Sketching for Data Streams

Michael Kapralov

EPFL

August 30, 2023

Streaming model (Alon, Matias, Szegedy'96)

Observe a (very long) stream of data, e.g. IP packets, tweets, search queries....

Task: maintain (approximate) statistics of the stream

Streaming model

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

Typically, assume N is known

- Small (sublinear) storage: typically N^α, α < 1 or log^{O(1)} N
 Units of storage: bits, words or 'data items' (e.g., points, nodes/edges)
- Fast processing time per element
- Mostly randomized algorithms Randomness often necessary

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: $i_1, i_2, ..., i_N$

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!



3 4 6 3

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)

Much better than storing all items!



3 4 6 3 2

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!



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, \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)



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)



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			



		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cior	n			
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



		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cio	n			
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



		de	esti	Lnat	cio	n			
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	



		de	esti	Lnat	cio	n			
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


		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cio	n			
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					
DATA						



		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cio	n			
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	DATA						



		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cio	n			
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	DATA						



		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cio	n			
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	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	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



		de	esti	Lnat	cio	n			
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						
DA	DATA						



		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cio	n			
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	DATA						



		de	esti	Lnat	cio	n		
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



		de	esti	Lnat	cio	n			
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	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	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



		de	esti	Lnat	cio	n			
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	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	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>	$O(\log N)$	

Streaming model

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

Streaming model

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

Are constants small?

Streaming model

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

Are constants small?

HyperLogLog: estimate Shakespeare's vocabulary using 128 bits of memory


Streaming model

Widely used in practice for scalable data analytics



most frequent searches on google.com over a time period



most frequent tweets

Heavy hitters problem

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)

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*)

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 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



1 4



146



1 4 6 1



1 4 6 1 2










































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

INITIALIZE UPDATE(C, i) $C \leftarrow 0$ $C \leftarrow C + s(i)$

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

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

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

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

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

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

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

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

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

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

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'

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

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

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

The mean is correct: our estimator is unbiased!

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

The mean is correct: our estimator is unbiased!

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

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

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

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

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

since

►
$$\mathbf{E}[s(j)s(j')] = \mathbf{E}[s(j)]\mathbf{E}[s(j')] = 0$$
 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

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

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

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

- 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

COUNTSKETCH algorithm (Charikar, Chen, Farach-Colton'02) Main ideas:

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

APPROXPOINTQUERY and COUNTSKETCH

COUNTSKETCH algorithm (Charikar, Chen, Farach-Colton'02) Main ideas:

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

APPROXPOINTQUERY and COUNTSKETCH

COUNTSKETCH algorithm (Charikar, Chen, Farach-Colton'02) 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,

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}





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

1







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

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 $r \in [1:t]$ $C[r, h_r(i)] \leftarrow C[r, h_r(i)] + s_r(i)$ 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 { $C[r, h_r(i)] \cdot s_r(i)$ }

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
$$B \ge 8 \max\left\{k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2}\right\}$$
 and $t = O(\log N)$, then for every $i \in [m]$
 $|ESTIMATE(C, i) - f_i| \le \varepsilon f_k$

at every point in the stream whp.

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
$$B \ge 8 \max\left\{k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2}\right\}$$
 and $t = O(\log N)$, then for every $i \in [m]$
 $|ESTIMATE(C, i) - f_i| \le \varepsilon f_k$

at every point in the stream whp.

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

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
$$B \ge 8 \max\left\{k, \frac{32\sum_{j \in TAIL} f_j^2}{(\varepsilon f_k)^2}\right\}$$
 and $t = O(\log N)$, then for every $i \in [m]$
 $|ESTIMATE(C, i) - f_i| \le \varepsilon f_k$

at every point in the stream whp.

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

How large is B?



