entfernt werden.
- Für die Gesamtkosten ergibt sich

$$\sum_{i=1}^n \mathsf{Kosten}(\mathsf{op}_i) \leq 2n$$

und somit

amortisierte Kosten $(op_i) \le 2 \in \mathcal{O}(1)$

326

32

Kontomethode

Modell

- Der Computer basiert auf Münzen: jede Elementaroperation der Maschine kostet eine Münze.
- Für jede Operation op_k einer Datenstruktur wird eine bestimme Anzahl Münzen a_k auf eine Konto A eingezahlt: $A_k = A_{k-1} + a_k$
- Die Münzen vom Konto A werden verwendet, um die anfallenden echten Kosten t_k zu bezahlen.
- Das Konto A muss zu jeder Zeit genügend Münzen aufweisen, um die laufende Operation op_k zu bezahlen: $A_k t_k \ge 0 \, \forall k$.
- $\Rightarrow a_k$ sind die amortisierten Kosten der Operation op_k .

Kontomethode (Stack)

- Aufruf von push: kostet 1 CHF und zusätzlich kommt 1 CHF auf das Bankkonto ($a_k=2$)
- Aufruf von pop: kostet 1 CHF, wird durch Rückzahlung vom Bankkonto beglichen. ($a_k = 0$)

Kontostand wird niemals negativ.

 $a_k \leq 2 \, \forall \, k$, also: konstante amortisierte Kosten.

Potentialmethode

Leicht anderes Modell

- Definiere ein Potential Φ_i , welches zum Zustand der betrachteten Datenstruktur zum Zeitpunkt i gehört.
- Das Potential soll zum Ausgleichen teurer Operationen verwendet werden und muss daher so gewählt sein, dass es bei (häufigen) günstigen Operationen erhöht wird, während es die (seltenen) teuren Operationen durch einen fallenden Wert bezahlt.

Potentialmethode (formal)

Bezeichne t_i die realen Kosten der Operation op_i .

Potential funktion $\Phi_i \geq 0$ zur Datenstruktur nach i Operationen.

Voraussetzung: $\Phi_i \geq \Phi_0 \ \forall i$.

Amortisierte Kosten der i-ten Operation:

$$a_i := t_i + \Phi_i - \Phi_{i-1}.$$

Es gilt nämlich

332

$$\sum_{i=1}^{n} a_i = \sum_{i=1}^{n} (t_i + \Phi_i - \Phi_{i-1}) = \left(\sum_{i=1}^{n} t_i\right) + \Phi_n - \Phi_0 \ge \sum_{i=1}^{n} t_i.$$

330

Beispiel Stack

Potentialfunktion Φ_i = Anzahl Elemente auf dem Stack.

- **push**(x,S): Reale Kosten $t_i=1$. $\Phi_i-\Phi_{i-1}=1$. Amortisierte Kosten $a_i=2$.
- **p**op(S): Reale Kosten $t_i = 1$. $\Phi_i \Phi_{i-1} = -1$. Amortisierte Kosten $a_i = 0$.
- multipop(k,S): Reale Kosten $t_i=k$. $\Phi_i-\Phi_{i-1}=-k$. Amortisierte Kosten $a_i=0$.

Alle Operationen haben konstante amortisierte Kosten! Im Durchschnitt hat also Multipop konstanten Zeitbedarf. ¹⁴

Beispiel Binärer Zähler

Binärer Zähler mit k bits. Im schlimmsten Fall für jede Zähloperation maximal k Bitflips. Also $\mathcal{O}(n \cdot k)$ Bitflips für Zählen von 1 bis n. Geht das besser?

Reale Kosten t_i = Anzahl Bitwechsel von 0 nach 1 plus Anzahl Bitwechsel von 1 nach 0.

$$...0\underbrace{1111111}_{l \text{ Einsen}} + 1 = ...1\underbrace{0000000}_{l \text{ Nullen}}.$$
 $\Rightarrow t_i = l+1$

¹⁴Achtung: es geht nicht um den probabilistischen Mittelwert sondern den (worst-case) Durchschnitt der Kosten.

Binärer Zähler: Aggregatmethode

Zähle die Anzahl Bitwechsel beim Zählen von 0 bis n-1. Beobachtung

- Bit 0 wechselt für jedes $k-1 \rightarrow k$
- Bit 1 wechselt für jedes $2k 1 \rightarrow 2k$
- Bit 2 wechselt für jedes $4k 1 \rightarrow 4k$

Gesamte Anzahl Bitwechsel $\sum_{i=0}^{n-1} \frac{n}{2^i} \le n \cdot \sum_{i=0}^{\infty} \frac{1}{2^i} = 2n$ Amortisierte Kosten für jede Erhöhung: $\mathcal{O}(1)$ Bitwechsel.

Binärer Zähler: Potentialmethode

$$...0\underbrace{1111111}_{l \text{ Einsen}} + 1 = ...1\underbrace{0000000}_{l \text{ Nullen}}$$

Potentialfunktion Φ_i : Anzahl der 1-Bits von x_i .

$$\Rightarrow \Phi_0 = 0 \le \Phi_i \,\forall i$$

$$\Rightarrow \Phi_i - \Phi_{i-1} = 1 - l,$$

$$\Rightarrow a_i = t_i + \Phi_i - \Phi_{i-1} = l + 1 + (1 - l) = 2.$$

Amortisiert konstante Kosten für eine Zähloperation.

Binärer Zähler: Kontomethode

Beobachtung: bei jedem Inkrementieren wird genau ein Bit auf 1 gesetzt, während viele Bits auf 0 gesetzt werden könnten. Nur ein vorgängig auf 1 gesetztes Bit kann wieder auf 0 zurückgesetzt werden.

 $a_i=2$: 1 CHF reale Kosten für das Setzen $0\to 1$ plus 1 CHF für das Konto. Jedes Zurücksetzen $1\to 0$ kann vom Konto beglichen werden.

4

13. Wörterbücher

Wörterbuch, Selbstandornung, Implementation Wörterbuch mit Array / Liste / Skipliste. [Ottman/Widmayer, Kap. 3.3,1.7, Cormen et al, Kap. Problem 17-5]

Wörterbuch (Dictionary)

ADT zur Verwaltung von Schlüsseln aus ${\mathcal K}$ mit Operationen

- insert(k, D): Hinzufügen von $k \in \mathcal{K}$ in Wörterbuch D. Bereits vorhanden \Rightarrow Fehlermeldung.
- delete(k, D): Löschen von k aus dem Wörterbuch D. Nicht vorhanden \Rightarrow Fehlermeldung.
- \blacksquare $\operatorname{search}(k,D)\text{:}$ Liefert true wenn $k\in D\text{,}$ sonst false.

Idee

Implementiere Wörterbuch als sortiertes Array. Anzahl Elementaroperationen im schlechtesten Fall

Suchen $\mathcal{O}(\log n)$ \bigcirc Einfügen $\mathcal{O}(n)$ \bigcirc Löschen $\mathcal{O}(n)$

8

Andere Idee

Implementiere Wörterbuch als verkettete Liste Anzahl Elementaroperationen im schlechtesten Fall

Suchen $\mathcal{O}(n)$ $\stackrel{\textstyle \smile}{ }$ Einfügen $\mathcal{O}(1)^{15}$ $\stackrel{\textstyle \smile}{ }$ Löschen $\mathcal{O}(n)$ $\stackrel{\textstyle \smile}{ }$

13.1 Selbstanordnung

¹⁵Unter der Voraussetzung, dass wir die Existenz nicht prüfen wollen.

Selbstanordnung

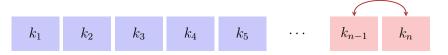
Problematisch bei der Verwendung verketteter Listen: lineare Suchzeit Idee: Versuche, die Listenelemente so anzuordnen, dass Zugriffe über die Zeit hinweg schneller möglich sind

Zum Beispiel

- Transpose: Bei jedem Zugriff auf einen Schlüssel wird dieser um eine Position nach vorne bewegt.
- Move-to-Front (MTF): Bei jedem Zugriff auf einen Schlüssel wird dieser ganz nach vorne bewegt.

Transpose

Transpose:

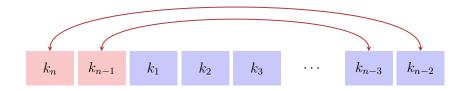


Worst case: n Wechselnde Zugriffe auf k_{n-1} und k_n . Laufzeit: $\Theta(n^2)$

342

Move-to-Front

Move-to-Front:



n Wechselnde Zugriffe auf k_{n-1} und k_n . Laufzeit: $\Theta(n)$ Man kann auch hier Folge mit quadratischer Laufzeit angeben, z.B. immer das letzte Element. Aber dafür ist keine offensichtliche Strategie bekannt, die viel besser sein könnte als MTF.

Analyse

Vergleichen MTF mit dem bestmöglichen Konkurrenten (Algorithmus) A. Wie viel besser kann A sein?

Annahmen:

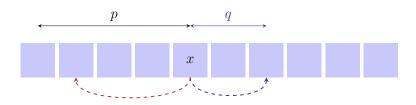
- MTF und A dürfen jeweils nur das zugegriffene Element x verschieben.
- MTF und A starten mit derselben Liste.

 M_k und A_k bezeichnen die Liste nach dem k-ten Schritt. $M_0 = A_0$.

Analyse

Kosten:

- Zugriff auf x: Position p von x in der Liste.
- Keine weiteren Kosten, wenn x vor p verschoben wird.
- Weitere Kosten q für jedes Element, das x von p aus nach hinten verschoben wird.



Amortisierte Analyse

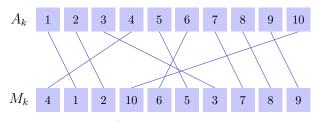
Sei eine beliebige Folge von Suchanfragen gegeben und seien $G_k^{(M)}$ und $G_k^{(A)}$ jeweils die Kosten im Schritt k für Move-to-Front und A. Suchen Abschätzung für $\sum_k G_k^{(M)}$ im Vergleich zu $\sum_k G_k^{(A)}$.

 \Rightarrow Amortisierte Analyse mit Potentialfunktion Φ .

346

Potentialfunktion

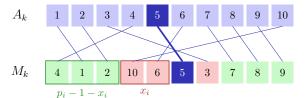
Potentialfunktion $\Phi=$ Anzahl der Inversionen von A gegen MTF. Inversion = Paar x, y so dass für die Positionen von x und y $\left(p^{(A)}(x) < p^{(A)}(y)\right) \neq \left(p^{(M)}(x) < p^{(M)}(y)\right)$

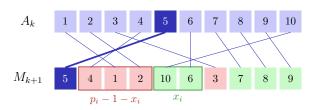


#Inversionen = #Kreuzungen

Abschätzung der Potentialfunktion: MTF

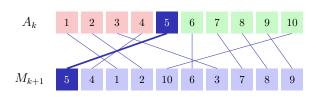
- Element i an Position $p_i := p^{(M)}(i)$.
- **Zugriffskosten** $C_k^{(M)} = p_i$.
- **a** x_i : Anzahl Elemente, die in M vor p_i und in A nach i stehen.
- lacktriangle MTF löst x_i Inversionen auf.
- $p_i x_i 1$: Anzahl Elemente, die in M vor p_i und in A vor i stehen.
- MTF erzeugt $p_i 1 x_i$ Inversionen.

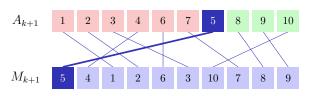




Abschätzung der Potentialfunktion: A

- Element i an Position $p^{(A)}(i)$.
- $X_k^{(A)}$: Anzahl Verschiebungen nach hinten (sonst 0).
- Zugriffskosten für i: $C_k^{(A)} = p^{(A)}(i) \ge p^{(M)}(i) x_i$.
- $\begin{tabular}{l} $\bf A$ erh\"{o}ht die Anzahl \\ Inversionen h\"{o}chstens um \\ $X_k^{(A)}$. \end{tabular}$





350

Abschätzung

$$\Phi_{k+1} - \Phi_k \le -x_i + (p_i - 1 - x_i) + X_k^{(A)}$$

Amortisierte Kosten von MTF im k-ten Schritt:

$$\begin{aligned} a_k^{(M)} &= C_k^{(M)} + \Phi_{k+1} - \Phi_k \\ &\leq p_i - x_i + (p_i - 1 - x_i) + X_k^{(A)} \\ &= (p_i - x_i) + (p_i - x_i) - 1 + X_k^{(A)} \\ &\leq C_k^{(A)} + C_k^{(A)} - 1 + X_k^{(A)} \leq 2 \cdot C_k^{(A)} + X_k^{(A)}. \end{aligned}$$

351

Abschätzung

Kosten Summiert

$$\sum_{k} G_{k}^{(M)} = \sum_{k} C_{k}^{(M)} \le \sum_{k} a_{k}^{(M)} \le \sum_{k} 2 \cdot C_{k}^{(A)} + X_{k}^{(A)}$$
$$\le 2 \cdot \sum_{k} C_{k}^{(A)} + X_{k}^{(A)}$$
$$= 2 \cdot \sum_{k} G_{k}^{(A)}$$

MTF führt im schlechtesten Fall höchstens doppelt so viele Operationen aus wie eine optimale Strategie.

13.2 Skiplisten

Sortierte Verkettete Liste



Element / Einfügeort suchen: worst-case n Schritte.

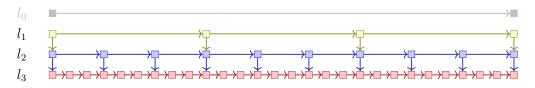
Sortierte Verkettete Liste mit 2 Ebenen

- **Anzahl Elemente:** $n_0 := n$
- Schrittweite Ebene 1: n_1
- Schrittweite Ebene 2: $n_2 = 1$
- \Rightarrow Element / Einfügeort suchen: worst-case $rac{n_0}{n_1}+rac{n_1}{n_2}$.
- \Rightarrow Beste Wahl für¹⁶ n_1 : $n_1=rac{n_0}{n_1}=\sqrt{n_0}$.

Element / Einfügeort suchen: worst-case $2\sqrt{n}$ Schritte.

354

Sortierte Verkettete Liste mit 3 Ebenen



- Anzahl Elemente: $n_0 := n$
- **Schrittweiten Ebenen** 0 < i < 3: n_i
- Schrittweite auf Ebene 3: $n_3 = 1$
- \Rightarrow Beste Wahl für (n_1,n_2) : $n_2=\frac{n_0}{n_1}=\frac{n_1}{n_2}=\sqrt[3]{n_0}$.

Element / Einfügeort suchen: worst-case $3 \cdot \sqrt[3]{n}$ Schritte.

Sortierte Verkettete Liste mit k Ebenen (Skipliste)

- Anzahl Elemente: $n_0 := n$
- Schrittweiten Ebenen 0 < i < k: n_i
- Schrittweite auf Ebene k: $n_k = 1$
- \Rightarrow Beste Wahl für (n_1,\ldots,n_k) : $n_{k-1}=rac{n_0}{n_1}=rac{n_1}{n_2}=\cdots=\sqrt[k]{n_0}$.

Element / Einfügeort suchen: worst-case $k \cdot \sqrt[k]{n}$ Schritte¹⁷.

Annahme: $n = 2^k$

 \Rightarrow worst case $\log_2 n \cdot 2$ Schritte und $\frac{n_i}{n_{i+1}} = 2 \, \forall \, 0 \leq i < \log_2 n$.

355

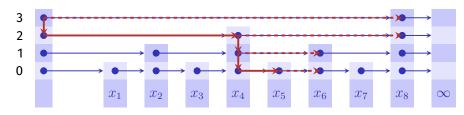
•

¹⁶Differenzieren und 0 setzen, siehe Anhang

¹⁷(Herleitung: Anhang)

Suche in Skipliste

Perfekte Skipliste



 $x_1 \leq x_2 \leq x_3 \leq \cdots \leq x_9$.

Beispiel: Suche nach einem Schlüssel x mit $x_5 < x < x_6$.

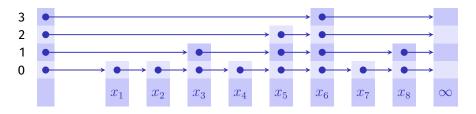
Analyse Perfekte Skipliste (schlechtester Fall)

Suchen in $\mathcal{O}(\log n)$. Einfügen in $\mathcal{O}(n)$.

358

Randomisierte Skipliste

Idee: Füge jeweils einen Knoten mit zufälliger Höhe H ein, wobei $\mathbb{P}(H=i)=\frac{1}{2^{i+1}}$.



Analyse Randomisierte Skipliste

Theorem 15

Die Anzahl an erwarteten Elementaroperationen für Suchen, Einfügen und Löschen eines Elements in einer randomisierten Skipliste ist $\mathcal{O}(\log n)$.

Der längliche Beweis, welcher im Rahmen dieser Vorlesung nicht geführt wird, betrachtet die Länge eines Weges von einem gesuchten Knoten zurück zum Startpunkt im höchsten Level.

13.3 Anhang

Mathematik zur Skipliste

[Mathematik zur k-Level Skipliste]

Vorige Folie
$$\Rightarrow \frac{n_t}{n_0} = \frac{n_t}{n_{t-1}} \frac{n_{t-1}}{n_{t-2}} \dots \frac{n_1}{n_0} = \left(\frac{n_1}{n_0}\right)^t$$
Insbesondere $1 = n_k = \frac{n_1^k}{n_0^{k-1}} \Rightarrow n_1 = \sqrt[k]{n_0^{k-1}}$

Also
$$n_{k-1} = \frac{n_0}{n_1} = \sqrt[k]{\frac{n_0^k}{n_0^{k-1}}} = \sqrt[k]{n_0}.$$

Maximale Anzahl Schritte in der Skipliste $f(\vec{n}) = k \cdot (\sqrt[k]{n_0})$ Angenommen $n_0 = 2^k$, dann $\frac{n_l}{n_{l+1}} = 2$ für alle $0 \le l < k$ (Skipliste halbiert die Daten in jedem Level), und $f(n) = k \cdot 2 = 2\log_2 n \in \Theta(\log n)$.

[Mathematik zur k-Level Skipliste]

Gegeben Anzahl Datenpunkte n_0 , Anzahl Level k>0 und Anzahl Elemente n_l die pro Level l übersprungen werden, $n_k=1$. Maximale Anzahl totaler Schritte in der Skipliste:

$$f(\vec{n}) = \frac{n_0}{n_1} + \frac{n_1}{n_2} + \dots \frac{n_{k-1}}{n_k}$$

363

Minimiere f für (n_1,\ldots,n_{k-1}) : $\frac{\partial f(\vec{n})}{\partial n_t} = 0$ für alle 0 < t < k, $\frac{\partial f(\vec{n})}{\partial n_t} = -\frac{n_{t-1}}{n_t^2} + \frac{1}{n_{t+1}} = 0 \Rightarrow n_{t+1} = \frac{n_t^2}{n_{t-1}}$ und $\frac{n_{t+1}}{n_t} = \frac{n_t}{n_{t-1}}$.

362

14. Hashing

Hashtabellen, Pre-Hashing, Hashing, Kollisionsauflösung durch Verketten, Einfaches gleichmässiges Hashing, Gebräuchliche Hashfunktionen, Tabellenvergrösserung, offene Addressierung: Sondieren, Gleichmässiges Hashing [Ottman/Widmayer, Kap. 4.1-4.3.2, 4.3.4, Cormen et al, Kap. 11-11.4]

Motivierendes Beispiel

Ziel: Effiziente Verwaltung einer Tabelle aller n ETH-Studenten. Mögliche Anforderung: Schneller Zugriff (Einfügen, Löschen, Finden) von Datensätzen nach Name.

Wörterbuch (Dictionary)

Abstrakter Datentyp (ADT) D zur Verwaltung einer Menge von Einträgen i mit Schlüsseln $k \in \mathcal{K}$. Operationen

- **D.** insert(i): Hinzufügen oder Überschreiben von i im Wörterbuch D.
- D.delete(i): Löschen von i aus dem Wörterbuch D. Nicht vorhanden \Rightarrow Fehlermeldung.
- D.search(k): Liefert Eintrag mit Schlüssel k, wenn er existiert.

366

367

Wörterbuch in C++

Assoziativer Container std::unordered_map<>

Motivation / Verwendung

Wahrscheinlich die gängigste Datenstruktur

- Unterstützt in vielen Programmiersprachen (C++, Java, Python, Ruby, Javascript, C# ...)
- Offensichtliche Verwendung
 - Datenbanken / Tabellenkalkulation
 - Symboltabellen in Compilern und Interpretern
- Weniger offensichtlich
 - Substring Suche (Google, grep)
 - Ähnlichkeit von Texten (Dokumentenvergleich, DNA)
 - Dateisynchronisation
 - Kryptographie: Filetransfer / Identifikation

 $^{^{\}rm 18}$ Schlüssel-Wert Paare (k,v) , im Folgenden betrachten wir hauptsächlich die Schlüssel.

1. Idee: Direkter Zugriff (Array)

Index	Eintrag		
0	-		
1	-		
2	-		
3	[3,wert(3)]		
4	-		
5	-		
:	:		
k	[k,wert(k)]		
÷	:		

Probleme

- 1. Schlüssel müssen nichtnegative ganze Zahlen sein
- 2. Grosser Schlüsselbereich ⇒ grosses Array

Lösung zum ersten Problem: Pre-hashing

Prehashing: Bilde Schlüssel ab auf positive Ganzzahlen mit einer Funktion $ph:\mathcal{K}\to\mathbb{N}$

- Theoretisch immer möglich, denn jeder Schlüssel ist als Bitsequenz im Computer gespeichert
- Theoretisch auch: $x = y \Leftrightarrow ph(x) = ph(y)$
- In der Praxis: APIs bieten Funktionen zum pre-hashing an. (Java: object.hashCode(), C++: std::hash<>, Python: hash(object))
- APIs bilden einen Schlüssel aus der Schlüsselmenge ab auf eine Ganzzahl mit beschränkter Grösse.¹9

370

Prehashing Beispiel: String

Zuordnung Name $s = s_1 s_2 \dots s_{l_a}$ zu Schlüssel

$$ph(s) = \left(\sum_{i=0}^{l_s-1} s_{l_s-i} \cdot b^i\right) \bmod 2^w$$

b so, dass verschiedene Namen möglichst verschiedene Schlüssel erhalten. w Wortgrösse des Systems (z.B. 32 oder 64).

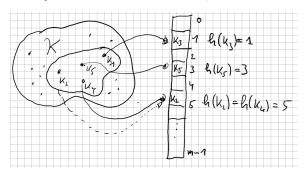
Beispiel (Java), mit
$$b=31$$
, $w=32$ Ascii-Werte s_i .
Anna $\mapsto 2045632$

Jacqueline $\mapsto 2042089953442505 \mod 2^{32} = 507919049$

Kollision: $h(k_i) = h(k_i)$.

Lösung zum zweiten Problem: Hashing

Reduziere des Schlüsseluniversum: Abbildung (Hash-Funktion) $h: \mathcal{K} \to \{0, ..., m-1\}$ ($m \approx n = \text{Anzahl Einträge in der Tabelle})$



373

¹⁹Somit gilt die Implikation $ph(x) = ph(y) \Rightarrow x = y$ **nicht** mehr für alle x,y.

Nomenklatur

Hashfunktion h: Abbildung aus der Menge der Schlüssel $\mathcal K$ auf die Indexmenge $\{0,1,\ldots,m-1\}$ eines Arrays (Hashtabelle).

$$h: \mathcal{K} \to \{0, 1, \dots, m-1\}.$$

Meist $|\mathcal{K}| \gg m$. Es gibt also $k_1, k_2 \in \mathcal{K}$ mit $h(k_1) = h(k_2)$ (Kollision).

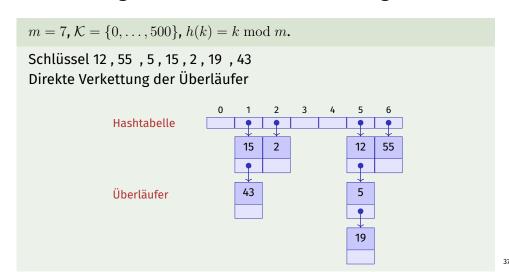
Eine Hashfunktion sollte die Menge der Schlüssel möglichst gleichmässig auf die Positionen der Hashtabelle verteilen.

374

Algorithmen zum Hashing mit Verkettung

- insert(i) Prüfe ob Schlüssel k vom Eintrag i in Liste an Position h(k). Falls nein, füge i am Ende der Liste ein; andernfalls ersetze das Element durch i.
- find(k) Prüfe ob Schlüssel k in Liste an Position h(k). Falls ja, gib die Daten zum Schlüssel k zurück. Andernfalls Rückgabe eines leeren Elements null.
- lacktriangle delete(k) Durchsuche die Liste an Position h(k) nach k. Wenn Suche erfolgreich, entferne das entsprechende Listenelement.

Behandlung von Kollisionen: Verkettung



Worst-case Analyse

Schlechtester Fall: alle Schlüssel werden auf den gleichen Index abgebildet.

 $\Rightarrow \Theta(n)$ pro Operation im schlechtesten Fall.

Einfaches Gleichmässiges Hashing

Starke Annahmen: Jeder beliebige Schlüssel wird

- mit gleicher Wahrscheinlichkeit (Uniformität)
- und unabhängig von den anderen Schlüsseln (Unabhängigkeit)

auf einen der m verfügbaren Slots abgebildet.

Einfaches Gleichmässiges Hashing

Unter der Voraussetzung von einfachem gleichmässigen Hashing: Erwartete Länge einer Kette, wenn n Elemente in eine Hashtabelle mit m Elementen eingefügt werden

$$\begin{split} \mathbb{E}(\text{L\"{a}nge Kette j}) &= \mathbb{E}\bigg(\sum_{i=0}^{n-1}\mathbb{I}(h(k_i)=j)\bigg) = \sum_{i=0}^{n-1}\mathbb{P}(h(k_i)=j) \\ &= \sum_{i=1}^{n}\frac{1}{m} = \frac{n}{m} \end{split}$$

 $\alpha=n/m$ heisst Belegungsfaktor oder Füllgrad der Hashtabelle.

378

379

Einfaches Gleichmässiges Hashing

Theorem 16

Sei eine Hashtabelle Verkettung gefüllt mit Füllgrad $\alpha=\frac{n}{m}<1$. Unter der Annahme vom einfachen gleichmässigen Hashing hat die nächste Operation erwartete Laufzeitkosten von $\Theta(1+\alpha)$.

Folgerung: ist die Anzahl der Slots m der Hashtabelle immer mindestens proportional zur Anzahl Elemente n in der Hashtabelle, $n \in \mathcal{O}(m) \Rightarrow$ Erwartete Laufzeit der Operationen Suchen, Einfügen und Löschen ist $\mathcal{O}(1)$.

Weitere Analyse (direkt verkettete Liste)

- 1. Erfolglose Suche. Durchschnittliche Listenlänge ist $\alpha=\frac{n}{m}$. Liste muss ganz durchlaufen werden.
 - ⇒ Durchschnittliche Anzahl betrachteter Einträge

$$C'_n = \alpha$$
.

- 2. Erfolgreiche Suche. Betrachten die Einfügehistorie: Schlüssel j sieht durchschnittliche Listenlänge (j-1)/m.
 - \Rightarrow Durchschnittliche Anzahl betrachteter Einträge

$$C_n = \frac{1}{n} \sum_{j=1}^{n} (1 + (j-1)/m) = 1 + \frac{1}{n} \frac{n(n-1)}{2m} \approx 1 + \frac{\alpha}{2}.$$

Vor und Nachteile der Verkettung

Vorteile der Strategie:

- Belegungsfaktoren $\alpha > 1$ möglich
- Entfernen von Schlüsseln einfach

Nachteile

Speicherverbrauch der Verkettung

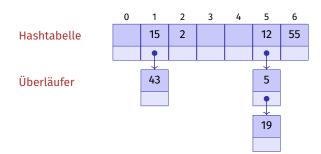
Beispiele gebräuchlicher Hashfunktionen

$$h(k) = k \mod m$$

Ideal: m Primzahl, nicht zu nahe bei Potenzen von 2 oder 10 Aber oft: $m=2^k-1$ ($k\in\mathbb{N}$)

[Variante:Indirekte Verkettung]

Beispiel m=7, $\mathcal{K}=\{0,\ldots,500\}$, $h(k)=k \bmod m$. Schlüssel 12, 55, 5, 15, 2, 19, 43 Indirekte Verkettung der Überläufer



382

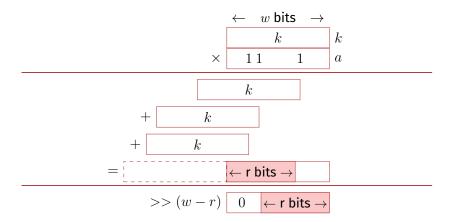
Beispiele gebräuchlicher Hashfunktionen

Multiplikationsmethode

$$h(k) = \left\lfloor (a \cdot k \bmod 2^w)/2^{w-r} \right\rfloor \bmod m$$

- Guter Wert für a: $\left\lfloor \frac{\sqrt{5}-1}{2} \cdot 2^w \right\rfloor$: Integer, der die ersten w Bits des gebrochenen Teils der irrationalen Zahl darstellt.
- \blacksquare Tabellengrösse $m=2^r$, w= Grösse des Maschinenworts in Bits.
- Multiplikation addiert k entlang aller Bits von a, Ganzzahldivision durch 2^{w-r} und $mod\ m$ extrahieren die oberen r Bits.
- Als Code geschrieben sehr einfach: a * k >> (w-r)

Illustration



Tabellenvergrösserung

In Abhängigkeit vom Füllgrad Tabellengrösse jeweils verdoppeln.

 \Rightarrow Amortisierte Analyse.

Jede Operation vom Hashing mit Verketten hat erwartet amortisierte Kosten $\Theta(1)$.

Tabellenvergrösserung

- \blacksquare Wissen nicht a priori, wie gross n sein wird.
- lacksquare Benötigen $m=\Theta(n)$ zu jeder Zeit.

Grösse der Tabelle muss angepasst werden. Hash-Funktion ändert sich ⇒ Rehashing

- Alloziere Array A' mit Grösse m' > m
- Füge jeden Eintrag von A erneut in A' ein (mit erneutem Hashing)
- Setze $A \leftarrow A'$.
- Kosten: $\mathcal{O}(n+m+m')$.

Wie wählt man m'?

386

Offene Addressierung

Speichere die Überläufer direkt in der Hashtabelle mit einer Sondierungsfunktion $s:\mathcal{K}\times\{0,1,\ldots,m-1\}\to\{0,1,\ldots,m-1\}$ Tabellenposition des Schlüssels entlang der Sondierungsfolge

$$S(k) := (s(k,0), s(k,1), \dots, s(k,m-1)) \mod m$$

Sondierungsfolge muss für jedes $k \in \mathcal{K}$ eine Permutation sein von $\{0,1,\ldots,m-1\}$

Begriffsklärung: Dieses Verfahren nutzt offene Addressierung (Positionen in der Hashtabelle nicht fixiert), heisst aber trotzdem ein geschlossenes Hashverfahren (Einträge bleiben in der Hashtabelle).

Algorithmen zur offenen Addressierung

- insert(i) Suche Schlüssel k von i in der Tabelle gemäss Sondierungssequenz S(k). Ist k nicht vorhanden, füge k an die erste freie Position in der Sondierungsfolge ein. Andernfalls Fehlermeldung.
- find(k) Durchlaufe Tabelleneinträge gemäss S(k). Wird k gefunden, gib die zu k gehörenden Daten zurück. Andernfalls Rückgabe eines leeres Elements null.
- delete(k) Suche k in der Tabelle gemäss S(k). Wenn k gefunden, ersetze k durch den speziellen Schlüssel removed.

Lineares Sondieren

$$s(k,j) = h(k) + j \Rightarrow S(k) = (h(k), h(k) + 1, \dots, h(k) + m - 1) \mod m$$

$$m = 7$$
, $\mathcal{K} = \{0, \dots, 500\}$, $h(k) = k \mod m$.

Schlüssel 12, 55, 5, 15, 2, 19

0	1	2	3	4	5	6
5	15	2	19		12	55

90

[Analyse Lineares Sondieren (ohne Herleitung)]

1. Erfolglose Suche. Durchschnittliche Anzahl betrachteter Einträge

$$C_n' \approx \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} \right)$$

2. Erfolgreiche Suche. Durchschnittliche Anzahl betrachteter Einträge

$$C_n \approx \frac{1}{2} \left(1 + \frac{1}{1 - \alpha} \right).$$

Diskussion

Beispiel $\alpha = 0.95$

Erfolglose Suche betrachtet im Durchschnitt 200 Tabelleneinträge! (Hier ohne Herleitung.).

Grund für die schlechte Performance?

Primäre Häufung: Ähnliche Hashaddressen haben ähnliche Sondierungsfolgen \Rightarrow lange zusammenhängende belegte Bereiche.

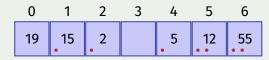
Quadratisches Sondieren

$$s(k,j) = h(k) + \lceil j/2 \rceil^2 (-1)^{j+1}$$

$$S(k) = (h(k), h(k) + 1, h(k) - 1, h(k) + 4, h(k) - 4, \dots) \mod m$$

$$m = 7$$
, $\mathcal{K} = \{0, \dots, 500\}$, $h(k) = k \mod m$.

Schlüssel 12, 55, 5, 15, 2, 19



[Analyse Quadratisches Sondieren (ohne Herleitung)]

1. Erfolglose Suche. Durchschnittliche Anzahl betrachteter Einträge

$$C'_n \approx \frac{1}{1-\alpha} - \alpha + \ln\left(\frac{1}{1-\alpha}\right)$$

2. Erfolgreiche Suche. Durchschnittliche Anzahl betrachteter Einträge

$$C_n \approx 1 + \ln\left(\frac{1}{1-\alpha}\right) - \frac{\alpha}{2}.$$

394

Diskussion

Beispiel $\alpha=0.95$

Erfolglose Suche betrachtet im Durchschnitt 22 Tabelleneinträge (Hier ohne Herleitung.)

Grund für die schlechte Performance?

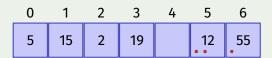
Sekundäre Häufung: Synonyme k und k' (mit h(k) = h(k')) durchlaufen dieselbe Sondierungsfolge.

Double Hashing

Zwei Hashfunktionen h(k) und h'(k). $s(k, j) = h(k) + j \cdot h'(k)$. $S(k) = (h(k), h(k) + h'(k), h(k) + 2h'(k), \dots, h(k) + (m-1)h'(k)) \mod m$

m = 7, $K = \{0, ..., 500\}$, $h(k) = k \mod 7$, $h'(k) = 1 + k \mod 5$.

Schlüssel 12, 55, 5, 15, 2, 19



Double Hashing

- Sondierungsreihenfolge muss Permutation aller Hashadressen bilden. Also $h'(k) \neq 0$ und h'(k) darf m nicht teilen, z.B. garantiert mit m prim.
- lacksquare h' sollte möglichst unabhängig von h sein (Vermeidung sekundärer Häufung).

Unabhängigkeit:

$$\mathbb{P}\big((h(k)=h(k')) \ \wedge \ (h'(k)=h'(k'))\big) = \mathbb{P}\big(h(k)=h(k')\big) \cdot \mathbb{P}\big(h'(k)=h'(k')\big).$$

Unabhängigkeit weitgehend erfüllt von $h(k) = k \mod m$ und $h'(k) = 1 + k \mod (m-2)$ (m prim).

[Analyse Double Hashing]

Sind h und h' unabhängig, dann:

1. Erfolglose Suche. Durchschnittliche Anzahl betrachteter Einträge

$$C_n' \approx \frac{1}{1-\alpha}$$

2. Erfolgreiche Suche. Durchschnittliche Anzahl betrachteter Einträge

$$C_n \approx \frac{1}{\alpha} \ln \left(\frac{1}{1 - \alpha} \right)$$

398

Gleichmässiges Hashing

Starke Annahme: Die Sondierungssequenz S(k) eines Schlüssels k ist mit gleicher Wahrscheinlichkeit eine der m! vielen Permutationssequenzen von $\{0,1,\ldots,m-1\}$.

(Double Hashing kommt dem am ehesten nahe)

Analyse gleichmässiges Hashing mit offener Addressierung

Theorem 17

Sei eine Hashtabelle mit offener Addressierung gefüllt mit Füllgrad $\alpha=\frac{n}{m}<1$. Unter der Annahme vom gleichmässigen Hashing hat die nächste Operation erwartete Laufzeitkosten von $\leq \frac{1}{1-\alpha}$.

Analyse: Beweis des Theorems

Zufallsvariable X: Anzahl Sondierungen bei einer erfolglosen Suche.

$$\mathbb{P}(X \ge i) \stackrel{*}{=} \frac{n}{m} \cdot \frac{n-1}{m-1} \cdot \frac{n-2}{m-2} \cdots \frac{n-i+2}{m-i+2}$$

$$\stackrel{**}{\le} \left(\frac{n}{m}\right)^{i-1} = \alpha^{i-1}. \qquad (1 \le i \le m)$$

* : Ereignis A_i : Slot beim j-ten Schritt belegt.

$$\mathbb{P}(A_1 \cap \cdots \cap A_{i-1}) = \mathbb{P}(A_1) \cdot \mathbb{P}(A_2 | A_1) \cdot \dots \cdot \mathbb{P}(A_{i-1} | A_1 \cap \cdots \cap A_{i-2}),$$

** : $\frac{n-1}{m-1} < \frac{n}{m}$ da n < m: $\frac{n-1}{m-1} < \frac{n}{m} \Leftrightarrow \frac{n-1}{n} < \frac{m-1}{m} \Leftrightarrow 1 - \frac{1}{n} < 1 - \frac{1}{m} \Leftrightarrow n < m$ (n > 0, m > 0)

Ausserdem $\mathbb{P}(x \geq i) = 0$ für $i \geq m$. Also

$$\mathbb{E}(X) \overset{\mathsf{Anhang}}{=} \sum_{i=1}^{\infty} \mathbb{P}(X \geq i) \leq \sum_{i=1}^{\infty} \alpha^{i-1} = \sum_{i=0}^{\infty} \alpha^{i} = \frac{1}{1-\alpha}.$$

Übersicht

	$\alpha = 0.50$		$\alpha = 0.90$		$\alpha = 0.95$	
	C_n	C'_n	C_n	C'_n	C_n	C'_n
(Direkte) Verkettung	1.25	0.50	1.45	0.90	1.48	0.95
Lineares Sondieren	1.50	2.50	5.50	50.50	10.50	200.50
Quadratisches Sondieren	1.44	2.19	2.85	11.40	3.52	22.05
Gleichmässiges Hashing	1.39	2.00	2.56	10.00	3.15	20.00

 α : Belegungsgrad.

 C_n : Anzahl Schritte erfolgreiche Suche,

 C'_n : Anzahl Schritte erfolglose Suche

[Erfolgreiche Suche beim gleichmässigen offenen Hashing]

Theorem 18

Sei eine Hashtabelle mit offener Addressierung gefüllt mit Füllgrad $\alpha=\frac{n}{m}<1$. Unter der Annahme vom gleichmässigen Hashing hat die erfolgreiche Suche erwartete Laufzeitkosten von $\leq \frac{1}{\alpha}\cdot\log\frac{1}{1-\alpha}$.

Beweis: Cormen et al, Kap. 11.4

402

14.7 Anhang

Mathematische Formeln

[Geburtstagsparadoxon]

Annahme: m Urnen, n Kugeln (oBdA $n \le m$). n Kugeln werden gleichverteilt in Urnen gelegt.



Wie gross ist die Kollisionswahrscheinlichkeit?

Geburtstagsparadoxon: Bei wie vielen Personen (n) ist die Wahrscheinlichkeit, dass zwei am selben Tag (m=365) Geburtstag haben grösser als 50%?

