

Algorithms for data streams

Irene Finocchi

finocchi@di.uniroma1.it

http://www.dsi.uniroma1.it/~finocchi/

May 16, 2012

Outline

Introduction

- Stream sources
- Data stream model
- 2 Data stream challenges
 - Puzzle 1: missing number
 - Puzzle 2: fishing
 - Puzzle 3: pointer & chaser
 - Lessons
- 3 Sampling
 - Reservoir sampling
 - Heavy hitters
 Sticky sampling
- 4 Sketches

- Distinct items
 - Probabilistic counting
- Frequency moments
 AMS sketches
- 5 Mining graphs
 - Triangle counting
 - Semi-streaming model
 - Maximum matching
- 6 Lower bounds
 - Communication complexity
- 🕡 Summing up
 - More streaming algorithms
 - What's next?

Goals:

- give a flavor of the theoretical results and techniques of data stream algorithmics
- only a representative sample of each topic: many other problems, algorithms, and techniques not covered in these lectures (non-exhaustive overview at the end of the talk)

Goals:

- give a flavor of the theoretical results and techniques of data stream algorithmics
- only a representative sample of each topic: many other problems, algorithms, and techniques not covered in these lectures (non-exhaustive overview at the end of the talk)

Math contents: some probability ahead (e.g., Chernoff bounds). Will introduce basic tools along the way.

Goals:

- give a flavor of the theoretical results and techniques of data stream algorithmics
- only a representative sample of each topic: many other problems, algorithms, and techniques not covered in these lectures (non-exhaustive overview at the end of the talk)

Math contents: some probability ahead (e.g., Chernoff bounds). Will introduce basic tools along the way.

RequestIf you get bored, ask questions

Goals:

- give a flavor of the theoretical results and techniques of data stream algorithmics
- only a representative sample of each topic: many other problems, algorithms, and techniques not covered in these lectures (non-exhaustive overview at the end of the talk)

Math contents: some probability ahead (e.g., Chernoff bounds). Will introduce basic tools along the way.

Request

- If you get bored, ask questions
- If you get lost, ask questions

Goals:

- give a flavor of the theoretical results and techniques of data stream algorithmics
- only a representative sample of each topic: many other problems, algorithms, and techniques not covered in these lectures (non-exhaustive overview at the end of the talk)

Math contents: some probability ahead (e.g., Chernoff bounds). Will introduce basic tools along the way.

Request

- If you get bored, ask questions
- If you get lost, ask questions
- If you'd like to ask questions, ask questions

Massive data

Data is growing faster than our ability to store and index it:

- networking: phone call networks, Internet, social networks
- scientific data: astronomical data, genome sequences, GIS geo-spatial data
- economic transactions: credit cards, online auctions



. . .

Network management

Monitoring flow of IP packets through the routers (Internet traffic):

- how many IP addresses used a given link in the last month?
- which are the top 100 IP addresses in terms of traffic?
- which destinations use most bandwidth?
- what's the average duration of an IP session?
- which hosts have similar usage patterns (clusters)?
- does traffic distribution change in different periods of time?



Network management

Monitoring flow of IP packets through the routers (Internet traffic):

- how many IP addresses used a given link in the last month?
- which are the top 100 IP addresses in terms of traffic?
- which destinations use most bandwidth?
- what's the average duration of an IP session?
- which hosts have similar usage patterns (clusters)?
- does traffic distribution change in different periods of time?

Up to 1 Billion packets per hour per router Many hundreds of routers per ISP



Many terabytes of data per hour!



Sensor data

Sensors with GPS unit deployed in the ocean:

- Each sensor reports surface height (4-byte real number) every tenth of second
- Base station receives 3.5 MB per day per sensor



Sensor data

Sensors with GPS unit deployed in the ocean:

- Each sensor reports surface height (4-byte real number) every tenth of second
- Base station receives 3.5 MB per day per sensor



What about a million sensors?

3.5 TB of data per day, coming at a high rate

A million sensors isn't very many: roughly one sensor per 150 square miles of ocean...

More streams...

Image data

- satellites send down to earth many TBs of images per day
- surveillance cameras produce roughly one image per second: London has about six millions such cameras



More streams...

Image data

- satellites send down to earth many TBs of images per day
- surveillance cameras produce roughly one image per second: London has about six millions such cameras





Web traffic

Google receives several hundreds million search queries per day

More streams...

Image data

- satellites send down to earth many TBs of images per day
- surveillance cameras produce roughly one image per second: London has about six millions such cameras





Web traffic

Google receives several hundreds million search queries per day

Economic trend analysis

 in online auction systems, users continuously submit bids for items and items for auction



Issues in data stream processing

Some features common to all these applications:

- huge volumes of data (terabytes, even petabytes)
- records arrive at a rapid rate
- need to monitor data continuously to support exploratory analyses and to detect correlations, patterns, rare events, fraud, intrusion, unusual activities

Issues in data stream processing

Some features common to all these applications:

- huge volumes of data (terabytes, even petabytes)
- records arrive at a rapid rate
- need to monitor data continuously to support exploratory analyses and to detect correlations, patterns, rare events, fraud, intrusion, unusual activities

Many problems about streaming data would be easy to solve if we had enough memory, but require new techniques for realistic data rates and sizes

What can be computed without even storing the input?

