

Häufigkeitsmomente

Wir betrachten Datenströme x_1, x_2, x_3, \dots ,
 deren Elemente x_i einem Universum U
 angehören. (z.B.: $U = \{0, 1, \dots, m-1\}$ für eine
 feste nat. Zahl m).

Zum Zeitpunkt $n \geq 1$ ist die
 Häufigkeit (oder Frequenz) eines Elements
 $u \in U$ im bisher gesehenen Strom x_1, \dots, x_n
 der Wert

$$f_u := |\{i \in \{1, \dots, n\} : x_i = u\}|.$$

Definition:

Für jede Zahl $k \in \mathbb{N}$ ist das
 k -te Häufigkeitsmoment des Stroms x_1, \dots, x_n
 definiert als die Zahl

$$F_k := \sum_{u \in U} (f_u)^k.$$

Insbesondere ist

$F_0 = \sum_{u \in U} (f_u)^0$ die Anzahl der verschiedenen Elemente,
 die im Strom vorkommen (hier wird $0^0 := 1$ und
 $a^0 := 1$ für alle $a \neq 0$ gesetzt)

$F_1 = \sum_{u \in U} (f_u)^1 = n$ die Länge des bisher gesehenen Stroms

$F_2 = \sum_{u \in U} (f_u)^2$ die "self-join-size" des bisher gesehenen Stroms.

Beispiel:

Für den Datenstrom

3, 1, 3, 3, 2, 1, 5, 2

gilt: $f_1 = 2$, $f_2 = 2$, $f_3 = 3$, $f_4 = 0$, $f_5 = 1$

$F_0 = 4$ (Anzahl verschiedener Elemente)

$F_1 = 8$

$$F_2 = f_1^2 + f_2^2 + f_3^2 + f_5^2$$

$$= 2^2 + 2^2 + 3^2 + 1^2$$

$$= 8 + 9 + 1$$

$$= 18$$

Die Häufigkeitsmomente liefern wichtige statistische Informationen über die Verteilung der Elemente aus U , die im Datenstrom vorkommen. Insbes. werden die Werte F_0 (und F_2 ?) z.B. bei der Optimierung von Datenbankabfragen genutzt.

Die Häufigkeitsmomente kann man natürlich leicht bestimmen, indem man für jedes $u \in U$ einen Zähler abspeichert, der besagt, wie oft u im bisher gesehenen Teil des Stroms bereits vorkam. Bei enorm großen U und enorm großen Datenströmen ist dies aber zu aufwändig. Daher wollen wir (unter Verwendung von deutlich weniger Speicherplatz) eine Approximation der Häufigkeitsmomente berechnen.

Der Einfachheit halber sei im Folgenden stets

$$U = \{0, 1, \dots, m-1\}$$

für eine natürliche Zahl m .

Das 0-te Häufigkeitsmoment:

Wie viele verschiedene Elemente enthält der Strom?

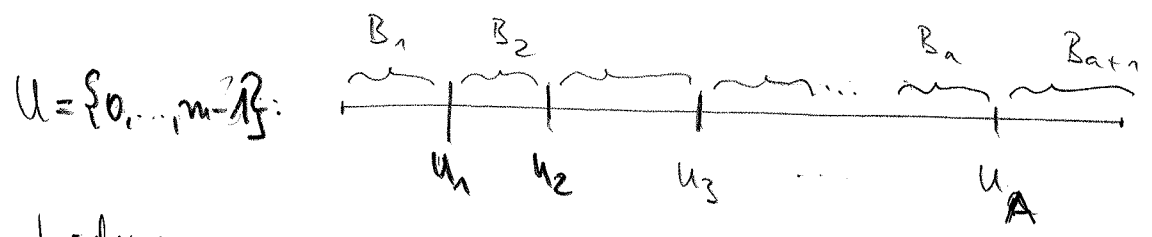
Intuition:

Betrachte zunächst den Fall, in dem die im Datenstrom vorkommenden Elemente zufällig, gleichverteilt aus $U = \{0, \dots, m-1\}$ gewählt wurden.

Sei m_{in} das kleinste Element aus U , das im Strom vorkommt.

Anzahl verschiedener Elemente im Strom	durchschnittlicher Wert von m_{in}
1	$\approx \frac{m}{2}$
2	$\approx \frac{m}{3}$
A	$\approx \frac{m}{A+1}$

Skizze:



Verschiedene Elemente im Strom: $u_1 < u_2 < \dots < u_A$

$\Rightarrow U$ wird in $A+1$ verschiedene Blöcke aufgeteilt; jeder Block hat erwartete Größe $\approx \frac{|U|}{A+1} = \frac{m}{A+1}$

Also:

In einem zufällig gewählten Strom mit genau A verschiedenen Elementen (d.h. $F_0 = A$) vermuten wir für das kleinste Element \min , das im Strom vorkommt, dass gilt:

$$\min \approx \frac{m}{A+1}$$

Somit: $A \approx \frac{m}{\min} - 1$ bzw. $A \approx \frac{m}{\min}$.

