

Höhere Algorithmik

Mitschrift

Max Wisniewski

WS 2011/2012
25. Oktober 2011

Inhaltsverzeichnis

1	Einleitung	2
1.1	Ziel	2
1.2	Algorithmus	2
2	Repräsentant einer Menge	3
2.1	Naiver Ansatz	3
2.2	K - SELECT	3
2.2.1	Das Problem	4
2.2.2	Algorithmus I in $\Theta(n \cdot \log n)$	4
2.2.3	SELECT in $O(n)$	4
2.2.4	Implementierung von SPLITTER	5
2.3	Bemerkung	8

Kapitel 1

Einleitung

1.1 Ziel

Bisher haben wir einfache solcher Algorithmen betrachtet (ALP1, ALP3, etc.). In dieser Vorlesung werden wir uns nun mit komplexeren Problemen beschäftigen. Wir wollen die Probleme unter folgenden Aspekten betrachten:

- Entwurf von Algorithmen
- Analyse dieser Algorithmen
- Bewertung dieser Algorithmen

1.2 Algorithmus

Def.: Ein Algorithmus ist ein endlich beschriebenes, effektives Verfahren, das eine Eingabe in eine Ausgabe überführt.

Zu Beginn betrachten wir ein einfaches Problem.

Kapitel 2

Repräsentant einer Menge

Gegeben sei folgendes statistisches Problem:

Es seien n Zahlen / Datensätze gegeben, wobei $n \gg 0$ gilt.

Gesucht ist ein Repräsentativer Wert für diese Menge.

2.1 Naiver Ansatz

Idee Wir verwenden den Durchschnitt / Mittelwert.

Die Laufzeit ist einfach, da wir nur einmal über alle Datensätze müssen. Setzen wir dabei eine konstante Zeit für Addition und Division voraus, ist die Laufzeit $O(n)$.

Problem: Der Mittelwert ist Anfällig für Außreißer und daher nicht sehr aussagekräftig.

Sind beispielsweise $n - 1$ Werte zwischen 0 und 10 und ein n ter liegt bei 10.000.000 so wird das ganze Ergebnis zu diesem Wert hin verfälscht.

Dieser Repräsentant ist leicht zu berechnen, aber nicht sehr schön.
Betrachten wir daher einen anderen Ansatz.

2.2 K - SELECT

Def.: Ein Element s einer total geordneten Menge S hat den Rang k
: \Leftrightarrow es gibt genau $(k - 1)$ Elemente in S , die kleiner sind als s .

Man schreibt dafür $rg(s)$.

Def.: Sei S total geordnet mit $n = |S|$ und $s \in S$.

$$s \text{ heißt Median} :\Leftrightarrow rg(s) = \left\lceil \frac{n+1}{2} \right\rceil.$$

2.2.1 Das Problem

Gegeben Sei S , $|S| = n$ paarweise verschiedene Zahlen.

Nun wollen wir den Median s von S möglichst effizient finden.

2.2.2 Algorithmus I in $\Theta(n \cdot \log n)$

Was die Laufzeit schon nahe legt, bedienen wir uns hier eines Sortieralgorithmuses.

1. Sortiere S . z. B. mit Heap - Sort .
Benötigt $\Theta(n \cdot \log n)$ Schritte.
2. Gib das Element an der Stelle $\lceil \frac{n+1}{2} \rceil$ aus.
Benötigt $\Theta(1)$ Schritte.

Laufzeit: $T(n) = \Theta(n \cdot \log n) + \Theta(1) = \Theta(n \cdot \log n)$.

Da für (vergleichsbasiertes) Sortieren jede Lösung mit $\Omega(n \cdot \log n)$ beschränkt ist, kann eine Lösung für das Medianproblem die Sortierung verwendet nicht schneller sein. Bleibt zu untersuchen, ob der Median ähnlich schwer ist, oder ob es einen Algorithmus gibt, der das Problem schneller lösen kann.

2.2.3 SELECT in $O(n)$

Angenommen es existiert eine Funktion SPLITTER(S), welche uns ein Element $q \in S$ liefert, so dass gilt:

$$rg(q) \geq \left\lfloor \frac{1}{4} n \right\rfloor \quad \wedge \quad rg(q) \leq \left\lceil \frac{3}{4} n \right\rceil.$$

Lemma: Angenommen wir können SPLITTER ohne weitere Kosten benutzen. Dann können wir den Median in $O(n)$ Zeit berechnen.

Beweis: Um diese Aussage zu beweisen lösen wir das allgemeinere Problem

$$\text{SELECT}(k, S)$$

finde Element mit Rang k . Dieses Problem wird "Auswahlproblem" genannt.