Basic data stream model

Data stream = sequence $\sigma = \langle a_1, a_2, ..., a_m \rangle$ of tokens drawn from universe $[n] = \{1, 2, ..., n\}$

Basic data stream model

Data stream = sequence $\sigma = \langle a_1, a_2, ..., a_m \rangle$ of tokens drawn from universe $[n] = \{1, 2, ..., n\}$

Input parameters: *m* and *n*

1 Stream σ is massively long. Stream length *m* is:

- typically unknown
- possibly infinite

Basic data stream model

Data stream = sequence $\sigma = \langle a_1, a_2, ..., a_m \rangle$ of tokens drawn from universe $[n] = \{1, 2, ..., n\}$

Input parameters: *m* and *n*

1 Stream σ is massively long. Stream length *m* is:

- typically unknown
- possibly infinite
- Universe size n is also typically very large (e.g., IP addresses, URLs, item prices)

Minimize space, passes, and processing time upon token arrivals

Minimize space, passes, and processing time upon token arrivals

Use a sublinear amount of space s:

 $s = o(\min\{n, m\})$

where s = bits of random-access working memory

Minimize space, passes, and processing time upon token arrivals

• Use a sublinear amount of space *s*:

 $s = o(\min\{n, m\})$

where s = bits of random-access working memory

Make p passes over the data, for some small integer p (no random access to tokens)

Minimize space, passes, and processing time upon token arrivals

• Use a sublinear amount of space *s*:

 $s = o(\min\{n, m\})$

where s = bits of random-access working memory

- Make p passes over the data, for some small integer p (no random access to tokens)

Minimize space, passes, and processing time upon token arrivals

• Use a sublinear amount of space *s*:

 $s = o(\min\{n, m\})$

where s = bits of random-access working memory

- Make p passes over the data, for some small integer p (no random access to tokens)
- **③** Use small per-item processing time t





"You're looking for the Holy Grail? Have you tried Ebay?"

Minimize space, passes, and processing time upon token arrivals

• Use a sublinear amount of space s:

 $s = o(\min\{n, m\})$

where s = bits of random-access working memory

- Make p passes over the data, for some small integer p (no random access to tokens)
- O Use small per-item processing time t



$$\begin{cases} s = O(\log m + \log n) \\ \text{Happy if } s = O(polylog(\min\{n, m\})) \\ p = 1 \\ t = O(1) \end{cases}$$

Token frequencies

Data stream = sequence $\sigma = \langle a_1, a_2, ..., a_m \rangle$ of tokens drawn from universe $[n] = \{1, 2, ..., n\}$

Token frequencies

Data stream = sequence $\sigma = \langle a_1, a_2, ..., a_m \rangle$ of tokens drawn from universe $[n] = \{1, 2, ..., n\}$

 σ represents a multiset of items and implicitly defines a frequency vector

$$f = \langle f_1, f_2, \dots f_n \rangle$$

where f_i = number of occurrences of item $i \in [n]$ in σ

Example

If $\sigma = \langle 2, 1, 2, 1, 5, 2, 3, 2 \rangle$ and n = 5, then $f = \langle 2, 4, 1, 0, 1 \rangle$

Token frequencies

Data stream = sequence $\sigma = \langle a_1, a_2, ..., a_m \rangle$ of tokens drawn from universe $[n] = \{1, 2, ..., n\}$

 σ represents a multiset of items and implicitly defines a frequency vector

$$f = \langle f_1, f_2, \dots f_n \rangle$$

where f_i = number of occurrences of item $i \in [n]$ in σ

Example If $\sigma = (2, 1, 2, 1, 5, 2, 3, 2)$ and n = 5, then f = (2, 4, 1, 0, 1)

In many streaming problems, wish to compute some statistical properties of the multiset: e.g., majority token (if any), most frequent items, or number of distinct items

Data stream = sequence of **tuples** $\sigma = \langle (a_1, c_1), (a_2, c_2), ... \rangle$ where $(a_i, c_i) \in [n] \times \{-F, ..., F\}$

Upon arrival of (a_i, c_i) , update frequency $f_{a_i} = f_{a_i} + c_i$

New role for *m*: $m = \sum_{j=1}^{n} f_j$

Data stream = sequence of **tuples** $\sigma = \langle (a_1, c_1), (a_2, c_2), ... \rangle$ where $(a_i, c_i) \in [n] \times \{-F, ..., F\}$

Upon arrival of (a_i, c_i) , update frequency $f_{a_i} = f_{a_i} + c_i$

New role for *m*: $m = \sum_{j=1}^{n} f_j$

Basic data stream model: $c_i = 1$ (m = stream length)

Data stream = sequence of **tuples** $\sigma = \langle (a_1, c_1), (a_2, c_2), ... \rangle$ where $(a_i, c_i) \in [n] \times \{-F, ..., F\}$

Upon arrival of (a_i, c_i) , update frequency $f_{a_i} = f_{a_i} + c_i$ New role for m: $m = \sum_{j=1}^n f_j$

Basic data stream model: $c_i = 1$ (m =stream length)

Cash register model: $c_i > 0$ (items can only arrive, their frequencies can be incremented by variable amounts)

Data stream = sequence of **tuples** $\sigma = \langle (a_1, c_1), (a_2, c_2), ... \rangle$ where $(a_i, c_i) \in [n] \times \{-F, ..., F\}$

Upon arrival of (a_i, c_i) , update frequency $f_{a_i} = f_{a_i} + c_i$ New role for m: $m = \sum_{j=1}^n f_j$

Basic data stream model: $c_i = 1$ (m = stream length)

Cash register model: $c_i > 0$ (items can only arrive, their frequencies can be incremented by variable amounts)

Turnstile model: generic c_i (items can arrive and depart from the multiset)

Historical remarks



Origin in the 70s (seminal paper by Munro & Paterson, STOC'78)

