

Lösung - Übung 6

Tutoraufgabe 1 (Gerrymandering):

Verschiedene Wahlsysteme haben verschiedene Stärken und Schwächen. Eine Anforderung an ein Wahlsystem ist, dass es nach der Wahl nicht nur einen Präsident für das ganze Land gibt, sondern auch einen Vertreter des lokalen Wahlkreises. Falls der Präsident jedoch nicht direkt, sondern von den Wahlkreisvertretern gewählt wird, so ergibt sich das Problem des Gerrymandering: Das Wahlergebnis hängt stark davon ab welche Personen welchem Wahlkreis zugeordnet sind. Die Person, welche bestimmt wer welchem Wahlkreis zugeordnet ist hat somit einen starken Einfluss auf den Wahlausgang. Wir betrachten dieses Problem nun in einem Zweiparteiensystem, das heißt jede Person kann entweder den Vertreter von Partei A oder den Vertreter von Partei B wählen. Die Wahlkreisvertreter können wiederum entweder den Präsidenten der Partei A oder den Präsidenten der Partei B wählen. Jeder Wahlkreisvertreter entscheidet sich für den Präsidenten der eigenen Partei.

- a) Angenommen jeder Wahlkreis besteht aus 5 Personen. Wie viel Prozent der Wählerstimmen müssen mindestens für Vertreter von Partei A abgegeben werden, damit der Präsident von Partei A durch die Wahlkreisvertreter gewählt werden kann?
- b) Angenommen jeder Wahlkreis besteht aus einer beliebigen aber festen Anzahl $k \in \mathbb{N}$ Personen.
 - i. Wie viel Prozent der Wählerstimmen müssen mindestens für Vertreter von Partei A abgegeben werden, damit der Präsident von Partei A durch die Wahlkreisvertreter gewählt werden kann?
 - ii. Wie muss k gewählt werden um mit einem möglichst geringen Prozentsatz trotzdem zu gewinnen.

Hinweise:

- Gehen Sie davon aus, dass in allen Wahlkreisen und während der Präsidentschaftswahl die Parteien nie gleich viele Stimmen bekommen.

Lösung

Sei $P = \{p_1, \dots, p_n\}$ die Menge der Wähler und $v : P \rightarrow \{A, B\}$ die Funktion welche jedem Wähler seine Stimme zuordnet. Weiter sein $K = \{K_1, \dots, K_m\}$ die Menge der Wahlkreise mit $K_1 \uplus \dots \uplus K_m = P$. Die Wahlkreise partitionieren also die Menge der Wähler.

Die Stimmen der Wahlkreisvertreter lassen sich nun wie folgt berechnen:

$$v_p : K \rightarrow \{A, B\}, K_i \mapsto \text{median}(\{\{v(p) \mid p \in K_i\}\})$$

Sei nun $x = \text{median}(\{\{v_p(K_i) \mid K_i \in K\}\})$. Der Präsident der Partei x wird von den Wahlkreisvertretern gewählt.

Um die Wahl zu gewinnen muss der Präsident der Partei A mehr als 50% der Stimmen der Wahlkreisvertreter erhalten.

Hinweis:

- Die Notation $\{\{1, 1, 1\}\}$ stellt eine Multimenge mit drei Mal dem Wert 1 dar. Eine normale Menge würde natürlich nur einmal den Wert 1 enthalten. Für die Medianberechnung ist die Anzahl jedoch entscheidend. Daher nutzen wir hier Multimengen.
- a) Um die Stimme eines Wahlkreisvertreters zu erhalten, muss in diesem Wahlkreis mindestens drei der fünf Stimmen (60%) für Partei A abgegeben worden sein.

Partei A muss also in mehr als 50% der Wahlkreise mindestens 60% der Stimmen erhalten um gewinnen zu können. Insgesamt ist es also ausreichend mehr als $50\% \cdot 60\% = 30\%$ der Stimmen zu erhalten. Wie viel mehr als 30% der Stimmen benötigt werden hängt von der Anzahl der Wahlkreise ab. Je mehr Wahlkreise es gibt, desto geringer ist die Anzahl an zusätzlich benötigten Stimmen. Mit höchstens 30% der Stimmen ist es also nicht möglich die Wahl zu gewinnen.