Ein erster Algorithmus zum Abschätzen von F_0 :

- Das Universum $U = \{0, \dots, m-1\}$ ist bekannt
- Beim Lesen des Stroms speichere das kleinste im Strom gesehene Element \min
- Gib den Wert $\frac{m}{\min} - 1$ aus und behaupte, dass dieser Wert eine Abschätzung für die Anzahl $A = F_0$ der im Strom vorkommenden verschiedenen Elemente ist.

Problem:

verschiedenen

Wenn die Datenstrom-Elemente zufällig, gleichverteilt sind, mag zumindest die Intuition stimmen. Aber tatsächliche Datenströme sind in der Regel nicht "zufällig und gleichverteilt".

Lösung:

Nimm eine hinreichend große Zahl M und eine geeignete Familie H von Hash-Funktionen von $U = \{0, 1, \dots, m-1\}$ nach $V := \{0, 1, \dots, M-1\}$.

Wähle zufällig, gleichverteilt ein $h \in H$.

An Stelle des Stroms x_1, x_2, \dots, x_n

betrachte den Strom $h(x_1), h(x_2), \dots, h(x_n)$,

bestimme das kleinste Element \min dieses Stroms und gib als Schätzer für F_0 den Wert

$$F_0^* := \frac{M}{\min + 1} \quad \text{aus, für } \min := \begin{cases} 1 & \text{falls } \min=0 \\ \min & \text{sonst.} \end{cases}$$

Fragen:

(1) Welches M und welches H sollen wir wählen?

(2) Wie nah liegt der Schätzwert F_0^* am tatsächlichen Wert F_0 dran?

Zu (1):

- Wähle eine Zahl $M > 9m$.
- Wähle eine streng 2-universelle Familie H von Hash-Funktionen von U nach $V = \{0, \dots, M-1\}$.

(Zur Erinnerung:

Wir haben bereits gesehen, dass wenn M eine Primzahl ist, die Menge $H = \{h_{a,b} : a, b \in V\}$ streng 2-universell ist, wobei

$$h_{a,b}(x) := a \cdot x + b \pmod{M} \quad \text{ist} \quad (\forall a, x \in U).$$

Diese Wahl von M und H stellt sicher, dass wir beweisen können, dass Frage (2) wie folgt beantwortet werden kann:

zu (2): Es gilt:

$$\Pr \left(\frac{F_0}{g} \leq F_0^* \leq g \cdot F_0 \right) \geq \frac{2}{3} \quad (\triangle)$$

d.h. mit Wk. mindestens $\frac{2}{3}$ liegt die Abschätzung F_0^* , die wir angeben, zwischen $\frac{1}{g} \cdot F_0$ und $g \cdot F_0$,

wobei F_0 die tatsächliche Anzahl der verschiedenen im Strom vorkommenden Elemente ist.

Beweis von (\triangle) :

$$\text{Behauptung 1: } \Pr (F_0^* > g \cdot F_0) \leq \frac{2}{g}$$

Beweis:

$$\begin{aligned} \Pr \left(\underbrace{F_0^*}_{\min} > g \cdot F_0 \right) &= \Pr \left(\min < \frac{M}{g \cdot F_0} \right) \\ &= \frac{M}{\min} \quad (\text{bzw } M, \text{ falls } \min = 0) \end{aligned}$$

$$\frac{M}{\min} > g \cdot F_0 \quad (\Rightarrow) \quad \frac{M}{g \cdot F_0} > \min$$

Sei $A := F_0$ und seien u_1, \dots, u_A die verschiedenen im Strom x_1, \dots, x_n vorkommenden Elemente. Dann ist \min das kleinste Element in $h(u_1), \dots, h(u_A)$, wobei h unsere zufällig gewählte Hash-Funktion ist.

Es gilt

$$A = F_0$$

(42)

$$\Pr\left(\min < \frac{M}{g \cdot F_0}\right) = \Pr\left(\exists i \in \{1, \dots, A\} : h(u_i) < \frac{M}{g \cdot A}\right)$$

$$\leq \sum_{i=1}^A \Pr\left(h(u_i) < \frac{M}{g \cdot A}\right)$$

$$\leq A \cdot \frac{\frac{M}{gA}}{M}$$

$$\leq A \cdot \frac{\frac{M}{gA} + 1}{M}$$

$$= A \cdot \left(\frac{\frac{M}{gA}}{M} + \frac{1}{M}\right) = A \cdot \left(\frac{1}{gA} + \frac{1}{M}\right) = \frac{1}{g} + \frac{A}{M} = \frac{1}{g} + \frac{F_0}{M}$$

Aus "H ist streng 2-universell" folgt dass für ein zufällig gewähltes l und festes u_i der Wert $h(u_i)$ wie zufällig aus $\{0, 1, \dots, M-1\}$ gewählt ist. D.h.: $\Pr(h(u_i) = v) = \frac{1}{M}$ für jedes $v \in \{0, 1, \dots, \frac{M}{gA}-1\}$.
 Und: $h(u_i) < \frac{M}{gA} \Leftrightarrow h(u_i) \in \{0, 1, \dots, \frac{M}{gA}-1\}$

Da jedes im Strom vorkommende Element aus $U = \{0, \dots, m-1\}$ ist, können im Strom nur maximal m verschiedene Elemente vorkommen, d.h.: $F_0 \leq m$.