Historical remarks



Origin in the 70s (seminal paper by Munro & Paterson, STOC'78)

Gained popularity in the last fifteen years:

- theoretical interest:
 - easy-to-state, but hard-to-solve problems
 - links to other theory areas and to novel computing paradigms (MapReduce)

Historical remarks



Origin in the 70s (seminal paper by Munro & Paterson, STOC'78)

Gained popularity in the last fifteen years:

- theoretical interest:
 - easy-to-state, but hard-to-solve problems
 - links to other theory areas and to novel computing paradigms (MapReduce)
- practical appeal: fast and effective solutions, wide applicability
Historical remarks



Origin in the 70s (seminal paper by Munro & Paterson, STOC'78)

Gained popularity in the last fifteen years:

- theoretical interest:
 - easy-to-state, but hard-to-solve problems
 - links to other theory areas and to novel computing paradigms (MapReduce)
- practical appeal: fast and effective solutions, wide applicability



Alon, Matias & Szegedy: Gödel prize (2005) for their paper on frequency moments approximation (STOC'96, JCSS'99), foundational work for streaming and sketching algorithms

Three puzzles Data stream challenges

 $\pi = \langle \pi_1, \pi_2, ... \pi_{n-1} \rangle$ is a permutation of [1, *n*] with one number missing



 $\pi = \langle \pi_1, \pi_2, ... \pi_{n-1} \rangle$ is a permutation of [1, n] with one number missing





What's the missing number?



 $\pi = \langle \pi_1, \pi_2, ... \pi_{n-1} \rangle$ is a permutation of [1, n] with one number missing





What's the missing number?

Constraint: Carole has limited memory: she can only use $O(\log n)$ bits



 $\pi = \langle \pi_1, \pi_2, ... \pi_{n-1} \rangle$ is a permutation of [1, n] with one number missing





What's the missing number?

Constraint: Carole has limited memory: she can only use $O(\log n)$ bits

$$\frac{n(n-1)}{2} - \sum_{i=1}^{n-1} \pi_i$$





Now π has two missing numbers, x and y: find them, but use only $O(\log n)$ bits!





Now π has two missing numbers, x and y: find them, but use only $O(\log n)$ bits!

Track
$$\begin{cases} S = \frac{n(n+1)}{2} - \sum_{i=1}^{n-2} \pi_i \\ P = n! - \prod_{i=1}^{n-2} \pi_i \end{cases}$$

Solve equations x + y = S and x y = P





Now π has two missing numbers, x and y: find them, but use only $O(\log n)$ bits!

Track
$$\begin{cases} S = \frac{n(n+1)}{2} - \sum_{i=1}^{n-2} \pi_i \\ P = n! - \prod_{i=1}^{n-2} \pi_i \end{cases}$$



Solve equations x + y = S and x y = P

How many bits?
$$\Omega(\log n!) = \Omega(n \log n)$$



Now π has two missing numbers, x and y: find them, but use only $O(\log n)$ bits!





Now π has two missing numbers, x and y: find them, but use only $O(\log n)$ bits!

Track
$$\begin{cases} S_1 = \frac{n(n-1)}{2} - \sum_{i=1}^{n-2} \pi_i \\ S_2 = \frac{n(n+1)(2n+1)}{6} - \sum_{i=1}^{n-2} \pi_i^2 \end{cases}$$



Solve equations $x + y = S_1$ and $x^2 + y^2 = S_2$



Now π has two missing numbers, x and y: find them, but use only $O(\log n)$ bits!

Track
$$\begin{cases} S_1 = \frac{n(n-1)}{2} - \sum_{i=1}^{n-2} \pi_i \\ S_2 = \frac{n(n+1)(2n+1)}{6} - \sum_{i=1}^{n-2} \pi_i^2 \end{cases}$$



Solve equations $x + y = S_1$ and $x^2 + y^2 = S_2$

How many bits?
$$O(\log n^3) = O(\log n)$$

Lesson 1

Some problems can be deterministically solved in:

- logarithmic space
- one pass

Lesson 1

Some problems can be deterministically solved in:

- logarithmic space
- one pass



Most of the times, we're not so lucky



$U = \{1, ...u\}$ fish species in the universe



 $U = \{1, \dots u\}$ fish species in the universe $a_t \in U$ fish species caught at time t



 $U = \{1, ...u\}$ fish species in the universe $a_t \in U$ fish species caught at time t $f_t[j] = |\{a_i \mid a_i = j, i \leq t\}|$ frequency of species j up to time t



 $U = \{1, \dots u\} \text{ fish species in the universe}$ $a_t \in U \text{ fish species caught at time } t$ $f_t[j] = |\{a_i \mid a_i = j, i \leq t\}| \text{ frequency of species } j \text{ up to time } t$

j is rare iff $f_t[j] = 1$



 $U = \{1, ...u\} \text{ fish species in the universe}$ $a_t \in U \text{ fish species caught at time } t$ $f_t[j] = |\{a_i | a_i = j, i \leq t\}| \text{ frequency of species } j \text{ up to time } t$ $j \text{ is rare iff } f_t[j] = 1$

Rarity of catch at time t:
$$\rho_t = \frac{|\{j \mid f_t[j] = 1\}|}{u} = \frac{R_t}{u}$$



 $U = \{1, ...u\} \text{ fish species in the universe}$ $a_t \in U \text{ fish species caught at time } t$ $f_t[j] = |\{a_i \mid a_i = j, i \leq t\}| \text{ frequency of species } j \text{ up to time } t$ $j \text{ is rare iff } f_t[j] = 1$