- b) Sei k die Anzahl der Personen in jedem Wahlkreis und m die Anzahl der Wahlkreise. Bei einer geraden Anzahl von Stimmen braucht es eine ganze Stimme mehr um zu gewinnen, bei einer ungeraden Anzahl von Stimmen reicht jedoch eine halbe Stimme mehr. Wir gehen daher nun davon aus, dass sowohl k als auch m ungerade ist, da so jeweils weniger Stimmen notwendig sind. Für gerade Stimmanzahlen funktioniert das Argument aber analog, mit einer etwas anderen Formel. Sei nun $k = 2k' + 1$ und $m = 2m' + 1$. Um die Wahl zu gewinnen sind mindestens $(k' + 1)(m' + 1)$ der $(2k' + 1)(2m' + 1)$ Stimmen notwendig, also ein Stimmanteil von

$$\frac{(k' + 1)(m' + 1)}{(2k' + 1)(2m' + 1)} > \frac{(k' + 1)(m' + 1)}{(2k' + 2)(2m' + 2)} = \frac{(k' + 1)(m' + 1)}{4(k' + 1)(m' + 1)} = \frac{1}{4}.$$

Dieser Wert ist immer größer als 25%. Mit ausreichend großen k' und m' ist es jedoch möglich beliebig nah an die 25% heranzukommen.

Tatsächlich konvergiert obige Formel gegen 30% für $k = 5$ und m gegen unendlich.

Wir stellen fest, dass diese Fragestellung stark verwandt ist mit der Frage wie weit der Median der Mediane vom tatsächlichen Median entfernt sein kann.

Hier einige Anmerkungen welche erklären, warum dieses Problem im deutschen Wahlsystem praktisch nicht auftritt, um das Vertrauen in das deutsche Wahlsystem nicht unnötig zu erschüttern:

- In Deutschland werden zwar Wahlkreisvertreter gewählt, welche auch im Parlament den Kanzler mitwählen. Das Stimmverhältnis im Parlament (und somit auch bei der Kanzlerwahl) wird jedoch per Verhältniswahlrecht über die Zweitstimmen festgelegt.
- In Deutschland sind die Wahlen geheim. Das bedeutet es gibt keine zentrale Stelle, welche vor der Wahl weiß welche Person welche Partei wählen wird. Somit ist eine Zuordnung der Stimmen zu den passenden Wahlkreisen praktisch unmöglich.
- In Deutschland bleiben die Wahlkreise über relativ lange Zeiträume unverändert. Es ist beispielsweise nicht üblich nach der Wahl die Wahlkreise neu festzulegen, um bei der nächsten Wahl eine bessere Gewinnchance zu haben.
- Deutschland hat kein Zweiparteiensystem.

Tutoraufgabe 2 (Selektionsalgorithmus mit 9er-Gruppen):

In der Vorlesung haben Sie gesehen, dass der Median der Mediane in linearer Zeit berechnet werden kann. Dazu wurde mithilfe des Medians der Mediane ein Pivot-Element gewählt. Bei der Berechnung des Medians der Mediane wurde zunächst die Eingabe in 5er-Gruppe unterteilt, anschließend jeweils der Median jeder 5er-Gruppe gebildet und schließlich der Median all dieser Mediane gebildet.

Angenommen die Eingabe wäre nicht in 5er-Gruppen eingeteilt worden, sondern in **9er-Gruppen**. Der restliche Selektionsalgorithmus bleibt hingegen unverändert.

Arbeitet der Selektionsalgorithmus weiterhin in linearer Zeit? Beweisen Sie Ihre Antwort.

Lösung

Der Selektionsalgorithmus benötigt mit 9er-Gruppen ebenfalls lineare Zeit.

Zunächst versuchen wir den Beweis der Linearität aus der Vorlesung zu übertagen.

Dafür haben wir, dass mindestens die Hälfte der $\lfloor \frac{n}{9} \rfloor$ Mediane, also $\lceil \frac{1}{2} \lfloor \frac{n}{9} \rfloor \rceil$ Mediane, größer gleich sind als der

Median of Medians. Mindestens die Hälfte der Elemente in einer Gruppe (das sind bei 9 genau 5) sind größer gleich deren Median. Also sind mindestens $5 \lceil \frac{1}{2} \lfloor \frac{n}{9} \rfloor \rceil$ Elemente größer gleich der Median of Medians. Damit ist $|S_{<}| < n - 5 \lceil \frac{1}{2} \lfloor \frac{n}{9} \rfloor \rceil < \lfloor \frac{7n}{9} \rfloor$ für $n \geq 50$. Symmetrisch ist damit auch $|S_{>}| < n - 5 \lceil \frac{1}{2} \lfloor \frac{n}{9} \rfloor \rceil < \lfloor \frac{7n}{9} \rfloor$ für $n \geq 50$. Folgende Laufzeitabschätzung lesen wir aus dem Algorithmus ab.

