

ALGORITHMS FOR DATA SCIENCE: LECTURE 4

VASILEIOS NAKOS

NATIONAL TECHNICAL UNIVERSITY OF ATHENS

APRIL 2, 2021

Algorithms for processing a stream of elements and maintaining statistics in *sublinear* space about the elements in the stream.

Algorithms for processing a stream of elements and maintaining statistics in *sublinear* space about the elements in the stream.

Alternatively, low-space data structures for big data computation.

Algorithms for processing a stream of elements and maintaining statistics in *sublinear* space about the elements in the stream.

Alternatively, low-space data structures for big data computation.

★ Compress data "on-the-fly": store a small piece of information sufficient to answer approximate answer for the data set.

- Traffic Monitoring in networks.

- Traffic Monitoring in networks.
- Very large Databases, i.e. estimating the size of a join.

APPLICATIONS

- Traffic Monitoring in networks.
- Very large Databases, i.e. estimating the size of a join.
- Manipulation of astronomical (satellite imagery), financial data.

- Traffic Monitoring in networks.
- Very large Databases, i.e. estimating the size of a join.
- Manipulation of astronomical (satellite imagery), financial data.
- GPS or seismometer readings to detect geological anomalies, sensor networks etc

APPLICATIONS

- Traffic Monitoring in networks.
- Very large Databases, i.e. estimating the size of a join.
- Manipulation of astronomical (satellite imagery), financial data.
- GPS or seismometer readings to detect geological anomalies, sensor networks etc
- Training Machine Learning models when the training data set is huge.



DISTINCT ELEMENTS

Input: A sequence of elements $x_1, x_2, \dots, \subseteq [n]$, where you can think of n as 2^{64} .

Output: Number of *distinct* elements in x

DISTINCT ELEMENTS

Input: A sequence of elements $x_1, x_2, \dots, \subseteq [n]$, where you can think of n as 2^{64} .

Output: Number of *distinct* elements in x

- Distinct users hitting a webpage.
- Distinct values in a column of a database
- Number of distinct queries to a search engine.
- Distinct patterns in DNA sequence.

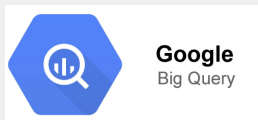
DISTINCT ELEMENTS

Input: A sequence of elements $x_1, x_2, \dots, \subseteq [n]$, where you can think of n as 2^{64} .

Output: Number of *distinct* elements in x

- Distinct users hitting a webpage.
- Distinct values in a column of a database
- Number of distinct queries to a search engine.
- Distinct patterns in DNA sequence.

Implementations used by Google (Sawzall, Dremel, PowerDrill), Yahoo, Twitter, Facebook Presto, etc.



Input: Vector $x \in \mathbb{R}^n$, updates (i, Δ) causing $x_i \leftarrow x_i + \Delta$.

Input: Vector $x \in \mathbb{R}^n$, updates (i, Δ) causing $x_i \leftarrow x_i + \Delta$.

Output: An approximate answer to $\sum_{i=1}^n x_i^2$.

- Anomaly detection in traffic monitoring.
- Detect DoS attacks.
- Database optimization engine to estimate *self join size*.
- Subroutine in many other streaming algorithms.

ARCHITECTURE OF A STREAMING ALGORITHM

Implements two routines, **Update()** and **Query()** using space S usually *sublinear* in the size of the input, i.e. does **not** store the whole input.

ARCHITECTURE OF A STREAMING ALGORITHM

Implements two routines, **Update()** and **Query()** using space S usually *sublinear* in the size of the input, i.e. does **not** store the whole input.

- Distinct elements: $O(\frac{\log \log n}{\epsilon^2} + \log n)$ bits of space to estimate number of distinct elements up to $(1 + \epsilon)$
- F_2 estimation: $O(\frac{\log n}{\epsilon^2})$ bits of space to estimate the F_2 moment up to $1 + \epsilon$, i.e. find V such that

$$(1 - \epsilon) \sum_{i=1}^n x_i^2 \leq V \leq (1 + \epsilon) \sum_{i=1}^n x_i^2.$$

Streaming algorithms are almost always *randomized* and *approximate*.