Rarity of catch at time t:
$$\rho_t = \frac{|\{j \mid f_t[j] = 1\}|}{u} = \frac{R_t}{u}$$

• George is curious and wants to compute rarity



 $U = \{1, ...u\} \text{ fish species in the universe}$ $a_t \in U \text{ fish species caught at time } t$ $f_t[j] = |\{a_i \mid a_i = j, i \leq t\}| \text{ frequency of species } j \text{ up to time } t$ $j \text{ is rare iff } f_t[j] = 1$

Rarity of catch at time t:
$$\rho_t = \frac{|\{j \mid f_t[j] = 1\}|}{u} = \frac{R_t}{u}$$

- George is curious and wants to compute rarity
- 2*u*-bit vector would suffice
- ... but George's suitcase has o(u) size

George cannot compute ρ_t precisely with a deterministic algorithm using only o(u) bits

By contradiction

George cannot compute ρ_t precisely with a deterministic algorithm using only o(u) bits

By contradiction

Let $S \subseteq U$ be a set of species: no duplicates, $|S| = \Theta(u)$

Need $\Omega(|S|) = \Omega(u)$ bits to represent S

If claim is false, could break information theoretic lower bound

George cannot compute ρ_t precisely with a deterministic algorithm using only o(u) bits

By contradiction

Let $S \subseteq U$ be a set of species: no duplicates, $|S| = \Theta(u)$

Need $\Omega(|S|) = \Omega(u)$ bits to represent S

If claim is false, could break information theoretic lower bound

To retrieve S, for each $i \in U$, stream $\langle S, i \rangle$ to George and compare ρ_t and ρ_{t+1} :

George cannot compute ρ_t precisely with a deterministic algorithm using only o(u) bits

By contradiction

Let $S \subseteq U$ be a set of species: no duplicates, $|S| = \Theta(u)$

Need $\Omega(|S|) = \Omega(u)$ bits to represent S

If claim is false, could break information theoretic lower bound

To retrieve S, for each $i \in U$, stream $\langle S, i \rangle$ to George and compare ρ_t and ρ_{t+1} :

• if
$$i
ot\in S$$
, then $R_{t+1} = R_t + 1$ and $ho_{t+1} >
ho_t$

George cannot compute ρ_t precisely with a deterministic algorithm using only o(u) bits

By contradiction

Let
$$S\subseteq U$$
 be a set of species: no duplicates, $|S|=\Theta(u)$

Need $\Omega(|S|) = \Omega(u)$ bits to represent S

If claim is false, could break information theoretic lower bound

To retrieve S, for each $i \in U$, stream $\langle S, i \rangle$ to George and compare ρ_t and ρ_{t+1} :

• if
$$i
ot\in S$$
, then $R_{t+1} = R_t + 1$ and $ho_{t+1} >
ho_t$

• if
$$i \in S$$
, then $R_{t+1} = R_t - 1$ and $\rho_{t+1} < \rho_t$

George cannot compute ρ_t precisely with a deterministic algorithm using only o(u) bits

By contradiction

Let
$$S \subseteq U$$
 be a set of species: no duplicates, $|S| = \Theta(u)$

Need $\Omega(|S|) = \Omega(u)$ bits to represent S

If claim is false, could break information theoretic lower bound

To retrieve S, for each $i \in U$, stream $\langle S, i \rangle$ to George and compare ρ_t and ρ_{t+1} :

• if
$$i
ot\in S$$
, then $R_{t+1} = R_t + 1$ and $ho_{t+1} >
ho_t$

• if
$$i \in S$$
, then $R_{t+1} = R_t - 1$ and $\rho_{t+1} < \rho_t$

Hence ρ decreases $\Leftrightarrow i \in S$

Randomized fish rarity (1/2)

George can approximate ρ_t using 2k = o(u) bits

Sampling:

- pick k random fish species
- maintain rarity $c_1[t], \dots c_k[t]$ of each sampled species (2 bits)

• Return
$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

Randomized fish rarity (1/2)

George can approximate ρ_t using 2k = o(u) bits

Sampling:

- pick k random fish species
- maintain rarity $c_1[t], \dots c_k[t]$ of each sampled species (2 bits)

• Return
$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

Claim: $E[\widetilde{\rho_t}] = \rho_t$

Randomized fish rarity (1/2)

George can approximate ρ_t using 2k = o(u) bits

Sampling:

- pick k random fish species
- maintain rarity $c_1[t], \dots c_k[t]$ of each sampled species (2 bits)

• Return
$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

Claim:
$$E[\widetilde{\rho_t}] = \rho_t$$

- If ρ_t large enough, $\tilde{\rho_t}$ is a good estimate for ρ_t with arbitrarily small precision and good probability
- Requires more advanced probabilistic tools: examples later

Randomized fish rarity (2/2)

$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

 $E[\widetilde{\rho_t}] = \rho_t$

Missing number Fishing Pointer & chaser Recap

Randomized fish rarity (2/2)

$$\widetilde{\rho_t} = \frac{|\{i \in [1, k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

$$E[\widetilde{\rho_t}] = \rho_t$$

$$Y_i$$
 indicator variable: $\begin{cases} Y_i = 1 & \text{if } c_i[t] = 1 \\ Y_i = 0 & \text{otherwise} \end{cases}$

Missing number Fishing Pointer & chaser Recap

Randomized fish rarity (2/2)

$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

$$E[\widetilde{\rho_t}] = \rho_t$$

$$Y_i$$
 indicator variable:

$$\begin{cases}
Y_i = 1 & \text{if } c_i[t] = 1 \\
Y_i = 0 & \text{otherwise}
\end{cases}$$