[Geburtstagsparadoxon]

 $\mathbb{P}(\mathsf{keine\ Kollision}) = \frac{m}{m} \cdot \frac{m-1}{m} \cdot \dots \cdot \frac{m-n+1}{m} = \frac{m!}{(m-n)! \cdot m^m}.$

Sei $a \ll m$. Mit $e^x = 1 + x + \frac{x^2}{2!} + \ldots$ approximiere $1 - \frac{a}{m} \approx e^{-\frac{a}{m}}$. Damit:

$$1 \cdot \left(1 - \frac{1}{m}\right) \cdot \left(1 - \frac{2}{m}\right) \cdot \dots \cdot \left(1 - \frac{n-1}{m}\right) \approx e^{-\frac{1 + \dots + n - 1}{m}} = e^{-\frac{n(n-1)}{2m}}.$$

Es ergibt sich

406

$$\mathbb{P}(\mathsf{Kollision}) = 1 - e^{-\frac{n(n-1)}{2m}}.$$

Auflösung zum Geburtstagsparadoxon: Bei 23 Leuten ist die Wahrscheinlichkeit für Geburstagskollision 50.7%. Zahl stammt von der leicht besseren Approximation via Stirling Formel. $n! \approx \sqrt{2\pi n} \cdot n^n \cdot e^{-n}$

407

[Erwartungswertformel]

 $X \geq 0$ diskrete Zufallsvariable mit $\mathbb{E}(X) < \infty$

$$\begin{split} \mathbb{E}(X) &\stackrel{(def)}{=} \sum_{x=0}^{\infty} x \mathbb{P}(X=x) \\ &\stackrel{\text{Aufz\"{a}hlen}}{=} \sum_{x=1}^{\infty} \sum_{y=x}^{\infty} \mathbb{P}(X=y) \\ &= \sum_{x=0}^{\infty} \mathbb{P}(X>x) \end{split}$$

15. C++ vertieft (III): Funktoren und Lambda

Was lernen wir heute?

- Funktoren: Objekte mit überladenem Funktionsoperator ().
- Closures
- Lambda-Ausdrücke: Syntaktischer Zucker
- Captures

Funktoren: Motivierung

```
template <typename T, typename Function>
void filter(const T& collection, Function f);

template <typename T>
bool even(T x){
   return x % 2 == 0;
}

std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
filter(a,even<int>); // output: 2,4,6,16
```

Funktoren: Motivierung

Ein simpler Ausgabefilter

```
template <typename T, typename Function>
void filter(const T& collection, Function f){
   for (const auto& x: collection)
      if (f(x)) std::cout << x << " ";
   std::cout << "\n";
}</pre>
```

filter funktioniert wenn das erste Argument einen Iterator anbietet und das zweite auf Elemente des Iterators angewendet werden kann und das Resultat zu bool konvertiertbar ist.

410

Funktor: Objekt mit überladenem Operator ()

```
class GreaterThan{
  int value; // state
  public:
    GreaterThan(int x):value{x}{}

  bool operator() (int par) const {
    return par > value;
  }
};
```

Ein Funktor ist ein aufrufbares Objekt. Kann verstanden werden als Funktion mit Zustand.

```
std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
int value=8;
filter(a,GreaterThan(value)); // 9,11,16,19
```

Funktor: Objekt mit überladenem Operator ()

```
template <typename T>
class GreaterThan{
   T value;
public:
   GreaterThan(T x):value{x}{}

   bool operator() (T par) const{
      return par > value;
   }
};

std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
int value=8;
filter(a,GreaterThan<int>(value)); // 9,11,16,19
```

Dasselbe mit Lambda-Expression

```
std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
int value=8;
filter(a, [value](int x) {return x > value;} );
```

14

Summe aller Elemente - klassisch

```
std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
int sum = 0;
for (auto x: a)
   sum += x;
std::cout << sum << std::endl; // 83</pre>
```

Summe aller Elemente - mit Funktor

```
template <typename T>
struct Sum{
    T value = 0;

    void operator() (T par){ value += par; }
};

std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
Sum<int> sum;
// for_each copies sum: we need to copy the result back
sum = std::for_each(a.begin(), a.end(), sum);
std::cout << sum.value << std::endl; // 83</pre>
```

Summe aller Flemente - mit Referenzen

```
template <typename T>
struct SumR{
   T& value;
   SumR (T& v):value{v} {}

   void operator() (T par){ value += par; }
};

std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
int s=0;
SumR<int> sum{s};
// cannot (and do not need to) assign to sum here
std::for_each(a.begin(), a.end(), sum);
std::cout << s << std::endl; // 83</pre>
Geht natürlich sehr ähnlich auch mit Zeigern
```

Sortieren, mal anders

```
// pre: i >= 0
// post: returns sum of digits of i
int q(int i){
    int res =0;
    for(;i>0;i/=10)
        res += i % 10;
    return res;
}

std::vector<int> v {10,12,9,7,28,22,14};
std::sort (v.begin(), v.end(),
    [] (int i, int j) { return q(i) < q(j);}
);

Jetzt v =10,12,22,14,7,9,28 (sortiert nach Quersumme)</pre>
```

Summe aller Elemente - mit Λ

```
std::vector<int> a {1,2,3,4,5,6,7,9,11,16,19};
int s=0;
std::for_each(a.begin(), a.end(), [&s] (int x) {s += x;} );
std::cout << s << std::endl;</pre>
```

418

Lambda-Expressions im Detail

```
[value] (int x) ->bool {return x > value;}
Capture Parameter Rück-
gabe-
typ
```

Closure

```
[value] (int x) ->bool {return x > value;}
```

- Lambda-Expressions evaluieren zu einem temporären Objekt einer closure
- Die closure erhält den Ausführungskontext der Funktion, also die captured Objekte.
- Lambda-Expressions können als Funktoren implementiert werden.

Simple Lambda-Expression

```
[]()->void {std::cout << "Hello World";}
Aufruf:
[]()->void {std::cout << "Hello World";}();
Zuweisung:
auto f = []()->void {std::cout << "Hello World";};</pre>
```

2

Minimale Lambda-Expression

[]{}

■ Rückgabetyp kann inferiert werden, wenn kein oder nur ein return:²⁰

```
[]() {std::cout << "Hello World";}</pre>
```

■ Keine Parameter und kein expliziter Rückgabetyp ⇒ () kann weggelassen werden

²⁰Seit C++14 auch mehrere returns, sofern derselbe Rückgabetyp deduziert wird

```
[]{std::cout << "Hello World";}
```

■ [...] kann nie weggelassen werden.

```
[](int x, int y) {std::cout << x * y;} (4,5);
Output: 20
```

Beispiele

Beispiele

```
int k = 8;
auto f = [](int& v) {v += v;};
f(k);
std::cout << k;
Output: 16</pre>
```

Beispiele

```
int k = 8;
auto f = [](int v) {v += v;};
f(k);
std::cout << k;
Output: 8</pre>
```

26

Capture - Lambdas

Für Lambda-Expressions bestimmt die capture-Liste über den zugreifbaren Teil des Kontextes

Syntax:

- [x]: Zugriff auf kopierten Wert von x (nur lesend)
- [&x]: Zugriff zur Referenz von x
- [&x,y]: Zugriff zur Referenz von x und zum kopierten Wert von y
- [&]: Default-Referenz-Zugriff auf alle Objekte im Kontext der Lambda-Expression
- [=]: Default-Werte-Zugriff auf alle Objekte im Kontext der Lambda-Expression

Capture - Lambdas

```
int elements=0;
int sum=0;
std::for_each(v.begin(), v.end(),
   [&] (int k) {sum += k; elements++;} // capture all by reference
)
```

Capture - Lambdas

```
template <typename T>
void sequence(vector<int> & v, T done){
  int i=0;
  while (!done()) v.push_back(i++);
}

vector<int> s;
sequence(s, [&] {return s.size() >= 5;} )
jetzt v = 0 1 2 3 4
```

Die capture liste bezieht sich auf den Kontext der Lambda Expression

Capture – Lambdas

Wann wird der Wert gelesen?

```
int v = 42;
auto func = [=] {std::cout << v << "\n"};
v = 7;
func();</pre>
```

Ausgabe: 42

Werte werden bei der Definition der (temporären) Lambda-Expression zugewiesen.

430

Capture – Lambdas

(Warum) funktioniert das?

```
class Limited{
  int limit = 10;
public:
  // count entries smaller than limit
  int count(const std::vector<int>& a){
    int c = 0;
    std::for_each(a.begin(), a.end(),
        [=,&c] (int x) {if (x < limit) c++;}
  );
  return c;
};</pre>
```

Der this pointer wird per default implizit kopiert

Capture – Lambdas

```
struct mutant{
  int i = 0;
  void do(){ [=] {i=42;}();}
};

mutant m;
m.do();
std::cout << m.i;</pre>
```

Ausgabe: 42

Der this pointer wird per default implizit kopiert

Lambda Ausdrücke sind Funktoren

```
[x, &y] () {y = x;}
kann implementiert werden als
unnamed {x,y};
mit

class unnamed {
  int x; int& y;
  unnamed (int x_, int& y_) : x (x_), y (y_) {}
  void operator () () {y = x;}
};
```

Lambda Ausdrücke sind Funktoren

```
[=] () {return x + y;}
kann implementiert werden als
unnamed {x,y};
mit

class unnamed {
  int x; int y;
  unnamed (int x_, int y_) : x (x_), y (y_) {}
  int operator () () const {return x + y;}
};
```

Polymorphic Function Wrapper std::function

```
#include <functional>
int k= 8;
std::function<int(int)> f;
f = [k](int i){ return i+k; };
std::cout << f(8); // 16

Kann verwendet werden, um Lambda-Expressions zu speichern.
Andere Beispiele
std::function<int(int,int)>; std::function<void(double)> ...
http://en.cppreference.com/w/cpp/utility/functional/function
```

Beispiel

```
template <typename T>
auto toFunction(std::vector<T> v){
  return [v] (T x) -> double {
    int index = (int)(x+0.5);
    if (index < 0) index = 0;
    if (index >= v.size()) index = v.size()-1;
    return v[index];
  };
}
```

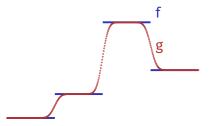
Beispiel

```
auto Gaussian(double mu, double sigma){
    return [mu,sigma](double x) {
        const double a = ( x - mu ) / sigma;
        return std::exp( -0.5 * a * a );
    };
}

template <typename F, typename Kernel>
auto smooth(F f, Kernel kernel){
    return [kernel,f] (auto x) {
        // compute convolution ...
        // and return result
    };
}
```

Beispiel

```
std::vector<double> v {1,2,5,3};
auto f = toFunction(v);
auto k = Gaussian(0,0.1);
auto g = smooth(f,k);
```



438

Zusammenfassung

- Funktoren erlauben die funktionale Programmierung mit C++. Lambdas sind syntaktischer Zucker, der das deutlich vereinfacht
- Mit Funktoren/Lambdas sind klassische Muster aus der funktionalen Programmierung (z.B. map / filter / reduce) auch in C++ möglich.
- In Kombination mit Templates und Typinferenz (auto) können sehr mächtige Funktionen in Variablen gespeichert werden, Funktionen können sogar Funktionen zurückgeben (sog. Funktionen höherer Ordnung).

16. Natürliche Suchbäume

[Ottman/Widmayer, Kap. 5.1, Cormen et al, Kap. 12.1 - 12.3]

Wörterbuchimplementationen

Hashing: Implementierung von Wörterbüchern mit erwartet sehr schnellen Zugriffszeiten.

Nachteile von Hashing: im schlechtesten Fall lineare Zugriffszeit. Manche Operationen gar nicht unterstützt:

- Aufzählen von Schlüssel in aufsteigender Anordnung
- Nächst kleinerer Schlüssel zu gegebenem Schlüssel
- Schlüssel k in vorgegebenem Intervall $k \in [l, r]$

Bäume

Bäume sind

- Verallgemeinerte Listen: Knoten können mehrere Nachfolger haben
- Spezielle Graphen: Graphen bestehen aus Knoten und Kanten. Ein Baum ist ein zusammenhängender, gerichteter, azyklischer Graph.

42

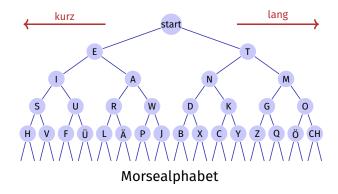
Bäume

Verwendung

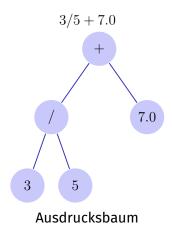
- Entscheidungsbäume: Hierarchische Darstellung von Entscheidungsregeln
- Syntaxbäume: Parsen und Traversieren von Ausdrücken, z.B. in einem Compiler
- Codebäume: Darstellung eines Codes, z.B. Morsealphabet, Huffmann Code
- Suchbäume: ermöglichen effizientes Suchen eines Elementes



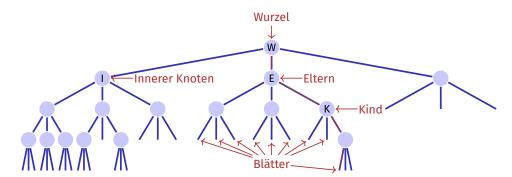
Beispiele



Beispiele



Nomenklatur



- Ordnung des Baumes: Maximale Anzahl Kindknoten, hier: 3
- Höhe des Baumes: maximale Pfadlänge Wurzel Blatt (hier: 4)

Binäre Bäume

Ein binärer Baum ist

- entweder ein Blatt, d.h. ein leerer Baum,
- oder ein innerer Knoten mit zwei Bäumen T_l (linker Teilbaum) und T_r (rechter Teilbaum) als linken und rechten Nachfolger.

In jedem inneren Knoten v wird gespeichert



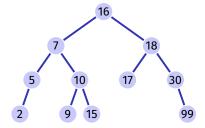
- ein Schlüssel v.key und
- zwei Zeiger v.left und v.right auf die Wurzeln der linken und rechten Teilbäume.

Ein Blatt wird durch den null-Zeiger repräsentiert

Binärer Suchbaum

Ein binärer Suchbaum ist ein binärer Baum, der die Suchbaumeigenschaft erfüllt:

- Jeder Knoten v speichert einen Schlüssel
- Schlüssel im linken Teilbaum v.left kleiner als v.key
- Schlüssel im rechten Teilbaum v.right grösser als v.key



447

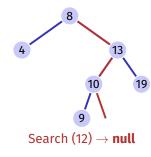
Suchen

 $\begin{array}{ll} \textbf{Input:} \ \, \text{Bin\"{a}rer Suchbaum mit Wurzel } r, \\ \text{Schl\"{u}ssel } k \end{array}$

 $\textbf{Output:} \ \mathsf{Knoten} \ v \ \mathsf{mit} \ v. \mathsf{key} = k \ \mathsf{oder} \ \mathbf{null}$

 $\begin{array}{c|c} \textbf{while } v \neq \textbf{null do} \\ & \textbf{if } k = v. \text{key then} \\ & & \textbf{return } v \\ & \textbf{else if } k < v. \text{key then} \\ & & & v \leftarrow v. \text{left} \\ & \textbf{else} \\ & & & & v \leftarrow v. \text{right} \end{array}$

return null



Höhe eines Baumes

Die Höhe h(T) eines binären Baumes T mit Wurzel r ist gegeben als

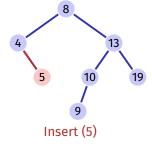
$$h(r) = \begin{cases} 0 & \text{falls } r = \text{null} \\ 1 + \max\{h(r.\text{left}), h(r.\text{right})\} & \text{sonst.} \end{cases}$$

Die Laufzeit der Suche ist somit im schlechtesten Fall $\mathcal{O}(h(T))$

Einfügen eines Schlüssels

Einfügen des Schlüssels k

- \blacksquare Suche nach k.
- Wenn erfolgreich: z.B. Fehlerausgabe
- Wenn erfolglos: Einfügen des Schlüssels am erreichten Blatt.



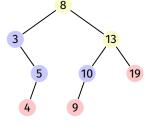
Drei Fälle möglich

■ Knoten hat keine Kinder

Knoten entfernen

- Knoten hat ein Kind
- Knoten hat zwei Kinder

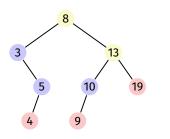
[Blätter zählen hier nicht]



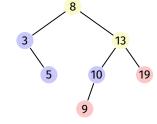
Knoten entfernen

Knoten hat keine Kinder

Einfacher Fall: Knoten durch Blatt ersetzen.





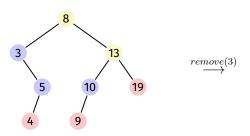


....

Knoten entfernen

Knoten hat ein Kind

Auch einfach: Knoten durch das einzige Kind ersetzen.



10 19

4

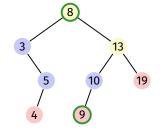
Knoten entfernen

Knoten v hat zwei Kinder

Beobachtung: Der kleinste Schlüssel im rechten Teilbaum v.right (der symmetrische Nachfolger von v)

- ist kleiner als alle Schlüssel in v.right
- ist grösser als alle Schlüssel in v.left
- und hat kein linkes Kind.

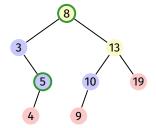
Lösung: ersetze v durch seinen symmetrischen Nachfolger



Aus Symmetriegründen...

Knoten v hat zwei Kinder

Auch möglich: ersetze v durch seinen symmetrischen Vorgänger



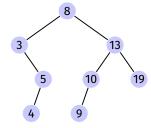
Algorithmus SymmetricSuccessor(v)

 $\begin{array}{l} \textbf{Input:} \ \, \textbf{Knoten} \ \, v \ \, \text{eines} \ \, \text{bin\"aren} \ \, \textbf{Suchbaumes} \\ \textbf{Output:} \ \, \textbf{Symmetrischer} \ \, \textbf{Nachfolger} \ \, \text{von} \ \, v \\ w \leftarrow v. \text{right} \\ x \leftarrow w. \text{left} \\ \textbf{while} \ \, x \neq \textbf{null} \ \, \textbf{do} \\ & | \ \, w \leftarrow x \\ & | \ \, x \leftarrow x. \text{left} \\ \end{array}$

return w

Traversierungsarten

- Hauptreihenfolge (preorder): v, dann $T_{\text{left}}(v)$, dann $T_{\text{right}}(v)$. 8, 3, 5, 4, 13, 10, 9, 19
- Nebenreihenfolge (postorder): $T_{\rm left}(v)$, dann $T_{\rm right}(v)$, dann v.
 4, 5, 3, 9, 10, 19, 13, 8
- Symmetrische Reihenfolge (inorder): $T_{\rm left}(v)$, dann v, dann $T_{\rm right}(v)$. 3, 4, 5, 8, 9, 10, 13, 19



Analyse

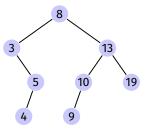
Löschen eines Elementes v aus einem Baum T benötigt $\mathcal{O}(h(T))$ Elementarschritte:

- Suchen von v hat Kosten $\mathcal{O}(h(T))$
- lacksquare Hat v maximal ein Kind ungleich **null**, dann benötigt das Entfernen $\mathcal{O}(1)$
- Das Suchen des symmetrischen Nachfolgers n benötigt $\mathcal{O}(h(T))$ Schritte. Entfernen und Einfügen von n hat Kosten $\mathcal{O}(1)$

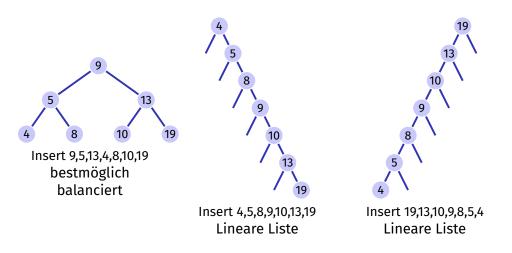
458

Weitere unterstützte Operationen

- Min(T): Auslesen des Minimums in $\mathcal{O}(h)$
- ExtractMin(T): Auslesen und Entfernen des Minimums in $\mathcal{O}(h)$
- List(*T*): Ausgeben einer sortierten Liste der Elemente von *T*
- Join(T_1, T_2): Zusammenfügen zweier Bäume mit $\max(T_1) < \min(T_2)$ in $\mathcal{O}(n)$.



Degenerierte Suchbäume



Probabilistisch

Ein Suchbaum, welcher aus einer zufälligen Sequenz von Zahlen erstellt wird hat erwartete Pfadlänge von $\mathcal{O}(\log n)$.

Achtung: das gilt nur für Einfügeoperation. Wird der Baum zufällig durch Einfügen und Entfernen gebildet, ist die erwartete Pfadlänge $\mathcal{O}(\sqrt{n})$. Balancierte Bäume stellen beim Einfügen und Entfernen (z.B. durch Rotationen) sicher, dass der Baum balanciert bleibt und liefern eine $\mathcal{O}(\log n)$ Worst-Case-Garantie.

17. Heaps

Datenstruktur optimiert zum schnellen Extrahieren von Minimum oder Maximum und Sortieren. [Ottman/Widmayer, Kap. 2.3, Cormen et al, Kap. 6]

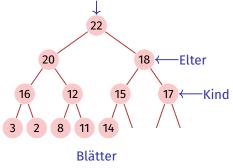
[Max-]Heap*

Binärer Baum mit folgenden Eigenschaften

 vollständig, bis auf die letzte Ebene

2. Lücken des Baumes in der letzten Ebene höchstens rechts.

3. Heap-Bedingung:
Max-(Min-)Heap: Schlüssel eines
Kindes kleiner (grösser) als der
des Elternknotens



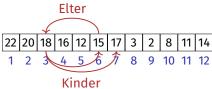
Wurzel

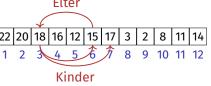
^{*}Heap (Datenstruktur), nicht: wie in "Heap und Stack" (Speicherallokation)

Heap als Array

Baum \rightarrow Array:

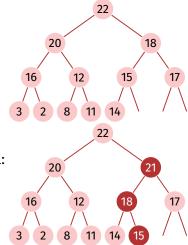
- **Kinder** $(i) = \{2i, 2i + 1\}$
- Elter(i) = |i/2|





Abhängig von Startindex!21

- Füge neues Element an erste freie Stelle ein. Verletzt Heap Eigenschaft potentiell.
- Stelle Heap Eigenschaft wieder her: Sukzessives Aufsteigen.
- $\mathcal{O}(\log n)$



[10] [11] [12]

Höhe eines Heaps

Welche Höhe H(n) hat ein Heap mit n Knoten? Auf der i-ten Ebene eines Binären Baumes befinden sich höchstens 2^i Knoten. Bis auf die letzte Ebene sind alle Ebenen eines Heaps aufgefüllt.

$$H(n) = \min\{h \in \mathbb{N} : \sum_{i=0}^{h-1} 2^i \ge n\}$$

Mit $\sum_{i=0}^{h-1} 2^i = 2^h - 1$:

$$H(n) = \min\{h \in \mathbb{N} : 2^h \ge n+1\},\$$

also

$$H(n) = \lceil \log_2(n+1) \rceil.$$

Einfügen

Anzahl Operationen im schlechtesten Fall:

Algorithmus Aufsteigen(A, m)

Array A mit mindestens m Elementen und Max-Heap-Struktur auf

$$A[1,\ldots,m-1]$$

Output: Array A mit Max-Heap-Struktur auf $A[1, \ldots, m]$.

 $v \leftarrow A[m] // \text{Wert}$

 $c \leftarrow m$ // derzeitiger Knoten (child)

 $p \leftarrow |c/2|$ // Elternknoten (parent)

while c > 1 and v > A[p] do

 $A[c] \leftarrow A[p]$ // Wert Elternknoten \rightarrow derzeitiger Knoten $c \leftarrow p$ // Elternknoten ightarrow derzeitiger Knoten

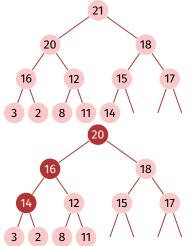
 $p \leftarrow |c/2|$

 $A[c] \leftarrow v \text{ // Wert} \rightarrow \text{Wurzel des (Teil-)Baumes}$

Für Arrays, die bei 0 beginnen: $\{2i, 2i+1\} \rightarrow \{2i+1, 2i+2\}, |i/2| \rightarrow |(i-1)/2|$

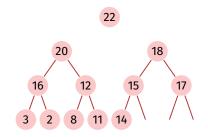
Maximum entfernen

- Ersetze das Maximum durch das unterste rechte Element.
- Stelle Heap Eigenschaft wieder her: Sukzessives Absinken (in Richtung des grösseren Kindes).
- Anzahl Operationen im schlechtesten Fall: $\mathcal{O}(\log n)$



Warum das korrekt ist: Rekursive Heap-Struktur

Ein Heap besteht aus zwei Teilheaps:



47

Algorithmus Versickern(A, i, m)

 $\textbf{Input} \colon \quad \mathsf{Array} \,\, A \,\, \mathsf{mit} \,\, \mathsf{Heapstruktur} \,\, \mathsf{f\"{u}r} \,\, \mathsf{die} \,\, \mathsf{Kinder} \,\, \mathsf{von} \,\, i. \,\, \mathsf{Letztes} \,\, \mathsf{Element} \,\, m.$

Output: Array A mit Heapstruktur für i mit letztem Element m.

while $2i \leq m$ do

Heap Sortieren

A[1,...,n] ist Heap. Solange n>1

- \blacksquare swap(A[1], A[n])
- \blacksquare Versickere(A, 1, n 1);
- $n \leftarrow n-1$

)					ı	
		7	6	4	5	1	2	
Tauschen	\Rightarrow	2	6	4	5	1	7	
Versickern	\Rightarrow	6	5	4	2	1	7	
Tauschen	\Rightarrow	1	5	4	2	6	7	
Versickern	\Rightarrow	5	4	2	1	6	7	
Tauschen	\Rightarrow	1	4	2	5	6	7	
Versickern	\Rightarrow	4	1	2	5	6	7	
Tauschen	\Rightarrow	2	1	4	5	6	7	
Versickern	\Rightarrow	2	1	4	5	6	7	
Tauschen	\Rightarrow	1	2	4	5	6	7	

Heap erstellen

Beobachtung: Jedes Blatt eines Heaps ist für sich schon ein korrekter Heap.

Folgerung: Induktion von unten!

Analyse: Sortieren eines Heaps

Versickere durchläuft maximal $\log n$ Knoten. An jedem Knoten 2 Schlüsselvergleiche. \Rightarrow Heap Sortieren kostet im schlechtesten Fall $2n\log n$ Vergleiche.

Anzahl der Bewegungen vom Heap Sortieren auch $\mathcal{O}(n \log n)$.

Algorithmus HeapSort(A, n)

```
\begin{array}{lll} \textbf{Input:} & \text{Array } A \text{ der L\"ange } n. \\ \textbf{Output:} & A \text{ sortiert.} \\ // & \text{Heap Bauen.} \\ \textbf{for } i \leftarrow n/2 \text{ downto } 1 \text{ do} \\ & & & \text{Versickere}(A,i,n); \\ // & \text{Nun ist } A \text{ ein Heap.} \\ \textbf{for } i \leftarrow n \text{ downto } 2 \text{ do} \\ & & & \text{swap}(A[1],A[i]) \\ & & & \text{Versickere}(A,1,i-1) \\ // & \text{Nun ist } A \text{ sortiert.} \end{array}
```

474

Analyse: Heap bauen

Aufrufe an Versickern: n/2.

Also Anzahl Vergleiche und Bewegungen $v(n) \in \mathcal{O}(n \log n)$.

Versickerpfade sind aber im Mittel viel kürzer:

Wir verwenden, dass $h(n) = \lceil \log_2 n + 1 \rceil = \lceil \log_2 n \rceil + 1$ für n > 0

$$\begin{split} v(n) &= \sum_{l=0}^{\lfloor \log_2 n \rfloor} \underbrace{2^l}_{\text{Anzahl Heaps auf Level l}} \cdot (\underbrace{\lfloor \log_2 n \rfloor + 1 - l}_{\text{H\"ohe Heaps auf Level l}} - 1) = \sum_{k=0}^{\lfloor \log_2 n \rfloor} 2^{\lfloor \log_2 n \rfloor - k} \cdot k \\ &= 2^{\lfloor \log_2 n \rfloor} \cdot \sum_{k=0}^{\lfloor \log_2 n \rfloor} \frac{k}{2^k} \leq n \cdot \sum_{k=0}^{\infty} \frac{k}{2^k} \leq n \cdot 2 \in \mathcal{O}(n) \end{split}$$

$$\text{mit } s(x) := \textstyle \sum_{k=0}^{\infty} kx^k = \frac{x}{(1-x)^2} \quad (0 < x < 1) \text{ und } s(\tfrac{1}{2}) = 2$$

(1 w)

Nachteile

Heapsort: $\mathcal{O}(n \log n)$ Vergleiche und Bewegungen.

Nachteile von Heapsort?

- Wenig Lokalität: per Definition springt Heapsort im sortierten Array umher (Negativer Cache Effekt).
- 🕨 Zwei Vergleiche vor jeder benötigten Bewegung.

18. AVL Bäume

Balancierte Bäume [Ottman/Widmayer, Kap. 5.2-5.2.1, Cormen et al, Kap. Problem 13-3]

478

Ziel

Suchen, Einfügen und Entfernen eines Schlüssels in Baum mit n Schlüsseln, welche in zufälliger Reihenfolge eingefügt wurden im Mittel in $\mathcal{O}(\log_2 n)$ Schritten.

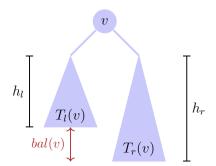
Schlechtester Fall jedoch: $\Theta(n)$ (degenerierter Baum). **Ziel:** Verhinderung der Degenerierung. Künstliches, bei jeder Update-Operation erfolgtes Balancieren eines Baumes Balancierung: garantiere, dass ein Baum mit n Knoten stets eine Höhe von $\mathcal{O}(\log n)$ hat.

Adelson-Venskii und Landis (1962): AVL-Bäume

Balance eines Knotens

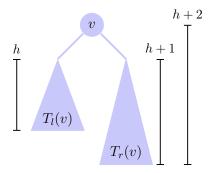
Die Balance eines Knotens v ist definiert als die Höhendifferenz seiner beiden Teilbäume $T_l(v)$ und $T_r(v)$

$$bal(v) := h(T_r(v)) - h(T_l(v))$$



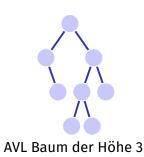
AVL Bedingung

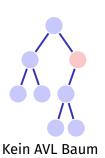
AVL Bedingung: für jeden Knoten v eines Baumes gilt $\mathrm{bal}(v) \in \{-1,0,1\}$



(Gegen-)Beispiele







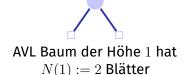
482

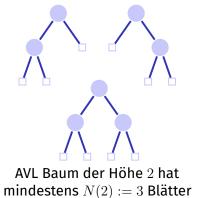
48

Anzahl Blätter

- lacksquare 1. Beobachtung: Ein Suchbaum mit n Schlüsseln hat genau n+1 Blätter. Einfaches Induktionsargument.
 - lacksquare Der Suchbaum mit n=0 Schlüsseln hat m=1 Blätter
 - Wird ein Schlüssel (Knoten) hinzugefügt ($n \to n+1$), so ersetzt er ein Blatt und fügt zwei Blätter hinzu ($m \to m-1+2=m+1$).
- 2. Beobachtung: untere Grenze für Anzahl Blätter eines Suchbaums zu gegebener Höhe erlaubt Abschätzung der maximalen Höhe eines Suchbaums zu gegebener Anzahl Schlüssel.

Untere Grenze Blätter

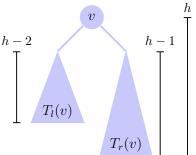




Untere Grenze Blätter für h > 2

- Höhe eines Teilbaums $\geq h-1$.
- Höhe des anderen Teilbaums $\geq h-2$. Minimale Anzahl Blätter N(h) ist

$$N(h) = N(h-1) + N(h-2)$$



Insgesamt gilt $N(h)=F_{h+2}$ mit Fibonacci-Zahlen $F_0:=0$, $F_1:=1$, $F_n:=F_{n-1}+F_{n-2}$ für n>1.

Es gilt²²

$$F_i = \frac{1}{\sqrt{5}} (\phi^i - \hat{\phi}^i)$$

mit den Wurzeln $\phi, \hat{\phi}$ der Gleichung vom goldenen Schnitt $x^2-x-1=0$:

$$\phi = \frac{1+\sqrt{5}}{2} \approx 1.618$$

$$\hat{\phi} = \frac{1-\sqrt{5}}{2} \approx -0.618$$

Fibonacci Zahlen, geschlossene Form

Fibonacci Zahlen, Induktiver Beweis

$$F_i \stackrel{!}{=} \frac{1}{\sqrt{5}} (\phi^i - \hat{\phi}^i) \qquad [*]$$

$$\left(\phi = \frac{1+\sqrt{5}}{2}, \hat{\phi} = \frac{1-\sqrt{5}}{2}\right).$$

- **1.** Klar für i = 0, i = 1.
- 2. Sei i > 2 und Behauptung [*] wahr für alle F_j , j < i.

$$F_i \stackrel{def}{=} F_{i-1} + F_{i-2} \stackrel{[*]}{=} \frac{1}{\sqrt{5}} (\phi^{i-1} - \hat{\phi}^{i-1}) + \frac{1}{\sqrt{5}} (\phi^{i-2} - \hat{\phi}^{i-2})$$

$$= \frac{1}{\sqrt{5}} (\phi^{i-1} + \phi^{i-2}) - \frac{1}{\sqrt{5}} (\hat{\phi}^{i-1} + \hat{\phi}^{i-2}) = \frac{1}{\sqrt{5}} \phi^{i-2} (\phi + 1) - \frac{1}{\sqrt{5}} \hat{\phi}^{i-2} (\hat{\phi} + 1)$$

$$(\phi, \hat{\phi} \text{ erfüllen } x+1=x^2)$$

$$= \frac{1}{\sqrt{5}}\phi^{i-2}(\phi^2) - \frac{1}{\sqrt{5}}\hat{\phi}^{i-2}(\hat{\phi}^2) = \frac{1}{\sqrt{5}}(\phi^i - \hat{\phi}^i).$$

Baumhöhe

Da $|\hat{\phi}| < 1$, gilt insgesamt

$$N(h) \in \Theta\left(\left(\frac{1+\sqrt{5}}{2}\right)^h\right) \subseteq \Omega(1.618^h)$$

und somit

$$N(h) \ge c \cdot 1.618^h$$

$$\Rightarrow h \le 1.44 \log_2 n + c'.$$

Ein AVL Baum ist asymptotisch nicht mehr als 44% höher als ein perfekt balancierter Baum. 23

²²Herleitung mit Erzeugendenfunktionen (Potenzreihen) im Anhang

 $^{^{23}}$ Ein perfekt balancierter Baum hat Höhe $\lceil \log_2 n + 1 \rceil$

Einfügen

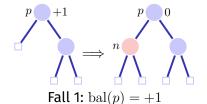
Balancieren

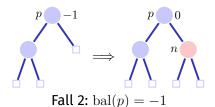
- Speichern der Balance für jeden Knoten
- Baum rebalancieren bei jeder Update-Operation

Neuer Knoten n wird eingefügt:

- Zuerst einfügen wie bei Suchbaum.
- lacksquare Prüfe die Balance-Bedingung für alle Knoten aufsteigend von n zur Wurzel.

Balance am Einfügeort



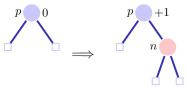


\- /

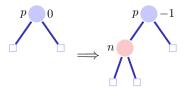
Fertig in beiden Fällen, denn der Teilbaum ist nicht gewachsen.

490

Balance am Einfügeort



Fall 3.1: bal(p) = 0 rechts



Fall 3.2: bal(p) = 0, links

In beiden Fällen noch nicht fertig. Aufruf von upin(p).

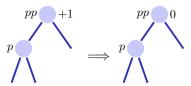
upin(p) - Invariante

Beim Aufruf von upin(p) gilt, dass

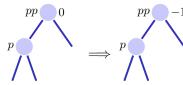
- lacktriangle der Teilbaum ab p gewachsen ist und
- $\blacksquare \ \operatorname{bal}(p) \in \{-1, +1\}$

upin(p)

Annahme: p ist linker Sohn von pp^{24}



Fall 1: bal(pp) = +1, fertig.



Fall 2: bal(pp) = 0, upin(pp)

In beiden Fällen gilt nach der Operation die AVL-Bedingung für den Teilbaum ab pp

upin(p)

Annahme: p ist linker Sohn von pp



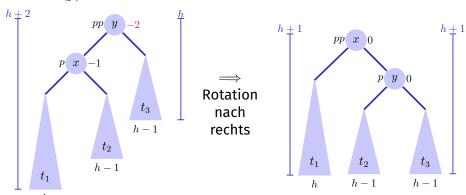
Fall 3: bal(pp) = -1,

Dieser Fall ist problematisch: das Hinzufügen von n im Teilbaum ab pp hat die AVL-Bedingung verletzt. Rebalancieren!

Zwei Fälle bal(p) = -1, bal(p) = +1

Rotationen

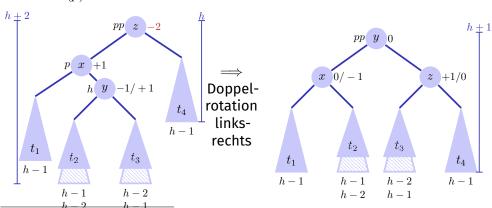
Fall 1.1 bal(p) = -1. ²⁵



 ^{25}p rechter Sohn $\Rightarrow \mathrm{bal}(pp) = \mathrm{bal}(p) = +1$, Linksrotation

Rotationen

Fall 1.2 bal(p) = +1. ²⁶



 ^{26}p rechter Sohn $\Rightarrow \text{bal}(pp) = +1$, bal(p) = -1, Doppelrotation rechts links

 $^{^{-24}}$ Ist p rechter Sohn: symmetrische Fälle unter Vertauschung von +1 und -1

Analyse

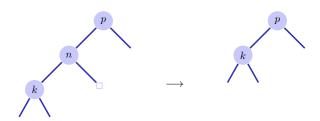
- Höhe des Baumes: $\mathcal{O}(\log n)$.
- Einfügen wie beim binären Suchbaum.
- Balancieren durch Rekursion vom Knoten zur Wurzel. Maximale Pfadlänge $\mathcal{O}(\log n)$.

Das Einfügen im AVL-Baum hat Laufzeitkosten von $\mathcal{O}(\log n)$.

Löschen

Fall 2: Knoten n hat einen inneren Knoten k als Kind

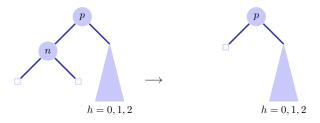
■ Ersetze n durch k. upout(k)



Löschen

Fall 1: Knoten n hat zwei Blätter als Kinder Sei p Elternknoten von $n. \Rightarrow$ Anderer Teilbaum hat Höhe h'=0, 1 oder 2

- h' = 1: bal(p) anpassen.
- h' = 0: bal(p) anpassen. Aufruf upout(p).
- h' = 2: Rebalancieren des Teilbaumes. Aufruf upout(p).



8

Löschen

Fall 3: Knoten n hat zwei inneren Knoten als Kinder

- lacktriangle Ersetze n durch symmetrischen Nachfolger. upout (k)
- Löschen des symmetrischen Nachfolgers wie in Fall 1 oder 2.

upout(p)

Sei pp der Elternknoten von p

(a) p linkes Kind von pp

1. $bal(pp) = -1 \Rightarrow bal(pp) \leftarrow 0$. upout (pp)

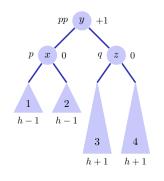
2. $\operatorname{bal}(pp) = 0 \Rightarrow \operatorname{bal}(pp) \leftarrow +1$.