Note: if $B \ge k$, can detect elements with counts above $O\left(\sqrt{\frac{1}{B} \cdot \sum_{j \in TAIL} f_j^2}\right)$

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

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

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

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

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

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

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

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

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

 $|\mathsf{ESTIMATE}(C, i) - f_i| \le \varepsilon f_k$

with high probability.

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

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

Variance of estimate for *i* from *r*-th row:

$$\sum_{j\neq i:\mathbf{h_r(j)}=\mathbf{h_r(i)}} f_j^2$$

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

Variance of estimate for *i* from *r*-th row:

$$\sum_{j \neq i: \mathbf{h}_{\mathbf{r}}(\mathbf{j}) = \mathbf{h}_{\mathbf{r}}(\mathbf{i})} f_j^2$$

Show that

$$\sum_{j \neq i: \mathbf{h_r(j)} = \mathbf{h_r(i)}} f_j^2 = O(1/B) \sum_{j \in \textit{TAIL}, j \neq i} f_j^2$$

with high constant probability.

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

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

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

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

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

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

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

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

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

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

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

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

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

Small variance from tail elements

SMALL-VARIANCE_r(i):=event that

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

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)

Small variance from tail elements

SMALL-VARIANCE_r(i):=event that

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

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

$$\mathbf{E}\left[\sum_{\substack{j\in \mathsf{TAIL}, j\neq i\\h_r(j)=h_r(i)}} f_j^2\right] = \sum_{j\in \mathsf{TAIL}, j\neq i} f_j^2 \cdot \mathbf{Pr}_{h_r}[h_r(i) = h_r(j)]$$
Small variance from tail elements

SMALL-VARIANCE_r(i):=event that

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

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

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

We proved that

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

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:

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

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

$$\frac{8}{B}\sum_{j\in TA/L}f_j^2$$

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 Markov'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 $\text{SMALL-DEVIATION}_{r}(i)] \ge 5/8$. For every $p \in [1 : N]$ let $f_i(p) :=$ frequency of *i* 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 with probability $\ge 1 - 1/N^3$ for every $i \in [m]$ $|\text{ESTIMATE}(C, i) - f_i(p)| \le \varepsilon f_k$

at the end of the stream.

Remarks, related results, open problems

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



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



1 4

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



1 4 6

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



1 4 6 1

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



1 4 6 1 2

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



1 4 6 1 2 10

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



1 4 6 1 2 10 1

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



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



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



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



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



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



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



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
$$r \in [1:t]$$

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



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
$$r \in [1:t]$$

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


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



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



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

Sketching: take (randomized) linear measurements of the input



Easy to maintain sketch in dynamic streams (insertions and deletions)



Let S be a COUNTSKETCH matrix with $O(k \log n)$ rows

Lemma

For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(Sx)$, then whp

$$\left\|x-\widehat{x}\right\|_{\infty} \leq \frac{1}{\sqrt{k}} \|x_{TA/L}\|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Let S be a COUNTSKETCH matrix with $O(k \log n)$ rows

Lemma

For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(Sx)$, then whp

$$\left\|x-\widehat{x}\right\|_{\infty} \leq \frac{1}{\sqrt{k}} \|x_{TAIL}\|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Let S be a COUNTSKETCH matrix with $O(k \log n)$ rows

Lemma

For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(Sx)$, then whp

$$\|x-\widehat{x}\|_{\infty} \leq \frac{1}{\sqrt{k}} \|x_{TA/L}\|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Observation 1: # of measurements is optimal for ℓ_∞/ℓ_2 guarantee above

(capacity of the Gaussian channel, i.e. Shannon-Hartley theorem, – see Do Ba, Indyk, Price, Woodruff'10)

Let S be a COUNTSKETCH matrix with $O(k \log n)$ rows

Lemma

For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(Sx)$, then whp

$$\|x-\widehat{x}\|_{\infty} \leq \frac{1}{\sqrt{k}} \|x_{TA/L}\|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Observation 1: # of measurements is optimal for ℓ_∞/ℓ_2 guarantee above