Daher ist $\frac{F_0}{M} \leq \frac{m}{M} < \frac{1}{g}$, da $M > gm$ ist.

Insgesamt ist also $\Pr(\min < \frac{M}{gF_0}) \leq \frac{1}{g} + \frac{1}{g} = \frac{2}{g}$

□ Beh 1

Behauptung 2: $\Pr\left(F_0^* < \frac{F_0}{g}\right) \leq \frac{1}{g}$

Beweis: Wegen $F_0^* = \frac{M}{\min}$ und $A := F_0$ gilt:

$$F_0^* < \frac{F_0}{g} \Leftrightarrow \frac{M}{\min} < \frac{A}{g} \Leftrightarrow \min > \frac{gM}{A}$$

Also: $\Pr\left(F_0^* < \frac{F_0}{g}\right) = \Pr\left(\min > \frac{gM}{A}\right) = \Pr\left(\forall i \in \{1, \dots, A\} : h(u_i) > \frac{gM}{A}\right)$

wobei u_1, \dots, u_A sämtliche verschiedene im Strom vorkommende Elemente sind.

Für $i \in \{1, \dots, A\}$ sei Y_i die Zufallsvariable

$$\text{mit } Y_i = \begin{cases} 0 & \text{falls } h(U_i) > \frac{gM}{A} \\ 1 & \text{sonst, d.h. } h(U_i) \leq \frac{gM}{A} \end{cases}$$

und sei $Y := \sum_{i=1}^A Y_i$. Es gilt: $Y=0 \Leftrightarrow \forall i \in \{1, \dots, A\}: h(U_i) > \frac{gM}{A}$

Ziel: Zeige, dass $\Pr(Y=0) < \frac{1}{9}$ ist.

Dazu nutzen wir die Tschebyscheff-Ungleichung, die folgendes besagt:

Tschebyscheff-Ungleichung (über einem endlichen Wahrscheinlichkeitsraum)

Für jede Zufallsvariable X mit Erwartungswert $E(X)$ und Varianz $\text{Var}(X)$ gilt für jede reelle Zahl $r > 0$:

$$\Pr(|X - E(X)| \geq r) \leq \frac{\text{Var}(X)}{r^2}$$

Zur Erinnerung:

$$\text{Erwartungswert } E(X) := \sum_i i \cdot \Pr(X=i)$$

$$\begin{aligned} \text{Varianz } \text{Var}(X) &:= E((X - E(X))^2) \\ &= E(X^2) - (E(X))^2 \end{aligned}$$

↑
ein Maß dafür, wie weit die tatsächlichen Werte von X um den Erwartungswert $E(X)$ verstreut sind

↑
Linearität des Erwartungswerts

Für uns: Je kleiner die Varianz, desto besser!

Es gilt:

- Für jedes $i \in \{1, \dots, A\}$ ist

$$E(Y_i) = P_i(Y_i = 1) \geq \frac{g}{A}$$

↑
 bei der Wahl des Wertes $h(i)$ gibt es
 $\geq \frac{gM}{A}$ "gute Werte" bei insgesamt M Werten
 nämlich: $0, 1, \dots, \frac{gM}{A} - 1$

- $E(Y) = E\left(\sum_{i=1}^A Y_i\right) = \sum_{i=1}^A E(Y_i) \geq A \cdot \frac{g}{A}$

D.h.: $E(Y) \geq g$ (1)

- Da H streng 2-universell ist,

gilt: $\text{Var}\left(\sum_{i=1}^A Y_i\right) = \sum_{i=1}^A \text{Var}(Y_i)$

(Beweis: Übung!)

Somit ist $\text{Var}(Y) = A \cdot \text{Var}(Y_1)$

- $\text{Var}(Y_1) = \underbrace{E(Y_1^2)} - \underbrace{\left(E(Y_1)\right)^2}_{\geq 0} \leq E(Y_1)$
 $= E(Y_1)$, da Y_1 nur den Wert 0 oder den Wert 1 annehmen kann

- Somit ist $\text{Var}(Y_1) \leq E(Y_1)$ und

$$\text{Var}(Y) = A \cdot \text{Var}(Y_1) \leq A \cdot E(Y_1) = E(Y)$$

Also $\text{Var}(Y) \leq E(Y)$ (2)

- Tschebyscheff liefert für $X := Y$ und $\mu := E(Y)$: (45)

$$\Pr \left(|Y - E(Y)| \geq E(Y) \right) \leq \frac{\text{Var}(Y)}{E(Y)^2}$$

$$\stackrel{\textcircled{2}}{\leq} \frac{E(Y)}{E(Y)^2}$$

$$= \frac{1}{E(Y)}$$

$$\stackrel{\textcircled{1}}{\leq} \frac{1}{9}$$

- Insgesamt gilt:

$$\Pr \left(F_0^* < \frac{F_0}{9} \right) = \Pr(Y = 0)$$

$$\begin{aligned} &\leq \Pr \left(|Y - E(Y)| \geq E(Y) \right) \\ &\leq \frac{1}{9} \end{aligned}$$