3. $bal(pp) = +1 \Rightarrow n\ddot{a}chste Folien.$

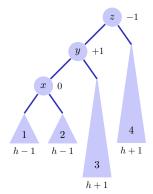
(b) p rechtes Kind von pp: Symmetrische Fälle unter Vertauschung von +1 und -1.

upout(p)

Fall (a).3: $\operatorname{bal}(pp) = +1$. Sei q Bruder von p (a).3.1: $\operatorname{bal}(q) = 0$.27



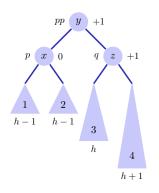
⇒ Linksrotation (y)



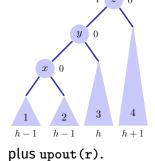
 $\overline{^{27}}$ (b).3.1: $\mathrm{bal}(pp) = -1$, $\mathrm{bal}(q) = -1$, Rechtsrotation.

upout(p)

Fall (a).3: bal(pp) = +1. (a).3.2: bal(q) = +1.²⁸



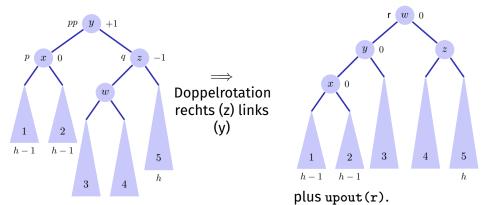
⇒ Linksrotation (y)



 $^{^{28}}$ (b).3.2: $\mathrm{bal}(pp) = -1$, $\mathrm{bal}(q) = +1$, Rechtsrotation+upout

upout(p)

Fall (a).3: bal(pp) = +1. (a).3.3: bal(q) = -1.²⁹



 $\overline{^{29}(\mathsf{b}).3.3: \mathrm{bal}(pp) = -1, \mathrm{bal}(q)} = -1$, Links-Rechts-Rotation + upout

Zusammenfassung

- AVL-Bäume haben asymptotische Laufzeit von $\mathcal{O}(\log n)$ (schlechtester Fall) für das Suchen, Einfügen und Löschen von Schlüsseln
- Einfügen und Löschen ist verhältnismässig aufwändig und für kleine Probleme relativ langsam.

18.5 Anhang

Herleitung einiger mathematischen Formeln

506

[Fibonacci Zahlen: geschlossene Form]

Geschlossene Form der Fibonacci Zahlen: Berechnung über erzeugende Funktionen:

1. Potenzreihenansatz

$$f(x) := \sum_{i=0}^{\infty} F_i \cdot x^i$$

[Fibonacci Zahlen: geschlossene Form]

2. Für Fibonacci Zahlen gilt $F_0=0$, $F_1=1$, $F_i=F_{i-1}+F_{i-2}\ \forall\ i>1$. Daher:

$$f(x) = x + \sum_{i=2}^{\infty} F_i \cdot x^i = x + \sum_{i=2}^{\infty} F_{i-1} \cdot x^i + \sum_{i=2}^{\infty} F_{i-2} \cdot x^i$$

$$= x + x \sum_{i=2}^{\infty} F_{i-1} \cdot x^{i-1} + x^2 \sum_{i=2}^{\infty} F_{i-2} \cdot x^{i-2}$$

$$= x + x \sum_{i=0}^{\infty} F_i \cdot x^i + x^2 \sum_{i=0}^{\infty} F_i \cdot x^i$$

$$= x + x \cdot f(x) + x^2 \cdot f(x).$$

[Fibonacci Zahlen: geschlossene Form]

3. Damit:

$$f(x) \cdot (1 - x - x^2) = x.$$

$$\Leftrightarrow f(x) = \frac{x}{1 - x - x^2} = -\frac{x}{x^2 + x - 1}$$

Mit den Wurzeln $-\phi$ und $-\hat{\phi}$ von x^2+x-1 ,

$$\phi = \frac{1 + \sqrt{5}}{2} \approx 1.6, \qquad \hat{\phi} = \frac{1 - \sqrt{5}}{2} \approx -0.6.$$

gilt $\phi \cdot \hat{\phi} = -1$ und somit

$$f(x) = -\frac{x}{(x+\phi)\cdot(x+\hat{\phi})} = \frac{x}{(1-\phi x)\cdot(1-\hat{\phi}x)}$$

[Fibonacci Zahlen: geschlossene Form]

4. Es gilt:

$$(1 - \hat{\phi}x) - (1 - \phi x) = \sqrt{5} \cdot x.$$

Damit:

$$f(x) = \frac{1}{\sqrt{5}} \frac{(1 - \hat{\phi}x) - (1 - \phi x)}{(1 - \phi x) \cdot (1 - \hat{\phi}x)}$$
$$= \frac{1}{\sqrt{5}} \left(\frac{1}{1 - \phi x} - \frac{1}{1 - \hat{\phi}x} \right)$$

510

[Fibonacci Zahlen: geschlossene Form]

5. Potenzreihenentwicklung von $g_a(x) = \frac{1}{1-a \cdot x}$ ($a \in \mathbb{R}$):

$$\frac{1}{1 - a \cdot x} = \sum_{i=0}^{\infty} a^i \cdot x^i.$$

Sieht man mit Taylor-Entwicklung von $g_a(x)$ um x=0 oder so: Sei $\sum_{i=0}^{\infty} G_i \cdot x^i$ eine Potenzreihenentwicklung von g. Mit der Identität $g_a(x)(1-a\cdot x)=1$ gilt für alle x (im Konvergenzradius)

$$1 = \sum_{i=0}^{\infty} G_i \cdot x^i - a \cdot \sum_{i=0}^{\infty} G_i \cdot x^{i+1} = G_0 + \sum_{i=1}^{\infty} (G_i - a \cdot G_{i-1}) \cdot x^i$$

Für x = 0 folgt $G_0 = 1$ und für $x \neq 0$ folgt dann $G_i = a \cdot G_{i-1} \Rightarrow G_i = a^i$.

[Fibonacci Zahlen: geschlossene Form]

6. Einsetzen der Potenzreihenentwicklung:

$$f(x) = \frac{1}{\sqrt{5}} \left(\frac{1}{1 - \phi x} - \frac{1}{1 - \hat{\phi} x} \right) = \frac{1}{\sqrt{5}} \left(\sum_{i=0}^{\infty} \phi^i x^i - \sum_{i=0}^{\infty} \hat{\phi}^i x^i \right)$$
$$= \sum_{i=0}^{\infty} \frac{1}{\sqrt{5}} (\phi^i - \hat{\phi}^i) x^i$$

Koeffizientenvergleich mit $f(x) = \sum_{i=0}^{\infty} F_i \cdot x^i$ liefert

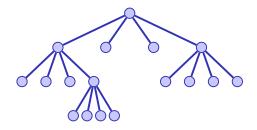
$$F_i = \frac{1}{\sqrt{5}} (\phi^i - \hat{\phi}^i).$$

19. Quadtrees

Quadtrees, Kollisionsdetektion, Bildsegmentierung

Quadtree

Ein Quadtree ist ein Baum der Ordnung 4.



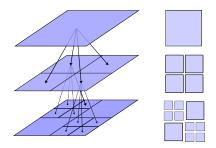
... und ist als solcher nicht besonders interessant, ausser man verwendet ihn zur...

514

-

Quadtree - Interpretation und Nutzen

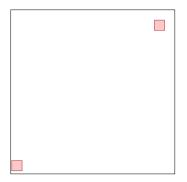
Partitionierung eines zweidimensionalen Bereiches in 4 gleich grosse Teile.



[Analog für drei Dimensionen mit einem Octtree (Baum der Ordnung 8)]

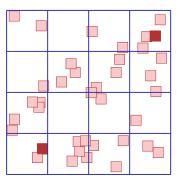
Beispiel 1: Erkennung von Kollisionen

- Objekte in der 2D-Ebene, z.B. Teilchensimulation auf dem Bildschirm.
- Ziel: Erkennen von Kollisionen



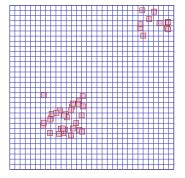
Idee

- Viele Objekte: n^2 Vergleiche (naiv)
- Verbesserung?
- Offensichtlich: keine Kollisionsdetektion für weit entfernte Objekte nötig.
- Was ist "weit entfernt"?
- Gitter ($m \times m$)
- Kollisionsdetektion pro Gitterzelle



Gitter

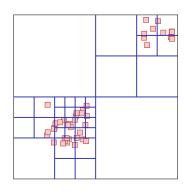
- Gitter hilft oft, aber nicht immer
- Verbesserung?
- Gitter verfeinern?
- Zu viele Gitterzellen!



118

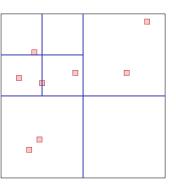
Adaptive Gitter

- Gitter hilft oft, aber nicht immer
- Verbesserung?
- Gitter adaptiv verfeinern!
- Quadtree!



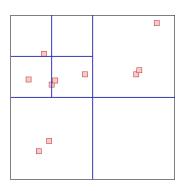
Algorithmus: Einfügen

- Quadtree startet mit einem einzigen Knoten
- Objekte werden zu dem Knoten hinzugefügt. Wenn in einem Knoten zu viele Objekte sind, wird der Knoten geteilt.
- Objekte, die beim Split auf dem Rand zu liegen kommen, werden im höher gelegenen Knoten belassen.

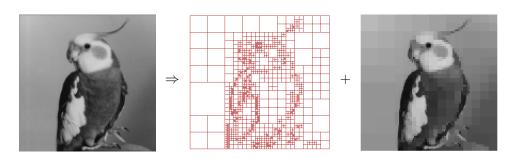


Algorithmus: Kollisionsdetektion

Durchlaufe den Quadtree rekursiv. Für jeden Knoten teste die Kollision der enthaltenen Objekte mit Objekten im selben Knoten oder (rekursiv) enthaltenen Knoten.



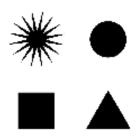
Beispiel 2: Bildsegmentierung

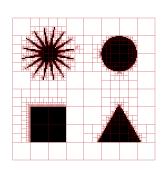


(Mögliche Anwendungen: Kompression, Entrauschen, Kantendetektion)

522

Quadtree auf Einfarbenbild

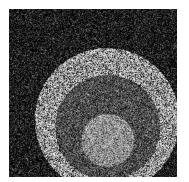


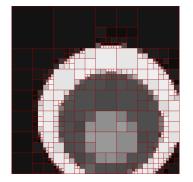


Erzeugung des Quadtree ähnlich wie oben: unterteile Knoten rekursiv bis jeder Knoten nur Pixel einer Farbe enthält.

Quadtree mit Approximation

Wenn mehr als zwei Farbewerte vorhanden sind, wird der Quadtree oft sehr gross. ⇒ Komprimierte Darstellung: *approximiere* das Bild stückweise konstant auf Rechtecken eines Quadtrees.





Stückweise konstante Approximation

(Graustufen-)Bild $\boldsymbol{y} \in \mathbb{R}^S$ auf den Pixelindizes S. ³⁰ Rechteck $r \subset S$.

Ziel: bestimme

$$\arg\min_{v\in\mathbb{R}}\sum_{s\in r}(y_s-v)^2$$

Lösung: das arithmetisches Mittel $\mu_r = \frac{1}{|r|} \sum_{s \in r} y_s$

Zwischenergebnis

Die im Sinne des mittleren quadratischen Fehlers beste Approximation

$$\mu_r = \frac{1}{|r|} \sum_{s \in r} y_s$$

und der dazugehörige Fehler

$$\sum_{s \in r} (y_s - \mu_r)^2 =: \| \boldsymbol{y}_r - \boldsymbol{\mu}_r \|_2^2$$

können nach einer $\mathcal{O}(|S|)$ Tabellierung schnell berechnet werden: Präfixsummen!

Welcher Ouadtree?

Konflikt

- Möglichst nahe an den Daten ⇒ kleine Rechtecke, grosser Quadtree. Extremer Fall: ein Knoten pro Pixel. Approximation = Original
- Möglichst wenige Knoten ⇒ Grosse Rechtecke, kleiner Quadtree Extremfall: ein einziges Rechteck. Approximation = ein Grauwert

Welcher Quadtree?

526

Idee: wähle zwischen Datentreue und Komplexität durch Einführung eines Regularisierungsparameters $\gamma \geq 0$

Wähle Quadtree T mit Blättern $^{\rm 31}$ L(T) so, dass T folgenden Funktion minimiert

$$H_{\gamma}(T, \boldsymbol{y}) := \gamma \cdot \underbrace{|L(T)|}_{\text{Anzahl Blätter}} + \underbrace{\sum_{r \in L(T)} \|y_r - \mu_r\|_2^2}_{\text{Summierter Approximationsfehler aller Blätter}}$$

 $^{^{30}}$ Wir nehmen an, dass S ein Quadrat ist mit Seitenlänge 2^k für ein $k \geq 0$

³¹hier: Blatt = Knoten mit Nullkindern

Regularisierung

Sei T ein Quadtree über einem Rechteck S_T und seien $T_{ll}, T_{lr}, T_{ul}, T_{ur}$ vier mögliche Unterbäume und

$$\widehat{H}_{\gamma}(T,y) := \min_{T} \gamma \cdot |L(T)| + \sum_{r \in L(T)} ||y_r - \mu_r||_2^2$$

Extremfälle:

 $\gamma = 0 \Rightarrow$ Originaldaten; $\gamma \to \infty \Rightarrow$ ein Rechteck

Algorithmus: Minimize(y,r,γ)

Input: Bilddaten $y \in \mathbb{R}^S$, Rechteck $r \subset S$, Regularisierung $\gamma > 0$ Output: $\min_T \gamma |L(T)| + ||y - \mu_{L(T)}||_2^2$

$$\begin{split} & \text{if } |r| = 0 \text{ then return } 0 \\ & m \leftarrow \gamma + \sum_{s \in r} (y_s - \mu_r)^2 \\ & \text{if } |r| > 1 \text{ then} \\ & | & \text{Split } r \text{ into } r_{ll}, r_{lr}, r_{ul}, r_{ur} \\ & m_1 \leftarrow \text{Minimize}(\boldsymbol{y}, r_{ll}, \gamma); \ m_2 \leftarrow \text{Minimize}(\boldsymbol{y}, r_{lr}, \gamma) \\ & m_3 \leftarrow \text{Minimize}(\boldsymbol{y}, r_{ul}, \gamma); \ m_4 \leftarrow \text{Minimize}(\boldsymbol{y}, r_{ur}, \gamma) \\ & m' \leftarrow m_1 + m_2 + m_3 + m_4 \end{split}$$

else $m' \leftarrow \infty$

if
$$m' < m$$
 then $m \leftarrow m'$

 ${\it return}\ m$

Beobachtung: Rekursion

Wenn der (Sub-)Quadtree T nur ein Pixel hat, so kann nicht aufgeteilt werden und es gilt

$$\widehat{H}_{\gamma}(T, \boldsymbol{y}) = \gamma$$

Andernfalls seien

$$M_1 := \gamma + \|\boldsymbol{y}_{S_T} - \boldsymbol{\mu}_{S_T}\|_2^2$$

$$M_2 := \widehat{H}_{\gamma}(T_{ll}, \boldsymbol{y}) + \widehat{H}_{\gamma}(T_{lr}, \boldsymbol{y}) + \widehat{H}_{\gamma}(T_{ul}, \boldsymbol{y}) + \widehat{H}_{\gamma}(T_{ur}, \boldsymbol{y})$$

Dann

$$\widehat{H}_{\gamma}(T,y) = \min\{\underbrace{M_1(T,\gamma,\boldsymbol{y})}_{\text{kein Solit}},\underbrace{M_2(T,\gamma,\boldsymbol{y})}_{\text{Solit}}\}$$

531

Analyse

530

Der Minimierungsalgorithmus über dyadische Partitionen (Quadtree) benötigt $\mathcal{O}(|S|\log |S|)$ Schritte.

Anwendung: Entrauschen (zusätzlich mit Wedgelets)











noised

 $\gamma = 0.003$

 $\gamma = 0.01$









 $\gamma = 0.3$

 $\gamma = 1$

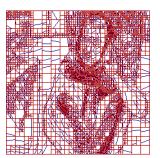
 $\gamma = 3$

 $\gamma = 10$

Erweiterungen: Affine Regression + Wedgelets



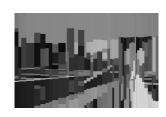


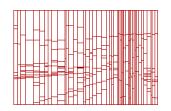


Andere Ideen

kein Quadtree: hierarchisch-eindimensionales Modell (benötigt Dynamic Programming)







19.1 Anhang

Lineare Regression

Das Lernproblem

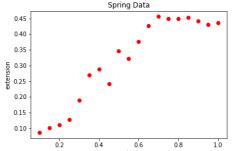
Ausgangslage

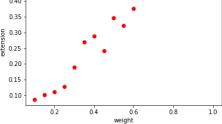
- \blacksquare Wir beobachten N Datenpunkte
- \blacksquare Eingaben: $\boldsymbol{X} = (\boldsymbol{X}_1, \dots, \boldsymbol{X}_N)^{\top}$
- Beobachtete Ausgaben: $\boldsymbol{y} = (y_1, \dots, y_N)^{\top}$
- Annahme: es gibt eine zugrundeliegende Wahrheit



Ziel: Finden einer approximativen Wahrheit $h \approx f$, um Prädiktionen h(x)für neue Datenpunkte x zu machen oder um die Daten zu erklären und beispielsweise komprimiert darzustellen.

Hier
$$\mathcal{X} = \mathbb{R}^d$$
. $\mathcal{Y} = \mathbb{R}$ (Regression).





538

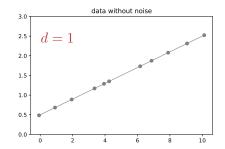
Modell: Lineare Regression

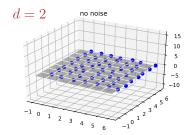
Annahme: Die zugrunde liegende Wahrheit lässt sich darstellen als

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x_1 + \dots + w_d x_d = w_0 + \sum_{i=1}^d w_i x_i.$$

 \Rightarrow Wir suchen w (manchmal auch d).







Trick für vereinfachte Notation

$$\boldsymbol{x} = (x_1, \dots, x_d) \to (\underbrace{x_0}_{\equiv 1}, x_1, \dots, x_d)$$

$$h_{\boldsymbol{w}}(\boldsymbol{x}) = w_0 x_0 + w_1 x_1 + \dots + w_d x_d$$
$$= \sum_{i=0}^d w_i x_i$$
$$= \boldsymbol{w}^\top \boldsymbol{x}$$

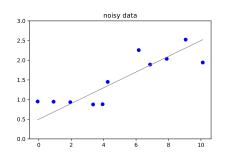
Datenmatrix

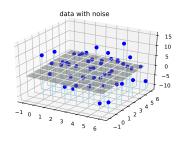
$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{X}_1 \\ \boldsymbol{X}_2 \\ \vdots \\ \boldsymbol{X}_n \end{bmatrix} = \begin{bmatrix} X_{1,0} & X_{1,1} & X_{1,2} & \dots & X_{1,d} \\ X_{2,0} & X_{2,1} & X_{2,2} & \dots & X_{2,d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ X_{n,0} & X_{n,1} & X_{n,2} & \dots & X_{n,d} \end{bmatrix}, \quad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \qquad \boldsymbol{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$$Xw \approx y$$
?

Ungenaue Beobachtungen

Realität: die Daten sind ungenau bzw. das Modell ist nur ein Modell.





Was tun?

Lösung aus der Linearen Algebra

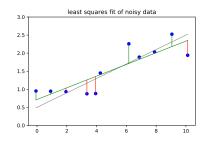
$$\widehat{\boldsymbol{w}} = \underbrace{\left(\boldsymbol{X}^{\top}\boldsymbol{X}\right)^{-1}\boldsymbol{X}^{\top}}_{=:\boldsymbol{X}^{\dagger}}\boldsymbol{y}$$

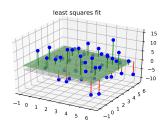
 X^{\dagger} : Moore-Penroe Pseudo-Inverse

Fehlerfunktion

$$E(\boldsymbol{w}) = \sum_{i=1}^{N} (h_{\boldsymbol{w}}(\boldsymbol{X}_i) - y_i)^2$$

Suchen ein $\widehat{\boldsymbol{w}}$ das E(w) minimiert. Linearität von h_w in $w \Rightarrow$ Lösung mit linearer Algebra.





542

Polynome fitten

Geht genauso mit linearer Regression.

$$h_{\mathbf{w}}(x) = w_0 + w_1 x^1 + w_2 x^2 + \dots + w_d x^d = w_0 + \sum_{i=1}^d w_i x^i.$$

Denn $h_{\boldsymbol{w}}(x)$ ist immer noch linear in \boldsymbol{w} !

$$\boldsymbol{X} = \begin{bmatrix} 1 & x_1 & (x_1)^2 & \dots & (x_1)^d \\ 1 & x_2 & (x_2)^2 & \dots & (x_2)^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & (x_n)^2 & \dots & (x_n)^d \end{bmatrix}, \quad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \qquad \boldsymbol{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}$$

Beispiel: Konstante Approximation

$$m{X} = egin{bmatrix} 1 \ 1 \ \vdots \ 1 \end{bmatrix}, \quad m{y} = egin{bmatrix} y_1 \ y_2 \ \vdots \ y_n \end{bmatrix}, \qquad m{w} = m{w} = m{w}_0 \end{bmatrix}$$

$$\widehat{\boldsymbol{w}} = (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} \boldsymbol{y} = \left[\frac{1}{n} \sum y_i\right].$$

Beispiel: Lineare Approximation

$$\boldsymbol{X} = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} \\ 1 & x_2^{(1)} & x_2^{(2)} \\ \vdots & & & \\ 1 & x_n^{(1)} & x_n^{(2)} \end{bmatrix}, \quad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \boldsymbol{w} = \begin{bmatrix} w_0 \end{bmatrix}$$

$$\widehat{\boldsymbol{w}} = \left(\boldsymbol{X}^{\top} \boldsymbol{X}\right)^{-1} \boldsymbol{X}^{\top} \boldsymbol{y} = \begin{bmatrix} N & \sum x_i^{(1)} & \sum x_i^{(2)} \\ \sum x_i^{(1)} & \sum \left(x_i^{(1)}\right)^2 & \sum x_i^{(1)} \cdot x_i^{(2)} \\ \sum x_i^{(2)} & \sum x_i^{(1)} \cdot x_i^{(2)} & \sum \left(x_i^{(2)}\right)^2 \end{bmatrix}^{-1} \cdot \begin{bmatrix} \sum y_i \\ \sum y_i \cdot x_i^{(1)} \\ \sum y_i \cdot x_i^{(2)} \end{bmatrix}$$

547

546

20. Dynamische Programmierung I

Memoisieren, Optimale Substruktur, Überlappende Teilprobleme, Abhängigkeiten, Allgemeines Vorgehen. Beispiele: Fibonacci, Schneiden von Eisenstangen, Längste aufsteigende Teilfolge, längste gemeinsame Teilfolge, Editierdistanz, Matrixkettenmultiplikation, Matrixmultiplikation nach Strassen

[Ottman/Widmayer, Kap. 1.2.3, 7.1, 7.4, Cormen et al, Kap. 15]

Fibonacci Zahlen



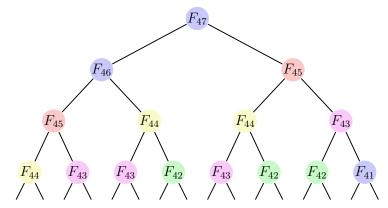
$$F_n := egin{cases} n & \mathsf{wenn} \ n < 2 \ F_{n-1} + F_{n-2} & \mathsf{wenn} \ n \geq 2. \end{cases}$$

Analyse: warum ist der rekursive Algorithmus so langsam.

Algorithmus FibonacciRecursive(n)

```
Input: n \geq 0 Output: n-te Fibonacci Zahl if n < 2 then | f \leftarrow n else | f \leftarrow FibonacciRecursive(n-1) + FibonacciRecursive(n-2) return f
```

Grund, visualisiert



Knoten mit denselben Werten werden (zu) oft ausgewertet.

Analyse

T(n): Anzahl der ausgeführten Operationen.

- n = 0, 1: $T(n) = \Theta(1)$
- $n \ge 2$: T(n) = T(n-2) + T(n-1) + c. $T(n) = T(n-2) + T(n-1) + c \ge 2T(n-2) + c \ge 2^{n/2}c' = (\sqrt{2})^n c'$

Algorithmus ist exponentiell (!) in n.

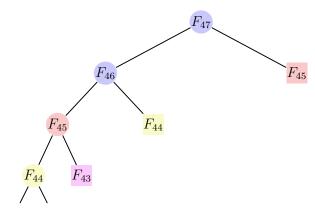
Memoization

550

Memoization (sic) Abspeichern von Zwischenergebnissen.

- Bevor ein Teilproblem gelöst wird, wird Existenz eines entsprechenden Zwischenergebnis geprüft.
- Existiert ein gespeichertes Zwischenergebnis bereits, so wird dieses verwendet.
- Andernfalls wird der Algorithmus ausgeführt und das Ergebnis wird entsprechend gespeichert.

Memoization bei Fibonacci



Rechteckige Knoten wurden bereits ausgewertet.

Analyse

Berechnungsaufwand:

$$T(n) = T(n-1) + c = \dots = \mathcal{O}(n).$$

denn nach dem Aufruf von f(n-1) wurde f(n-2) bereits berechnet. Das lässt sich auch so sehen: Für jedes n wird f(n) maximal einmal rekursiv berechnet. Laufzeitkosten: n Aufrufe mal $\Theta(1)$ Kosten pro Aufruf $n\cdot c\in\Theta(n)$. Die Rekursion verschwindet aus der Berechnung der Laufzeit. Algorithmus benötigt $\Theta(n)$ Speicher. $\Omega(n)$

Algorithmus FibonacciMemoization(n)

```
\begin{array}{l} \textbf{Input: } n \geq 0 \\ \textbf{Output: } n\text{-te Fibonacci Zahl} \\ \textbf{if } n \leq 2 \textbf{ then} \\ \mid f \leftarrow 1 \\ \textbf{else if } \exists \mathsf{memo}[n] \textbf{ then} \\ \mid f \leftarrow \mathsf{memo}[n] \textbf{ else} \\ \mid f \leftarrow \mathsf{FibonacciMemoization}(n-1) + \mathsf{FibonacciMemoization}(n-2) \\ \mid \mathsf{memo}[n] \leftarrow f \\ \textbf{return } f \end{array}
```

554

555

Genauer hingesehen ...

... berechnet der Algorithmus der Reihe nach die Werte F_1 , F_2 , F_3 , ... verkleidet im Top-Down Ansatz der Rekursion.

Man kann den Algorithmus auch gleich Bottom-Up hinschreiben. Das ist charakteristisch für die dynamische Programmierung.

 $[\]overline{\ \ \ }^{32}$ Allerdings benötigt der naive Algorithmus auch $\Theta(n)$ Speicher für die Rekursionsverwaltung.

Algorithmus FibonacciBottomUp(n)

Input: $n \ge 0$

Output: n-te Fibonacci Zahl

$$\begin{split} F[1] &\leftarrow 1 \\ F[2] &\leftarrow 1 \\ \text{for } i \leftarrow 3, \dots, n \text{ do} \end{split}$$

 $|F[i] \leftarrow F[i-1] + F[i-2]$

return F[n]

Dynamische Programmierung: Konsequenz

Identische Teilprobleme werden nur einmal gerechnet

⇒ Resultate werden zwischengespeichert



Laufzeit Wir tauschen Wir Speicherplatz gegen Speicherplatz

Dynamische Programmierung: Idee

- Aufteilen eines komplexen Problems in eine vernünftige Anzahl kleinerer Teilprobleme
- Die Lösung der Teilprobleme wird zur Lösung des komplexeren Problems verwendet
- Identische Teilprobleme werden nur einmal gerechnet

Dynamic Programming: Beschreibung

- 1. Verwalte DP-Tabelle mit Information zu den Teilproblemen. Dimension der Tabelle? Bedeutung der Einträge?
- 2. Berechnung der Randfälle. Welche Einträge hängen nicht von anderen ab?
- 3. Berechnungsreihenfolge bestimen. In welcher Reihenfolge können Einträge berechnet werden, so dass benötigte Einträge jeweils vorhanden sind?
- 4. Auslesen der Lösung.
 Wie kann sich Lösung aus der Tabelle konstruieren lassen?

Laufzeit (typisch) = Anzahl Einträge der Tabelle mal Aufwand pro Eintrag.

Dynamic Programming: Beschreibung am Beispiel

Dimension der Tabelle? Bedeutung der Einträge?
 Tabelle der Grösse n × 1. n-ter Eintrag enthält n-te Fibonacci Zahl.
 Welche Einträge hängen nicht von anderen ab?
 Werte F₁ und F₂ sind unabhängig einfach "berechenbar".
 Berechnungsreihenfolge?
 F_i mit aufsteigenden i.
 Rekonstruktion einer Lösung?
 F_n ist die n-te Fibonacci-Zahl.

Dynamic Programming = Divide-And-Conquer?

- In beiden Fällen ist das Ursprungsproblem (einfacher) lösbar, indem Lösungen von Teilproblemen herangezogen werden können. Das Problem hat **optimale Substruktur**.
- Bei Divide-And-Conquer Algorithmen (z.B. Mergesort) sind Teilprobleme unabhängig; deren Lösungen werden im Algorithmus nur einmal benötigt.
- Beim DP sind Teilprobleme nicht unabhängig. Das Problem hat überlappende Teilprobleme, welche im Algorithmus mehrfach gebraucht werden.
- Damit sie nur einmal gerechnet werden müssen, werden Resultate tabelliert. Dafür darf es zwischen Teilproblemen keine zirkulären Abhängigkeiten geben.

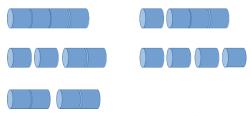
562

Schneiden von Eisenstäben

- Metallstäbe werden zerschnitten und verkauft.
- lacktriangle Metallstäbe der Länge $n\in\mathbb{N}$ verfügbar. Zerschneiden kostet nichts.
- lacksquare Für jede Länge $l\in\mathbb{N}$, $l\leq n$ bekannt: Wert $v_l\in\mathbb{R}^+$
- lacksquare Ziel: Zerschneide die Stange so (in $k\in\mathbb{N}$ Stücke), dass

 $\sum_{i=1}^k v_{l_i}$ maximal unter $\sum_{i=1}^k l_i = n$.

Schneiden von Eisenstäben: Beispiel



Arten, einen Stab der Länge 4 zu zerschneiden (ohne Permutationen)

Länge	0	1	2	3	4	⇒ Bester Schnitt: 3 + 1 mit Wert 10.
Preis	0	2	3	8	9	Dester Schiller 3 . Thirt were 10.

564

Wie findet man den DP Algorithmus

- 0. Genaue Formulierung der gesuchten Lösung
- 1. Definiere Teilprobleme (und bestimme deren Anzahl)
- 2. Raten / Aufzählen (und bestimme die Laufzeit für das Raten)
- 3. Rekursion: verbinde die Teilprobleme
- 4. Memoisieren / Tabellieren. Bestimme die Abhängigkeiten der Teilprobleme
- 5. Lösung des Problems Laufzeit = #Teilprobleme × Zeit/Teilproblem

Algorithmus RodCut(v,n)

$$\overline{^{33}}T(n) = T(n-1) + \sum_{i=0}^{n-2} \overline{T(i)} + c = T(n-1) + (T(n-1) - c) + c = 2T(n-1) \quad (n > 0)$$

Struktur des Problems

- 0. Gesucht: r_n = maximal erreichbarer Wert von (ganzem oder geschnittenem) Stab mit Länge n.
- 1. Teilprobleme: maximal erreichbarer Wert r_k für alle $0 \le k < n$
- 2. Rate Länge des ersten Stückes
- 3. Rekursion

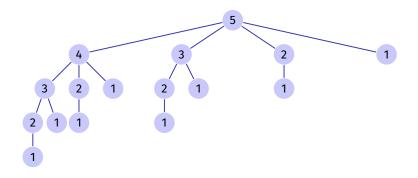
$$r_k = \max\{v_i + r_{k-i} : 0 < i \le k\}, \quad k > 0$$

 $r_0 = 0$

- 4. Abhängigkeit: r_k hängt (nur) ab von den Werten v_i , $l \le i \le k$ und den optimalen Schnitten r_i , i < k
- 5. Lösung in r_n

566

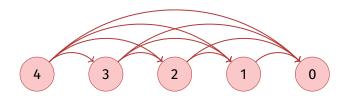
Rekursionsbaum



Algorithmus RodCutMemoized(m, v, n)

Teilproblem-Graph

beschreibt die Abhängigkeiten der Teilprobleme untereinander



und darf keine Zyklen enthalten

570

Konstruktion des optimalen Schnittes

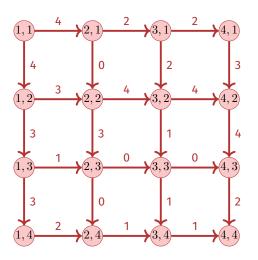
- Während der (rekursiven) Berechnung der optimalen Lösung für jedes $k \leq n$ bestimmt der rekursive Algorithmus die optimale Länge des ersten Stabes
- Speichere die Länge des ersten Stabes für jedes $k \le n$ in einer Tabelle mit n Einträgen.

Bottom-Up Beschreibung am Beispiel

1.	Dimension der Tabelle? Bedeutung der Einträge?
	Tabelle der Grösse $n \times 1$. n -ter Eintrag enthält besten Wert eines Stabes der Länge n .
2.	Welche Einträge hängen nicht von anderen ab?
	Wert r_0 ist 0 .
3.	Berechnungsreihenfolge?
	r_i , $i=1,\ldots,n$.
4.	Rekonstruktion einer Lösung?
	r_n ist der beste Wert für eine Stange der Länge n

Kaninchen!

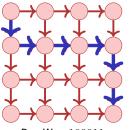
Ein Kaninchen sitzt auf Platz (1,1) eines $n \times n$ Gitters. Es kann nur nach Osten oder nach Süden gehen. Auf jedem Wegstück liegt eine Anzahl Rüben. Wie viele Rüben sammelt das Kaninchen maximal ein?



Kaninchen!

Anzahl mögliche Pfade?

- Auswahl von n-1 Wegen nach Süden aus 2n-2 Wegen insgesamt.
- $\binom{2n-2}{n-1} \in \Omega(2^n)$
- \Rightarrow Naiver Algorithmus hat keine Chance



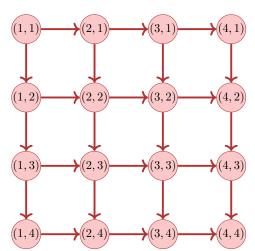
Der Weg 100011 (1:nach Süden, 0:nach Osten)

Rekursion

Gesucht: $T_{1,1}$ = Maximale Anzahl Rüben von (1,1) nach (n,n). Sei $w_{(i,j)-(i',j')}$ Anzahl Rüben auf Kante von (i,j) nach (i',j'). Rekursion (maximale Anzahl Rüben von (i,j) nach (n,n))

$$T_{ij} = \begin{cases} \max\{w_{(i,j)-(i,j+1)} + T_{i,j+1}, w_{(i,j)-(i+1,j)} + T_{i+1,j}\}, & i < n, j < n \\ w_{(i,j)-(i,j+1)} + T_{i,j+1}, & i = n, j < n \\ w_{(i,j)-(i+1,j)} + T_{i+1,j}, & i < n, j = n \\ 0 & i = j = n \end{cases}$$

Teilproblemabhängigkeitsgraph



3/3

Bottom-Up Beschreibung am Beispiel

Dimension der Tabelle? Bedeutung der Einträge?

- Tabelle T der Grösse $n \times n$. Eintrag bei i, j enthält die maximale Anzahl Rüben von (i, j) nach (n, n).
 - Welche Einträge hängen nicht von anderen ab?

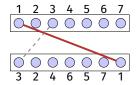
Wert $T_{n,n}$ ist 0.

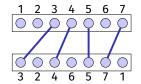
Berechnungsreihenfolge?

- $T_{i,j}$ mit $i=n \setminus 1$ und für jedes $i: j=n \setminus 1$, (oder umgekehrt: $j=n \setminus 1$ und für jedes $i: i = n \setminus 1$).
 - Rekonstruktion einer Lösung?

 $T_{1,1}$ enthält die maximale Anzahl Rüben

Längste aufsteigende Teilfolge (LAT)



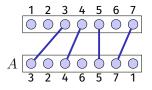


Verbinde so viele passende Anschlüsse wie möglich, ohne dass sich die Anschlüsse kreuzen.

Formalisieren

4.

- Betrachte Folge $A_n = (a_1, \dots, a_n)$.
- Suche eine längste aufsteigende Teilfolge von
- \blacksquare Beispiele aufsteigender Teilfolgen: (3, 4, 5), (2,4,5,7), (3,4,5,7), (3,7).



Verallgemeinerung: Lasse Zahlen ausserhalb von 1, ..., n zu, auch mit Mehrfacheinträgen. (Weitehrhin aber nur strikt aufsteigende Teilfolgen) Beispiel: (2,3,3,3,5,1) mit aufsteigender Teilfolge (2,3,5).

Erster Entwurf

Sei L_i = längste Teilfolge von A_i , ($1 \le i \le n$).

Annahme: LAT L_k von A_k bekannt. Wollen nun LAT L_{k+1} für A_{k+1} berechnen.

Wenn a_{k+1} zu L_k passt, dann $L_{k+1} = L_k \oplus a_{k+1}$?

Gegenbeispiel: $A_5 = (1, 2, 5, 3, 4)$. Sei $A_3 = (1, 2, 5)$ mit $L_3 = A_3$ und $L_4 = A_3$.

Bestimme L_5 aus L_4 ?

So kommen wir nicht weiter: können nicht von L_k auf L_{k+1} schliessen.

Zweiter Entwurf

Sei L_i = längste Teilfolge von A_i , $(1 \le i \le n)$.

Annahme: eine LAT L_j für alle $j \leq k$ bekannt. Wollen nun LAT L_{k+1} für k+1 berechnen.

Betrachte alle passenden $L_{k+1} = L_j \oplus a_{k+1}$ ($j \leq k$) und wähle eine längste solche Folge.

Gegenbeispiel: $A_5 = (1, 2, 5, 3, 4)$. Sei $A_4 = (1, 2, 5, 3)$ mit $L_1 = (1)$, $L_2 = (1, 2)$, $L_3 = (1, 2, 5)$, $L_4 = (1, 2, 5)$. Bestimme L_5 aus L_1, \ldots, L_4 ?

So kommen wir nicht weiter: können nicht von jeweils nur einer beliebigen Lösung L_j auf L_{k+1} schliessen. Wir müssten alle möglichen LAT betrachten. Zu viel!

Dritter Entwurf

Sei $M_{n,i}$ = längste Teilfolge von A_n der Länge i ($1 \le i \le n$)

Annahme: die LAT $M_{k,j}$ für A_k , welche mit kleinstem Element enden seien für alle Längen $1 \le j \le k$ bekannt.

Betrachte nun alle passenden $M_{k,j} \oplus a_{k+1}$ ($j \leq k$) und aktualisiere die Tabelle der längsten aufsteigenden Folgen, welche mit kleinstem Element enden.

32

Dritter Entwurf Beispiel

Beispiel: A = (1, 1000, 1001, 4, 5, 2, 6, 7)

	(-,, -, -, -, -, -, -, -, -,
A	LAT $M_{k,\cdot}$
1	(1)
+ 1000	(1), (1, 1000)
+ 1001	(1), (1, 1000), (1, 1000, 1001)
+4	(1), (1, 4), (1, 1000, 1001)
+ 5	(1), (1,4), (1,4,5)
+2	(1), (1, 2), (1, 4, 5)
+6	(1), (1, 2), (1, 4, 5), (1, 4, 5, 6)
+ 7	(1), (1, 2), (1, 4, 5), (1, 4, 5, 6), (1, 4, 5, 6, 7)

DP Table

- Idee: speichere jeweils nur das letzte Element der aufsteigenden Folge $M_{k,j}$ am Slot j.
- Beispielfolge: 3 2 5 1 6 4
- Problem: Tabelle enthält zum Schluss nicht die Folge, nur den letzten Wert.
- Lösung: Zweite Tabelle mit den Vorgängern.

Index	1	2	3	4	5	6
Wert	3	2	5	1	6	4
Vorgänger	$-\infty$	$-\infty$	2	$-\infty$	5	1

Index	0	1	2	3	4	•••
$(L_j)_j$	-∞	1	4	6	∞	

Dynamic Programming Algorithmus LAT

Dimension der Tabelle? Bedeutung der Einträge?

Zwei Tabellen $T[0,\ldots,n]$ und $V[1,\ldots,n]$. T[j]: letztes Element der aufsteigenden Folge $M_{n,j}$ V[j]: Wert des Vorgängers von a_j . Zu Beginn $T[0] \leftarrow -\infty$, $T[i] \leftarrow \infty \ \forall i>1$

Berechnung eines Eintrags

Einträge in T aufsteigend sortiert. Für jeden Neueintrag a_k binäre Suche nach l, so dass $T[l] < a_k < T[l+1]$. Setze $T[l+1] \leftarrow a_k$. Setze V[k] = T[l].

