

# Approximative Speicherung von Daten: "Sketching"

## Bloom-Filter

Ziel: Eine große Menge  $S \subseteq U$  kompakt repräsentieren, so dass bei Eingabe eines Elements  $u \in U$  schnell getestet werden kann, ob das Element  $u$

- garantiert nicht zu  $S$  gehört
- oder
- möglicherweise zu  $S$  gehört

Hierbei wollen wir natürlich möglichst wenige "falsche positive Antworten" (d.h. Antworten der Form "gehört möglicherweise zu  $S$ ", obwohl das betrachtete Element tatsächlich nicht zu  $S$  gehört) bekommen.

### Anwendungsbeispiel:

$S$  ist eine "white list" von Email-Adressen, deren Emails der Spam-Filter an die Adressaten anliefern soll.

Szenario: Cans Leskovec, Rajaraman, Ullman "Mining of Massive Datasets" 2014)

- $S$  enthält  $10^9$  Email-Adressen
  - zum Speichern einer einzelnen Email-Adresse werden ca 20 Bytes benötigt
- $S$  kann nicht im Hauptspeicher abgelegt werden

(60)

Aber: Bei Ankunft einer Email soll schnell (ohne Zugriffe auf externen Speicher) entschieden werden, ob die Mail vom Spam-Filter aussortiert oder ob sie an den Adressaten weitergeleitet werden soll.

### Lösung:

- Lege im Hauptspeicher eine geeignete "Skizze" (engl: "Sketch") von  $S$  ab
- Nimm in Kauf, dass es "falsche positive Entscheidungen" gibt — aber Sorge dafür, dass das möglichst selten vorkommt.
- Gestalte die "Skizze" so, dass es garantiert nicht zu "falschen negativen Entscheidungen" kommt (dh: dass Mails von Adressen, die auf der "white list"  $S$  stehen, stets an den Adressaten weitergeleitet werden).

Genau dies leisten Bloom-Filter!

# Der Bloom-Filter

(61)

- Szenario:
- Im Hauptspeicher steht ein Bit-Array  $A[1..n]$  der Länge  $n$  zur Verfügung (für eine feste Zahl  $n$ ).
  - Wie wollen eine "Skizze" einer Menge  $S = \{s_1, \dots, s_N\} \subseteq U$  speichern.

## Speicherung der "Skizze":

- Das Bit-Array wird auf 0 gesetzt:  
 $A[j] := 0$  für alle  $j \in \{1, \dots, n\}$
- Wähle zufällig und unabhängig voneinander  $k$  Hash-Funktionen  $h_1, \dots, h_k$ , wobei jedes  $h_i$  eine Funktion von  $U$  nach  $\{1, \dots, n\}$  ist. Die Zahl  $k$  wird dazu geeignet gewählt (siehe unten).
- Für jedes  $s \in S$  tue Folgendes:  
Für jedes  $i \in \{1, \dots, k\}$   
berechne  $h_i(s)$  und  
setze  $A[h_i(s)] := 1$ .

## Abfragen der "Skizze":

Bei Eingabe eines  $u \in U$ :

- berechne  $h_1(u), \dots, h_k(u)$
- Teste, ob  $A[h_i(u)] = 1$  für alle  $i \in \{1, \dots, k\}$  gilt.
- Falls "ja": gib aus "u gehört möglicherweise zu S"  
sonst: gib aus "u gehört nicht zu S".

Man sieht leicht, dass die Aussage korrekt ist. (62)

Frage: Wie sollen wir  $k$  wählen?

Ansichtlich scheint klar zu sein:

- Je kleiner  $k$ , desto weniger Array-Einträge sind auf 1 gesetzt, und desto geringer scheint die Wahrscheinlichkeit dafür zu sein, dass für ein  $u \notin S$  die Array-Einträge an allen Positionen  $h_1(u), \dots, h_k(u)$  auf 1 gesetzt sind (und daher eine "falsche positive Antwort" gegeben wird).
- Je größer  $k$ , desto geringer scheint die Wahrscheinlichkeit dafür zu sein, dass für ein  $u \notin S$  die Array-Einträge an allen Positionen  $h_1(u), \dots, h_k(u)$  auf 1 gesetzt sind. (und daher eine "falsche positive Antwort" gegeben wird).

Was ist also ein guter Kompromiss zwischen "kleinem  $k$ " und "großem  $k$ "?

Ziel: Für festes  $n$  und  $N$  wähle  $k$  so, dass die erwartete Anzahl falscher positiver Antworten möglichst klein ist.

Um die Analyse zu ermöglichen gehen wir im Folgenden davon aus, dass die von uns genutzten Hash-Funktionen die Werte aus  $U$  zufällig, gleichverteilt und unabhängig voneinander auf die Bits  $1 \dots n$  verteilen.

(In der Praxis genutzte Hash-Funktionen werden diese idealisierte Eigenschaft i.d.R. nicht haben; unsere Analyse wird dadurch aber überhaupt erst möglich. Bisher sind keine Ergebnisse bekannt, die theoretische Garantien für Bloom-Filter geben, wenn (strenge)  $k$ -universelle Familien von Hash-Funktionen genutzt werden.

