Lecture 2: Frequency Moments, Heavy Hitters

Michael Kapralov

EPFL

May 24, 2017

Linear sketching



Sketching algorithms for basic statistics, and then graph sketching

In this lecture:

- Frequency moments (AMS sketch)
- Heavy hitters (CountSketch)

In this lecture:

Frequency moments (AMS sketch)

Heavy hitters (CountSketch)

Goal: approximate

$$||x||_2 = \sqrt{\sum_{i \in [n]} x_i^2}$$

from a stream of increments/decrements to x_i .

$x \in \mathbb{R}^n$ 1 2 3 4 5 6 7 8 9 10

- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



3 4

- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- 3 4 3
- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



3 4 3 2

- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



3 4 3 2 10

- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



3 4 3 2 10 3

- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



3 4 3 2 10 3 1

- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- ► Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$



- Initially, x = 0
- Insertion of *i* interpreted as

$$x_i := x_i + 1$$

Goal: approximate

$$||x||_2 = \sqrt{\sum_{i \in [n]} x_i^2}$$

from a stream of increments/decrements to x_i .

Goal: approximate

$$||x||_2 = \sqrt{\sum_{i \in [n]} x_i^2}$$

from a stream of increments/decrements to x_i .

Choose r_1, \ldots, r_n to be i.i.d. r.v., with

$$\Pr[r_i = +1] = \Pr[r_i = -1] = 1/2.$$

Goal: approximate

$$||x||_2 = \sqrt{\sum_{i \in [n]} x_i^2}$$

from a stream of increments/decrements to x_i .

Choose r_1, \ldots, r_n to be i.i.d. r.v., with

$$\Pr[r_i = +1] = \Pr[r_i = -1] = 1/2.$$

Maintain

$$Z = \sum_{i=1}^{n} r_i x_i$$

under increments/decrements of x.

Goal: approximate

$$||x||_2 = \sqrt{\sum_{i \in [n]} x_i^2}$$

from a stream of increments/decrements to x_i .

Choose r_1, \ldots, r_n to be i.i.d. r.v., with

$$\Pr[r_i = +1] = \Pr[r_i = -1] = 1/2.$$

Maintain

$$Z = \sum_{i=1}^{n} r_i x_i$$

under increments/decrements of x.

Basic algorithm: output Z^2

Goal: approximate

$$||x||_2 = \sqrt{\sum_{i \in [n]} x_i^2}$$

from a stream of increments/decrements to x_i .

Choose r_1, \ldots, r_n to be i.i.d. r.v., with

$$\Pr[r_i = +1] = \Pr[r_i = -1] = 1/2.$$

Maintain

$$Z = \sum_{i=1}^{n} r_i x_i$$

under increments/decrements of x.

Basic algorithm: output Z^2

Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability'

Alon-Matias-Szegedy – analysis (expectation) Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability'

Alon-Matias-Szegedy – analysis (expectation)

Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability'

Compute expectation of Z^2 , then bound the variance
Alon-Matias-Szegedy – analysis (expectation)

Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability'

Compute expectation of Z^2 , then bound the variance

Expectation:

$$\mathbf{E}[Z^{2}] = \mathbf{E}\left[\left(\sum_{i=1}^{n} r_{i} x_{i}\right)^{2}\right]$$

= $\sum_{i=1}^{n} \sum_{j=1}^{n} \mathbf{E}[r_{i} r_{j} x_{i} x_{j}]$
= $\sum_{i=1}^{n} \sum_{j=1}^{n} \mathbf{E}[r_{i} r_{j}] x_{i} x_{j}$
= $\sum_{i=1}^{n} x_{i}^{2} + \sum_{i,j:i \neq j}^{n} \mathbf{E}[r_{i}] \mathbf{E}[r_{j}] x_{i} x_{j}$
= $\sum_{i=1}^{n} x_{i}^{2}$
= $||x||_{2}^{2}$

Alon-Matias-Szegedy – analysis (expectation)

Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability'

Compute expectation of Z^2 , then bound the variance

Expectation:

$$\mathbf{E}[Z^2] = \mathbf{E}\left[\left(\sum_{i=1}^n r_i x_i\right)^2\right]$$
$$= \sum_{i=1}^n \sum_{j=1}^n \mathbf{E}[r_i r_j x_i x_j]$$
$$= \sum_{i=1}^n \sum_{j=1}^n \mathbf{E}[r_i r_j] x_i x_j$$
$$= \sum_{i=1}^n x_i^2 + \sum_{i,j:i\neq j}^n \mathbf{E}[r_i] \mathbf{E}[r_j] x_i x_j$$
$$= \sum_{i=1}^n x_i^2$$
$$= ||x||_2^2 \quad (\text{our estimator is unbiased!})$$

Alon-Matias-Szegedy – analysis (variance)

Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability' Bound the variance $Var[Z^2] = E[Z^4] - (E[Z^2])^2$?

Alon-Matias-Szegedy – analysis (variance)

Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability' Bound the variance $Var[Z^2] = E[Z^4] - (E[Z^2])^2$? Compute

 $\mathbf{E}[Z^4] = \mathbf{E}\left[(\sum_{i=1}^n r_i x_i) (\sum_{j=1}^n r_j x_j) (\sum_{k=1}^n r_k x_k) (\sum_{l=1}^n r_l x_l) \right]$

Alon-Matias-Szegedy – analysis (variance)

Want to claim that Z^2 is 'close' to $||x||_2^2$ with 'high probability' Bound the variance $Var[Z^2] = E[Z^4] - (E[Z^2])^2$? Compute

 $\mathbf{E}[Z^{4}] = \mathbf{E}\left[\left(\sum_{i=1}^{n} r_{i} x_{i}\right)\left(\sum_{i=1}^{n} r_{j} x_{j}\right)\left(\sum_{k=1}^{n} r_{k} x_{k}\right)\left(\sum_{l=1}^{n} r_{l} x_{l}\right)\right]$

Can be decomposed as follows:

- $\sum_{i=1}^{n} (r_i x_i)^4$ expectation $\sum_{i=1}^{n} x_i^4$
- $6\sum_{i < j} (r_i r_j x_i x_j)^2$ expectation $6\sum_{i < j} x_i^2 x_j^2$
- ► Terms involving a single *r_ix_i* − expectation zero.

In total: $\sum_{i=1}^{n} x_i^4 + 6 \sum_{i < j} x_i^2 x_j^2$

Alon-Matias-Szegedy – analysis (variance) Bound the variance $Var[Z^2] = E[Z^4] - (E[Z^2])^2$? Computed

$$\mathbf{E}[Z^4] = \sum_{i=1}^{n} x_i^4 + 6 \sum_{i < j} x_i^2 x_j^2$$

Alon-Matias-Szegedy – analysis (variance) Bound the variance $Var[Z^2] = E[Z^4] - (E[Z^2])^2$? Computed

$$\mathbf{E}[Z^4] = \sum_{i=1}^n x_i^4 + 6 \sum_{i < j} x_i^2 x_j^2$$

So

$$Var[Z^{2}] = \mathbf{E}[Z^{4}] - (\mathbf{E}[Z^{2}])^{2}$$

$$= \sum_{i=1}^{n} x_{i}^{4} + 6 \sum_{i < j} x_{i}^{2} x_{j}^{2} - (\sum_{i=1}^{n} x_{i}^{2})^{2}$$

$$= \sum_{i=1}^{n} x_{i}^{4} + 6 \sum_{i < j} x_{i}^{2} x_{j}^{2} - \sum_{i=1}^{n} x_{i}^{4} - 2 \sum_{i < j} x_{i}^{2} x_{j}^{2}$$

$$\leq 4 \sum_{i < j} x_{i}^{2} x_{j}^{2}$$

$$\leq 2(\sum_{i=1}^{n} x_{i}^{2})^{2}$$

We showed that

•
$$\mathbf{E}[Z^2] = ||x||_2^2$$

• $\sigma^2 = \mathbf{Var}[Z^2] \le 2||x||_2^4$

We showed that

•
$$\mathbf{E}[Z^2] = ||x||_2^2$$

• $\sigma^2 = \mathbf{Var}[Z^2] \le 2||x||_2^4$

So by Chebyshev's inequality for $t \ge 1$

$$\mathbf{Pr}[|Z^2 - \mathbf{E}[Z^2]| \ge t\sigma] \le 1/t^2$$

We showed that

•
$$\mathbf{E}[Z^2] = ||x||_2^2$$

• $\sigma^2 = \mathbf{Var}[Z^2] \le 2||x||_2^4$

So by Chebyshev's inequality for $t \ge 1$

$$\mathbf{Pr}[|Z^2 - \mathbf{E}[Z^2]| \ge t\sigma] \le 1/t^2$$

and we get

$$\mathbf{Pr}[|Z^2 - ||x||_2^2] \ge \sqrt{2t}||x||_2^2] \le 1/t^2$$

We showed that

•
$$\mathbf{E}[Z^2] = ||x||_2^2$$

• $\sigma^2 = \mathbf{Var}[Z^2] \le 2||x||_2^4$

So by Chebyshev's inequality for $t \ge 1$

$$\mathbf{Pr}[|Z^2 - \mathbf{E}[Z^2]| \ge t\sigma] \le 1/t^2$$

and we get

$$\mathbf{Pr}[|Z^2 - ||x||_2^2] \ge \sqrt{2t}||x||_2^2] \le 1/t^2$$

Not good...

