

Basic algorithmic techniques for data streams

Mariano Zelke

Goethe-Universität Frankfurt am Main
Institut für Informatik
Arbeitsgruppe Theorie Komplexer Systeme

DEIS'10, Dagstuhl, Nov. 2010

Contents

Basic Definitions

Sampling

- General Idea

- Reservoir Sampling

- Sampling Applications

Advanced Sampling

- AMS Sampling

- Sliding Window Sampling

Count-Min Sketch

Data Stream Model

Data Stream Model

- ▶ **Stream**: m elements of some universe of size n

Data Stream Model

- ▶ **Stream**: m elements of some universe of size n

$$a_1, a_2, a_3, \dots, a_m \qquad a_i \in \{1, \dots, n\}$$

Data Stream Model

- ▶ **Stream:** m elements of some universe of size n
 $a_1, a_2, a_3, \dots, a_m$ $a_i \in \{1, \dots, n\}$
- ▶ **Goal:** gain information about stream

Data Stream Model

- ▶ **Stream:** m elements of some universe of size n

$$a_1, a_2, a_3, \dots, a_m \quad a_i \in \{1, \dots, n\}$$

- ▶ **Goal:** gain information about stream
statistical information (median, frequency moments, ...),
longest increasing subsequence,...

Data Stream Model

- ▶ **Stream**: m elements of some universe of size n

$$a_1, a_2, a_3, \dots, a_m \quad a_i \in \{1, \dots, n\}$$

- ▶ **Goal**: gain information about stream
statistical information (median, frequency moments, ...),
longest increasing subsequence,...
- ▶ **But**: algorithms are restricted to

Data Stream Model

- ▶ **Stream:** m elements of some universe of size n

$$a_1, a_2, a_3, \dots, a_m \quad a_i \in \{1, \dots, n\}$$

- ▶ **Goal:** gain information about stream
statistical information (median, frequency moments, ...),
longest increasing subsequence,...
- ▶ **But:** algorithms are restricted to
 - ▶ sequential access to items in stream

Data Stream Model

- ▶ **Stream:** m elements of some universe of size n

$$a_1, a_2, a_3, \dots, a_m \quad a_i \in \{1, \dots, n\}$$

- ▶ **Goal:** gain information about stream
statistical information (median, frequency moments, ...),
longest increasing subsequence,...
- ▶ **But:** algorithms are restricted to
 - ▶ sequential access to items in stream
 - ▶ limited memory, sublinear in m and n

Sampling

General idea:

Sampling

General idea:

- ▶ sample (= select) items from stream according to some rule

Sampling

General idea:

- ▶ sample (= select) items from stream according to some rule
- ▶ rule is randomized or deterministic

Sampling

General idea:

- ▶ sample (= select) items from stream according to some rule
- ▶ rule is randomized or deterministic

Sampling

General idea:

- ▶ sample (= select) items from stream according to some rule
- ▶ rule is randomized or deterministic
- ▶ use sampled items to get information about whole stream

Sampling

General idea:

- ▶ sample (= select) items from stream according to some rule
- ▶ rule is randomized or deterministic
- ▶ use sampled items to get information about whole stream
- ▶ only need to store sampled items
⇒ good for streaming

Easy Starter

Sample out a single item uniformly at random from the stream

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

↓

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

↓

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$

↓

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$


↓

- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$




- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item a_r

Easy Starter

Sample out a single item uniformly at random from the stream

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_r, \dots, a_m$



- ▶ Easy if m is known in advance:
- ▶ Pick a random number $r \in \{1, 2, \dots, m\}$
- ▶ Go over the stream and snatch out sampled item a_r

But what if m is not known in advance?

Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream
without knowing its length

Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream
without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream
without knowing its length

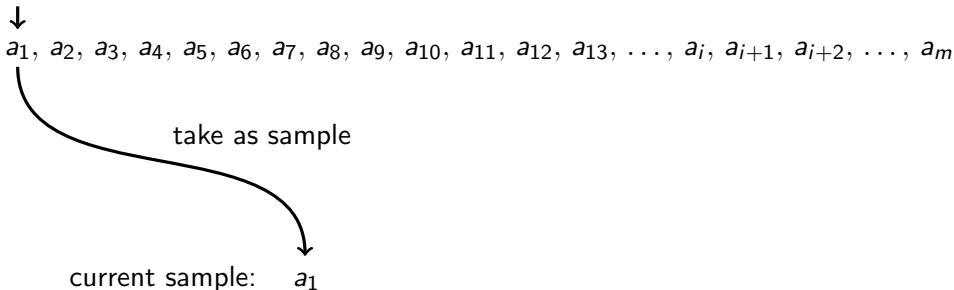
$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

current sample:

Reservoir Sampling

[J. Vitter '85]

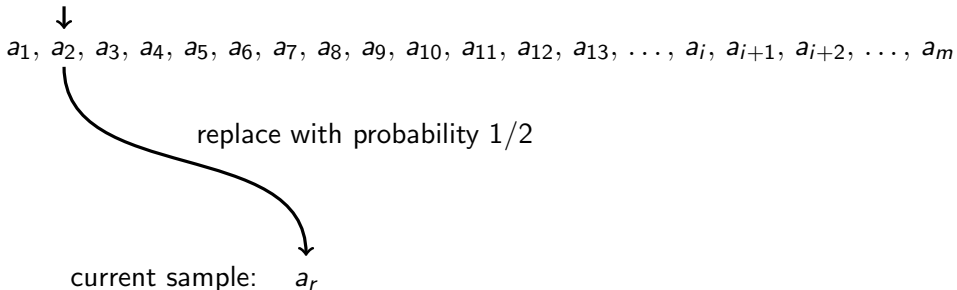
Sample out a single item uniformly at random from the stream without knowing its length



Reservoir Sampling

[J. Vitter '85]

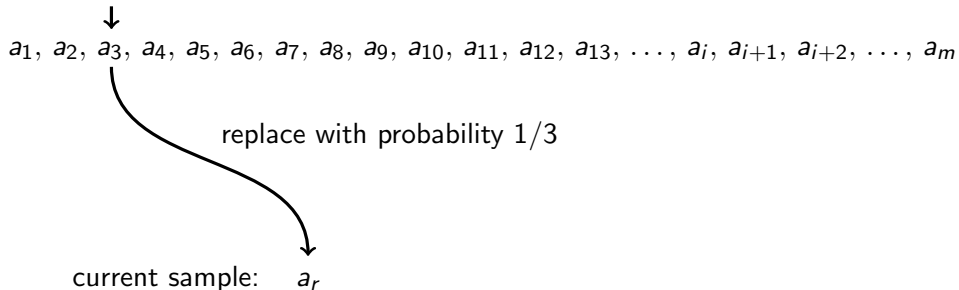
Sample out a single item uniformly at random from the stream without knowing its length



Reservoir Sampling

[J. Vitter '85]

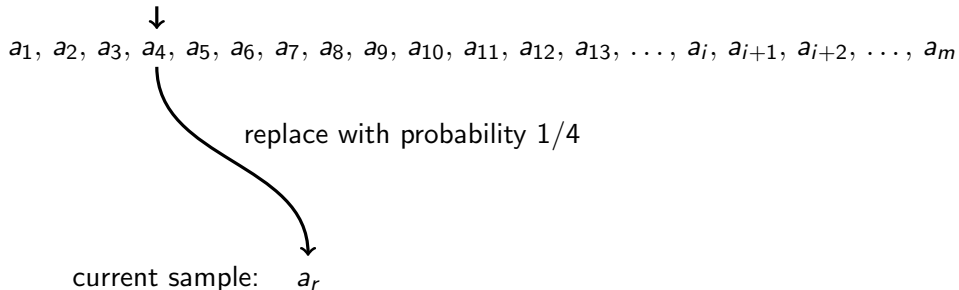
Sample out a single item uniformly at random from the stream without knowing its length



Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream without knowing its length



Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

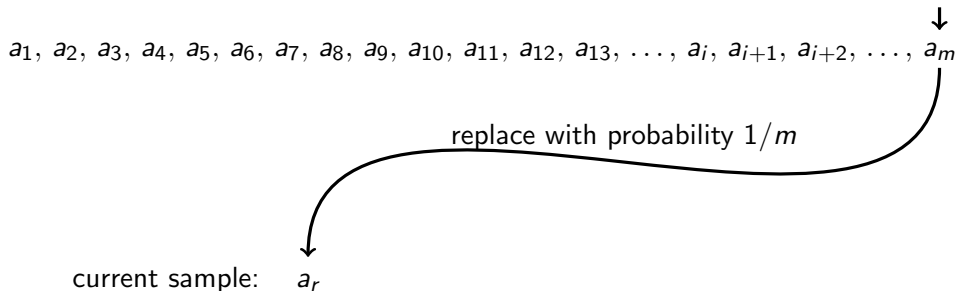
replace with probability $1/i$

current sample: a_r

Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream without knowing its length



Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream
without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

final sample: a_r

Sample out a single item uniformly at random from the stream
without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

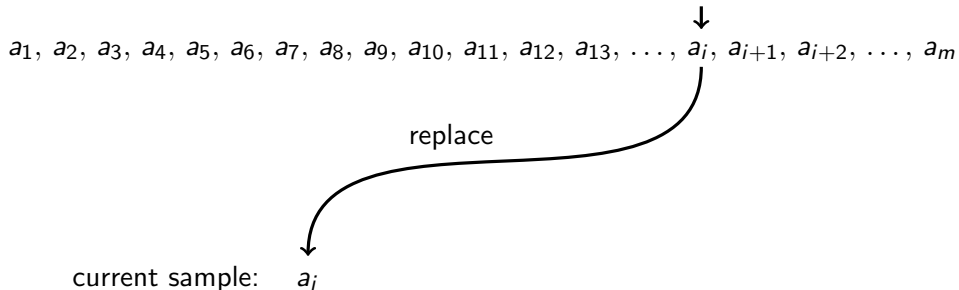
current sample:

$Pr[\text{final sample } a_i] =$

Reservoir Sampling

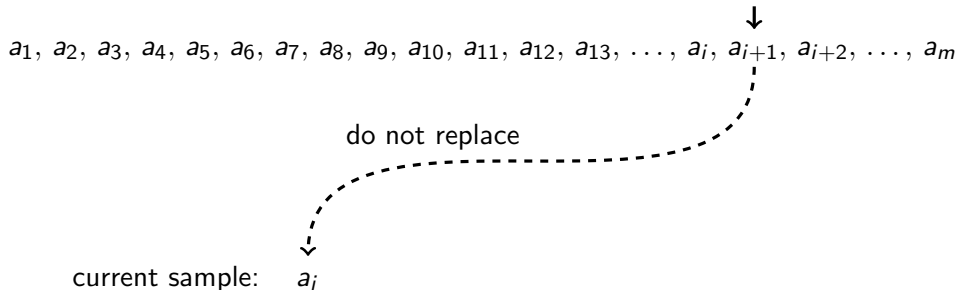
[J. Vitter '85]

Sample out a single item uniformly at random from the stream without knowing its length



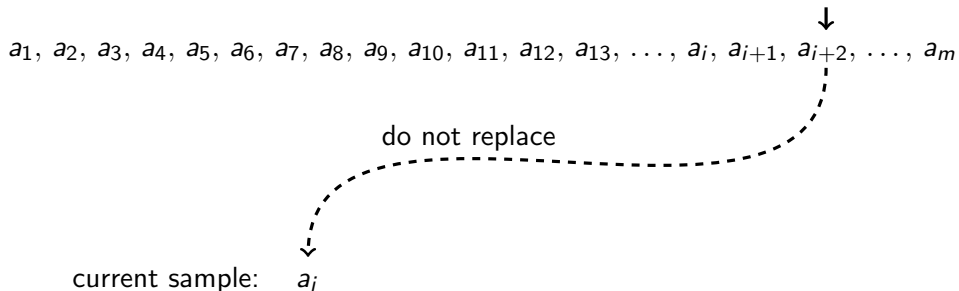
$$Pr[\text{final sample } a_i] = \frac{1}{i}$$

Sample out a single item uniformly at random from the stream without knowing its length



$$Pr[\text{final sample } a_i] = \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right)$$

Sample out a single item uniformly at random from the stream without knowing its length

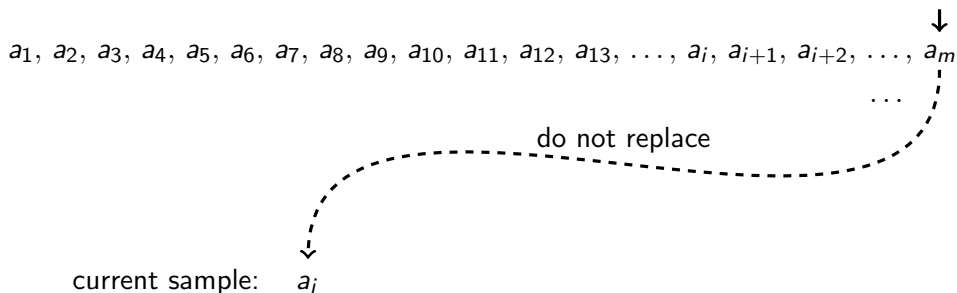


$$Pr[\text{final sample } a_i] = \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right) \times \left(1 - \frac{1}{i+2}\right)$$

Reservoir Sampling

[J. Vitter '85]

Sample out a single item uniformly at random from the stream without knowing its length



$$Pr[\text{final sample } a_i] = \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right) \times \left(1 - \frac{1}{i+2}\right) \times \dots \times \left(1 - \frac{1}{m}\right)$$

Sample out a single item uniformly at random from the stream
without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

final sample: a_i

$$\begin{aligned} Pr[\text{final sample } a_i] &= \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right) \times \left(1 - \frac{1}{i+2}\right) \times \dots \times \left(1 - \frac{1}{m}\right) \\ &= \frac{1}{i} \times \frac{i}{i+1} \times \frac{i+1}{i+2} \times \dots \times \frac{m-1}{m} \end{aligned}$$

Sample out a single item uniformly at random from the stream
without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

final sample: a_i

$$\begin{aligned} Pr[\text{final sample } a_i] &= \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right) \times \left(1 - \frac{1}{i+2}\right) \times \dots \times \left(1 - \frac{1}{m}\right) \\ &= \frac{1}{i} \times \frac{i}{i+1} \times \frac{i+1}{i+2} \times \dots \times \frac{m-1}{m} \end{aligned}$$

Sample out a single item uniformly at random from the stream
without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

final sample: a_i

$$\begin{aligned} Pr[\text{final sample } a_i] &= \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right) \times \left(1 - \frac{1}{i+2}\right) \times \dots \times \left(1 - \frac{1}{m}\right) \\ &= \frac{1}{m} \end{aligned}$$

Sample out **a single item** uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, a_{i+1}, a_{i+2}, \dots, a_m$

final sample: a_i

$$\begin{aligned} Pr[\text{final sample } a_i] &= \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right) \times \left(1 - \frac{1}{i+2}\right) \times \dots \times \left(1 - \frac{1}{m}\right) \\ &= \frac{1}{m} \end{aligned}$$

Reservoir Sampling

[J. Vitter '85]

Sample out **several items** uniformly at random from the stream without knowing its length

Sample out **several items** uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Reservoir Sampling

[J. Vitter '85]

Sample out **several items** uniformly at random from the stream
without knowing its length

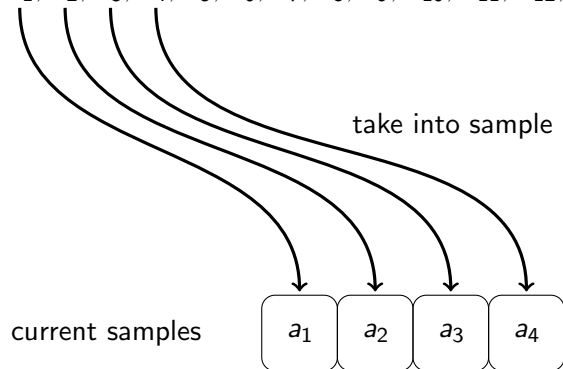
$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current samples



Sample out **several items** uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

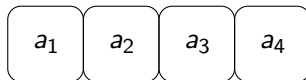


Sample out **several items** uniformly at random from the stream without knowing its length

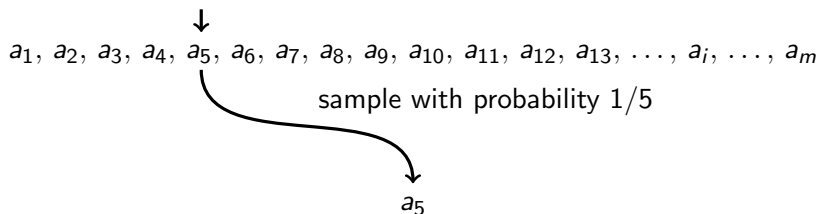


$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

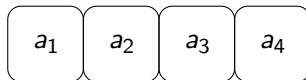
current samples



Sample out **several items** uniformly at random from the stream without knowing its length



current samples



Reservoir Sampling

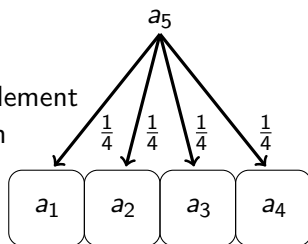
[J. Vitter '85]

Sample out **several items** uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$
sample with probability $1/5$

replace a sampled element uniformly at random

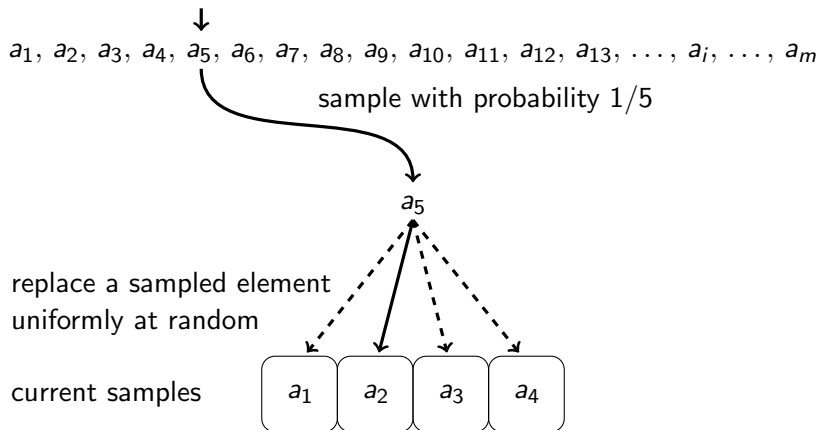
current samples



Reservoir Sampling

[J. Vitter '85]

Sample out **several items** uniformly at random from the stream without knowing its length

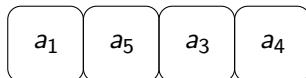


Sample out **several items** uniformly at random from the stream without knowing its length



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current samples

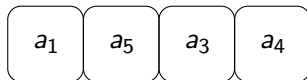


Sample out **several items** uniformly at random from the stream without knowing its length

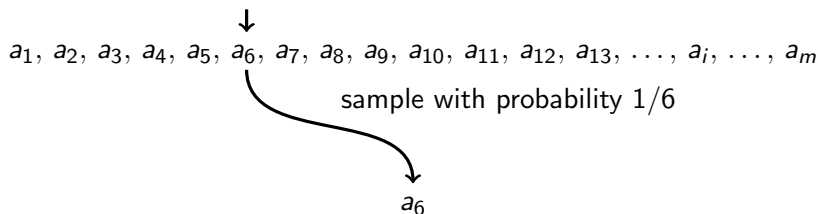


$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

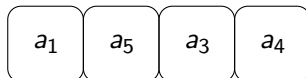
current samples



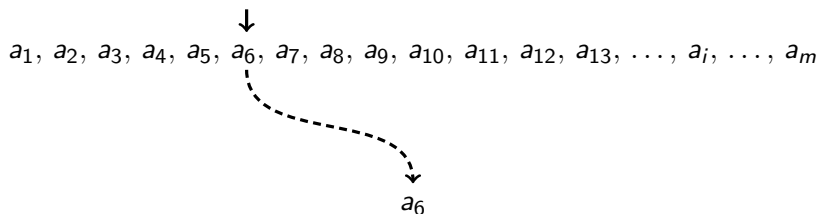
Sample out **several items** uniformly at random from the stream without knowing its length