Let X be a discrete random variable that takes values on $\{\dots, -1, 0, 1, \dots\}$. Then

- (expectation) $\mathbb{E}(X) := \sum_{i=-\infty}^{\infty} i \cdot Pr[X = i]$
- (variance) $\text{Var}(X) := \mathbb{E}(X^2) - \mathbb{E}(X)^2$

Let X be a discrete random variable that takes values on $\{\dots, -1, 0, 1, \dots\}$. Then

- (expectation) $\mathbb{E}(X) := \sum_{i=-\infty}^{\infty} i \cdot Pr[X = i]$
- (variance) $\text{Var}(X) := \mathbb{E}(X^2) - \mathbb{E}(X)^2$
- (linearity of expectation) $\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y)$.

Let X be a discrete random variable that takes values on $\{\dots, -1, 0, 1, \dots\}$. Then

- (expectation) $\mathbb{E}(X) := \sum_{i=-\infty}^{\infty} i \cdot \Pr[X = i]$
- (variance) $\text{Var}(X) := \mathbb{E}(X^2) - \mathbb{E}(X)^2$
- (linearity of expectation) $\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y)$.
- (Markov's inequality) For a variable X that takes only positive values we have $\Pr[X \geq x] \leq \frac{\mathbb{E}(X)}{x}$.

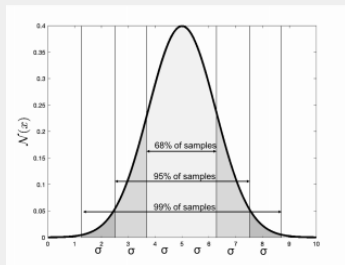
Let X be a discrete random variable that takes values on $\{\dots, -1, 0, 1, \dots\}$. Then

- (expectation) $\mathbb{E}(X) := \sum_{i=-\infty}^{\infty} i \cdot \Pr[X = i]$
- (variance) $\text{Var}(X) := \mathbb{E}(X^2) - \mathbb{E}(X)^2$
- (linearity of expectation) $\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y)$.
- (Markov's inequality) For a variable X that takes only positive values we have $\Pr[X \geq x] \leq \frac{\mathbb{E}(X)}{x}$.
- (Chebyshev's inequality) $\Pr[|X - \mathbb{E}(X)| \geq \lambda] \leq \frac{\text{Var}(X)}{\lambda^2}$.

Chebyshev's inequality is very useful in the design of randomized algorithms, showing that an estimator concentrates around its expected value.

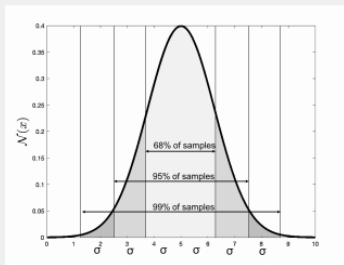
Question: How tight is Chebyshev's inequality?

Question: How tight is Chebyshev's inequality?



Follow the "68-95-99 rule" for Gaussian bell-curve $\mathcal{N}(0, \sigma^2)$.

Question: How tight is Chebyshev's inequality?



Follow the "68-95-99 rule" for Gaussian bell-curve $\mathcal{N}(0, \sigma^2)$.

- Chebyshev's inequality versus true value.
- $Pr[|X - \mathbb{E}(X)| \geq 1\sigma] \leq 100\%$ vs $Pr[|X - \mathbb{E}(X)| \geq 1\sigma] \approx 32\%$
- $Pr[|X - \mathbb{E}(X)| \geq 2\sigma] \leq 25\%$ vs $Pr[|X - \mathbb{E}(X)| \geq 1\sigma] \approx 5\%$
- $Pr[|X - \mathbb{E}(X)| \geq 3\sigma] \leq 11\%$ vs $Pr[|X - \mathbb{E}(X)| \geq 1\sigma] \approx 1\%$
- $Pr[|X - \mathbb{E}(X)| \geq 4\sigma] \leq 6\%$ vs $Pr[|X - \mathbb{E}(X)| \geq 1\sigma] \approx 0.01\%$

BACK TO DISTINCT ELEMENTS

Input: A sequence of elements $x_1, x_2, \dots, \subseteq [n]$, where you can think of n as 2^{64} .