We showed that

•
$$\mathbf{E}[Z^2] = ||x||_2^2$$

• $\sigma^2 = \mathbf{Var}[Z^2] \le 2||x||_2^4$

So by Chebyshev's inequality for $t \ge 1$

$$\mathbf{Pr}[|Z^2 - \mathbf{E}[Z^2]| \ge t\sigma] \le 1/t^2$$

and we get

$$\mathbf{Pr}[|Z^2 - ||x||_2^2] \ge \sqrt{2t}||x||_2^2] \le 1/t^2$$

Not good...but can reduce variance by averaging!

Actual algorithm:

• Maintain
$$Z_1, \ldots, Z_k, Z_i = \sum_{j=1}^n r_j^i x_j$$

• Output
$$A := \frac{1}{k} \sum_{i=1}^{k} Z_i^2$$

Now

$$\operatorname{Var}[A] = \operatorname{Var}\left[\frac{1}{k}\sum_{i=1}^{k}Z_{i}^{2}\right] = \frac{1}{k}\operatorname{Var}\left[Z^{2}\right] \leq (2/k)||x||_{2}^{4},$$

Actual algorithm:

• Maintain
$$Z_1, \ldots, Z_k, Z_j = \sum_{j=1}^n r_j^j x_j$$

• Output
$$A := \frac{1}{k} \sum_{i=1}^{k} Z_i^2$$

Now

$$\operatorname{Var}[A] = \operatorname{Var}\left[\frac{1}{k}\sum_{i=1}^{k}Z_{i}^{2}\right] = \frac{1}{k}\operatorname{Var}\left[Z^{2}\right] \leq (2/k)||x||_{2}^{4},$$

and by Chebyshev's inequality

$$\Pr\left[|A - ||x||_2^2| \ge t \cdot (2/k)^{1/2} ||x||_2^2\right] \le 1/t^2,$$

so setting $k = O(1/\epsilon^2)$ and t = 10 suffices for a $(1 \pm \epsilon)$ -approximation with probability $\ge 99/100$!

How much space do we need to store r_i 's?

4-wise independence suffices, hence $O(\log n)$ space

Can we reduce failure probability to $1 - \delta$?

Can we reduce failure probability to $1 - \delta$?

Use O(1/(ε²δ)) repetitions – bad dependence on δ

Can we reduce failure probability to $1 - \delta$?

- Use O(1/(ε²δ)) repetitions bad dependence on δ
- ► Median trick: keep T = O(log(1/δ)) copies of the estimator, output the median

Can we reduce failure probability to $1 - \delta$?

- Use O(1/(ε²δ)) repetitions bad dependence on δ
- ► Median trick: keep T = O(log(1/δ)) copies of the estimator, output the median

Let $Y_t = 1$ if *t*-th algorithm fails, and 0 otherwise.

We have $\mathbf{E}[Y_t] \le 1/100$, so by the Chernoff bound

$$\begin{aligned} & \mathbf{Pr} \left[|\text{median}_{i=1,...,T}(A_i) - ||x||_2^2 | > \varepsilon ||x||_2^2 \right] \\ & \leq \mathbf{Pr} [\text{at least half of } A_i \text{ fail, } i = 1,...,T] \\ & \leq \mathbf{Pr} [\sum_{t=1}^T Y_i \ge T/2] \\ & \leq e^{-\Omega(T)}. \end{aligned}$$

So setting $T = O(\log(1/\delta))$ suffices.

Downside of the median trick: nonlinear embedding

Median trick not needed if we have enough independence

Johnson-Lindenstrauss transform (see Ilya's lecture)

Take (randomized) linear measurements of the input



Take (randomized) linear measurements of the input



Can get $(1 \pm \varepsilon)$ -approximation to $||x||^2$ with $O(\frac{1}{\varepsilon^2} \log(1/\delta))$ rows

Take (randomized) linear measurements of the input



Can get $(1 \pm \varepsilon)$ -approximation to $||x||^2$ with $O(\frac{1}{\varepsilon^2} \log(1/\delta))$ rows Easy to maintain linear sketches in the (dynamic) streaming model In this lecture:

- Frequency moments (AMS sketch)
- Heavy hitters (CountSketch)

In this lecture:

- Frequency moments (AMS sketch)
- Heavy hitters (CountSketch)

Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)

Single pass over the data: i_1, i_2, \dots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)

Single pass over the data: i_1, i_2, \dots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)

Single pass over the data: i_1, i_2, \dots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \dots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \dots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



3 4 6 3 2 10

Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \dots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: $i_1, i_2, ..., i_N$

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)


Single pass over the data: $i_1, i_2, ..., i_N$

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: $i_1, i_2, ..., i_N$

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: $i_1, i_2, ..., i_N$

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: $i_1, i_2, ..., i_N$

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)



Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)





Estimate the dominant IP flows through a router

	destination										
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



	destination											
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
Δ	Ο	Ο	Ο	Ο	Ο	Ο	Ο	Ο				



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	1	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst
DA	TA



	destination											
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	1	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
Δ	Ο	Ο	Δ	Δ	Δ	Δ	Ο	Ο				



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	1	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst
DA	TA



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	1	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	2	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst
DA	TA



	destination											
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	2	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	1	0	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	2	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	1	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst				
DA	ТА				



	destination										
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	2	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	1	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination										
0	0	0	0	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		
0	0	0	2	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		
0	1	0	0	0	0	0	0	0		
0	0	0	0	0	1	0	0	0		
0	0	1	0	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		

Src	Dst
DA	ТА



	destination										
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	2	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	1	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	3	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	1	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst
DA	TA



	destination										
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	3	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	1	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination										
0	0	0	0	0	0	0	0	0		
0	0	0	0	0	0	1	0	0		
0	0	0	0	0	0	0	0	0		
0	0	0	3	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		
0	1	0	0	0	0	0	0	0		
0	0	0	0	0	1	0	0	0		
0	0	1	0	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		

Src	Dst
DA	TA



	destination										
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	3	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	1	0	0	0	0	0	0	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination										
0	0	0	0	0	0	0	0	0		
0	0	0	0	0	0	1	0	0		
0	0	0	0	0	0	0	0	0		
0	0	0	3	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		
0	2	0	0	0	0	0	0	0		
0	0	0	0	0	1	0	0	0		
0	0	1	0	0	0	0	0	0		
0	0	0	0	0	0	0	0	0		

Src	Dst				
DATA					



destination								
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0
0	0	0	3	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	2	0	0	0	0	0	0	0
0	0	0	0	0	1	0	0	0
0	0	1	0	0	0	0	0	0
Δ	Ο	Ο	Ο	Ο	Ο	Ο	Ο	Ο



destination									
0	0	0	0	0	0	0	0	0	
0	0	0	0	0	0	1	0	0	
0	0	0	0	0	0	0	0	0	
0	0	0	3	0	0	0	0	0	
0	0	0	0	0	0	0	0	0	
0	2	0	0	0	0	0	1	0	
0	0	0	0	0	1	0	0	0	
0	0	1	0	0	0	0	0	0	
0	0	0	0	0	0	0	0	0	