Idee: Nehme SPLITTER als PIVOT Element und teile die Menge der Daten daran auf.

Pseudocode:

```

SELECT( k , S )
  IF |S| < 100 THEN
    RETURN BRUTFORCE( k , S )  // z. B. Algorithmus I
  q ← SPLITTER( S )
  S< ← { s ∈ S | s < q }
  S> ← { s ∈ S | s > q }
  IF |S<| ≥ k THEN
    RETURN SELECT( k , S< )
  ELSE IF |S<| = k - 1 THEN
    RETURN q
  ELSE
    RETURN SELECT( k - |S<| - 1 , S> )

```

Laufzeitanalyse:

Da $rg(q) \in [\lfloor \frac{1}{4} n \rfloor, \lceil \frac{3}{4} n \rceil]$ gilt $|S_{<}|, |S_{>}| \leq \frac{3}{4} n$.

Also gilt:

$$T(n) \leq \begin{cases} O(1) & , n < 100 \\ O(n) + T(\frac{3}{4} n) & , \text{sonst} \end{cases}$$

Behauptung:

$$T(n) \in O(n)$$

Beweis:

$$\begin{aligned}
T(n) &\leq c \cdot n + T\left(\frac{3}{4} n\right) \\
&\leq c \cdot n + c \left(\frac{3}{4} n\right) + T\left(\left(\frac{3}{4}\right)^2 n\right) \\
&\leq c \cdot n \sum_{i=0}^{\infty} \left(\frac{3}{4}\right)^i + O(1) \\
&\leq (4c) \cdot n + O(1) \\
&= O(n)
\end{aligned}$$

□

2.2.4 Implementierung von SPLITTER

Damit k-SELECT die versprochene lineare Laufzeit erreicht, müssen wir uns als nächstes die Implementierung von Splitter ansehen. Da wir uns in jedem Schritt einen neuen Splitter besorgen, muss die Laufzeit sehr gering gehalten werden.

Randomisierte Lösung

Die erste Idee ist es, statt dem SPLITTER mit den gewünschten Eigenschaften einfach einen zufällig Gewählten zu nehmen. Wenn man diese Laufzeit berechnen wird man auch auf eine lineare kommen.

Um dieses Problem zu lösen werden wir später Randomisierte Algorithmen betrachten und wie man Laufzeiten aus Erwartungswerten bestimmt.

BFRPT - Algorithm

Der Algorithmus wurde nach seinen Entdeckern Blum¹, Floyd², Pratt, Rivest³, Tarijan⁴ benannt.

Grundlegend funktioniert der Algorithmus folgender Maßen:

Man wählt zufällig eine Stichprobe $S' \subseteq S$ mit $|S'| = \lfloor \frac{n}{5} \rfloor$, so dass der Median von S' ein guter Splitter von S ist. Bestimme rekursiv den Media von S' .

Wählen von S'

Die Idee ist S in 5er Gruppen zu unterteilen. Innerhalb dieser Gruppen können wir den Median in konstanter Zeit findet. Baue aus den Medianen der 5er Gruppen die Menge S' und nimm deren Median.

Lemma: Der Median von S' ist ein guter SPLITTER von S , wenn n groß genug ist.

Anschaung: HIER WIRD NOCH EIN BILD UND ERKLÄRUNG EINGEFÜGT!!

Beweis:

Wir wollen prüfen, ob der Median g , den wir finden, wirklich SPLITTER Eigenschaften besitzt. Das heißt wir wollen wissen, ob min. $\frac{1}{4}$ kleiner und $\frac{1}{4}$ größer ist.

Größer:

Es sind $\lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil$ Elemente aus S' größer als g . Da alle Elemente aus S' Mediane ihrer 5er Gruppen sind, wissen wir, das in jeder dieser Gruppen 3 Elemente größer sind als g . Dies gilt für alle Gruppen, außer die Gruppe von g selber und die mögliche letzte Gruppe.

Dies führt zu $3 \cdot \lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil - 3$

Kleiner:

Es gibt ebenso $\lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil$ Gruppen, deren Meridiane kleiner sind als g . In jeder

¹Turing Award 1995

²Turing Award 1978

³Turing Award 2002

⁴Turing Award 1986

dieser Gruppe, wissen wir von 3 Elementen die kleiner sind, bis auf die Gruppe von g und die letzte Gruppe, die in diese Klasse fallen könnte. Dies führt zu mindestens $3 \lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil - 3$ Elementen, die kleiner sind als q .