\nearrow für $Y=0$ ist $|Y - E(Y)| = E(Y)$

□ Beh 2

Zusammenfassung:

$$\begin{aligned} &\Pr \left(F_0^* < \frac{F_0}{9} \text{ oder } F_0^* > 9 \cdot F_0 \right) \\ &\leq \Pr \left(F_0^* < \frac{F_0}{9} \right) + \Pr \left(F_0^* > 9 \cdot F_0 \right) \\ &\stackrel{\text{Beh 2+1}}{\leq} \frac{1}{9} + \frac{2}{9} = \frac{3}{9} = \frac{1}{3} \end{aligned}$$

$$\text{Somit: } \Pr \left(\frac{F_0}{9} \leq F_0^* \leq 9 \cdot F_0 \right) \geq 1 - \frac{1}{3} = \frac{2}{3}$$

□ Beweis von \textcircled{A}

Insgesamt haben wir gesehen, dass Folgendes gilt:

Satz (Abschätzung von F_0)

Für $m \geq 1$ sei $U := \{0, \dots, m-1\}$

Das folgende Verfahren liefert bei Eingabe eines Stroms x_1, \dots, x_n von Elementen aus U einen Schätzwert F_0^* , für den gilt:

$$\Pr \left(\frac{F_0}{3} \leq F_0^* \leq 3 \cdot F_0 \right) \geq \frac{2}{3}$$

und benutzt dazu $O(\log m)$ Speicherbits.

Verfahren:

- Wähle eine Primzahl M mit $3m < M \leq 18m$ (eine solche Primzahl gibt es gemäß Bertrand'schem Postulat).
- Wähle zufällig, gleichverteilt und unabhängig voneinander zwei Zahlen $a, b \in \{0, 1, \dots, M-1\}$ und nutze im Folgenden die Hash-Funktion $h: U \rightarrow \{0, 1, \dots, M-1\}$ mit $h(x) := ax + b \pmod{M}$ (v.a. $x \in U$).
- Beim Lesen des ersten Strom-Elements x_1 initialisiere die Variable min durch $min := h(x_1)$
- Für jedes $i \geq 1$ tue beim Lesen des $(i+1)$ -ten Strom-Elements x_{i+1} folgendes:
 - berechne $h(x_{i+1})$
 - falls $h(x_{i+1}) < min$, so setze $min := h(x_{i+1})$
- Nach dem Lesen des n -ten Strom-Elements x_n gib den Wert $F_0^* := \frac{M}{min'}$ aus, wobei $min' := \begin{cases} 1, & \text{falls } min = 0 \\ min, & \text{sonst.} \end{cases}$

(Quellenachweis: Die Darstellung auf den Seiten 37-46 orientiert sich an dem Abschnitt "Number of distinct elements in a data stream" im Buch "Foundations of Data Science" von Blum, Hopcraft, Kannan (Version: May 14, 2015).)

Anmerkungen zur Abschätzung von F_0 :

(47)

Mit dem auf Seite 46 vorgestellten Verfahren bekommen wir eine Abschätzung F_0^* , für die gilt:

$$\Pr\left(\frac{1}{9} \cdot F_0 \leq F_0^* \leq 9 \cdot F_0\right) \geq \frac{2}{3} = 1 - \frac{1}{3}$$

Für jede feste Zahl δ mit $0 < \delta < 1$ erhalten wir durch Anwenden des Median-Tricks (siehe Seite 50) ein Verfahren, das einen Schätzwert \hat{F}_0 für F_0 liefert, so dass gilt:

$$\Pr\left(\frac{1}{9} \cdot F_0 \leq \hat{F}_0 \leq 9 \cdot F_0\right) \geq 1 - \delta$$

Dieses Verfahren nutzt Speicherplatz $O(s \cdot \log m)$

für $s = \lceil c \cdot \ln\left(\frac{1}{\delta}\right) \rceil$ mit

$$c = \frac{6}{(1 - 2 \cdot \frac{1}{3})^2} = \frac{6}{(\frac{1}{3})^2} = \frac{6}{\frac{1}{9}} = 6 \cdot 9 = 54.$$

Gut an dem Verfahren ist:

- hohe Erfolgswahrscheinlichkeit: $1 - \delta$
- geringer Speicherplatzverbrauch: $O(54 \cdot \ln\left(\frac{1}{\delta}\right) \cdot \log m)$

Schlecht ist:

- Der Schätzwert \hat{F}_0 kann noch recht weit vom tatsächlichen Wert F_0 wegliegen.