Src	Dst
DA	TA



destination								
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0
0	0	0	3	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	2	0	0	0	0	0	1	0
0	0	0	0	0	1	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0



destination									
0	0	0	0	0	0	0	0	0	
0	0	0	0	0	0	1	0	0	
0	0	0	0	0	0	0	0	0	
0	0	0	4	0	0	0	0	0	
0	0	0	0	0	0	0	0	0	
0	2	0	0	0	0	0	1	0	
0	0	0	0	0	1	0	0	0	
0	0	1	0	0	0	0	0	0	
0	0	0	0	0	0	0	0	0	

Src	Dst
DA	TA


destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	4	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	2	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	4	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	3	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst
DA	TA



destination											
0	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	4	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	3	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination											
1	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	4	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	3	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst
DA	ТА



destination											
1	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	4	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	3	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination											
1	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	4	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	3	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination												
1	0	0	0	0	0	0	0	0				
0	0	0	0	0	0	1	0	0				
0	0	0	0	0	0	0	0	0				
0	0	0	4	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				
0	4	0	0	0	0	0	1	0				
0	0	0	0	0	1	0	0	0				
0	0	1	0	0	0	0	0	0				
0	0	0	0	0	0	0	0	0				

Src	Dst
DA	ТА



destination											
1	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	4	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	4	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			



destination											
1	0	0	0	0	0	0	0	0			
0	0	0	0	0	0	1	0	0			
0	0	0	0	0	0	0	0	0			
0	0	0	5	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			
0	4	0	0	0	0	0	1	0			
0	0	0	0	0	1	0	0	0			
0	0	1	0	0	0	0	0	0			
0	0	0	0	0	0	0	0	0			

Src	Dst
DA	TA



destination								
1	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0
0	0	0	5	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	4	0	0	0	0	0	1	0
0	0	0	0	0	1	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0



destination								
1	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0
0	0	0	5	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	4	0	0	0	0	0	1	0
0	0	0	0	0	1	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0



Estimate the dominant IP flows through a router

destination								
1	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	0
0	0	0	5	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	4	0	0	0	0	0	1	0
0	0	0	0	0	1	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0



Estimate the dominant IP flows through a router





Estimate the dominant IP flows through a router



Trivial: store all distinct IP pairs Space complexity: $\Theta(N)$



Estimate the dominant IP flows through a router



Trivial: store all distinct IP pairs Space complexity: $\Theta(N)$

This lecture: solve in space $O(\log N)$

Exponential improvement!

Given a set of items as a stream (e.g. queries on google.com over a period of time)

Geneva to NYC, coffee in Geneva, Geneva to NYC

Given a set of items as a stream (e.g. queries on google.com over a period of time)

Geneva to NYC, coffee in Geneva, Geneva to NYC

Find the most frequent items in the set

Geneva to NYC, coffee in Geneva

Given a set of items as a stream (e.g. queries on google.com over a period of time)

Geneva to NYC, coffee in Geneva, Geneva to NYC

Find the most frequent items in the set

Geneva to NYC, coffee in Geneva

Given a set of items as a stream (e.g. queries on google.com over a period of time)

Geneva to NYC, coffee in Geneva, Geneva to NYC

Find the most frequent items in the set

Geneva to NYC, coffee in Geneva

	Trivial	This lecture
Solution	hash <string> h;</string>	COUNTSKETCH
Space	<pre># of distinct items</pre>	$O(\log N)$

Heavy hitters problem

Single pass over the data: i_1, i_2, \ldots, i_N

Assume N is known

Output k most frequent items

(Heavy hitters)

Small storage: will get O(k log N)

Much better than storing all items!

Goal: design a small space data structure

FINDTOP(S, k): returns top k most frequent items seen so far

Goal: design a small space data structure

FINDTOP(S, k): returns top k most frequent items seen so far

Useful to first design

POINTQUERY(*S*, *i*): processes stream, then for any query item *i* can return f_i =number of times item *i* appeared

Denote the number of times item *i* appears in the stream by f_i (frequency of *i*)

Assume elements are ordered by frequency: $f_1 \ge f_2 \ge ... \ge f_m$

Denote the number of times item *i* appears in the stream by f_i (frequency of *i*)

Assume elements are ordered by frequency: $f_1 \ge f_2 \ge ... \ge f_m$

POINTQUERY(S, i) in space $O(k \log N)$?

Denote the number of times item *i* appears in the stream by f_i (frequency of *i*)

Assume elements are ordered by frequency: $f_1 \ge f_2 \ge ... \ge f_m$

POINTQUERY(S, i) in space $O(k \log N)$?

Impossible in general...

Imagine a stream where all elements occur with about the same frequency

FINDAPPROXTOP(S, k, ε): returns set of k items such that $f_i \ge (1 - \varepsilon)f_k$ for all reported i

APPROXPOINTQUERY(S, i, ε): processes stream, then for any query item *i* can return approximation $\hat{f}_i \in [f_i - \varepsilon f_k, f_i + \varepsilon f_k]$

FINDAPPROXTOP(S, k, ε): returns set of k items such that $f_i \ge (1 - \varepsilon)f_k$ for all reported i

APPROXPOINTQUERY(S, i, ε): processes stream, then for any query item *i* can return approximation $\hat{f}_i \in [f_i - \varepsilon f_k, f_i + \varepsilon f_k]$

In this lecture: find most frequent (head) items if they contribute the bulk of the stream under some measure

- 1. Finding top k elements via (APPROX)POINTQUERY
- 2. Basic version of APPROXPOINTQUERY
- 3. APPROXPOINTQUERY and the COUNTSKETCH algorithm

- 1. Finding top k elements via (APPROX)POINTQUERY
- 2. Basic version of APPROXPOINTQUERY
- 3. APPROXPOINTQUERY and the COUNTSKETCH algorithm

Given: stream $S = (i_1, \ldots, i_N)$

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for POINTQUERY currMax=NULL; currFreq=0;

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for POINTQUERY
 currMax=NULL; currFreq=0;
For p = 1,...,N
 Compute frequency f ← POINTQUERY(ip)

Given: stream $S = (i_1, \ldots, i_N)$

```
end if
```

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for POINTQUERY

```
currMax=NULL; currFreq=0;
```

For *p* = 1,...,*N*

Compute frequency $f \leftarrow \text{POINTQUERY}(i_p)$

```
If currMax==ip
currFreq=f; continue;
end if
```

```
If currMax==NULL
```

currMax=*i*_p; currFreq=1;

else

```
If currFreq<f
    currMax=ip; currFreq=f;
    end if
    end if
end for</pre>
```

Why does this work?

At each point in the stream currMax is either NULL or the most frequent element so far...

At each point in the stream currMax is either NULL or the most frequent element so far...

What about finding k most frequent elements for k > 1?
Given: stream $S = (i_1, \ldots, i_N)$

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for **POINTQUERY** a heap *H* of at most *k* items, by count

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for POINTQUERY a heap H of at most k items, by count For p = 1,...,NCompute frequency $f \leftarrow POINTQUERY(i_p)$

Given: stream $S = (i_1, \ldots, i_N)$

```
Maintain: data structure for POINTQUERY
a heap H of at most k items, by count
For p = 1,...,N
Compute frequency f \leftarrow POINTQUERY(i_p)
If H contains i_p
update i_p's key to f; continue;
end if
```

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for **POINTQUERY** a heap *H* of at most *k* items, by count

For *p* = 1,...,*N*

```
Compute frequency f \leftarrow \text{POINTQUERY}(i_p)
```

```
If H contains ip
update ip's key to f; continue;
end if
```

```
If H contains < k elements,
```

```
add < f, i_p > to H
```

else

 $< f_{min}, i_{min} > \leftarrow$ element in *H* with smallest key

```
If fmin < f, insert < f, ip > and evict < fmin, imin >
end if
end for
```

Why does this work?

(Assume the set of top *k* items is unique for simplicity)

For every item *i* in top *k* let p_i denote the last position where *i* occurs

Why does this work?

