

Script generated by TTT

Title: Seidl: GAD (19.05.2015)

Date: Tue May 19 13:45:10 CEST 2015

Duration: 149:57 min

Pages: 55

Familie für universelles Hashing

Definiere für jeden Vektor

$$\mathbf{a} = (a_1, \dots, a_k) \in \{0, \dots, m-1\}^k$$

mittels Skalarprodukt

$$\mathbf{a} \cdot \mathbf{x} = \sum_{i=1}^k a_i x_i$$

eine Hashfunktion von der Schlüsselmenge
in die Menge der Zahlen $\{0, \dots, m-1\}$

$$h_{\mathbf{a}}(\mathbf{x}) = \mathbf{a} \cdot \mathbf{x} \pmod{m}$$

Familie für universelles Hashing

Satz

Wenn m eine Primzahl ist, dann ist

$$H = \{h_{\mathbf{a}} : \mathbf{a} \in \{0, \dots, m-1\}^k\}$$

eine **[1]-universelle** Familie von Hashfunktionen.

Oder anders:

das Skalarprodukt zwischen einer Tupeldarstellung des Schlüssels
und einem Zufallsvektor modulo m definiert eine gute Hashfunktion.

Familie für universelles Hashing

Beispiel

- 32-Bit-Schlüssel, Hashtabellengröße $m = 269$

Familie für universelles Hashing

Beispiel

- 32-Bit-Schlüssel, Hashtabellengröße $m = 269$
- ⇒ Schlüssel unterteilt in $k = 4$ Teile mit $w = \lfloor \log_2 m \rfloor = 8$ Bits
- Schlüssel sind also 4-Tupel von Integers aus dem Intervall $[0, 2^8 - 1] = \{0, \dots, 255\}$, z.B. $\mathbf{x} = (11, 7, 4, 3)$

Familie für universelles Hashing

Beispiel

- 32-Bit-Schlüssel, Hashtabellengröße $m = 269$
- ⇒ Schlüssel unterteilt in $k = 4$ Teile mit $w = \lfloor \log_2 m \rfloor = 8$ Bits
- Schlüssel sind also 4-Tupel von Integers aus dem Intervall $[0, 2^8 - 1] = \{0, \dots, 255\}$, z.B. $\mathbf{x} = (11, 7, 4, 3)$
- Die Hashfunktion wird auch durch ein 4-Tupel von Integers, aber aus dem Intervall $[0, 269 - 1] = \{0, \dots, 268\}$, spezifiziert z.B. $\mathbf{a} = (2, 4, 261, 16)$

Familie für universelles Hashing

Beispiel

- 32-Bit-Schlüssel, Hashtabellengröße $m = 269$
- ⇒ Schlüssel unterteilt in $k = 4$ Teile mit $w = \lfloor \log_2 m \rfloor = 8$ Bits
- Schlüssel sind also 4-Tupel von Integers aus dem Intervall $[0, 2^8 - 1] = \{0, \dots, 255\}$, z.B. $\mathbf{x} = (11, 7, 4, 3)$
- Die Hashfunktion wird auch durch ein 4-Tupel von Integers, aber aus dem Intervall $[0, 269 - 1] = \{0, \dots, 268\}$, spezifiziert z.B. $\mathbf{a} = (2, 4, 261, 16)$

⇒ Hashfunktion:

$$h_{\mathbf{a}}(\mathbf{x}) = (2x_1 + 4x_2 + 261x_3 + 16x_4) \bmod 269$$

$$h_{\mathbf{a}}(\mathbf{x}) = (2 \cdot 11 + 4 \cdot 7 + 261 \cdot 4 + 16 \cdot 3) \bmod 269 = 66$$

Eindeutiges a_j

Beweis

- Betrachte zwei beliebige verschiedene Schlüssel $\mathbf{x} = \{x_1, \dots, x_k\}$ und $\mathbf{y} = \{y_1, \dots, y_k\}$
- Wie groß ist $\Pr[h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y})]$?
- Sei j ein Index (von evt. mehreren möglichen) mit $x_j \neq y_j$ (muss es geben, sonst wäre $\mathbf{x} = \mathbf{y}$)
- ⇒ $(x_j - y_j) \not\equiv 0 \pmod m$
d.h., es gibt genau ein multiplikatives Inverses $(x_j - y_j)^{-1}$
- ⇒ gegebene Primzahl m und Zahlen $x_j, y_j, b \in \{0, \dots, m-1\}$ hat jede Gleichung der Form

$$a_j (x_j - y_j) \equiv b \pmod m$$

eine eindeutige Lösung: $a_j \equiv (x_j - y_j)^{-1} b \pmod m$

Wann wird $h(\mathbf{x}) = h(\mathbf{y})$?

Beweis.

Wenn man alle Variablen a_i außer a_j festlegt, gibt es **exakt eine Wahl für a_j** , so dass $h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y})$, denn

$$\begin{aligned}
 h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y}) &\Leftrightarrow \sum_{i=1}^k a_i x_i \equiv \sum_{i=1}^k a_i y_i \pmod{m} \\
 &\Leftrightarrow a_j(x_j - y_j) \equiv \sum_{i \neq j} a_i(y_i - x_i) \pmod{m} \\
 &\Leftrightarrow a_j \equiv (x_j - y_j)^{-1} \sum_{i \neq j} a_i(y_i - x_i) \pmod{m}
 \end{aligned}$$

Wie oft wird $h(\mathbf{x}) = h(\mathbf{y})$?