Zusammensetzen:

Es gilt $3 \lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil - 3 \geq 3 \cdot \frac{1}{2} \cdot \frac{1}{5}n - 3 = \frac{3}{10}n - 3$.

Für einen guten SPLITTER, muss die Anzahl der außerhalb liegenden Elemente (sowohl größer als auch kleiner) größer als $\frac{1}{4}n$ sein.

$$\begin{aligned} \Rightarrow \frac{3}{10}n - 3 \geq \frac{1}{4} &\Leftrightarrow \left(\frac{3}{10} - \frac{1}{4}\right)n \geq 3 \\ &\Leftrightarrow n \geq 60 \end{aligned}$$

Wir sehen hier, dass wir mit dem Verfahren garantiert einen guten SPLITTER finden, wenn wir mehr als 60 Elemente haben. Mit diesem Problem haben wir uns aber schon im Algorithmus beschäftigt. Dort haben wir gesagt, dass wir bei Listen bestimmter Größe das ganze Problem mit BRUTEFORCE lösen wollen. Damit werden die Listen in denen wir einen SPLITTER suchen immer garantiert 60 Elemente besitzen.

Nachdem wir nun wissen, dass wir mit diesem Verfahren einen SPLITTER erhalten, müssen wir prüfen, ob dieses Verfahren den Algorithmus asymptotisch langsamer macht oder ob wir bei einer linearen Laufzeit bleiben.

Laufzeit K-SELECT mit BFRPT Algorithmus

Betrachten wir noch einmal, wie unser Algorithmus nun nach dem Einsetzen des SPLITTER Codes aussieht.

```

SELECT (S, K)
  IF |S| < 100 THEN
    RETURN BRUTEFORCE(S, K)
  Unterteile S in 5er Gruppen
  S' ← {Median jeder 5er Gruppe}
  q ← SELECT(S', ⌈(|S'|+1)/2⌉)
  S< ← {s ∈ S | s < q}
  S> ← {s ∈ S | s > q}
  IF |S<| ≥ k THEN
    RETURN SELECT(S<, K)
  ELSE IF |S<| = K - 1 THEN
    RETURN q
  ELSE
    RETURN SELECT(S>, K - |S<| - 1)

```

Nun können wir aus dem Programm die Anzahl der Vergleiche ablesen:

$$T(n) \leq \begin{cases} O(1) & , n < 100 \\ O(n) + T\left(\lceil \frac{n}{5} \rceil\right) + T\left(\frac{3}{4}n\right) & , \text{sonst} \end{cases}$$

Behauptung: $T(n) = O(n)$

Beweis: Induktion über n mit $\exists \alpha > 0 : T(n) \leq \alpha \cdot n$

I.A.: $n < 100$ klar, da $T(n)$ konstant.

I.S.: $n \rightarrow n + 1$

$$\begin{aligned} T(n) &\leq c \cdot n + T\left(\left\lceil \frac{n}{5} \right\rceil\right) + T\left(\frac{3}{4}n\right) \\ &\stackrel{\text{I.A.}}{\leq} c \cdot n + \alpha \left\lceil \frac{n}{5} \right\rceil + \alpha \left(\frac{3}{4}n\right) \\ &\leq c \cdot n + \alpha \left(\frac{n}{5} + 1\right) + \alpha \left(\frac{3}{4}n\right) \\ &= n \left(c + \frac{19}{20}\right) \alpha + \alpha \\ &\leq \alpha \cdot n \end{aligned}$$

Den letzten Schritt darf jeder für sich selbst nachvollziehen.
So sehen wir, dass die Laufzeit immer noch lineare ist.

□

2.3 Bemerkung

Was wollten wir an diesem Beispiel sehen?

Was wir erreicht haben, ist ein Algorithmus, der kurz, elegant und optimal ist. Um diesen zu gewinnen mussten wir nicht-triviale Strukturen benutzen, auf die man nicht mehr so leicht kommt, wie auf Quicksort oder Mergesort.

Die Laufzeit des Algorithmus war nicht sofort offensichtlich und brauchte eine Analyse samt Beweis.

Mit solchen Algorithmen werden wir uns im folgenden in der Vorlesung beschäftigen. Ein jeder hat mindestens einen kleinen Kniff dabei.

Zu Bemerkungen bleibt, dass der oben genannte und analysierte Algorithmus zwar linear läuft, die Konstanten sind allerdings so hoch, dass bis zu einer Listengröße von 2^{25} Quicksort und anschließendes Nehmen des k -ten Elements schneller ist.