

Big Data Analytics in Theorie und Praxis

Prof. Johann-Christoph Freytag, Ph.D. und Prof. Dr. Nicole Schweikardt
Sommersemester 2015

Theorie-Übungsblatt 4

Zu bearbeiten bis 9. Juli 2015

Aufgabe 1:

(40 Punkte)

Für einen Datenstrom x_1, \dots, x_n mit $x_i \in \{1, \dots, m\}$, $m > n$ und $x_i \neq x_j$ für $i, j \in \{1, \dots, n\}$ mit $i \neq j$ soll eine Näherung für den Median ermittelt werden. Diese Näherung des Medians soll durch eine zufällig, gleichverteilte Stichprobe S der Größe s aus dem Datenstrom bestimmt werden.

Sei $x_{\pi_1}, \dots, x_{\pi_n}$ der Datenstrom, wenn man ihn nach den Werten aufsteigend sortieren würde. Dies wird natürlich nicht durchgeführt, aber ist für die nun folgende Definition hilfreich.

Sei $\text{Pos}(x_{\pi_i}) := i$ für $i \in \{1, \dots, n\}$ die *Position* von x_{π_i} im sortierten Datenstrom. Für ungerade n ist offensichtlich der Median *Med* des Datenstroms das Element $x_{\pi_{\frac{n+1}{2}}}$, d.h. das Element mit Position $\frac{n+1}{2}$ im sortierten Datenstrom. Für gerade n wäre der Median *Med* das Element $x_{\pi_{\frac{n}{2}}}$ oder $x_{\pi_{\frac{n}{2}+1}}$. Zur Vereinfachung nehmen wir für all diese Fälle $\text{Pos}(\text{Med}) = \lceil \frac{n}{2} \rceil$ an, und in Rechnungen können Sie für $\text{Pos}(\text{Med})$ einfach $\frac{n}{2}$ einsetzen.

Bestimmen Sie eine möglichst kleine Zahl s für die Größe der zu ziehenden Stichprobe S , so dass s nur von ε und δ (für gegebene ε, δ mit $0 < \varepsilon < 1$ und $0 < \delta < 1$) und einer Konstanten c abhängt. Hierbei hängt c nicht von $\varepsilon, \delta, x_1, \dots, x_n, n, m$ oder s ab. Außerdem soll für den Median *med* der Stichprobe S Folgendes gelten:

$$\Pr \left(\text{Pos}(\text{med}) \in \left(\frac{1}{2} \pm \varepsilon \right) \cdot n \right) \geq 1 - \delta.$$

D.h., die Position vom Median *med* der Stichprobe im sortierten Datenstrom weicht von der Position des Medians des gesamten Datenstroms mit einer Wahrscheinlichkeit von mindestens $1 - \delta$ höchstens um $\varepsilon \cdot n$ ab.

Hinweis: Überlegen Sie, in welchen Fällen die Position des Medians der Stichprobe im sortierten Datenstrom um mehr als $\varepsilon \cdot n$ von $\frac{n}{2}$ abweicht und benutzen Sie dann die Chernoff-Schranke aus der Vorlesung für die einzelnen Fälle. Beachten Sie, dass die Chernoff-Schranke auch gilt, wenn die Betragsstriche weggelassen werden. Außerdem ist das ε in der Chernoff-Schranke nicht notwendigerweise dasselbe ε wie in dieser Aufgabenstellung.

Aufgabe 2:**(30 Punkte)**

Sei $m \in \mathbb{N}$ und sei M eine Primzahl mit $M > m$. Für alle $a, b \in \{0, 1, \dots, M-1\}$ sei $h_{a,b}$ die Funktion $h_{a,b} : \{1, \dots, m\} \rightarrow \{0, \dots, M-1\}$ mit

$$h_{a,b}(x) = ax + b \pmod{M}, \quad \text{für alle } x \in \{1, \dots, m\}.$$

Sei H die Familie all dieser Funktionen $h_{a,b}$, d.h.

$$H := \{ h_{a,b} : a, b \in \{0, \dots, M-1\} \}$$

Beweisen Sie, dass H paarweise unabhängig ist.

Zur Erinnerung: Eine Familie H von Funktionen von $\{1, \dots, m\}$ nach $\{0, \dots, M-1\}$ heißt *paarweise unabhängig* (oder *2-universell*), falls für alle $x_1, x_2 \in \{1, \dots, m\}$ mit $x_1 \neq x_2$ und alle $y_1, y_2 \in \{0, \dots, M-1\}$ gilt: Für ein zufällig, gleichverteilt aus H gewähltes h ist

$$\Pr(h(x_1) = y_1 \text{ und } h(x_2) = y_2) = \frac{1}{M^2}.$$

Aufgabe 3:**(20 Punkte)**

Beweisen Sie, dass für $A \geq 2$ paarweise unabhängige Zufallsvariablen Y_1, \dots, Y_A gilt, dass

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

Erläutern Sie Ihre Zwischenschritte.

Hinweis: Seien X und Y zwei unabhängige Zufallsvariablen und sei a eine beliebige Zahl. Dann gelten folgende Regeln:

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}((X - \mathbb{E}(X))^2), \\ \text{Var}(X) &= \mathbb{E}(X^2) - \mathbb{E}(X)^2, \\ \mathbb{E}(X + Y) &= \mathbb{E}(X) + \mathbb{E}(Y), \\ \mathbb{E}(X \cdot Y) &= \mathbb{E}(X) \cdot \mathbb{E}(Y), \\ \mathbb{E}(a \cdot X) &= a \cdot \mathbb{E}(X). \end{aligned}$$

Aufgabe 4:**(10 Punkte)**

Beweisen Sie, dass für jede Zufallsvariable Y und jede Zahl a gilt:

$$\text{Var}(a \cdot Y) = a^2 \cdot \text{Var}(Y).$$