Beweis.

- Es gibt m^{k-1} Möglichkeiten, Werte für die Variablen a_i mit $i \neq j$ zu wählen.
- Für jede solche Wahl gibt es genau eine Wahl für a_j , so dass $h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y})$.
- Für \mathbf{a} gibt es insgesamt m^k Auswahlmöglichkeiten.
- Also

$$\Pr[h_{\mathbf{a}}(\mathbf{x}) = h_{\mathbf{a}}(\mathbf{y})] = \frac{m^{k-1}}{m^k} = \frac{1}{m}$$



Familie für k -universelles Hashing

Definiere für $a \in \{1, \dots, m-1\}$ die Hashfunktion

$$h'_a(\mathbf{x}) = \sum_{i=1}^k a^{i-1} x_i \pmod{m}$$

(mit $x_i \in \{0, \dots, m-1\}$)

Satz

Für jede Primzahl m ist

$$H' = \{h'_a : a \in \{1, \dots, m-1\}\}$$

eine **k -universelle** Familie von Hashfunktionen.

Familie für k -universelles Hashing

Beweisidee:

Für Schlüssel $\mathbf{x} \neq \mathbf{y}$ ergibt sich folgende Gleichung:

$$\begin{aligned}
 h'_a(\mathbf{x}) &= h'_a(\mathbf{y}) \\
 h'_a(\mathbf{x}) - h'_a(\mathbf{y}) &\equiv 0 \pmod{m} \\
 \sum_{i=1}^k a^{i-1} (x_i - y_i) &\equiv 0 \pmod{m}
 \end{aligned}$$

Anzahl der Nullstellen des Polynoms in a ist durch den Grad des Polynoms beschränkt (Fundamentalsatz der Algebra), also durch $k-1$.

Falls $k \leq m$ können also höchstens $k-1$ von $m-1$ möglichen Werten für a zum gleichen Hashwert für \mathbf{x} und \mathbf{y} führen.

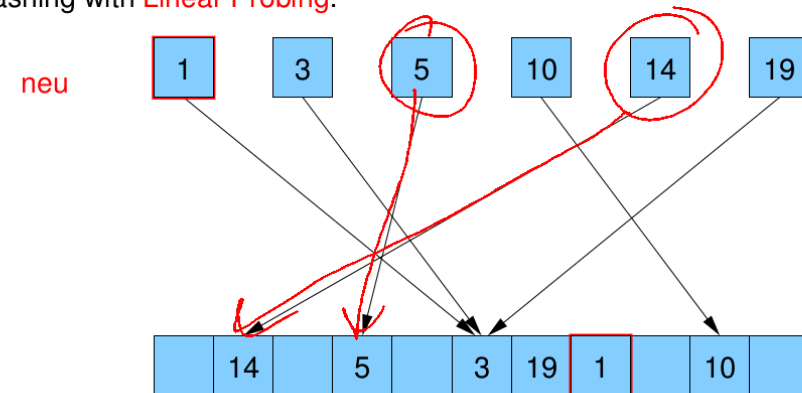
Aus $\Pr[h(\mathbf{x}) = h(\mathbf{y})] \leq \frac{\min\{k-1, m-1\}}{m-1} \leq \frac{k}{m}$ folgt, dass H' k -universell ist.

Übersicht

5 Hashing

- Hashtabellen
- Hashing with Chaining
- Universelles Hashing
- Hashing with Linear Probing
- Anpassung der Tabellengröße
- Perfektes Hashing
- Diskussion / Alternativen

Dynamisches Wörterbuch

Hashing with **Linear Probing**:

Speichere Element e im ersten freien
Ort $T[i]$, $T[i + 1]$, $T[i + 2]$, ... mit $i == h(\text{key}(e))$
(Ziel: Folgen besetzter Positionen möglichst kurz)

Hashing with Linear Probing

Elem[m] T; // Feld sollte genügend groß sein

```
insert(Elem e) {
  i = h(key(e));
  while (T[i] ≠ null ∧ T[i] ≠ e)
    i = (i+1) % m;
  T[i] = e;
}
```

```
find(Key k) {
  i = h(k);
  while (T[i] ≠ null ∧ key(T[i]) ≠ k)
    i = (i+1) % m;
  return T[i];
}
```

Hashing with Linear Probing

Vorteil:

Es werden im Gegensatz zu Hashing with Chaining
(oder auch im Gegensatz zu anderen Probing-Varianten)
nur **zusammenhängende** Speicherzellen betrachtet.

⇒ Cache-Effizienz!