 $Pr{Y_i = 1} = Pr{\text{the i-th sampled species is rare}} = \frac{R_t}{n} = \rho_t$

Randomized fish rarity (2/2)

$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

$$E[\widetilde{\rho_t}] = \rho_t$$

$$Y_i$$
 indicator variable: $\begin{cases} Y_i = 1 & \text{if } c_i[t] = 1 \\ Y_i = 0 & \text{otherwise} \end{cases}$

$$Pr{Y_i = 1} = Pr{\text{the i-th sampled species is rare}} = \frac{R_t}{u} = \rho_t$$

 $\Rightarrow E[Y_i] = \rho_t$

Missing number Fishing Pointer & chaser Recap

Randomized fish rarity (2/2)

$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

$$E[\widetilde{\rho_t}] = \rho_t$$

$$Y_i$$
 indicator variable:
$$\begin{cases} Y_i = 1 & \text{if } c_i[t] = 1 \\ Y_i = 0 & \text{otherwise} \end{cases}$$

 $Pr{Y_i = 1} = Pr{\text{the i-th sampled species is rare}} = \frac{R_t}{u} = \rho_t$

$$\Rightarrow E[Y_i] = \rho_t$$

$$\Rightarrow E[\widetilde{R_t}] = \sum_{i=1}^k E[Y_i] = k\rho_t$$

Missing number Fishing Pointer & chaser Recap

Randomized fish rarity (2/2)

$$\widetilde{\rho_t} = \frac{|\{i \in [1,k] \mid c_i[t] = 1\}|}{k} = \frac{\widetilde{R_t}}{k}$$

$$E[\widetilde{\rho_t}] = \rho_t$$

$$Y_i$$
 indicator variable:
$$\begin{cases} Y_i = 1 & \text{if } c_i[t] = 1 \\ Y_i = 0 & \text{otherwise} \end{cases}$$

 $Pr{Y_i = 1} = Pr{\text{the i-th sampled species is rare}} = \frac{R_t}{u} = \rho_t$

$$\Rightarrow E[Y_i] = \rho_t$$

$$\Rightarrow E[\widetilde{R_t}] = \sum_{i=1}^k E[Y_i] = k\rho_t$$

$$\Rightarrow E[\widetilde{\rho_t}] = \frac{E[\widetilde{R_t}]}{k} = \rho_t$$
Lesson 2

It is often impossible to solve problems precisely and deterministically in small (sublinear) space



Randomization and approximation greatly help:

- find an answer correct within some factor (guarantee that $\tilde{\rho}$ is within 10% of ρ)
- allow a small probability of failure (answer is correct, except with probability 1 in 10,000)

Pointer and chaser



Paul has n + 1 pointers

For each pointer *i*, he points to a position $P[i] \in [1, n]$



Pointer and chaser



Paul has n + 1 pointers

For each pointer *i*, he points to a position $P[i] \in [1, n]$



Carole has to guess any duplicate pointer

Constraints:

- $O(\log n)$ bits
- O(n) queries
- cannot move items



Repeated scans



• Trivial solution

- for each i, count how many j are such that P[j]=i
- $O(\log n)$ bits, but $O(n^2)$ queries

Repeated scans



• Trivial solution

- for each i, count how many j are such that P[j]=i
- $O(\log n)$ bits, but $O(n^2)$ queries

Ø Better solution

- if # of items below n/2 > # of items above n/2 then search for duplicates < n/2 else search for duplicates ≥ n/2
- $O(\log n)$ bits and passes, $O(n \log n)$ queries
- With $O(\log n)$ bits, $\Omega(\log n / \log \log n)$ passes are needed



- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!





- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!





- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!





- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!





- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!





- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!





- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!



$$\begin{cases} a+b=t\\ a+k(b+c)+b=2t \end{cases}$$

t and k known



- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!



$$\left\{\begin{array}{l} a+b=t\\ a+k(b+c)+b=2t \end{array} \Rightarrow \left\{\begin{array}{l} a+b=t\\ b+c=t/k \end{array}\right.$$

t and k known

31 / 99



- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!



t and k known



- Chase pointers, starting from position n + 1
- Problem equivalent to finding a loop in a linked list
- Can be solved in O(n) time with just 2 pointers!



t and k known



Tokens come as a stream: no random access



Sometimes impossible to achieve the same bounds as in the RAM model

Recap on lessons



Typically impossible to solve problems precisely and deterministically in small (sublinear) space



Randomize and approximate!

Sequential access 1 2 3 4 5 6 7 8Random access

Sequential data access makes things harder

Sampling Working with less

Why sampling?

- Basic problem: sample *s* items uniformly from a stream
- Answer queries (e.g., compute fish species rarity) on the sample
- Utility depends on the problem: in some cases, sampling-based approaches not effective unless taking large (almost linear) samples

Why sampling?

- Basic problem: sample *s* items uniformly from a stream
- Answer queries (e.g., compute fish species rarity) on the sample
- Utility depends on the problem: in some cases, sampling-based approaches not effective unless taking large (almost linear) samples

How can we sample uniformly if we don't know in advance how long is the stream? When do we sample a stream token?