Dynamic Programming Algorithmus LAT

Berechnungsreihenfolge

Beim Traversieren der Liste werden die Einträge T[k] und V[k] mit aufsteigendem k berechnet.

Rekonstruktion einer Lösung?

4. Suche das grösste l mit $T[l]<\infty$. l ist der letzte Index der LAT. Suche von l ausgehend den Index i< l, so dass $V[l]=a_i$, i ist der Vorgänger von l. Repetiere mit $l\leftarrow i$ bis $T[l]=-\infty$

586

Analyse

■ Berechnung Tabelle:

- Initialisierung: $\Theta(n)$ Operationen
- Berechnung k-ter Eintrag: Binäre Suche auf Positionen $\{1,\ldots,k\}$ plus konstante Anzahl Zuweisungen.

$$\sum_{k=1}^{n} (\log k + \mathcal{O}(1)) = \mathcal{O}(n) + \sum_{k=1}^{n} \log(k) = \Theta(n \log n).$$

Rekonstruktion: Traversiere A von rechts nach links: $\mathcal{O}(n)$.

Somit Gesamtlaufzeit

$$\Theta(n \log n)$$
.

Minimale Editierdistanz

Editierdistanz von zwei Zeichenketten $A_n=(a_1,\ldots,a_n)$, $B_m=(b_1,\ldots,b_m)$. Editieroperationen:

- Einfügen eines Zeichens
- Löschen eines Zeichens
- Änderung eines Zeichens

Frage: Wie viele Editieroperationen sind mindestens nötig, um eine gegebene Zeichenkette ${\cal A}$ in eine Zeichenkette ${\cal B}$ zu überführen.

TIGER ZIGER ZIEGER ZIEGE

Minimale Editierdistanz

Gesucht: Günstigste zeichenweise Transformation $A_n \to B_m$ mit Kosten

Operation	Levenshtein	LGT ³⁴	allgemein
c einfügen	1	1	ins(c)
c löschen	1	1	del(c)
Ersetzen $c \rightarrow c'$	$1(c \neq c')$	$\infty \cdot \mathbb{1}(c \neq c')$	repl(c,c')

Beispiel

Τ	1	G	Ε	R	Т	1	_	G	Ε	R	T	≻Z	+E	-R
Ζ	ı	Ε	G	Ε	Ζ	1	Ε	G	Ε		Z	₹	-E	+R

³⁴Längste gemeinsame Teilfolge – Spezialfall des Editierproblems

DP

4. Abhängigkeiten



- \Rightarrow Berechnung von links oben nach rechts unten. Zeilen- oder Spaltenweise.
- 5. Lösung steht in E(n, m)

DP

- 0. E(n,m) = minimale Anzahl Editieroperationen (ED Kosten) für $a_{1...n} \to b_{1...m}$
- 1. Teilprobleme E(i, j) = ED von $a_{1...i}$, $b_{1...j}$.

 $\#\mathsf{TP} = n \cdot m$ $\mathsf{Kosten}\Theta(1)$

- 2. Raten/Probieren
 - $\blacksquare a_{1..i} \rightarrow a_{1...i-1}$ (löschen)
 - $lacksquare a_{1...i}
 ightarrow a_{1...i} b_j$ (einfügen)
 - $\blacksquare a_{1..i} \rightarrow a_{1...i-1}b_j$ (ersetzen)
- 3. Rekursion

$$E(i, j) = \min \begin{cases} \mathsf{del}(a_i) + E(i - 1, j), \\ \mathsf{ins}(b_j) + E(i, j - 1), \\ \mathsf{repl}(a_i, b_j) + E(i - 1, j - 1) \end{cases}$$

Beispiel (Levenshteinabstand)

$$E[i,j] \leftarrow \min \left\{ E[i-1,j] + 1, E[i,j-1] + 1, E[i-1,j-1] + \mathbb{1}(a_i \neq b_j) \right\}$$

	Ø	Z	I	Ε	G	Ε
Ø	0	1	2	3	4	5
T	1	1	2	3	4	5
- 1	2	2	1	2	3	4
G	3	3	2	2	2	3
Ε	4	4	3	2	3	2
R	5	5	4	3	4 4 3 2 3 3	3

Editierschritte: von rechts unten nach links oben, der Rekursion folgend. Bottom-Up Beschreibung des Algorithmus: Übung

Bottom-Up DP Algorithmus ED

Dimension der Tabelle? Bedeutung der Einträge?

Tabelle $E[0,\ldots,m][0,\ldots,n]$. E[i,j]: Minimaler Editierabstand der Zeichenketten (a_1,\ldots,a_i) und (b_1,\ldots,b_j)

Berechnung eines Eintrags

2. $E[0,i] \leftarrow i \ \forall 0 \leq i \leq m$, $E[j,0] \leftarrow i \ \forall 0 \leq j \leq n$. Berechnung von E[i,j] sonst mit $E[i,j] = \min\{ \operatorname{del}(a_i) + E(i-1,j), \operatorname{ins}(b_j) + E(i,j-1), \operatorname{repl}(a_i,b_j) + E(i-1,j-1) \}$

Bottom-Up DP Algorithmus ED

Berechnungsreihenfolge

Abhängigkeiten berücksichtigen: z.B. Zeilen aufsteigend und innerhalb von Zeilen Spalten aufsteigend.

Rekonstruktion einer Lösung?

Beginne bei j=m, i=n. Falls $E[i,j]=\operatorname{repl}(a_i,b_j)+E(i-1,j-1)$ gilt, gib $a_i\to b_j$ aus und fahre fort mit $(j,i)\leftarrow (j-1,i-1)$; sonst, falls $E[i,j]=\operatorname{del}(a_i)+E(i-1,j)$ gib $\operatorname{del}(a_i)$ aus fahre fort mit $j\leftarrow j-1$; sonst, falls $E[i,j]=\operatorname{ins}(b_j)+E(i,j-1)$, gib $\operatorname{ins}(b_j)$ aus und fahre fort mit $i\leftarrow i-1$. Terminiere für i=0 und j=0.

14

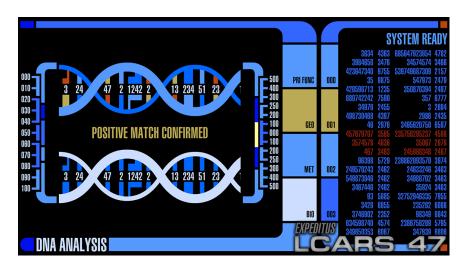
Analyse ED

- lacksquare Anzahl Tabelleneinträge: $(m+1)\cdot (n+1)$.
- lacktriangle Berechnung jeweils mit konstanter Anzahl Zuweisungen und Vergleichen. Anzahl Schritte $\mathcal{O}(mn)$
- Bestimmen der Lösung: jeweils Verringerung von i oder j. Maximal $\mathcal{O}(n+m)$ Schritte.

Laufzeit insgesamt:

 $\mathcal{O}(mn)$.

DNA - Vergleich (Star Trek)



DNA - Vergleich

- DNA besteht aus Sequenzen von vier verschiedenen Nukleotiden Adenin Guanin Thymin Cytosin
- DNA-Sequenzen (Gene) werden mit Zeichenketten aus A, G, T und C beschrieben.
- Ein möglicher Vergleich zweier Gene: Bestimme Längste gemeinsame Teilfolge

Das Problem, die längste gemeinsame Teilfolge zu finden ist ein Spezialfall der minimalen Editierdistanz.

Längste Gemeiname Teilfolge

Teilfolgen einer Zeichenkette:

Teilfolgen(KUH): (), (K), (U), (H), (KU), (KH), (UH), (KUH)

Problem:

- Eingabe: Zwei Zeichenketten $A=(a_1,\ldots,a_m)$, $B=(b_1,\ldots,b_n)$ der Längen m>0 und n>0.
- Gesucht: Eine längste gemeinsame Teilfolge (LGT) von A und B.

370

Längste Gemeiname Teilfolge

Beispiele:

LGT(IGEL,KATZE)=E, LGT(TIGER,ZIEGE)=IGE

Ideen zur Lösung?

Rekursives Vorgehen

Annahme: Lösungen L(i,j) bekannt für $A[1,\ldots,i]$ und $B[1,\ldots,j]$ für alle $1 \le i \le m$ und $1 \le j \le n$, jedoch nicht für i=m und j=n.

Betrachten Zeichen a_m , b_n . Drei Möglichkeiten:

- 1. A wird um ein Leerzeichen erweitert. L(m,n) = L(m,n-1)
- 2. B wird um ein Leerzeichen erweitert. L(m,n) = L(m-1,n)
- 3. $L(m,n)=L(m-1,n-1)+\delta_{mn}$ mit $\delta_{mn}=1$ wenn $a_m=b_n$ und $\delta_{mn}=0$ sonst

Rekursion

 $L(m,n) \leftarrow \max\{L(m-1,n-1) + \delta_{mn}, L(m,n-1), L(m-1,n)\}$ für m,n>0 und Randfälle $L(\cdot,0)=0$, $L(0,\cdot)=0$.

	Ø	Z	1	Ε	0 0 1 2 2	Ε
Ø	0	0	0	0	0	0
Τ	0	0	0	0	0	0
I	0	0	1	1	1	1
G	0	0	1	1	2	2
Ε	0	0	1	2	2	3
R	0	0	1	2	2	3

Dynamic Programming Algorithmus LGT

Dimension der Tabelle? Bedeutung der Einträge?

Tabelle $L[0,\ldots,m][0,\ldots,n]$. L[i,j]: Länge einer LGT der Zeichenketten (a_1,\ldots,a_i) und (b_1,\ldots,b_j)

Berechnung eines Eintrags

2. $L[0,i] \leftarrow 0 \ \forall 0 \le i \le m$, $L[j,0] \leftarrow 0 \ \forall 0 \le j \le n$. Berechnung von L[i,j] sonst mit $L[i,j] = \max(L[i-1,j-1] + \delta_{ij}, L[i,j-1], L[i-1,j])$.

602

Dynamic Programming Algorithmus LGT

Berechnungsreihenfolge

 Abhängigkeiten berücksichtigen: z.B. Zeilen aufsteigend und innerhalb von Zeilen Spalten aufsteigend.

Rekonstruktion einer Lösung?

4. Beginne bei j=m, i=n. Falls $a_i=b_j$ gilt, gib a_i aus und fahre fort mit $(j,i)\leftarrow (j-1,i-1)$; sonst, falls L[i,j]=L[i,j-1] fahre fort mit $j\leftarrow j-1$; sonst, falls L[i,j]=L[i-1,j] fahre fort mit $i\leftarrow i-1$. Terminiere für i=0 oder j=0.

Analyse LGT

- \blacksquare Anzahl Tabelleneinträge: $(m+1)\cdot (n+1)$.
- \blacksquare Berechnung jeweils mit konstanter Anzahl Zuweisungen und Vergleichen. Anzahl Schritte $\mathcal{O}(mn)$
- \blacksquare Bestimmen der Lösung: jeweils Verringerung von i oder j. Maximal $\mathcal{O}(n+m)$ Schritte.

Laufzeit insgesamt:

 $\mathcal{O}(mn)$.

Matrix-Kettenmultiplikation

Aufgabe: Berechnung des Produktes $A_1 \cdot A_2 \cdot \ldots \cdot A_n$ von Matrizen A_1, \ldots, A_n . Matrizenmultiplikation ist assoziativ, d.h. Klammerung kann beliebig gewählt werden.

Ziel: möglichst effiziente Berechnung des Produktes.

Annahme: Multiplikation einer $(r \times s)$ -Matrix mit einer $(s \times u)$ -Matrix hat Kosten $r \cdot s \cdot u$.

606

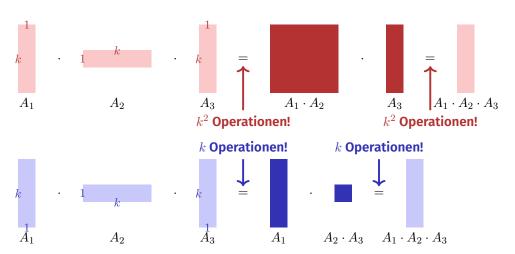
Rekursion

- Annahme, dass die bestmögliche Berechnung von $(A_1 \cdot A_2 \cdots A_i)$ und $(A_{i+1} \cdot A_{i+2} \cdots A_n)$ für jedes i bereits bekannt ist.
- Bestimme bestes *i*, fertig.

 $n \times n$ -Tabelle M. Eintrag M[p,q] enthält Kosten der besten Klammerung von $(A_p \cdot A_{p+1} \cdots A_q)$.

 $M[p,q] \leftarrow \min_{p \leq i < q} (M[p,i] + M[i+1,q] + \text{Kosten letzte Multiplikation})$

Macht das einen Unterschied?



Berechnung der DP-Tabelle

- Randfälle: $M[p,p] \leftarrow 0$ für alle $1 \le p \le n$.
- Berechnung von M[p,q] hängt ab von M[i,j] mit $p \le i \le j \le q$, $(i,j) \ne (p,q)$.

Insbesondere hängt M[p,q] höchstens ab von Einträgen M[i,j] mit i-j < q-p.

Folgerung: Fülle die Tabelle von der Diagonale ausgehend.

Analyse

DP-Tabelle hat n^2 Einträge. Berechung eines Eintrages bedingt Betrachten von bis zu n-1 anderen Einträgen. Gesamtlaufzeit $\mathcal{O}(n^3)$.

Auslesen der Reihenfolge aus M: Übung!

Exkurs: Matrixmultiplikation

Betrachten Multiplikation zweier $n \times n$ -Matrizen. Seien

$$A = (a_{ij})_{1 \le i,j \le n}, B = (b_{ij})_{1 \le i,j \le n}, C = (c_{ij})_{1 \le i,j \le n},$$

 $C = A \cdot B$

dann

$$c_{ij} = \sum_{k=1}^{n} a_{ik} b_{kj}.$$

Naiver Algorithmus benötigt $\Theta(n^3)$ elementare Multiplikationen.

610

611

Divide and Conquer

	В
A	C = AB

		a	b
		c	d
e	f	ea + fc	eb + fd
g	h	ga + hc	gb + hd

Divide and Conquer

- Annahme $n=2^k$.
- Anzahl elementare Multiplikationen: M(n) = 8M(n/2), M(1) = 1.
- Ergibt $M(n) = 8^{\log_2 n} = n^{\log_2 8} = n^3$. Kein Gewinn

a	ь
c	d

e	f	ea + fc	eb + fd
g	h	ga + hc	gb + hd

Strassens Matrixmultiplikation

- Nichttriviale Beobachtung von Strassen (1969): Es genügt die Berechnung der sieben Produkte $A = (e + h) \cdot (a + d), B = (g + h) \cdot a, C = e \cdot (b - d),$ $D = h \cdot (c - a)$, $E = (e + f) \cdot d$, $F = (g - e) \cdot (a + b)$, $G = (f - h) \cdot (c + d)$. Denn: ea + fc = A + D - E + G, eb + fd = C + E,ga + hc = B + D, gb + hd = A - B + C + F.
- Damit ergibt sich M'(n) = 7M(n/2), M'(1) = 1. Also $M'(n) = 7^{\log_2 n} = n^{\log_2 7} \approx n^{2.807}.$
- Schnellster bekannter Algorithmus: $\mathcal{O}(n^{2.37})$

		a	b
		c	d
e	f	ea + fc	eb + fd
g	h	ga + hc	gb + hd

21. Dynamic Programming II

Subset Sum Problem, Rucksackproblem, Greedy Algorithmus vs dynamische Programmierung [Ottman/Widmayer, Kap. 7.2, 7.3, 5.7, Cormen et al, Kap. 15,35.5]

Aufgabe













Teile obige "Gegenstände" so auf zwei Mengen auf, dass beide Mengen den gleichen Wert haben.

Eine Lösung:











Subset Sum Problem

Seien $n \in \mathbb{N}$ Zahlen $a_1, \ldots, a_n \in \mathbb{N}$ gegeben.

Ziel: Entscheide, ob eine Auswahl $I \subseteq \{1, \dots, n\}$ existiert mit

$$\sum_{i \in I} a_i = \sum_{i \in \{1, \dots, n\} \setminus I} a_i.$$

Naiver Algorithmus

Prüfe für jeden Bitvektor $b = (b_1, \dots, b_n) \in \{0, 1\}^n$, ob

$$\sum_{i=1}^{n} b_i a_i \stackrel{?}{=} \sum_{i=1}^{n} (1 - b_i) a_i$$

Schlechtester Fall: n Schritte für jeden der 2^n Bitvektoren b. Anzahl Schritte: $\mathcal{O}(n\cdot 2^n)$.

Algorithmus mit Aufteilung

- Zerlege Eingabe in zwei gleich grosse Teile: $a_1, \ldots, a_{n/2}$ und $a_{n/2+1}, \ldots, a_n$.
- Iteriere über alle Teilmengen der beiden Teile und berechne Teilsummen $S_1^k, \dots, S_{2^{n/2}}^k$ (k = 1, 2).
- Sortiere die Teilsummen: $S_1^k \leq S_2^k \leq \cdots \leq S_{2n/2}^k$.
- Prüfe ob es Teilsummen gibt, so dass $S_i^1 + S_i^2 = \frac{1}{2} \sum_{i=1}^n a_i =: h$
 - Beginne mit $i = 1, j = 2^{n/2}$.

 - Gilt $S_i^1 + S_j^2 = h$ dann fertig Gilt $S_i^1 + S_j^2 > h$ dann $j \leftarrow j 1$ Gilt $S_i^1 + S_j^2 < h$ dann $i \leftarrow i + 1$

Beispiel

Menge $\{1, 6, 2, 3, 4\}$ mit Wertesumme 16 hat 32 Teilmengen. Aufteilung in $\{1, 6\}$, $\{2, 3, 4\}$ ergibt folgende 12 Teilmengen mit Wertesummen:

 \Leftrightarrow Eine Lösung: $\{1, 3, 4\}$

Analyse

- Teilsummegenerierung in jedem Teil: $\mathcal{O}(2^{n/2} \cdot n)$.
- Sortieren jeweils: $\mathcal{O}(2^{n/2}\log(2^{n/2})) = \mathcal{O}(n2^{n/2})$.
- **Z**usammenführen: $\mathcal{O}(2^{n/2})$

Gesamtlaufzeit

$$\mathcal{O}(n \cdot 2^{n/2}) = \mathcal{O}(n(\sqrt{2})^n).$$

Wesentliche Verbesserung gegenüber ganz naivem Verfahren aber immer noch exponentiell!

Dynamische Programmierung

Aufgabe: sei $z=\frac{1}{2}\sum_{i=1}^n a_i$. Suche Auswahl $I\subset\{1,\ldots,n\}$, so dass $\sum_{i\in I} a_i=z$.

DP-Tabelle: $[0,\ldots,n] \times [0,\ldots,z]$ -Tabelle T mit Wahrheitseinträgen. T[k,s] gibt an, ob es eine Auswahl $I_k \subset \{1,\ldots,k\}$ gibt, so dass $\sum_{i \in I_k} a_i = s$. Initialisierung: $T[0,0] = \text{true.} \ T[0,s] = \text{false für } s > 1$. Berechnung:

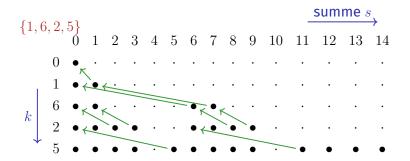
$$T[k,s] \leftarrow \begin{cases} T[k-1,s] & \text{falls } s < a_k \\ T[k-1,s] \lor T[k-1,s-a_k] & \text{falls } s \ge a_k \end{cases}$$

für aufsteigende k und innerhalb k dann s.

Rätselhaftes

Der Algorithmus benötigt $\mathcal{O}(n \cdot z)$ Elementaroperationen. Was ist denn jetzt los? Hat der Algorithmus plötzlich polynomielle Laufzeit?

Beispiel



Auslesen der Lösung: wenn T[k,s]=T[k-1,s] dann a_k nicht benutzt und bei T[k-1,s] weiterfahren, andernfalls a_k benutzt und bei $T[k-1,s-a_k]$ weiterfahren.

623

Aufgelöst

Zahl und keine Anzahl!

622

Der Algorithmus hat nicht unbedingt eine polynomielle Laufzeit. z ist eine

Eingabelänge des Algorithmus \cong Anzahl Bits zur *vernünftigen* Repräsentation der Daten. Bei der Zahl z wäre das $\zeta = \log z$.

Also: Algorithmus benötigt $\mathcal{O}(n\cdot 2^{\zeta})$ Elementaroperationen und hat exponentielle Laufzeit in ζ .

Sollte z allerdings polynomiell sein in n, dann hat der Algorithmus polynomielle Laufzeit in n. Das nennt man pseudopolynomiell.

NP

Man weiss, dass der Subset-Sum Algorithmus zur Klasse der NP-vollständigen Probleme gehört (und somit *NP-schwer* ist).

P: Menge aller in Polynomialzeit lösbarer Probleme.

NP: Menge aller Nichtdeterministisch in Polynomialzeit lösbarer Probleme. Implikationen:

- NP enthält P.
- Probleme in Polynomialzeit verifizierbar.
- Unter der (noch?) unbewiesenen³⁵ Annahme, dass NP ≠ P, gibt es für das Problem keinen Algorithmus mit polynomieller Laufzeit.

Das Rucksackproblem

Wir packen unseren Koffer und nehmen mit ...

- Zahnbürste
- Zahnbürste

Zahnbürste

Hantelset

Luftballon

■ Kaffemaschine

- Kaffemaschine
- Taschenmesser
- Taschenmesser

- Oh jeh zu schwer.
- Ausweis

Ausweis

Hantelset

- Oh jeh zu schwer.
- Oh jeh zu schwer.

Wollen möglichst viel mitnehmen. Manche Dinge sind uns aber wichtiger als andere.

626

626

Rucksackproblem (engl. Knapsack problem)

Gegeben:

- Menge von $n \in \mathbb{N}$ Gegenständen $\{1, \ldots, n\}$.
- Jeder Gegenstand i hat Nutzwert $v_i \in \mathbb{N}$ und Gewicht $w_i \in \mathbb{N}$.
- lacktriangle Maximalgewicht $W \in \mathbb{N}$.
- Bezeichnen die Eingabe mit $E = (v_i, w_i)_{i=1,...,n}$.

Gesucht:

eine Auswahl $I \subseteq \{1, \dots, n\}$ die $\sum_{i \in I} v_i$ maximiert unter $\sum_{i \in I} w_i \leq W$.

Gierige (engl. greedy) Heuristik

Sortiere die Gegenstände absteigend nach Nutzen pro Gewicht v_i/w_i :

Permutation p mit $v_{p_i}/w_{p_i} \ge v_{p_{i+1}}/w_{p_{i+1}}$

Füge Gegenstände in dieser Reihenfolge hinzu ($I \leftarrow I \cup \{p_i\}$), sofern das zulässige Gesamtgewicht dadurch nicht überschritten wird.

Das ist schnell: $\Theta(n\log n)$ für Sortieren und $\Theta(n)$ für die Auswahl. Aber ist es auch gut?

³⁵Die bedeutenste ungelöste Frage der theoretischen Informatik!

Gegenbeispiel zur greedy strategy

$$v_1 = 1$$
 $w_1 = 1$ $v_1/w_1 = 1$ $v_2 = W - 1$ $w_2 = W$ $v_2/w_2 = \frac{W-1}{W}$

Greedy Algorithmus wählt $\{v_1\}$ mit Nutzwert 1. Beste Auswahl: $\{v_2\}$ mit Nutzwert W-1 und Gewicht W. Greedy kann also beliebig schlecht sein.

Dynamic Programming

Unterteile das Maximalgewicht.

Dreidimensionale Tabelle m[i,w,v] ("machbar") aus Wahrheitswerten. m[i,w,v]= true genau dann wenn

- Auswahl der ersten i Teile existiert ($0 \le i \le n$)
- lacktriangle deren Gesamtgewicht höchstens w ($0 \le w \le W$) und
- Nutzen mindestens v ($0 \le v \le \sum_{i=1}^{n} v_i$) ist.

Berechnung der DP Tabelle

Initial

- \blacksquare $m[i, w, 0] \leftarrow$ true für alle $i \ge 0$ und alle $w \ge 0$.
- \blacksquare $m[0, w, v] \leftarrow$ false für alle $w \ge 0$ und alle v > 0.

Berechnung

$$m[i,w,v] \leftarrow \begin{cases} m[i-1,w,v] \lor m[i-1,w-w_i,v-v_i] & \text{falls } w \ge w_i \text{ und } v \ge v_i \\ m[i-1,w,v] & \text{sonst.} \end{cases}$$

aufsteigend nach i und für festes i aufsteigend nach w und für festes i und w aufsteigend nach v.

Lösung: Grösstes v, so dass $m[i,w,v]={\rm true}\;{\rm für}\;{\rm ein}\;i\;{\rm und}\;w$.

Beobachtung

Nach der Definition des Problems gilt offensichtlich, dass

- für m[i,w,v]= true gilt: m[i',w,v]= true $\forall i'\geq i$, m[i,w',v]= true $\forall w'\geq w$, m[i,w,v']= true $\forall v'\leq v$.
- für m[i, w, v] = false gilt: m[i', w, v] = false $\forall i' \leq i$, m[i, w', v] = false $\forall w' \leq w$, m[i, w, v'] = false $\forall v' \geq v$.

Das ist ein starker Hinweis darauf, dass wir keine 3d-Tabelle benötigen.

DP Tabelle mit 2 Dimensionen

Tabelleneintrag t[i, w] enthält statt Wahrheitswerten das jeweils grösste v, das erreichbar ist³⁶ mit

- den Gegenständen 1, ..., i ($0 \le i \le n$)
- bei höchstem zulässigen Gewicht w ($0 \le w \le W$).

Berechnung

Initial

 \bullet $t[0,w] \leftarrow 0$ für alle $w \geq 0$.

Berechnung

$$t[i,w] \leftarrow \begin{cases} t[i-1,w] & \text{falls } w < w_i \\ \max\{t[i-1,w],t[i-1,w-w_i]+v_i\} & \text{sonst.} \end{cases}$$

aufsteigend nach i und für festes i aufsteigend nach w. Lösung steht in t[n,w]

Beispiel

Auslesen der Lösung: wenn t[i,w]=t[i-1,w] dann Gegenstand i nicht benutzt und bei t[i-1,w] weiterfahren, andernfalls benutzt und bei $t[i-1,s-w_i]$ weiterfahren.

Analyse

634

Die beiden Algorithmen für das Rucksackproblem haben eine Laufzeit in $\Theta(n\cdot W\cdot \sum_{i=1}^n v_i)$ (3d-Tabelle) und $\Theta(n\cdot W)$ (2d-Tabelle) und sind beide damit pseudopolynomiell, liefern aber das bestmögliche Resultat. Der greedy Algorithmus ist sehr schnell, liefert aber unter Umständen beliebig schlechte Resultate.

Im folgenden beschäftigen wir uns mit einer Lösung dazwischen.

,

³⁶So etwas ähnliches hätten wir beim Subset Sum Problem auch machen können, um die dünnbesetzte Tabelle etwas zu verkleinern

22. Dynamic Programming III

FPTAS [Ottman/Widmayer, Kap. 7.2, 7.3, Cormen et al, Kap. 15,35.5], Optimale Suchbäume [Ottman/Widmayer, Kap. 5.7]

 $\sum_{i \in I} v_i \geq (1 - \varepsilon) \sum_{i \in I_{\mathsf{opt}}} v_i.$

Sei ein $\varepsilon \in (0,1)$ gegeben. Sei I_{opt} eine bestmögliche Auswahl.

Summe der Gewichte darf W natürlich in keinem Fall überschreiten.

Andere Formulierung des Algorithmus

Bisher: Gewichtsschranke $w \to \text{maximaler Nutzen } v$ Umkehrung Nutzen $v \to \text{minimales Gewicht } w$

- \Rightarrow **Alternative Tabelle:** g[i, v] gibt das minimale Gewicht an, welches
- lacksquare eine Auswahl der ersten i Gegenstände ($0 \le i \le n$) hat, die
- lacksquare einen Nutzen von genau v aufweist ($0 \le v \le \sum_{i=1}^n v_i$).

Berechnung

Approximation

Suchen eine gültige Auswahl I mit

Initial

- $\blacksquare g[0,0] \leftarrow 0$
- $\blacksquare g[0,v] \leftarrow \infty$ (Nutzen v kann mit 0 Gegenständen nie erreicht werden.).

Berechnung

$$g[i,v] \leftarrow \begin{cases} g[i-1,v] & \text{falls } v < v_i \\ \min\{g[i-1,v], g[i-1,v-v_i] + w_i\} & \text{sonst.} \end{cases}$$

aufsteigend nach i und für festes i aufsteigend nach v. Lösung ist der grösste Index v mit $g[n,v] \leq w$.

Beispiel

$$E = \{(2,3), (4,5), (1,1)\}$$

$$0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9$$

$$\emptyset \quad 0 \leftarrow \infty \quad \infty$$

$$\downarrow (2,3) \quad 0 \leftarrow \infty \quad \infty \quad 2 \leftarrow \infty \quad \infty \quad \infty \quad \infty$$

$$\downarrow (4,5) \quad 0 \leftarrow \infty \quad \infty \quad 2 \leftarrow \infty \quad 4 \leftarrow \infty \quad \infty \quad 6 \leftarrow \infty$$

$$\downarrow (1,1) \quad 0 \quad 1 \quad \infty \quad 2 \quad 3 \quad 4 \quad 5 \quad \infty \quad 6 \quad 7$$

Auslesen der Lösung: wenn g[i,v]=g[i-1,v] dann Gegenstand i nicht benutzt und bei g[i-1,v] weiterfahren, andernfalls benutzt und bei $g[i-1,b-v_i]$ weiterfahren.

Der Approximationstrick

Pseudopolynomielle Laufzeit wird polynomiell, wenn vorkommenden Werte in Polynom der Eingabelänge beschränkt werden können. Sei K>0 geeignet gewählt. Ersetze die Nutzwerte v_i durch "gerundete Werte" $\tilde{v_i}=\lfloor v_i/K \rfloor$ und erhalte eine neue Eingabe $E'=(w_i,\tilde{v_i})_{i=1\dots n}$. Wenden nun den Algorithmus auf Eingabe E' mit derselben Gewichtsschranke W an.

642

Idee

Beispiel K=5

Eingabe Nutzwerte

$$1, 2, 3, 4, 5, 6, 7, 8, 9, 10, \dots, 98, 99, 100$$
 \rightarrow
 $0, 0, 0, 0, 1, 1, 1, 1, 1, 2, \dots, 19, 19, 20$

Offensichtlich weniger unterschiedliche Nutzwerte

Eigenschaften des neuen Algorithmus

- lacksquare Auswahl von Gegenständen aus E' ist genauso gültig wie die aus E. Gewicht unverändert!
- Laufzeit des Algorithmus ist beschränkt durch $\mathcal{O}(n^2 \cdot v_{\max}/K)$ $(v_{\max} := \max\{v_i | 1 \le i \le n\})$

Wie gut ist die Approximation?

Es gilt

$$v_i - K \le K \cdot \left| \frac{v_i}{K} \right| = K \cdot \tilde{v_i} \le v_i$$

Sei I'_{opt} eine optimale Lösung von E'. Damit

$$\left(\sum_{i \in I_{\mathsf{Opt}}} v_i \right) - n \cdot K \overset{|I_{\mathsf{Opt}}| \le n}{\le} \sum_{i \in I_{\mathsf{Opt}}} (v_i - K) \le \sum_{i \in I_{\mathsf{Opt}}} (K \cdot \tilde{v_i}) = K \sum_{i \in I_{\mathsf{Opt}}} \tilde{v_i}$$

$$\le K \sum_{I_{\mathsf{Opt}}' \mathsf{Optimal}} K \sum_{i \in I_{\mathsf{Opt}}'} \tilde{v_i} = \sum_{i \in I_{\mathsf{Opt}}'} K \cdot \tilde{v_i} \le \sum_{i \in I_{\mathsf{Opt}}'} v_i.$$

Wahl von K

Forderung:

$$\sum_{i \in I'} v_i \ge (1 - \varepsilon) \sum_{i \in I_{\mathsf{opt}}} v_i.$$

Ungleichung von oben:

$$\sum_{i \in I_{\mathsf{opt}}'} v_i \ge \left(\sum_{i \in I_{\mathsf{opt}}} v_i\right) - n \cdot K$$

Also:
$$K = \varepsilon \frac{\sum_{i \in I_{\mathsf{opt}}} v_i}{n}$$
.

646

Wahl von K

Wähle $K=arepsilon rac{\sum_{i\in I_{\mathrm{opt}}}v_i}{n}$. Die optimale Summe ist aber unbekannt, daher wählen wir $K'=arepsilon rac{v_{\mathrm{max}}}{n}$. 37

Es gilt $v_{\max} \leq \sum_{i \in I_{\text{opt}}} v_i$ und somit $K' \leq K$ und die Approximation ist sogar etwas besser.

Die Laufzeit des Algorithmus ist beschränkt durch

$$\mathcal{O}(n^2 \cdot v_{\text{max}}/K') = \mathcal{O}(n^2 \cdot v_{\text{max}}/(\varepsilon \cdot v_{\text{max}}/n)) = \mathcal{O}(n^3/\varepsilon).$$

FPTAS

Solche Familie von Algorithmen nennt man Approximationsschema: die Wahl von ε steuert Laufzeit und Approximationsgüte.

Die Laufzeit $\mathcal{O}(n^3/\varepsilon)$ ist ein Polynom in n und in $\frac{1}{\varepsilon}$. Daher nennt man das Verfahren auch ein voll polynomielles Approximationsschema FPTAS - Fully Polynomial Time Approximation Scheme

 $[\]overline{\ \ \ }^{37}$ Wir können annehmen, dass vorgängig alle Gegenstände i mit $w_i>W$ entfernt wurden.

22.2 Optimale Suchbäume

Optimale binäre Suchbäume

Gegeben: Suchwahrscheinlichkeiten p_i zu jedem Schlüssel k_i ($i=1,\ldots,n$) und q_i zu jedem Intervall d_i ($i=0,\ldots,n$) zwischen Suchschlüsseln eines binären Suchbaumes. $\sum_{i=1}^n p_i + \sum_{i=0}^n q_i = 1$.

Gesucht: Optimaler Suchbaum T mit Schlüsseltiefen $\operatorname{depth}(\cdot)$, welcher die erwarteten Suchkosten

$$C(T) = \sum_{i=1}^{n} p_i \cdot (\operatorname{depth}(k_i) + 1) + \sum_{i=0}^{n} q_i \cdot (\operatorname{depth}(d_i) + 1)$$
$$= 1 + \sum_{i=1}^{n} p_i \cdot \operatorname{depth}(k_i) + \sum_{i=0}^{n} q_i \cdot \operatorname{depth}(d_i)$$

minimiert.

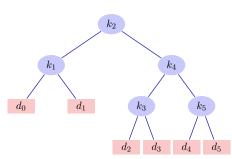
650

Beispiel

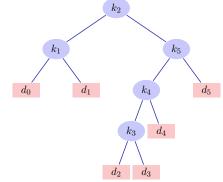
Erwartete Häufigkeiten

i	0	1	2	3	4	5	
				0.05			
q_i	0.05	0.10	0.05	0.05	0.05	0.10	

Beispiel



Suchbaum mit erwarteten Kosten 2.8

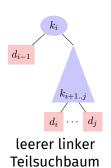


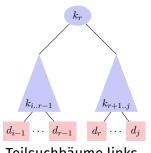
Suchbaum mit erwarteten Kosten 2.75

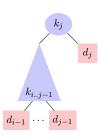
Struktur eines optimalen Suchbaumes

- Teilsuchbaum mit Schlüsseln k_i, \ldots, k_j und Intervallschlüsseln d_{i-1}, \ldots, d_j muss für das entsprechende Teilproblem optimal sein. ³⁸
- Betrachten aller Teilsuchbäume mit Wurzel k_r , $i \le r \le j$ und optimalen Teilbäumen k_i, \ldots, k_{r-1} und k_{r+1}, \ldots, k_j

Teilsuchbäume







Teilsuchbäume links und rechts nichtleer

leerer rechter Teilsuchbaum

654

Erwartete Suchkosten

Sei depth $_T(k)$ die Tiefe des Knotens im Teilbaum T. Sei k_r die Wurzel eines Teilbaumes T_r und T_{L_r} und T_{R_r} der linke und rechte Teilbaum von T_r . Dann

$$depth_T(k_i) = depth_{T_{L_r}}(k_i) + 1, (i < r)$$

$$depth_T(k_i) = depth_{T_{R_r}}(k_i) + 1, (i > r)$$

Erwartete Suchkosten

Seien e[i,j] die Kosten eines optimalen Suchbaumes mit Knoten k_i, \ldots, k_j . Basisfall: e[i,i-1], erwartete Suchkosten d_{i-1}

Sei $w(i,j) = \sum_{l=i}^{j} p_l + \sum_{l=i-1}^{j} q_l$.

Wenn k_r die Wurzel eines optimalen Teilbaumes mit Schlüsseln k_i,\ldots,k_j , dann

$$\begin{split} e[i,j] &= p_r + (e[i,r-1] + w(i,r-1)) + (e[r+1,j] + w(r+1,j)) \\ \text{mit } w(i,j) &= w(i,r-1) + p_r + w(r+1,j) \text{:} \\ e[i,j] &= e[i,r-1] + e[r+1,j] + w(i,j). \end{split}$$

033

³⁸Das übliche Argument: wäre er nicht optimal, könnte er durch eine bessere Lösung ersetzt werden, welche die Gesamtlösung verbessert.

Dynamic Programming

$$e[i,j] = \begin{cases} q_{i-1} & \text{falls } j=i-1, \\ \min_{i \leq r \leq j} \{e[i,r-1] + e[r+1,j] + w[i,j]\} & \text{falls } i \leq j \end{cases}$$

Berechnung

Tabellen $e[1\ldots n+1,0\ldots n], w[1\ldots n+1,0\ldots m], r[1\ldots n,1\ldots n]$ Initial $e[i,i-1]\leftarrow q_{i-1}, w[i,i-1]\leftarrow q_{i-1}$ für alle $1\leq i\leq n+1$.

Berechnung

658

$$w[i,j] = w[i,j-1] + p_j + q_j$$

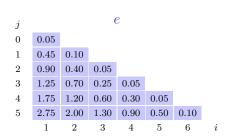
$$e[i,j] = \min_{i \le r \le j} \{e[i,r-1] + e[r+1,j] + w[i,j]\}$$

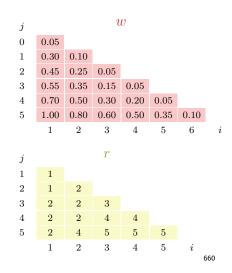
$$r[i,j] = \arg\min_{i \le r \le j} \{e[i,r-1] + e[r+1,j] + w[i,j]\}$$

für Intervalle [i,j] mit ansteigenden Längen $l=1,\ldots,n$, jeweils für $i=1,\ldots,n-l+1$. Resultat steht in e[1,n], Rekonstruktion via r. Laufzeit $\Theta(n^3)$.

Beispiel

i	0	1	2	3	4	5
p_i	0.05	0.15	0.10	0.05	0.10	0.20
q_i	0.05	0.10	0.05	0.05	0.05	0.10





23. Gierige (Greedy) Algorithmen

Gebrochenes Rucksack Problem, Huffman Coding [Cormen et al, Kap. 16.1, 16.3]

659

Gierige Auswahl

Ein rekursiv lösbares Optimierungsproblem kann mit einem **gierigen** (greedy) Algorithmus gelöst werden, wenn es die folgende Eigenschaften hat:

- Das Problem hat **optimale Substruktur**: die Lösung eines Problems ergibt sich durch Kombination optimaler Teillösungen.
- Es gilt die greedy choice property: Die Lösung eines Problems kann konstruiert werden, indem ein lokales Kriterium herangezogen wird, welches nicht von der Lösung der Teilprobleme abhängt.

Beispiele: Gebrochenes Rucksackproblem, Huffman-Coding (s.u.) Gegenbeispiele: Rucksackproblem. Optimaler binärer Suchbaum.