Hashing with Linear Probing

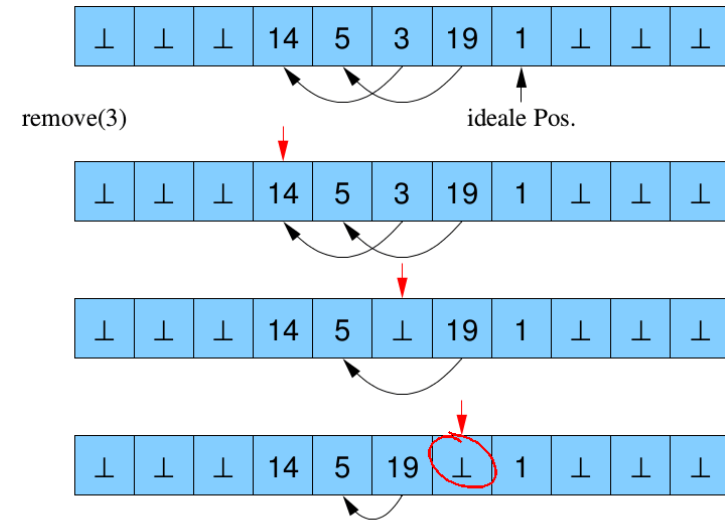
Problem: **Löschen** von Elementen

- 1 Löschen verbieten
 - 2 Markiere Positionen als gelöscht (mit speziellem Zeichen \perp)
Suche endet bei \perp , aber nicht bei markierten Zellen
- Problem: Anzahl echt freier Zellen sinkt monoton
 \Rightarrow Suche wird evt. langsam oder periodische Reorganisation

- 3 Invariante sicherstellen:
Für jedes $e \in S$ mit idealer Position $i = h(key(e))$ und aktueller Position j gilt:
 $T[i], T[i+1], \dots, T[j]$ sind besetzt

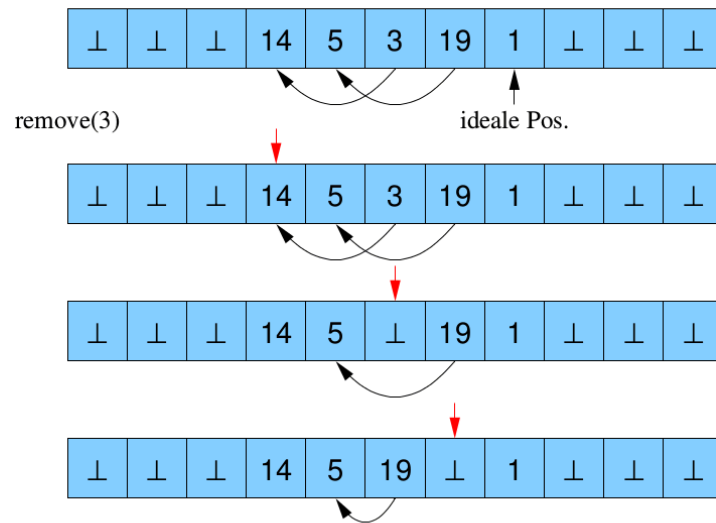
Hashing with Linear Probing

Löschen / Aufrechterhaltung der Invariante



Hashing with Linear Probing

Löschen / Aufrechterhaltung der Invariante



Übersicht

- 5 Hashing
 - Hashtabellen
 - Hashing with Chaining
 - Universelles Hashing
 - Hashing with Linear Probing
 - Anpassung der Tabellengröße
 - Perfektes Hashing
 - Diskussion / Alternativen

Dynamisches Wörterbuch

Problem: Hashtabelle ist **zu groß** oder **zu klein**
 (sollte nur um konstanten Faktor von Anzahl der Elemente abweichen)

Lösung: Reallokation

- Wähle geeignete Tabellengröße
- Wähle neue Hashfunktion
- Übertrage Elemente auf die neue Tabelle

Dynamisches Wörterbuch

Problem: Tabellengröße m sollte **prim** sein
 (für eine gute Verteilung der Schlüssel)

Lösung:

- Für jedes k gibt es eine Primzahl in $[k^3, (k+1)^3]$
- Jede Zahl $z \leq (k+1)^3$, die nicht prim ist, muss einen Teiler $t \leq \sqrt{(k+1)^3} = (k+1)^{3/2}$ haben.
- Für eine gewünschte ungefähre Tabellengröße m' (evt. nicht prim) bestimme k so, dass $k^3 \leq m' \leq (k+1)^3$

Dynamisches Wörterbuch

Problem: Tabellengröße m sollte **prim** sein
 (für eine gute Verteilung der Schlüssel)

Lösung:

- Für jedes k gibt es eine Primzahl in $[k^3, (k+1)^3]$
- Jede Zahl $z \leq (k+1)^3$, die nicht prim ist, muss einen Teiler $t \leq \sqrt{(k+1)^3} = (k+1)^{3/2}$ haben.
- Für eine gewünschte ungefähre Tabellengröße m' (evt. nicht prim) bestimme k so, dass $k^3 \leq m' \leq (k+1)^3$
- Größe des Intervalls:
 $(k+1)^3 - k^3 + 1 = \cancel{k^3} + 3k^2 + 3k + 1 - \cancel{k^3} + 1 = 3k^2 + 3k + 2$
- Für jede Zahl $j = 2, \dots, (k+1)^{3/2}$:
 streiche die Vielfachen von j in $[k^3, (k+1)^3]$
- Für jedes j kostet das Zeit $((k+1)^3 - k^3 + 1) / j \in O(k^2 / j)$

