

# 14. Hashing

---

Hashtabellen, Pre-Hashing, Hashing, Kollisionsauflösung durch Verketteten, 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<sup>18</sup>  $i$  mit Schlüsseln  $k \in \mathcal{K}$ . Operationen

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

---

<sup>18</sup>Schlüssel-Wert Paare  $(k, v)$ , im Folgenden betrachten wir hauptsächlich die Schlüssel.

# Wörterbuch in C++

Assoziativer Container `std::unordered_map<>`

```
// Create an unordered_map of strings that map to strings
std::unordered_map<std::string, std::string> u = {
    {"RED", "#FF0000"}, {"GREEN", "#00FF00"}
};

u["BLUE"] = "#0000FF"; // Add

std::cout << "The HEX of color RED is: " << u["RED"] << "\n";

for( const auto& n : u ) // iterate over key-value pairs
    std::cout << n.first << ":" << n.second << "\n";
```

# 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

# 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  $\Rightarrow$  grosses Array

# Lösung zum ersten Problem: Pre-hashing

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

- 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.<sup>19</sup>

---

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

# Prehashing Beispiel: String

Zuordnung Name  $s = s_1s_2 \dots s_{l_s}$  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

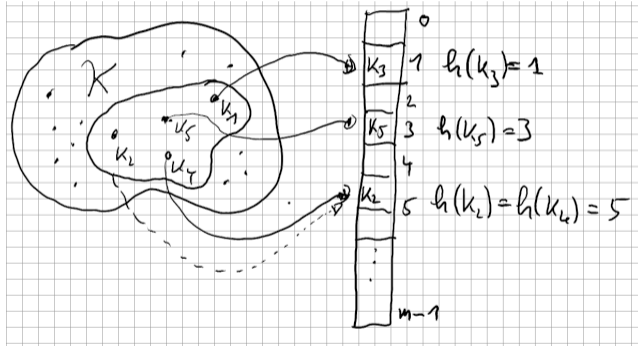
Jacqueline  $\mapsto$  2042089953442505  $\bmod 2^{32} = 507919049$



# Lösung zum zweiten Problem: Hashing

Reduziere des Schlüsseluniversum: Abbildung (Hash-Funktion)

$h : \mathcal{K} \rightarrow \{0, \dots, m - 1\}$  ( $m \approx n =$  Anzahl Einträge in der Tabelle)



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

# Nomenklatur

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

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

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

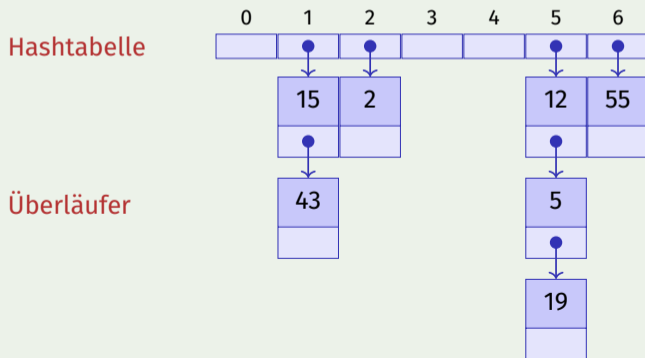
Eine Hashfunktion sollte die Menge der Schlüssel möglichst gleichmässig auf die Positionen der Hashtabelle verteilen.

# Behandlung von Kollisionen: Verkettung

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

Schlüssel 12, 55, 5, 15, 2, 19, 43

Direkte Verkettung der Überläufer



# 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`.
- `delete( $k$ )` Durchsuche die Liste an Position  $h(k)$  nach  $k$ . Wenn Suche erfolgreich, entferne das entsprechende Listenelement.

# Worst-case Analyse

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

⇒  $\Theta(n)$  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

$$\begin{aligned}\mathbb{E}(\text{Länge Kette } j) &= \mathbb{E}\left(\sum_{i=0}^{n-1} \mathbb{1}(h(k_i) = j)\right) = \sum_{i=0}^{n-1} \mathbb{P}(h(k_i) = j) \\ &= \sum_{i=1}^n \frac{1}{m} = \frac{n}{m}\end{aligned}$$

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

# 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. Erfolgreiche 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 Einfügeschicht: Schlüssel  $j$  sieht durchschnittliche Listenlänge  $(j - 1)/m$ .

⇒ 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

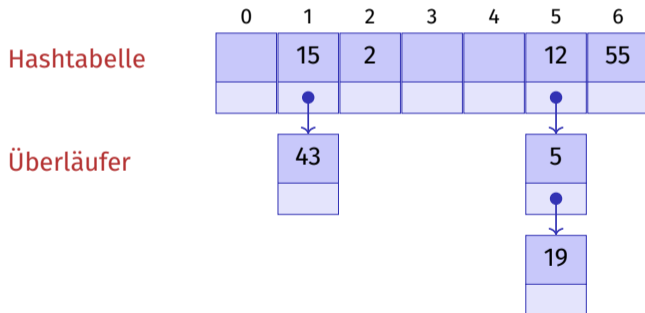
- Speicherverbrauch der Verkettung

# [Variante:Indirekte Verkettung]

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

Schlüssel 12, 55, 5, 15, 2, 19, 43

Indirekte Verkettung der Überläufer



# Beispiele gebräuchlicher Hashfunktionen

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

Ideal:  $m$  Primzahl, nicht zu nahe bei Potenzen von 2 oder 10

Aber oft:  $m = 2^k - 1$  ( $k \in \mathbb{N}$ )