$$T(n) \leq \begin{cases} c_1, & \text{falls } n < 50 \\ c_2 n + T(\lfloor \frac{n}{9} \rfloor) + T(\lfloor \frac{7n}{9} \rfloor), & \text{falls } n \geq 50 \end{cases}$$

Wir benötigen $c_2 n$ Schritte für die Partitionierung, $T(\lfloor \frac{n}{9} \rfloor)$ Schritte um den Median der Mediane exakt zu bestimmen und müssen auf maximal $\frac{7n}{9}$ Werten im rekursiven Aufruf arbeiten.

Wir wollen zeigen, dass $T(n) \leq g(n)$ mit $g(n) = cn = O(n)$.

Induktionsanfang Angenommen $c_1 \leq c$, dann gilt $T(n) = c_1 \leq cn = g(n)$ für $n < 50$.

Induktionshypothese Angenommen es gilt $T(m) \leq cm$ für $m < n$.

Induktionsschritt Sei $n \geq 50$.

$$\begin{aligned} T(n) &\leq c_2 n + T(\lfloor \frac{n}{9} \rfloor) + T(\lfloor \frac{7n}{9} \rfloor) \\ &\leq c_2 n + c \lfloor \frac{n}{9} \rfloor + c \lfloor \frac{7n}{9} \rfloor && \text{(Induktionshypothese)} \\ &\leq c_2 n + c \frac{n}{9} + c \frac{7n}{9} \\ &= c_2 n + cn - c \frac{n}{9} \\ &= cn + n(c_2 - c \frac{1}{9}) \\ &\leq cn && \text{falls } c_2 - c \frac{1}{9} \leq 0 \end{aligned}$$

Die Zusatzbedingung lässt sich wie folgt in eine untere Grenze für c umformen.

$$c_2 - c \frac{1}{9} \leq 0 \iff c_2 \leq c \frac{1}{9} \iff 9c_2 \leq c$$

Die Induktion funktioniert also für $\max(c_1, 9c_2) \leq c$. Damit benötigt der Selektionsalgorithmus für 9er-Gruppen ebenfalls lineare Zeit.

Tutoraufgabe 3 (Geschlossenes Hashing mit großen Hashtabellen):

In dieser Aufgabe beschäftigen wir uns mit der erwarteten Anzahl an Sondierungen in einer Hashtabelle bei Nutzung von geschlossenem Hashing. Dafür möchten wir zeigen, dass bei Nutzung einer Hashtabelle, die ausreichend groß ist (genauer, dass $n \leq \frac{m}{2}$ mit n die Anzahl der Eingaben und m die Größe der Hashtabelle) die Länge einer solche Sondierungssequenz im Averagecase lediglich logarithmisch groß in der Größe der Eingabe ist.

Dazu nehmen wir an, dass wir *uniformes hashing* nutzen, d.h. für eine zufällige Eingabe (Schlüssel) x und eine uniforme Hashfunktion mit j ter Sondierung $h(x, j)$ auf eine Menge der Größe m ist jede Sequenz $h(x, 1) \dots h(x, m)$ aus unterschiedlichen Hashwerten gleich wahrscheinlich. Insbesondere ist damit auch jeder Wert für $h(x, j)$ gleich wahrscheinlich.

Sei hierfür X_i die Anzahl der Sondierungen, die beim i ten Einfügen in die Hashtabelle vorgenommen werden müssen.

Hinweis:

- Aufgrund der Annahme $n \leq \frac{m}{2}$ können wir etwas bessere Abschätzungen angeben als in der Vorlesung.
- a) Wir nehmen an, wir nutzen uniformes Hashing. Zeigen Sie, dass die Wahrscheinlichkeit, dass wir für die i te Eingabe mit $i \leq n$ in eine Hashtabelle der Größe m mit $n \leq \frac{m}{2}$ echt mehr als k viele Sondierungen vornehmen müssen, maximal $\frac{1}{2^k}$ ist. Formal ausgedrückt, dass $\mathbb{P}(X_i > k) < \frac{1}{2^k}$ gilt.
- b) Zeigen Sie, dass die Wahrscheinlichkeit, dass die i te Eingabe mit $i \leq n$ in eine Hashtabelle der Größe m mit $n \leq \frac{m}{2}$ mehr als $2 \log_2(n)$ viele Sondierungen benötigt, in $O(\frac{1}{n^2})$ ist. Formal ausgedrückt, dass $\mathbb{P}(X_i > 2 \log_2(n)) = O(\frac{1}{n^2})$ gilt.