(capacity of the Gaussian channel, i.e. Shannon-Hartley theorem, – see Do Ba, Indyk, Price, Woodruff'10)

Observation 2: ℓ_2/ℓ_2 sparse recovery guarantee follows:

$$||x - \hat{x}||_2 = O(1) \cdot ||x_{TAIL}||_2.$$

Let S be a COUNTSKETCH matrix with $O(k \log n)$ rows

Lemma

For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(Sx)$, then whp

$$\left\| x - \widehat{x} \right\|_{\infty} \leq \frac{1}{\sqrt{k}} \| x_{TAIL} \|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Observation 1: # of measurements is optimal for ℓ_∞/ℓ_2 guarantee above

(capacity of the Gaussian channel, i.e. Shannon-Hartley theorem, – see Do Ba, Indyk, Price, Woodruff'10)

Let S be a COUNTSKETCH matrix with $O(k \log n)$ rows

Lemma

For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(Sx)$, then whp

$$\left\| x - \widehat{x} \right\|_{\infty} \leq \frac{1}{\sqrt{k}} \| x_{TAIL} \|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Observation 1: # of measurements is optimal for ℓ_∞/ℓ_2 guarantee above

(capacity of the Gaussian channel, i.e. Shannon-Hartley theorem, – see Do Ba, Indyk, Price, Woodruff'10)

Observation 2: ℓ_2/ℓ_2 sparse recovery guarantee follows:

$$||x - \hat{x}||_2 = O(1) \cdot \min_{k-sparse \ x'} ||x - x'||_2.$$

Let S be a COUNTSKETCH matrix with $O(\frac{1}{s^2} k \log n)$ rows

Lemma

For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(Sx)$, then whp

$$\left\| x - \widehat{x} \right\|_{\infty} \leq \frac{1}{\sqrt{k}} \| x_{TAIL} \|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Observation 1: # of measurements is optimal for ℓ_∞/ℓ_2 guarantee above

(capacity of the Gaussian channel, i.e. Shannon-Hartley theorem, – see Do Ba, Indyk, Price, Woodruff'10)

Observation 2: ℓ_2/ℓ_2 sparse recovery guarantee follows:

$$\|x-\widehat{x}\|_{2} = (1+\varepsilon) \cdot \min_{k-sparse \ x'} \|x-x'\|_{2}.$$

Let **S** be a COUNTSKETCH matrix with $O(k \log n)$ rows

Lemma For every $x \in \mathbb{R}^n$ if $\hat{x} = \text{EST}(\frac{S}{x})$, then whp

$$\|x-\widehat{x}\|_{\infty} \leq \frac{1}{\sqrt{k}} \|x_{TAIL}\|_2.$$

 $(x_{TAIL} - x \text{ with largest } k \text{ elements zeroed out})$

Observation 3: if $||x||_0 \le k$, then $x_{TAIL} = 0$ and EST(Sx) = x whp

Exact sparse recovery: a *k*-sparse vector can be recovered from O(k) linear measurements

Input: *r* parties hold vectors $x_1, \ldots, x_r \in \mathbb{R}^n$

each party sends $O(B \log n)$ bits to coordinator

(assume shared randomness)

Input: *r* parties hold vectors $x_1, \ldots, x_r \in \mathbb{R}^n$

each party sends $O(B \log n)$ bits to coordinator (assume shared randomness)

(assume shared randomness)

Output: find largest entries in $A = \sum_{i=1}^{r} x_i x_i^T$

more precisely, output approximation \widehat{A}

$$\left|\widehat{A_{ij}}-A_{ij}\right|=\frac{O(1)}{\sqrt{B}}\|A\|_F.$$

Input: *r* parties hold vectors $x_1, ..., x_r \in \mathbb{R}^n$

each party sends $O(B \log n)$ bits to coordinator (assume shared randomness)