# 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.
- Tabellengröße  $m = 2^r$ ,  $w =$  Grösse des Maschinenworts in Bits.
- Multiplikation addiert  $k$  entlang aller Bits von  $a$ , Ganzzahldivision durch  $2^{w-r}$  und  $\bmod m$  extrahieren die oberen  $r$  Bits.
- Als Code geschrieben sehr einfach: `a * k >> (w-r)`

# Illustration

$$\begin{array}{r} \leftarrow w \text{ bits} \rightarrow \\ \times \begin{array}{|c|c|} \hline k \\ \hline 11 & 1 \\ \hline \end{array} \begin{array}{l} k \\ a \end{array} \end{array}$$

---

$$\begin{array}{r} \begin{array}{|c|} \hline k \\ \hline \end{array} \\ + \begin{array}{|c|} \hline k \\ \hline \end{array} \\ + \begin{array}{|c|} \hline k \\ \hline \end{array} \\ = \begin{array}{|c|c|} \hline \phantom{0} & \leftarrow r \text{ bits} \rightarrow \\ \hline \end{array} \end{array}$$

---

$$\gg (w - r) \begin{array}{|c|c|} \hline 0 & \leftarrow r \text{ bits} \rightarrow \\ \hline \end{array}$$

# Tabellenvergrößerung

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

Grösse der Tabelle muss angepasst werden.

Hash-Funktion ändert sich  $\Rightarrow$  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  $m'$ ?

# Tabellenvergrößerung

In Abhängigkeit vom Füllgrad Tabellengröße jeweils verdoppeln.

⇒ Amortisierte Analyse.

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



# Offene Adressierung

Speichere die Überläufer direkt in der Hashtabelle mit einer Sondierfunktion  $s : \mathcal{K} \times \{0, 1, \dots, m - 1\} \rightarrow \{0, 1, \dots, m - 1\}$   
Tabellenposition des Schlüssels entlang der Sondierfolge

$$S(k) := (s(k, 0), s(k, 1), \dots, s(k, m - 1)) \pmod{m}$$

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

---

Begriffsklärung: Dieses Verfahren nutzt offene Adressierung (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

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

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

Schlüssel 12, 55, 5, 15, 2, 19

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

# [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 Hashadressen haben ähnliche Sondierungsfolgen  $\Rightarrow$  lange zusammenhängende belegte Bereiche.

# Quadratisches Sondieren

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

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

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

Schlüssel 12, 55, 5, 15, 2, 19

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

# [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

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

# 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  $h'(k)$ .  $s(k, j) = h(k) + j \cdot h'(k)$ .

$S(k) = (h(k), h(k) + h'(k), h(k) + 2h'(k), \dots, h(k) + (m - 1)h'(k)) \pmod m$

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

Schlüssel 12, 55, 5, 15, 2, 19

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

# 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.
- $h'$  sollte möglichst unabhängig von  $h$  sein (Vermeidung sekundärer Häufung).

## Unabhängigkeit:

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

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

# [Analyse Double Hashing]

Sind  $h$  und  $h'$  unabhängig, dann:

1. Erfolgreiche Suche. Durchschnittliche Anzahl betrachteter Einträge

$$C'_n \approx \frac{1}{1 - \alpha}$$

2. Erfolgreiche Suche. Durchschnittliche Anzahl betrachteter Einträge

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

# Gleichmässiges Hashing

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

(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.

$$\begin{aligned}\mathbb{P}(X \geq i) &\stackrel{*}{=} \frac{n}{m} \cdot \frac{n-1}{m-1} \cdot \frac{n-2}{m-2} \cdots \frac{n-i+2}{m-i+2} \\ &\stackrel{**}{\leq} \left(\frac{n}{m}\right)^{i-1} = \alpha^{i-1}. \quad (1 \leq i \leq m)\end{aligned}$$

**\*** : Ereignis  $A_j$ : Slot beim  $j$ -ten Schritt belegt.

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

**\*\*** :  $\frac{n-1}{m-1} < \frac{n}{m}$  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

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

# [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

# 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

$\alpha$ : Belegungsgrad.

$C_n$ : Anzahl Schritte erfolgreiche Suche,

$C'_n$ : Anzahl Schritte erfolglose Suche



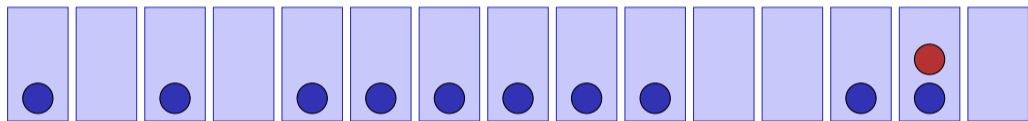
## 14.7 Anhang

---

Mathematische Formeln

# [Geburtstagsparadoxon]

Annahme:  $m$  Urnen,  $n$  Kugeln (oBdA  $n \leq m$ ).  
 $n$  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]

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

Sei  $a \ll m$ . Mit  $e^x = 1 + x + \frac{x^2}{2!} + \dots$  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

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

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

# [Erwartungswertformel]

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

$$\begin{aligned}\mathbb{E}(X) &\stackrel{(def)}{=} \sum_{x=0}^{\infty} x\mathbb{P}(X = x) \\ &\stackrel{\text{Aufzählen}}{=} \sum_{x=1}^{\infty} \sum_{y=x}^{\infty} \mathbb{P}(X = y) \\ &= \sum_{x=0}^{\infty} \mathbb{P}(X > x)\end{aligned}$$