BACK TO DISTINCT ELEMENTS

Input: A sequence of elements $x_1, x_2, \dots, \subseteq [n]$, where you can think of n as 2^{64} .

Output: Number of *distinct* elements D in x . Example:
 $0, 1, 0, 0, 0, 1, 2 \rightarrow 3$

BACK TO DISTINCT ELEMENTS

Input: A sequence of elements $x_1, x_2, \dots, \subseteq [n]$, where you can think of n as 2^{64} .

Output: Number of *distinct* elements D in x . Example:
 $0, 1, 0, 0, 0, 1, 2 \rightarrow 3$

The (idealized) Flajolet-Martin algorithm:

- Choose a random hash function $h : [n] \rightarrow [0, 1]$
- $S = \infty$
- For every element e set $S \leftarrow \min \{S, h(e)\}$.
- Output $\frac{1}{S} - 1$.

BACK TO DISTINCT ELEMENTS

Input: A sequence of elements $x_1, x_2, \dots, \subseteq [n]$, where you can think of n as 2^{64} .

Output: Number of *distinct* elements D in x . Example:
 $0, 1, 0, 0, 0, 1, 2 \rightarrow 3$

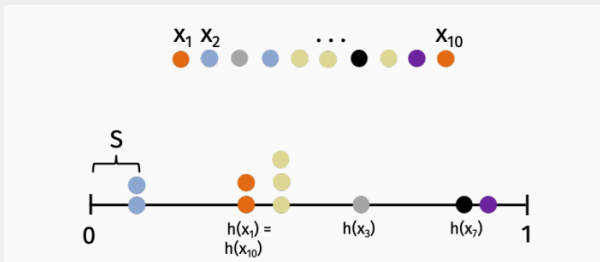
The (idealized) Flajolet-Martin algorithm:

- Choose a random hash function $h : [n] \rightarrow [0, 1]$
- $S = \infty$
- For every element e set $S \leftarrow \min \{S, h(e)\}$.
- Output $\frac{1}{S} - 1$.

We must store S and a description of $h \dots$

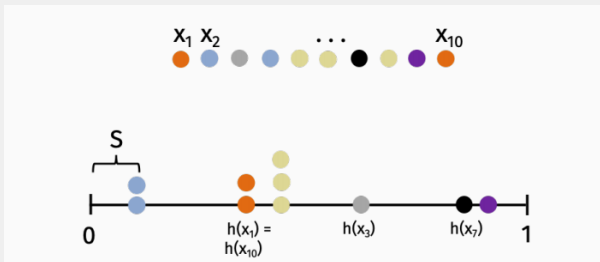
SO WHAT IS V?

S is the minimum hash value ever seen so far.



SO WHAT IS V?

S is the minimum hash value ever seen so far.



If an element comes many times, it will be always mapped to the same position. We return $\frac{1}{S} - 1$.

Lemma

$$\mathbb{E}(S) = \frac{1}{D+1}.$$

Lemma

$$\mathbb{E}(S) = \frac{1}{D+1}.$$

$$\begin{aligned}\mathbb{E}(S) &= \int_0^1 \Pr[S \geq \lambda] d\lambda = \\ &= \int_0^1 (1 - \lambda)^D d\lambda = \frac{1}{D+1}\end{aligned}$$

We know that $\mathbb{E}(S) = \frac{1}{D+1}$. **Estimate** $\tilde{D} = \frac{1}{S} - 1$.

We know that $\mathbb{E}(S) = \frac{1}{D+1}$. **Estimate** $\tilde{D} = \frac{1}{S} - 1$.
Note that it **does not** hold that $\mathbb{E}(\tilde{D}) = D$.

We know that $\mathbb{E}(S) = \frac{1}{D+1}$. **Estimate** $\tilde{D} = \frac{1}{S} - 1$.

Note that it **does not** hold that $\mathbb{E}(\tilde{D}) = D$.

However,

$$|S - \mathbb{E}(S)| \leq \epsilon \cdot \mathbb{E}(S) \Rightarrow (1 - 4\epsilon)D \leq \tilde{D} \leq (1 + 4\epsilon)D.$$

We know that $\mathbb{E}(S) = \frac{1}{D+1}$. **Estimate** $\tilde{D} = \frac{1}{S} - 1$.