- c) Sei $X = \max_{1 \leq i \leq n} X_i$ die maximale Anzahl an Sondierung aller n Eingaben in eine Hashtabelle der Größe m mit $n \leq \frac{m}{2}$. Zeigen Sie, dass die Wahrscheinlichkeit, dass X größer ist als $2 \log_2(n)$ in $O(\frac{1}{n})$ ist. Formal ausgedrückt, dass $\mathbb{P}(X > 2 \log_2(n)) = O(\frac{1}{n})$.
- d) Sei $X = \max_{1 \leq i \leq n} X_i$ die maximale Anzahl an Sondierung aller n Eingaben in eine Hashtabelle der Größe m mit $n \leq \frac{m}{2}$. Zeigen Sie, dass die maximalen Sondierungen pro Eingabe X im Averagecase in $O(\log(n))$ ist. Formal ausgedrückt, dass $\mathbb{E}(X) = O(\log(n))$.

Lösung

- a) Da $n \leq \frac{m}{2}$ gilt, ist stets mindestens die Hälfte aller Einträge frei. Damit eine Kollision entsteht, müssen wir aber bereits auf ein belegtes Feld landen. Da $h(x, j)$ aber uniform verteilt ist (sowohl für die Sequenz, als auch für jedes i), ist die Wahrscheinlichkeit eines der belegten Felder zu landen maximal $\frac{1}{2}$. Die Wahrscheinlichkeit, dass dies für eine Eingabe x aber nun k oft passiert ist damit auch maximal $\frac{1}{2^k}$.
- b) Mithilfe der Ungleichung der vorherigen Teilaufgabe haben wir mit $k = 2 \log_2(n)$ dass $\mathbb{P}(X_i > 2 \log_2(n)) < \frac{1}{2^{2 \log_2(n)}} = \frac{1}{2^{\log_2(n^2)}} = \frac{1}{n^2} = O(\frac{1}{n^2})$.
- c) Wir haben nach Definition von X und dem Maximum, dass

$$\mathbb{P}(X > 2 \log_2(n)) = \mathbb{P}(\max_{1 \leq i \leq n} X_i > 2 \log_2(n)) = \mathbb{P}\left(\bigvee_{1 \leq i \leq n} X_i > 2 \log_2(n)\right).$$

Wir können die Veroderung nach oben durch die Summe der einzel Wahrscheinlichkeiten abschätzen als

$$\mathbb{P}\left(\bigvee_{1 \leq i \leq n} X_i > 2 \log_2(n)\right) \leq \sum_{1 \leq i \leq n} \mathbb{P}(X_i > 2 \log_2(n)).$$

Und schließlich mit der Lösung der vorhergehenden Aufgabe haben wir

$$\sum_{1 \leq i \leq n} \mathbb{P}(X_i > 2 \log_2(n)) \leq \sum_{1 \leq i \leq n} \frac{1}{n^2} = \frac{n}{n^2} = \frac{1}{n} = O\left(\frac{1}{n}\right).$$

- d) Die längste Anzahl an Sondierungen möglich ist n , da wir auch nur n viele Einträge in der Hashtabelle haben werden.

Weiterhin wissen wir aus der vorherigen Aufgabe, dass $\mathbb{P}(X > 2 \log_2(n)) \leq \frac{1}{n}$ gilt und da Wahrscheinlichkeiten maximal 1 sind, dass $\mathbb{P}(X \leq 2 \log_2(n)) \leq 1$ gilt. Somit ist:

$$\begin{aligned} \mathbb{E}(X) &\leq \mathbb{P}(X \leq 2 \log_2(n)) \cdot 2 \log_2(n) + \mathbb{P}(X > 2 \log_2(n)) \cdot n \\ &\leq 1 \cdot 2 \log_2(n) + \frac{1}{n} \cdot n \\ &= 2 \log_2(n) + 1 = O(\log(n)). \end{aligned}$$

Aufgabe 4 (Selektion des zweitkleinsten Elements):

7 Punkte

Entwerfen Sie einen Algorithmus, welcher aus n Elementen das zweitkleinste auswählt und dafür höchstens $n + \lceil \log_2(n) \rceil$ Vergleiche von Elementen benötigt. Dabei kann es hilfreich sein zusätzlich das kleinste Element zu finden. Begründen Sie Ihre Antwort.