(Assume the set of top *k* items is unique for simplicity)

For every item *i* in top *k* let p_i denote the last position where *i* occurs

Note that

- 1. at position *p_i* element *i* is inserted into heap *H* if it was not in *H* at that time
- 2. element *i* is never evicted after p_i

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for **POINTQUERY** a heap *H* of at most *k* items, by count

For *p* = 1,...,*N*

```
Compute frequency f \leftarrow \text{POINTQUERY}(i_p)
```

```
If H contains ip
update ip's key to f; continue;
end if
```

```
If H contains < k elements,
```

```
add < f, i_p > to H
```

else

 $< f_{min}, i_{min} > \leftarrow$ element in *H* with smallest key

```
If fmin < f, insert < f, ip > and evict < fmin, imin >
end if
end for
```

Why is this useful? We know that POINTQUERY requires essentially storing the entire stream...

Why is this useful? We know that POINTQUERY requires essentially storing the entire stream...

A similar reduction shows that APPROXPOINTQUERY implies FINDAPPROXTOP! APPROXPOINTQUERY implies FINDAPPROXTOP

Given: stream $S = (i_1, \ldots, i_N)$

Maintain: data structure for APPROXPOINTQUERY a heap *H* of at most *k* items, by count

For *p* = 1,...,*N*

Compute frequency $\hat{f} \leftarrow \text{APPROXPOINTQUERY}(i_p)$

If *H* contains i_p update i_p 's key to \hat{f} ; continue; end if

```
If H contains < k elements,
```

```
add <\hat{f}, i_p > \text{to } H
```

else

 $\langle \hat{f}_{min}, i_{min} \rangle \leftarrow$ element in *H* with smallest key If $\hat{f}_{min} \langle \hat{f}, \text{ insert } \langle \hat{f}, i_p \rangle$ and evict $\langle \hat{f}_{min}, i_{min} \rangle$ end if end for FINDAPPROXTOP(S, k, ε): returns set S of k items such that $f_i \ge (1 - \varepsilon)f_k$ for all $i \in S$

APPROXPOINTQUERY(S, i, ε): returns $\hat{f}_i \in [f_i - \varepsilon f_k, f_i + \varepsilon f_k]$

In what follows: APPROXPOINTQUERY in small space

Observe a stream of updates, maintain small space data structure

Task: after observing the stream, given $i \in \{1, 2, ..., m\}$, compute estimate \hat{f}_i of f_i

In what follows: APPROXPOINTQUERY in small space

Observe a stream of updates, maintain small space data structure

Task: after observing the stream, given $i \in \{1, 2, ..., m\}$, compute estimate \hat{f}_i of f_i

To be specified:

- space complexity?
- quality of approximation?
- success probability?

- 1. Finding top k elements via (APPROX)POINTQUERY
- 2. Basic version of APPROXPOINTQUERY
- 3. APPROXPOINTQUERY and the COUNTSKETCH algorithm

- 1. Finding top k elements via (APPROX)POINTQUERY
- 2. Basic version of APPROXPOINTQUERY
- 3. APPROXPOINTQUERY and the COUNTSKETCH algorithm

1 2 3 4 5 6 7 8 9 10









1 4 6 1





1 4 6 1 2 10



1 4 6 1 2 10 1





1 4 6 1 2 10 1 5 1


































Will design a basic estimate with O(1) space complexity, analyze precision

Will design a basic estimate with O(1) space complexity, analyze precision

Choose a hash function $s: [m] \rightarrow \{-1, +1\}$ uniformly at random

NITIALIZE	Update(C, i)
<i>C</i> ← 0	$C \leftarrow C + s(i)$

Will design a basic estimate with O(1) space complexity, analyze precision

Choose a hash function $s: [m] \rightarrow \{-1, +1\}$ uniformly at random

NITIALIZEUPDATE(C, i)
$$C \leftarrow 0$$
 $C \leftarrow C + s(i)$

for every p = 1, ..., N (every element in the stream) UPDATE(C, i_p) end for

Will design a basic estimate with O(1) space complexity, analyze precision

Choose a hash function $s: [m] \rightarrow \{-1, +1\}$ uniformly at random

NITIALIZE	Update(C, i)
<i>C</i> ← 0	$C \leftarrow C + s(i)$

```
for every p = 1, ..., N (every element in the stream)
UPDATE(C, i_p)
end for
```

ESTIMATE(C, i) return $C \cdot s(i)$

How does one argue that a randomized estimate works?



How does one argue that a randomized estimate works?



How does one argue that a randomized estimate works?

Want to show that $C \cdot s(i)$ is close to f_i 'with high probability'



How does one argue that a randomized estimate works?

Want to show that $C \cdot s(i)$ is close to f_i 'with high probability'

Typically show this in two steps:

• show that $\mathbf{E}_{s}[C \cdot s(i)] = f_{i}$

(so $C \cdot s(i)$ is an unbiased estimate of f_i)

• show that $Var_s[C \cdot s(i)]$ is 'small'



How does one argue that a randomized estimate works?

Want to show that $C \cdot s(i)$ is close to f_i 'with high probability'

Typically show this in two steps:

- Show that E_s[C ⋅ s(i)] = f_i (so C ⋅ s(i) is an unbiased estimate of f_i)
- show that $Var_s[C \cdot s(i)]$ is 'small'

It then follows that $|C \cdot s(i) - f_i|$ is 'small' with high probability (essentially law of large numbers)

UPDATE(C, i) $C \leftarrow C + s(i)$

$$C \cdot s(i) = \sum_{p=1}^{N} s(i_p)s(i)$$

UPDATE(C, i) $C \leftarrow C + s(i)$

$$C \cdot s(i) = \sum_{p=1}^{N} s(i_p) s(i) = \sum_{j \in [m]} f_j \cdot s(j) s(i)$$

UPDATE(C, i) $C \leftarrow C + s(i)$

$$C \cdot s(i) = \sum_{p=1}^{N} s(i_p) s(i) = \sum_{j \in [m]} f_j \cdot s(j) s(i)$$
$$= f_i s(i)^2 + \sum_{j \in [m] \setminus i} f_j \cdot s(j) s(i)$$

UPDATE(C, i) $C \leftarrow C + s(i)$

$$C \cdot s(i) = \sum_{p=1}^{N} s(i_p) s(i) = \sum_{j \in [m]} f_j \cdot s(j) s(i)$$
$$= f_i s(i)^2 + \sum_{j \in [m] \setminus i} f_j \cdot s(j) s(i)$$
$$= f_i + \sum_{j \in [m] \setminus i} f_j \cdot s(j) s(i) \quad \longleftarrow \text{ random } \pm 1\text{'s}$$

UPDATE(C, i) $C \leftarrow C + s(i)$

$$C \cdot s(i) = f_i + \sum_{j \in [m] \setminus i} f_j \cdot s(j)s(i)$$

UPDATE(C, i) $C \leftarrow C + s(i)$

$$\mathbf{E}[C \cdot \mathbf{s}(i)] = f_i + \mathbf{E}[\sum_{j \in [m] \setminus i} f_j \cdot \mathbf{s}(j)\mathbf{s}(i)]$$

UPDATE(C, i) $C \leftarrow C + s(i)$

$$\mathbf{E}[C \cdot s(i)] = f_i + \mathbf{E}[\sum_{j \in [m] \setminus i} f_j \cdot s(j)s(i)]$$

= $f_i + \sum_{j \in [m] \setminus i} f_j \cdot \mathbf{E}[s(j)]\mathbf{E}[s(i)]$ (by independence of $s(i)$)

UPDATE(C, i) $C \leftarrow C + s(i)$

$$\mathbf{E}[C \cdot \mathbf{s}(i)] = f_i + \mathbf{E}[\sum_{j \in [m] \setminus i} f_j \cdot \mathbf{s}(j)\mathbf{s}(i)]$$

= $f_i + \sum_{j \in [m] \setminus i} f_j \cdot \mathbf{E}[\mathbf{s}(j)]\mathbf{E}[\mathbf{s}(i)]$ (by independence of $\mathbf{s}(i)$)
= f_i

UPDATE(C, i) $C \leftarrow C + s(i)$ ESTIMATE(C, i) return C · s(i)

$$\begin{split} \mathbf{E}[C \cdot \mathbf{s}(i)] &= f_i + \mathbf{E}[\sum_{j \in [m] \setminus i} f_j \cdot \mathbf{s}(j)\mathbf{s}(i)] \\ &= f_i + \sum_{j \in [m] \setminus i} f_j \cdot \mathbf{E}[\mathbf{s}(j)]\mathbf{E}[\mathbf{s}(i)] \text{ (by independence of } \mathbf{s}(i)) \\ &= f_i \end{split}$$

The mean is correct: our estimator is unbiased!