Note that it **does not** hold that $\mathbb{E}(\tilde{D}) = D$.

However,

$$|S - \mathbb{E}(S)| \leq \epsilon \cdot \mathbb{E}(S) \Rightarrow (1 - 4\epsilon)D \leq \tilde{D} \leq (1 + 4\epsilon)D.$$

Thus, precise estimation of D reduces to showing that S concentrates around its value \rightarrow Chebyshev's inequality!

LET'S COMPUTE THE VARIANCE.

Lemma

$$\text{Var}(S) = \mathbb{E}(S^2) - \mathbb{E}(S)^2 = \frac{2}{(D+1)(D+2)} - \frac{1}{(D+1)^2} \leq \frac{1}{(D+1)^2}.$$

Similarly as before,

$$\begin{aligned}\mathbb{E}(S^2) &= \int_0^1 \Pr[S^2 \geq \lambda] d\lambda = \\ &= \int_0^1 \Pr[S \geq \sqrt{\lambda}] d\lambda = \\ &= \int_0^1 (1 - \sqrt{\lambda})^D d\lambda = \frac{2}{(D+1)(D+2)}\end{aligned}$$

LET'S SEE WHAT WE GET

- $\mathbb{E}(S) = \frac{1}{D+1}$
- $\text{Var}(S) \leq \frac{1}{(D+1)^2}$

LET'S SEE WHAT WE GET

- $\mathbb{E}(S) = \frac{1}{D+1}$
- $\text{Var}(S) \leq \frac{1}{(D+1)^2}$
- $\Pr[|S - \mathbb{E}(S)| \geq \epsilon \sqrt{\text{Var}(S)}] \leq \frac{1}{\epsilon^2}$ (too large!)

THE BOOSTING OF INDEPENDENT REPETITIONS

Consider identically distributed random variables S_1, S_2, \dots, S_r , and let

$$\bar{S} := \frac{1}{r}(S_1 + S_2 + \dots + S_r).$$

■ $\mathbb{E}(S) = \mathbb{E}\left(\frac{1}{r}(S_1 + S_2 + \dots + S_r)\right) =$
 $\frac{1}{r}(\mathbb{E}(S)_1 + \mathbb{E}(S)_2 + \dots + \mathbb{E}(S)_r) = \mathbb{E}(S).$

THE BOOSTING OF INDEPENDENT REPETITIONS

Consider identically distributed random variables S_1, S_2, \dots, S_r , and let

$$\bar{S} := \frac{1}{r}(S_1 + S_2 + \dots + S_r).$$

- $\mathbb{E}(\bar{S}) = \mathbb{E}\left(\frac{1}{r}(S_1 + S_2 + \dots + S_r)\right) = \frac{1}{r}(\mathbb{E}(S)_1 + \mathbb{E}(S)_2 + \dots + \mathbb{E}(S)_r) = \mathbb{E}(S).$
- $\text{Var}(\bar{S}) = \frac{1}{r} \cdot \text{Var}(S)$, since
 1. For random variables X, Y we have $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$, and

THE BOOSTING OF INDEPENDENT REPETITIONS

Consider identically distributed random variables S_1, S_2, \dots, S_r , and let

$$\bar{S} := \frac{1}{r}(S_1 + S_2 + \dots + S_r).$$

- $\mathbb{E}(\bar{S}) = \mathbb{E}\left(\frac{1}{r}(S_1 + S_2 + \dots + S_r)\right) = \frac{1}{r}(\mathbb{E}(S)_1 + \mathbb{E}(S)_2 + \dots + \mathbb{E}(S)_r) = \mathbb{E}(S)$.
- $\text{Var}(\bar{S}) = \frac{1}{r} \cdot \text{Var}(S)$, since
 1. For random variables X, Y we have $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$, and
 2. $\text{Var}(aX) = a^2 \text{Var}(X)$.