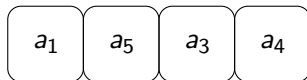
current samples



Sample out **several items** uniformly at random from the stream without knowing its length



current samples

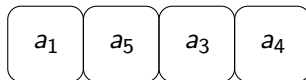


Sample out **several items** uniformly at random from the stream without knowing its length



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current samples

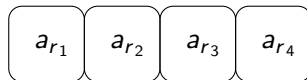


Sample out **several items** uniformly at random from the stream without knowing its length

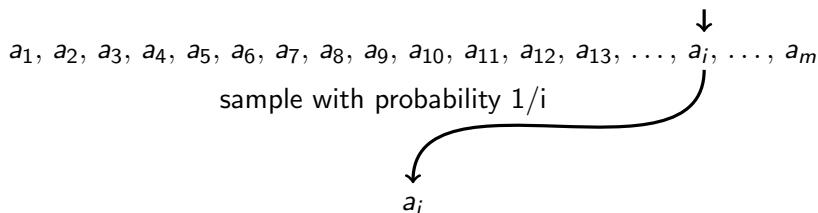
$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

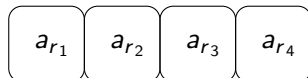
current samples



Sample out **several items** uniformly at random from the stream without knowing its length



current samples



Reservoir Sampling

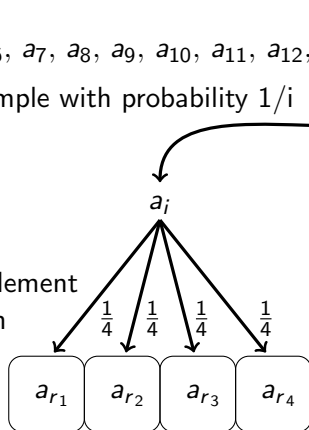
[J. Vitter '85]

Sample out **several items** uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$
sample with probability $1/i$

replace a sampled element
uniformly at random

current samples

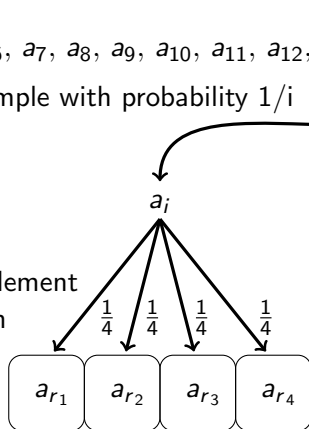


Sample out **several items** uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$
sample with probability $1/i$

replace a sampled element uniformly at random

current samples



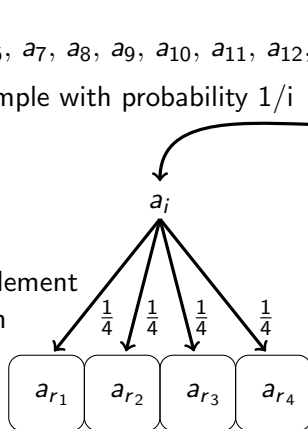
Memory usage for sampling k items:

Sample out **several items** uniformly at random from the stream without knowing its length

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$
sample with probability $1/i$

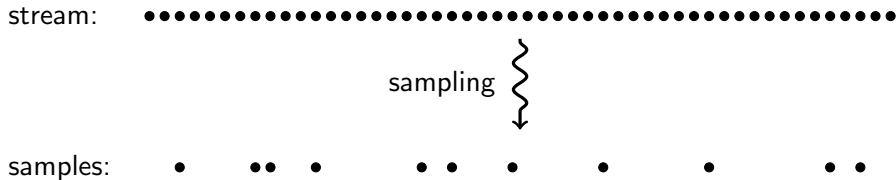
replace a sampled element uniformly at random

current samples



Memory usage for sampling k items: $k \cdot \log n$

Sampling Applications

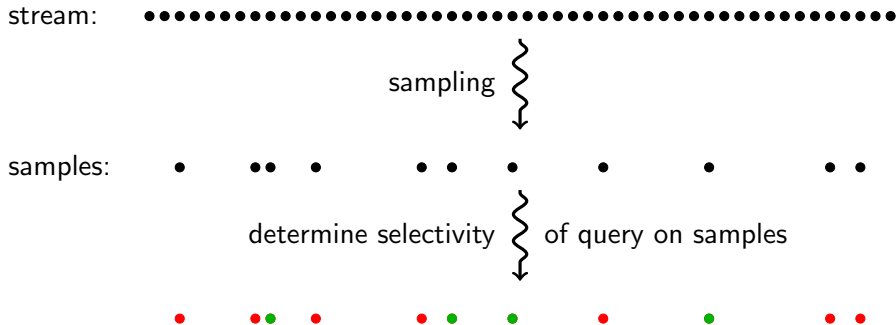


Sampling Applications

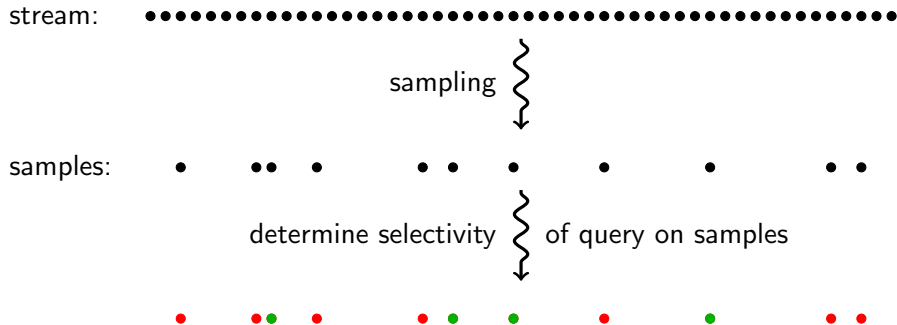


Goal: determine query selectivity on items of stream
(assume query to be invariant of stream order)

Sampling Applications

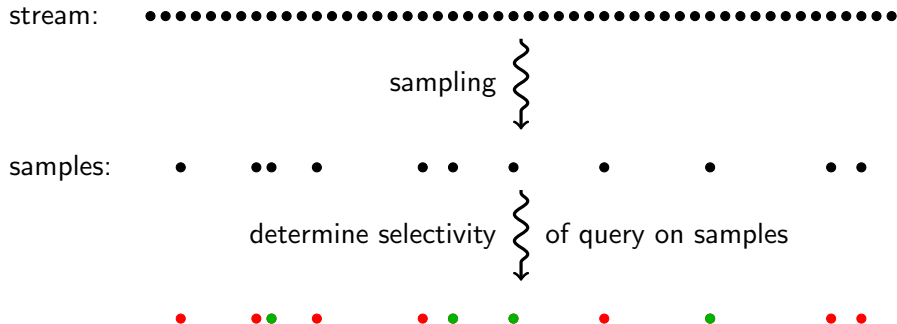


Sampling Applications



To get $(1 \pm \epsilon)$ -estimate with probability $1 - \delta$: What sample size s ?

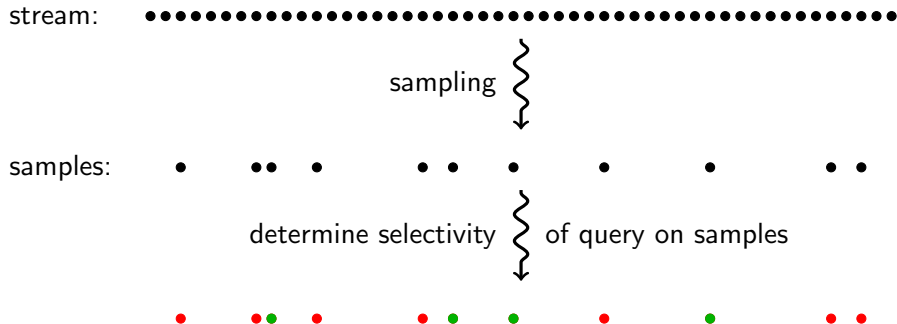
Sampling Applications



To get $(1 \pm \epsilon)$ -estimate with probability $1 - \delta$: What sample size s ?

Query selects $\frac{m}{c}$ stream items $\Rightarrow \text{Exp}[s^+] = \frac{s}{c}$

Sampling Applications

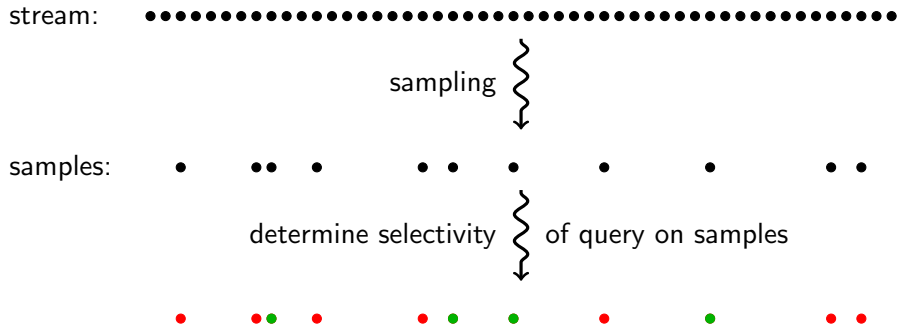


To get $(1 \pm \epsilon)$ -estimate with probability $1 - \delta$: What sample size s ?

Query selects $\frac{m}{c}$ stream items $\Rightarrow \text{Exp}[s^+] = \frac{s}{c}$

Chernoff-Hoeffding-Ineq.: $Pr[|s^+ - \text{Exp}[s^+]| > \epsilon \cdot s^+] \leq e^{-\Theta(\epsilon^2 s)}$

Sampling Applications



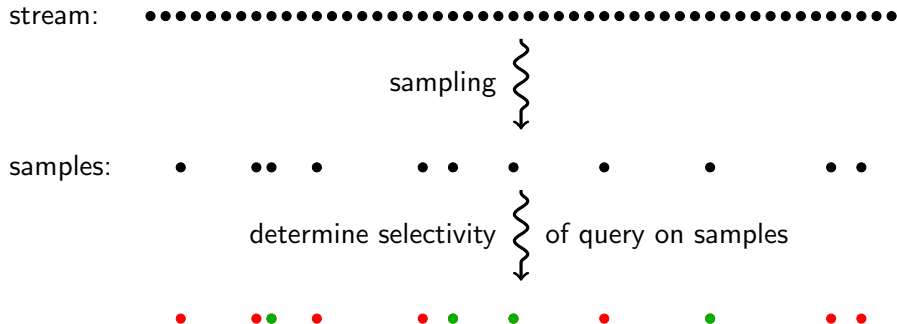
To get $(1 \pm \varepsilon)$ -estimate with probability $1 - \delta$: What sample size s ?