Hinweis:

- Sie dürfen beliebig viele Vergleiche zwischen Booleans nutzen.

Lösung

Wir berechnen zunächst das Minimum indem wir unten stehenden Vergleichsbaum aufbauen. Die erste Ebene entsteht indem wir jeweils zwei direkt nebeneinander liegende Elemente miteinander vergleichen und jeweils das kleinere Element in die neue Ebene schreiben. Zusätzlich merken wir uns auf welcher Seite das kleinere Element stand. Außerdem merken wir uns ob wir gerade die erste Baumebene bauen. Die zweite Ebene entsteht aus der ersten Ebene in genau derselben Art und Weise. Damit haben wir $\frac{n}{2}$ Vergleiche auf der ersten Ebene, $\frac{n}{4}$ auf der zweiten, und so weiter. Zusammen ergibt das

$$\sum_{i=1}^{\lceil \log_2(n) \rceil} \frac{n}{2^i} = -n + \sum_{i=0}^{\lceil \log_2(n) \rceil} \frac{n}{2^i} < n \cdot \frac{1}{0.5} - n = n.$$

Anschließend können wir mit folgender Funktion das zweitkleinste Element im Array finden, welche das Minimum der unten rot hervorgehobenen Werte berechnet, wobei wir Folgendes annehmen:

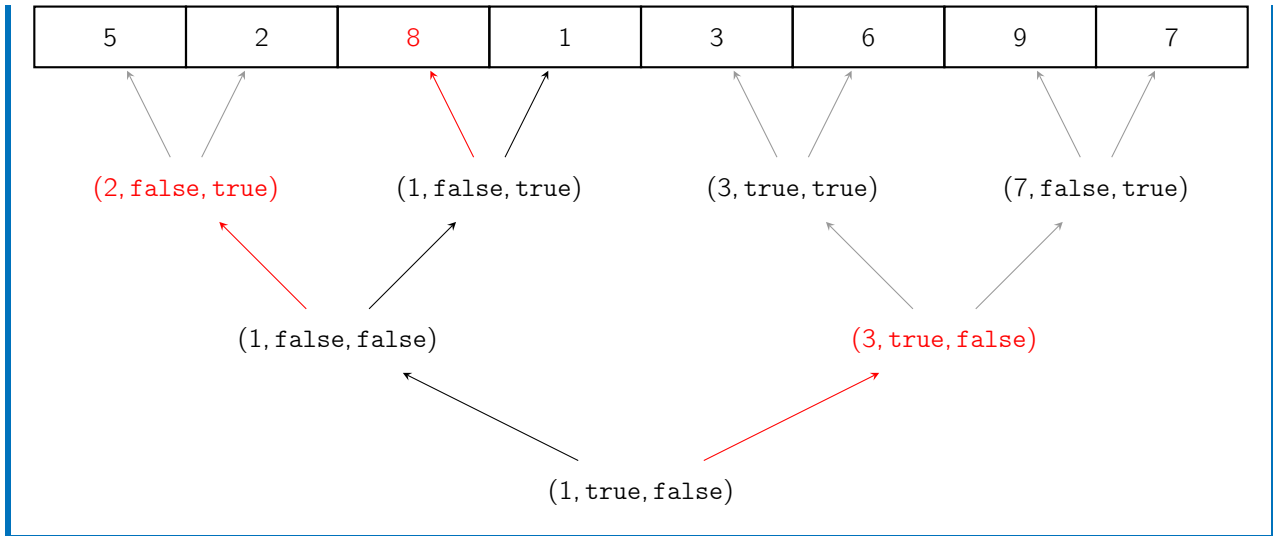
- `(x,y,z).value == x`
- `(x,y,z).minIsLeft == y`
- `(x,y,z).isLastLayer == z`

```
def min2(a):
    if a.isLastLayer:
        if a.minIsLeft:
            return a.right
        else:
            return a.left
    else:
        if a.minIsLeft:
            return min(a.right.value, min2(a.left))
        else:
            return min(a.left.value, min2(a.right))
```

Ohne Beschränkung der Allgemeinheit nehmen wir an, dass das Minimum im linken Teilbaum ist. Wir wissen, dass die Wurzel im rechten Teilbaum das kleinste Element des rechten Teilbaums ist. Rekursiv können wir nun noch das zweit kleinste Element im linken Teilbaum berechnen. Das kleinere dieser beiden muss nun auch das zweit kleinste Element sein.