THE BOOSTING OF INDEPENDENT REPETITIONS

Consider identically distributed random variables S_1, S_2, \dots, S_r , and let

$$\bar{S} := \frac{1}{r}(S_1 + S_2 + \dots + S_r).$$

- $\mathbb{E}(\bar{S}) = \mathbb{E}(\frac{1}{r}(S_1 + S_2 + \dots + S_r)) = \frac{1}{r}(\mathbb{E}(S)_1 + \mathbb{E}(S)_2 + \dots + \mathbb{E}(S)_r) = \mathbb{E}(S)$.
- $\text{Var}(\bar{S}) = \frac{1}{r} \cdot \text{Var}(S)$, since
 1. For random variables X, Y we have $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$, and
 2. $\text{Var}(aX) = a^2 \text{Var}(X)$.

Thus, in our case take $r = O(\frac{1}{\epsilon^2})$ different instantiations of the algorithm, and output the average!

The (idealized) Flajolet-Martin algorithm:

- $r := \frac{3}{\epsilon^2}$
- Choose random hash function $h_r : [n] \rightarrow [0, 1]$ for all $j \in [r]$.
- For every $j \in [r]$ set $S_j = \infty$
- For every element e set and every $j \in [r]$ set $S_j \leftarrow \min \{S_j, h_j(e)\}$.
- $S := \frac{1}{r}(S_1 + S_2 + \dots + S_r)$.
- Output $\frac{1}{S} - 1$.

The (idealized) Flajolet-Martin algorithm:

- $r := \frac{3}{\epsilon^2}$
- Choose random hash function $h_r : [n] \rightarrow [0, 1]$ for all $j \in [r]$.
- For every $j \in [r]$ set $S_j = \infty$
- For every element e set and every $j \in [r]$ set $S_j \leftarrow \min \{S_j, h_j(e)\}$.
- $S := \frac{1}{r}(S_1 + S_2 + \dots + S_r)$.
- Output $\frac{1}{S} - 1$.

WHAT HAPPENS

We have $\mathbb{E}(S) = \frac{1}{D+1}$, $\text{Var}(S) < \frac{1}{(D+1)^2 r}$, so applying Chebyshev's inequality yields

$$\Pr[|S - \mathbb{E}(S)| \geq \frac{\epsilon}{D+1}] \leq \frac{1}{3}.$$

Thus, with probability $\frac{2}{3}$ our estimator will satisfy

$$(1 - 4\epsilon)D \leq \tilde{D} \leq (1 + 4\epsilon)D.$$

So we need to run the algorithm with $\epsilon' := \frac{\epsilon}{4}$.

HOW TO BOOST THE SUCCESS PROBABILITY?

Idea I: For target probability δ , we can set $r = \frac{16}{\epsilon^2 \delta}$ and obtain an estimate satisfying

$$(1 - \epsilon)D \leq \tilde{D} \leq (1 + \epsilon)D.$$

with probability $1 - \delta$.

HOW TO BOOST THE SUCCESS PROBABILITY?

Idea I: For target probability δ , we can set $r = \frac{16}{\epsilon^2 \delta}$ and obtain an estimate satisfying

$$(1 - \epsilon)D \leq \tilde{D} \leq (1 + \epsilon)D.$$

with probability $1 - \delta$.

Why? $\text{Var}(S)$ becomes $\frac{\epsilon^2 \delta}{16 \cdot (D+1)^2}$, so the same analysis yields

$$\Pr[|S - \mathbb{E}(S)| \geq \frac{\epsilon}{4(D+1)}] \leq \delta.$$

But we can do waaay better!

TRICK OF THE MEDIANS

For each $t \in [2 \log(1/\delta)]$, keep a distinct elements data structures D_t with $O(\frac{1}{\epsilon^2})$ counters, and let $S^{(t)}$ be the estimate produced by each data structure.

TRICK OF THE MEDIANS

For each $t \in [2 \log(1/\delta)]$, keep a distinct elements data structures D_t with $O(\frac{1}{\epsilon^2})$ counters, and let $S^{(t)}$ be the estimate produced by each data structure.

Then the quantity

$$S := \text{median}_t S^{(t)}$$

satisfies the desired inequality with probability $1 - \delta$. Good, huh?

WHAT HAPPENED?

We began with an *unbiased* estimator which has large variance, and

- for a multiplicative cost of $O(\frac{1}{\epsilon^2})$ we reduced the variance to the desired one (trick of the means).

WHAT HAPPENED?

