

Datenströme

Prof. Dr. Nicole Schweikardt

Arbeitsgruppe Logik in der Informatik
Humboldt-Universität zu Berlin

Vorlesung „Big Data Analytics in Theorie und Praxis“

Datenströme



Situation:

- riesige Mengen von Daten
- automatisch generiert
- ständig kommen neue Daten hinzu

Beispiele:

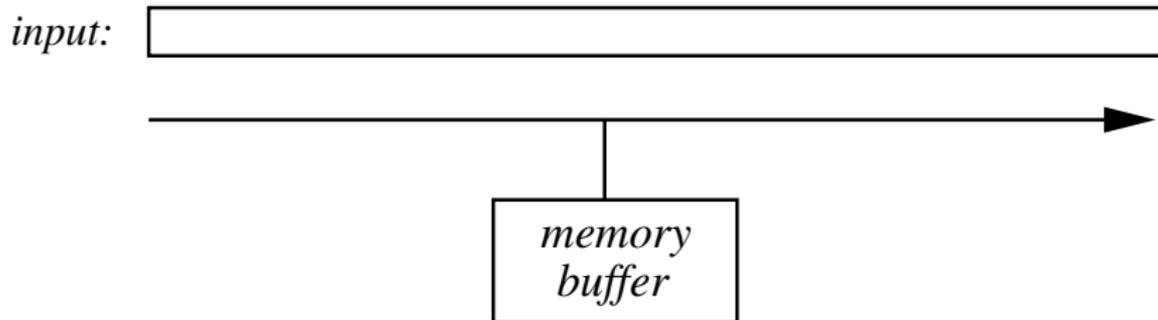
- Meteorologische Daten (Sensor-Netzwerke)
- Astronomische Daten
- Börsen-Transaktionen
- Netzwerk-Monitoring

Herausforderungen:

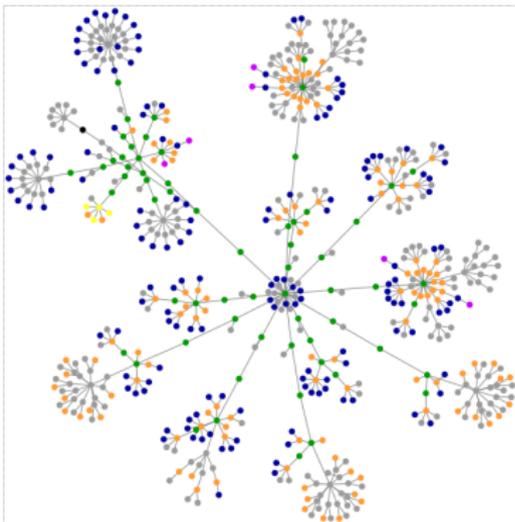
- Warten, bis alle Daten da sind, ist nicht möglich
 ~> Daten müssen “on-the-fly” verarbeitet werden
- Speichern der gesamten Daten ist zu teuer
 ~> speichere eine “Skizze”
- Daten können so schnell aufeinander folgen, dass man keine Gelegenheit hat, sich jeden einzelnen Datensatz anzusehen
 ~> ziehe “Stichproben”

Einmaliges Lesen des Datenstroms

Szenario:



Ein Beispiel aus dem Bereich Netzwerk-Monitoring



A sei ein Knoten im Internet.

Als Eingabe erhält A einen Strom von "Paketen"

$$p_1, p_2, p_3, p_4, \dots$$

Jedes Paket p_i enthält Informationen über

- ▶ die IP-Adresse, von der das Paket losgeschickt wurde,
- ▶ die IP-Adresse, an die das Paket geschickt werden soll,
- ▶ die eigentlichen Daten, die verschickt werden sollen

Frage: Wie viele verschiedene IP-Adressen haben seit heute morgen ein oder mehrere Pakete über Knoten A geleitet?

Problem: A will/kann nicht den gesamten Strom p_1, p_2, p_3, \dots zwischenspeichern.

Lösung: Nutze randomisierte Algorithmen, die eine gute Approximation für die Antwort auf obige Frage berechnen.