Das Gebrochene Rucksackproblem

Menge von $n\in\mathbb{N}$ Gegenständen $\{1,\ldots,n\}$ gegeben. Jeder Gegenstand i hat Nutzwert $v_i\in\mathbb{N}$ und Gewicht $w_i\in\mathbb{N}$. Das Maximalgewicht ist gegeben als $W\in\mathbb{N}$. Bezeichnen die Eingabe mit $E=(v_i,w_i)_{i=1,\ldots,n}$.

Gesucht: Anteile $0 \le q_i \le 1$ ($1 \le i \le n$) die die Summe $\sum_{i=1}^n q_i \cdot v_i$ maximieren unter $\sum_{i=1}^n q_i \cdot w_i \le W$.

62

Gierige (Greedy) Heuristik

Sortiere die Gegenstände absteigend nach Nutzen pro Gewicht v_i/w_i .

Annahme
$$v_i/w_i \geq v_{i+1}/w_{i+1}$$
 Sei $j = \max\{0 \leq k \leq n: \sum_{i=1}^k w_i \leq W\}$. Setze

- $\blacksquare q_i = 1$ für alle $1 \le i \le j$.
- $q_{j+1} = \frac{W \sum_{i=1}^{j} w_i}{w_{i+1}}$.
- $q_i = 0$ für alle i > j + 1.

Das ist schnell: $\Theta(n \log n)$ für Sortieren und $\Theta(n)$ für die Berechnung der q_i .

Korrektheit

Annahme: Optimale Lösung (r_i) ($1 \le i \le n$).

Der Rucksack wird immer ganz gefüllt: $\sum_i r_i \cdot w_i = \sum_i q_i \cdot w_i = W$.

Betrachte k: kleinstes i mit $r_i \neq q_i$. Die gierige Heuristik nimmt per

Definition so viel wie möglich: $q_k > r_k$. Sei $x = q_k - r_k > 0$.

Konstruiere eine neue Lösung (r_i') : $r_i' = r_i \forall i < k$. $r_k' = q_k$. Entferne Gewicht $\sum_{i=k+1}^n \delta_i = x \cdot w_k$ von den Gegenständen k+1 bis n. Das geht, denn

 $\sum_{i=k}^{n} r_i \cdot w_i = \sum_{i=k}^{n} q_i \cdot w_i.$

Korrektheit

$$\sum_{i=k}^{n} r_i' v_i = r_k v_k + x w_k \frac{v_k}{w_k} + \sum_{i=k+1}^{n} (r_i w_i - \delta_i) \frac{v_i}{w_i}$$

$$\geq r_k v_k + x w_k \frac{v_k}{w_k} + \sum_{i=k+1}^{n} r_i w_i \frac{v_i}{w_i} - \delta_i \frac{v_k}{w_k}$$

$$= r_k v_k + x w_k \frac{v_k}{w_k} - x w_k \frac{v_k}{w_k} + \sum_{i=k+1}^{n} r_i w_i \frac{v_i}{w_i} = \sum_{i=k}^{n} r_i v_i.$$

Also ist (r'_i) auch optimal. Iterative Anwendung dieser Idee erzeugt die Lösung (q_i) .

Huffman-Codierungen

- Betrachten Präfixcodes: kein Codewort kann mit einem anderen Codewort beginnen.
- Präfixcodes können im Vergleich mit allen Codes die optimale Datenkompression erreichen (hier ohne Beweis).
- Codierung: Verkettung der Codewörter ohne Zwischenzeichen (Unterschied zum Morsen!)

$$affe \rightarrow 0 \cdot 1100 \cdot 1100 \cdot 1101 \rightarrow 0110011001101$$

■ Decodierung einfach da Präfixcode $0110011001101 \rightarrow 0 \cdot 1100 \cdot 1100 \cdot 1101 \rightarrow affe$

Huffman-Codierungen

Ziel: Speicherplatzeffizientes Speichern einer Folge von Zeichen mit einem binären Zeichencode aus Codewörtern.

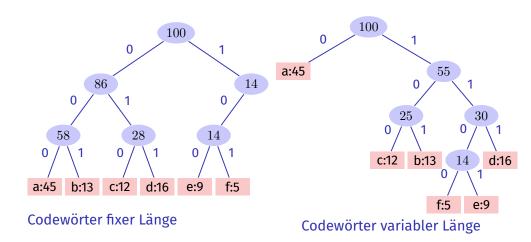
Beispiel

File aus 100.000 Buchstaben aus dem Alphabet $\{a, \ldots, f\}$

	a	b	С	d	е	f
Häufigkeit (Tausend)	45	13	12	16	9	5
Codewort fester Länge	000	001	010	011	100	101
Codewort variabler Länge	0	101	100	111	1101	1100

Speichergrösse (Code fixe Länge): 300.000 bits. Speichergrösse (Code variabler Länge): 224.000 bits.

Codebäume



Eigenschaften der Codebäume

- Optimale Codierung eines Files wird immer durch vollständigen binären Baum dargestellt: jeder innere Knoten hat zwei Kinder.
- Sei C die Menge der Codewörter, f(c) die Häufigkeit eines Codeworts c und $d_T(c)$ die Tiefe eines Wortes im Baum T. Definieren die Kosten eines Baumes als

$$B(T) = \sum_{c \in C} f(c) \cdot d_T(c).$$

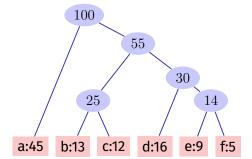
(Kosten = Anzahl Bits des codierten Files)

Bezeichnen im folgenden einen Codebaum als optimal, wenn er die Kosten minimiert.

Algorithmus Idee

Baum Konstruktion von unten nach oben

- Starte mit der Menge C der Codewörter
- Ersetze iterativ die beiden Knoten mit kleinster Häufigkeit durch ihren neuen Vaterknoten.



670

671

Algorithmus Huffman(C)

 $\mathbf{Input} \colon \quad \mathsf{Codew\"{o}rter} \ c \in C$

return ExtractMin(Q)

Output: Wurzel eines optimalen Codebaums

$$\begin{array}{l} n \leftarrow |C| \\ Q \leftarrow C \\ \text{for } i = 1 \text{ to } n-1 \text{ do} \\ & \text{Alloziere neuen Knoten } z \\ & z.\text{left} \leftarrow \text{ExtractMin}(Q) \\ & z.\text{right} \leftarrow \text{ExtractMin}(Q) \\ & z.\text{freq} \leftarrow z.\text{left.freq} + z.\text{right.freq} \\ & \text{Insert}(Q,z) \end{array}$$

Analyse

Verwendung eines Heaps: Heap bauen in $\mathcal{O}(n)$. Extract-Min in $O(\log n)$ für n Elemente. Somit Laufzeit $O(n \log n)$.

Das gierige Verfahren ist korrekt

Theorem 19

Seien x,y zwei Symbole mit kleinsten Frequenzen in C und sei T'(C') der optimale Baum zum Alphabet $C'=C-\{x,y\}+\{z\}$ mit neuem Symbol z mit f(z)=f(x)+f(y). Dann ist der Baum T(C) der aus T'(C') entsteht, indem der Knoten z durch einen inneren Knoten mit Kindern x und y ersetzt wird, ein optimaler Codebaum zum Alphabet C.

24. Graphen

Notation, Repräsentation, Traversieren (DFS, BFS), Topologisches Sortieren , Reflexive transitive Hülle, Zusammenhangskomponenten [Ottman/Widmayer, Kap. 9.1 - 9.4,Cormen et al, Kap. 22]

Beweis

Es gilt

$$f(x) \cdot d_T(x) + f(y) \cdot d_T(y) = (f(x) + f(y)) \cdot (d_{T'}(z) + 1) = f(z) \cdot d_{T'}(x) + f(x) + f(y).$$

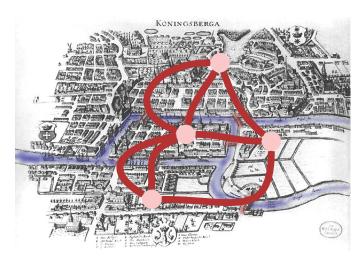
Also $B(T') = B(T) - f(x) - f(y).$

Annahme: T sei nicht optimal. Dann existiert ein optimaler Baum T'' mit B(T'') < B(T). Annahme: x und y Brüder in T''. T''' sei der Baum T'' in dem der innere Knoten mit Kindern x und y gegen z getauscht wird. Dann gilt B(T''') = B(T'') - f(x) - f(y) < B(T) - f(x) - f(y) = B(T'). Widerspruch zur Optimalität von T'.

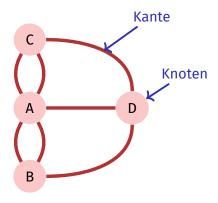
Die Annahme, dass x und y Brüder sind in T'' kann man rechtfertigen, da ein Tausch der Elemente mit kleinster Häufigkeit auf die unterste Ebene den Wert von B höchstens verkleinern kann.

4

Königsberg 1736



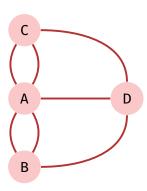
[Multi]Graph



Zyklen

- Gibt es einen Rundweg durch die Stadt (den Graphen), welcher jede Brücke (jede Kante) genau einmal benutzt?
- Euler (1736): nein.
- Solcher Rundweg (*Zyklus*) heisst Eulerscher Kreis.
- Eulerzyklus ⇔ jeder Knoten hat gerade Anzahl Kanten (jeder Knoten hat einen *geraden Grad*).

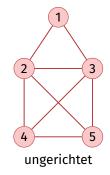
" \Rightarrow " ist sofort klar, " \Leftarrow " ist etwas schwieriger, aber auch elementar.



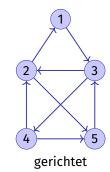
678

679

Notation



$$\begin{split} V = & \{1, 2, 3, 4, 5\} \\ E = & \{\{1, 2\}, \{1, 3\}, \{2, 3\}, \{2, 4\}, \\ & \{2, 5\}, \{3, 4\}, \{3, 5\}, \{4, 5\}\} \end{split}$$



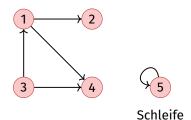
$$V = \{1, 2, 3, 4, 5\}$$

$$E = \{(1, 3), (2, 1), (2, 5), (3, 2),$$

$$(3, 4), (4, 2), (4, 5), (5, 3)\}$$

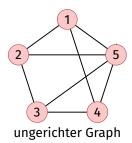
Notation

Ein gerichteter Graph besteht aus einer Menge $V = \{v_1, \dots, v_n\}$ von Knoten (*Vertices*) und einer Menge $E \subseteq V \times V$ von Kanten (*Edges*). Gleiche Kanten dürfen nicht mehrfach enthalten sein.



Notation

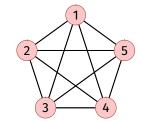
Ein ungerichteter Graph besteht aus einer Menge $V = \{v_1, \dots, v_n\}$ von Knoten und einer Menge $E \subseteq \{\{u, v\} | u, v \in V\}$ von Kanten. Kanten dürfen nicht mehrfach enthalten sein.³⁹



³⁹Im Gegensatz zum Eingangsbeispiel – dann Multigraph genannt.

Notation

Ein ungerichteter Graph G=(V,E) ohne Schleifen in dem jeder Knoten mit jedem anderen Knoten durch eine Kante verbunden ist, heisst vollständig.

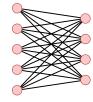


ein vollständiger ungerichter Graph

682

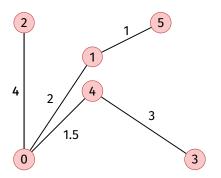
Notation

Ein Graph, bei dem V so in disjunkte U und W aufgeteilt werden kann, dass alle $e \in E$ einen Knoten in U und einen in W haben heisst bipartit.



Notation

Ein gewichteter Graph G=(V,E,c) ist ein Graph G=(V,E) mit einer Kantengewichtsfunktion $c:E\to\mathbb{R}.\ c(e)$ heisst Gewicht der Kante e.



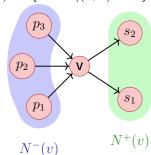
683

i e

Notation

Für gerichtete Graphen G = (V, E)

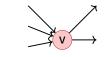
- $\mathbf{w} \in V$ heisst adjazent zu $v \in V$, falls $(v, w) \in E$
- Vorgängermenge von $v \in V$: $N^-(v) := \{u \in V | (u, v) \in E\}$. Nachfolgermenge: $N^+(v) := \{u \in V | (v, u) \in E\}$



Notation

Für gerichtete Graphen G = (V, E)

■ Eingangsgrad: $\deg^-(v) = |N^-(v)|$, Ausgangsgrad: $\deg^+(v) = |N^+(v)|$



 $\deg^-(v) = 3$, $\deg^+(v) = 2$



 $\deg^-(w) = 1$, $\deg^+(w) = 1$

Notation

Für ungerichtete Graphen G = (V, E):

- $lacksquare w \in V$ heisst adjazent zu $v \in V$, falls $\{v,w\} \in E$
- Nachbarschaft von $v \in V$: $N(v) = \{w \in V | \{v, w\} \in E\}$
- Grad von v: deg(v) = |N(v)| mit Spezialfall Schleifen: erhöhen Grad um 2.



deg(v) = 5



 $\deg(w) = 2$

Beziehung zwischen Knotengraden und Kantenzahl

In jedem Graphen G = (V, E) gilt

- 1. $\sum_{v \in V} \deg^-(v) = \sum_{v \in V} \deg^+(v) = |E|$, falls G gerichtet
- 2. $\sum_{v \in V} \deg(v) = 2|E|$, falls G ungerichtet.

00/

Wege

- Weg: Sequenz von Knoten $\langle v_1, \ldots, v_{k+1} \rangle$ so dass für jedes $i \in \{1 \ldots k\}$ eine Kante von v_i nach v_{i+1} existiert.
- lacktriangle Länge des Weges: Anzahl enthaltene Kanten k.
- Gewicht des Weges (in gewichteten Graphen): $\sum_{i=1}^k c((v_i, v_{i+1}))$ (bzw. $\sum_{i=1}^k c(\{v_i, v_{i+1}\})$)
- Pfad (auch: einfacher Pfad): Weg der keinen Knoten mehrfach verwendet.

Zusammenhang

- Ungerichteter Graph heisst zusammenhängend, wenn für jedes Paar $v, w \in V$ ein verbindender Weg existiert.
- Gerichteter Graph heisst stark zusammenhängend, wenn für jedes Paar $v, w \in V$ ein verbindender Weg existiert.
- Gerichteter Graph heisst schwach zusammenhängend, wenn der entsprechende ungerichtete Graph zusammenhängend ist.

0

Einfache Beobachtungen

- Allgemein: $0 \le |E| \in \mathcal{O}(|V|^2)$
- lacksquare Zusammenhängender Graph: $|E| \in \Omega(|V|)$
- \blacksquare Vollständiger Graph: $|E|=\frac{|V|\cdot(|V|-1)}{2}$ (ungerichtet)
- lacksquare Maximal $|E|=|V|^2$ (gerichtet), $|E|=rac{|V|\cdot(|V|+1)}{2}$ (ungerichtet)

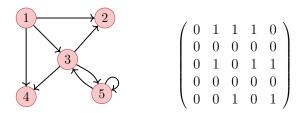
Zyklen

- **Zyklus:** Weg $\langle v_1, \dots, v_{k+1} \rangle$ mit $v_1 = v_{k+1}$
- Kreis: Zyklus mit paarweise verschiedenen v_1, \ldots, v_k , welcher keine Kante mehrfach verwendet.
- Kreisfrei (azyklisch): Graph ohne jegliche Kreise.

Eine Folgerung: Ungerichtete Graphen können keinen Kreis der Länge 2 enthalten (Schleifen haben Länge 1).

Repräsentation mit Matrix

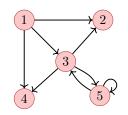
Graph G=(V,E) mit Knotenmenge v_1,\ldots,v_n gespeichert als Adjazenzmatrix $A_G=(a_{ij})_{1\leq i,j\leq n}$ mit Einträgen aus $\{0,1\}$. $a_{ij}=1$ genau dann wenn Kante von v_i nach v_j .

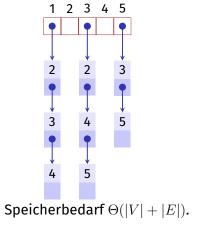


Speicherbedarf $\Theta(|V|^2)$. A_G ist symmetrisch, wenn G ungerichtet.

Repräsentation mit Liste

Viele Graphen G=(V,E) mit Knotenmenge v_1,\ldots,v_n haben deutlich weniger als n^2 Kanten. Repräsentation mit Adjazenzliste: Array $A[1],\ldots,A[n]$, A_i enthält verkettete Liste aller Knoten in $N^+(v_i)$.



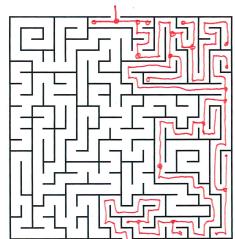


94

Laufzeiten einfacher Operationen

Operation	Matrix	Liste	
${\bf Nachbarn/Nachfolger\ von\ }v\in V\ {\bf finden}$	$\Theta(n)$	$\Theta(\deg^+ v)$	
$v \in V$ ohne Nachbar/Nachfolger finden	$\Theta(n^2)$	$\Theta(n)$	
$(v,u) \in E$?	$\Theta(1)$	$\Theta(\deg^+ v)$	
Kante einfügen	$\Theta(1)$	$\Theta(1)$	
Kante (v,u) löschen	$\Theta(1)$	$\Theta(\deg^+ v)$	

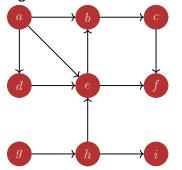
Tiefensuche



596

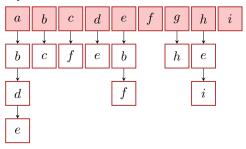
Graphen Traversieren: Tiefensuche

Verfolge zuerst Pfad in die Tiefe, bis nichts mehr besucht werden kann.



Reihenfolge a, b, c, f, d, e, g, h, i

Adjazenzliste



Farben

Konzeptuelle Färbung der Knoten

- Weiss: Knoten wurde noch nicht entdeckt.
- Grau: Knoten wurde entdeckt und zur Traversierung vorgemerkt / in Bearbeitung.
- Schwarz: Knoten wurde entdeckt und vollständig bearbeitet

Algorithmus Tiefensuche DFS-Visit(G, v)

```
Input: Graph G = (V, E), Knoten v.
v.color \leftarrow \mathsf{grey}
foreach w \in N^+(v) do
    if w.color = white then
          \mathsf{DFS}\text{-Visit}(G,w)
v.color \leftarrow \mathsf{black}
```

Tiefensuche ab Knoten v. Laufzeit (ohne Rekursion): $\Theta(\deg^+ v)$

Algorithmus Tiefensuche DFS-Visit(G)

```
Input: Graph G = (V, E)
foreach v \in V do
   v.color \leftarrow \mathsf{white}
foreach v \in V do
    if v.color = white then
        DFS-Visit(G,v)
```

Tiefensuche für alle Knoten eines Graphen. Laufzeit $\Theta(|V| + \sum_{v \in V} (\deg^+(v) + 1)) = \Theta(|V| + |E|).$

Iteratives DFS-Visit(G, v)

```
Input: Graph G = (V, E), v \in V mit v.color = white
Stack S \leftarrow \emptyset
v.color \leftarrow \mathsf{grey}; S.\mathsf{push}(v)
                                                      // invariant: grey nodes always on stack
while S \neq \emptyset do
    w \leftarrow \mathsf{nextWhiteSuccessor}(v)
                                                                                 // code: next slide
    if w \neq \text{null then}
         w.color \leftarrow \mathsf{grey}; S.\mathsf{push}(w)
                                                  // work on w. parent remains on the stack
    else
                                                       // no grey successors, v becomes black
          v.color \leftarrow \mathsf{black}
         if S \neq \emptyset then
                                                                        // visit/revisit next node
               v \leftarrow S.\mathsf{pop}()
              if v.color = grey then S.push(v)
                                                                    Speicherbedarf Stack \Theta(|V|)
```

nextWhiteSuccessor(v)

```
\textbf{Input:} \ \mathsf{Knoten} \ v \in V
```

Output: Nachfolgeknoten u von v mit u.color = white, null sonst

return null

Interpretation der Farben

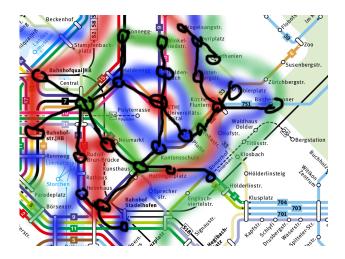
Beim Traversieren des Graphen wird ein Baum (oder Wald) aufgebaut. Beim Entdecken von Knoten gibt es drei Fälle

■ Weisser Knoten: neue Baumkante

■ Grauer Knoten: Zyklus ("Rückwärtskante")

■ Schwarzer Knoten: Vorwärts-/Seitwärtskante

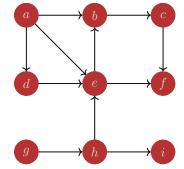
Breitensuche



/(

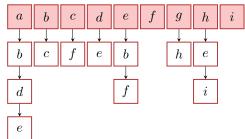
Graphen Traversieren: Breitensuche

Verfolge zuerst Pfad in die Breite, gehe dann in die Tiefe.





Adjazenzliste



Reihenfolge a, b, d, e, c, f, g, h, i

(Iteratives) BFS-Visit(G, v)

```
Input: Graph G = (V, E)
Queue Q \leftarrow \emptyset
v.color \leftarrow \mathsf{grey}
enqueue(Q, v)
while Q \neq \emptyset do
     w \leftarrow \mathsf{dequeue}(Q)
     foreach c \in N^+(w) do
          if c.color = white then
                c.color \leftarrow \mathsf{grey}
                enqueue(Q, c)
     w.color \leftarrow \mathsf{black}
```

Algorithmus kommt mit $\mathcal{O}(|V|)$ Extraplatz aus.

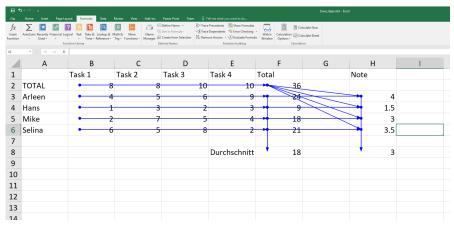
706

Rahmenprogramm BFS-Visit(G)

```
Input: Graph G = (V, E)
foreach v \in V do
   v.color \leftarrow \mathsf{white}
foreach v \in V do
    if v.color = white then
        BFS-Visit(G,v)
```

Breitensuche für alle Knoten eines Graphen. Laufzeit $\Theta(|V| + |E|)$.

Topologisches Sortieren



Auswertungsreihenfolge?

Topologische Sortierung

Topologische Sortierung eines azyklischen gerichteten Graphen G=(V,E):

Bijektive Abbildung

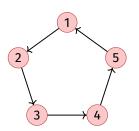
ord :
$$V \to \{1, ..., |V|\}$$

so dass

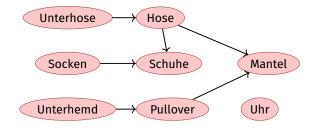
$$\operatorname{ord}(v) < \operatorname{ord}(w) \ \forall \ (v, w) \in E.$$

Identifizieren Wert i mit dem Element $v_i := \operatorname{ord}^{-1}(i)$. Topologische Sortierung $\triangleq \langle v_1, \dots, v_{|V|} \rangle$.

(Gegen-)Beispiele



Zyklischer Graph: kann nicht topol- Graphen: ogisch sortiert werden. Unterhem



Eine mögliche topologische Sortierung des Graphen: Unterhemd, Pullover, Unterhose, Uhr, Hose, Mantel,

711

710

Socken, Schuhe

Beobachtung

Theorem 20

Ein gerichteter Graph G=(V,E) besitzt genau dann eine topologische Sortierung, wenn er kreisfrei ist

Beweis "⇒"

Wenn G einen Kreis besitzt, so besitzt er keine topologische Sortierung. Denn in einem Kreis $\langle v_{i_1}, \dots, v_{i_m} \rangle$ gälte $v_{i_1} < \dots < v_{i_m} < v_{i_1}$.

Beweis "←"

- Anfang (n = 1): Graph mit einem Knoten ohne Schleife ist topologisch sortierbar. Setze $ord(v_1) = 1$.
- \blacksquare Hypothese: Graph mit n Knoten kann topologisch sortiert werden.
- \blacksquare Schritt ($n \rightarrow n+1$):
 - 1. G enthält einen Knoten v_q mit Eingangsgrad $\deg^-(v_q)=0$. Andernfalls verfolge iterativ Kanten rückwärts nach spätestens n+1 Iterationen würde man einen Knoten besuchen, welcher bereits besucht wurde. Widerspruch zur Zyklenfreiheit.
 - 2. Graph ohne Knoten v_q und ohne dessen Eingangskanten kann nach Hypothese topologisch sortiert werden. Verwende diese Sortierung, setze $\operatorname{ord}(v_i) \leftarrow \operatorname{ord}(v_i) + 1$ für alle $i \neq q$ und setze $\operatorname{ord}(v_q) \leftarrow 1$.

Algorithmus Topological-Sort(G)

```
Input: Graph G = (V, E).
```

Output: Topologische Sortierung ord

```
Stack S \leftarrow \emptyset
```

```
foreach v \in V do A[v] \leftarrow 0
```

foreach $(v, w) \in E$ **do** $A[w] \leftarrow A[w] + 1 //$ Eingangsgrade berechnen

foreach $v \in V$ with A[v] = 0 **do** push(S, v) // Merke Nodes mit Eingangsgrad 0 $i \leftarrow$ 1

while $S \neq \emptyset$ do

```
v \leftarrow \mathsf{pop}(S); \ \mathrm{ord}[v] \leftarrow i; \ i \leftarrow i+1 \ // \ \mathsf{W\"{a}hle} \ \mathsf{Knoten} \ \mathsf{mit} \ \mathsf{Eingangsgrad} \ \mathsf{0} foreach (v,w) \in E \ \mathsf{do} \ // \ \mathsf{Verringere} \ \mathsf{Eingangsgrad} \ \mathsf{der} \ \mathsf{Nachfolger}  \ | \ A[w] \leftarrow A[w] - 1  if A[w] = 0 \ \mathsf{then} \ \mathsf{push}(S,w)
```

if i = |V| + 1 then return ord else return "Cycle Detected"

714

Algorithmus Korrektheit

Theorem 21

Sei G=(V,E) ein gerichteter, kreisfreier Graph. Der Algorithmus **TopologicalSort**(G) berechnet in Zeit $\Theta(|V|+|E|)$ eine topologische Sortierung ord für G.

Beweis: folgt im wesentlichen aus vorigem Theorem:

- 1. Eingangsgrad verringern entspricht Knotenentfernen.
- 2. Im Algorithmus gilt für jeden Knoten v mit A[v] = 0 dass entweder der Knoten Eingangsgrad 0 hat oder dass zuvor alle Vorgänger einen Wert $\operatorname{ord}[u] \leftarrow i$ zugewiesen bekamen und somit $\operatorname{ord}[v] > \operatorname{ord}[u]$ für alle Vorgänger u von v. Knoten werden nur einmal auf den Stack gelegt.
- 3. Laufzeit: Inspektion des Algorithmus (mit Argumenten wie beim Traversieren).

Algorithmus Korrektheit

Theorem 22

Sei G=(V,E) ein gerichteter, nicht kreisfreier Graph. Der Algorithmus $\mathtt{TopologicalSort}(G)$ terminiert in Zeit $\Theta(|V|+|E|)$ und detektiert Zyklus.

Beweis: Sei $\langle v_{i_1},\dots,v_{i_k} \rangle$ ein Kreis in G. In jedem Schritt des Algorithmus bleibt $A[v_{i_j}] \geq 1$ für alle $j=1,\dots,k$. Also werden k Knoten nie auf den Stack gelegt und somit ist zum Schluss $i \leq V+1-k$.

Die Laufzeit des zweiten Teils des Algorithmus kann kürzer werden, jedoch kostet die Berechnung der Eingangsgrade bereits $\Theta(|V| + |E|)$.

715

, , ,

Alternative: Algorithmus DFS-Topsort(G, v)

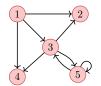
Rufe Algorithmus für jeden noch nicht besuchten Knoten auf. Asymptotische Laufzeit $\Theta(|V|+|E|)$.

Interpretation

Theorem 23

Sei G=(V,E) ein Graph und $k\in\mathbb{N}$. Dann gibt das Element $a_{i,j}^{(k)}$ der Matrix $(a_{i,j}^{(k)})_{1\leq i,j\leq n}=(A_G)^k$ die Anzahl der Wege mit Länge k von v_i nach v_i an.

Adjazenzmatrizen multipliziert



$$B := A_G^2 = \begin{pmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 \end{pmatrix}^2 = \begin{pmatrix} 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 2 \end{pmatrix}$$

718

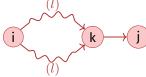
Beweis

720

Per Induktion.

Anfang: Klar für k=1. $a_{i,j}=a_{i,j}^{(1)}$. **Hypothese:** Aussage wahr für alle $k \leq l$

Schritt (
$$l o l+1$$
):
$$a_{i,j}^{(l+1)} = \sum^n a_{i,k}^{(l)} \cdot a_{k,j}$$



$$a_{k,j}=1$$
 g.d.w. Kante von k nach j , 0 sonst. Summe zählt die Anzahl Wege

721

der Länge l vom Knoten v_i zu allen Knoten v_k welche direkte Verbindung zu Knoten v_j haben, also alle Wege der Länge l+1.

Relation

Gegeben: endliche Menge $\,V\,$

(Binäre) **Relation** R auf V: Teilmenge des kartesischen Produkts

 $V \times V = \{(a,b)|a \in V, b \in V\}$ Relation $R \subseteq V \times V$ heisst

lacksquare reflexiv, wenn $(v,v)\in R$ für alle $v\in V$

 \blacksquare symmetrisch, wenn $(v,w) \in R \Rightarrow (w,v) \in R$

Transitiv, wenn $(v,x) \in R$, $(x,w) \in R \Rightarrow (v,w) \in R$

Die (Reflexive) Transitive Hülle R^* von R ist die kleinste Erweiterung $R \subseteq R^* \subseteq V \times V$ von R_* so dass R^* reflexiv und transitiv ist.

Graphen und Relationen

Graph G = (V, E)

Adjazenzen $A_G \cong \text{Relation } E \subseteq V \times V \text{ auf } V$

- reflexiv $\Leftrightarrow a_{i,i} = 1$ für alle i = 1, ..., n. (Schleifen)
- **symmetrisch** $\Leftrightarrow a_{i,j} = a_{j,i}$ für alle $i, j = 1, \dots, n$ (ungerichtet)
- transitiv \Leftrightarrow $(u,v) \in E$, $(v,w) \in E \Rightarrow (u,w) \in E$. (Erreichbarkeit)

722

724

Reflexive Transitive Hülle

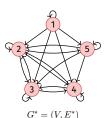
Reflexive transitive Hülle von $G \Leftrightarrow \text{Erreichbarkeits relation } E^*: (v, w) \in E^*$ gdw. \exists Weg von Knoten v zu w.

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ \end{bmatrix}$$



G=(V,E)





Algorithmus $A \cdot A$

Input: (Adjazenz-)Matrix $A=(a_{ij})_{i,j=1\dots n}$

Output: Matrixprodukt $B = (b_{ij})_{i,j=1...n} = A \cdot A$

$$B \leftarrow 0$$

for $r \leftarrow 1$ to n do

for
$$c \leftarrow 1$$
 to n do for $k \leftarrow 1$ to n do
$$b_{rc} \leftarrow b_{rc} + a_{rk} \cdot a_{kc}$$

// Anzahl Pfade

return B

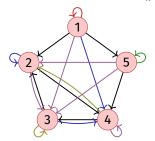
Berechnet Anzahl Pfade der Länge 2

Algorithmus $A \otimes A$

Berechnet Existenz von Pfaden der Längen $1\ \mathrm{und}\ 2$

Verbesserung: Algorithmus von Warshall (1962)

Induktiver Ansatz: Alle Wege bekannt über Knoten aus $\{v_i : i < k\}$. Hinzunahme des Knotens v_k .



$$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \end{bmatrix}$$

Berechnung Reflexive Transitive Hülle

Ziel: Berechnung von $B = (b_{ij})_{1 \le i,j \le n}$ mit $b_{ij} = 1 \Leftrightarrow (v_i, v_j) \in E^*$ Erste Idee:

- Starte mit $B \leftarrow A$ und setze $b_{ii} = 1$ für alle i (Reflexivität).
- Berechne

$$B_n = \bigotimes_{i=1}^n B$$

727

mit Potenzen von 2 $B_2:=B\otimes B$, $B_4:=B_2\otimes B_2$, $B_8=B_4\otimes B_4$... \Rightarrow Laufzeit $n^3\lceil\log_2 n\rceil$

726

Algorithmus TransitiveClosure(A_G)

```
Input: Adjazenzmatrix A_G = (a_{ij})_{i,j=1...n}

Output: Reflexive Transitive Hülle B = (b_{ij})_{i,j=1...n} von G
```

$$\begin{array}{c|c} B \leftarrow A_G \\ \textbf{for } k \leftarrow 1 \textbf{ to } n \textbf{ do} \\ & b_{kk} \leftarrow 1 \\ \textbf{for } r \leftarrow 1 \textbf{ to } n \textbf{ do} \\ & \textbf{for } c \leftarrow 1 \textbf{ to } n \textbf{ do} \\ & & b_{rc} \leftarrow \max\{b_{rc}, b_{rk} \cdot b_{kc}\} \end{array} \qquad // \textbf{ Alle Wege "uber } v_k$$

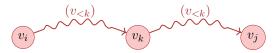
return B

Laufzeit des Algorithmus $\Theta(n^3)$.

Korrektheit des Algorithmus (Induktion)

Invariante (k**):** alle Wege über Knoten mit maximalem Index < k berücksichtigt

- Anfang (k = 1): Alle direkten Wege (alle Kanten) in A_G berücksichtigt.
- **Hypothese**: Invariante (k) erfüllt.
- Schritt $(k \to k+1)$: Für jeden Weg von v_i nach v_j über Knoten mit maximalen Index k: nach Hypothese $b_{ik}=1$ und $b_{kj}=1$. Somit im k-ten Schleifendurchlauf: $b_{ij} \leftarrow 1$.



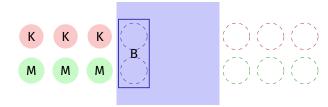
25. Kürzeste Wege

Motivation, Universeller Algorithmus, Dijkstras Algorithmus auf Distanzgraphen, Algorithmus von Bellman-Ford, Algorithmus von Floyd-Warshall, Johnson Algorithmus [Ottman/Widmayer, Kap. 9.5 Cormen et al, Kap. 24.1-24.3, 25.2-25.3]

0

Flussüberquerung (Missionare und Kannibalen)

Problem: Drei Kannibalen und drei Missionare stehen an einem Ufer eines Flusses. Ein dort bereitstehendes Boot fasst maximal zwei Personen. Zu keiner Zeit dürfen an einem Ort (Ufer oder Boot) mehr Kannibalen als Missionare sein. Wie kommen die Missionare und Kannibalen möglichst schnell über den Fluss? 40



⁴⁰Es gibt leichte Variationen dieses Problems, es ist auch äquivalent zum Problem der eifersüchtigen Ehemänner

Formulierung als Graph

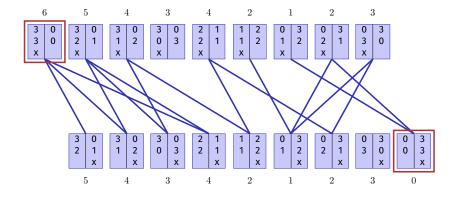
Zähle alle erlaubten Konfigurationen als Knoten auf und verbinde diese mit einer Kante, wenn Überfahrt möglich ist. Das Problem ist dann ein Problem des kürzesten Pfades Beispiel

	links	rechts			links	rechts
Missionare	3	0	Überfahrt möglich	Missionare	2	1
Kannibalen	3	0		Kannibalen	2	1
Boot	Х			Boot		Х

6 Personen am linken Ufer

4 Personen am linken Ufer

Das ganze Problem als Graph



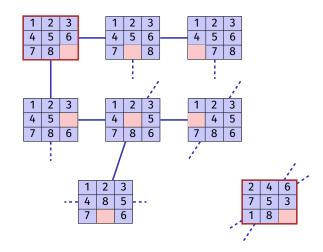
Anderes Beispiel: Schiebepuzzle

Wollen die schnelleste Lösung finden für

	2	4	6	1	2	3
Ì	7	5	3	 4	5	6
Ì	1	8		7	8	

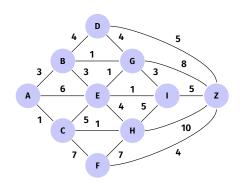
734

Problem als Graph



Routenfinder

Gegeben Städte A - Z und Distanzen zwischen den Städten.

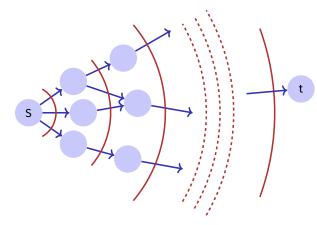


Was ist der kürzeste Weg von A nach Z?

Einfachster Fall

Konstantes Kantengewicht 1 (oBdA)

Lösung: Breitensuche



Kürzeste Wege

Notation: Wir schreiben

$$u \stackrel{p}{\leadsto} v$$
 oder $p: u \leadsto v$

und meinen einen Weg p von u nach v

Notation: $\delta(u, v)$ = Gewicht eines kürzesten Weges von u nach v:

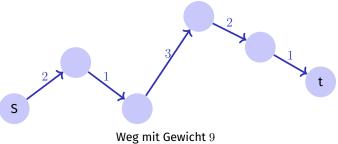
$$\delta(u,v) = \begin{cases} \infty & \text{kein Weg von } u \text{ nach } v \\ \min\{c(p) : u \stackrel{p}{\leadsto} v\} & \text{sonst} \end{cases}$$

Gewichtete Graphen

Gegeben: G=(V,E,c), $c:E\to\mathbb{R}$, $s,t\in V$.

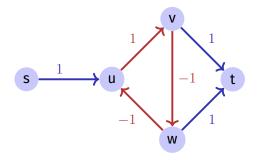
Gesucht: Länge (Gewicht) eines kürzesten Weges von s nach t.

Weg: $p = \langle s = v_0, v_1, \dots, v_k = t \rangle$, $(v_i, v_{i+1}) \in E$ ($0 \le i < k$) Gewicht: $c(p) := \sum_{i=0}^{k-1} c((v_i, v_{i+1}))$.



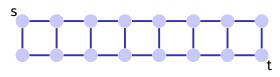
Beobachtungen (1)

Es gibt Situationen, in denen kein kürzester Weg existiert: negative Zyklen könnten auftreten.



Beobachtungen (2)

Es kann exponentiell viele Wege geben.



(mindestens $2^{|V|/2}$ Wege von s nach t)

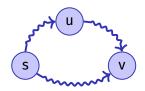
 \Rightarrow Alle Wege probieren ist zu ineffizient.

Beobachtungen (3)

Dreiecksungleichung

Für alle $s, u, v \in V$:

 $\delta(s, v) \le \delta(s, u) + \delta(u, v)$



Ein kürzester Weg von s nach v (ohne weitere Einschränkungen) kann nicht länger sein als ein kürzester Weg von s nach v, der u enthalten muss.

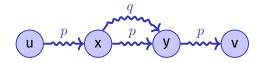
742

743

Beobachtungen (4)

Optimale Substruktur

Teilpfade von kürzesten Pfaden sind kürzeste Pfade: Sei $p=\langle v_0,\ldots,v_k\rangle$ ein kürzester Pfad von v_0 nach v_k . Dann ist jeder der Teilpfade $p_{ij}=\langle v_i,\ldots,v_j\rangle$ ($0\leq i< j\leq k$) ein kürzester Pfad von v_i nach v_j .



Wäre das nicht so, könnte man einen der Teilpfade kürzen, Widerspruch zur Voraussetzung.