Output: find largest entries in $A = \sum_{i=1}^{r} x_i x_i^T$

more precisely, output approximation \widehat{A}



Input: *r* parties hold vectors $x_1, \ldots, x_r \in \mathbb{R}^n$

each party sends $O(B \log n)$ bits to coordinator (assume shared randomness)

Output: find largest entries in $A = \sum_{i=1}^{r} x_i x_i^T$

more precisely, output approximation \widehat{A}



Every party *i* sends COUNTSKETCH($x_i x_j^T$) into O(B) buckets (slow)

$$\boldsymbol{A} = \sum_{i=1}^{r} x_i x_i^{T} \in \mathbb{R}^{n \times n}.$$

$$\boldsymbol{A} = \sum_{i=1}^{r} x_i x_i^{T} \in \mathbb{R}^{n \times n}.$$

Hash function

$$h: [n] \times [n] \rightarrow [B]$$

and random signs

$$s:[n]\times[n]\to\{-1,+1\}.$$

$$\boldsymbol{A} = \sum_{i=1}^{r} \boldsymbol{X}_{i} \boldsymbol{X}_{i}^{T} \in \mathbb{R}^{n \times n}.$$

Hash function

$$h: [n] \times [n] \rightarrow [B]$$

and random signs

$$s:[n]\times[n]\to\{-1,+1\}.$$

$$(Sx)_b = \sum_{i,j\in[n]:h(i,j)=b} s(i,j) \cdot x_i x_j.$$

$$\boldsymbol{A} = \sum_{i=1}^{r} x_i x_i^{T} \in \mathbb{R}^{n \times n}.$$

Hash function

$$h: [n] \times [n] \rightarrow [B]$$

and random signs

$$s:[n]\times[n]\to\{-1,+1\}.$$

$$(Sx)_b = \sum_{i,j\in[n]:h(i,j)=b} s(i,j)\cdot x_i x_j.$$

COUNTSKETCH($x_i x_i^T$) takes n^2 time to compute...

$$\boldsymbol{A} = \sum_{i=1}^{r} x_i x_i^{T} \in \mathbb{R}^{n \times n}.$$

Hash function

$$h:[n]\times [n]\to [B]$$

and random signs

$$s:[n]\times[n]\to\{-1,+1\}.$$

$$(Sx)_b = \sum_{i,j\in[n]:h(i,j)=b} s(i,j)\cdot x_i x_j.$$

COUNTSKETCH $(x_i x_i^T)$ takes n^2 time to compute...

Make hash functions 'separable'?

Input: *r* parties hold vectors $x_1, ..., x_r \in \mathbb{R}^n$

each party sends $O(B \log n)$ bits to coordinator (assume shared randomness)

Output: find largest entries in $A = \sum_{i=1}^{r} x_i x_i^T$

more precisely, output approximation \widehat{A}



Input: *r* parties hold vectors $x_1, \ldots, x_r \in \mathbb{R}^n$

each party sends $O(B \log n)$ bits to coordinator (assume shared randomness)

Output: find largest entries in $A = \sum_{i=1}^{r} x_i x_i^T$

more precisely, output approximation \widehat{A}



Every party *i* sends SOMESKETCH(*x_i*) into *B* buckets? (fast)



Take two independent instances of COUNTSKETCH: hash functions

 $h_1, h_2: [n] \rightarrow [B],$

random signs

$$s_1, s_2 : [n] \to \{-1, +1\}$$

Tensor COUNTSKETCH₁ and COUNTSKETCH₂!



Define tensoring of COUNTSKETCH₁ and COUNTSKETCH₂:

 $h(i,j) = h_1(i) + h_2(j) \pmod{B}.$



Define tensoring of COUNTSKETCH1 and COUNTSKETCH2:

$$h(i,j) = h_1(i) + h_2(j) \pmod{B}.$$

and $s(i,j) = s_1(i) \cdot s_2(j)$.