Beispiel: Suche nach der fehlenden Zahl

SUCHE NACH DER FEHLENDEN ZAHL

Eingabe: Strom $x_1, x_2, x_3, \dots, x_{n-1}$ von
 $n-1$ verschiedenen Zahlen aus $\{1, \dots, n\}$

Frage: Welche Zahl aus $\{1, \dots, n\}$ fehlt?

Naive Lösung: 2 5 1 3 4 8 6 \dots n benötigt n Speicher-Bits

1	2	3	4	5	6	7	8	\dots	n
✓	✓	✓	✓	✓	✓		✓	✓	✓

Bessere Lösung: Speichere Zwischensumme $O(\log n)$ Bits reichen aus

$$s := x_1 + x_2 + x_3 + x_4 + \dots + x_{n-1}$$

$$\text{fehlende Zahl} = \frac{n \cdot (n+1)}{2} - s$$

Untere Schranke: mindestens $\log n$ Bits sind nötig

Das MULTISSET-EQUALITY Problem

MULTISSET-EQUALITY

Eingabelänge: $N = O(m \cdot \log n)$ Bits

Eingabe: Zwei Multimengen $\{x_1, \dots, x_m\}$ und $\{y_1, \dots, y_m\}$ von Zahlen x_i, y_j aus $\{1, \dots, n\}$

Frage: Ist $\{x_1, \dots, x_m\} = \{y_1, \dots, y_m\}$?

D.h.: Kommen dieselben Zahlen mit denselben Vielfachheiten vor?

Satz:

Jeder *deterministische* Datenstrom-Algorithmus zur Lösung dieses Problems benötigt $\Omega(N)$ Speicher-Bits.

Beweismethode:

- Nutze Resultate aus der **Kommunikationskomplexität:**

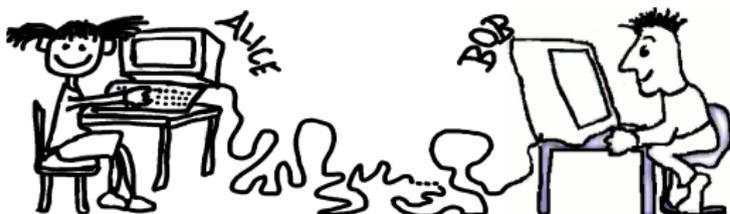
Zu entscheiden, ob zwei m -elementige Teilmengen von $\{1, \dots, n\}$ gleich sind, benötigt man mindestens $\log \binom{n}{m}$ Kommunikations-Bits.

- Falls $n = m^2$, so sind $\log \binom{n}{m} \geq m \cdot \log m$ Kommuniktions-Bits nötig, und die Eingabelänge der MULTISSET-EQUALITY Instanz ist $N = \Theta(m \cdot \log m)$.

Kommunikationskomplexität

Yaos 2-Parteien Kommunikations-Modell:

- 2 Spieler: Alice & Bob
- beide kennen eine Abbildung $f : A \times B \rightarrow \{0, 1\}$
- Alice sieht nur die Eingabe $a \in A$, Bob sieht nur die Eingabe $b \in B$
- beide wollen gemeinsam den Wert $f(a, b)$ berechnen
- Ziel: Tausche so wenig Kommunikations-Bits wie möglich aus



Satz: Zu entscheiden, ob zwei m -elementige Mengen

$$a = \{x_1, \dots, x_m\} \subseteq \{1, \dots, n\} \quad \text{und} \quad b = \{y_1, \dots, y_m\} \subseteq \{1, \dots, n\}$$

gleich sind, erfordert mindestens $\log \binom{n}{m}$ Kommunikations-Bits.

Das MULTISSET-EQUALITY Problem

MULTISSET-EQUALITY

Eingabelänge: $N = O(m \cdot \log n)$ Bits

Eingabe: Zwei Multimengen $\{x_1, \dots, x_m\}$ und $\{y_1, \dots, y_m\}$ von Zahlen x_i, y_j aus $\{1, \dots, n\}$

Frage: Ist $\{x_1, \dots, x_m\} = \{y_1, \dots, y_m\}$?

D.h.: Kommen dieselben Zahlen mit denselben Vielfachheiten vor?

Satz:

Jeder *deterministische* Datenstrom-Algorithmus zur Lösung dieses Problems benötigt $\Omega(N)$ Speicher-Bits.