Beobachtungen (5)

Kürzeste Wege enthalten keine Zyklen

- 1. Kürzester Weg enthält negativen Zyklus: es existiert kein kürzester Weg. Widerspruch.
- 2. Weg enthält positiven Zyklus: Weglassen des positiven Zyklus kann den Weg verkürzen: Widerspruch
- 3. Weg enthält Zyklus vom Gewicht 0: Weglassen des Zyklus verändert das Pfadgewicht nicht. Weglassen (Konvention).

Zutaten für einen Algorithmus

Gesucht: Kürzeste Wege von einem Startknoten s aus.

■ Gewicht des kürzesten bisher gefundenen Pfades

$$d_s:V\to\mathbb{R}$$

Zu Beginn: $d_s[v] = \infty$ für alle Knoten $v \in V$. Ziel: $d_s[v] = \delta(s,v)$ für alle $v \in V$.

■ Vorgänger eines Knotens

$$\pi_s: V \to V$$

Zu Beginn $\pi_s[v]$ undefiniert für jeden Knoten $v \in V$

Relaxieren ist sicher

Zu jeder Zeit gilt in obigem Algorithmus

$$d_s[v] \ge \delta(s, v) \quad \forall v \in V$$

Im Relaxierschritt:

$$\begin{split} \delta(s,v) & \leq \delta(s,u) + \delta(u,v) \\ \delta(s,u) & \leq d_s[u] \\ \delta(u,v) & \leq c(u,v) \\ \Rightarrow & d_s[u] + c(u,v) \geq \delta(s,v) \end{split} \qquad \begin{aligned} & \text{[Induktionsvorraussetzung]}. \\ & \text{[Minimalität von δ]} \end{aligned}$$

$$\Rightarrow \min\{d_s[v], d_s[u] + c(u, v)\} \ge \delta(s, v)$$

Allgemeiner Algorithmus

- 1. Initialisiere d_s und π_s : $d_s[v] = \infty$, $\pi_s[v] = \text{null für alle } v \in V$
- 2. Setze $d_s[s] \leftarrow 0$
- 3. Wähle eine Kante $(u, v) \in E$

Relaxiere
$$(u,v)$$
: if $d_s[v] > d_s[u] + c(u,v)$ then
$$d_s[v] \leftarrow d_s[u] + c(u,v)$$

$$\pi_s[v] \leftarrow u$$

4. Wiederhole 3 bis nichts mehr relaxiert werden kann.

(bis
$$d_s[v] \le d_s[u] + c(u,v) \quad \forall (u,v) \in E$$
)

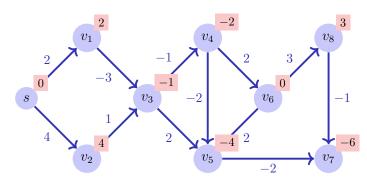
Zentrale Frage

746

Wie / in welcher Reihenfolge wählt man die Kanten in obigem Algorithmus?

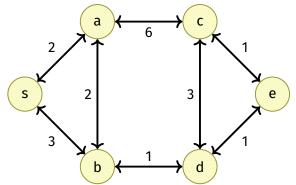
Spezialfall: Gerichteter Azyklischer Graph (DAG)

DAG ⇒ Topologische Sortierung liefert optimale Besuchsreihenfolge



Top. Sortieren: \Rightarrow Reihenfolge $s, v_1, v_2, v_3, v_4, v_6, v_5, v_8, v_7$.

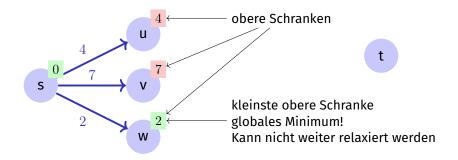
Annahme (vorübergehend)



Alle Gewichte von G sind positiv.

50

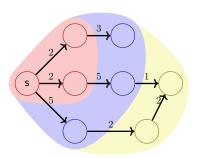
Beobachtung (Dijkstra)



Grundidee

Menge V aller Knoten wird unterteilt in

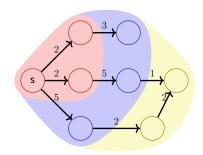
- die Menge *M* von Knoten, für die schon ein kürzester Weg von *s* bekannt ist
- die Menge $R = \bigcup_{v \in M} N^+(v) \setminus M$ von Knoten, für die kein kürzester Weg bekannt ist, die jedoch von M direkt erreichbar sind.
- die Menge $U = V \setminus (M \cup R)$ von Knoten, die noch nicht berücksichtigt wurden.



Induktion

Induktion über |M|: Wähle Knoten aus R mit kleinster oberer Schranke. Nimm r zu M hinzu, und update R und U.

Korrektheit: Ist innerhalb einer "Wellenfront" einmal ein Knoten mit minimalem Pfadgewicht w gefunden, kann kein Pfad über später gefundene Knoten (mit Gewicht $\geq w$) zu einer Verbesserung führen.



754

Algorithmus Dijkstra(G, s)

Input: Positiv gewichteter Graph G = (V, E, c), Startpunkt $s \in V$

Output: Minimale Gewichte d der kürzesten Pfade und Vorgängerknoten für jeden Knoten.

foreach
$$u \in V$$
 do

$$d_s[s] \leftarrow 0; R \leftarrow \{s\}$$

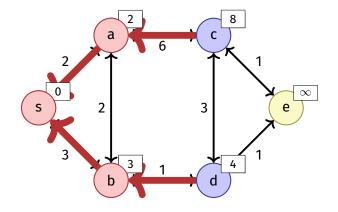
while
$$R \neq \emptyset$$
 do

$$u \leftarrow \mathsf{ExtractMin}(R)$$
 for each $v \in N^+(u)$ do
$$| \text{ if } d_s[u] + c(u,v) < d_s[v] \text{ then }$$

$$| d_s[v] \leftarrow d_s[u] + c(u,v)$$

$$| \pi[v] \leftarrow u$$

Beispiel



$$M = \{s, a, b\}$$
$$R = \{c, d\}$$
$$U = \{e\}$$

Zur Implementation: Datenstruktur für R?

Benötigte Operationen:

- Insert (Hinzunehmen zu R)
- ExtractMin (über R) und DecreaseKey (Update in R)

MinHeap!

DecreaseKey

- DecreaseKey: Aufsteigen im MinHeap in $\mathcal{O}(\log |V|)$
- Position im Heap?
 - Möglichkeit (a): Speichern am Knoten
 - Möglichkeit (b): Hashtabelle über Knoten
 - Möglichkeit (c): Knoten nach erfolgreichem Relaxieren erneut einfügen. Knoten beim Entnehmen als "deleted" kennzeichnen (Lazy Deletion).⁴¹

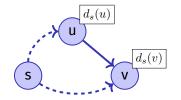
758

Allgemeine Bewertete Graphen

Verbesserungsschritt wie bisher, aber mit Rückgabewert:

$$\begin{aligned} & \mathsf{Relax}(u,v) \ \textbf{(}u,v \in V \textbf{,} \ (u,v) \in E \textbf{)} \\ & \mathsf{if} \ d_s[u] + c(u,v) < d_s[v] \ \mathsf{then} \\ & d_s[v] \leftarrow d_s[u] + c(u,v) \\ & \pi_s[v] \leftarrow u \\ & \mathsf{return} \ \mathsf{true} \end{aligned}$$

return false



Problem: Zyklen mit negativen Gewichten können Weg verkürzen: es muss keinen kürzesten Weg mehr geben

Laufzeit

- $|V| \times \text{ExtractMin: } \mathcal{O}(|V| \log |V|)$
- \blacksquare $|E| \times$ Insert oder DecreaseKey: $\mathcal{O}(|E| \log |V|)$
- $1 \times Init: \mathcal{O}(|V|)$
- Insgesamt⁴²: $\mathcal{O}((|V| + |E|) \log |V|)$

Kann verbessert werden unter Verwendung einer für ExtractMin und DecreaseKey optimierten Datenstruktur (Fibonacci Heap), dann Laufzeit $\mathcal{O}(|E|+|V|\log|V|)$.

Dynamic Programming Ansatz (Bellman)

Induktion über Anzahl Kanten. $d_s[i,v]$: Kürzeste Weglänge von s nach v über maximal i Kanten.

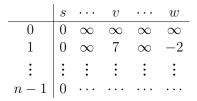
$$\begin{split} d_s[i,v] &= \min\{d_s[i-1,v], \min_{(u,v) \in E} (d_s[i-1,u] + c(u,v)) \\ d_s[0,s] &= 0, d_s[0,v] = \infty \ \forall v \neq s. \end{split}$$

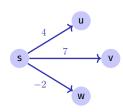
759

⁴¹Für die lazy deletion benötigt man ein Paar von Kante (oder Zielknoten) und Distanz

 $^{^{42}\}mathcal{O}(|E|\log|V|)$ für zusammenhängende Graphen

Dynamic Programming Ansatz (Bellman)





Algorithmus: Iteriere über letzte Zeile bis die Relaxationsschritte keine Änderung mehr ergeben, maximal aber n-1 mal. Wenn dann noch Änderungen, dann gibt es keinen kürzesten Pfad.

Alle kürzesten Pfade

Ziel: Berechne das Gewicht eines kürzesten Pfades für jedes Knotenpaar.

- $|V| \times \text{Anwendung von Dijkstras ShortestPath: } \mathcal{O}(|V| \cdot (|E| + |V|) \cdot \log |V|)$ (Mit Fibonacci-Heap:: $\mathcal{O}(|V|^2 \log |V| + |V| \cdot |E|)$)
- $|V| \times \text{Anwendung von Bellman-Ford: } \mathcal{O}(|E| \cdot |V|^2)$
- Es geht besser!

Algorithmus Bellman-Ford(G, s)

Input: Graph G = (V, E, c), Startpunkt $s \in V$

Output: Wenn Rückgabe true, Minimale Gewichte d der kürzesten Pfade zu jedem Knoten, sonst kein kürzester Pfad.

```
foreach u \in V do
   d_s[u] \leftarrow \infty; \ \pi_s[u] \leftarrow \mathsf{null}
d_s[s] \leftarrow 0;
for i \leftarrow 1 to |V| do
     f \leftarrow \mathsf{false}
     foreach (u, v) \in E do
       f \leftarrow f \vee \text{Relax}(u, v)
     if f = false then return true
```

return false:

Induktion über Knotennummer.

Betrachte die Gewichte aller kürzesten Wege S^k mit Zwischenknoten in 43 $V^k := \{v_1, \dots, v_k\}$, wenn Gewichte zu allen kürzesten Wegen S^{k-1} mit Zwischenknoten in V^{k-1} gegeben sind.

- \mathbf{v}_k kein Zwischenknoten eines kürzesten Pfades von $v_i \rightsquigarrow v_i$ in V^k : Gewicht eines kürzesten Pfades $v_i \leadsto v_j$ in S^{k-1} dann auch das Gewicht eines kürzesten Pfades in S^k .
- \mathbf{v}_k Zwischenknoten eines kürzesten Pfades $v_i \rightsquigarrow v_i$ in V^k : Teilpfade $v_i \leadsto v_k$ und $v_k \leadsto v_i$ enthalten nur Zwischenknoten aus S^{k-1} .

⁴³wie beim Algorithmus für die reflexive transitive Hülle von Warshall

DP Induktion

 $d^k(u,v)$ = Minimales Gewicht eines Pfades $u \leadsto v$ mit Zwischenknoten aus V^k

Induktion

$$d^{k}(u,v) = \min\{d^{k-1}(u,v), d^{k-1}(u,k) + d^{k-1}(k,v)\}(k \ge 1)$$

$$d^{0}(u,v) = c(u,v)$$

DP Algorithmus Floyd-Warshall(G)

Input: Graph G = (V, E, c) ohne Tyklen mit negativem Gewicht.

Output: Minimale Gewichte aller Pfade d

$$d^0 \leftarrow c$$

$$\begin{array}{c|c} \text{for } k \leftarrow 1 \text{ to } |V| \text{ do} \\ & \text{for } i \leftarrow 1 \text{ to } |V| \text{ do} \\ & \text{for } j \leftarrow 1 \text{ to } |V| \text{ do} \\ & & L d^k(v_i,v_j) = \min\{d^{k-1}(v_i,v_j), d^{k-1}(v_i,v_k) + d^{k-1}(v_k,v_j)\} \end{array}$$

Laufzeit: $\Theta(|V|^3)$

Bemerkung: Der Algorithmus kann auf einer einzigen Matrix d (in place) ausgeführt werden.

766

Umgewichtung

Idee: Anwendung von Dijkstras Algorithmus auf Graphen mit negativen Gewichten durch Umgewichtung

Das folgende geht nicht. Die Graphen sind nicht äquivalent im Sinne der kürzesten Pfade.



Umgewichtung

Andere Idee: "Potentialfunktion" (Höhe) auf den Knoten

- \blacksquare G = (V, E, c) ein gewichteter Graph.
- Funktion $h: V \to \mathbb{R}$
- Neue Gewichte

$$\tilde{c}(u,v) = c(u,v) + h(u) - h(v), (u,v \in V)$$

Umgewichtung

Beobachtung: Ein Pfad p ist genau dann kürzester Pfad in G=(V,E,c), wenn er in $\tilde{G}=(V,E,\tilde{c})$ kürzester Pfad ist.

$$\tilde{c}(p) = \sum_{i=1}^{k} \tilde{c}(v_{i-1}, v_i) = \sum_{i=1}^{k} c(v_{i-1}, v_i) + h(v_{i-1}) - h(v_i)$$

$$= h(v_0) - h(v_k) + \sum_{i=1}^{k} c(v_{i-1}, v_i) = c(p) + h(v_0) - h(v_k)$$

Also $\tilde{c}(p)$ minimal unter allen $v_0 \leadsto v_k \Longleftrightarrow c(p)$ minimal unter allen $v_0 \leadsto v_k$. Zyklengewichte sind invariant: $\tilde{c}(v_0,\ldots,v_k=v_0)=c(v_0,\ldots,v_k=v_0)$

Johnsons Algorithmus

Hinzunahme eines neuen Knotens $s \notin V$:

$$G' = (V', E', c')$$

$$V' = V \cup \{s\}$$

$$E' = E \cup \{(s, v) : v \in V\}$$

$$c'(u, v) = c(u, v), \ u \neq s$$

$$c'(s, v) = 0(v \in V)$$

Johnsons Algorithmus

Falls keine negativen Zyklen: wähle für Höhenfunktion Gewicht der kürzesten Pfade von s,

$$h(v) = d(s, v).$$

Für minimales Gewicht d eines Pfades gilt generell folgende Dreiecksungleichung:

$$d(s, v) < d(s, u) + c(u, v).$$

Einsetzen ergibt $h(v) \leq h(u) + c(u, v)$. Damit

$$\tilde{c}(u,v) = c(u,v) + h(u) - h(v) \ge 0.$$

Algorithmus Johnson(G)

Input: Gewichteter Graph G = (V, E, c)

Output: Minimale Gewichte aller Pfade D.

Neuer Knoten s. Berechne G' = (V', E', c')

if BellmanFord(G',s)=false then return "graph has negative cycles"

foreach $v \in V'$ do

 $h(v) \leftarrow d(s,v) \; / / \; d$ aus BellmanFord Algorithmus

foreach $(u,v) \in E'$ do

foreach $u \in V$ do

$$\tilde{d}(u,\cdot) \leftarrow \mathsf{Dijkstra}(\tilde{G}',u)$$

$$D(u,v) \leftarrow \tilde{d}(u,v) + h(v) - h(u)$$

771

Analyse

Laufzeiten

■ Berechnung von G': $\mathcal{O}(|V|)$

■ Bellman Ford G': $\mathcal{O}(|V| \cdot |E|)$

■ $|V| \times \text{Dijkstra } \mathcal{O}(|V| \cdot |E| \cdot \log |V|)$ (Mit Fibonacci-Heap:: $\mathcal{O}(|V|^2 \log |V| + |V| \cdot |E|)$)

Insgesamt $\mathcal{O}(|V| \cdot |E| \cdot \log |V|)$ $(\mathcal{O}(|V|^2 \log |V| + |V| \cdot |E|))$

25.8 A*-Algorithmus

Disclaimer

Diese Folien beinhalten die wichtigsten Formalien zum A*-Algorithmus und dessen Korrektheit. In der Vorlesung wird der Algorithmus motiviert und mit Beispielen unterlegt.

Eine sehr schöne Motivation des Algorithmus findet sich zum Beispiel hier: https://www.youtube.com/watch?v=bRvs8r0QU-Q

A*-Algorithmus

Voraussetzungen

- Positiv gewichteter Graph G = (V, E, c)
- G endlich oder δ -Graph: $\exists \delta > 0 : c(e) \geq \delta$ für alle $e \in E$
- $\blacksquare \ s \in V$, $t \in V$
- Abstandsschätzung $\hat{h}_t(v) \leq h_t(v) := \delta(v,t) \ \forall \ v \in V$.
- lacktriangle Gesucht: kürzester Pfad $p:s\leadsto t$

A*-Algorithmus(G, s, t, \hat{h})

Input: Positiv gewichteter Graph G=(V,E,c), Startpunkt $s\in V$, Endpunkt $t\in V$, Schätzung $\widehat{h}(v)<\delta(v,t)$

 $\textbf{Output:} \ \, \mathsf{Existenz} \ \, \mathsf{und} \ \, \mathsf{Wert} \ \, \mathsf{eines} \ \, \mathsf{k\"{u}rzesten} \ \, \mathsf{Pfades} \ \, \mathsf{von} \ \, s \ \, \mathsf{nach} \ \, t$

foreach $u \in V$ do

return failure

Warum der Algorithmus funktioniert

Lemma 24

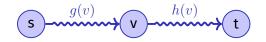
Sei $u \in V$ und, zu einem Zeitpunkt des A*-Algorithmus, $u \notin M$. Sei p ein kürzester Pfad von s nach u. Dann existiert ein $u' \in p$ mit $\widehat{g}(u') = g(u')$ und $u' \in R$.

Das Lemma besagt, dass es immer einen Knoten in der offenen Menge R gibt, dessen wahre Entfernung von s schon berechnet wurde und der zum kürzesten Pfad gehört (sofern ein solcher existiert).

Notation

Sei f(v) die Distanz eines kürzesten Weges von s nach t über v, also

$$f(v) := \underbrace{\delta(s,v)}_{g(v)} + \underbrace{\delta(v,t)}_{h(v)}$$



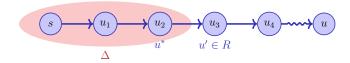
Sei p ein kürzester Weg von s nach t.

Dann gilt $f(s) = \delta(s, t)$ und f(v) = f(s) für alle $v \in p$.

Sei $\widehat{g}(v):=d[v]$ die Schätzung von g(v) in obigem Algorithms. Es gilt, dass $\widehat{g}(v)\geq g(v)$.

 $\widehat{h}(v)$ ist eine Schätzung von h(v) mit $\widehat{h}(v) \leq h(v)$.

Illustration und Beweis



Beweis: Wenn $s \in R$, dann $\widehat{g}(s) = g(s) = 0$. Sei also $s \notin R$.

Sei $p=\langle s=u_0,u_1,\ldots,u_k=u\rangle$ und $\Delta=\{u_i\in p,u_i\in M,\widehat{g}(u_i)=g(u_i)\}$. $\Delta\neq\emptyset$, denn $s\in\Delta$.

Sei $m = \max\{i : u_i \in \Delta\}$, $u^* = u_m$. Dann $u^* \neq u$, da $u \notin M$. Sei $u' = u_{m+1}$.

- 1. $\widehat{g}(u') \leq \widehat{g}(u^*) + c(u^*, u')$ weil u' schon relaxiert wurde
- **2.** $\widehat{g}(u^*) = g(u^*)$ (da $u^* \in \Delta$)
- 3. $\hat{g}(u') \ge g(u')$ (Konstruktion von \hat{g})
- **4.** $g(u') = g(u^*) + c(u^*, u')$ (da *p* optimal)

Also: $\widehat{g}(u') = g(u')$ und somit auch $u' \in R$ da $u' \notin \Delta$.

779

Folgerung

Corollary 25

Wenn $\hat{h}(u) \leq h(u)$ für alle $u \in V$ und A*- Algorithmus hat noch nicht terminiert. Dann existiert für jeden kürzesten Pfad p von s nach t ein Knoten $u' \in p$ mit $\hat{f}(u') \leq \delta(s,t) = f(t)$.

Wenn es einen kürzesten Weg p von s nach t gibt, steht also stets ein Knoten in der offenen Menge bereit, der die Gesamtentfernung maximal unterschätzt und der auf dem kürzesten Weg liegt.

Beweis des Corollars

Beweis:

Nach Lemma $\exists u' \in p \text{ mit } \widehat{g}(u') = g(u').$ Also:

$$\widehat{f}(u') = \widehat{g}(u') + \widehat{h}(u')$$

$$= g(u') + \widehat{h}(u')$$

$$\leq g(u') + h(u') = f(u')$$

Da p optimal: $f(u') = \delta(s, t)$.

782

78

Zulässigkeit

Theorem 26

Wenn es einen kürzesten Weg von s nach t gibt und $\hat{h}(u) \leq h(u) \ \forall \ u \in V$, dann terminiert der A*-Algorithmus mit $\ \hat{g}(t) = \delta(s,t)$

Beweis: Wenn der Algorithmus terminiert, dann terminiert er in t mit $f(t)=\widehat{g}(t)+0=g(t)$. Denn \widehat{g} überschätzt g höchstens und nach obigem Korrolar findet der Algorithmus stets ein Element $v\in R$ mit $f(v)\leq \delta(s,t)$. Der Algorithmus terminiert in endlichen vielen Schritten. Für endliche Graphen ist die maximale Anzahl an Relaxierschritten beschränkt.

Erneutes Besuchen von Knoten

- Der A*-Algorithmus kann Knoten mehrfach aus der Menge R entnehmen und sie später wieder einfügen.
- Das kann zu suboptimalem Verhalten im Sinne der Laufzeit des Algorithmus führen.
- Wenn \hat{h} zusätzlich zur Zulässigkeit $(\hat{h}(v) \le h(v))$ für alle $v \in V$) auch noch monoton ist, d.h. wenn für alle $(u, u') \in E$:

$$\hat{h}(u') \le \hat{h}(u) + c(u', u)$$

dann ist der A* Algorithmus äquivalent zum Dijkstra-Algorithmus mit Kantengewichten $\tilde{c}(u,v)=c(u,v)+\hat{h}(u)-\hat{h}(v)$ und kein Knoten wird aus R entnommen und wieder eingefügt.

■ Es ist allerdings nicht immer möglich, eine monotone Heuristik zu finden.

⁴⁴Für einen δ -Graphen ist die maximale Anzahl an Relaxierschritten bevor R nur noch Knoten mit $\hat{f}(s) > \delta(s,t)$ enthält, auch beschränkt. Das genaue Argument findet sich im Originalartikel Hart, P. E.; Nilsson, N. J.; Raphael, B. (1968). "A Formal Basis for the Heuristic Determination of Minimum Cost Paths".

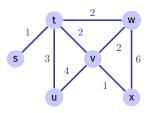
26. Minimale Spannbäume

Motivation, Greedy, Algorithmus von Kruskal, Allgemeine Regeln, Union-Find Struktur, Algorithmus von Jarnik, Prim, Dijkstra , Fibonacci Heaps [Ottman/Widmayer, Kap. 9.6, 6.2, 6.1, Cormen et al, Kap. 23, 19]

Problem

Gegeben: Ungerichteter, zusammenhängender, gewichteter Graph G=(V,E,c).

Gesucht: Minimaler Spannbaum T=(V,E'): zusammenhängender, zyklenfreier Teilgraph $E'\subset E$, so dass $\sum_{e\in E'}c(e)$ minimal.



86

Beispiele von Anwendungen

- Netzwerk-Design: finde das billigste / kürzeste Netz oder Leitungssystem, welches alle Knoten miteinander verbindet.
- Approximation einer Lösung des Travelling-Salesman Problems: finde einen möglichst kurzen Rundweg, welcher jeden Knoten einmal besucht.

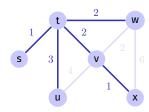
Greedy Verfahren

Zur Erinnerung:

- Gierige Verfahren berechnen eine Lösung schrittweise, indem lokal beste Lösungen gewählt werden.
- Die meisten Probleme sind nicht mit einer greedy Strategie lösbar.
- Das Problem des Minimalen Spannbaumes kann mit einem gierigen Verfahren effizient gelöst werden.

Greedy Idee (Kruskal, 1956)

Konstruiere ${\cal T}$ indem immer die billigste Kante hinzugefügt wird, welche keinen Zyklus erzeugt.



(Lösung ist nicht eindeutig.)

Korrektheit

Zu jedem Zeitpunkt ist (V,A) ein Wald, eine Menge von Bäumen. MST-Kruskal betrachtet jede Kante e_k einmal und wählt e_k oder verwirft e_k Notation (Momentaufnahme im Algorithmus)

- A: Menge gewählte Kanten
- R: Menge verworfener Kanten
- *U*: Menge der noch unentschiedenen Kanten

Algorithmus MST-Kruskal(G)

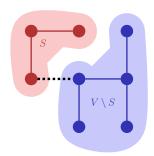
 $\label{eq:local_continuous_cont$

return (V, A, c)

790

Schnitt

Ein Schnitt von G ist eine Partition $S, V \setminus S$ von V. ($S \subseteq V$). Eine Kante kreuzt einen Schnitt, wenn einer Ihrer Endpunkte in S und der andere in $V \setminus S$ liegt.



Regeln

- 1. Auswahlregel: Wähle einen Schnitt, den keine gewählte Kante kreuzt. Unter allen unentschiedenen Kanten, welche den Schnitt kreuzen, wähle die mit minimalem Gewicht.
- 2. Verwerfregel: Wähle einen Kreis ohne verworfene Kanten. Unter allen unentschiedenen Kanten im Kreis verwerfe die mit maximalem Gewicht.

Regeln

Kruskal wendet beide Regeln an:

- 1. Ein gewähltes e_k verbindet zwei Zusammenhangskomponenten, sonst würde ein Kreis erzeugt werden. e_k ist beim Verbinden minimal, man kann also einen Schnitt wählen, den e_k mit minimalem Gewicht kreuzt.
- 2. Ein verworfenes e_k ist Teil eines Kreises. Innerhalb des Kreises hat e_k maximales Gewicht.

794

7/

Korrektheit

Theorem 27

Jeder Algorithmus, welcher schrittweise obige Regeln anwendet bis $U=\emptyset$ ist korrekt.

Folgerung: MST-Kruskal ist korrekt.

Auswahlinvariante

Invariante: Es gibt stets einen minimalen Spannbaum, der alle gewählten und keine der verworfenen Kanten enthält.

Wenn die beiden Regeln die Invariante erhalten, dann ist der Algorithmus sicher korrekt. Induktion:

- Zu Beginn: U = E, $R = A = \emptyset$. Invariante gilt offensichtlich.
- Invariante bleibt nach jedem Schritt des Algorithmus erhalten.
- Am Ende: $U = \emptyset$, $R \cup A = E \Rightarrow (V, A)$ ist Spannbaum.

Beweis des Theorems: zeigen nun, dass die beiden Regeln die Invariante erhalten.

Auswahlregel erhält Invariante

Es gibt stets einen minimalen Spannbaum T, der alle gewählten und keine der verworfenen Kanten enthält.

Wähle einen Schnitt, den keine gewählte Kante kreuzt. Unter allen unentschiedenen Kanten, welche den Schnitt kreuzen, wähle eine Kante e mit minimalem Gewicht.

- Fall 1: $e \in T$ (fertig)
- Fall 2: $e \notin T$. Dann hat $T \cup \{e\}$ einen Kreis, der e enthält. Kreis muss einen zweite Kante e' enthalten, welche den Schnitt auch kreuzt. Da $e' \notin R$ ist $e' \in U$. Somit $c(e) \le c(e')$ und $T' = T \setminus \{e'\} \cup \{e\}$ ist auch minimaler Spannbaum (und c(e) = c(e')).

Verwerfregel erhält Invariante

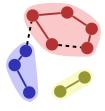
Es gibt stets einen minimalen Spannbaum ${\cal T}$, der alle gewählten und keine der verworfenen Kanten enthält.

Wähle einen Kreis ohne verworfene Kanten. Unter allen unentschiedenen Kanten im Kreis verwerfe die Kante e mit maximalem Gewicht.

- Fall 1: $e \notin T$ (fertig)
- Fall 2: $e \in T$. Entferne e von T, Das ergibt einen Schnitt. Diesen Schnitt muss eine weitere Kante e' aus dem Kreis kreuzen. Da $c(e') \le c(e)$ ist $T' = T \setminus \{e\} \cup \{e'\}$ auch minimal (und c(e) = c(e')).

Zur Implementation

Gegeben eine Menge von Mengen $i\equiv A_i\subset V$. Zur Identifikation von Schnitten und Kreisen: Zugehörigkeit der beiden Endpunkte einer Kante zu einer der Mengen.



Zur Implementation

Allgemeines Problem: Partition (Menge von Teilmengen) z.B. $\{\{1,2,3,9\},\{7,6,4\},\{5,8\},\{10\}\}$

Benötigt: Abstrakter Datentyp "Union-Find" mit folgenden Operationen

- Make-Set(*i*): Hinzufügen einer neuen Menge *i*.
- Find(e): Name i der Menge, welche e enthält.
- Union(i, j): Vereingung der Mengen mit Namen i und j.

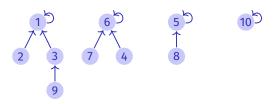
 $^{^{45}}$ Ein solcher Kreis enthält mindestens einen Knoten in S und einen in $V\setminus S$ und damit mindestens zwei Kanten zwischen S und $V\setminus S$.

Union-Find Algorithmus MST-Kruskal(G)

```
Input: Gewichteter Graph G = (V, E, c) Output: Minimaler Spannbaum mit Kanten A. Sortiere Kanten nach Gewicht c(e_1) \leq ... \leq c(e_m) A \leftarrow \emptyset for k = 1 to |V| do \bigcup MakeSet(k) for k = 1 to m do \bigcup (u, v) \leftarrow e_k if Find(u) \neq \operatorname{Find}(v) then \bigcup Union(Find(u), Find(v)) \bigcup A \leftarrow A \cup e_k else // konzeptuell: R \leftarrow R \cup e_k return (V, A, c)
```

Implementation Union-Find

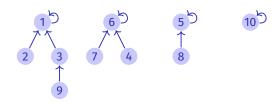
Idee: Baum für jede Teilmenge in der Partition, z.B. $\{\{1,2,3,9\},\{7,6,4\},\{5,8\},\{10\}\}$



803

Baumwurzeln = Namen (Stellvertreter) der Mengen, Bäume = Elemente der Mengen

Implementation Union-Find



Repräsentation als Array:

Index 1 2 3 4 5 6 7 8 9 10 Parent 1 1 1 6 5 6 5 5 3 10

Implementation Union-Find

802

Index 1 2 3 4 5 6 7 8 9 10 Parent 1 1 1 6 5 6 5 5 3 10

Make-Set(i)	$p[i] \leftarrow i$; return i	
Find(i)	while $(p[i] \neq i)$ do $i \leftarrow p[i]$ return i	
Union(i,j) ⁴⁶	$p[j] \leftarrow i;$	

 $^{^{46}}i$ und j müssen Namen (Wurzeln) der Mengen sein. Andernfalls verwende Union(Find(i),Find(j))

Optimierung der Laufzeit für Find

Baum kann entarten: Beispiel Union(8,7), Union(7,6), Union(6,5), ...

Index 1 2 3 4 5 6 7 8 .. Parent 1 1 2 3 4 5 6 7 ...

Laufzeit von Find im schlechtesten Fall in $\Theta(n)$.

Optimierung der Laufzeit für Find

Idee: Immer kleineren Baum unter grösseren Baum hängen. Benötigt zusätzliche Grösseninformation (Array) g

$$\label{eq:make-Set} \begin{array}{ll} \text{Make-Set(i)} & p[i] \leftarrow i; \ g[i] \leftarrow 1; \ \text{return } i \\\\ \hline\\ \text{Union(i, j)} & \begin{array}{ll} \text{if} \ g[j] > g[i] \ \text{then swap(i, j)} \\\\ p[j] \leftarrow i \\\\ \text{if} \ g[i] = g[j] \ \text{then } g[i] \leftarrow g[i] + 1 \end{array}$$

 \Rightarrow Baumtiefe (und schlechteste Laufzeit für Find) in $\Theta(\log n)$

Theorem 28

Beobachtung

Obiges Verfahren Vereinigung nach Grösse konserviert die folgende Eigenschaft der Bäume: ein Baum mit Höhe h hat mindestens 2^h Knoten.

Unmittelbare Folgerung: Laufzeit Find = $O(\log n)$.

Beweis

Induktion: nach Vorraussetzung haben Teilbäume jeweils mindestens 2^{h_i} Knoten. ObdA: $h_2 \leq h_1$.

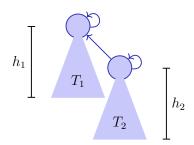
■
$$h_2 < h_1$$
:

$$h(T_1 \oplus T_2) = h_1 \Rightarrow g(T_1 \oplus T_2) \ge 2^h$$

 $h_2 = h_1$:

$$g(T_1) \ge g(T_2) \ge 2^{h_2}$$

 $\Rightarrow g(T_1 \oplus T_2) = g(T_1) + g(T_2) \ge 2 \cdot 2^{h_2} = 2^{h(T_1 \oplus T_2)}$



Weitere Verbesserung

Bei jedem Find alle Knoten direkt an den Wurzelknoten hängen. Find(i):

$$\begin{split} j &\leftarrow i \\ \text{while } (p[i] \neq i) \text{ do } i \leftarrow p[i] \\ \text{while } (j \neq i) \text{ do} \\ \mid & t \leftarrow j \\ j \leftarrow p[j] \\ p[t] \leftarrow i \end{split}$$

return i

Laufzeit: amortisiert fast konstant (Inverse der Ackermannfunktion).⁴⁷

Laufzeit des Kruskal Algorithmus

- Sortieren der Kanten: $\Theta(|E|\log|E|) = \Theta(|E|\log|V|)$. ⁴⁸
- Initialisieren der Union-Find Datenstruktur $\Theta(|V|)$
- $|E| \times \text{Union(Find(}x\text{),Find(}y\text{))}$: $\mathcal{O}(|E|\log|E|) = \mathcal{O}(|E|\log|V|)$. Insgesamt $\Theta(|E|\log|V|)$.

Algorithmus von Jarnik (1930), Prim, Dijkstra (1959)

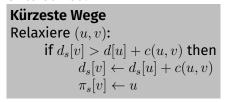
Idee: Starte mit einem $v \in V$ und lasse von dort einen Spannbaum wachsen:

$$\begin{array}{l} A \leftarrow \emptyset \\ S \leftarrow \{v_0\} \\ \text{for } i \leftarrow 1 \text{ to } |V| \text{ do} \\ & \quad \text{W\"{a}hle billigste } (u,v) \text{ mit } u \in S, \, v \not \in S \\ & \quad A \leftarrow A \cup \{(u,v)\} \\ & \quad S \leftarrow S \cup \{v\} \; // \; (\text{F\"{a}rbung}) \end{array}$$

Anmerkung: man benötigt keine Union-Find Datenstruktur. Es genügt, Knoten zu färben, sobald sie zu ${\cal S}$ hinzugenommen werden.

Implementation and Laufzeit

Implementation wie bei Dijkstra's Kürzeste Wege Algorithmus. Einziger Unterschied:



- Mit Min-Heap, Kosten $\mathcal{O}(|E| \cdot \log |V|)$:
 - lacktriang Initialisierung (Knotenfärbung) $\mathcal{O}(|V|)$
 - $|V| \times \text{ExtractMin} = \mathcal{O}(|V| \log |V|)$,
 - lacksquare |E| imes Insert oder DecreaseKey: $\mathcal{O}(|E|\log|V|)$,
- Mit Fibonacci-Heap: $\mathcal{O}(|E| + |V| \cdot \log |V|)$.

811

⁴⁷Wird hier nicht vertieft.

 $^{^{48}}$ da G zusammenhängend: $|V| \leq |E| \leq |V|^2$

Fibonacci Heaps

Datenstruktur zur Verwaltung von Elementen mit Schlüsseln. Operationen

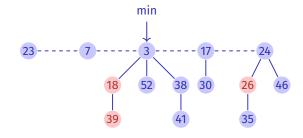
- MakeHeap(): Liefere neuen Heap ohne Elemente
- Insert(H, x): Füge x zu H hinzu
- lacktriangle Minimum(H): Liefere Zeiger auf das Element m mit minimalem Schlüssel
- **ExtractMin**(H): Liefere und entferne (von H) Zeiger auf das Element m
- Union (H_1, H_2) : Liefere Verschmelzung zweier Heaps H_1 und H_2
- **DecreaseKey**(H, x, k): Verkleinere Schlüssel von x in H zu k
- Delete (H, x): Entferne Element x von H

Vorteil gegenüber Binary Heap?

	Binary Heap (worst-Case)	Fibonacci Heap (amortisiert)
MakeHeap	$\Theta(1)$	$\Theta(1)$
Insert	$\Theta(\log n)$	$\Theta(1)$
Minimum	$\Theta(1)$	$\Theta(1)$
ExtractMin	$\Theta(\log n)$	$\Theta(\log n)$
Union	$\Theta(n)$	$\Theta(1)$
DecreaseKey	$\Theta(\log n)$	$\Theta(1)$
Delete	$\Theta(\log n)$	$\Theta(\log n)$

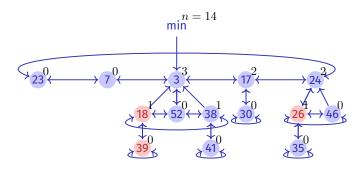
Struktur

Menge von Bäumen, welche der Min-Heap Eigenschaft genügen. Markierbare Knoten.



Implementation

Doppelt verkettete Listen von Knoten mit marked-Flag und Anzahl Kinder. Zeiger auf das minimale Element und Anzahl Knoten.



Einfache Operationen

- MakeHeap (trivial)
- Minimum (trivial)
- Insert(H, e)
 - 1. Füge neues Element in die Wurzelliste ein
 - 2. Wenn Schlüssel kleiner als Minimum, min-pointer neu setzen.
- Union (H_1, H_2)
 - 1. Wurzellisten von H_1 und H_2 aneinander hängen
 - 2. Min-Pointer neu setzen.
- \blacksquare Delete(H, e)
 - 1. DecreaseKey($H, e, -\infty$)
 - 2. ExtractMin(H)

ExtractMin

- 1. Entferne Minimalknoten m aus der Wurzelliste
- 2. Hänge Liste der Kinder von m in die Wurzelliste
- 3. Verschmelze solange heapgeordnete Bäume gleichen Ranges, bis alle Bäume unterschiedlichen Rang haben: Rangarray $a[0,\ldots,n]$ von Elementen, zu Beginn leer. Für jedes Element e der Wurzelliste:
 - a Sei g der Grad von e.
 - b Wenn $a[g] = nil: a[g] \leftarrow e$.
 - c Wenn $e':=a[g] \neq nil$: Verschmelze e mit e' zu neuem e'' und setze $a[g] \leftarrow nil$. Setze e'' unmarkiert Iteriere erneut mit $e \leftarrow e''$ vom Grad g+1.

818

DecreaseKey (H, e, k)

- 1. Entferne e von seinem Vaterknoten p (falls vorhanden) und erniedrige den Rang von p um eins.
- 2. Insert(H, e)
- 3. Vermeide zu dünne Bäume:
 - a Wenn p = nil , dann fertig
 - b Wenn p unmarkiert: markiere p, fertig.
 - c Wenn p markiert: unmarkiere p, trenne p von seinem Vater pp ab und Insert(H,p). Iteriere mit $p \leftarrow pp$.

Abschätzung für den Rang

Theorem 29

Sei p Knoten eines F-Heaps H. Ordnet man die Söhne von p in der zeitlichen Reihenfolge, in der sie an p durch Zusammenfügen angehängt wurden, so gilt: der i-te Sohn hat mindestens Rang i-2

Beweis: p kann schon mehr Söhne gehabt haben und durch Abtrennung verloren haben. Als der ite Sohn p_i angehängt wurde, müssen p und p_i jeweils mindestens Rang i-1 gehabt haben. p_i kann maximal einen Sohn verloren haben (wegen Markierung), damit bleibt mindestens Rang i-2.