An Stelle von " $\frac{1}{9} \cdot F_0 \leq \hat{F}_0 \leq 9 \cdot F_0$ " hätten wir lieber:

$$(1 - \epsilon) \cdot F_0 \leq \hat{F}_0 \leq (1 + \epsilon) \cdot F_0 \quad (\text{also: } \hat{F}_0 \in (1 \pm \epsilon) \cdot F_0).$$

Lösung:

(48)

An Stelle des auf Seite 46 beschriebenen Verfahrens
("Speichere in m den kleinsten Wert $h(x_i)$
für x_1, \dots, x_n und gib den Schätzwert $\tilde{F}_0 := \frac{M}{\min}$ aus")

gehe wie folgt vor:

- Wähle $M := m^3$
- Wähle zufällig, gleichverteilt und unabhängig voneinander zwei Zahlen $a, b \in \{0, 1, \dots, M-1\}$ und nutze im Folgenden die Hash-Funktion
 $h: U \rightarrow \{0, 1, \dots, M-1\}$ mit $h(x) = ax + b \pmod{M}$
(f.a. $x \in U$)
- Wähle $t := \lceil \frac{96}{\epsilon^2} \rceil$ und speichere während des Lesens des Stroms x_1, x_2, \dots die t kleinsten paarweise verschiedenen Werte in $h(x_1), h(x_2), \dots$
- Nach dem Lesen des n -ten Strom-Elements x_n betrachte den t -kleinsten gespeicherten Wert \min_t
z.B. für $t=4$ wäre das der viert-kleinste Wert
und gib den Wert $\tilde{F}_0 := \frac{t \cdot M}{\min_t}$ aus.

Speicherverbrauch: $O\left(\frac{1}{\epsilon^2} \cdot \log m\right)$ Bits

Satz (Bar-Yossef et al, 2002)

Für den mit diesem Verfahren berechneten Schätzwert \tilde{F}_0 gilt:

$$\Pr\left((1-\epsilon) \cdot F_0 \leq \tilde{F}_0 \leq (1+\epsilon) F_0\right) \geq \frac{2}{3} - \frac{1}{m}$$

Hier ohne Beweis — die Details sind ähnlich wie beim Beweis des Satzes auf Seite 46, aber etwas schwieriger. Siehe Bar-Yossef, Jayram, Kumar, Sivakumar, Trevisan "Counting Distinct Elements in a Data Stream", RANDOM 2002.

Durch Anwenden des Median-Tricks (siehe Seite 50) erhalten wir dann ein Verfahren, das einen Schätzwert \hat{f}_0 liefert, so dass gilt:

$$\Pr \left(\hat{f}_0 \in (1 \pm \epsilon) \cdot f_0 \right) \geq 1 - \delta$$

und dabei $O \left(\frac{1}{\epsilon^2} \cdot \ln \left(\frac{1}{\delta} \right) \cdot \log m \right)$ Speicherbits nutzt:

Für eine geeignete Zahl $s = \Theta \left(\ln \left(\frac{1}{\delta} \right) \right)$

führe unabhängig voneinander s -mal das auf Seite 4.8 vorgestellte Verfahren durch und

gib als \hat{f}_0 den Median der s berechneten Schätzwerte aus.

Der Median-Trick

... zum Erhöhen der Erfolgswahrscheinlichkeit auf $1-\delta$
für $0 < \delta < 1$.

Wenn wir ein Verfahren V haben, das einen Schätzwert F^* für einen Zielwert F berechnet, so dass gilt

$$Pr(aF \leq F^* \leq bF) = 1-p,$$

für Werte a, b, p mit $0 \leq a \leq b$ und $0 \leq p < \frac{1}{2}$,

so liefert beliebiges Verfahren \hat{V} einen Schätzwert \hat{F} für F , so dass gilt

$$Pr(aF \leq \hat{F} \leq bF) \geq 1-\delta.$$

Bsp: $\delta = \frac{1}{1.000.000} \Rightarrow \frac{1}{\delta} = 10^6 \Rightarrow \ln(\frac{1}{\delta}) \approx 13,81...$
 $p = \frac{1}{3} \Rightarrow (1-2p)^2 = (1-\frac{2}{3})^2 = \frac{1}{9} \Rightarrow \frac{6}{(1-2p)^2} = 6 \cdot 9 = 54$ } \Rightarrow Wähle $s \geq 750$

Verfahren \hat{V} :

- Wähle ein beliebiges $c \geq \frac{6}{(1-2p)^2}$ und ein beliebiges $s \geq c \cdot \ln(\frac{1}{\delta})$.
- Führe unabhängig voneinander s -mal das Verfahren V durch, um s Schätzwerte $F_1^*, F_2^*, \dots, F_s^*$ zu erhalten
- Berechne den Median \hat{F} der Werte F_1^*, \dots, F_s^* und gib diesen aus.

Um zu beweisen, dass tatsächlich $Pr(aF \leq \hat{F} \leq bF) \geq 1-\delta$

gilt, wenden wir die Chernoff-Schranke an:

Die Chernoff-Schranke

Satz B.20 (Chernoff-Schranke).

Sei p eine reelle Zahl mit $0 < p < 1$. Wir werfen eine Münze, die mit Wahrscheinlichkeit p auf „Kopf“ und mit Wahrscheinlichkeit $1 - p$ auf „Zahl“ landet. Nach dem i -ten Münzwurf setzen wir

$$X_i := \begin{cases} 1 & \text{falls die Münze beim } i\text{-ten Wurf auf „Kopf“ gefallen ist} \\ 0 & \text{sonst.} \end{cases}$$

Wir machen insgesamt $s \in \mathbb{N}_{\geq 1}$ Münzwürfe und setzen

$$X := \sum_{i=1}^s X_i,$$

d.h. X gibt an, bei wie vielen der s Münzwürfe die Münze auf „Kopf“ gefallen ist. Dann gilt:

(a) Der Erwartungswert von X ist $E(X) = p \cdot s$.