Dynamisches Wörterbuch

- Hilfsmittel: Wachstum der harmonischen Reihe

$$\ln n \leq H_n = \sum_{i=1}^n \frac{1}{i} \leq 1 + \ln n$$

- insgesamt:

$$\begin{aligned} \sum_{2 \leq j \leq (k+1)^{3/2}} O\left(\frac{k^2}{j}\right) &\leq k^2 \sum_{2 \leq j \leq (k+1)^{3/2}} O\left(\frac{1}{j}\right) \\ &\in k^2 \cdot O\left(\ln\left((k+1)^{3/2}\right)\right) \\ &\in O(k^2 \ln k) \subseteq \Theta(k^3) \\ &\in o(m) \end{aligned}$$

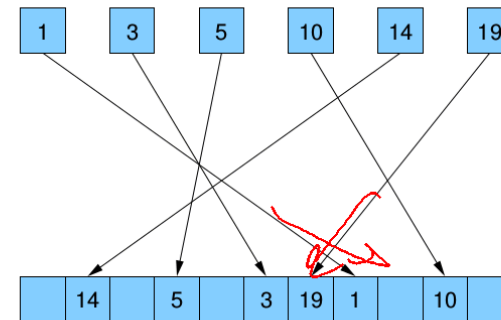
⇒ Kosten zu vernachlässigen im Vergleich zur Initialisierung der Tabelle der Größe m (denn m ist kubisch in k)

Übersicht

- 5 Hashing
 - Hashtabellen
 - Hashing with Chaining
 - Universelles Hashing
 - Hashing with Linear Probing
 - Anpassung der Tabellengröße
 - Perfektes Hashing
 - Diskussion / Alternativen

Perfektes Hashing für statisches Wörterbuch

- bisher: konstante *erwartete* Laufzeit
falls m im Vergleich zu n genügend groß gewählt wird
(nicht ausreichend für Real Time Szenario)
- Ziel: konstante Laufzeit im **worst case** für find()
durch perfekte Hashtabelle ohne Kollisionen
- Annahme: statische Menge S von n Elementen



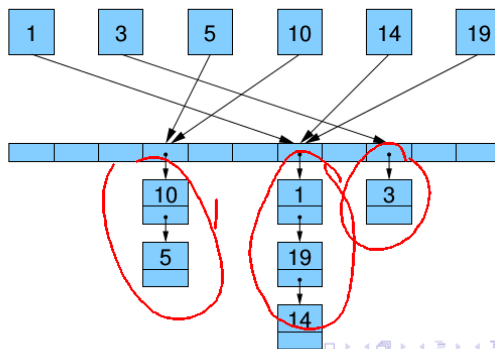
Statisches Wörterbuch

- S : feste Menge von n Elementen mit Schlüssel k_1 bis k_n
- H_m : c -universelle Familie von Hashfunktionen auf $\{0, \dots, m-1\}$
(Hinweis: 2-universelle Familien existieren für alle m)
- $C(h)$ für $h \in H_m$: Anzahl Kollisionen in S für h , d.h.

$$C(h) = |\{(x, y) : x, y \in S, x \neq y, h(x) = h(y)\}|$$

Beispiel:

$$C(h) = 2 + 6 = 8$$



Erwartete Anzahl von Kollisionen

Lemma

Für die Anzahl der Kollisionen gilt: $\mathbb{E}[C(h)] \leq cn(n-1)/m$.

Erwartete Anzahl von Kollisionen

Lemma

Für die Anzahl der Kollisionen gilt: $\mathbb{E}[C(h)] \leq cn(n-1)/m$.

Beweis.

- Definiere $n(n-1)$ Indikator-Zufallsvariablen $X_{ij}(h)$:
Für $i \neq j$ sei $X_{ij}(h) = 1 \Leftrightarrow h(k_i) = h(k_j)$.
- Dann ist $C(h) = \sum_{i \neq j} X_{ij}(h)$

$$\mathbb{E}[C] = \mathbb{E}\left[\sum_{i \neq j} X_{ij}\right] = \sum_{i \neq j} \mathbb{E}[X_{ij}] = \sum_{i \neq j} \Pr[X_{ij} = 1] \leq n(n-1) \cdot c/m$$

⇒ Für **quadratische Tabellengröße** ist die erwartete Anzahl von Kollisionen (und damit die erwartete worst-case-Laufzeit für find) eine **Konstante**. □

Markov-Ungleichung

Satz (Markov-Ungleichung)

Für jede nichtnegative Zufallsvariable X und Konstante $k > 0$ gilt:

$$\Pr[X \geq k] \leq \frac{\mathbb{E}[X]}{k} \quad \text{und für } \mathbb{E}[X] > 0: \Pr[X \geq k \cdot \mathbb{E}[X]] \leq \frac{1}{k}$$

Erwartete Anzahl von Kollisionen

Lemma

Für die Anzahl der Kollisionen gilt: $\mathbb{E}[C(h)] \leq cn(n-1)/m$.

Beweis.