Abschätzung für den Rang

Theorem 30

Jeder Knoten p vom Rang k eines F-Heaps ist Wurzel eines Teilbaumes mit mindestens F_{k+1} Knoten. (F: Fibonacci-Folge)

Beweis: Sei S_k Minimalzahl Nachfolger eines Knotens vom Rang k in einem F-Heap plus 1 (der Knoten selbst). Klar: $S_0=1$, $S_1=2$. Nach vorigem Theorem $S_k \geq 2 + \sum_{i=0}^{k-2} S_i, k \geq 2$ (p und Knoten p_1 jeweils 1). Für Fibonacci-Zahlen gilt (Induktion) $F_k \geq 2 + \sum_{i=2}^k F_i, k \geq 2$ und somit (auch Induktion) $S_k \geq F_{k+2}$. Fibonacci-Zahlen wachsen exponentiell ($\mathcal{O}(\varphi^k)$) Folgerung: Maximaler Grad eines beliebigen Knotens im Fibonacci-Heap mit n Knoten ist $\mathcal{O}(\log n)$.

Amortisierte Worst-case-Analyse Fibonacci Heap

t(H): Anzahl Bäume in der Wurzelliste von H, m(H): Anzahl markierte Knoten in H ausserhalb der Wurzelliste, Potentialfunktion $\Phi(H)=t(H)+2\cdot m(H)$. Zu Anfang $\Phi(H)=0$. Potential immer nichtnegativ. Amortisierte Kosten:

- Insert(H,x): t'(H) = t(H) + 1, m'(H) = m(H), Potentialerhöhung 1, Amortisierte Kosten $\Theta(1) + 1 = \Theta(1)$
- Minimum(H): Amortisierte Kosten = tatsächliche Kosten = $\Theta(1)$
- Union (H_1, H_2) : Amortisierte Kosten = tatsächliche Kosten = $\Theta(1)$

2

Amortisierte Kosten ExtractMin

- \blacksquare Anzahl der Bäume in der Wurzelliste t(H).
- Tatsächliche Kosten der ExtractMin Operation: $\mathcal{O}(\log n + t(H))$.
- Nach dem Verschmelzen noch $\mathcal{O}(\log n)$ Knoten.
- Anzahl der Markierungen kann beim Verschmelzen der Bäume maximal kleiner werden.
- Amortisierte Kosten von ExtractMin also maximal

$$\mathcal{O}(\log n + t(H)) + \mathcal{O}(\log n) - \mathcal{O}(t(H)) = \mathcal{O}(\log n).$$

Amortisierte Kosten DecreaseKey

- Annahme: DecreaseKey führt zu c Abtrennungen eines Knotens von seinem Vaterknoten, tatsächliche Kosten $\mathcal{O}(c)$
- c Knoten kommen zur Wurzelliste hinzu
- lacktriang Löschen von (c-1) Markierungen, Hinzunahme maximal einer Markierung
- Amortisierte Kosten von DecreaseKey:

$$\mathcal{O}(c) + (t(H) + c) + 2 \cdot (m(H) - c + 2)) - (t(H) + 2m(H)) = \mathcal{O}(1)$$

27. Flüsse in Netzen

Flussnetzwerk, Maximaler Fluss, Schnitt, Restnetzwerk, Max-flow Min-cut Theorem, Ford-Fulkerson Methode, Edmonds-Karp Algorithmus, Maximales Bipartites Matching [Ottman/Widmayer, Kap. 9.7, 9.8.1], [Cormen et al, Kap. 26.1-26.3]

27.1 Flussnetzwerk

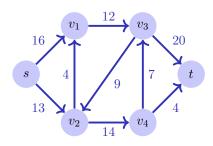
Motivation

- Modelliere Fluss von Flüssigkeiten, Bauteile auf Fliessbändern, Strom in elektrischen Netwerken oder Information in Kommunikationsnetzwerken.
- Konnektivität von Kommunikationsnetzwerken, Bipartites Matching, Zirkulationen, Scheduling, Image Segmentation, Baseball Eliminination...

6

Flussnetzwerk

- Flussnetzwerk G = (V, E, c): gerichteter Graph mit Kapazitäten
- Antiparallele Kanten verboten: $(u, v) \in E \Rightarrow (v, u) \notin E$.
- Fehlen einer Kante (u, v) auch modelliert durch c(u, v) = 0.
- Quelle s und Senke t: spezielle Knoten. Jeder Knoten v liegt auf einem Pfad zwischen s und t: $s \leadsto v \leadsto t$

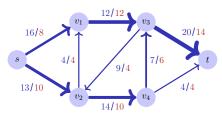


Fluss

Ein Fluss $f:V\times V\to\mathbb{R}$ erfüllt folgende Bedingungen:

- Kapazitätsbeschränkung: Für alle $u, v \in V$: $f(u, v) \le c(u, v)$.
- $\begin{tabular}{l} \blacksquare & Schiefsymmetrie: \\ & F\"{u}r alle \ u,v \in V\hbox{:} \ f(u,v) = -f(v,u). \\ \end{tabular}$
- Flusserhaltung: Für alle $u \in V \setminus \{s, t\}$:

$$\sum_{v \in V} f(u, v) = 0.$$



Wert w des Flusses: $|f| = \sum_{v \in V} f(s, v)$.

Hier
$$|f| = 18$$
.

Wie gross kann ein Fluss sein?

Begrenzende Faktoren: Schnitte

- s von t trennender Schnitt: Partitionierung von V in S und T mit $s \in S$, $t \in T$.
- Kapazität eines Schnittes: $c(S,T) = \sum_{v \in S, v' \in T} c(v,v')$
- Minimaler Schnitt: Schnitt mit minimaler Kapazität.
- Fluss über Schnitt: $f(S,T) = \sum_{v \in S, v' \in T} f(v,v')$

830

831

Implizites Summieren

Notation: Seien $U, U' \subseteq V$

$$f(U, U') := \sum_{\substack{u \in U \\ u' \in U'}} f(u, u'), \qquad f(u, U') := f(\{u\}, U')$$

Somit

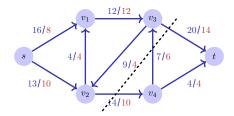
- $\blacksquare |f| = f(s, V)$
- f(U,U) = 0

- f(R,V) = 0 wenn $R \cap \{s,t\} = \emptyset$. [Flusserhaltung!]

Wie gross kann ein Fluss sein?

Es gilt für jeden Fluss und jeden Schnitt, dass f(S,T) = |f|:

$$f(S,T) = f(S,V) - \underbrace{f(S,S)}_{0} = f(S,V)$$
$$= f(s,V) + \underbrace{f(S-\{s\},V)}_{\not\ni t,\not\ni s} = |f|.$$



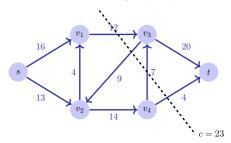
27.3 Maximaler Fluss

Maximaler Fluss?

Es gilt insbesondere für alle Schnitte (S,T) von V.

$$|f| \le \sum_{v \in S, v' \in T} c(v, v') = c(S, T)$$

Werden sehen, dass Gleicheit gilt für $\min_{S,T} c(S,T)$.

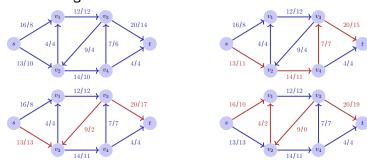


835

834

Maximaler Fluss?

Naives Vorgehen:



Folgerung: Greedy Flusserhöhung löst das Problem nicht.

Die Ford-Fulkerson Methode

- Starte mit f(u, v) = 0 für alle $u, v \in V$
- lacksquare Bestimme Restnetzwerk* G_f und Erweiterungspfad in G_f
- Erhöhe Fluss über den Erweiterungspfad*
- Wiederholung bis kein Erweiterungspfad mehr vorhanden.

$$G_f := (V, E_f, c_f)$$

 $c_f(u, v) := c(u, v) - f(u, v) \quad \forall u, v \in V$
 $E_f := \{(u, v) \in V \times V | c_f(u, v) > 0\}$

*Wird im Folgenden erklärt

Flusserhöhung, negativ

Sei ein Fluss f im Netzwerk gegeben.

Erkenntnis:

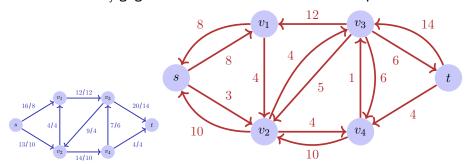
- Flusserhöhung in Richtung einer Kante möglich, wenn Fluss entlang der Kante erhöht werden kann, also wenn f(u,v) < c(u,v). Restkapazität $c_f(u,v) = c(u,v) f(u,v) > 0$.
- Flusserhöhung entgegen der Kantenrichtung möglich, wenn Fluss entlang der Kante verringert werden kann, also wenn f(u,v)>0. Restkapazität $c_f(v,u)=f(u,v)>0$.

27.4 Restnetzwerk

38

Restnetzwerk

Restnetzwerk G_f gegeben durch alle Kanten mit Restkapazität:



Restnetzwerke haben dieselben Eigenschaften wie Flussnetzwerke, ausser dass antiparallele Kapazitäten-Kanten zugelassen sind.

Beobachtung

Theorem 31

Sei G=(V,E,c) ein Flussnetzwerk mit Quelle s und Senke t und f ein Fluss in G. Sei G_f das dazugehörige Restnetzwerk und sei f' ein Fluss in G_f . Dann definiert $f\oplus f'$ mit

$$(f \oplus f')(u,v) = f(u,v) + f'(u,v)$$

einen Fluss in G mit Wert |f| + |f'|.

Beweis

 $f \oplus f'$ ist ein Fluss in G:

■ Kapazitätsbeschränkung

$$(f \oplus f')(u,v) = f(u,v) + \underbrace{f'(u,v)}_{\leq c(u,v) - f(u,v)} \leq c(u,v)$$

Schiefsymmetrie

$$(f \oplus f')(u,v) = -f(v,u) + -f'(v,u) = -(f \oplus f')(v,u)$$

■ Flusserhaltung $u \in V - \{s, t\}$:

$$\sum_{v \in V} (f \oplus f')(u, v) = \sum_{v \in V} f(u, v) + \sum_{v \in V} f'(u, v) = 0$$

Beweis

Wert von $f \oplus f'$

$$|f \oplus f'| = (f \oplus f')(s, V)$$

$$= \sum_{u \in V} f(s, u) + f'(s, u)$$

$$= f(s, V) + f'(s, V)$$

$$= |f| + |f'|$$

842

Erweiterungspfade

Erweiterungspfad p: einfacher Pfad von s nach t im Restnetzwerk G_f . Restkapazität $c_f(p) = \min\{c_f(u,v): (u,v) \text{ Kante in } p\}$

Fluss in G_f

Theorem 32

Die Funktion $f_p: V \times V \to \mathbb{R}$,

$$f_p(u,v) = \begin{cases} c_f(p) & \text{wenn } (u,v) \text{ Kante in } p \\ -c_f(p) & \text{wenn } (v,u) \text{ Kante in } p \\ 0 & \text{sonst} \end{cases}$$

ist ein Fluss in G_f mit dem Wert $|f_p| = c_f(p) > 0$.

 f_p ist ein Fluss (leicht nachprüfbar). Es gibt genau einen Knoten $u \in V$ mit $(s,u) \in p$. Somit $|f_p| = \sum_{v \in V} f_p(s,v) = f_p(s,u) = c_f(p)$.

Folgerung

Strategie für den Algorithmus:

Mit einem Erweiterungspfad p in G_f definiert $f \oplus f_p$ einen neuen Fluss mit Wert $|f \oplus f_p| = |f| + |f_p| > |f|$.

27.5 Max-Flow Min-Cut

846

Max-Flow Min-Cut Theorem

Theorem 33

Wenn f ein Fluss in einem Flussnetzwerk G=(V,E,c) mit Quelle s und Senke t is, dann sind folgende Aussagen äquivalent:

- 1. f ist ein maximaler Fluss in G
- 2. Das Restnetzwerk G_f enthält keine Erweiterungspfade
- 3. Es gilt |f| = c(S, T) für einen Schnitt (S, T) von G.

Beweis

- $(3) \Rightarrow (1)$: Es gilt $|f| \le c(S,T)$ für alle Schnitte S,T. Aus |f| = c(S,T) folgt also |f| maximal.
- (1) \Rightarrow (2): f maximaler Fluss in G. Annahme: G_f habe einen Erweiterungsfad. Dann gilt $|f \oplus f_p| = |f| + |f_p| > |f|$. Widerspruch.

Beweis $(2) \Rightarrow (3)$

Annahme: G_f habe keinen Erweiterungsfad Definiere $S=\{v\in V: \text{ es existiert Pfad } s\leadsto v \text{ in } G_f\}.$ $(S,T):=(S,V\setminus S)$ ist ein Schnitt: $s\in S,t\in T.$ Sei $u\in S$ und $v\in T.$ Dann $c_f(u,v)=0$, also $c_f(u,v)=c(u,v)-f(u,v)=0$. Somit f(u,v)=c(u,v).

$$|f| = f(S,T) = \sum_{u \in S} \sum_{v \in T} f(u,v) = \sum_{u \in S} \sum_{v \in T} c(u,v) = C(S,T).$$

27.6 Ford-Fulkerson Algorithmus

850

Algorithmus Ford-Fulkerson(G, s, t)

```
\label{eq:local_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_continuous_cont
```

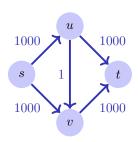
Praktische Anmerkung

In einer Implementation des Ford-Fulkerson Algorithmus werden die negativen Flusskanten normalerweise nicht gespeichert, da ihr Wert sich stets als der negierter Wert der Gegenkante ergibt.

$$\begin{split} f(u,v) &\leftarrow f(u,v) + c_f(p) \\ f(v,u) &\leftarrow f(v,u) - c_f(p) \\ \text{wird dann zu} \\ \text{if } (u,v) &\in E \text{ then} \\ \mid f(u,v) \leftarrow f(u,v) + c_f(p) \\ \text{else} \\ \mid f(v,u) \leftarrow f(v,u) - c_f(p) \end{split}$$

Analyse

- Der Ford-Fulkerson Algorithmus muss für irrationale Kapazitäten nicht einmal terminieren!
 Für ganze oder rationale Zahlen terminiert der Algorithmus.
- Für ganzzahligen Fluss benötigt der Algorithmus maximal $|f_{\max}|$ Durchläufe der While-Schleife (denn der Fluss erhöht sich mindestens um 1). Suche einzelner zunehmender Weg (z.B. Tiefensuche oder Breitensuche) $\mathcal{O}(|E|)$. Also $\mathcal{O}(f_{\max}|E|)$.



Bei schlecht gewählter Strategie benötigt der Algorithmus hier bis zu 2000 Iterationen.

27.7 Edmonds-Karp Algorithmus

Edmonds-Karp Algorithmus

Wähle in der Ford-Fulkerson-Methode zum Finden eines Pfades in G_f jeweils einen Erweiterungspfad kürzester Länge (z.B. durch Breitensuche).

Edmonds-Karp Algorithmus

Theorem 34

Wenn der Edmonds-Karp Algorithmus auf ein ganzzahliges Flussnetzwerk G=(V,E) mit Quelle s und Senke t angewendet wird, dann ist die Gesamtanzahl der durch den Algorithmus angewendete Flusserhöhungen in $\mathcal{O}(|V|\cdot|E|)$.

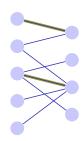
 \Rightarrow Gesamte asymptotische Laufzeit: $\mathcal{O}(|V| \cdot |E|^2)$

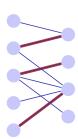
[Ohne Beweis]

27.8 Maximales Bipartites Matching

Anwendung: Maximales bipartites Matching

Gegeben: bipartiter ungerichteter Graph G=(V,E). Matching M: $M\subseteq E$ so dass $|\{m\in M:v\in m\}|\leq 1$ für alle $v\in V$. Maximales Matching M: Matching M, so dass $|M|\geq |M'|$ für jedes Matching M'.

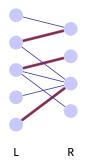


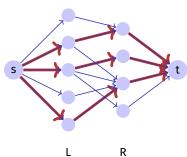


8

Korrespondierendes Flussnetzwerk

Konstruiere zur einer Partition L,R eines bipartiten Graphen ein korrespondierendes Flussnetzwerk mit Quelle s und Senke t, mit gerichteten Kanten von s nach L, von L nach R und von R nach t. Jede Kante bekommt Kapazität 1.





Ganzzahligkeitstheorem

Theorem 35

Wenn die Kapazitäten eines Flussnetzwerks nur ganzzahlige Werte anehmen, dann hat der durch Ford-Fulkerson erzeugte maximale Fluss die Eigenschaft, dass der Wert von f(u,v) für alle $u,v\in V$ eine ganze Zahl ist.

[ohne Beweis]

Folgerung: Ford Fulkerson erzeugt beim zum bipartiten Graph gehörenden Flussnetzwerk ein maximales Matching $M = \{(u, v) : f(u, v) = 1\}$.

27.9 Push-Relabel Algorithmus

Disclaimer

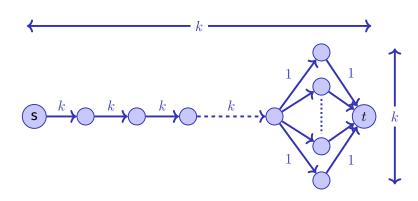
Die folgenden Folien beinhalten die wichtigsten Formalien zum Push-Relabel Algorithmus und dessen Korrektheit. Es fehlt noch ein Beispiel. In der Vorlesung wird der Algorithmus motiviert und mit Beispielen unterlegt.

Die Konzeption dieser Vorlesung ist übernommen von Tim Roughgarden (Stanford)

https://www.youtube.com/watch?v=0h189H39USg

52

Beispiel



Der Ford-Fulkerson Algorithmus (und Edmonds Karp) führt hier $\Omega(k^2)$ viele Schritte aus.

Pre-Flow

Ein Pre-Flow $f:V\times V\to\mathbb{R}$ ist ein Fluss mit einer abgeschwächten Flusserhaltung:

- Kapazitätsbeschränkung: Für alle $u, v \in V$: $f(u, v) \le c(u, v)$.
- Schiefsymmetrie: Für alle $u, v \in V$: f(u, v) = -f(v, u).
- Abgeschwächte Flusserhaltung: Für alle $u \in V \setminus \{s, t\}$:

$$\alpha_f(u) := \sum_{v \in V} f(v, u) \ge 0.$$



Knoten mit Exzess $\alpha_f(u) = 3 + 2 - 1 - 2 = 2.$

Der Wert $\alpha_f(u)$ heisst Exzess von f in u

Algorithmus Push(u, v)

Das Erweiterungsnetzwerk $G_f = (V, E_f, c_f)$ ist auf einem Pre-Flow f definiert wie vorher auf einem Fluss.

$$\begin{array}{c|c} \textbf{if} \ \alpha_f(u) > 0 \ \textbf{then} \\ & \textbf{if} \ c_f(u,v) > 0 \ \text{in} \ G_f \ \textbf{then} \\ & \ \ \, \Delta \leftarrow \min\{c_f(u,v),\alpha_f(u)\} \\ & \ \ \, f(u,v) \leftarrow f(u,v) + \Delta. \end{array}$$

Höhenfunktion

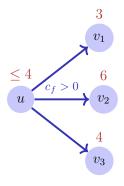
Eine Höhenfunktion $h:V\to\mathbb{N}_0$ auf G wird dafür sorgen, dass der Fluss nicht unendlich oft im Kreis weitergegeben wird. Ausserdem sorgen die folgenden Invarianten dafür, dass im Erweiterungsnetzwerk s von t getrennt bleibt.

Invarianten der Höhenfunktion

- 1. h(s) = n
- 2. h(t) = 0
- Für alle $u, v \in V$ mit $c_f(u, v) > 0$ gilt $h(u) \le h(v) + 1$.

66

Beispiel



Kanten im Erweiterungsnetzwerk gehen höchstens einen Schritt bergab (oder bleiben gleich oder gehen bergauf).

Kein Erweiterungspfad

Die Länge eines Pfades von s nach t im Erweiterungsnetzwerk ist maximal n-1. Da für jede Kante (u,v) mit $c_f(u,v)>0$ gilt, dass $h(v)\geq h(u)-1$ und da h(s)=n und h(t)=0 (der Pfad von s nach t müsste also mindestens Länge n haben), kann es bei Einhalten der Höheninvarianten keinen Erweiterungspfad geben.

Strategien

Ford-Fulkerson (konservativ)

Invariante: FlusserhaltungSchritte: Erweiterungspfade

Tiel: Trenne s von t im Erweiterungsnetzwerk G_f .

Push-Relabel

Invariante: Höheninvariante (kein Erweiterungspfad!)

Schritte: Push Flow

Ziel: Stelle Flusserhaltung her.

Push-Relabel-Algorithmus

```
\begin{array}{l} \text{Input: Flussnetzwerk } G=(V,E,c)\text{, Quelle } s \text{ und Senke } t. \ n:=|v| \\ h(s) \leftarrow n \\ \text{foreach } v \neq s \text{ do } h(v) \leftarrow 0 \\ \text{foreach } (u,v) \in E \text{ do } f(u,v) \leftarrow 0 \\ \text{foreach } (s,v) \in E \text{ do } f(s,v) \leftarrow c(s,v) \\ \\ \text{while } \exists u \in V \setminus \{s,t\} : \alpha_f(u) > 0 \text{ do} \\ \text{wähle } u \text{ mit } \alpha_f(u) > 0 \text{ und maximalem } h(u) \\ \text{if } \exists v \in V : c_f(u,v) > 0 \ \land \ h(v) = h(u) - 1 \text{ then} \\ \text{| } \text{push}(u,v) \\ \text{else} \\ \text{| } h(u) \leftarrow h(u) + 1 \\ \end{array} \right. // \text{ relabel} \\ \\ \end{array}
```

870

Korrektheit: Invarianten-Lemma

Lemma 36

Während der Ausführung des Push-Relabel Algorithmus bleiben die Invarianten zur Höhenfunktion erhalten.

Unmittelbare Folgerung: Wenn der Push-Relabel-Algorithmus terminiert, terminiert er mit einem maximalen Fluss

Invarianten-Lemma: Beweis

Beweis:

Nach der Initialisierung sind die Invarianten erhalten, denn nur für Kanten (s,u) ist die Höhendifferenz kleiner als -1, dort gilt $c_f(s,u)=0$ Invarianten für s und t bleiben erhalten, da die Höhe von s und t nicht verändert wird.

Ausführung von push(u, v) erzeugt im Restnetzwerk hochstens die Kante (v, u), mit h(v) > h(u)

Ausführung von Relabel findet nur statt, wenn keine Abwärtskante vorhanden ist. Damit gilt danach $h(u) \ge h(v) - 1$ für alle Kanten (u, v)

Terminierung und Laufzeit

Theorem 37

Der Push-Relabel Algorithmus terminiert nach

- \square $\mathcal{O}(n^2)$ Relabel-Operationen und
- \square $\mathcal{O}(n^3)$ Push-Operationen.

Der Beweis wird im Folgenden für Relabel und Push-Operationen separat geführt.

Hauptlemma

Lemma 38

Sei f ein Pre-Flow in G. Wenn $\alpha_f(u) > 0$ gilt für einen Knoten $u \in V - \{s,t\}$, dann gibt es einen Pfad $p:u \leadsto s$ im Erweiterungsnetzwerk G_f

74

Hauptlemma: Beweis

Beweis: Sei $A:=\{u\in V:\exists p:s\leadsto u \text{ with } f(e)>0\ \forall\ e\in p\}$ und $B:=V\setminus A$. Für jedes $u\in A$ gibt es einen Pfad von s mit positiven Fluss. Daher gibt es im Erweiterungsnetzwerk einen Pfad von u nach s.

Sei $u \in B$. Dann $\sum_{v \in V} f(v, u) \ge 0$, denn f ist pre-flow.

Aber auch
$$\sum_{v \in V} \sum_{u \in B} f(v, u) = \underbrace{\sum_{v \in A} \sum_{u \in B} f(v, u)}_{\leq 0} + \underbrace{\sum_{v \in B} \sum_{u \in B} f(v, u)}_{=0} \leq 0$$
 denn es

kann keine Kante mit positivem Fluss von A nach B geben und für jede Kante innerhalb von B gilt f(u,v)=-f(v,u). $\Rightarrow \alpha_f(u)=0 \ \forall \ u \in B$. Also impliziert $\alpha_f(u)>0$ dass $u \in A$.

Maximale Knotenhöhe

Corollary 39

Während der Ausführung des Push-Relabel Algorithmus gilt stets h(u) < 2n für alle $u \in V$.

Beweis:

Hauptlemma: für jeden Knoten u mit $\alpha_f(u)>0$ gibt es Pfad $p:u\leadsto s$ in Restnetzwerk.

Höheninvarianten: Kanten in G_f gehen maximal einen Schritt abwärts , $h(s)=n. \label{eq:hohen}$

Maximale Länge von $p:u\leadsto s$ (zyklenfrei!) ist $n-1.\Rightarrow$ Maximale Höhe Knoten ist n+n-1=2n-1.

Anzahl Relabels

Aus dem vorherigen Korollar folgt direkt

Corollary 40

Der Push-Relabel Algorithmus führt $\mathcal{O}(n^2)$ Relabel-Operationen aus.

(Nicht) sättigende Push-Operationen

push(u, v) heisst

 \blacksquare sättigend, wenn $c_f(u,v) \leq \alpha_f(u)$

$$\alpha_f = 3 \qquad \alpha_f = 1$$

$$u \xrightarrow{c_f = 2} v \Rightarrow u \leftarrow v$$

■ nicht sättigend, wenn $c_f(u,v) > \alpha_f(u)$

$$\begin{array}{ccc}
\alpha_f = 3 & & \alpha_f = 0 \\
u & & c_f = 4
\end{array}$$

$$\Rightarrow & \begin{array}{c}
\alpha_f = 0 \\
v \\
\end{array}$$

78

Anzahl sättigende Push-Operationen

Lemma 41

Zwischen zwei sättigenden Push-Operationen auf derselben Kante (u,v) führt der Push-Relabel Algorithmus auf u und v jeweils mindestens zwei Relabel- Operationen aus

Unmittelbare Folgerung: es finden insgesamt $\mathcal{O}(n^3)$ sättigende Push-Operationen statt, denn für jeden Knoten werden nach Korollar 39 $\mathcal{O}(n)$ Relabel-Operationen ausgeführt.

Beweis: Anzahl sättigende Push-Operationen

Beweis:

Nach sättigendem push(u, v) (mit h(u) = h(v) + 1) verschwindet Kante (u, v) vom Restnetzwerk.

Damit Kante (u,v) wieder im Restnetzwerk erscheint, muss $\mathtt{push}(v,u)$ (Gegenkante) ausgeführt werden. Dafür muss h(v) = h(u) + 1 gelten, also sind mindestens zwei Relabel Operationen auf v nötig.

Nachfolgend sind mindestens zwei Relabel-Operationen auf u nötig, damit $\mathtt{push}(u,v)$ ausgeführt werden kann.

Anzahl nichtsättigende Push-Operationen

Lemma 42

Zwischen zwei Relabel-Operationen führt der Push-Relabel Algorithmus maximal n nicht-sättigende Push-Operationen aus.

Unmittelbare Folgerung: es finden insgesamt $\mathcal{O}(n^3)$ nichtsättigende Push-Operationen statt, denn nach Korollar 40 finden insgesamt $\mathcal{O}(n^2)$ Relabel-Operationen statt

882

28. Parallel Programming I

Moore's Law und The Free Lunch, Hardware Architekturen, Parallele Ausführung Multi-Threading, Parallelität und Nebenläufigkeit, C++ Threads, Skalierbarkeit: Amdahl und Gustafson, Daten- und Taskparallelität, Scheduling

[Task-Scheduling: Cormen et al, Kap. 27] [Concurrency, Scheduling: Williams. Kap. 1.1 – 1.2]

Beweis: Anzahl nichtsättigende Push-Operationen

Beweis:

Sei $A_f := \{ v \in V : \alpha_f(v) > 0 \}$

Wahl von u bei push: $u \in A_f$ mit $h(u) \ge h(v)$ für alle $v \in A_f$.

Bei nichstättigender push-Operation verschwindet u von A_f . Bei diesem und jedem weiteren Push kommen höchstens $v \in A_f$ hinzu mit h(v) < h(u). Bis eine erneute Relabel-Operation stattgefunden hat, gilt also stets $u \not\in A_f$.

Dieses Argument gilt für jedes gewählte u und es können bis zur nächster Relabel-Operation höchstens n nichtsättigende Push-Operationen ausgeführt werden.

The Free Lunch

The free lunch is over 49

⁴⁹"The Free Lunch is Over", a fundamental turn toward concurrency in software, Herb Sutter, Dr. Dobb's Journal, 2005

Beobachtung von Gordon E. Moore: Die Anzahl Transistoren in integrierten Schaltkreisen verdoppelt sich ungefähr alle zwei Jahre.

Gordon E. Moore (1929)

Für eine lange Zeit...

- wurde die sequentielle Ausführung schneller ("Instruction Level Parallelism", "Pipelining", Höhere Frequenzen)
- mehr und kleinere Transistoren = mehr Performance
- Programmierer warteten auf die nächste schnellere Generation

Heute

5,000,000,000 1,000,000,000

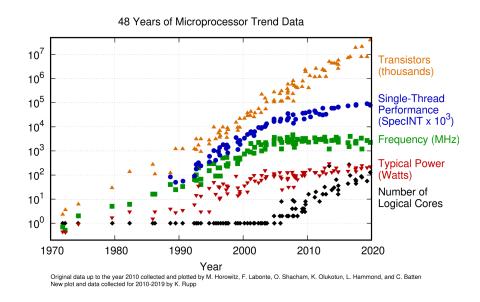
> 100,000,000 50,000,000

10.000.000 5,000,000

1,000,000

500,000

- steigt die Frequenz der Prozessoren kaum mehr an (Kühlproblematik)
- steigt die Instruction-Level Parallelität kaum mehr an
- ist die Ausführungsgeschwindigkeit in vielen Fällen dominiert von Speicherzugriffszeiten (Caches werden aber immer noch grösser und schneller)



Multicore

- Verwende die Transistoren für mehr Rechenkerne
- Parallelität in der Software
- Implikation: Programmierer müssen parallele Programme schreiben, um die neue Hardware vollständig ausnutzen zu können

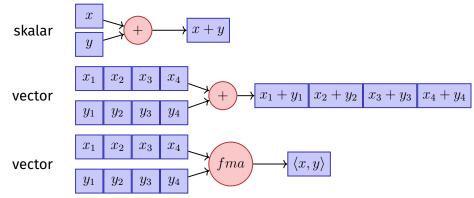
890

Formen der Parallelen Ausführung

- Vektorisierung
- Pipelining
- Instruction Level Parallelism
- Multicore / Multiprocessing
- Verteiltes Rechnen

Vektorisierung

Parallele Ausführung derselben Operation auf Elementen eines Vektor(Register)s



891

Pipelining in CPUs

Fetch

Decode

Execute

Data Fetch

Writeback

Mehrere Stufen

- Jede Instruktion dauert 5 Zeiteinheiten (Zyklen)
- Im besten Fall: 1 Instruktion pro Zyklus, nicht immer möglich ("stalls")

Parallelität (mehrere funktionale Einheiten) führt zu **schnellerer Ausführung.**

ILP - Instruction Level Parallelism

Moderne CPUs führen unabhängige Instruktionen intern auf mehreren Einheiten parallel aus

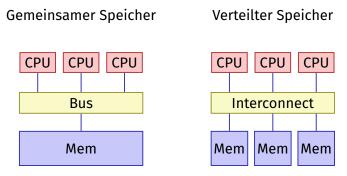
- Pipelining
- Superskalare CPUs (Mehrere Instruktionen pro Zyklus)
- Out-Of-Order Execution (Programmierer sieht die sequentielle Ausführung)
- Speculative Execution (Instruktionen werden spekulativ ausgeführt und unterbrochen, wenn die Bedingung zu deren Ausführung nicht erfüllt ist.)

894

89

Gemeinsamer vs. verteilter Speicher

28.2 Hardware Architekturen



Architekturen mit gemeinsamen Speicher

- Multicore (Chip Multiprocessor CMP)
- Symmetric Multiprocessor Systems (SMP)
- Non-Uniform Memory Access (NUMA)
- Simultaneous Multithreading (SMT = Hyperthreading)
 - nur ein physischer Kern, Mehrere Instruktionsströme/Threads: mehrere virtuelle Kerne
 - Zwischen ILP (mehrere Units für einen Strom) und Multicore (mehrere Units für mehrere Ströme). Limitierte parallele Performance

Gleiches Programmierinterface!

Shared vs. Distributed Memory Programming

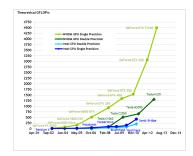
- Kategorien des Programmierinterfaces
 - Kommunikation via Message Passing
 - Kommunikation via geteiltem Speicher
- Es ist möglich:
 - Systeme mit gemeinsamen Speicher als verteilte Systeme zu programmieren (z.B. mit Message Passing Interface MPI)
 - Systeme mit verteiltem Speicher als System mit gemeinsamen Speicher zu programmieren (z.B. Partitioned Global Address Space PGAS)

8

Massiv Parallele Hardware

[General Purpose] Graphical Processing Units ([GP]GPUs)

- Revolution im High Performance Computing
 - Calculation 4.5 TFlops vs. 500 GFlops
 - Memory Bandwidth 170 GB/s vs. 40 GB/s
- Single Instruction Multiple Data (SIMD)
 - Hohe Datenparallelität
 - Benötigt eigenes Programmiermodell.Z.B. CUDA / OpenCL



28.3 Multi-Threading, Parallelität und Nebenläufigkeit

Prozesse und Threads

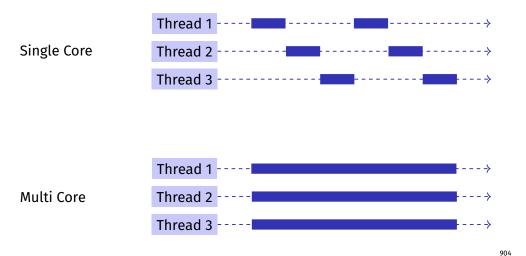
- Prozess: Instanz eines Programmes
 - jeder Prozess hat seinen eigenen Kontext, sogar eigenen Addresraum
 - OS verwaltet Prozesse (Resourcenkontrolle, Scheduling, Synchronisierung)
- Threads: Ausführungsfäden eines Programmes
 - Threads teilen sich einen Addressraum
 - Schneller Kontextwechsel zwischen Threads

Warum Multithreading?

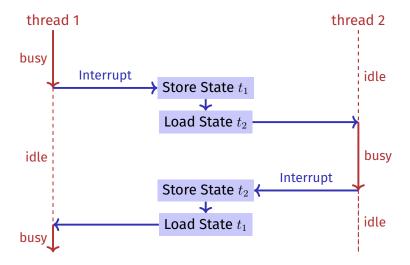
- Verhinderung vom "Polling" auf Resourcen (Files, Netwerkzugriff, Tastatur)
- Interaktivität (z.B. Responsivität von GUI Programmen)
- Mehrere Applikationen / Clients gleichzeitig instanzierbar
- Parallelität (Performanz!)

02

Multithreading konzeptuell



Threadwechsel auf einem Core (Preemption)



Parallelität vs. Nebenläufigkeit (Concurrency)

- Parallelität: Verwende zusätzliche Resourcen (z.B. CPUs), um ein Problem schneller zu lösen
- **Nebenläufigkeit:** Vewalte gemeinsam genutzte Resourcen (z.B. Speicher) korrekt und effizient
- Begriffe überlappen offensichtlich. Bei parallelen Berechnungen besteht fast immer Synchronisierungsbedarf.

Parallelität Nebenläufigkeit Arbeit Anfragen Resourcen Resourcen

Thread-Sicherheit

Thread-Sicherheit bedeutet, dass in der nebenläufigen Anwendung eines Programmes dieses sich immer wie gefordert verhält.

Viele Optimierungen (Hardware, Compiler) sind darauf ausgerichtet, dass sich ein sequentielles Programm korrekt verhält.

Nebenläufige Programme benötigen für ihre Synchronisierungen auch eine Annotation, welche gewisse Optimierungen selektiv abschaltet

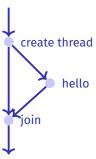
906

C++11 Threads

```
#include <iostream>
#include <thread>

void hello(){
   std::cout << "hello\n";
}

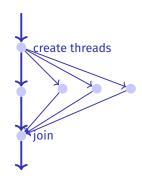
int main(){
   // create and launch thread t
   std::thread t(hello);
   // wait for termination of t
   t.join();
   return 0;
}</pre>
```



28.4 C++ Threads

C++11 Threads

```
void hello(int id){
  std::cout << "hello from " << id << "\n";
}
int main(){
  std::vector<std::thread> tv(3);
  int id = 0;
  for (auto & t:tv)
    t = std::thread(hello, ++id);
  std::cout << "hello from main \n";
  for (auto & t:tv)
    t.join();
  return 0;
}</pre>
```



Nichtdeterministische Ausführung!

Eine Ausführung:

hello from main hello from 2 hello from 1 hello from 0

Andere Ausführung:

hello from 1 hello from main hello from 0 hello from 2

Andere Ausführung:

hello from main hello from 0 hello from hello from 1

911

10

Technisches Detail

void background();

Um einen Thread als Hintergrundthread weiterlaufen zu lassen:

```
void someFunction(){
    ...
    std::thread t(background);
    t.detach();
    ...
} // no problem here, thread is detached
```

Mehr Technische Details

- Beim Erstellen von Threads werden auch Referenzparameter kopiert, ausser man gibt explizit std::ref bei der Konstruktion an.
- Funktoren oder Lambda-Expressions können auch auf einem Thread ausgeführt werden
- In einem Kontext mit Exceptions sollte das join auf einem Thread im catch-Block ausgeführt werden

Noch mehr Hintergründe im Kapitel 2 des Buches *C++ Concurrency in Action*, Anthony Williams, Manning 2012. Auch online bei der ETH Bibliothek erhältlich.

28.5 Task- und Datenparallelität

Paradigmen der Parallelen Programmierung

Arbeitsteilung: Aufteilung der Arbeit in parallele Tasks

- Task / Thread Parallel: Programmierer legt parallele Tasks manuell / explizit fest.
- **Daten-Parallel:** Dieselben Operationen finden auf einer Menge von Datenobjekten parallel statt. Programmierer legt die Operationen fest und das System erledigt den Rest.

914

Beispiel Data Parallel (OMP)

```
double sum = 0, A[MAX];
#pragma omp parallel for reduction (+:ave)
for (int i = 0; i < MAX; ++i)
   sum += A[i];
return sum;</pre>
```

Beispiel Task Parallel (C++11 Threads/Futures)

```
double sum(Iterator from, Iterator to)
{
  auto len = from - to;
  if (len > threshold){
   auto future = std::async(sum, from, from + len / 2);
   return sumS(from + len / 2, to) + future.get();
  }
  else
   return sumS(from, to);
}
```

Partitionierung und Scheduling

- Aufteilung der Arbeit in parallele Tasks (Programmierer oder System)
 - Ein Task ist eine Arbeitseinheit
 - Frage: welche Granularität?
- Scheduling (Laufzeitsystem)
 - Zuweisung der Tasks zu Prozessoren
 - Ziel: volle Resourcennutzung bei wenig Overhead

Granularität: Wie viele Tasks?

- #Tasks = #Cores?
- Problem: wenn ein Core nicht voll ausgelastet werden kann
- Beispiel: 9 Einheiten Arbeit. 3 Cores. Scheduling von 3 sequentiellen Tasks.



Exklusive Auslastung:

P1	s 1
P2	s2
P3	s3

Ausführungszeit: 3 Einheiten

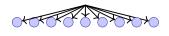
Fremder Thread "stört":

P1	s 1		
P2		s2	s1
Р3		s3	

Ausführungszeit: 5 Einheiten

Granularität: Wie viele Tasks?

- #Tasks = Maximum?
- Beispiel: 9 Einheiten Arbeit. 3 Cores. Scheduling von 9 sequentiellen Tasks.