UPDATE(C, i) $C \leftarrow C + s(i)$ ESTIMATE(C, i) return C · s(i)

$$\begin{split} \mathbf{E}[C \cdot \mathbf{s}(i)] &= f_i + \mathbf{E}[\sum_{j \in [m] \setminus i} f_j \cdot \mathbf{s}(j)\mathbf{s}(i)] \\ &= f_i + \sum_{j \in [m] \setminus i} f_j \cdot \mathbf{E}[\mathbf{s}(j)]\mathbf{E}[\mathbf{s}(i)] \text{ (by independence of } \mathbf{s}(i)) \\ &= f_i \end{split}$$

The mean is correct: our estimator is unbiased!

Is the estimate $C \cdot s(i)$ close to f_i with high probability?

Theorem

For every random variable X with mean μ and variance σ^2 , and every $t \ge 1$ one has

$$\mathbf{Pr}[|X - \mu| \ge t \cdot \sigma] \le 1/t^2$$



Theorem

For every random variable X with mean μ and variance σ^2 , and every $t \ge 1$ one has

 $\mathbf{Pr}[|X - \mu| \ge t \cdot \sigma] \le 1/t^2$



Apply Chebyshev's inequality with $X = C \cdot s(i)$ and $\mu = \mathbf{E}[C \cdot s(i)]?$

Theorem

For every random variable X with mean μ and variance σ^2 , and every $t \ge 1$ one has

 $\mathbf{Pr}[|X - \mu| \ge t \cdot \sigma] \le 1/t^2$



Apply Chebyshev's inequality with $X = C \cdot s(i)$ and $\mu = \mathbf{E}[C \cdot s(i)]?$

A quantitative form of the 'law of large numbers'

Theorem

For every random variable X with mean μ and variance σ^2 , and every $t \ge 1$ one has

 $\mathbf{Pr}[|X - \mu| \ge t \cdot \sigma] \le 1/t^2$



Apply Chebyshev's inequality with $X = C \cdot s(i)$ and $\mu = \mathbf{E}[C \cdot s(i)]?$

A quantitative form of the 'law of large numbers'

Need to compute the variance $\sigma^2 = \mathbf{E}[(C \cdot s(i) - \mu)^2]$

UPDATE(C, i) $C \leftarrow C + s(i)$

We have

ESTIMATE(C, i) return $C \cdot s(i)$

$$C \cdot s(i) = f_i + \sum_{j \in [m] \setminus i} f_j \cdot s(j) s(i)$$

and

 $\mathbf{E}[C \cdot \mathbf{s}(i)] = f_i.$

UPDATE(C, i) $C \leftarrow C + s(i)$

We have

ESTIMATE(C, i) return $C \cdot s(i)$

$$C \cdot s(i) = f_i + \sum_{j \in [m] \setminus i} f_j \cdot s(j)s(i)$$

and

$$\mathbf{E}[C \cdot \mathbf{s}(i)] = f_i.$$

We need to bound

$$\mathsf{Var}(C \cdot s(i)) = \mathsf{E}[(C \cdot s(i) - \mathsf{E}[C \cdot s(i)])^2]$$
$$= \mathsf{E}[(C \cdot s(i) - f_i)^2]$$
$$= \mathsf{E}\left[\left(\sum_{j \in [m] \setminus i} f_j \cdot s(j)s(i)\right)^2\right]$$

UPDATE(C, i) $C \leftarrow C + s(i)$

$$(C \cdot s(i) - f_i)^2 = \left(\sum_{j \in [m] \setminus i} f_j \cdot s(j) s(i)\right)^2$$
$$= \sum_{j \in [m] \setminus i} \sum_{j' \in [m] \setminus i} f_j f_{j'} \cdot s(j) s(j') \cdot s^2(i)$$
$$= \sum_{j \in [m] \setminus i} \sum_{j' \in [m] \setminus i} f_j f_{j'} \cdot s(j) s(j')$$

UPDATE(C, i) $C \leftarrow C + s(i)$ ESTIMATE(C, i) return C · s(i)

$$\mathbf{E}[(C \cdot s(i) - f_i)^2] = \mathbf{E}[\sum_{j \in [m] \setminus i} \sum_{j' \in [m] \setminus i} f_j f_{j'} \cdot s(j) s(j')]$$
$$= \sum_{j \in [m] \setminus i} \sum_{j' \in [m] \setminus i} f_j f_{j'} \cdot \mathbf{E}[s(j) s(j')]$$
$$= \sum_{j \in [m] \setminus i} f_j^2$$

since

•
$$\mathbf{E}[s(j)s(j')] = \mathbf{E}[s(j)]\mathbf{E}[s(j')] = 0 \text{ for } j \neq j'.$$

UPDATE(C, i) $C \leftarrow C + s(i)$

UPDATE(C, i) $C \leftarrow C + s(i)$ ESTIMATE(C, i) return C · s(i)

We have proved that

$$\operatorname{Var}(C \cdot s(i)) = \operatorname{E}[(C \cdot s(i) - f_i)^2] = \sum_{j \in [m] \setminus i} f_j^2$$

UPDATE(C, i) $C \leftarrow C + s(i)$ ESTIMATE(C, i) return C · s(i)

We have proved that

$$\operatorname{Var}(C \cdot s(i)) = \operatorname{E}[(C \cdot s(i) - f_i)^2] = \sum_{j \in [m] \setminus i} f_j^2$$

By Chebyshev's inequality

$$\Pr\left[|C \cdot s(i) - f_i| \ge 8 \cdot \sqrt{\sum_{j \in [m] \setminus i} f_j^2}\right] \le 1/64$$
Basic estimate: variance

UPDATE(C, i) $C \leftarrow C + s(i)$ ESTIMATE(C, i) return C · s(i)

We have proved that

$$\operatorname{Var}(C \cdot s(i)) = \operatorname{E}[(C \cdot s(i) - f_i)^2] = \sum_{j \in [m] \setminus i} f_j^2$$

By Chebyshev's inequality

$$\Pr\left[|C \cdot s(i) - f_i| \ge 8 \cdot \sqrt{\sum_{j \in [m] \setminus i} f_j^2}\right] \le 1/64$$

So $C \cdot s(i)$ is close (?) to f_i with high probability

Basic estimate: variance

UPDATE(C, i) $C \leftarrow C + s(i)$ ESTIMATE(C, i) return C · s(i)

We have proved that

$$\operatorname{Var}(C \cdot s(i)) = \operatorname{E}[(C \cdot s(i) - f_i)^2] = \sum_{j \in [m] \setminus i} f_j^2$$

By Chebyshev's inequality

$$\Pr\left[|C \cdot s(i) - f_i| > \mathbf{8} \cdot \sqrt{\sum_{j \in [\mathbf{m}] \setminus i} f_j^2}\right] \le 1/64$$

So $C \cdot s(i)$ is close (?) to f_i with high probability

Basic estimate: summary

UPDATE(C, i) $C \leftarrow C + s(i)$

Estimate f_i up to

ESTIMATE(C, i) return $C \cdot s(i)$

$$8 \cdot \sqrt{\sum_{j \in [m] \setminus i} f_j^2}$$



Pro: works well for most frequent item, if other items are small

Basic estimate: summary

UPDATE(C, i) $C \leftarrow C + s(i)$

Estimate f_i up to

ESTIMATE(C, i) return $C \cdot s(i)$

$$8 \cdot \sqrt{\sum_{j \in [m] \setminus i} f_j^2}$$



Pro: works well for most frequent item, if other items are small Con: estimate for a small items contaminated by large items

Next: final APPROXPOINTQUERY and the COUNTSKETCH algorithm

COUNTSKETCH algorithm: find top k elements (approximately)

- hash items into O(k) buckets (i.e. substreams)
- run simple estimate on every bucket
- repeat O(log N) times independently, take median as answer