We began with an *unbiased* estimator which has large variance, and

- for a multiplicative cost of $O(\frac{1}{\epsilon^2})$ we reduced the variance to the desired one (trick of the means).
- for a multiplicative cost of $O(\log(1/\delta))$ we reduced the failure probability to δ (trick of the medians).

WHAT HAPPENED?

We began with an *unbiased* estimator which has large variance, and

- for a multiplicative cost of $O(\frac{1}{\epsilon^2})$ we reduced the variance to the desired one (trick of the means).
- for a multiplicative cost of $O(\log(1/\delta))$ we reduced the failure probability to δ (trick of the medians).

Space complexity (besides storing the hash functions):

$O(\log(1/\delta) \cdot \frac{1}{\epsilon^2} \cdot \log n)$ bits or $O(\log(1/\delta)/\epsilon^2)$ words.

WHAT HAPPENED?

We began with an *unbiased* estimator which has large variance, and

- for a multiplicative cost of $O(\frac{1}{\epsilon^2})$ we reduced the variance to the desired one (trick of the means).
- for a multiplicative cost of $O(\log(1/\delta))$ we reduced the failure probability to δ (trick of the medians).

Space complexity (besides storing the hash functions):

$O(\log(1/\delta) \cdot \frac{1}{\epsilon^2} \cdot \log n)$ bits or $O(\log(1/\delta)/\epsilon^2)$ words. Compare with the trivial space of $O(D)$ words.

AVOIDING TAKING $h : [n] \rightarrow [0, 1]$

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
	⋮
$h(x_n)$	1011000

In the practical version of Flajolet-Martin (HyperLogLog) we estimate distinct elements based on maximum number of trailing zeros.

AVOIDING TAKING $h : [n] \rightarrow [0, 1]$

$h(x_1)$	1010010
$h(x_2)$	1001100
$h(x_3)$	1001110
	⋮
$h(x_n)$	1011000

In the practical version of Flajolet-Martin (HyperLogLog) we estimate distinct elements based on maximum number of trailing zeros.

$$\Pr[h(x) \text{ has } \log D \text{ trailing zeros}] = \frac{1}{D}.$$

WHAT HAPPENS FOR THE PRACTICAL FLAJOLET-MARTIN?

Total space: $O(\log \log D/\epsilon^2 + \log D)$ for an ϵ approximation with constant probability.

Quote from "Loglog Counting of Large Cardinalities"

"Using an auxiliary memory smaller than the size of this abstract, the LogLog algorithm makes it possible to estimate in a single pass and within a few percents the number of different words in the whole of Shakespeare's works." - Flajolet, Durand.

WHAT HAPPENS FOR THE PRACTICAL FLAJOLET-MARTIN?

Total space: $O(\log \log D/\epsilon^2 + \log D)$ for an ϵ approximation with constant probability.

Quote from "Loglog Counting of Large Cardinalities"

"Using an auxiliary memory smaller than the size of this abstract, the LogLog algorithm makes it possible to estimate in a single pass and within a few percents the number of different words in the whole of Shakespeare's works." - Flajolet, Durand.

Using HyperLogLog to approximate 1 billion distinct items to 2% accuracy can be in approximately $1.6\text{KB} = 12800$ bits!

DISTINCT ELEMENTS IN THE DISTRIBUTED SETTING

The Flajolet-Martin algorithm is totally distributed: why share lists of distinct elements when you can only share a bunch of minimum hash values seen?

DISTINCT ELEMENTS IN THE DISTRIBUTED SETTING

The Flajolet-Martin algorithm is totally distributed: why share lists of distinct elements when you can only share a bunch of minimum hash values seen?



Estimate spam rate: Count number of distinct subject lines in emails sent by users that have registered in the last week, in comparison to number of emails sent overall.

DISTINCT ELEMENTS IN THE DISTRIBUTED SETTING

The Flajolet-Martin algorithm is totally distributed: why share lists of distinct elements when you can only share a bunch of minimum hash values seen?



Estimate spam rate: Count number of distinct subject lines in emails sent by users that have registered in the last week, in comparison to number of emails sent overall.

Good news: Answering the above query can be done in 2 seconds in Google's distributed implementations!

Input: A vector $x \in \mathbb{R}^n$ and updates (i, Δ) causing $x_i \leftarrow x_i + \Delta$.