Reservoir sampling

- Add to S the first s stream items
- **2** Upon seeing x_i at time, sample x_i with probability s/i
- **③** If x_i added to S, evict a random item from S (other than x_i)

Sample is uniform

At any time t and for each $i \leq t$, it holds: $Pr\{x_i \in_t S\} = \frac{s}{t}$

Reservoir sampling

- Add to S the first s stream items
- **2** Upon seeing x_i at time, sample x_i with probability s/i
- **③** If x_i added to S, evict a random item from S (other than x_i)

Sample is uniform

At any time t and for each $i \le t$, it holds: $Pr\{x_i \in_t S\} = \frac{s}{t}$

Warmup analysis: s = 1

 $Pr\{x_i \in_t S\} =$

 $= Pr\{x_i \text{ sampled at time } i\} \times Pr\{x_i \text{ survives up to time } t\} =$

$$=\frac{1}{i}\times\frac{i}{i+1}\times\frac{i+1}{i+2}\times\ldots\times\frac{t-2}{t-1}\times\frac{t-1}{t}=\frac{1}{t}$$

Sample is uniform:
$$Pr\{x_i \in t S\} = \frac{s}{t}$$

Sample is uniform: $Pr\{x_i \in t S\} = \frac{s}{t}$

By induction on t (base step: $t \leq s$)

Sample is uniform: $Pr\{x_i \in t S\} = \frac{s}{t}$

By induction on t (base step: $t \leq s$)

How does S change at time t when x_t arrives?

Sample is uniform:
$$Pr\{x_i \in t S\} = \frac{s}{t}$$

By induction on t (base step: $t \leq s$)

How does S change at time t when x_t arrives?

• Pr
$$\{x_t \text{ added to } S\} = \frac{s}{t}$$

Sample is uniform:
$$Pr\{x_i \in t S\} = \frac{s}{t}$$

By induction on t (base step: $t \leq s$)

How does S change at time t when x_t arrives?

•
$$Pr\{x_t \text{ added to } S\} = \frac{s}{t}$$

2 Inductive hypothesis:
$$Pr\{x_i \in t-1} S\} = \frac{s}{t-1}$$

Sample is uniform:
$$Pr\{x_i \in t S\} = \frac{s}{t}$$

By induction on t (base step: $t \leq s$)

How does S change at time t when x_t arrives?

•
$$Pr\{x_t \text{ added to } S\} = \frac{s}{t}$$

② Inductive hypothesis: $Pr\{x_i \in t-1} S\} = \frac{s}{t-1}$

• Pr{x_i $\in_t S | x_t \text{ added to } S$ } = Pr{x_i $\in_{t-1} S \text{ and not evicted}} =$ $= <math>\frac{s}{t-1} \left(1 - \frac{1}{s}\right)$

Sample is uniform:
$$Pr\{x_i \in t S\} = \frac{s}{t}$$

By induction on t (base step: $t \leq s$)

How does S change at time t when x_t arrives?

•
$$Pr\{x_t \text{ added to } S\} = \frac{s}{t}$$

2 Inductive hypothesis: $Pr\{x_i \in t-1} S\} = \frac{s}{t-1}$

 $Pr\{x_i \in_t S \mid x_t \text{ added to } S\} = Pr\{x_i \in_{t-1} S \text{ and not evicted}\} = \frac{s}{t-1} \left(1 - \frac{1}{s}\right)$

 $Pr\{x_i \in_t S \mid x_t \text{ not added to } S\} = Pr\{x_i \in_{t-1} S\} = \frac{s}{t-1}$

Sample is uniform:
$$Pr\{x_i \in t S\} = \frac{s}{t}$$

By induction on t (base step: $t \leq s$)

How does S change at time t when x_t arrives?

•
$$Pr\{x_t \text{ added to } S\} = \frac{s}{t}$$

② Inductive hypothesis: $Pr\{x_i \in t-1} S\} = \frac{s}{t-1}$

• Pr{x_i $\in_t S | x_t \text{ added to } S$ } = Pr{x_i $\in_{t-1} S \text{ and not evicted}} =$ $= <math>\frac{s}{t-1} \left(1 - \frac{1}{s}\right)$

 $Pr\{x_i \in_t S \mid x_t \text{ not added to } S\} = Pr\{x_i \in_{t-1} S\} = \frac{s}{t-1}$

By combining conditional probabilities:

$$Pr\{x_i \in t S\} = \frac{s}{t} \frac{s}{t-1} \left(1 - \frac{1}{s}\right) + \left(1 - \frac{s}{t}\right) \frac{s}{t-1} = \frac{s}{t}$$

Optimizations and drawbacks

Skip numbers

Instead of flipping a coin at each stream element, generate number of elements to be skipped before the next element is added to S [Vitter 85]

Optimizations and drawbacks

Skip numbers

Instead of flipping a coin at each stream element, generate number of elements to be skipped before the next element is added to S [Vitter 85]

Other issues:

- Frequently occurring values are a wasteful use of the available sample space: concise sampling [Gibbons and Matias '98]
- Runs into difficulties in the presence of data deletions: [Babcock *et al.* '02]
- Hard to parallelize on multiple streams: how do we sample if more than one item comes at any time? Min-wise sampling [Nath *et al.* '04]

The Britney Spears problem...



Algorithms for data streams

... tracking who's hot and who's not



"... can't just pay attention to a few popular subjects, because you can't know in advance which ones are going to rank near the top. To be certain of catching every new trend as it unfolds, you have to monitor *all* the incoming queries – and their variety is unbounded. "



Given a stream of n items, find those that appear "most frequently"



E.g., items occurring more than 1% of the time

- Formally "hard" in small space, so allow approximation
- No false negatives: return all items with count $\geq \varphi n$
- "Good" false positives: no item with count $< (\varphi \varepsilon)n$ is returned (error $\varepsilon \in (0, 1), \varepsilon \ll \varphi$)
- Related problem: estimate each frequency with error $\pm arepsilon n$

Heavy hitters

Why heavy hitters?