Falls das Minimum im rechten Teilbaum ist, ist die Korrektheit symmetrisch begründet.

Für dieses Verfahren brauchen wir höchstens $\lceil \log_2(n) \rceil$ Vergleiche, wodurch wir insgesamt mit maximal $n + \lceil \log_2(n) \rceil$ Vergleichen zwischen Elementen auskommen.



Aufgabe 5 (Selektionsalgorithmus mit 7er-Gruppen):

6 Punkte

In der Vorlesung haben Sie gesehen, dass der Median der Mediane in linearer Zeit berechnet werden kann. Dazu wurde mithilfe des Medians der Mediane ein Pivot-Element gewählt. Bei der Berechnung des Medians der Mediane wurde zunächst die Eingabe in 5er-Gruppe unterteilt, anschließend jeweils der Median jeder 5er-Gruppe gebildet und schließlich der Median all dieser Mediane gebildet.

Angenommen die Eingabe wäre nicht in 5er-Gruppen eingeteilt worden, sondern in **7er-Gruppen**. Der restliche Selektionsalgorithmus bleibt hingegen unverändert.

Arbeitet der Selektionsalgorithmus weiterhin in linearer Zeit? Beweisen Sie Ihre Antwort.

Lösung

Der Selektionsalgorithmus benötigt mit 7er-Gruppen ebenfalls lineare Zeit.

Zunächst versuchen wir den Beweis der Linearität aus der Vorlesung zu übertragen.

Dafür haben wir, dass mindestens die Hälfte der $\lfloor \frac{n}{7} \rfloor$ Mediane, also $\lceil \frac{1}{2} \lfloor \frac{n}{7} \rfloor \rceil$ Mediane, größer gleich sind als der Median of Medians. Mindestens die Hälfte der Elemente in einer Gruppe (das sind bei 7 genau 4) sind größer gleich deren Median. Also sind mindestens $4 \lceil \frac{1}{2} \lfloor \frac{n}{7} \rfloor \rceil$ Elemente größer gleich der Median of Medians. Damit ist $|S_{<}| < n - 4 \lceil \frac{1}{2} \lfloor \frac{n}{7} \rfloor \rceil < \lfloor \frac{11n}{14} \rfloor$ für $n \geq 50$. Symmetrisch ist damit auch $|S_{>}| < n - 4 \lceil \frac{1}{2} \lfloor \frac{n}{7} \rfloor \rceil < \lfloor \frac{11n}{14} \rfloor$ für $n \geq 50$. Folgende Laufzeitabschätzung lesen wir aus dem Algorithmus ab.

$$T(n) \leq \begin{cases} c_1, & \text{falls } n < 50 \\ c_2 n + T(\lfloor \frac{n}{7} \rfloor) + T(\lfloor \frac{11n}{14} \rfloor), & \text{falls } n \geq 50 \end{cases}$$

Wir benötigen $c_2 n$ Schritte für die Partitionierung, $T(\lfloor \frac{n}{7} \rfloor)$ Schritte um den Median der Mediane exakt zu bestimmen und müssen auf maximal $\frac{11n}{14}$ Werten im rekursiven Aufruf arbeiten.

Wir wollen zeigen, dass $T(n) \leq g(n)$ mit $g(n) = cn = O(n)$.

Induktionsanfang Angenommen $c_1 \leq c$, dann gilt $T(n) = c_1 \leq cn = g(n)$ für $n < 50$.

Induktionshypothese Angenommen es gilt $T(m) \leq cm$ für $m < n$.

Induktionsschritt Sei $n \geq 50$.

$$\begin{aligned}
 T(n) &\leq c_2 n + T(\lfloor \frac{n}{7} \rfloor) + T(\lfloor \frac{11n}{14} \rfloor) \\
 &\leq c_2 n + c \lfloor \frac{n}{7} \rfloor + c \lfloor \frac{11n}{14} \rfloor && \text{(Induktionshypothese)} \\
 &\leq c_2 n + c \frac{n}{7} + c \frac{11n}{14} \\
 &= c_2 n + cn - c \frac{n}{14} \\
 &= cn + n(c_2 - c \frac{1}{14}) \\
 &\leq cn && \text{falls } c_2 - c \frac{1}{14} \leq 0
 \end{aligned}$$

Die Zusatzbedingung lässt sich wie folgt in eine untere Grenze für c umformen.

$$c_2 - c \frac{1}{14} \leq 0 \iff c_2 \leq c \frac{1}{14} \iff 14c_2 \leq c$$

Die Induktion funktioniert also für $\max(c_1, 14c_2) \leq c$. Damit benötigt der Selektionsalgorithmus für 7er-Gruppen ebenfalls lineare Zeit.