(b) Für jede reelle Zahl ε mit $0 < \varepsilon < 1$ ist

$$P\left(X < (1 - \varepsilon) \cdot E(X)\right) < e^{-\frac{1}{2} \cdot \varepsilon^2 \cdot ps}.$$

(c) Für jede reelle Zahl $\beta > 0$ ist

$$P\left(X > (1 + \beta) \cdot E(X)\right) < e^{-\frac{1}{3} \cdot \min(\beta, \beta^2) \cdot ps}.$$

(d) Für jede reelle Zahl ε mit $0 < \varepsilon < 1$ ist

$$P\left(X \notin (1 \pm \varepsilon) \cdot E(X)\right) < 2 \cdot e^{-\frac{1}{3} \cdot \varepsilon^2 \cdot ps}.$$

Hierbei verwenden wir „ $X \notin (1 \pm \varepsilon) \cdot E(X)$ “ als Kurzschreibweise für

$$„X < (1 - \varepsilon) \cdot E(X) \text{ oder } X > (1 + \varepsilon) \cdot E(X)“$$

Wir betrachten nun folgendes Szenario:

Die bei der Chernoff-Schranke betrachtete Münze fällt beim i -ten Wurf auf "Kopf" (mit Wk p), falls nicht $aT \leq F_i^* \leq bT$ gilt.

Zur Anwendung der Chernoff-Schranke setzen wir

$$X_i = \begin{cases} 1 & \text{falls nicht } aT \leq F_i^* \leq bT \text{ gilt} \\ 0 & \text{sonst} \end{cases}$$

und $X = \sum_{i=1}^s X_i$. D.h. X gibt an, bei wie vielen der s Durchläufe des Verfahrens V der Schätzwert zu weit vom Zielwert T weg ist.

Gemäß Chernoff-Schranke gilt

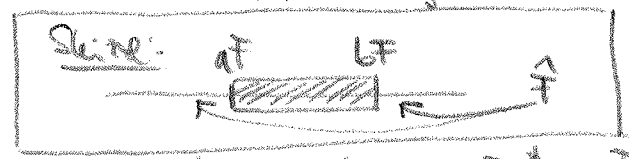
(1) $E(X) = p \cdot s$, und

(3) für jede Zahl $\beta > 0$ gilt:
$$\Pr(X > (1+\beta)E(X)) < e^{-\frac{1}{2} \cdot \min(\beta, \beta^2) \cdot p \cdot s} \quad (*)$$

• Wir wollen zeigen, dass $\Pr(\text{es gilt nicht } aT \leq \hat{T} \leq bT) \leq \delta$ ist.
Dann beachte, dass \hat{T} der Median von F_1^*, \dots, F_s^* ist.

Daher gilt Folgendes:

$aT \leq \hat{T} \leq bT$ gilt nicht



\Rightarrow Für mehr als die Hälfte der Werte F_1^*, \dots, F_s^* ist $F_i^* < aT$
oder
für mehr als die Hälfte der Werte F_1^*, \dots, F_s^* ist $F_i^* > bT$

$\Rightarrow X > \frac{s}{2}$

Somit gilt:

(52)

$$\Pr\left(\text{es gilt nicht } aT \leq \hat{T} \leq bT\right) \leq \Pr\left(X > \frac{s}{2}\right)$$

Es reicht also, zu zeigen, dass $\Pr\left(X > \frac{s}{2}\right) \leq \delta$ ist.

• Wegen $E(X) = p \cdot s$ ist $s = \frac{E(X)}{p}$. Daher ist

$$X > \frac{s}{2} \quad (\Leftrightarrow) \quad X > \frac{1}{2p} \cdot E(X)$$

• Um $\textcircled{*}$ anzuwenden, wählen wir β so, dass

$$1 + \beta = \frac{1}{2p} \quad \text{ist.}$$

D.h. wir wählen $\beta := \frac{1}{2p} - 1$.

Beachte: $\beta > 0$, denn:

$$\frac{1}{2p} - 1 > 0 \quad (\Leftrightarrow) \quad \frac{1}{2p} > 1 \quad (\Leftrightarrow) \quad 1 > 2p \quad (\Leftrightarrow) \quad p < \frac{1}{2}, \text{ was laut Voraussetzung erfüllt ist.}$$

Gemäß $\textcircled{*}$ gilt also:

$$\Pr\left(X > \frac{s}{2}\right) = \Pr\left(X > (1+\beta) \cdot E(X)\right) < e^{-\frac{1}{3} \cdot \min(\beta, \beta^2) \cdot p \cdot s}$$

Es reicht also, zu zeigen, dass $e^{-\frac{1}{3} \cdot \min(\beta, \beta^2) \cdot p \cdot s} \leq \delta$ ist.