Input: A vector $x \in \mathbb{R}^n$ and updates (i, Δ) causing $x_i \leftarrow x_i + \Delta$.

Output: A value V satisfying $(1 - \epsilon) \cdot \sum_{i=1}^n x_i^2 \leq V \leq (1 + \epsilon) \cdot \sum_{i=1}^n x_i^2$.

Alon-Mattias-Szegedy (AMS) sketch

There exists an algorithm which uses $O(\frac{\log n}{\epsilon^2})$ bits of space and returns an estimator as above with constant probability.

The algorithm and proof is just a couple of lines, yet the authors received the Goedel prize for that!

- $V \leftarrow \mathbf{0}$
- Pick hash function $\sigma : [n] \rightarrow \{-1, 1\}$ (random signs)

- $V \leftarrow 0$
- Pick hash function $\sigma : [n] \rightarrow \{-1, 1\}$ (random signs)
- Upon update (i, Δ) set $V \leftarrow V + \sigma(i) \cdot \Delta$

- $V \leftarrow 0$
- Pick hash function $\sigma : [n] \rightarrow \{-1, 1\}$ (random signs)
- Upon update (i, Δ) set $V \leftarrow V + \sigma(i) \cdot \Delta$
- Output V^2

- $V \leftarrow 0$
- Pick hash function $\sigma : [n] \rightarrow \{-1, 1\}$ (random signs)
- Upon update (i, Δ) set $V \leftarrow V + \sigma(i) \cdot \Delta$
- Output V^2

Note that $V = \sum_{i=1}^n \sigma(i)x_i$ and

- $V \leftarrow 0$
- Pick hash function $\sigma : [n] \rightarrow \{-1, 1\}$ (random signs)
- Upon update (i, Δ) set $V \leftarrow V + \sigma(i) \cdot \Delta$
- Output V^2

Note that $V = \sum_{i=1}^n \sigma(i)x_i$ and

$$V^2 = \sum_{i,j} \sigma(i)\sigma(j)x_i x_j.$$

Let $X := V^2$.

SO WHAT ABOUT $\mathbb{E}(X)$?

It holds that

$$\begin{aligned}\mathbb{E}(V^2) &= \mathbb{E}\left(\sum_{i,j} \sigma(i)\sigma(j)x_i x_j\right) = \\ &= \sum_{i,j} \mathbb{E}(\sigma(i) \cdot \sigma(j)) \cdot x_i x_j = \\ &= \sum_{i \neq j} \mathbb{E}(\sigma(i)) \cdot \mathbb{E}(\sigma(j)) \cdot x_i x_j + \sum_i x_i^2 = \\ &= \sum_{i \neq j} 0 \cdot 0 \cdot x_i x_j + \sum_i x_i^2 = \sum_i x_i^2.\end{aligned}$$

AND WHAT ABOUT $\text{Var}(X^2)$?

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = \mathbb{E}(V^4) - (\mathbb{E}(V^2))^2 = \\ &\mathbb{E}\left(\sum_{i_1, i_2, i_3, i_4} \sigma(i_1)\sigma(i_2)\sigma(i_3)\sigma(i_4)X_{i_1}X_{i_2}X_{i_3}X_{i_4}\right) = \\ &\left(\sum_{i_1, i_2, i_3, i_4} \mathbb{E}(\sigma(i_1)\sigma(i_2)\sigma(i_3)\sigma(i_4))X_{i_1}X_{i_2}X_{i_3}X_{i_4}\right) = \\ &\sum_{i=1}^n X_i^4 + \sum_{i,j} 6X_i^2X_j^2 \leq 3 \cdot \sum_{i=1}^n X_i^2 = 3\mathbb{E}(X).\end{aligned}$$

AND WHAT ABOUT $\text{Var}(X^2)$?