Beweismethode:

- Nutze Resultate aus der **Kommunikationskomplexität:**

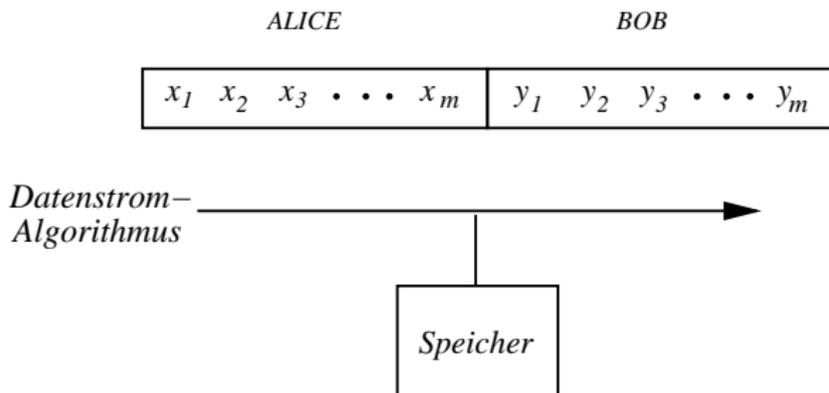
Zu entscheiden, ob zwei m -elementige Teilmengen von $\{1, \dots, n\}$ gleich sind, benötigt man mindestens $\log \binom{n}{m}$ Kommunikations-Bits.

- Falls $n = m^2$, so sind $\log \binom{n}{m} \geq m \cdot \log m$ Kommuniktions-Bits nötig, und die Eingabelänge der MULTISSET-EQUALITY Instanz ist $N = \Theta(m \cdot \log m)$.

Eine untere Schranke für MULTISSET-EQUALITY

Beweis (Fortsetzung):

- **Wir haben eben gesehen:** Eingabelänge: $N = \Theta(m \cdot \log m)$;
zum Lösen von MULTISSET-EQUALITY benötigte Kommunikations-Bits: $\geq m \cdot \log m$
- Ein deterministischer Datenstrom-Algorithmus, der MULTISSET-EQUALITY löst und dabei B Speicher-Bits nutzt, führt zu einem Kommunikations-Protokoll, das B Kommunikations-Bits nutzt:



- Somit: Untere Schranke für Kommunikationskomplexität \rightsquigarrow untere Schranke für Speicherplatz des Datenstrom-Algorithmus

Ein randomisierter Datenstrom-Algorithmus

Satz: (Grohe, Hernich, Schweikardt, PODS'06)

Das MULTISSET-EQUALITY Problem kann im folgenden Sinn durch einen **randomisierten** Datenstrom-Algorithmus, der $O(\log N)$ Speicher-Bits nutzt, gelöst werden: Sind m, n und ein Strom von Zahlen $a_1, \dots, a_m, b_1, \dots, b_m$ aus $\{1, \dots, n\}$ gegeben, so

- akzeptiert der Algo mit Wk 1 falls $\{\{a_1, \dots, a_m\}\} = \{\{b_1, \dots, b_m\}\}$
- verwirft der Algo mit Wk ≥ 0.9 falls $\{\{a_1, \dots, a_m\}\} \neq \{\{b_1, \dots, b_m\}\}$.

Grundprinzip des Algorithmus: Benutze "Fingerprinting"-Techniken:

- Repräsentiere $\{\{a_1, \dots, a_m\}\}$ durch ein Polynom $f(x) := \sum_{i=1}^m x^{a_i}$.
- Repräsentiere $\{\{b_1, \dots, b_m\}\}$ durch ein Polynom $g(x) := \sum_{i=1}^m x^{b_i}$.
- Wähle eine Zufallszahl z und überprüfe, ob $f(z) = g(z)$.
- Akzeptiere, falls $f(z) = g(z)$. Ansonsten verwirfe.



Falls $\{\{a_1, \dots, a_m\}\} = \{\{b_1, \dots, b_m\}\}$, so ist $f(x) = g(x)$, und daher akzeptiert der Algorithmus stets. Falls $\{\{a_1, \dots, a_m\}\} \neq \{\{b_1, \dots, b_m\}\}$, dann gibt es nur $\text{Grad}(f-g)$ viele Zahlen z mit $f(z) = g(z)$, und daher verwirft der Algorithmus mit hoher Wk.