• Es gilt:

$$e^{-\frac{1}{3} \cdot \min(\beta, \beta^2) \cdot p \cdot s} \leq \delta$$

$$(\Leftrightarrow) \quad -\frac{1}{3} \cdot \min(\beta, \beta^2) \cdot p \cdot s \leq \ln(\delta)$$

$$(\Leftrightarrow) \quad \frac{1}{3} \cdot \min(\beta, \beta^2) \cdot p \cdot s \geq \underbrace{-\ln(\delta)}_{= \ln\left(\frac{1}{\delta}\right)}$$

$$(\Leftrightarrow) \quad s \geq \frac{3}{p \cdot \min(\beta, \beta^2)} \cdot \ln\left(\frac{1}{\delta}\right)$$

• Was ist $\frac{3}{p \cdot \min(\beta, \beta^2)}$?

Es gilt:

$$* \min(\beta, \beta^2) = \begin{cases} \beta^2 & \text{falls } 0 < \beta < 1 \\ \beta & \text{falls } \beta \geq 1 \end{cases}$$

$$* \beta \stackrel{\text{Def}}{=} \frac{1}{2p} - 1 > 0, \text{ da } p < \frac{1}{2} \text{ laut Voraussetzung}$$

$$* \beta \geq 1 \Leftrightarrow \frac{1}{2p} - 1 \geq 1 \Leftrightarrow \frac{1}{2p} \geq 2 \Leftrightarrow 1 \geq 4p$$

$$\Leftrightarrow p \leq \frac{1}{4}$$

Somit gilt:

* Falls $p \leq \frac{1}{4}$, so ist

$$\min(\beta, \beta^2) = \beta = \frac{1}{2p} - 1 = \frac{1-2p}{2p}, \text{ und daher}$$

$$p \cdot \min(\beta, \beta^2) = \frac{1-2p}{2}, \text{ und } 0 \leq 1-2p \leq 1, \text{ da } 0 \leq p \leq \frac{1}{2}$$

$$\frac{3}{p \cdot \min(\beta, \beta^2)} = \frac{3 \cdot 2}{1-2p} = \frac{6}{1-2p} \leq \frac{6}{(1-2p)^2} \leq c$$

* Falls $p > \frac{1}{4}$, so ist $\beta \geq 1$ und daher

$$\min(\beta, \beta^2) = \beta^2 = \left(\frac{1}{2p} - 1\right)^2 = \left(\frac{1-2p}{2p}\right)^2 = \frac{(1-2p)^2}{4p^2},$$

$$\text{und daher}$$

$$p \cdot \min(\beta, \beta^2) = \frac{(1-2p)^2}{4p}, \text{ und } 0 \leq p \leq \frac{1}{2}$$

$$\frac{3}{p \cdot \min(\beta, \beta^2)} = \frac{3 \cdot 4p}{(1-2p)^2} = \frac{12p}{(1-2p)^2} \leq \frac{6}{(1-2p)^2} \leq c$$

• Durch Wählen von $S \geq c \cdot \ln\left(\frac{1}{\delta}\right)$ erreichen wir also, dass

$S \geq \frac{3}{p \cdot \min(\beta, \beta^2)} \cdot \ln\left(\frac{1}{\delta}\right)$ ist und erhalten somit, dass $e^{-\frac{1}{3} \cdot \min(\beta, \beta^2) \cdot p \cdot S} \leq \delta$ ist. Dies beendet den Beweis dafür, dass der Median-Trick funktioniert. \square

Der BJKST-Algorithmus zur Abschätzung von F_0

In der Arbeit "Counting distinct elements in a data stream" bei RANDOM 2002 haben die Autoren Bar-Yossef, Jayram, Kumar, Sivakumar und Trevisan (kurz: BJKST) den folgenden Algorithmus zur Abschätzung von F_0 vorgestellt, der in vielerlei Hinsicht optimal (bzgl. Speicherplatzverbrauch und Approximationsgüte) ist.

Zur Darstellung des Algorithmus nutzen wir folgenden Begriff:

Für eine natürliche Zahl n setzen wir

$$\text{Nullen}(n) := \max \{ i : n \text{ ist ein Vielfaches von } 2^i \}$$

D.h.: $\text{Nullen}(n)$ ist die Anzahl der Nullen am Ende der Binärdarstellung von n

Bsp:

n	Binärdarstellung von n	Nullen(n)
0	0	1
1	1	0
2	10	1
3	11	0
4	100	2
5	101	0
6	110	1
7	111	0
8	1000	3
16	10000	4

Satz (BKST-Algorithmus)

Für $m \geq 1$ sei $U := \{0, 1, \dots, m-1\}$.
Sei ϵ mit $0 < \epsilon < 1$, so dass $\frac{1}{\epsilon^2} = o(m)$

Der folgende so genannte BKST-Algorithmus liefert bei Eingabe eines Stroms x_1, \dots, x_n von Elementen aus U einen Schätzwert T_0^* , für den gilt

$$\Pr (T_0^* \in (1 \pm \epsilon) T_0) \geq \frac{2}{3}$$

und benutzt dazu $O\left(\log m + \frac{1}{\epsilon^2} \cdot (\log(\frac{1}{\epsilon}) + \log \log m)\right)$ Speicherbits.