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = \mathbb{E}(V^4) - (\mathbb{E}(V^2))^2 = \\ &= \mathbb{E} \left(\sum_{i_1, i_2, i_3, i_4} \sigma(i_1)\sigma(i_2)\sigma(i_3)\sigma(i_4)X_{i_1}X_{i_2}X_{i_3}X_{i_4} \right) = \\ &= \left(\sum_{i_1, i_2, i_3, i_4} \mathbb{E}(\sigma(i_1)\sigma(i_2)\sigma(i_3)\sigma(i_4))X_{i_1}X_{i_2}X_{i_3}X_{i_4} \right) = \\ &= \sum_{i=1}^n X_i^4 + \sum_{i,j} 6X_i^2X_j^2 \leq 3 \cdot \sum_{i=1}^n X_i^2 = 3\mathbb{E}(X).\end{aligned}$$

If we apply Chebyshev's inequality, we run into the same issue as before (too large variance), so let's take $\frac{1}{\epsilon^2}$ different estimator and average them!

AMS SKETCH

- $r := \Theta(\frac{1}{\epsilon^2})$
- $V_j \leftarrow 0$ for all $j \in [r]$
- Pick hash functions $\sigma_j : [n] \rightarrow \{-1, 1\}$ (random signs) for all $j \in [r]$

AMS SKETCH

- $r := \Theta(\frac{1}{\epsilon^2})$
- $V_j \leftarrow 0$ for all $j \in [r]$
- Pick hash functions $\sigma_j : [n] \rightarrow \{-1, 1\}$ (random signs) for all $j \in [r]$
- Upon update (i, Δ) set $V_j \leftarrow V_j + \sigma_j(i) \cdot \Delta$ for all $j \in [r]$

- $r := \Theta(\frac{1}{\epsilon^2})$
- $V_j \leftarrow 0$ for all $j \in [r]$
- Pick hash functions $\sigma_j : [n] \rightarrow \{-1, 1\}$ (random signs) for all $j \in [r]$
- Upon update (i, Δ) set $V_j \leftarrow V_j + \sigma_j(i) \cdot \Delta$ for all $j \in [r]$
- Output $\frac{1}{r} \sum_{j=1}^r V_j^2$

AMS SKETCH

- $r := \Theta(\frac{1}{\epsilon^2})$
- $V_j \leftarrow 0$ for all $j \in [r]$
- Pick hash functions $\sigma_j : [n] \rightarrow \{-1, 1\}$ (random signs) for all $j \in [r]$
- Upon update (i, Δ) set $V_j \leftarrow V_j + \sigma_j(i) \cdot \Delta$ for all $j \in [r]$
- Output $\frac{1}{r} \sum_{j=1}^r V_j^2$

To store each V_j we need $O(\log n)$ bits of space, for a total of all $O(\frac{\log n}{\epsilon^2})$ overall. Storing a compact representations of the hash functions can also be done in the same space.

- $r := \Theta(\frac{1}{\epsilon^2})$
- $V_j \leftarrow 0$ for all $j \in [r]$
- Pick hash functions $\sigma_j : [n] \rightarrow \{-1, 1\}$ (random signs) for all $j \in [r]$
- Upon update (i, Δ) set $V_j \leftarrow V_j + \sigma_j(i) \cdot \Delta$ for all $j \in [r]$
- Output $\frac{1}{r} \sum_{j=1}^r V_j^2$

To store each V_j we need $O(\log n)$ bits of space, for a total of all $O(\frac{\log n}{\epsilon^2})$ overall. Storing a compact representations of the hash functions can also be done in the same space.

Update time: $O(\frac{1}{\epsilon^2})$ assuming operations in happen constant time.

Query time: $O(\frac{1}{\epsilon^2})$ assuming operations happen in constant time.

WHAT WE'VE SEEN IN THIS LECTURE.

- The abstract architecture of a streaming algorithm.
- Applications of streaming algorithms.
- Distinct elements and the power of randomness.
- Estimating the F_2 moment.

WHAT WE'VE SEEN IN THIS LECTURE.

- The abstract architecture of a streaming algorithm.
- Applications of streaming algorithms.
- Distinct elements and the power of randomness.
- Estimating the F_2 moment.

Next lecture: Sketching and the foundations of Dimensionality Reduction.

WHAT WE'VE SEEN IN THIS LECTURE.

- The abstract architecture of a streaming algorithm.
- Applications of streaming algorithms.
- Distinct elements and the power of randomness.
- Estimating the F_2 moment.

Next lecture: Sketching and the foundations of Dimensionality Reduction.

Thank you!