Query selects $\frac{m}{c}$ stream items $\Rightarrow \text{Exp}[s^+] = \frac{s}{c}$

Chernoff-Hoeffding-Ineq.: $Pr[|s^+ - \text{Exp}[s^+]| > \varepsilon \cdot s^+] \leq e^{-\Theta(\varepsilon^2 s)}$

sample $\mathcal{O}\left(\frac{1}{\varepsilon^2} \cdot \log \frac{1}{\delta}\right)$ items

Sampling Applications



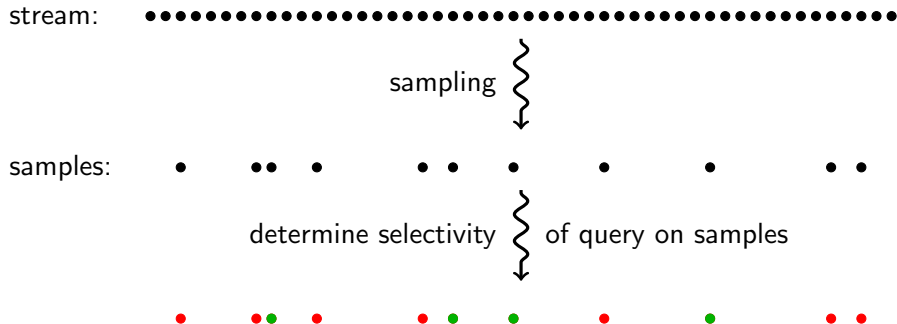
To get $(1 \pm \varepsilon)$ -estimate with probability $1 - \delta$: What sample size s ?

Query selects $\frac{m}{c}$ stream items $\Rightarrow \text{Exp}[s^+] = \frac{s}{c}$

Chernoff-Hoeffding-Ineq.: $Pr[|s^+ - \text{Exp}[s^+]| > \varepsilon \cdot s^+] \leq e^{-\Theta(\varepsilon^2 s)}$

Memory usage: $\mathcal{O}\left(\frac{1}{\varepsilon^2} \cdot \log \frac{1}{\delta}\right)$ items $\times \log n$ bits

Sampling Applications



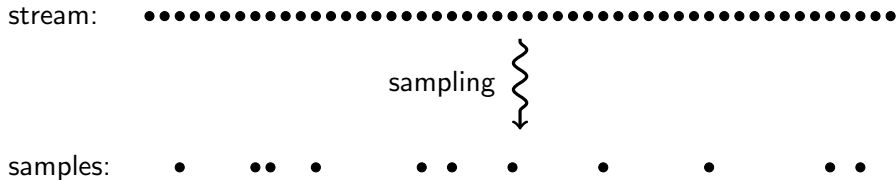
To get $(1 \pm \varepsilon)$ -estimate with probability $1 - \delta$: What sample size s ?

Query selects $\frac{m}{c}$ stream items $\Rightarrow \text{Exp}[s^+] = \frac{s}{c}$

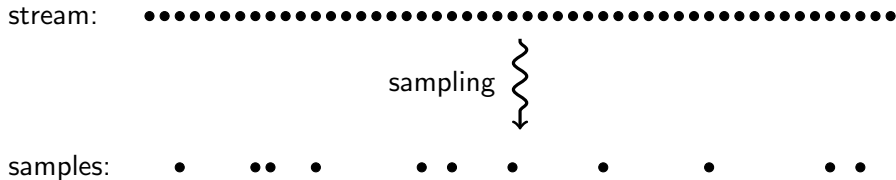
Chernoff-Hoeffding-Ineq.: $Pr[|s^+ - \text{Exp}[s^+]| > \varepsilon \cdot s^+] \leq e^{-\Theta(\varepsilon^2 s)}$

Memory usage: $\mathcal{O}\left(\frac{1}{\varepsilon^2} \cdot \log \frac{1}{\delta} \cdot \log n\right)$ bits

Sampling Applications

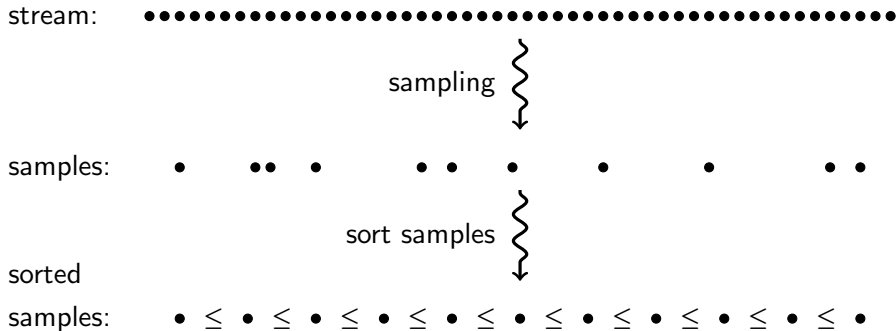


Sampling Applications

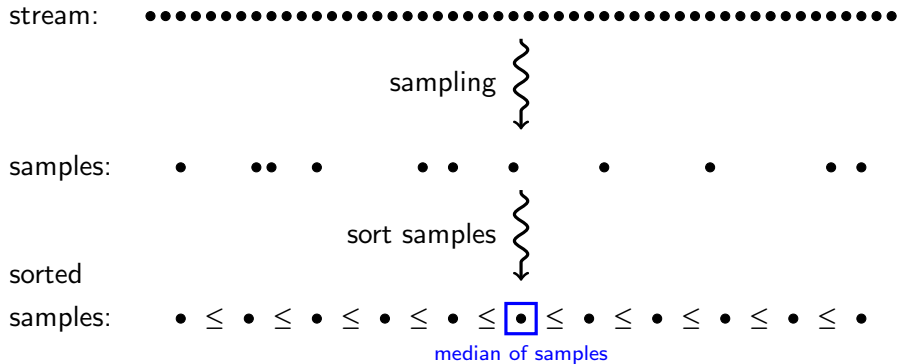


Goal: find the median of the stream

Sampling Applications

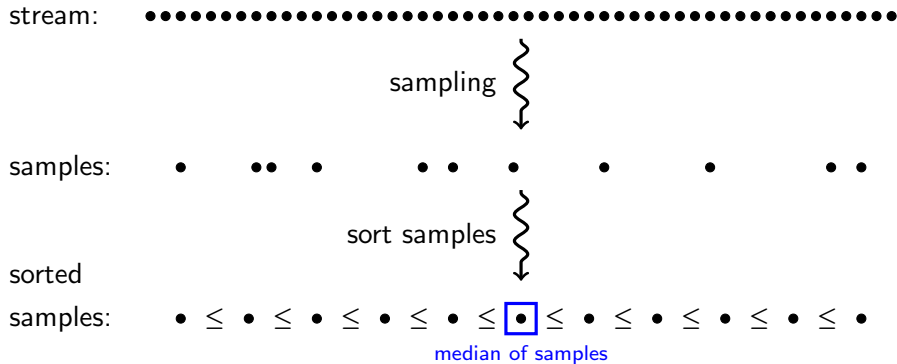


Sampling Applications



Pick median of samples as an estimate of stream's median

Sampling Applications



Pick median of samples as an estimate of stream's median

To get $(1 \pm \epsilon)$ -estimate with probability $1 - \delta$:

$$\text{sample } \mathcal{O}\left(\frac{1}{\epsilon^2} \cdot \log \frac{1}{\delta}\right) \text{ items}$$

AMS Sampling

[Alon, Matias, Szegedy '96]

AMS Sampling

[Alon, Matias, Szegedy '96]

Stream $a_1, a_2, a_3, \dots, a_m$

AMS Sampling

[Alon, Matias, Szegedy '96]

Stream $a_1, a_2, a_3, \dots, a_m$

Frequency of an item: $f_i = |\{j : a_j = i\}|$

Stream $a_1, a_2, a_3, \dots, a_m$

Frequency of an item: $f_i = |\{j : a_j = i\}|$

Example: stream 2, 3, 3, 2, 3, 1, 2, 3

$$f_1 = 1, \quad f_2 = 3, \quad f_3 = 4$$

Stream $a_1, a_2, a_3, \dots, a_m$

Frequency of an item: $f_i = |\{j : a_j = i\}|$

*k*th frequency moment $F_k = \sum_{i=1}^n f_i^k$

Stream $a_1, a_2, a_3, \dots, a_m$

Frequency of an item: $f_i = |\{j : a_j = i\}|$

k th frequency moment $F_k = \sum_{i=1}^n f_i^k$

Frequency moments provide useful statistics:

- ▶ F_0 : number of distinct elements in stream
- ▶ F_1 : length of stream, m
- ▶ F_2 : size of self join
- ▶ $F_k, k \geq 2$: skew of distribution

Stream $a_1, a_2, a_3, \dots, a_m$

Frequency of an item: $f_i = |\{j : a_j = i\}|$

*k*th frequency moment $F_k = \sum_{i=1}^n f_i^k$

Frequency moments provide useful statistics:

- ▶ F_0 : number of distinct elements in stream
- ▶ F_1 : length of stream, m
- ▶ F_2 : size of self join
- ▶ $F_k, k \geq 2$: skew of distribution

Trivial determination of F_k : maintain counters for each f_i

Stream $a_1, a_2, a_3, \dots, a_m$

Frequency of an item: $f_i = |\{j : a_j = i\}|$

k th frequency moment $F_k = \sum_{i=1}^n f_i^k$

Frequency moments provide useful statistics:

- ▶ F_0 : number of distinct elements in stream
- ▶ F_1 : length of stream, m
- ▶ F_2 : size of self join
- ▶ $F_k, k \geq 2$: skew of distribution

Trivial determination of F_k : maintain counters for each f_i

$\Rightarrow \Omega(n)$ bits needed

AMS Sampling (for estimating F_k)