- Many practical applications: mining of search logs, analysis of network data, DBMS optimization...
- Core streaming problem: connections with entropy estimation, itemsets mining, compressed sensing
- Extensive research: scores of streaming papers on frequent items and its variations

We'll see a counter-based algorithm named Sticky sampling:

- probabilistic, sampling-based approach
- ② correct with probability $\geq 1 \delta$, with $\delta \in (0, 1)$ user-specified probability of failure
Sticky sampling

Intuition

It should be possible to estimate frequent items by a good sample

Data structure S: set of pairs $\langle x, f_e(x) \rangle$, where

- $f_e(x)$ estimated frequency of x
- f(x) true frequency

Query algorithm: at time *n* report items $x \in S$ such that $f_e(x) \ge (\varphi - \varepsilon)n$

Update algorithm works in rounds:

- each round distinguished by a (fixed) sampling rate r
- sampling rate adjusted between rounds so that probability of sampling a stream item decreases as stream gets longer

Update algorithm

Structure of *r*-rate round

For each stream item x:

- if $x \in S$, then increase $f_e(x)$ by 1
- if x ∉ S, sample x with probability ¹/_r: if x sampled, add pair ⟨x,1⟩ to S

At the end of a round:

- double sampling rate *r* (*r* increases geometrically)
- adjust estimated frequencies so that S is transformed into exactly the state it would have been in, if new rate 2r had been used from the beginning

Adjusting frequencies

Assume x sampled at time k with probability $\frac{1}{r}$:

- $f_e(x) = \text{exact number of occurrences of } x \text{ after time } k$
- with smaller sampling probability $(\frac{1}{2r})$, x will be sampled at one of the later occurrences
- simulate all coin tosses not done with sampling rate r

For each $\langle x, f_e(x) \rangle \in S$ repeatedly toss a coin:

- first coin toss unbiased $(\frac{1}{2})$, makes probability of sampling x at time $k = \frac{1}{2r}$
- 2 next coin tosses biased with probability $\frac{1}{2r}$
- **③** for each unsuccessful coin toss, decrease $f_e(x)$ by 1
- stop when coin toss successful or f_e(x) = 0 (in this case remove x from S)

Round length

Recall:
$$\begin{cases} \varphi = \text{frequency threshold} \\ \varepsilon = \text{frequency error} \\ \delta = \text{algorithm failure probability} \end{cases}$$

Let
$$t=rac{1}{arepsilon}\lograc{1}{arphi\delta}$$



r-rate round has length *rt* (except for r = 1) expected sample size: 2*t* (we'll prove)

Irene Finocchi

A technical lemma

For each rate $r \ge 2$, let *n* be the number of stream items considered up to the *r*-rate round. It holds:

$$\frac{1}{r} \ge \frac{t}{n}$$

A technical lemma

For each rate $r \ge 2$, let *n* be the number of stream items considered up to the *r*-rate round. It holds:

$$\frac{1}{r} \ge \frac{t}{n}$$

By induction, at the beginning of *r*-rate round n = rt:

	n=	rt n'=n+	-rt=2rt
		rt	
rate	2	r	2r

Hence during the round $n \ge rt$

A technical lemma

For each rate $r \ge 2$, let *n* be the number of stream items considered up to the *r*-rate round. It holds:

$$\frac{1}{r} \ge \frac{t}{n}$$

By induction, at the beginning of *r*-rate round n = rt:

	n=	rt n'=n+	rt=2rt
		rt	
rate	e	r	2r

Hence during the round $n \ge rt$



For any $\varphi, \varepsilon, \delta \in (0, 1)$, with $\varepsilon < \varphi$, Sticky Sampling computes the heavy hitters with probability $\geq 1 - \delta$

For any $\varphi, \varepsilon, \delta \in (0, 1)$, with $\varepsilon < \varphi$, Sticky Sampling computes the heavy hitters with probability $\ge 1 - \delta$

Good false positives: items with frequency < (φ − ε)n are not returned
 f(x) < (φ − ε)n ⇒ f_e(x) < (φ − ε)n, since f_e(x) < f(x)

For any $\varphi, \varepsilon, \delta \in (0, 1)$, with $\varepsilon < \varphi$, Sticky Sampling computes the heavy hitters with probability $\ge 1 - \delta$

- Good false positives: items with frequency < (φ − ε)n are not returned
 f(x) < (φ − ε)n ⇒ f_e(x) < (φ − ε)n, since f_e(x) ≤ f(x)
- **2** No false negatives: all items with frequency $\geq \varphi n$ are returned

For any $\varphi, \varepsilon, \delta \in (0, 1)$, with $\varepsilon < \varphi$, Sticky Sampling computes the heavy hitters with probability $\geq 1 - \delta$