Define tensoring of COUNTSKETCH₁ and COUNTSKETCH₂:

 $h(i,j) = h_1(i) + h_2(j) \pmod{B}.$

and $s(i,j) = s_1(i) \cdot s_2(j)$.

$$(Sx)_{b} = \sum_{i,j \in [n]: h(i,j)=b} S(i,j) \cdot x_{i}x_{j}$$

=
$$\sum_{i,j \in [n]: h_{1}(i)+h_{2}(j)=b} S_{1}(i) \cdot S_{2}(j) \cdot x_{i}x_{j}$$

$$(Sx)_b = \sum_{i,j \in [n]: h_1(i) + h_2(j) = b} S_1(i) \cdot S_2(j) \cdot x_i x_j.$$

Can find Sx from COUNTSKETCH₁(x) and COUNTSKETCH₂(x) fast! (exercise)

Stronger analysis of COUNTSKETCH

The bound

$$\left\| x - \widehat{x} \right\|_{\infty} \le \frac{1}{\sqrt{k}} \| x_{TAIL} \|_2$$

is optimal for sketches with $O(k \log n)$ rows, for worst case x

Stronger analysis of COUNTSKETCH

The bound

$$\left\| \boldsymbol{x} - \hat{\boldsymbol{x}} \right\|_{\infty} \le \frac{1}{\sqrt{k}} \|\boldsymbol{x}_{\mathsf{TAIL}}\|_2$$

is optimal for sketches with $O(k \log n)$ rows, for worst case x

If *x* is drawn from a distribution (e.g., power law, Zipfian), one can do better by about log *n* factor: Minton-Price'14

Stronger analysis of COUNTSKETCH

The bound

$$\|x - \hat{x}\|_{\infty} \le \frac{1}{\sqrt{k}} \|x_{TA}\|_{2}$$

is optimal for sketches with $O(k \log n)$ rows, for worst case x

If *x* is drawn from a distribution (e.g., power law, Zipfian), one can do better by about log *n* factor: Minton-Price'14

Minton-Price'14 assumes uniformly random hashing. A very recent improvement:

Pseudorandom Hashing for Space-bounded Computation with Applications in Streaming

Praneeth Kacham^{*} Rasmus Pagh[†] Mikkel Thorup[‡] David P. Woodruff[§]

Abstract

We revisit Nisan's classical pseudorandom generator (PRG) for space-bounded computation (STOC 1990) and its applications in streaming algorithms. We describe a new generator, Hash-PRG, that can be thought of as a symmetric version of Nisan's generator over larger alphabets.

Good constants are achieved by ℓ_1 -minimization and related (non-sublinear) methods. Get best of both worlds?

Good constants are achieved by ℓ_1 -minimization and related (non-sublinear) methods. Get best of both worlds?

Ex: in LDPC codes, the Tanner graphs needs to be irregular to achieve capacity – similar effects here?

Good constants are achieved by ℓ_1 -minimization and related (non-sublinear) methods. Get best of both worlds?

Ex: in LDPC codes, the Tanner graphs needs to be irregular to achieve capacity – similar effects here?

Information Theoretic Limits of Cardinality Estimation: Fisher Meets Shannon*

Seth Pettie pettie@umich.edu University of Michigan Computer Science and Engineering Ann Arbor, MI, USA

ABSTRACT

Estimating the *cardinality* (number of distinct elements) of a large multiset is a classic problem in streaming and sketching, dating back to Flajolet and Martin's classic Probabilistic Counting (PCSA) algorithm from 1983.