Beobachtung 1:

Nach dem Einfügen von  $S$  ins Array  $A[1 \dots n]$  gilt für jedes feste Bit  $j \in \{1 \dots n\}$ :

$$Pr(A[j] = 0) = \left(1 - \frac{1}{n}\right)^{k \cdot N}$$

Zur Erinnerung:  
 $N = |S|$

Die Wahrscheinlichkeit läuft hier darüber, dass die Hash-Funktionen  $h_1, \dots, h_k$  zufällig und unabhängig voneinander gewählt werden.

Begründung: Es gilt:  $A[j] = 0 \Leftrightarrow \forall s \in S \text{ gilt: } h_1(s) \neq j \text{ und } \dots \text{ und } h_k(s) \neq j.$   
Für jedes  $s \in S$  und jedes  $i \in \{1, \dots, k\}$  ist

$$Pr(h_i(s) = j) = \frac{1}{n} \quad \text{und daher} \quad Pr(h_i(s) \neq j) = 1 - \frac{1}{n}.$$

Also gilt für jedes feste  $s \in S$ :

$$Pr(h_1(s) \neq j \text{ und } \dots \text{ und } h_k(s) \neq j) = \left(1 - \frac{1}{n}\right)^k$$

$h_1, \dots, h_k$  sind unabhängig voneinander gewählt

Auf Grund unserer idealisierenden Annahme, dass die Hash-Funktionen die Werte aus  $U$  zufällig, gleichverteilt und unabhängig voneinander auf die Bits  $1, \dots, n$  verteilen, erhalten wir:

$$\Pr \left( \forall s \in S \text{ gilt: } h_1(s) \neq j \text{ und } \dots \text{ und } h_k(s) \neq j \right) = \left(1 - \frac{1}{n}\right)^{k \cdot |S|}$$

$$\text{Somit gilt: } \Pr(A[j] = 0) = \left(1 - \frac{1}{n}\right)^{k \cdot N}$$

□ Beobachtung 1

Wir setzen  $p(k) := \left(1 - \frac{1}{n}\right)^{k \cdot N}$  Beob. 1  $\Pr(A[j] = 0)$  (für jedes feste  $j \in \{1, \dots, n\}$ )

Wir gehen im Folgenden vereinfachend davon aus, dass für paarweise verschiedene Werte  $j, i, i' \in \{1, \dots, n\}$  die Ereignisse " $A[j_i] \neq 0$ ", ..., " $A[j_{i'}] \neq 0$ " unabhängig voneinander sind. (Unter Verwendung der Chernoff-Schranke kann man zeigen, dass diese Annahme das Ergebnis nicht allzu sehr verfälscht.)

Beobachtung 2:

Für jedes feste  $u \in U \setminus S$  gilt:

$$\Pr \left( \underbrace{\text{"u liefert eine falsche positive Antwort"}} \right) = \left(1 - p(k)\right)^k$$

d.h.:  $A[h_i(u)] = 1$  für alle  $i \in \{1, \dots, k\}$

Begründung: Für jedes  $i \in \{1, \dots, k\}$  setze  $j_i := h_i(u)$ . Gemäß

$$\text{Beobachtung 1 gilt } \Pr(A[j_i] = 1) = 1 - p(k).$$

Auf Grund unserer vereinfachenden Annahme, dass die Ereignisse "A{j\_1} ≠ 0", ..., "A{j\_k} ≠ 0" alle unabhängig voneinander sind, erhalten wir:

$$Pr \left( \forall i \in \{1, \dots, k\} \text{ gilt: } A\{i\} = 1 \right) = (1 - p(k))^k$$

↳ Beobachtung 2

Ziel: Wähle k so, dass  $(1 - p(k))^k$  möglichst klein ist.

Lösungsansatz:

"Kurvendiskussion" für die Funktion

$$\tilde{f}: \mathbb{R} \rightarrow \mathbb{R} \quad \text{mit} \quad \tilde{f}(x) := (1 - p(x))^x \\ = \left( 1 - \left( 1 - \frac{1}{n} \right)^{N \cdot x} \right)^x$$

Typischer Trick, um die Rechnung zu vereinfachen, ohne das Ergebnis allzu sehr zu verfälschen:

Ersetze  $\left( 1 - \frac{1}{n} \right)$  durch  $e^{-\frac{1}{n}}$ .

Rechtfertigung: Es ist bekannt, dass

$$e = \lim_{n \rightarrow \infty} \left( 1 + \frac{1}{n} \right)^n \quad \text{und} \quad e^{-1} = \lim_{n \rightarrow \infty} \left( 1 - \frac{1}{n} \right)^n$$

Für hinreichend große n ist also daher

$e \approx \left( 1 + \frac{1}{n} \right)^n$

und

$e^{-1} \approx \left( 1 - \frac{1}{n} \right)^n$

Anstatt  $\tilde{f}(x) := \left(1 - \underbrace{\left(1 - \frac{1}{n}\right)^{N \cdot x}}_{\approx e^{-\frac{1}{n}}}\right)^x$  betrachten

wir daher die Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$  mit

$$f(x) := \left(1 - e^{-\frac{N}{n} \cdot x}\right)^x$$

Als unser gesuchtes  $k$  wählen wir dasjenige  $x$ , für das  $f(x)$  kleinstmöglich ist.