Main intuition: estimate large items from substreams like



Main intuition: estimate large items from substreams like



and small items from substreams like



- 1. Finding top k elements via (APPROX)POINTQUERY
- 2. Basic version of APPROXPOINTQUERY
- 3. APPROXPOINTQUERY and the COUNTSKETCH algorithm

- 1. Finding top k elements via (APPROX)POINTQUERY
- 2. Basic version of APPROXPOINTQUERY
- 3. APPROXPOINTQUERY and the COUNTSKETCH algorithm

APPROXPOINTQUERY and COUNTSKETCH

Main ideas:

1. run basic estimate on subsampled/hashed stream (reduces variance)

APPROXPOINTQUERY and COUNTSKETCH

Main ideas:

- 1. run basic estimate on subsampled/hashed stream (reduces variance)
- 2. aggregate independent estimates to boost confidence (take medians)

APPROXPOINTQUERY and COUNTSKETCH

Main ideas:

- 1. run basic estimate on subsampled/hashed stream (reduces variance)
- 2. aggregate independent estimates to boost confidence (take medians)



Hashing the items



Hashed into b = 8 buckets, get 8 subsampled streams

For item *i* its stream consists of $j \in [m]$ such that h(j) = h(i)

Hashing the items



Hashed into b = 8 buckets, get 8 subsampled streams

For item *i* its stream consists of $j \in [m]$ such that h(j) = h(i)

For example,

- subsampled stream of item 1 is {1, 6}
- subsampled stream of item 5 is {5, 7}





1

E.x. the subsampled stream of item 1 is $\{1, 6\}$







E.x. the subsampled stream of item 5 is $\{5, 7\}$



5

Final ApproxPointQuery

Choose

- t random hash functions h₁, h₂,..., h_t from items [m] to b ≈ k buckets {1,2,...,b}
- ► *t* random hash functions $s_1, s_2, ..., s_t$ from items [*m*] to $\{-1, +1\}$



Final ApproxPointQuery

Choose

- t random hash functions h₁, h₂,..., h_t from items [m] to b ≈ k buckets {1,2,...,b}
- ► *t* random hash functions $s_1, s_2, ..., s_t$ from items [*m*] to $\{-1, +1\}$



The algorithm runs *t* independent copies of basic estimate:

UPDATE(C, i) for $r \in [1:t]$ $C[r, h_r(i)] \leftarrow C[r, h_r(i)] + s_r(i)$ ESTIMATE(C, i) return median_r { $C[r, h_r(i)] \cdot s_r(i)$ } end for

Final ApproxPointQuery

Choose

- t random hash functions h₁, h₂,..., h_t from items [m] to b ≈ k buckets {1,2,...,b}
- ► *t* random hash functions $s_1, s_2, ..., s_t$ from items [*m*] to $\{-1, +1\}$



The algorithm runs *t* independent copies of basic estimate:

UPDATE(C, i) for $r \in [1:t]$ $C[r, h_r(i)] \leftarrow C[r, h_r(i)] + s_r(i)$ ESTIMATE(C, i) return median_r { $C[r, h_r(i)] \cdot s_r(i)$ } end for

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)]=f_i$

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

and

$$\mathbf{E}_{s}[(C[r,h_{r}(i)] \cdot \mathbf{s}_{r}(i) - f_{i})^{2}] = \sum_{j \neq i: \mathbf{h}_{r}(\mathbf{j}) = \mathbf{h}_{r}(\mathbf{i})} f_{j}^{2}$$

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

and

$$\mathbf{E}_{s}[(C[r, h_{r}(i)] \cdot s_{r}(i) - f_{i})^{2}] = \sum_{j \neq i: \mathbf{h}_{r}(\mathbf{j}) = \mathbf{h}_{r}(\mathbf{i})} f_{j}^{2}$$

How large can the variance be? Can be be reduced by making number of buckets *b* large?

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

and

$$\mathbf{E}_{s}[(C[r, h_{r}(i)] \cdot s_{r}(i) - f_{i})^{2}] = \sum_{j \neq i: \mathbf{h}_{r}(j) = \mathbf{h}_{r}(\mathbf{i})} f_{j}^{2}$$

How large can the variance be? Can be be reduced by making number of buckets *b* large?

YES: hashing into a sufficiently large number of buckets reduces estimation error to below $\varepsilon \cdot f_k$

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

and

$$\mathbf{E}_{s}[(C[r, h_{r}(i)] \cdot s_{r}(i) - f_{i})^{2}] = \sum_{j \neq i: \mathbf{h}_{r}(j) = \mathbf{h}_{r}(\mathbf{i})} f_{j}^{2}$$

How large can the variance be? Can be be reduced by making number of buckets *b* large?

YES: hashing into a sufficiently large number of buckets reduces estimation error to below $\varepsilon \cdot f_k$

O(log N) repetitions ensure estimates are correct for all *i* with high probability

UPDATE(C, i) for $r \in [1:t]$ $C[r, h_r(i)] \leftarrow C[r, h_r(i)] + s_r(i)$ end for

ESTIMATE(C, i)
return median_r {
$$C[r, h_r(i)] \cdot s_r(i)$$
}

UPDATE(C, i)
for
$$r \in [1:t]$$

 $C[r, h_r(i)] \leftarrow C[r, h_r(i)] + s_r(i)$
end for
ESTIMATE(C, i)
return median_r { $C[r, h_r(i)] \cdot s_r(i)$ }

Lemma If $b \ge 8 \max \left\{ k, \frac{32\sum_{j \in \text{TAIL}} f_j^2}{(\varepsilon f_k)^2} \right\}$ and $t = O(\log N)$, then for every $i \in [m]$

 $|\text{median}_r \{C[r, h_r(i)] \cdot s_r(i)\} - f_i(p)| \le \varepsilon f_k$

at every point $p \in [1 : N]$ in the stream.

 $(f_i(p) \text{ is the frequency of } i \text{ up to position } p)$

UPDATE(C, i)
for
$$r \in [1:t]$$

 $C[r, h_r(i)] \leftarrow C[r, h_r(i)] + s_r(i)$
end for
ESTIMATE(C, i)
return median_r { $C[r, h_r(i)] \cdot s_r(i)$ }

Lemma If $b \ge 8 \max \left\{ k, \frac{32\sum_{j \in \text{TAIL}} f_j^2}{(\varepsilon f_k)^2} \right\}$ and $t = O(\log N)$, then for every $i \in [m]$

 $|\text{median}_r \{C[r, h_r(i)] \cdot s_r(i)\} - f_i(p)| \le \varepsilon f_k$

at every point $p \in [1 : N]$ in the stream.

 $(f_i(p) \text{ is the frequency of } i \text{ up to position } p)$

Space complexity is O(blog N)

UPDATE(C, i)
for
$$r \in [1:t]$$

 $C[r, h_r(i)] \leftarrow C[r, h_r(i)] + s_r(i)$
ESTIMATE(C, i)
return median_r { $C[r, h_r(i)] \cdot s_r(i)$ }
end for

Lemma If $b \ge 8 \max \left\{ k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2} \right\}$ and $t = O(\log N)$, then for every $i \in [m]$

 $|\text{median}_r \{C[r, h_r(i)] \cdot s_r(i)\} - f_i(p)| \le \varepsilon f_k$

at every point $p \in [1 : N]$ in the stream.

 $(f_i(p) \text{ is the frequency of } i \text{ up to position } p)$

Space complexity is $O(b \log N)$

How large is *b*?



In practice, choose *b* subject to space constraints, detect elements with counts above $O\left(\epsilon \sqrt{\frac{1}{k}\sum_{j \in TAIL} f_j^2}\right)$

Then
$$f_1 = \sqrt{N}$$
, $f_i = 1$ for $i = 2, N - \sqrt{N}$.
Set $b = 8 \max \left\{ 1, \frac{32\sum_{j \in TA/L} f_j^2}{(\varepsilon f_1)^2} \right\}$

Then
$$f_1 = \sqrt{N}$$
, $f_i = 1$ for $i = 2, N - \sqrt{N}$.

Set
$$b = 8 \max\left\{1, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_1)^2}\right\}$$

We have
$$\sum_{j \in TAIL} f_j^2 = N - \sqrt{N} \le N$$
, and $f_1^2 = N$

Then
$$f_1 = \sqrt{N}$$
, $f_i = 1$ for $i = 2$, $N - \sqrt{N}$.