- Definiere $n(n-1)$ Indikator-Zufallsvariablen $X_{ij}(h)$:
Für $i \neq j$ sei $X_{ij}(h) = 1 \Leftrightarrow h(k_i) = h(k_j)$.
- Dann ist $C(h) = \sum_{i \neq j} X_{ij}(h)$

$$\mathbb{E}[C] = \mathbb{E}\left[\sum_{i \neq j} X_{ij}\right] = \sum_{i \neq j} \mathbb{E}[X_{ij}] = \sum_{i \neq j} \Pr[X_{ij} = 1] \leq n(n-1) \cdot c/m$$

⇒ Für **quadratische Tabellengröße** ist die erwartete Anzahl von Kollisionen (und damit die erwartete worst-case-Laufzeit für find) eine **Konstante**. □

Markov-Ungleichung

Satz (Markov-Ungleichung)

Für jede nichtnegative Zufallsvariable X und Konstante $k > 0$ gilt:

$$\Pr[X \geq k] \leq \frac{\mathbb{E}[X]}{k} \quad \text{und für } \mathbb{E}[X] > 0: \Pr[X \geq k \cdot \mathbb{E}[X]] \leq \frac{1}{k}$$

Beweis.

$$\begin{aligned} \mathbb{E}[X] &= \sum_{z \in \mathbb{R}} z \cdot \Pr[X = z] \\ &\geq \sum_{z \geq k \cdot \mathbb{E}[X]} z \cdot \Pr[X = z] \geq \sum_{z \geq k \cdot \mathbb{E}[X]} k \cdot \mathbb{E}[X] \cdot \Pr[X = z] \\ &\geq k \cdot \mathbb{E}[X] \cdot \Pr[X \geq k \cdot \mathbb{E}[X]] \end{aligned}$$

Hashfunktionen mit wenig Kollisionen

Lemma

Für mindestens **die Hälfte** der Funktionen $h \in H_m$ gilt:

$$C(h) \leq 2cn(n-1)/m$$

Erwartete Anzahl von Kollisionen

Lemma

Für die Anzahl der Kollisionen gilt: $\mathbb{E}[C(h)] \leq cn(n-1)/m$

Beweis.

- Definiere $n(n-1)$ Indikator-Zufallsvariablen $X_{ij}(h)$:
Für $i \neq j$ sei $X_{ij}(h) = 1 \Leftrightarrow h(k_i) = h(k_j)$.
- Dann ist $C(h) = \sum_{i \neq j} X_{ij}(h)$

$$\mathbb{E}[C] = \mathbb{E}\left[\sum_{i \neq j} X_{ij}\right] = \sum_{i \neq j} \mathbb{E}[X_{ij}] = \sum_{i \neq j} \Pr[X_{ij} = 1] \leq n(n-1) \cdot c/m$$

\Rightarrow Für **quadratische Tabellengröße** ist die erwartete Anzahl von Kollisionen (und damit die erwartete worst-case-Laufzeit für find) eine **Konstante**. \square

Hashfunktionen mit wenig Kollisionen

Lemma

Für mindestens **die Hälfte** der Funktionen $h \in H_m$ gilt:

$$C(h) \leq 2cn(n-1)/m$$

Beweis.

- Aus Lemma $\mathbb{E}[C(h)] \leq cn(n-1)/m$ und Markov-Ungleichung $\Pr[X \geq k \cdot \mathbb{E}[X]] \leq \frac{1}{k}$ folgt für $n \geq 2$:

$$\Pr[C(h) \geq 2cn(n-1)/m] \leq \Pr[C(h) \geq 2\mathbb{E}[C(h)]] \leq \frac{1}{2}$$

\Rightarrow Für höchstens die Hälfte der Funktionen ist $C(h) \geq \frac{2cn(n-1)}{m}$

\Rightarrow Für mindestens die Hälfte der Funktionen ist $C(h) \leq \frac{2cn(n-1)}{m}$ ($n \geq 1$) \square

Hashfunktionen ohne Kollisionen

Lemma

Wenn $m \geq cn(n-1) + 1$, dann bildet mindestens die Hälfte der Funktionen $h \in H_m$ die Schlüssel **injektiv** in die Indexmenge der Hashtabelle ab.

Beweis.

- Für mindestens die Hälfte der Funktionen $h \in H_m$ gilt $C(h) < 2$.
 - Da $C(h)$ immer eine gerade Zahl sein muss, folgt aus $C(h) < 2$ direkt $C(h) = 0$.
- \Rightarrow keine Kollisionen (bzw. injektive Abbildung) \square

- Wähle zufällig $h \in H_m$ mit $m \geq cn(n-1) + 1$
 - Prüfe auf Kollision \Rightarrow behalten oder erneut wählen
- \Rightarrow Nach durchschnittlich 2 Versuchen erfolgreich

Hashfunktionen mit wenig Kollisionen

Lemma

Für mindestens **die Hälfte** der Funktionen $h \in H_m$ gilt:

$$C(h) \leq 2cn(n-1)/m$$

Beweis.

- Aus Lemma $\mathbb{E}[C(h)] \leq cn(n-1)/m$ und Markov-Ungleichung $\Pr[X \geq k \cdot \mathbb{E}[X]] \leq \frac{1}{k}$ folgt für $n \geq 2$:

$$\Pr[C(h) \geq 2cn(n-1)/m] \leq \Pr[C(h) \geq 2\mathbb{E}[C(h)]] \leq \frac{1}{2}$$

- ⇒ Für höchstens die Hälfte der Funktionen ist $C(h) \geq \frac{2cn(n-1)}{m}$
- ⇒ Für mindestens die Hälfte der Funktionen ist $C(h) \leq \frac{2cn(n-1)}{m}$ ($n \geq 1$)

□

Hashfunktionen ohne Kollisionen

Lemma

Wenn $m \geq cn(n-1) + 1$, dann bildet mindestens die Hälfte der Funktionen $h \in H_m$ die Schlüssel **injektiv** in die Indexmenge der Hashtabelle ab.