Usw:  $f(x)$  minimal  $\Leftrightarrow g(x) := \ln(f(x))$  minimal.

Wir minimieren also

$$g(x) := \ln(f(x)) = x \cdot \ln\left(1 - e^{-\frac{N}{n} \cdot x}\right).$$

Dazu suchen wir Nullstellen der ersten Ableitung von  $g(x)$ , also der Funktion

$$g'(x) = \ln\left(1 - e^{-\frac{N}{n} \cdot x}\right) + \frac{N}{n} \cdot x \cdot \frac{e^{-\frac{N}{n} \cdot x}}{1 - e^{-\frac{N}{n} \cdot x}}$$

Standard-Ableitungsregeln

Man kann (leicht) nachprüfen, dass Folgendes gilt:

- 1)  $g'(x) = 0 \Leftrightarrow x = (\ln 2) \cdot \frac{n}{N}$ , und
- 2) Der Wert  $x = (\ln 2) \cdot \frac{n}{N}$  liefert ein globales Minimum der Funktion  $g(x)$  und somit auch ein globales Minimum der Funktion  $f(x)$ .

(6)

Wir wählen also  $k := \underbrace{(\ln 2)}_{\approx 0,7} \cdot \frac{n}{N} \approx 0,7 \cdot \frac{n}{N}$

Für jedes feste  $n \in \mathbb{N}$  gilt dann:

$$\begin{aligned}
 & \Pr(u \text{ liefert eine falsche positive Antwort}) \\
 &= (1 - p(k))^k \\
 \text{Beob. 2} & \\
 &= \tilde{f}(k) \\
 &\approx f(k) \\
 &\approx f\left(\ln 2 \cdot \frac{n}{N}\right) \\
 &= \left(1 - \underbrace{e^{-\frac{n}{N} \cdot (\ln 2) \cdot \frac{n}{N}}}_{= e^{-\ln 2} = \frac{1}{e^{\ln 2}} = \frac{1}{2}}\right)^{(\ln 2) \cdot \frac{n}{N}} \\
 &= \left(\frac{1}{2}\right)^{(\ln 2) \cdot \frac{n}{N}}
 \end{aligned}$$

Wegen  $\left(\frac{1}{2}\right)^{\ln 2} \approx 0,6185$  gilt mit unserer Wahl von  $k$  also:

$$\Pr(\text{falsche positive Antwort}) \approx \left(0,6185\right)^{\frac{n}{N}}$$

Es gilt:

$y = \frac{n}{N}$	$\frac{1}{2}$	1	2	3	4	5	...	8	...	16
0,6185	$\approx 0,786$	0,6185	$\approx 0,3825$	$\approx 0,2366$	$\approx 0,1463$	$\approx 0,0905$		$\approx 0,0214$		$\approx 0,00045$
	$\approx 79\%$	$\approx 62\%$	$\approx 38\%$	$\approx 24\%$	$\approx 15\%$	$\approx 9\%$	...	$\approx 2\%$	...	$\approx 0,05\%$

Insbesondere:

Um für ein gegebenes  $S$  der Größe  $N$  eine Fehlerwahrscheinlichkeit von  $\approx 2\%$  zu bekommen, sollte man ein Bit-Array

der Länge  $n = 8 \cdot N$  verwenden und

$$k := \lceil (\ln 2) \cdot \frac{n}{N} \rceil = \lceil 8 \cdot \ln 2 \rceil = 6 \text{ Hash-Funktionen nutzen.}$$

Wenn man stattdessen ein Bit-Array

der Länge  $n = 16 \cdot N$  verwendet und

$$k := \lceil (\ln 2) \cdot \frac{n}{N} \rceil = \lceil 16 \cdot \ln 2 \rceil = 12 \text{ Hash-Funktionen nutzt,}$$

verringert sich die Fehlerwahrscheinlichkeit sogar auf  $\approx 0,05\%$ .

Zurück zu unserem Anwendungsbeispiel:

Hier ist  $S$  eine "white list", die aus  $N := 10^9$  Email-Adressen besteht, wobei zur Speicherung jeder einzelnen Email-Adresse ca 20 Bytes nötig sind.

Zur exakten Speicherung von  $S$  benötigen wir also  $20 \cdot 10^9$  Byte = 20 GB.

(60)  
Wenn wir stattdessen einen Bloom-Filter  
mit einem Bit-Array der Länge  $m = 8 \cdot N$   
und  $k = 6$  Hash-Funktionen verwenden, erhalten  
wir falsche positive Antworten mit ca 2%.

Zur Speicherung des Arrays  $A[1 \dots n]$  benötigen wir  
"nur"  $8 \cdot N$  Bits =  $8 \cdot 10^9$  Bits =  $1 \cdot 10^9$  Bytes = 1 GB.