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$

Example: stream 2, 3, 3, 2, 3, 1, 2, 3

$$a_j = 3, r = 3$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$: $Exp[m(r^k - (r - 1)^k)]$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$: $Exp[m(r^k - (r - 1)^k)]$

$$= \left[m(f_1^k - (f_1 - 1)^k) \right]$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$:

$$\text{Exp}[m(r^k - (r - 1)^k)]$$
$$= \left[m(f_1^k - (f_1 - 1)^k) + m((f_1 - 1)^k - (f_1 - 2)^k) \right]$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$: $Exp[m(r^k - (r - 1)^k)]$

$$= \left[m(f_1^k - (f_1 - 1)^k) + m((f_1 - 1)^k - (f_1 - 2)^k) + \dots + m(2^k - 1^k) \right]$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$: $Exp[m(r^k - (r - 1)^k)]$

$$= \left[m(f_1^k - (f_1 - 1)^k) + m((f_1 - 1)^k - (f_1 - 2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right]$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$: $Exp[m(r^k - (r - 1)^k)]$

$$= \left[m(f_1^k - (f_1 - 1)^k) + m((f_1 - 1)^k - (f_1 - 2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right. \\ \left. + m(f_2^k - (f_2 - 1)^k) + m((f_2 - 1)^k - (f_2 - 2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right]$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$: $Exp[m(r^k - (r - 1)^k)]$

$$= \left[m(f_1^k - (f_1 - 1)^k) + m((f_1 - 1)^k - (f_1 - 2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right. \\ \left. + m(f_2^k - (f_2 - 1)^k) + m((f_2 - 1)^k - (f_2 - 2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right. \\ \dots \\ \left. + m(f_n^k - (f_n - 1)^k) + m((f_n - 1)^k - (f_n - 2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right]$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r-1)^k)$$

for $k \geq 1$: $Exp[m(r^k - (r-1)^k)]$

$$\begin{aligned}
 &= \left[m(f_1^k - (f_1-1)^k) + m((f_1-1)^k - (f_1-2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right. \\
 &\quad + m(f_2^k - (f_2-1)^k) + m((f_2-1)^k - (f_2-2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \\
 &\quad \dots \\
 &\quad \left. + m(f_n^k - (f_n-1)^k) + m((f_n-1)^k - (f_n-2)^k) + \dots + m(2^k - 1^k) + m \cdot 1^k \right] \frac{1}{m}
 \end{aligned}$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$:

$$\text{Exp}[m(r^k - (r - 1)^k)]$$

$$= f_1^k + f_2^k + f_3^k + \dots + f_n^k$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

for $k \geq 1$:

$$\text{Exp}[m(r^k - (r - 1)^k)]$$

$$= f_1^k + f_2^k + f_3^k + \dots + f_n^k = F_k$$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

To get $(1 \pm \varepsilon)$ -estimate of F_k with probability $1 - \delta$:

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

To get $(1 \pm \varepsilon)$ -estimate of F_k with probability $1 - \delta$:

Run $\mathcal{O}\left(\frac{n^{1-1/k}}{\varepsilon^2}\right)$ parallel instances, take average A

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

To get $(1 \pm \varepsilon)$ -estimate of F_k with probability $1 - \delta$:

Run $\mathcal{O}\left(\frac{n^{1-1/k}}{\varepsilon^2}\right)$ parallel instances, take average A

\Rightarrow Chebyshev-Ineq.: $Pr[|A - F_k| > \varepsilon \cdot F_k] \leq p < \frac{1}{2}$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

To get $(1 \pm \varepsilon)$ -estimate of F_k with probability $1 - \delta$:

Run $\mathcal{O}\left(\frac{n^{1-1/k}}{\varepsilon^2}\right)$ parallel instances, take average A

\Rightarrow Chebyshev-Ineq.: $Pr[|A - F_k| > \varepsilon \cdot F_k] \leq p < \frac{1}{2}$

Take median M over $\mathcal{O}\left(\log \frac{1}{\delta}\right)$ such averages

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

To get $(1 \pm \varepsilon)$ -estimate of F_k with probability $1 - \delta$:

Run $\mathcal{O}\left(\frac{n^{1-1/k}}{\varepsilon^2}\right)$ parallel instances, take average A

\Rightarrow Chebyshev-Ineq.: $Pr[|A - F_k| > \varepsilon \cdot F_k] \leq p < \frac{1}{2}$

Take median M over $\mathcal{O}\left(\log \frac{1}{\delta}\right)$ such averages

\Rightarrow Chernoff-Ineq.: $Pr[|M - F_k| > \varepsilon \cdot F_k] \leq \delta$

AMS Sampling (for estimating F_k)

1. Pick random item a_j from stream
2. Compute $r = |\{j' : j' \geq j, a_{j'} = a_j\}|$
3. At the end of the stream calculate

$$m(r^k - (r - 1)^k)$$

To get $(1 \pm \varepsilon)$ -estimate of F_k with probability $1 - \delta$:

Run $\mathcal{O}\left(\frac{n^{1-1/k}}{\varepsilon^2}\right)$ parallel instances, take average A

\Rightarrow Chebyshev-Ineq.: $Pr[|A - F_k| > \varepsilon \cdot F_k] \leq p < \frac{1}{2}$

Take median M over $\mathcal{O}\left(\log \frac{1}{\delta}\right)$ such averages

\Rightarrow Chernoff-Ineq.: $Pr[|M - F_k| > \varepsilon \cdot F_k] \leq \delta$

Memory consumption: $\mathcal{O}\left(\frac{n^{1-1/k}}{\varepsilon^2} \cdot \log \frac{1}{\delta} \cdot (\log n + \log m)\right)$ bits

Sliding Window Sampling

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

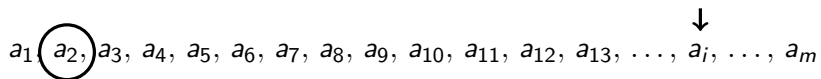
Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

↓

Sliding Window Sampling

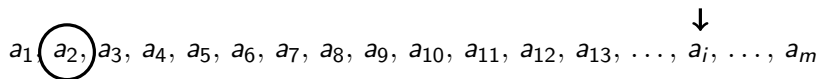
$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



current sample: a_2

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

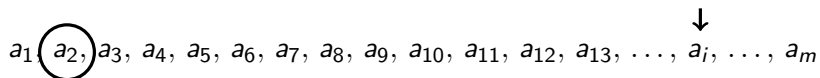


current sample: a_2

Current sample might be "far away" from actual item

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



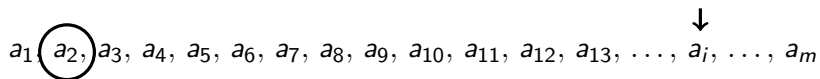
current sample: a_2

Current sample might be "far away" from actual item

- fine for some applications
- for others only w recent items matter

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$



current sample: a_2

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- \Rightarrow only consider items in a sliding window of size w

Sliding Window Sampling

↓
 a_1 , $a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

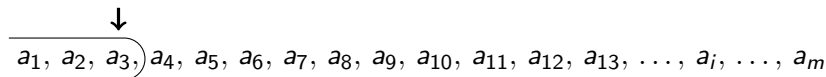
Sliding Window Sampling

↓
 $a_1, a_2,$ $a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

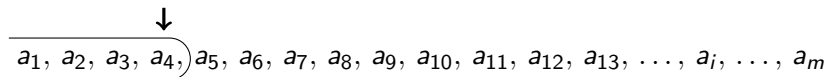
Sliding Window Sampling



Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

Sliding Window Sampling



Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- \Rightarrow only consider items in a sliding window of size w

Sliding Window Sampling



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- \Rightarrow only consider items in a sliding window of size w

Sliding Window Sampling

↓

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

Sliding Window Sampling

↓

$a_1, a_2, (a_3, a_4, a_5, a_6, a_7), a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

Sliding Window Sampling

↓

$a_1, a_2, a_3, (a_4, a_5, a_6, a_7, a_8), a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

Sliding Window Sampling

$a_1, a_2, a_3, a_4, \underbrace{a_5, a_6, a_7, a_8, a_9}, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, \underbrace{a_6, a_7, a_8, a_9, a_{10}}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

Current sample might be "far away" from actual item

- fine for some applications
 - for others only w recent items matter
- ⇒ only consider items in a sliding window of size w

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, \underbrace{a_6, a_7, a_8, a_9, a_{10}}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

Goal: sample an item uniformly from sliding window of size w

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, \underbrace{a_6, a_7, a_8, a_9, a_{10}}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

Goal: sample an item uniformly from sliding window of size w

Trivial: - memorize whole window content $\Rightarrow w \cdot \log n$ bits

- impractical if w is large

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, \underbrace{a_6, a_7, a_8, a_9, a_{10}}_{\downarrow}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

Goal: sample an item uniformly from sliding window of size w

Trivial: - memorize whole window content $\Rightarrow w \cdot \log n$ bits

- impractical if w is large

Better idea:

Algorithm:

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, \underbrace{a_6, a_7, a_8, a_9, a_{10}}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

Goal: sample an item uniformly from sliding window of size w

Trivial: - memorize whole window content $\Rightarrow w \cdot \log n$ bits

- impractical if w is large

Better idea:

Algorithm: 1: For each a_i pick random value $r_i \in (0, 1)$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, \underbrace{a_6, a_7, a_8, a_9, a_{10}}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

Goal: sample an item uniformly from sliding window of size w

Trivial: - memorize whole window content $\Rightarrow w \cdot \log n$ bits

- impractical if w is large

Better idea:

Algorithm: 1: For each a_i pick random value $r_i \in (0, 1)$

2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, \underbrace{a_6, a_7, a_8, a_9, a_{10}}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

↓

Goal: sample an item uniformly from sliding window of size w

Trivial: - memorize whole window content $\Rightarrow w \cdot \log n$ bits

- impractical if w is large

Better idea:

- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 a_1 , $a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_1 = 0.2$

- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling



$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_1 = 0.2$



- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 a_1 , $a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_1 = 0.2$

current
sample



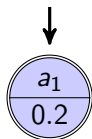
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_2 = 0.4$

current
sample



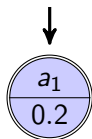
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_2 = 0.4$

current
sample



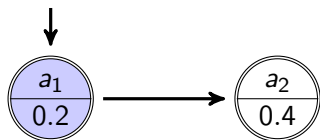
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_2 = 0.4$

current
sample



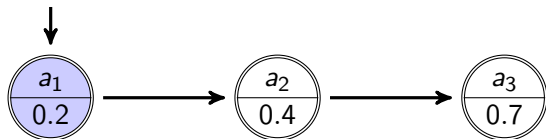
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_3 = 0.7$

current
sample



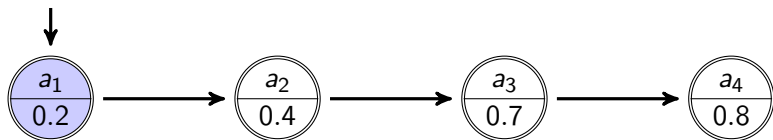
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_4 = 0.8$

current
sample



- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_5 = 0.6$

current
sample



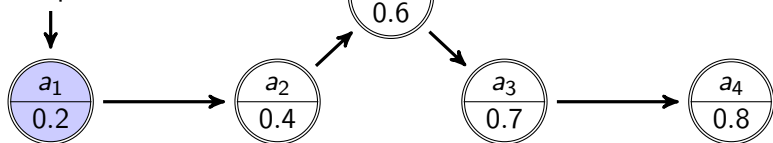
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_5 = 0.6$

current
sample



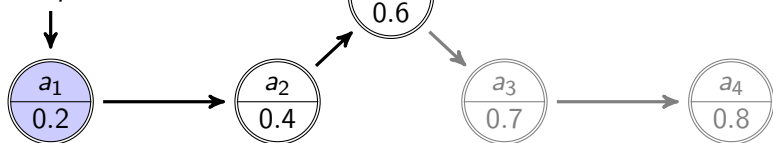
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_5 = 0.6$

current
sample



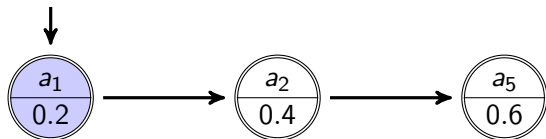
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

↓
 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_5 = 0.6$

current
sample



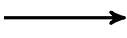
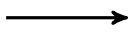
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_6 = 0.5$

current
sample

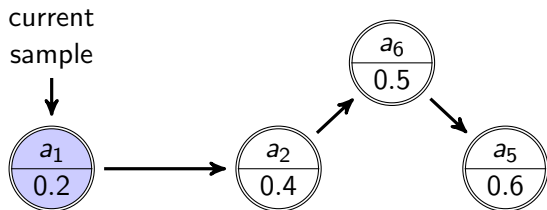


- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_6 = 0.5$



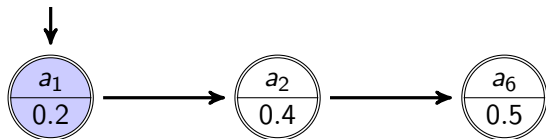
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_6 = 0.5$

current
sample



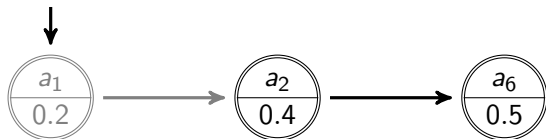
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_6 = 0.5$

current
sample



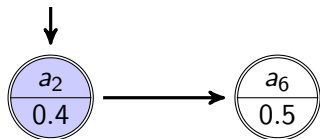
- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

random value $r_6 = 0.5$

current
sample

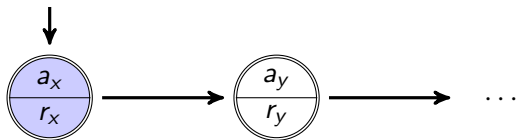


- Algorithm:**
- 1: For each a_i pick random value $r_i \in (0, 1)$
 - 2: In window (a_{i-w+1}, \dots, a_i) choose a_j with smallest r_j
 - 3: Only maintain items in window whose r -value is minimal among subsequent r -values

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

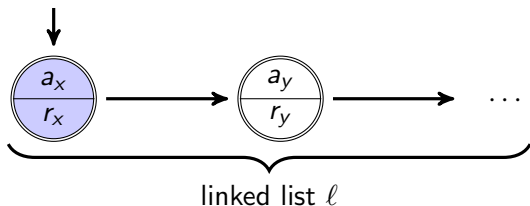
current
sample



Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample

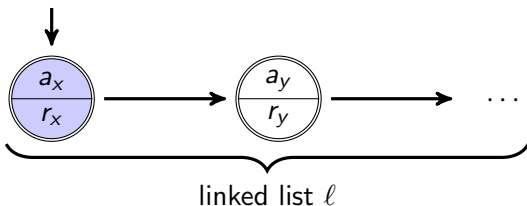


Analysis: What is the length of l ?

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



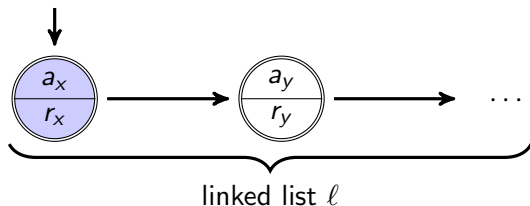
Analysis: What is the length of l ?

Worst case: $|\ell| = w$ But that is very unlikely.

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



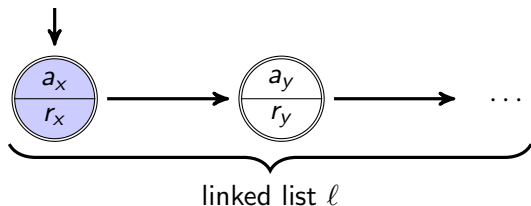
Analysis: What is the length of l ?

$$\text{Exp}[|l|] =$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



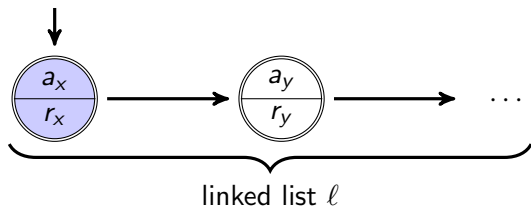
Analysis: What is the length of l ?

$$\text{Exp}[|l|] = \text{Pr}[a_i \in l] +$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



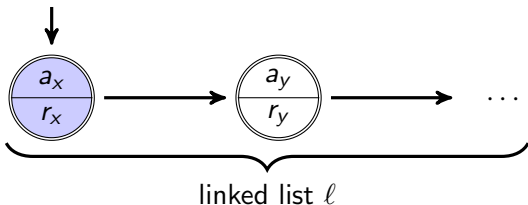
Analysis: What is the length of l ?

$$\text{Exp}[|l|] = \text{Pr}[a_i \in l] + \text{Pr}[a_{i-1} \in l] +$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



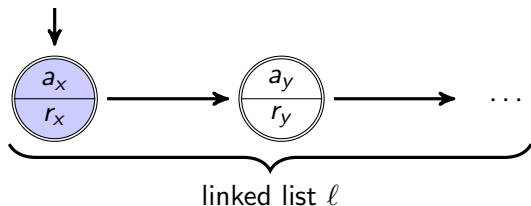
Analysis: What is the length of ℓ ?

$$\text{Exp}[|\ell|] = \text{Pr}[a_i \in \ell] + \text{Pr}[a_{i-1} \in \ell] + \dots + \text{Pr}[a_{i-w+1} \in \ell]$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



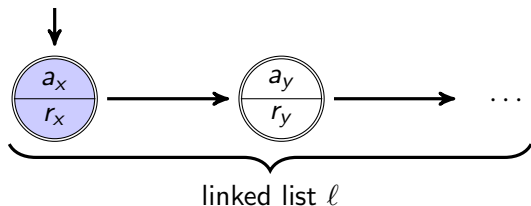
Analysis: What is the length of ℓ ?

$$\begin{aligned} \text{Exp}[|\ell|] &= \text{Pr}[a_i \in \ell] + \text{Pr}[a_{i-1} \in \ell] + \dots + \text{Pr}[a_{i-w+1} \in \ell] \\ &= 1 + \end{aligned}$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



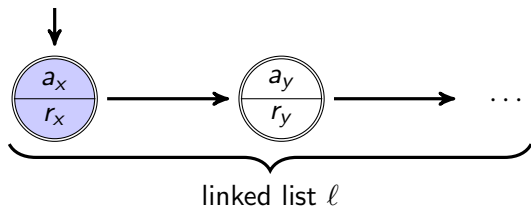
Analysis: What is the length of ℓ ?

$$\begin{aligned} \text{Exp}[|\ell|] &= \text{Pr}[a_i \in \ell] + \text{Pr}[a_{i-1} \in \ell] + \dots + \text{Pr}[a_{i-w+1} \in \ell] \\ &= 1 + \frac{1}{2} + \end{aligned}$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



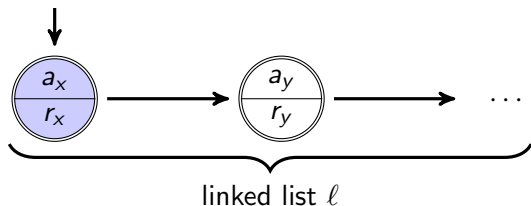
Analysis: What is the length of l ?

$$\begin{aligned} \text{Exp}[|l|] &= \text{Pr}[a_i \in l] + \text{Pr}[a_{i-1} \in l] + \dots + \text{Pr}[a_{i-w+1} \in l] \\ &= 1 + \frac{1}{2} + \frac{1}{3} + \end{aligned}$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



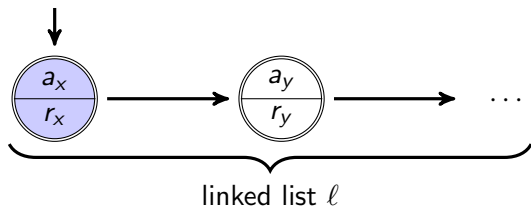
Analysis: What is the length of l ?

$$\begin{aligned} \text{Exp}[|l|] &= \text{Pr}[a_i \in l] + \text{Pr}[a_{i-1} \in l] + \dots + \text{Pr}[a_{i-w+1} \in l] \\ &= 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots + \frac{1}{w} \end{aligned}$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



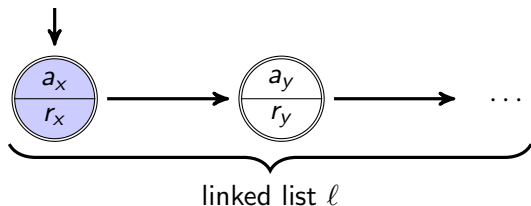
Analysis: What is the length of ℓ ?