Beweis.

- Für mindestens die Hälfte der Funktionen $h \in H_m$ gilt $C(h) < 2$.
- Da $C(h)$ immer eine gerade Zahl sein muss, folgt aus $C(h) < 2$ direkt $C(h) = 0$.

⇒ keine Kollisionen (bzw. injektive Abbildung)

□

- Wähle zufällig $h \in H_m$ mit $m \geq cn(n-1) + 1$
 - Prüfe auf Kollision ⇒ behalten oder erneut wählen
- ⇒ Nach durchschnittlich 2 Versuchen erfolgreich

Hashfunktionen ohne Kollisionen

Lemma

Wenn $m \geq cn(n-1) + 1$, dann bildet mindestens die Hälfte der Funktionen $h \in H_m$ die Schlüssel **injektiv** in die Indexmenge der Hashtabelle ab.

Beweis.

- Für mindestens die Hälfte der Funktionen $h \in H_m$ gilt $C(h) < 2$.
- Da $C(h)$ immer eine gerade Zahl sein muss, folgt aus $C(h) < 2$ direkt $C(h) = 0$.

⇒ keine Kollisionen (bzw. injektive Abbildung)

□

- Wähle zufällig $h \in H_m$ mit $m \geq cn(n-1) + 1$
 - Prüfe auf Kollision ⇒ behalten oder erneut wählen
- ⇒ Nach durchschnittlich 2 Versuchen erfolgreich

Statisches Wörterbuch

Ziel: lineare Tabellengröße

Idee: **zweistufige** Abbildung der Schlüssel

- Wähle Hashfunktion h mit wenig Kollisionen (≈ 2 Versuche)
- ⇒ 1. Stufe bildet Schlüssel auf Buckets von konstanter durchschnittlicher Größe ab
- Wähle Hashfunktionen h_ℓ ohne Kollisionen (≈ 2 Versuche pro h_ℓ)
- ⇒ 2. Stufe benutzt quadratisch viel Platz für jedes Bucket, um alle Kollisionen aus der 1. Stufe aufzulösen

Statisches Wörterbuch

- B_ℓ^h : Menge der Elemente in S , die h auf ℓ abbildet, $\ell \in \{0, \dots, m-1\}$ und $h \in H_m$
- b_ℓ^h : Kardinalität von B_ℓ^h , also $b_\ell^h := |B_\ell^h|$
- Für jedes ℓ führen die Schlüssel in B_ℓ^h zu $b_\ell^h(b_\ell^h - 1)$ Kollisionen
- Also ist die Gesamtzahl der Kollisionen

$$C(h) = \sum_{\ell=0}^{m-1} b_\ell^h(b_\ell^h - 1)$$

Perfektes statisches Hashing: 1. Stufe

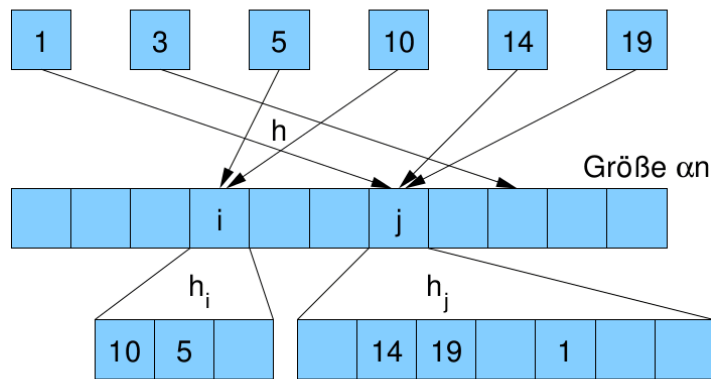
- 1. Stufe der Hashtabelle soll **linear** viel Speicher verwenden, also $\lceil \alpha n \rceil$ Adressen, wobei wir Konstante α später festlegen.
- Wähle Funktion $h \in H_{\lceil \alpha n \rceil}$, um S in Teilmengen B_ℓ aufzuspalten. Wähle h dabei so lange zufällig aus $H_{\lceil \alpha n \rceil}$ aus, bis gilt:

$$C(h) \leq \frac{2cn(n-1)}{\lceil \alpha n \rceil} \leq \frac{2cn(n-1)}{\alpha n} \leq \frac{2cn}{\alpha}$$

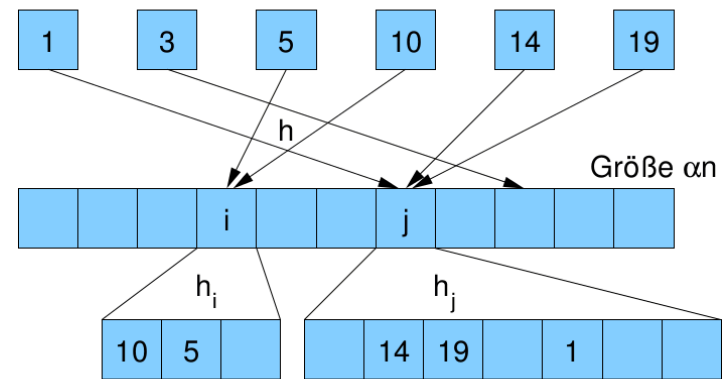
Da das für mindestens die Hälfte der Hashfunktionen in $H_{\lceil \alpha n \rceil}$ gilt (vorletztes Lemma), erwarten wir dafür ≤ 2 Versuche.