Exklusive Auslastung:

P1	s1	S 4	s7
P2	s2	s5	s8
Р3	s3	s6	s9

Ausführungszeit: $3 + \varepsilon$ Einheiten

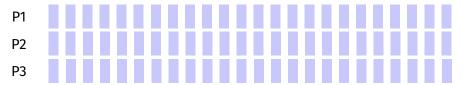
Fremder Thread "stört":

P1	s1				
P2	s2	S 4	s5	s8	
P3	s3	s6	s7	s9	

Ausführungszeit: 4 Einheiten. Volle Auslastung.

Granularität: Wie viele Tasks?

- #Tasks = Maximum?
- Beispiel: 10⁶ kleine Einheiten Arbeit.



Ausführungszeit: dominiert vom Overhead

Granularität: Wie viele Tasks?

Antwort: so viele Tasks wie möglich mit sequentiellem Cut-off, welcher den Overhead vernachlässigen lässt.

28.6 Skalierbarkeit: Amdahl und Gustafson

922

Skalierbarkeit

In der parallelen Programmierung:

- lacktriangle Geschwindigkeitssteigerung bei wachsender Anzahl p Prozessoren
- Was passiert, wenn $p \to \infty$?
- Linear skalierendes Programm: Linearer Speedup

Parallele Performanz

Gegeben fixierte Rechenarbeit W (Anzahl Rechenschritte)

 T_1 : Sequentielle Ausführungszeit sei

 T_p : Parallele Ausführungszeit auf p CPUs

■ Perfektion: $T_p = T_1/p$

■ Performanzverlust: $T_p > T_1/p$ (üblicher Fall)

■ Hexerei: $T_p < T_1/p$

Paralleler Speedup

Paralleler Speedup S_p auf p CPUs:

$$S_p = \frac{W/T_p}{W/T_1} = \frac{T_1}{T_p}.$$

- Perfektion: Linearer Speedup $S_p = p$
- Verlust: sublinearer Speedup $S_p < p$ (der übliche Fall)
- Hexerei: superlinearer Speedup $S_p > p$

Effizienz: $E_p = S_p/p$

Erreichbarer Speedup?

Paralleles Programm

Paralleler Teil Seq. Teil 80% 20%

$$T_1 = 10$$

$$T_8 = \frac{10 \cdot 0.8}{8} + 10 \cdot 0.2 = 1 + 2 = 3$$

$$S_8 = \frac{T_1}{T_8} = \frac{10}{3} \approx 3.3 < 8 \quad (!)$$

Amdahl's Law: Zutaten

Zu leistende Rechenarbeit W fällt in zwei Kategorien

- \blacksquare Parallelisierbarer Teil W_p
- \blacksquare Nicht paralleliserbarer, sequentieller Teil W_s

Annahme: W kann mit **einem** Prozessor in W Zeiteinheiten sequentiell erledigt werden $(T_1 = W)$:

$$T_1 = W_s + W_p$$
$$T_p \ge W_s + W_p/p$$

Amdahl's Law

$$S_p = \frac{T_1}{T_p} \le \frac{W_s + W_p}{W_s + \frac{W_p}{p}}$$

Amdahl's Law

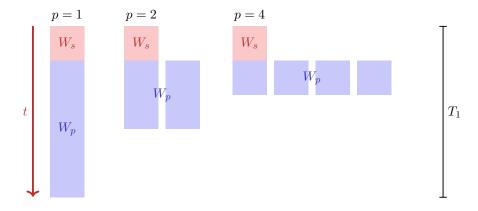
Mit seriellem, nicht parallelisierbaren Anteil λ : $W_s = \lambda W$, $W_p = (1 - \lambda)W$:

$$S_p \le \frac{1}{\lambda + \frac{1-\lambda}{p}}$$

Somit

$$S_{\infty} \le \frac{1}{\lambda}$$

Illustration Amdahl's Law



10

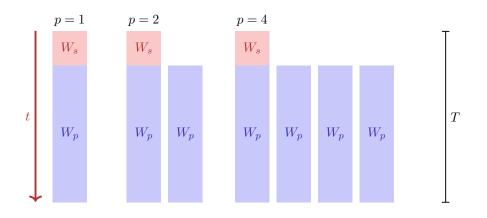
Amdahl's Law ist keine gute Nachricht

Alle nicht parallelisierbaren Teile können Problem bereiten und stehen der Skalierbarkeit entgegen.

Gustafson's Law

- Halte die Ausführungszeit fest.
- Variiere die Problemgrösse.
- Annahme: Der sequentielle Teil bleibt konstant, der parallele Teil wird grösser.

Illustration Gustafson's Law



Gustafson's Law

Arbeit, die mit einem Prozessor in der Zeit T erledigt werden kann:

$$W_s + W_p = T$$

Arbeit, die mit p Prozessoren in der Zeit T erledigt werden kann:

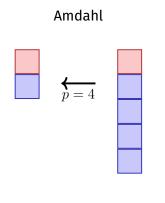
$$W_s + p \cdot W_p = \lambda \cdot T + p \cdot (1 - \lambda) \cdot T$$

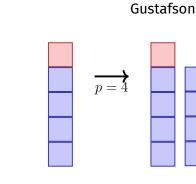
Speedup:

934

$$S_p = \frac{W_s + p \cdot W_p}{W_s + W_p} = p \cdot (1 - \lambda) + \lambda$$
$$= p - \lambda(p - 1)$$

Amdahl vs. Gustafson





Amdahl vs. Gustafson

Die Gesetze von Amdahl und Gustafson sind Modelle der Laufzeitverbesserung bei Parallelisierung.

Amdahl geht von einem festen **relativen** sequentiellen Anteil der Arbeit aus, während Gustafson von einem festen **absoluten** sequentiellen Teil ausgeht (der als Bruchteil der Arbeit W_1 ausgedrückt wird und bei Zunahme der Arbeit nicht wächst).

Die beiden Modelle widersprechen sich nicht, sondern beschreiben die Laufzeitverbesserung verschiedener Probleme und Algorithmen.

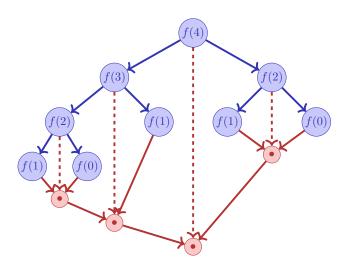
28.7 Scheduling

Beispiel: Fibonacci

```
int fib_task(int x){
   if (x < 2) {
      return x;
   } else {
      auto f1 = std::async(fib_task, x-1);
      auto f2 = std::async(fib_task, x-2);
      return f1.get() + f2.get();
   }
}</pre>
```

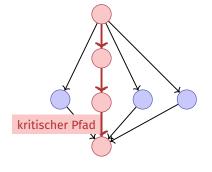
938

Task-Graph



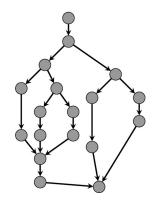
Frage

- Jeder Knoten (Task) benötigt 1 Zeiteinheit.
- Pfeile bezeichnen Abhängigkeiten.
- Minimale Ausführungseinheit wenn Anzahl Prozessoren = ∞ ?



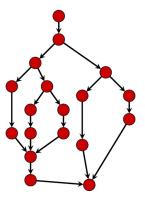
Performanzmodell

- p Prozessoren
- Dynamische Zuteilung
- \blacksquare T_p : Ausführungszeit auf p Prozessoren



Performanzmodell

- T_p : Ausführungszeit auf p Prozessoren
- *T*₁: **Arbeit:** Zeit für die gesamte Berechnung auf einem Prozessor
- T_1/T_p : Speedup

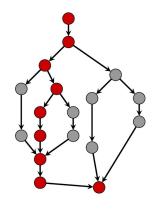


942

Performanzmodell

- T_{∞} : **Zeitspanne**: Kritischer Pfad. Ausführungszeit auf ∞ Prozessoren. Längster Pfad von der Wurzel zur Senke.
- T_1/T_∞ : Parallelität: breiter ist besser
- Untere Grenzen

 $T_p \geq T_1/p$ Arbeitsgesetz $T_p \geq T_\infty$ Zeitspannengesetz



Greedy Scheduler

Greedy Scheduler: teilt zu jeder Zeit so viele Tasks zu Prozessoren zu wie möglich.

Theorem 43

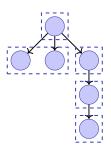
Auf einem idealen Parallelrechner mit p Prozessoren führt ein Greedy-Scheduler eine mehrfädige Berechnung mit Arbeit T_1 und Zeitspanne T_∞ in Zeit

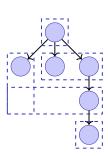
$$T_p \le T_1/p + T_\infty$$

aus.

Beispiel

Annahme p=2.





$$T_p = 5$$

$$T_p = 4$$

Beweis des Theorems

Annahme, dass alle Tasks gleich viel Arbeit aufweisen.

- lacktriangle Vollständiger Schritt: p Tasks stehen zur Berechnung bereit
- lacktriangle Unvollständiger Schritt: weniger als p Tasks bereit.

Annahme: Anzahl vollständige Schritte grösser als $\lfloor T_1/p \rfloor$. Ausgeführte Arbeit $\geq \lfloor T_1/p \rfloor \cdot p + p = T_1 - T_1 \mod p + p > T_1$. Widerspruch. Also maximal $\lfloor T_1/p \rfloor$ vollständige Schritte.

Betrachten nun den Graphen der ausstehenden Tasks. Jeder maximale (kritische) Pfad beginnt mit einem Knoten t mit $\deg^-(t)=0$. Jeder unvollständige Schritt führt zu jedem Zeitpunkt alle vorhandenen Tasks t mit $\deg^-(t)=0$ aus und verringert also die Länge der Zeitspanne. Anzahl unvollständige Schritte also begrenzt durch T_∞ .

946

Konsequenz

Wenn $p \ll T_1/T_{\infty}$, also $T_{\infty} \ll T_1/p$, dann

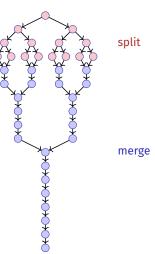
$$T_p \le T_1/p + T_\infty \quad \Rightarrow \quad T_p \lesssim T_1/p$$

Fibonacci

 $T_1(n)/T_\infty(n)=\Theta(\phi^n/n)$. Für moderate Grössen von n können schon viele Prozessoren mit linearem Speedup eingesetzt werden.

Beispiel: Parallelität von Mergesort

- Arbeit (sequentielle Laufzeit) von Mergesort $T_1(n) = \Theta(n \log n)$.
- Span $T_{\infty}(n) = \Theta(n)$
- Parallelität $T_1(n)/T_\infty(n) = \Theta(\log n)$ (Maximal erreichbarer Speedup mit $p = \infty$ Prozessoren)



29. Parallel Programming II

Gemeinsamer Speicher, Nebenläufigkeit, Gegenseitiger Ausschluss, Race Conditions [C++ Threads: Williams, Kap. 2.1-2.2], [C++ Race Conditions: Williams, Kap. 3.1] [C++ Mutexes: Williams, Kap. 3.2.1, 3.3.3]

29.1 Gemeinsamer Speicher, Nebenläufigkeit

950

Gemeinsam genutzte Resourcen (Speicher)

- Bis hier: fork-join Algorithmen: Datenparallel oder Divide und Conquer
- Einfache Struktur (Datenunabhängigkeit der Threads) zum Vermeiden von Wettlaufsituationen (race conditions)
- Funktioniert nicht mehr, wenn Threads gemeinsamen Speicher nutzen müssen.

Konsistenz des Zustands

Gemeinsamer Zustand: Hauptschwierigkeit beim nebenläufigen Programmieren.

Ansätze:

- Unveränderbarkeit, z.B. Konstanten
- Isolierte Veränderlichkeit, z.B. Thread-lokale Variablen, Stack.
- Gemeinsame veränderliche Daten, z.B. Referenzen auf gemeinsamen Speicher, globale Variablen

Schütze den gemeinsamen Zustand

- Methode 1: Locks, Garantiere exklusiven Zugriff auf gemeinsame Daten.
- Methode 2: lock-freie Datenstrukturen, garantiert exklusiven Zugriff mit sehr viel feinerer Granularität.
- Methode 3: Transaktionsspeicher (hier nicht behandelt)

Kanonisches Beispiel

```
class BankAccount {
  int balance = 0;
public:
  int getBalance(){ return balance; }
  void setBalance(int x) { balance = x; }
  void withdraw(int amount) {
    int b = getBalance();
    setBalance(b - amount);
  }
  // deposit etc.
};
```

(korrekt bei Einzelthreadausführung)

·

Ungünstige Verschachtelung (Bad Interleaving)

Paralleler Aufruf von withdraw (100) auf demselben Konto

Verlockende Fallen

FALSCH:

```
void withdraw(int amount) {
  int b = getBalance();
  if (b==getBalance())
    setBalance(b - amount);
}
```

Bad interleavings lassen sich fast nie mit wiederholtem Lesen lösen

Verlockende Fallen

Auch FALSCH:

```
void withdraw(int amount) {
    setBalance(getBalance() - amount);
}
```

Annahmen über Atomizität von Operationen sind fast immer falsch

Gegenseitiger Ausschluss (Mutual Exclusion)

Wir benötigen ein Konzept für den gegenseitigen Ausschluss Nur ein Thread darf zu einer Zeit die Operation withdraw auf demselben Konto ausführen.

Der Programmierer muss den gegenseitigen Ausschlus sicherstellen.

958

Mehr verlockende Fallen

```
class BankAccount {
  int balance = 0;
  bool busy = false;
public:
  void withdraw(int amount) {
    while (busy); // spin wait
    busy = true;
    int b = getBalance();
    setBalance(b - amount);
    busy = false;
}

// deposit would spin on the same boolean
};
```

Das Problem nur verschoben!

Wie macht man das richtig?

- Wir benutzen ein **Lock** (eine Mutex) aus Bibliotheken
- Eine Mutex verwendet ihrerseits Hardwareprimitiven, Read-Modify-Write (RMW) Operationen, welche atomar lesen und abhängig vom Leseergebis schreiben können.
- Ohne RMW Operationen ist der Algorithmus nichttrivial und benötigt zumindest atomaren Zugriff auf Variablen von primitivem Typ.

29.2 Gegenseitiger Ausschluss

962

Kritische Abschnitte und Gegenseitiger Ausschluss

Kritischer Abschnitt (Critical Section)

Codestück, welches nur durch einen einzigen Thread zu einer Zeit ausgeführt werden darf.

Gegenseitiger Ausschluss (Mutual Exclusion)

Algorithmus zur Implementation eines kritischen Abschnitts

```
acquire_mutex();  // entry algorithm\\
...  // critical section
release_mutex();  // exit algorithm
```

Anforderung an eine Mutex.

Korrektheit (Safety)

Maximal ein Prozess in der kritischen Region



Fortschritt (Liveness)

Das Betreten der kritischen Region darf nur endliche Zeit dauern, wenn kein Thread in der kritischen Region verweilt.



Fast Korrekt

```
class BankAccount {
  int balance = 0;
  std::mutex m; // requires #include <mutex>
public:
    ...
  void withdraw(int amount) {
    m.lock();
    int b = getBalance();
    setBalance(b - amount);
    m.unlock();
  }
};
```

Was, wenn eine Exception auftritt?

Reentrante Locks

Reentrantes Lock (rekursives Lock)

- merkt sich den betroffenen Thread;
- hat einen Zähler
 - Aufruf von lock: Zähler wird inkrementiert
 - Aufruf von unlock: Z\u00e4hler wird dekrementiert. Wenn Z\u00e4hler = 0, wird das Lock freigeben

RAII Ansatz

```
class BankAccount {
  int balance = 0;
  std::mutex m;
public:
    ...
  void withdraw(int amount) {
    std::lock_guard<std::mutex> guard(m);
    int b = getBalance();
    setBalance(b - amount);
  } // Destruction of guard leads to unlocking m
};
```

Was ist mit getBalance / setBalance?

966

Konto mit reentrantem Lock

```
class BankAccount {
  int balance = 0;
  std::recursive_mutex m;
  using guard = std::lock_guard<std::recursive_mutex>;
public:
  int getBalance(){ guard g(m); return balance;
  }
  void setBalance(int x) { guard g(m); balance = x;
  }
  void withdraw(int amount) { guard g(m);
   int b = getBalance();
   setBalance(b - amount);
  }
};
```

29.3 Race Conditions

Wettlaufsituation (Race Condition)

- Eine **Wettlaufsituation** (Race Condition) tritt auf, wenn das Resultat einer Berechnung vom Scheduling abhängt.
- Wir unterscheiden bad interleavings und data races
- Bad Interleavings können auch unter Verwendung einer Mutex noch auftreten.

770

Beispiel: Stack

Stack mit korrekt synchronisiertem Zugriff:

```
template <typename T>
class stack{
    ...
    std::recursive_mutex m;
    using guard = std::lock_guard<std::recursive_mutex>;
public:
    bool isEmpty(){ guard g(m); ... }
    void push(T value){ guard g(m); ... }
    T pop(){ guard g(m); ...}
};
```

Peek

Peek Funktion vergessen. Dann so?

```
template <typename T>
T peek (stack<T> &s){
  T value = s.pop();
  s.push(value);
  return value;
}
```

Code trotz fragwürdigem Stil in sequentieller Welt korrekt. Nicht so in nebenläufiger Programmierung!

Bad Interleaving!

Initial leerer Stack s, nur von Threads 1 und 2 gemeinsam genutzt. Thread 1 legt einen Wert auf den Stack und prüft, dass der Stack nichtleer ist. Thread 2 liest mit peek() den obersten Wert.

Die Lösung

Peek muss mit demselben Lock geschützt werden, wie die anderen Zugriffsmethoden.

974

Bad Interleavings

Race Conditions in Form eines Bad Interleavings können also auch auf hoher Abstraktionsstufe noch auftreten.

Betrachten nachfolgend andere Form der Wettlaufsitation: Data Race.

Wie ist es damit?

```
class counter{
  int count = 0;
  std::recursive_mutex m;
  using guard = std::lock_guard<std::recursive_mutex>;
public:
  int increase(){
    return ++count;
  }
  int get(){
    return count;
  }
}
```

Warum falsch?

Es sieht so aus, als könne hier nichts schiefgehen, da der Update von count in einem "winzigen Schritt" geschieht.

Der Code ist trotzdem falsch und von Implementationsdetails der Programmiersprache und unterliegenden Hardware abhängig.

Das vorliegende Problem nennt man **Data-Race**

Moral: Vermeide Data-Races, selbst wenn jede denkbare Form von Verschachtelung richtig aussieht. Mache keine Annahmen über die Anordnung von Speicheroperationen.

Ftwas formaler

Data Race (low-level Race-Conditions) Fehlerhaftes Programmverhalten verursacht durch ungenügend synchronisierten Zugriff zu einer gemeinsam genutzten Resource, z.B. gleichzeitiges Lesen/Schreiben oder Schreiben/Schreiben zum gleichen Speicherbereich.

Bad Interleaving (High Level Race Condition) Fehlerhaftes Programmverhalten verursacht durch eine unglückliche Ausführungsreihenfolge eines Algorithmus mit mehreren Threads, selbst dann wenn die gemeinsam genutzten Resourcen anderweitig gut synchronisiert sind.

78

Genau hingeschaut

Es gibt keine Verschachtelung zweier f und g aufrufender Threads die die Bedingung in der Assertion falsch macht:

- ABCD ✓
- ACBD ✓
- ACDB ✓
- CABD ✓
- CCDB√
- CDAB ✓

Es kann trotzdem schiefgehen!

Ein Grund: Memory Reordering

Daumenregel: Compiler und Hardware dürfen die Ausführung des Codes so ändern, dass die *Semantik einer sequentiellen Ausführung* nicht geändert wird.

Die Software-Perspektive

Moderne Compiler geben keine Garantie, dass die globale Anordnung aller Speicherzugriffe der Ordnung im Quellcode entsprechen

- Manche Speicherzugriffe werden sogar komplett wegoptimiert
- Grosses Potential für Optimierungen und Fehler in der nebenläufigen Welt, wenn man falsche Annahmen macht

Beispiel: Selbstgemachtes Rendevouz

```
int x; // shared

void wait(){
    x = 1;
    while(x == 1);
}

void arrive(){
    x = 2;
}
```

Angenommen Thread 1 ruft wait auf, später ruft Thread 2 arrive auf. Was passiert?

```
thread 1 wait hread 2 arrive
```

982

983

Kompilation

Source int x; // shared void wait(){ x = 1; while(x == 1); } void arrive(){ x = 2; } void arrive(){ arri movl

```
Ohne Optimierung
                        Mit Optimierung
wait:
                        wait:
movl $0x1, x
                        movl $0x1, x
test: ←
mov x, %eax
                        jmp test
                 if equa
cmp $0x1, %eax
ie test -
arrive:
                        arrive
movl $0x2, x
                        movl $0x2, x
```

Hardware Perspektive

Moderne Prozessoren erzwingen nicht die globale Anordnung aller Instruktionen aus Gründen der Performanz:

- Die meisten Prozessoren haben einen Pipeline und können Teile von Instruktionen simultan oder in anderer Reihenfolge ausführen.
- Jeder Prozessor(kern) hat einen lokalen Cache, der Effekt des Speicherns im gemeinsamen Speicher kann bei anderen Prozessoren u.U. erst zu einem späteren Zeitpunkt sichtbar werden.

Speicherhierarchie

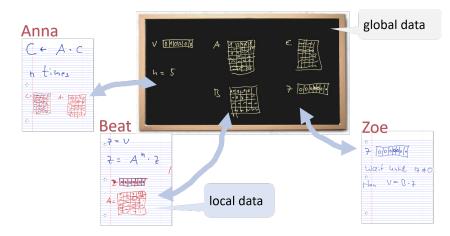
Registers L1 Cache L2 Cache

System Memory

schnell,kleine Latenz,hohe Kosten, geringe Kapazität

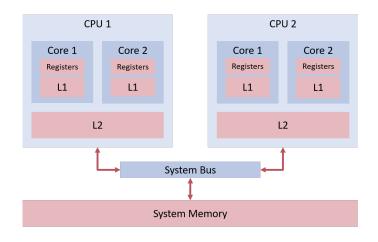
langsam, hohe Latenz, geringe Kosten, hohe Kapazität

Eine Analogie



986

Schematisch



Speichermodelle

Wann und ob Effekte von Speicheroperationen für Threads sichtbar werden, hängt also von Hardware, Laufzeitsystem und der Programmiersprache ab. Ein Speichermodell (z.B. das von C++) gibt Minimalgarantien für den Effekt von Speicheroperationen.

- Lässt Möglichkeiten zur Optimierung offen
- Enthält Anleitungen zum Schreiben Thread-sicherer Programme

C++ gibt zum Beispiel Garantien, wenn Synchronisation mit einer Mutex verwendet wird.

Repariert

```
class C {
   int x = 0;
   int y = 0;
   std::mutex m;
public:
   void f() {
      m.lock(); x = 1; m.unlock();
      m.lock(); y = 1; m.unlock();
   }
   void g() {
      m.lock(); int a = y; m.unlock();
      m.lock(); int b = x; m.unlock();
      assert(b >= a); // cannot fail
   }
};
```

Atomic

```
Hier auch möglich:
```

```
class C {
    std::atomic_int x{0}; // requires #include <atomic>
    std::atomic_int y{0};
public:
    void f() {
        x = 1;
        y = 1;
    }
    void g() {
        int a = y;
        int b = x;
        assert(b >= a); // cannot fail
    }
};
```

30. Parallel Programming III

Verklemmung (Deadlock) und Verhungern (Starvation), Producer-Consumer, Konzept des Monitors, Condition Variables [Deadlocks: Williams, Kap. 3.2.4-3.2.5] [Condition Variables: Williams, Kap. 4.1]

Verklemmung (Deadlock) Motivation

```
class BankAccount {
  int balance = 0;
  std::recursive_mutex m;
  using guard = std::lock_guard<std::recursive_mutex>;
public:
    ...
  void withdraw(int amount) { guard g(m); ... }
  void deposit(int amount){ guard g(m); ... }

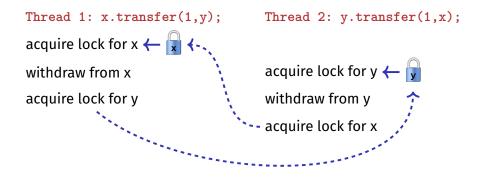
  void transfer(int amount, BankAccount& to){
      guard g(m);
      withdraw(amount);
      to.deposit(amount);
  }
};
```

992

990

Verklemmung (Deadlock) Motivation

Betrachte BankAccount Instanzen x und y



Deadlock

Deadlock: zwei oder mehr Prozesse sind gegenseitig blockiert, weil jeder Prozess auf einen anderen Prozess warten muss, um fortzufahren.



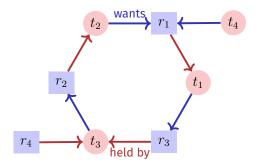
4

Threads und Resourcen

Grafisch: Threads t und Resourcen (Locks) r
 Thread t versucht Resource a zu bekommen: t a
 Resource b wird von Thread q gehalten: b

Deadlock - Erkennung

Ein Deadlock für Threads t_1, \ldots, t_n tritt auf, wenn der gerichtete Graph, der die Beziehung der n threads und Resourcen r_1, \ldots, r_m beschreibt, einen Kreis enthält.



Techniken

- **Deadlock Erkennung** findet die Zyklen im Abhängigkeitsgraph. Deadlock kann normalerweise nicht geheilt werden: Freigeben von Locks resultiert in inkonsistentem Zustand.
- **Deadlock Vermeidung** impliziert, dass Zyklen nicht auftreten können
 - Grobere Granularität "one lock for all"
 - Zwei-Phasen-Locking mit Retry-Mechanismus
 - Lock-Hierarchien
 - **...**
 - Anordnen der Resourcen

998

C++11 Stil

```
class BankAccount {
    ...
    std::recursive_mutex m;
    using guard = std::lock_guard<std::recursive_mutex>;
public:
    ...
    void transfer(int amount, BankAccount& to){
        std::lock(m,to.m); // lock order done by C++
        // tell the guards that the lock is already taken:
        guard g(m,std::adopt_lock); guard h(to.m,std::adopt_lock);
        withdraw(amount);
        to.deposit(amount);
}
```

Zurück zum Beispiel

```
class BankAccount {
  int id; // account number, also used for locking order
  std::recursive_mutex m; ...
public:
  ...
  void transfer(int amount, BankAccount& to){
    if (id < to.id){
       guard g(m); guard h(to.m);
       withdraw(amount); to.deposit(amount);
    } else {
       guard g(to.m); guard h(m);
       withdraw(amount); to.deposit(amount);
    }
}
};</pre>
```

Übrigens...

```
class BankAccount {
  int balance = 0;
  std::recursive_mutex m;
  using guard = std::lock_guard<std::recursive_mutex>;
public:
    ...
  void withdraw(int amount) { guard g(m); ... }
  void deposit(int amount) { guard g(m); ... }

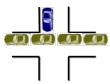
  void transfer(int amount, BankAccount& to) {
     withdraw(amount);
     to.deposit(amount);
   }
};

Das hätte auch funktioniert. Allerdings verschwindet dann kurz das Geld, was inakzeptabel ist (kurzzeitige Inkonsistenz!)
```

1000

Starvation und Livelock

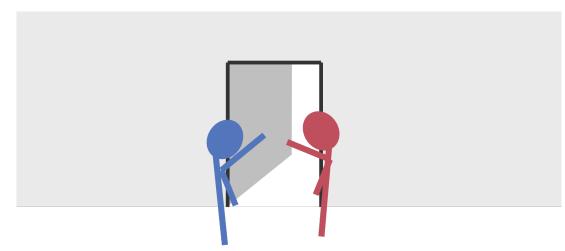
Starvation: der wiederholte, erfolglose Versuch eine zwischenzeitlich freigegebene Resource zu erhalten, um die Ausführung fortzusetzen.



Livelock: konkurrierende Prozesse erkennen einen potentiellen Deadlock, machen aber keinen Fortschritt beim Auflösen des Problems.



Politelock



1002

Produzenten-Konsumenten Problem

Zwei (oder mehr) Prozesse, Produzenten und Konsumenten von Daten, sollen mit Hilfe einer Datenstruktur entkoppelt werden. Fundamentale Datenstruktur für den Bau von Software-Pipelines!



Sequentielle Implementation (unbeschränkter Buffer)

```
class BufferS {
  std::queue<int> buf;
public:
    void put(int x) {
       buf.push(x);
    }
    int get() {
       while (buf.empty()) {} // wait until data arrive
       int x = buf.front();
       buf.pop();
       return x;
    }
};
```

Wie wärs damit?

```
class Buffer {
 std::recursive_mutex m;
 using guard = std::lock_guard<std::recursive_mutex>;
 std::queue<int> buf;
public:
   void put(int x){ guard g(m);
       buf.push(x);
   }
                               Deadlock
   int get(){ guard g(m);
       while (buf.empty()){}
       int x = buf.front();
       buf.pop();
       return x;
   }
};
```

Besser?

```
void put(int x){
  guard g(m);
  buf.push(x);
}
int get(){
  m.lock();
  while (buf.empty()){
    m.unlock();
    std::this_thread::sleep_for(std::chrono::milliseconds(10));
    m.lock();
}
int x = buf.front(); buf.pop();
  m.unlock();
return x;
```

Ok, so?

```
void put(int x){
   guard g(m);
   buf.push(x);
int get(){
   m.lock();
   while (buf.empty()){
                           Ok, das geht, verschwendet aber CPU
       m.unlock();
                           Zeit!
       m.lock();
   }
   int x = buf.front();
   buf.pop();
   m.unlock();
   return x;
}
```

1007

Moral

1006

Wir wollen das Warten auf eine Bedingung nicht selbst implementieren müssen.

Dafür gibt es bereits einen Mechanismus: **Bedingungsvariablen (condition variables)**.

Das zugrunde liegende Konzept nennt man Monitor.

Monitor

Monitor Abstrakte Datenstruktur, die mit einer Menge Operationen ausgestattet ist, die im gegenseitigen Ausschluss arbeiten und synchronisiert werden können.

Erfunden von C.A.R. Hoare und Per Brinch Hansen (cf. Monitors – An Operating System Structuring Concept, C.A.R. Hoare 1974)

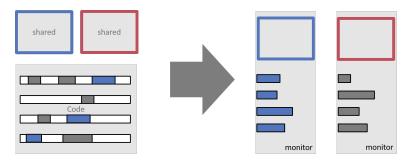


C.A.R. Hoare, *1934



Per Brinch Hansen (1938-2007)

Monitors vs. Locks



1011

1010

Monitor und Bedingungen

Ein Monitor stellt, zusätzlich zum gegenseitigen Ausschluss, folgenden Mechanismus bereit:

Warten auf Bedingungen: Ist eine Bedingung nicht erfüllt, dann

- Gib das Lock auf
- Warte auf die Erfüllung der Bedingung
- Prüfe die Erfüllung der Bedingung wenn ein Signal gesendet wird

Signalisieren: Thread, der die Bedingung wahr machen könnte:

Sende Signal zu potentiell wartenden Threads

Bedingungsvariablen

```
#include <mutex>
#include <condition_variable>
...

class Buffer {
   std::queue<int> buf;

   std::mutex m;
   // need unique_lock guard for conditions
   using guard = std::unique_lock<std::mutex>;
   std::condition_variable cond;
public:
   ...
};
```

Bedingungsvariablen

```
class Buffer {
...
public:
    void put(int x){
        guard g(m);
        buf.push(x);
        cond.notify_one();
}
    int get(){
        guard g(m);
        cond.wait(g, [&]{return !buf.empty();});
        int x = buf.front(); buf.pop();
        return x;
}
};
```

Technische Details

- Ein Thread, der mit cond.wait wartet, läuft höchstens sehr kurz auf einem Core. Danach belastet er das System nicht mehr und "schläft".
- Der Notify (oder Signal-) Mechanismus weckt schlafende Threads auf, welche nachfolgend ihre Bedingung prüfen.
 - cond.notify_one signalisiert einen wartenden Threads.
 - cond.notify_all signalisiert *alle* wartende Threads. Benötigt, wenn wartende Threads potentiell auf *verschiedene* Bedingungen warten.

1015

Technische Details

In vielen anderen Sprachen gibt es denselben Mechanismus. Das Prüfen von Bedingungen (in einem Loop!) muss der Programmierer dort oft noch selbst implementieren.

Java Beispiel

```
synchronized long get() {
    long x;
    while (isEmpty())
        try {
            wait ();
        } catch (InterruptedException e)
        x = doGet();
        return x;
}

synchronized put(long x){
        doPut(x);
        notify ();
}
```

1014

31. Parallel Programming IV

Futures, Read-Modify-Write Instruktionen, Atomare Variablen, Idee der lockfreien Programmierung

[C++ Futures: Williams, Kap. 4.2.1-4.2.3] [C++ Atomic: Williams, Kap. 5.2.1-5.2.4, 5.2.7] [C++ Lockfree: Williams, Kap. 7.1.-7.2.1]

Futures: Motivation

Threads waren bisher Funktionen ohne Resultat:

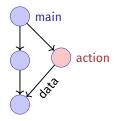
```
void action(some parameters){
    ...
}
std::thread t(action, parameters);
...
t.join();
// potentially read result written via ref-parameters
```

Futures: Motivation

Wir wollen nun etwas in dieser Art:

```
T action(some parameters){
    ...
    return value;
}

std::thread t(action, parameters);
    ...
value = get_value_from_thread();
```



1018

Wir können das schon!

- Wir verwenden das Producer/Consumer Pattern (implementiert mit Bedingungsvariablen)
- Starten einen Thread mit Referenz auf den Buffer
- Wenn wir das Resultat brauchen, holen wir es vom Buffer
- Synchronisation ist ja bereits implementiert

Zur Erinnerung

```
template <typename T>
class Buffer {
  std::queue<T> buf;
  std::mutex m;
  std::condition_variable cond;
public:
  void put(T x){ std::unique_lock<std::mutex> g(m);
    buf.push(x);
    cond.notify_one();
}
T get(){ std::unique_lock<std::mutex> g(m);
    cond.wait(g, [&]{return (!buf.empty());});
    T x = buf.front(); buf.pop(); return x;
}
}.
```

1020 };

Einfacher: nur ein einziger Wert

```
template <typename T>
class Buffer {
   T value; bool received = false;
   std::mutex m;
   std::condition_variable cond;
public:
   void put(T x){ std::unique_lock<std::mutex> g(m);
     value = x; received = true;
     cond.notify_one();
}
T get(){ std::unique_lock<std::mutex> g(m);
     cond.wait(g, [&]{return received;});
     return value;
}
};
```

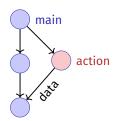
Anwendung

```
void action(Buffer<int>& c){
    // some long lasting operation ...
    c.put(42);
}
int main(){
    Buffer<int> c;
    std::thread t(action, std::ref(c));
    t.detach(); // no join required for free running thread
    // can do some more work here in parallel
    int val = c.get();
    // use result
    return 0;
}
```

1023

Mit C++11 Bordmitteln

```
int action(){
    // some long lasting operation
    return 42;
}
int main(){
    std::future<int> f = std::async(action);
    // can do some work here in parallel
    int val = f.get();
    // use result
    return 0;
}
```



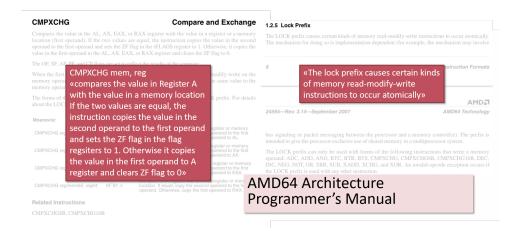
1022

Disclaimer

Die Darstellung oben ist vereinfacht. Die echte Implementation einer Future kennt Timeouts, Speicherallokatorenm kann mit Exceptions umgehen und ist näher am System geschrieben.

31.2 Read-Modify-Write

Beispiel: Atomare Operationen in Hardware



1026

Read-Modify-Write

Konzept von Read-Modify-Write: Der Effekt von Lesen, Verändern und Zurückschreiben, wird zu einem Zeitpunkt sichtbar (geschieht atomar).

Pseudo-Code für CAS – Compare-And-Swap

```
bool CAS(int& variable, int& expected, int desired){
   if (variable == expected){
     variable = desired;
     return true;
   }
   else{
     expected = variable;
     return false;
}
```

Verwendungsbeispiel CAS in C++11

Wir bauen unser eigenes (Spin-)Lock:

```
class Spinlock{
  std::atomic<bool> taken {false};
public:
  void lock(){
   bool old = false;
   while (!taken.compare_exchange_strong(old=false, true)){}
}
  void unlock(){
   bool old = true;
   assert(taken.compare_exchange_strong(old, false));
}
};
```

31.3 Lock-Freie Programmierung

Ideen

1030

Lock-freie Programmierung

Datenstruktur heisst

- **lock-frei**: zu jeder Zeit macht mindestens ein Thread in beschränkter Zeit Fortschritt, selbst dann, wenn viele Algorithmen nebenläufig ausgeführt werden. Impliziert systemweiten Fortschritt aber nicht Starvationfreiheit.
- wait-free: jeder Thread macht zu jeder Zeit in beschränkter Zeit Fortschritt, selbst dann wenn andere Algorithmen nebenläufig ausgeführt werden.

Fortschrittsbedingungen

	Lock-frei	Blockierend
Jeder macht Fortschritt	Wait-frei	Starvation-frei
Mindestens einer macht Fortschritt	Lock-frei	Deadlock-frei

Implikation

- Programmieren mit Locks: jeder Thread kann andere Threads beliebig blockieren.
- Lockfreie Programmierung: der Ausfall oder das Aufhängen eines Threads kann nicht bewirken, dass andere Threads blockiert werden

Wie funktioniert lock-freie Programmierung?

Beobachtung:

- RMW-Operationen sind in Hardware Wait-Free implementiert.
- Jeder Thread sieht das Resultat eines CAS in begrenzter Zeit.

Idee der lock-freien Programmierung: lese Zustand der Datenstruktur und verändere die Datenstruktur *atomar* dann und nur dann, wenn der gelesene Zustand unverändert bleibt.

1034

Beispiel: lock-freier Stack

Nachfolgend vereinfachte Variante eines Stacks

- pop prüft nicht, ob der Stack leer ist
- pop gibt nichts zurück

(Node)

```
Nodes:
struct Node {
   T value;

Node<T>* next;
   Node(T v, Node<T>* nxt): value(v), next(nxt) {}
};
```

value
next
value
next
value
value

next

value next 1035

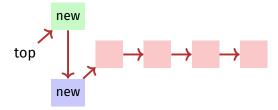
(Blockierende Version)

```
template <typename T>
class Stack {
                                                     top \rightarrow value
   Node<T> *top=nullptr;
                                                             next
   std::mutex m;
public:
                                                             value
   void push(T val){ guard g(m);
                                                             next
       top = new Node<T>(val, top);
   }
                                                             value
   void pop(){ guard g(m);
                                                             next
       Node<T>* old_top = top;
       top = top->next;
                                                             value
       delete old_top;
   }
                                                             next
};
```

Push

```
void push(T val){
  Node<T>* new_node = new Node<T> (val, top);
  while (!top.compare_exchange_weak(new_node->next, new_node));
}
```

2 Threads:



Lock-Frei

```
template <typename T>
class Stack {
   std::atomic<Node<T>**> top {nullptr};
public:
   void push(T val){
     Node<T>* new_node = new Node<T> (val, top);
     while (!top.compare_exchange_weak(new_node->next, new_node));
}
void pop(){
   Node<T>* old_top = top;
   while (!top.compare_exchange_weak(old_top, old_top->next));
   delete old_top;
}
};
```

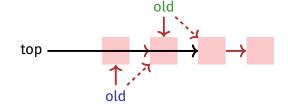
1038

1039

Pop

```
void pop(){
  Node<T>* old_top = top;
  while (!top.compare_exchange_weak(old_top, old_top->next));
  delete old_top;
}
```

2 Threads:



Lockfreie Programmierung – Grenzen

- Lockfreie Programmierung ist kompliziert.
- Wenn mehr als ein Wert nebenläufig angepasst werden muss (Beispiel: Queue), wird es schwieriger. Damit Algorithmen lock-frei bleiben, müssen Threads sich "gegenseitig helfen".
- Bei Speicherwiederverwendung kann das ABA Problem auftreten. Die Lösung dieses Problems ist aufwändig.