Aufgabe 6 (Hashing):

2 + 3 + 2 = 7 Punkte

- a) Wir speichern n Objekte in einer einfach verketteten Liste mit Schlüssel $\text{key}[o]$ und Hashwert $h(\text{key}[o])$ für Objekte o . Die Schlüssel sind beliebig lange Zeichenketten mit minimal Länge d und die Hashwerte haben maximal Länge m mit $m \ll d$.

Wie können Sie die Hashwerte ausnutzen, um ein Objekt o in der Liste zu suchen? Begründen Sie ihre Antwort.

- b) In dieser Aufgabe nehmen wir einfaches uniformes Hashing an. Das heißt für eine zufällige Eingabe x und eine Hashfunktion h ist jeder Hashwert für $h(x)$ gleich wahrscheinlich.

Berechnen Sie die erwartete Anzahl an Kollisionen beim Hashen der verschiedenen Werte k_1, \dots, k_n , das heißt berechnen Sie $\mathbb{E}(|\{ \{k_i, k_j\} \mid k_i \neq k_j \text{ und } h(k_i) = h(k_j) \}|)$.

Hinweise:

- Beachten Sie, dass wir z.B. bei vier Elementen k_1, k_2, k_3, k_4 mit $h(k_1) = h(k_2) = h(k_3) = h(k_4)$ die 6 Kollisionen $\{k_1, k_2\}, \{k_1, k_3\}, \{k_1, k_4\}, \{k_2, k_3\}, \{k_2, k_4\}, \{k_3, k_4\}$ haben.
 - Wegen einfachen uniformen Hashings einer Hashfunktion $h : A \rightarrow B$ folgt, dass für $x \neq y$ die Wahrscheinlichkeit, dass x und y gleich gehasht werden, uniform ist, also dass $\mathbb{P}(h(x) = h(y)) = \frac{1}{|B|}$.
 - Nutzen Sie aus, dass der Erwartungswert linear ist, also dass $\mathbb{E}(\sum_{i=0}^n X_i) = \sum_{i=0}^n \mathbb{E}(X_i)$ gilt.
- c) Angenommen wir benutzen offenes Hashing und wollen verhindern, dass bei einem zu hohen Füllgrad (= Anzahl der eingefügten Elemente geteilt durch die Größe der Hash-Tabelle) die Listen zu lang werden und dadurch die Suche in den Listen den Zugriffsaufwand dominiert.

Wir könnten dies erreichen, indem wir die Größe der Tabelle dynamisch anpassen, d.h. jedesmal, wenn der Füllgrad der Tabelle z.B. den Wert 2 übersteigt, erzeugen wir eine neue Tabelle mit doppelter Größe und überführen alle Element aus der alten Tabelle in die neue.

Ist dies eine sinnvolle Lösung?

Lösung

- a) Da der Hashwert der Objekte in der Liste bereits gespeichert sind, können wir einmalig den Hashwert $h(\text{key}[o])$ der Eingabe o mit Schlüssel $\text{key}[o]$ berechnen. Nun vergleichen wir die Hashwerte der Objekte in der Liste mit dem Hashwert der Eingabe. Sind die Hashwerte gleich, vergleichen wir anschließend noch deren Schlüssel. Damit reduzieren wir die Worstcase Laufzeitkomplexität der Suche in der Liste, da wir statt für jedes Objekt ein Vergleich mit Komplexität $O(d)$ ein Vergleich der Komplexität $O(m)$ in vielen Fällen vornehmen können.

- b) Wir berechnen zuerst, mit wie vielen Elementen der Wert k_i kollidieren wird, das heißt wir berechnen $\mathbb{E}(|\{k_j, k_i\} \mid j > i \text{ und } h(k_j) = h(k_i)|)|$ und addieren anschließend alle Erwartungswerte zusammen für den gewünschten Erwartungswert. Dies können wir tun, da $\{k_j, k_i\} \mid n \geq j > i \text{ und } h(k_j) = h(k_i)\}$ für jedes i eine Partition der Menge $\{k_i, k_j\} \mid k_i \neq k_j \text{ und } h(k_i) = h(k_j)\}$ darstellt.