In this paper we study the intrinsic tradeoff between the space complexity of the sketch and its estimation error in the RANDOM ONACLE model. We define a new measure of efficiency for cardinality estimators called the Fisher-Shannon (Fish) number \mathcal{H}/I . It captures the tension between the limiting Shannon entropy (\mathcal{H}) of the sketch and its normalized Fisher information (I), which characterizes the variance of a statistically efficient, asymptotically unbiased estimator. Dingyu Wang wangdy@umich.edu University of Michigan Computer Science and Engineering Ann Arbor, MI, USA

KEYWORDS

cardinality estimation, streaming algorithm

ACM Reference Format:

Seth Pettie and Dingyu Wang. 2021. Information Theoretic Limits of Carionality Estimation: Fisher Meets Shannon. In Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing (STOC '21), June 21–25, 2021, Virtual, Italy, ACM, New York, NY, USA, 14 pages. https: //doi.org/10.1145/3408325.3451032

1 INTRODUCTION

Cardinality Estimation (aka *Distinct Elements* or *F*₀-*estimation*) is a fundamental problem in streaming/sketching, with widespread industrial deployments in databases, networking, and sensing.

Good constants are achieved by ℓ_1 -minimization and related (non-sublinear) methods. Get best of both worlds?

Ex: in LDPC codes, the Tanner graphs needs to be irregular to achieve capacity – similar effects here?

Information Theoretic Limits of Cardinality Estimation: Fisher Meets Shannon*

Seth Pettie pettie@umich.edu University of Michigan Computer Science and Engineering Ann Arbor, MI, USA

Abstract

Dingyu Wang wangdy@umich.edu University of Michigan Computer Science and Engineering Ann Arbor, MI, USA

ABSTRACT

Estimating the ca multiset is a class back to Flajolet ar algorithm from 1! In this paper w complexity of the oracLE model. W nality estimators captures the tens of the sketch an characterizes the unbiased estimato

KEYWORDS

Peeling Close to the Orientability Threshold – Spatial Coupling in Hashing-Based Data Structures

Stefan Walzer*

1 Introduction

Good constants are achieved by ℓ_1 -minimization and related (non-sublinear) methods. Get best of both worlds?

Ex: in LDPC codes, the Tanner graphs needs to be irregular to achieve capacity – similar effects here?

Information Theoretic Limits of Cardinality Estimation: Fisher Meets Shannon*

Seth Pettie pettie@umich.edu University of Michigan Computer Science and Engineering Ann Arbor, MI, USA

Abstract

Dingyu Wang wangdy@umich.edu University of Michigan Computer Science and Engineering Ann Arbor, MI, USA

ABSTRACT

Estimating the ca multiset is a class back to Flajolet ar algorithm from 1! In this paper w complexity of the oracLE model. W nality estimators captures the tens of the sketch an characterizes the unbiased estimato

KEYWORDS

Peeling Close to the Orientability Threshold – Spatial Coupling in Hashing-Based Data Structures

Stefan Walzer*

Simple Set Sketching

Learning-augmented sketching: learn the hash function *h* in COUNTSKETCH (and more!) from data

Adversarially robust sketching: what if x is chosen by an adversary with (partial) knowledge of the data structure?
Learning-augmented sketching: learn the hash function *h* in COUNTSKETCH (and more!) from data

Adversarially robust sketching: what if x is chosen by an adversary with (partial) knowledge of the data structure?

ON THE ROBUSTNESS OF COUNTSKETCH TO ADAPTIVE INPUTS

Edith Cohen* Xin Lyu[†]

Jelani Nelson[‡]

Tamás Sarlós[§]

Moshe Shechner[¶] Uri Stemmer[∥]

March 1, 2022

ABSTRACT

CountSketch is a popular dimensionality reduction technique that maps vectors to a lower dimension using randomized linear measurements. The sketch supports recovering ℓ_2 -heavy hitters of a vector (entries with $\psi[i]^2 \ge \frac{1}{k} \|v\|_2^2$). We study the robustness of the sketch in *adaptive* settings where input vectors may depend on the output from prior inputs. Adaptive settings arise in processes with feedback or with adversarial attacks. We show that the classic estimator is not robust, and can be attacked with a number of queries of the order of the sketch size. We propose a robust estimator (for a slightly modified sketch) that allows for quadratic number of queries in the sketch size, which is an improvement factor of \sqrt{k} (for k heavy hitters) over prior work.