Set
$$b = 8 \max\left\{1, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_1)^2}\right\}$$

We have
$$\sum_{j \in TAIL} f_j^2 = N - \sqrt{N} \le N$$
, and $f_1^2 = N$

So
$$b = 8 \max \left\{ 1, \frac{32 \sum_{j \in TAIL} f_j^2}{(\varepsilon f_1)^2} \right\} = O(1/\varepsilon^2)$$
 suffices

Set k = 1. Suppose that 1 appears \sqrt{N} times in the stream, and other $N - \sqrt{N}$ elements are distinct

Then
$$f_1 = \sqrt{N}$$
, $f_i = 1$ for $i = 2$, $N - \sqrt{N}$.

Set
$$b = 8 \max\left\{1, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_1)^2}\right\}$$

We have
$$\sum_{j \in TAIL} f_j^2 = N - \sqrt{N} \le N$$
, and $f_1^2 = N$

So
$$b = 8 \max \left\{ 1, \frac{32 \sum_{j \in TAIL} f_j^2}{(\varepsilon f_1)^2} \right\} = O(1/\varepsilon^2)$$
 suffices

Remarkable, as 1 appears only in \sqrt{N} positions out of *N*: a vanishingly small fraction of positions!
Final algorithm: COUNTSKETCH

FINDAPPROXTOP(S, k, ε): returns set of k items such that $f_i \ge (1 - \varepsilon)f_k$ for all returned i

(In fact also every *i* with $f_i \ge (1 - \varepsilon)f_k$ is reported)

APPROXPOINTQUERY(*S*, *i*, ε): returns $\hat{f}_i \in [f_i - \varepsilon f_k, f_i + \varepsilon f_k]$

Find head items if they contribute the bulk of the stream in ℓ_2 sense

CountSketch: proof details

Lemma If $b \ge 8 \max\left\{k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2}\right\}$ and $t \ge A \log N$ for an absolute constant A > 0, then for every $i \in [m]$

 $|\text{median}_r \{C[r, h_r(i)] \cdot s_r(i)\} - f_i(p)| \le \varepsilon f_k$

at every point $p \in [1 : N]$ in the stream with high probability.

 $(f_i(p) \text{ is the frequency of } i \text{ up to position } p)$

Lemma If $b \ge 8 \max \left\{ k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2} \right\}$ and $t \ge A \log N$ for an absolute constant A > 0, then for every $i \in [m]$

|*median*_r { $C[r, h_r(i)] \cdot s_r(i)$ } - f_i $| \le \varepsilon f_k$

with high probability.

 $(f_i \text{ is the frequency of } i)$

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

and

$$\mathbf{E}_{s}[(C[r, h_{r}(i)] \cdot s_{r}(i) - f_{i})^{2}] = \sum_{j \neq i: \mathbf{h}_{r}(\mathbf{j}) = \mathbf{h}_{r}(\mathbf{i})} f_{j}^{2}$$

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

and

$$\mathbf{E}_{s}[(C[r,h_{r}(i)] \cdot s_{r}(i) - f_{i})^{2}] = \sum_{j \neq i: \mathbf{h}_{r}(\mathbf{j}) = \mathbf{h}_{r}(\mathbf{i})} f_{j}^{2}$$

How large can the variance be? Does it reduce by about a factor of *b*?

By basic estimate analysis for every $r \in [1:t]$

 $\mathbf{E}[C[r,h_r(i)]\cdot s_r(i)] = f_i$

and

$$\mathbf{E}_{s}[(C[r,h_{r}(i)] \cdot s_{r}(i) - f_{i})^{2}] = \sum_{j \neq i:\mathbf{h}_{r}(\mathbf{j}) = \mathbf{h}_{r}(\mathbf{i})} f_{j}^{2}$$

How large can the variance be? Does it reduce by about a factor of *b*?

Consider contribution of head and tail items separately:

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in H \in AD, j \neq i \\ \mathbf{h}_r(j) = \mathbf{h}_r(i)}} f_j^2 + \sum_{\substack{j \in TAIL, j \neq i \\ \mathbf{h}_r(j) = \mathbf{h}_r(i)}} f_j^2$$

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in H \in AD, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2 + \sum_{\substack{j \in TAIL, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2$$

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in HEAD, j \neq i \\ \mathbf{h_r(j)} = \mathbf{h_r(i)}}} f_j^2 + \sum_{\substack{j \in TAIL, j \neq i \\ \mathbf{h_r(j)} = \mathbf{h_r(i)}}} f_j^2$$

For each $r \in [1 : t]$ and each item $i \in [m]$ define three events:

NO-COLLISIONS_r(i) – i does not collide with any of the head items under hashing r

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in H \in AD, j \neq i \\ \mathbf{h}_r(j) = \mathbf{h}_r(\mathbf{i})}} f_j^2 + \sum_{\substack{j \in T A \mid L, j \neq i \\ \mathbf{h}_r(j) = \mathbf{h}_r(\mathbf{i})}} f_j^2$$

For each $r \in [1:t]$ and each item $i \in [m]$ define three events:

- NO-COLLISIONS_r(i) i does not collide with any of the head items under hashing r
- SMALL-VARIANCE_r(i) i does not collide with too many of tail items under hashing r

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in HEAD, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2 + \sum_{\substack{j \in TAIL, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2$$

For each $r \in [1:t]$ and each item $i \in [m]$ define three events:

- NO-COLLISIONS_r(i) i does not collide with any of the head items under hashing r
- SMALL-VARIANCE_r(i) i does not collide with too many of tail items under hashing r
- SMALL-DEVIATION_r(i) success event from basic analysis

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in HEAD, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2 + \sum_{\substack{j \in TAIL, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2$$

For each $r \in [1:t]$ and each item $i \in [m]$ define three events:

- NO-COLLISIONS_r(i) i does not collide with any of the head items under hashing r
- SMALL-VARIANCE_r(i) i does not collide with too many of tail items under hashing r
- SMALL-DEVIATION_r(i) success event from basic analysis

Show that all three events hold simultaneously with probability strictly bigger than 1/2 – so median gives good estimate

(No) collisions with head items

NO-COLLISIONS_r(i):=event that

$$\{j \in HEAD \setminus i : h_r(j) = h_r(i)\} = \emptyset,$$

i.e. that *i* collides with none of top *k* elements under h_r .

(No) collisions with head items

NO-COLLISIONS_r(i):=event that

$$\{j \in HEAD \setminus i : h_r(j) = h_r(i)\} = \emptyset,$$

i.e. that *i* collides with none of top *k* elements under h_r .

For every $j \neq i$ and every $r \in [1:t]$

 $\mathbf{Pr}[h_r(i) = h_r(j)] \le 1/b$

(No) collisions with head items

NO-COLLISIONS_{r(i)}:=event that

$$\{j \in HEAD \setminus i : h_r(j) = h_r(i)\} = \emptyset,$$

i.e. that *i* collides with none of top *k* elements under h_r .

For every $j \neq i$ and every $r \in [1:t]$

$$\mathbf{Pr}[h_r(i) = h_r(j)] \le 1/b$$

Suppose that $b \ge 8k$. Then by the union bound

$$\Pr[\text{NO-COLLISIONS}_r(i)] \ge 1 - k/b \\ \ge 1 - 1/8$$

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in H \in AD, j \neq i \\ \mathbf{h}_r(j) = \mathbf{h}_r(\mathbf{i})}} f_j^2 + \sum_{\substack{j \in T A \mid L, j \neq i \\ \mathbf{h}_r(j) = \mathbf{h}_r(\mathbf{i})}} f_j^2$$

For each $r \in [1:t]$ and each item $i \in [m]$ define three events:

- NO-COLLISIONS_r(i) i does not collide with any of the head items under hashing r
- SMALL-VARIANCE_r(i) i does not collide with too many of tail items under hashing r
- SMALL-DEVIATION_r(i) success event from basic analysis

Show that all three events hold simulaneously with probability strictly bigger than 1/2 – so median gives good estimate

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in HEAD, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2 + \sum_{\substack{j \in TAIL, j \neq i \\ \mathbf{h}_r(\mathbf{j}) = \mathbf{h}_r(\mathbf{i})}} f_j^2$$

For each $r \in [1:t]$ and each item $i \in [m]$ define three events:

- NO-COLLISIONS_r(i) i does not collide with any of the head items under hashing r
- SMALL-VARIANCE_r(i) i does not collide with too many of tail items under hashing r
- SMALL-DEVIATION_r(i) success event from basic analysis

Show that all three events hold simulaneously with probability strictly bigger than 1/2 – so median gives good estimate

Small variance from tail elements

SMALL-VARIANCE_r(i):=event that

$$\sum_{\substack{j \in \mathsf{TAIL}, j \neq i \\ h_r(j) = h_r(i)}} f_j^2 \le \frac{8}{b} \sum_{j \in \mathsf{TAIL}} f_j^2$$

Small variance from tail elements

SMALL-VARIANCE_r(i):=event that

$$\sum_{\substack{j \in \mathsf{TAIL}, j \neq i \\ h_r(j) = h_r(i)}} f_j^2 \le \frac{8}{b} \sum_{j \in \mathsf{TAIL}} f_j^2$$

For every $i, j \in [m], i \neq j$ and $r \in [1:t]$

 $\mathbf{Pr}_{h_r}[h_r(i) = h_r(j)] = 1/b$ (*b* is the number of buckets)

Small variance from tail elements

SMALL-VARIANCE_r(i):=event that

$$\sum_{\substack{j \in \mathsf{TAIL}, j \neq i \\ h_r(j) = h_r(i)}} f_j^2 \leq \frac{8}{b} \sum_{j \in \mathsf{TAIL}} f_j^2$$

For every $i, j \in [m], i \neq j$ and $r \in [1 : t]$

 $\mathbf{Pr}_{h_r}[h_r(i) = h_r(j)] = 1/b$ (*b* is the number of buckets)

So by linearity of expectation

$$\mathbf{E}\left[\sum_{\substack{j\in \mathsf{TA}|\mathcal{L},j\neq i\\h_r(j)=h_r(i)}} f_j^2\right] = \sum_{j\in \mathsf{TA}|\mathcal{L},j\neq i} f_j^2 \cdot \mathbf{Pr}_{h_r}[h_r(i) = h_r(j)]$$
$$\leq \frac{1}{b} \sum_{j\in \mathsf{TA}|\mathcal{L}} f_j^2$$

Markov's inequality

Theorem

For every non-negative random variable X with mean $\mu \ge 0$, and every $k \ge 1$ one has

$$\Pr[X \ge k \cdot \mu] \le 1/k$$



We proved that

$$\mathbf{E}\left[\sum_{\substack{j\in \mathsf{TA}|L, j\neq i\\h_r(j)=h_r(i)}} f_j^2\right] \leq \frac{1}{b} \sum_{j\in \mathsf{TA}|L} f_j^2$$

By Markov's inequality one has, for every *i* and every *r*,

 $\Pr[\text{SMALL-VARIANCE}_r(i)] \ge 1 - 1/8$

NO-COLLISIONS_r(i) and SMALL-VARIANCE_r(i): recap

Consider contribution of head and tail items separately:

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in H \in AD, j \neq i \\ h_r(j) = h_r(i)}} f_j^2 + \sum_{\substack{j \in \mathsf{TAIL}, j \neq i \\ h_r(j) = h_r(i)}} f_j^2$$

Conditioned on NO-COLLISIONS_r(i) and SMALL-VARIANCE_r(i)

NO-COLLISIONS_r(i) and SMALL-VARIANCE_r(i): recap

Consider contribution of head and tail items separately:

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in H \in AD, j \neq i \\ h_r(j) = h_r(i)}} f_j^2 + \sum_{\substack{j \in \mathsf{TAIL}, j \neq i \\ h_r(j) = h_r(i)}} f_j^2$$

Conditioned on NO-COLLISIONS_r(i) and SMALL-VARIANCE_r(i)

first term is zero

NO-COLLISIONS_r(i) and SMALL-VARIANCE_r(i): recap

Consider contribution of head and tail items separately:

$$\sum_{\substack{j \neq i: h_r(j) = h_r(i)}} f_j^2 = \sum_{\substack{j \in H \in AD, j \neq i \\ h_r(j) = h_r(i)}} f_j^2 + \sum_{\substack{j \in \mathsf{TAIL}, j \neq i \\ h_r(j) = h_r(i)}} f_j^2$$

Conditioned on NO-COLLISIONS_r(i) and SMALL-VARIANCE_r(i)

- first term is zero
- second term is at most

$$\frac{8}{b}\sum_{j\in TAIL}f_j^2$$

Small deviation event

SMALL-DEVIATION $_r(i)$ =event that

$$(C[r,h_r(i)] \cdot s_r(i) - f_i)^2 \leq 8 \operatorname{Var}(C[r,h_r(i)] \cdot s_r(i)).$$

Small deviation event

SMALL-DEVIATION $_r(i)$ =event that

$$(C[r,h_r(i)] \cdot s_r(i) - f_i)^2 \leq 8 \operatorname{Var}(C[r,h_r(i)] \cdot s_r(i)).$$

By Chebyshev's inequality one has, for every i and every r,

 $\Pr[\text{SMALL-DEVIATION}_r(i)] \ge 1 - 1/8$

$$\Pr[\text{SMALL-VARIANCE}_r(i)] \ge 1 - 1/8$$

$$Pr[NO-COLLISIONS_r(i)] \ge 1 - 1/8$$

$$\Pr[\text{SMALL-DEVIATION}_r(i)] \ge 1 - 1/8$$

So by the union bound

 $\Pr[\text{SMALL-VARIANCE}_{r}(i) \text{ and } \text{NO-COLLISIONS}_{r}(i)$ and $\text{SMALL-DEVIATION}_{r}(i)] \ge 5/8$.

$$\gamma := \sqrt{\frac{1}{b} \sum_{j \in TAIL} f_j^2}$$

Lemma If $b \ge 8k$, then for every *i*, every $r \in [1:t]$,

 $\mathbf{Pr}[|C[r,h_r(i)] \cdot s_r(i) - f_i| \le 8\gamma] \ge 5/8$

$$\gamma := \sqrt{\frac{1}{b} \sum_{j \in \mathsf{TAIL}} f_j^2}$$

Lemma If $b \ge 8k$ and $t \ge A \log N$ for an absolute constant A > 0, then for every *i*, with probability $\ge 1 - 1/N^4$

$$\left| median_r \left\{ C[r, h_r(i)] \cdot s_r(i) \right\} - f_i \right| \le 8\gamma$$

at the end of the stream.

Proof. Chernoff bounds.

$$\gamma := \sqrt{\frac{1}{b} \sum_{j \in TAIL} f_j^2}$$

Lemma If $b \ge 8k$ and $t \ge A \log N$ for an absolute constant A > 0, then with probability $\ge 1 - 1/N^3$ for every $i \in [m]$

$$\left| median_r \left\{ C[r, h_r(i)] \cdot s_r(i) \right\} - f_i(p) \right| \le 8\gamma$$

at the end of the stream.

Proof. Chernoff bounds.

$$\gamma := \sqrt{\frac{1}{b} \sum_{j \in \mathsf{TA}/\mathsf{I}} f_j^2}$$

Lemma If $b \ge 8 \max \left\{ k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2} \right\}$ and $t \ge A \log N$ for an absolute constant A > 0, then with probability $\ge 1 - 1/N^3$ for every $i \in [m]$ $\left| median_r \left\{ C[r, h_r(i)] \cdot s_r(i) \right\} - f_i(p) \right| \le \varepsilon f_k$

at the end of the stream.

$$\gamma := \sqrt{\frac{1}{b} \sum_{j \in \mathsf{TA}/\mathsf{I}} f_j^2}$$

Lemma If $b \ge 8 \max \left\{ k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2} \right\}$ and $t \ge A \log N$ for an absolute constant A > 0, then with probability $\ge 1 - 1/N^3$ for every $i \in [m]$ $\left| median_r \left\{ C[r, h_r(i)] \cdot s_r(i) \right\} - f_i(p) \right| \le \varepsilon f_k$

at the end of the stream.

Proof.

Substitute value of *b* into definition of γ :

$$\gamma = \sqrt{\frac{1}{b} \sum_{j \in TAIL} f_j^2} \le \varepsilon f_k / 8$$