$$\begin{aligned} \text{Exp}[|\ell|] &= \text{Pr}[a_i \in \ell] + \text{Pr}[a_{i-1} \in \ell] + \dots + \text{Pr}[a_{i-w+1} \in \ell] \\ &= 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots + \frac{1}{w} = \mathcal{O}(\log w) \end{aligned}$$

Sliding Window Sampling

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_i, \dots, a_m$

current
sample



Analysis:

$$\text{Exp}[\text{memory usage}] = \mathcal{O}(\log w \cdot \log n)$$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

Point query: $f_i = ?$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

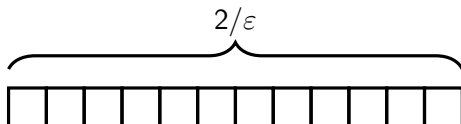
Point query: $f_i = ?$

Trivial solution: maintain n counters

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

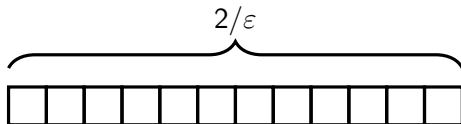


Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

rand. hash funct. h_1

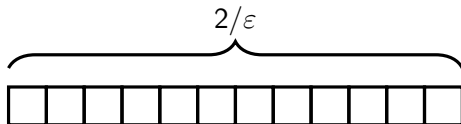


Count-Min Sketch

[Cormode, Muthukrishnan '04]

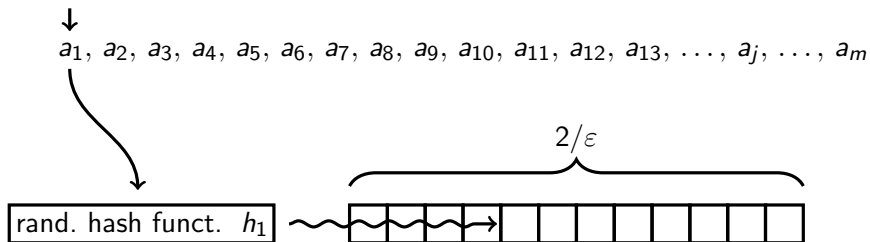
↓
 $a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

rand. hash funct. h_1



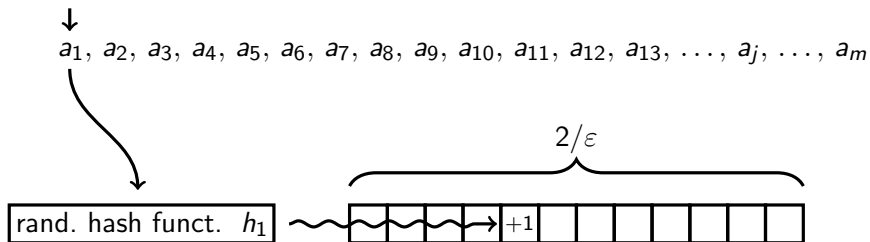
Count-Min Sketch

[Cormode, Muthukrishnan '04]



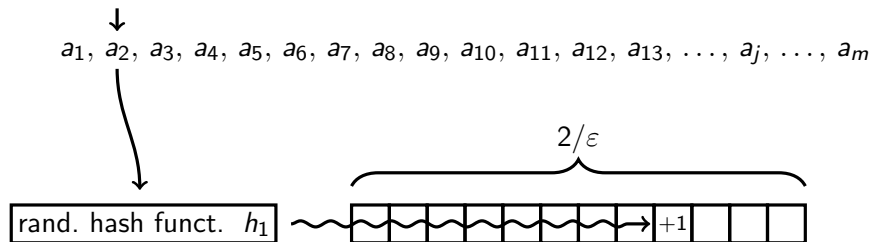
Count-Min Sketch

[Cormode, Muthukrishnan '04]



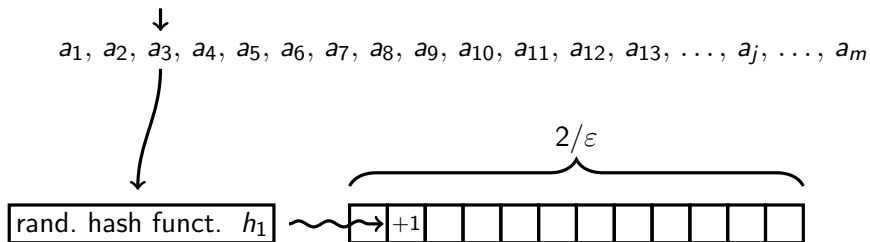
Count-Min Sketch

[Cormode, Muthukrishnan '04]



Count-Min Sketch

[Cormode, Muthukrishnan '04]

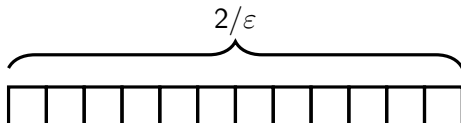


Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

rand. hash funct. h_1

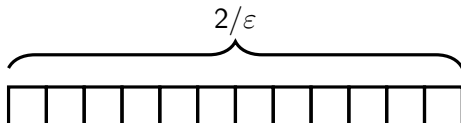


Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

rand. hash funct. h_1

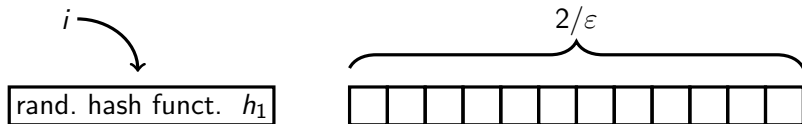


Point query: $f_i = ?$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

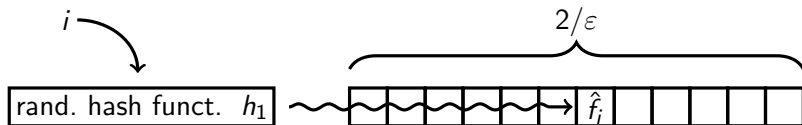


Point query: $f_i = ?$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

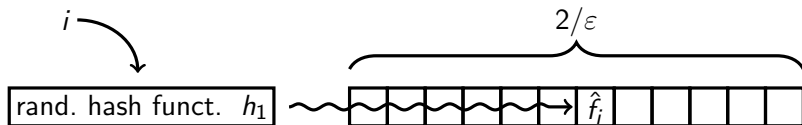


Point query: $f_i = ?$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

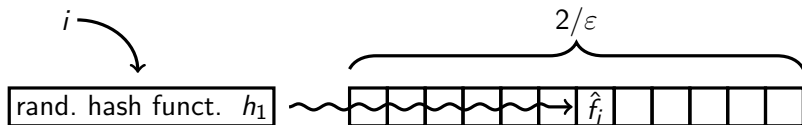
$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



Point query: $f_i = ?$

Give \hat{f}_i as estimate for f_i

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



Point query: $f_i = ?$

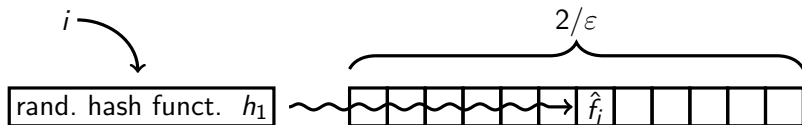
Give \hat{f}_i as estimate for f_i

Analysis:

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

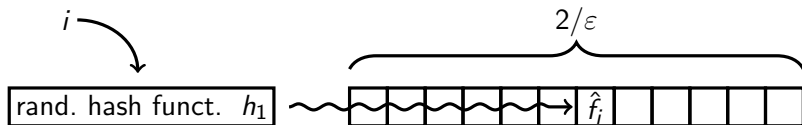


Point query: $f_i = ?$

Give \hat{f}_i as estimate for f_i

Analysis: $\hat{f}_i \geq f_i$

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



Point query: $f_i = ?$

Give \hat{f}_i as estimate for f_i

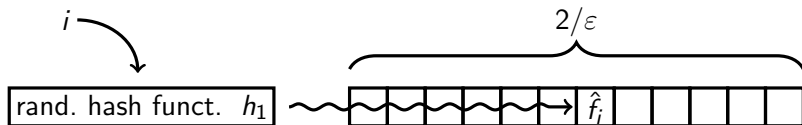
Analysis: $\hat{f}_i \geq f_i$

$$\text{Exp}[\hat{f}_i] \leq f_i + \epsilon \cdot m/2$$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



Point query: $f_i = ?$

Give \hat{f}_i as estimate for f_i

Analysis: $\hat{f}_i \geq f_i$

$$\text{Exp}[\hat{f}_i] \leq f_i + \epsilon \cdot m/2$$

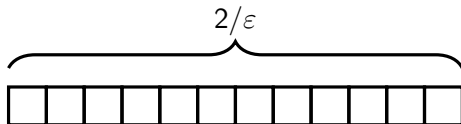
$$\text{Markov-Ineq.: } \Pr[\hat{f}_i > f_i + \epsilon \cdot m] \leq \frac{1}{2}$$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

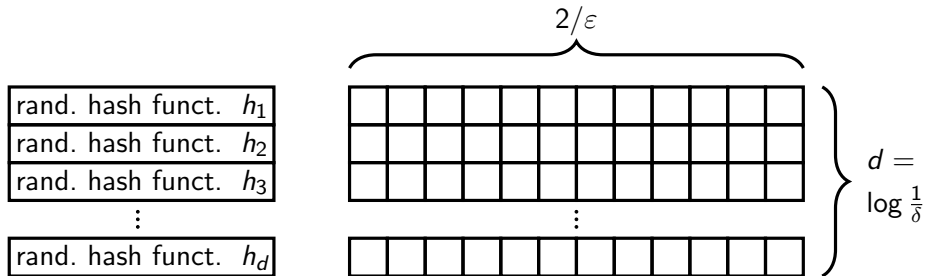
rand. hash funct. h_1



Count-Min Sketch

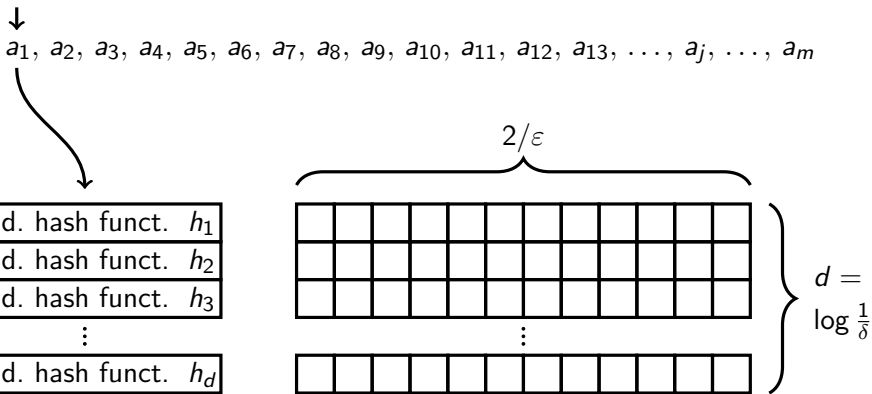
[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