1 Introduction

Take (randomized) linear measurements of the input



Distribution of the sketching matrix?



Bernoulli(±1) or Gaussian linear measurements on random subsets of the universe

The nonzeros are specified by the hash function $h: [n] \rightarrow [B]$



Bernoulli(±1) or Gaussian linear measurements on random subsets of the universe

The nonzeros are specified by the hash function $h: [n] \rightarrow [B]$



Bernoulli(±1) or Gaussian linear measurements on random subsets of the universe

The nonzeros are specified by the hash function $h: [n] \rightarrow [B]$



Bernoulli(±1) or Gaussian linear measurements on random subsets of the universe

The nonzeros are specified by the hash function $h: [n] \rightarrow [B]$



Bernoulli(±1) or Gaussian linear measurements on random subsets of the universe

The nonzeros are specified by the hash function $h: [n] \rightarrow [B]$



Bernoulli(±1) or Gaussian linear measurements on random subsets of the universe

The nonzeros are specified by the hash function $h: [n] \rightarrow [B]$

Random restrictions (hashing)

What can we learn from *Sx*, where *S* is just random restrictions?



Can learn $||x||_0$, i.e. number of nonzeros in x

Sketching matrix S = a row of i.i.d. Gaussians of unit variance

Sketching matrix S = a row of i.i.d. Gaussians of unit variance

Measures ℓ_2^2 norm of *x* in expectation:

$$\mathbb{E}\left[\left\|\frac{\mathbf{S}x}{\mathbf{S}}\right\|_{2}^{2}\right] = \|\mathbf{X}\|_{2}^{2}$$

Sketching matrix S = m rows of i.i.d. Gaussians of unit variance 1/m

Sketching matrix S = m rows of i.i.d. Gaussians of unit variance 1/m

Measures ℓ_2^2 norm of x with high probability: $\mathbb{P}\left[\|Sx\|_2^2 \neq \|x\|_2^2\right] = 1 - \exp(-\Omega(\varepsilon^2 m))$

Sketching matrix S = m rows of i.i.d. Gaussians of unit variance 1/m

Measures ℓ_2^2 norm of x with high probability: $\mathbb{P}\left[\|Sx\|_2^2 \neq \|x\|_2^2\right] = 1 - \exp(-\Omega(\varepsilon^2 m))$

Downside: Sx takes $m \cdot n$ time to compute

Subsampled randomized Hadamard transform

Subsampled randomized Hadamard transform Sketching matrix $S = P \cdot H \cdot D$

Subsampled randomized Hadamard transform

Sketching matrix $S = P \cdot H \cdot D$

D=diagonal random sign matrix, *H*=Hadamard transform, *P*=random sampling matrix

Subsampled randomized Hadamard transform

Sketching matrix $S = P \cdot H \cdot D$

D=diagonal random sign matrix, *H*=Hadamard transform, *P*=random sampling matrix

Sx can be computed in $O(m + n \log n)$ time

Frequency moments

The *p*-th frequency moment

$$F_p = \sum_{i \in [n]} f_i^p$$

Frequency moments

The *p*-th frequency moment

$$F_{p} = \sum_{i \in [n]} f_{i}^{p}$$

Theorem Can approximate F_p for all $p \in [0,2]$ in polylogarithmic space, but need $\Omega(n^{1-2/p})$ space for all p > 2

Frequency moments

The *p*-th frequency moment

$$F_{p} = \sum_{i \in [n]} f_{i}^{p}$$

Theorem Can approximate F_p for all $p \in [0,2]$ in polylogarithmic space, but need $\Omega(n^{1-2/p})$ space for all p > 2

Bar-Yossef et al. Information complexity approach to data stream lower bounds