BKST-Algorithmus:

- Wähle zufällig, gleichverteilt und unabhängig voneinander zwei Zahlen $a, b \in \{0, 1, \dots, m-1\}$ und nutze im Folgenden die Hash-Funktion $h := h_{a,b}$, dh $h: U \rightarrow U$ mit $h(x) := ax + b \pmod{m}$ für $a \neq 0$.
- Setze $M := \frac{4}{\epsilon^4} \cdot (\log m)^{27}$, wähle zufällig, gleichverteilt und unabhängig voneinander zwei Zahlen $a', b' \in \{0, 1, \dots, M-1\}$ und nutze im Folgenden die Hash-Funktion $g := g_{a',b'}$ mit $g: U \rightarrow \{0, 1, \dots, M-1\}$ mit $g(x) := a'x + b' \pmod{M}$ für $a' \neq 0$.
- Setze $N := 0$ und $S := \emptyset$
- Beim Lesen eines Strom-Elements x_i tue Folgendes:
Falls $N \leq h(x_i) < M$, dann
 - Füge das Tupel $(g(x_i), h(x_i))$ in S ein
 - Solange $|S| \geq \frac{600}{\epsilon^2}$ setze $N := N+1$ und entferne aus S alle Tupel (α, β) mit $\beta < N$.
- Nach dem Lesen des n -ten Strom-Elements gib $T_0^* := |S| \cdot 2^N$ aus.

Wir werden hier keinen formalen Beweis des Satzes angeben, wollen uns aber eine Intuition dazu ersarbeiten, wie der Algorithmus funktioniert. (56)

1) Speicherplatzbedarf:

- $\leq 6 \cdot \log m$ Bits zum Speichern von h, g, M, N
- Gemäß Algorithmus ist stets $|S| \leq \left\lfloor \frac{600}{\epsilon^2} \right\rfloor + 1$.

zum Speichern eines einzelnen Tupels $(g(x_i), \underbrace{\text{Nullen}(h(x_i))}_{\text{Bits}})$ in S reichen

$$\underbrace{\log M}_{\log\left(\frac{4}{\epsilon^4}\right)} + \log \log m \text{ Bits,}$$

$$\log\left(\frac{4}{\epsilon^4}\right) + \log((\log m)^2)$$

also $O\left(\log\left(\frac{1}{\epsilon}\right) + \log \log m\right)$ Bits.

Zum Speichern der gesamten Menge S werden also $O\left(\frac{1}{\epsilon^2} \cdot \left(\log\left(\frac{1}{\epsilon}\right) + \log \log m\right)\right)$ Bits genutzt.

2) Arbeitsweise:

- Die Hash-Funktion g nutzen wir, um Speicherplatz zu sparen. Man kann nachrechnen, dass g mit hoher Wahrscheinlichkeit kollisionsfrei ist, d.h. $|S| = |S'|$ für

$$S' := \left\{ (x_i, \text{Nullen}(h(x_i))) : (g(x_i), \text{Nullen}(h(x_i))) \in S \right\}$$

(Beachte dazu, dass $|S| \leq \frac{600}{\epsilon^2} + 1$ und dass insgesamt g auf maximal $\left(\frac{600}{\epsilon^2} + 1\right) \cdot \log m$ Elemente angewendet wird.)

Die Zahl $M := \frac{4}{\epsilon^2} (\log m)^2$ ist so gewählt, dass ⁽⁵⁷⁾ die Wahrscheinlichkeit hoch genug und gleichzeitig der Speicherplatzbedarf möglichst gering ist.

Im Folgenden gehen wir davon aus, dass $|S| = |S'|$ ist.

• Warum liefert $|S'| \cdot 2^N$ einen guten Schätzwert für T_0 ?

* Die Hash-Funktion h wird — wie bei unserem ersten Algorithmus zur Abschätzung von T_0 — genutzt, um die verschiedenen Datenstrom-Elemente möglichst zufällig auf Werte aus $\{0, 1, \dots, m-1\}$ zu verteilen, d.h. auf Bit-Strings der Länge $\log m$ (indem wir die Binärrépräsentation der Werte betrachten).

* Die Zahl N ist so gewählt, dass S' aus genau denjenigen $(x_i, \text{Nullen}(h(x_i)))$ besteht, für die $\text{Nullen}(h(x_i)) \geq N$ ist. Beachte: S' ist eine Menge und keine Multimenge, d.h. für $x_i = x_j$ mit $i \neq j$ enthält S' das Tupel $(x_i, \text{Nullen}(h(x_i)))$ nur einmal, nicht zweimal.

* Wenn der Strom aus $A := T_0$ verschiedenen Elementen besteht und h diese Elemente zufällig auf Bitstrings verteilt, ist zu erwarten, dass bei etwa $\frac{A}{2^N}$ dieser Bitstrings die letzten N Bits gleich 0 sind. D.h.: $|S'| \approx \frac{A}{2^N}$, und

somit ist $A \approx |S'| \cdot 2^N$ — und das ist der Schätzwert, den der BJKST-Algorithmus ausgibt. //

(58)
Unter Verwendung des Median-Tricks können wir natürlich wieder die Erfolgswahrscheinlichkeit des BJKST-Algorithmus von $\frac{2}{3}$ auf $1-\delta$ verstärken, indem wir parallel zueinander und unabhängig voneinander $s = O(\ln(\frac{1}{\delta}))$ mal den BJKST-Algorithmus durchführen und den Median der berechneten Schätzwerte ausgeben.