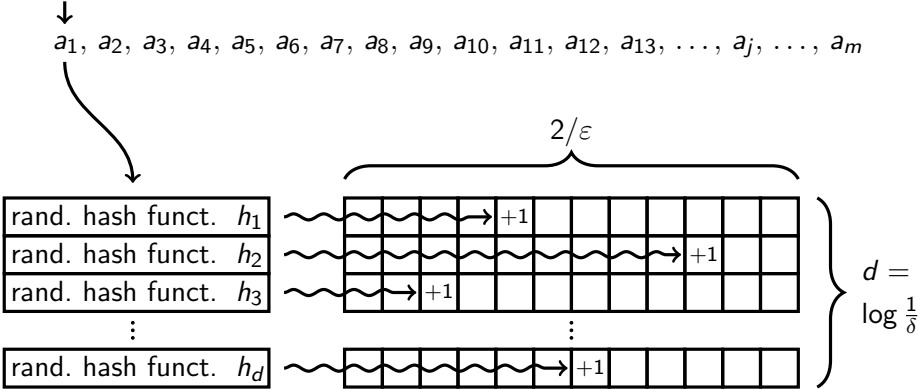
Count-Min Sketch

[Cormode, Muthukrishnan '04]



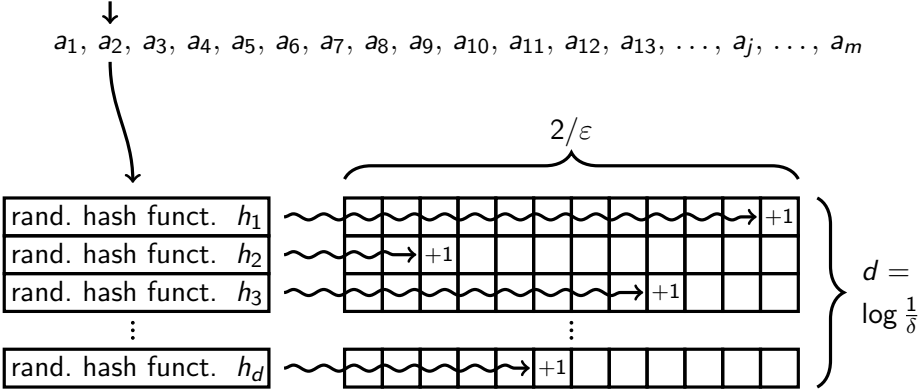
Count-Min Sketch

[Cormode, Muthukrishnan '04]



Count-Min Sketch

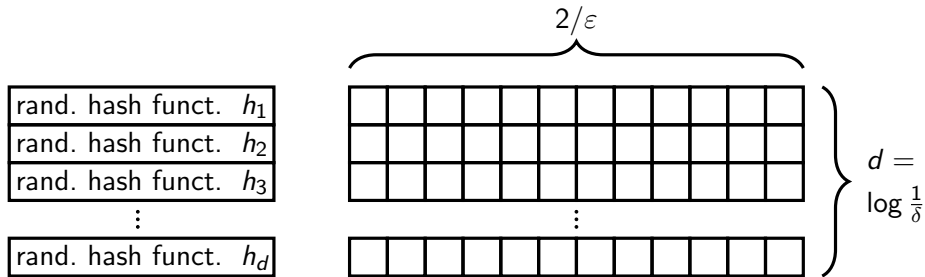
[Cormode, Muthukrishnan '04]



Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

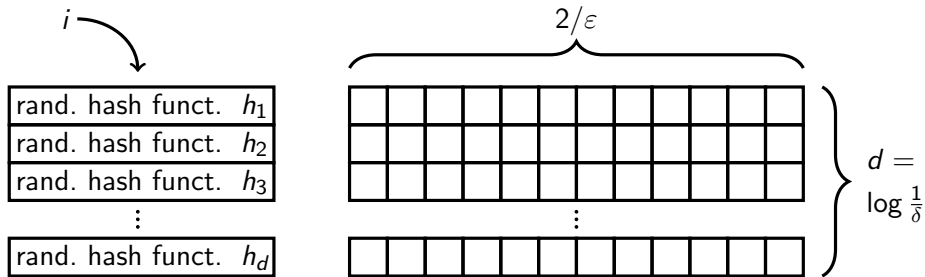


Point query: $f_i = ?$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

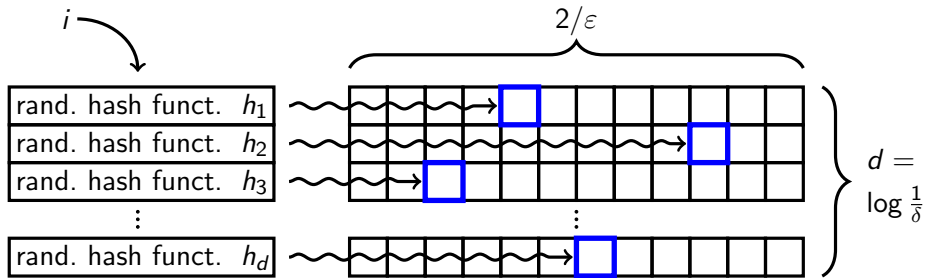


Point query: $f_i = ?$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

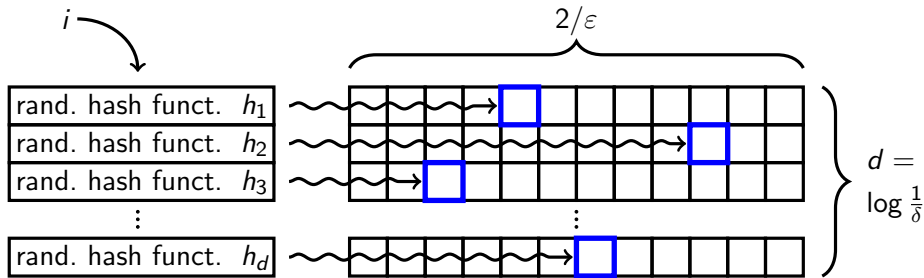


Point query: $f_i = ?$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



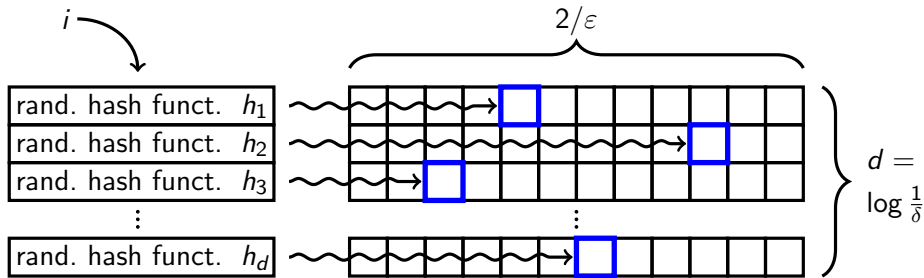
Point query: $f_i = ?$

Give $\hat{f}_i :=$ minimum of \square -values as estimate for f_i

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$

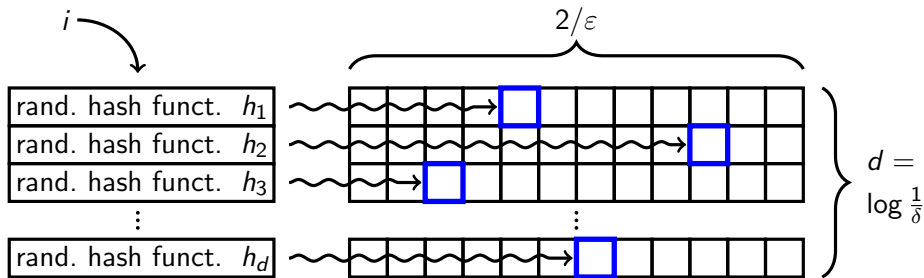


Analysis: $Pr[\hat{f}_i > f_i + \epsilon \cdot m] \leq \left(\frac{1}{2}\right)^d = \delta$

Count-Min Sketch

[Cormode, Muthukrishnan '04]

$a_1, a_2, a_3, a_4, a_5, a_6, a_7, a_8, a_9, a_{10}, a_{11}, a_{12}, a_{13}, \dots, a_j, \dots, a_m$



Analysis: $Pr[\hat{f}_i > f_i + \epsilon \cdot m] \leq \left(\frac{1}{2}\right)^d = \delta$

Memory consumption: $\mathcal{O}\left(\frac{1}{\epsilon} \cdot \log \frac{1}{\delta} \cdot \log m + \log \frac{1}{\delta} \cdot \log n\right)$

Recap

- ▶ Reservoir sampling for uniform selection
- ▶ AMS sampling for frequency moments
- ▶ Sliding window sampling
- ▶ Count-Min sketch for point queries

References

- ▶ N. Alon, Y. Matias, and M. Szegedy: The space complexity of approximating the frequency moments. *J. Comput. Syst. Sci.* 58, pp. 137-147, 1999.
- ▶ G. Cormode and S. Muthukrishnan: An Improved Data Stream Summary: The Count-Min Sketch and its Applications. *Journal of Algorithms* 55(1), pp. 58-75, 2005.
- ▶ B. Lahiri and S. Tirthapura: Stream Sampling. Pages 2838-2842 in: Ling Liu, M. Tamer zsu (Eds.): *Encyclopedia of Database Systems*. Springer US, 2009.
- ▶ S. Muthukrishnan: *Data Streams: Algorithms and Applications*. *Foundations and Trends in Theoretical Computer Science*, 1(2), 2005.
- ▶ J. Vitter: Random Sampling with a Reservoir. *ACM Transactions on Mathematical Software*, 11(1), pp. 37-57, 1985.