- Für jedes $\ell \in \{0, \dots, \lceil \alpha n \rceil - 1\}$ seien B_ℓ die Elemente, die durch h auf Adresse ℓ abgebildet werden und $b_\ell = |B_\ell|$ deren Anzahl.

Perfektes statisches Hashing



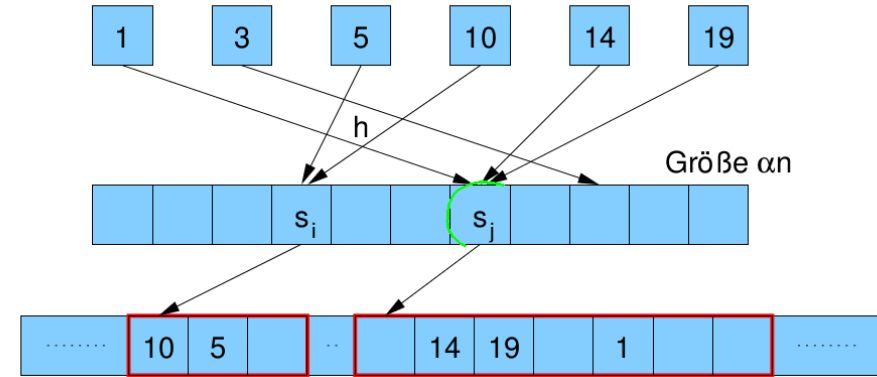
Perfektes statisches Hashing



Perfektes statisches Hashing: 2. Stufe

- Für jedes B_ℓ :
 Berechne $m_\ell = cb_\ell(b_\ell - 1) + 1$.
 Wähle zufällig Funktion $h_\ell \in H_{m_\ell}$, bis h_ℓ die Menge B_ℓ injektiv in $\{0, \dots, m_\ell - 1\}$ abbildet (also ohne Kollisionen).
 Mindestens die Hälfte der Funktionen in H_{m_ℓ} tut das.
- Hintereinanderreihung der einzelnen Tabellen ergibt eine Gesamtgröße der Tabelle von $\sum_\ell m_\ell$
- Teiltabelle für B_ℓ beginnt an Position $s_\ell = m_0 + m_1 + \dots + m_{\ell-1}$ und endet an Position $s_\ell + m_\ell - 1$
- Für gegebenen Schlüssel x , berechnen die Anweisungen
 $\ell = h(x); \text{ return } s_\ell + h_\ell(x);$
 dann eine injektive Funktion auf der Menge S .

Perfektes statisches Hashing



Perfektes statisches Hashing

- Die Funktion ist beschränkt durch:

$$\sum_{\ell=0}^{\lceil \alpha n \rceil - 1} m_\ell = \sum_{\ell=0}^{\lceil \alpha n \rceil - 1} (c \cdot b_\ell (b_\ell - 1) + 1) \quad (\text{siehe Def. der } m_\ell\text{'s})$$

$$\leq c \cdot C(h) + \lceil \alpha n \rceil$$

$$\leq c \cdot 2cn/\alpha + \alpha n + 1$$

$$\leq (2c^2/\alpha + \alpha)n + 1$$

- Zur Minimierung der Schranke betrachte die Ableitung

$$f(\alpha) = (2c^2/\alpha + \alpha)n + 1$$

$$f'(\alpha) = (-2c^2/\alpha^2 + 1)n$$

$\Rightarrow f'(\alpha) = 0$ liefert $\alpha = \sqrt{2c}$
 \Rightarrow Adressbereich: $0 \dots 2\sqrt{2}cn$

Perfektes statisches Hashing

- Die Funktion ist beschränkt durch:

$$\sum_{\ell=0}^{\lceil \alpha n \rceil - 1} m_\ell = \sum_{\ell=0}^{\lceil \alpha n \rceil - 1} (c \cdot b_\ell (b_\ell - 1) + 1) \quad (\text{siehe Def. der } m_\ell\text{'s})$$

$$\leq c \cdot C(h) + \lceil \alpha n \rceil$$

$$\leq c \cdot 2cn/\alpha + \alpha n + 1$$

$$\leq (2c^2/\alpha + \alpha)n + 1 = (\sqrt{2} \cdot c + \sqrt{2} \cdot c) \alpha n + 1$$

- Zur Minimierung der Schranke betrachte die Ableitung

$$f(\alpha) = (2c^2/\alpha + \alpha)n + 1$$

$$f'(\alpha) = (-2c^2/\alpha^2 + 1)n$$

$\Rightarrow f'(\alpha) = 0$ liefert $\alpha = \sqrt{2c}$
 \Rightarrow Adressbereich: $0 \dots 2\sqrt{2}cn$

Perfektes statisches Hashing

Satz

Für eine beliebige Menge von n Schlüsseln kann eine perfekte Hashfunktion mit Zielmenge $\{0, \dots, 2\sqrt{2}cn\}$ in linearer erwarteter Laufzeit konstruiert werden.