Um $\mathbb{E}(|\{k_j, k_i\} \mid j > i \text{ und } h(k_j) = h(k_i)|)|$ zu berechnen, addieren wir die Wahrscheinlichkeiten, dass k_i und k_j kollidieren, multipliziert mit der Anzahl der Kollisionen die k_i produziert, das heißt wir haben

$$\mathbb{E}(|\{k_j, k_i\} \mid n \geq j > i \text{ und } h(k_j) = h(k_i)|)| = \sum_{n \geq j > i} \mathbb{P}(h(k_j) = h(k_i)) \cdot 1 = \sum_{n \geq j > i} \mathbb{P}(h(k_j) = h(k_i)) .$$

Die Wahrscheinlichkeit für $\mathbb{P}(h(k_i) = h(k_j))$ ist wegen der Annahme einfacher uniformem Hashings genau $\frac{1}{m}$, wobei m die Größe des Bildes von h ist, daher haben wir

$$\sum_{n \geq j > i} \mathbb{P}(h(k_j) = h(k_i)) = \sum_{n \geq j > i} \frac{1}{m} = \frac{n-i}{m} .$$

Addieren wir nun alle Erwartungswerte zusammen haben wir

$$\begin{aligned} & \mathbb{E}(|\{k_i, k_j\} \mid k_i \neq k_j \text{ und } h(k_i) = h(k_j)|)| \\ &= \sum_{i=1}^n \mathbb{E}(|\{k_j, k_i\} \mid n \geq j > i \text{ und } h(k_j) = h(k_i)|)| \\ &= \sum_{i=1}^n \frac{n-i}{m} = \frac{n^2 - \frac{n(n+1)}{2}}{m} = \frac{n^2 - n}{2m} . \end{aligned}$$

- c) In der Tat ist dies eine sinnvolle Idee. Durch das Einfügen würden wir nun eine Laufzeit von $2m$ für das $2m$ te Element haben. Dies ist jedoch in Ordnung, da wir in den letzten m vielen Einfügeoperationen (so viele gab es bevor wir das letzte mal die Hashtabelle verdoppelt haben) immer sehr schnellen Zugriff hatten, kann diese Effizienzverbesserung der letzten m Einfügeoperationen genutzt werden, um nun beim $2m$ ten Element mehr Laufzeitkomplexität zu haben. Tatsächlich würde sich das ausgleichen: m viele Operationen in $O(1)$ und 1 Operation in $O(m)$ ergibt miteinander „verrechnet“ eine „verrechnete“ Laufzeit von $O(1)$ pro Operation.

Man kann das vorherige Argument auch formal in der so genannten Amortisierten Analyse^a nutzen und bekommt bei diesem Verfahren tatsächlich eine erwartete amortisierte Laufzeit von $O(1)$ pro Operation.

Aufgabe 7 (Programmierung in Python - Binäre Suchbäume):

4 + 6 + 3 + 3 + 4 = 20 Punkte

Bearbeiten Sie die Python Programmieraufgaben. In dieser praktischen Aufgabe werden Sie sich mit Bäumen, genauer mit binären Bäumen und binären Suchbäumen, auseinandersetzen. Diese Aufgabe dient dazu einige Konzepte der Vorlesung zu wiederholen.

Zum Bearbeiten der Programmieraufgabe können Sie einfach den Anweisungen des Notebooks *blatt06-python.ipynb* folgen. Das Notebook steht in der .zip-Datei zum Übungsblatt im Lernraum zur Vergütung.

Ihre Implementierung soll immer nach dem `# YOUR CODE HERE` Statement kommen. Ändern Sie keine weiteren Zellen.

Laden Sie spätestens bis zur Deadline dieses Übungsblatts auch Ihre Lösung der Programmieraufgabe im Lernraum hoch. Die Lösung des Übungsblatts und die Lösung der Programmieraufgabe muss im Lernraum an derselben Stelle hochgeladen werden. Die Lösung des Übungsblatts muss dazu als .pdf-Datei hochgeladen werden. Die Lösung der Programmieraufgabe muss als .ipynb-Datei hochgeladen werden.

Übersicht der zu bearbeitenden Aufgaben:

- a) Binäre Suchbäume

- `insert()`
- `is_binary_search_tree()`
- `bin_tree_2_list()`
- `list_2_bin_tree()`
- `bin_tree_2_bin_search_tree()`

Lösung

Die Lösung der Programmieraufgaben finden Sie im Lernraum. Die Datei trägt den Namen *blatt06-python-solution.ipynb*.