- Da wir wissen, dass wir für beliebige m eine 2-universelle Familie von Hashfunktionen finden können, kann man also z.B. eine Hashfunktion mit Adressmenge $\{0, \dots, 4\sqrt{2}n\}$ in linearer erwarteter Laufzeit konstruieren. (Las Vegas-Algorithmus)
 - Unsere Minimierung hat nicht den Platz berücksichtigt, der benötigt wird, um die Werte s_ℓ , sowie die ausgewählten Hashfunktionen h_ℓ zu speichern.
- ⇒ Berechnung von α sollte eigentlich angepasst werden (entsprechend Speicherplatz pro Element, pro m_ℓ und pro h_ℓ)

Perfektes dynamisches Hashing

Kann man perfekte Hashfunktionen auch **dynamisch** konstruieren?

ja, z.B. mit **Cuckoo** Hashing

- 2 Hashfunktionen h_1 und h_2
- 2 Hashtabellen T_1 und T_2
- bei find und remove jeweils in beiden Tabellen nachschauen
- bei insert abwechselnd beide Tabellen betrachten, das zu speichernde Element an die Zielposition der aktuellen Tabelle schreiben und wenn dort schon ein anderes Element stand, dieses genauso in die andere Tabelle verschieben usw.
- evt. Anzahl Verschiebungen durch $2 \log n$ beschränken, um Endlosschleife zu verhindern (ggf. kompletter Rehash mit neuen Funktionen h_1, h_2)

Probleme beim linearen Sondieren

- Offene Hashverfahren allgemein:
Erweiterte Hashfunktion $h(k, i)$ gibt an, auf welche Adresse ein Schlüssel k abgebildet werden soll, wenn bereits i Versuche zu einer Kollision geführt haben
- Lineares Sondieren (Linear Probing):

$$h(k, i) = (h(k) + i) \bmod m$$

Probleme beim linearen Sondieren

- Offene Hashverfahren allgemein:
Erweiterte Hashfunktion $h(k, i)$ gibt an, auf welche Adresse ein Schlüssel k abgebildet werden soll, wenn bereits i Versuche zu einer Kollision geführt haben
 - Lineares Sondieren (Linear Probing):
- $$h(k, i) = (h(k) + i) \bmod m$$
- **Primäre Häufung** (primary clustering):
tritt auf, wenn für Schlüssel k_1, k_2 mit unterschiedlichen Hashwerten $h(k_1) \neq h(k_2)$ ab einem bestimmten Punkt i_1 bzw. i_2 die gleiche Sondierfolge auftritt:

$$\exists i_1, i_2 \quad \forall j: \quad h(k_1, i_1 + j) = h(k_2, i_2 + j)$$

Probleme beim quadratischen Sondieren

- Quadratisches Sondieren (Quadratic Probing):

$$h(k, i) = (h(k) + c_1i + c_2i^2) \bmod m \quad (c_2 \neq 0)$$

oder: $h(k, i) = (h(k) - (-1)^i \lceil i/2 \rceil^2) \bmod m$

Probleme beim quadratischen Sondieren

- Quadratisches Sondieren (Quadratic Probing):

$$h(k, i) = (h(k) + c_1i + c_2i^2) \bmod m \quad (c_2 \neq 0)$$

oder: $h(k, i) = (h(k) - (-1)^i \lceil i/2 \rceil^2) \bmod m$

- $h(k, i)$ soll möglichst **surjektiv** auf die Adressmenge $\{0, \dots, m - 1\}$ abbilden, um freie Positionen auch **immer** zu finden. Bei

$$h(k, i) = (h(k) - (-1)^i \lceil i/2 \rceil^2) \bmod m$$

z.B. durch Wahl von $m \text{ prim} \wedge m \equiv 3 \pmod 4$

Probleme beim quadratischen Sondieren

- Quadratisches Sondieren (Quadratic Probing):

$$h(k, i) = (h(k) + c_1i + c_2i^2) \bmod m \quad (c_2 \neq 0)$$

oder: $h(k, i) = (h(k) - (-1)^i \lceil i/2 \rceil^2) \bmod m$

- $h(k, i)$ soll möglichst **surjektiv** auf die Adressmenge $\{0, \dots, m - 1\}$ abbilden, um freie Positionen auch **immer** zu finden. Bei

$$h(k, i) = (h(k) - (-1)^i \lceil i/2 \rceil^2) \bmod m$$

z.B. durch Wahl von $m \text{ prim} \wedge m \equiv 3 \pmod 4$

- **Sekundäre Häufung** (secondary clustering): tritt auf, wenn für Schlüssel k_1, k_2 mit gleichem Hashwert $h(k_1) = h(k_2)$ auch die nachfolgende Sondierfolge gleich ist:

$$\forall i: h(k_1, i) = h(k_2, i)$$

Double Hashing

- Auflösung der Kollisionen der Hashfunktion h durch eine zweite Hashfunktion h' :

$$h(k, i) = [h(k) + i \cdot h'(k)] \bmod m$$

wobei für alle k gelten soll, dass $h'(k)$ teilerfremd zu m ist,

z.B. $h'(k) = 1 + k \bmod m - 1$
 oder $h'(k) = 1 + k \bmod m - 2$

für Primzahl m

- primäre und sekundäre Häufung werden weitgehend vermieden, aber nicht komplett ausgeschlossen