- Good false positives: items with frequency < (φ − ε)n are not returned
 f(x) < (φ − ε)n ⇒ f_e(x) < (φ − ε)n, since f_e(x) ≤ f(x)
- **2** No false negatives: all items with frequency $\geq \varphi n$ are returned

$$\begin{array}{ll} y_1 \ \dots \ y_k \ \text{frequent items:} & f(y_i) \geq \varphi n \quad \forall i \\ & \Rightarrow k \leq \frac{1}{\varphi} \end{array}$$

$$Pr\{\exists \ \text{false negative}\} = & Pr\{\exists y_i : y_i \ \text{not returned}\} \leq \\ & \sum_{i=1}^k Pr\{y_i \ \text{not returned}\} \end{array}$$



$$Pr\{y_i \text{ not returned}\} = Pr\{f_e(y_i) < (\varphi - \varepsilon)n\} =$$

 $Pr\{ at | east \in n | unsuccessful | coin | tosses \} \leq$

$$\left(1-\frac{1}{r}\right)^{\varepsilon n} \leq \left(1-\frac{t}{n}\right)^{\varepsilon n} \leq e^{-t\varepsilon}$$



$$Pr\{y_i \text{ not returned}\} = Pr\{f_e(y_i) < (\varphi - \varepsilon)n\} = Pr\{\text{at least } \varepsilon n \text{ unsuccessful coin tosses}\} \le \left(1 - \frac{1}{r}\right)^{\varepsilon n} \le \left(1 - \frac{t}{n}\right)^{\varepsilon n} \le e^{-t\varepsilon}$$

Hence:

$$\begin{array}{ll} \Pr\{\exists \ \mathsf{false} \ \mathsf{negative}\} & \leq \sum_{i=1}^{k} \Pr\{y_i \ \mathsf{not} \ \mathsf{returned}\} \leq \\ & \leq k e^{-t\varepsilon} \leq \frac{e^{-t\varepsilon}}{\varphi} = \delta \quad \ \mathsf{by} \ \mathsf{definition} \ \mathsf{of} \ t \end{array}$$

Sketching streams



Sketches

- Not every problem can be solved with sampling
 E.g., counting distinct items in a stream: need to sample a large fraction of items to know if they are all same or different
- Sketches take advantage that the algorithm can "see" all the data even if it can't "remember" it all

Sketch = linear transform of the input (exploit hashing)

Sampling and sketching ideas at the heart of stream mining:

- A sample is a quite general representative of the data set
- Sketches tend to be tailored to a specific problem (e.g., distinct items)

Warmup example

Problem: test if two asynchronous binary streams are equal

100101110100 100100110100

To test in small space: pick a random hash function h and test $h(\sigma_1) = h(\sigma_2)$:

- no false negatives: if $\sigma_1 = \sigma_2$ then $h(\sigma_1) = h(\sigma_2)$
- small chance of false positive: it may be $h(\sigma_1) = h(\sigma_2)$ for $\sigma_1 \neq \sigma_2$ with very small probability

Compute $h(\sigma_1)$ and $h(\sigma_2)$ incrementally as new bits arrive (Karp-Rabin fingerprints)

Distinct items

Count of the number of distinct items seen in the stream

Trivial solution: maintain set of encountered items through its characteristic vector

O(1) processing time but $\Theta(u)$ space, where u = universe size

- Exact/deterministic algorithms need $\Omega(u)$ bits of space
- Approximate randomized algorithms use $O(\log u)$ bits of space

FM-sketch [Flajolet & Martin '85]

Sampling not appropriate here: we'll build a data summary (sketch)

Universal hashing

- $\bullet\,$ Idea: select a hash function at random from a family ${\cal H}$ of hash functions with a certain mathematical property
- Guarantee: low number of collisions in expectation, even if the data is chosen by an adversary

2-universal hashing

 \mathcal{H} is a 2-universal family (set) of hash functions $h: U \rightsquigarrow D$ if, for all $x, y \in U, x \neq y$: $Pr_{h \in \mathcal{H}} \{h(x) = h(y)\} \leq \frac{1}{|D|}$

Strongly 2-universal hashing

 \mathcal{H} is strongly 2-universal if, for all $x \neq y \in U$ and $a, b \in D$: $Pr_{h \in \mathcal{H}}\{h(x) = a \& h(y) = b\} = \frac{1}{|D|^2}$

FM skecth: probabilistic counter

Two useful functions:

h: *U* → [0, *u* − 1] drawn from a family of strongly 2-universal hash functions

Transforms values of the universe into integers uniformly distributed over the set of binary strings of length $\log u$

t: [0, u − 1] → [1, log u] gives the number t(i) in the binary representation of i
 E.g., t(5₁₀) = t(00101₂) = 2

FM sketch: counter C of log u bits

Counter update: upon seeing stream item *x*, set C[t(h(x))] = 1

Query algorithm: return 2^R , where $R \in [1, \log u]$ is the position of the rightmost 1 in C

E.g., if C = 1110100, then R = 5: returns 32

Intuition

h distributes items of the universe U uniformly on [0, u - 1]: important to avoid adversarial streams

- How many values in [0, u 1] have exactly 0 trailing 0s? u/2
- How many values have exactly 1 trailing 0? u/4
- How many values have exactly 2 trailing 0s? $u/8 \dots$

Hence, if the stream contains D distinct values:

- D/2 will be mapped to the first bit of C
- D/4 to the second bit
- D/8 to the third bit ...

We expect the first log D counter bits will be set to 1 Hence $R \approx \log D$ and $2^R \approx D$

Geometric distribution over counter bits

- |values with exactly j trailing $0s| = \frac{u}{2^{j+1}}$
- |values with $\geq j$ trailing 0s | = 1 + $\sum_{i=j}^{\log u-1} \frac{u}{2^{j+1}} = 2^{\log u-j}$
- W_x indicator random variable: $W_x = 1$ iff $t(h(x)) \ge j$ $Pr\{W_x = 1\} = Pr\{t(h(x)) \ge j\} = \frac{2^{\log u - j}}{u} = 2^{-j}$

since h distributes items uniformly over [0, u - 1]

- $E[W_x] = 2^{-j}$
- $Var[W_x] = E[W_x^2] E[W_x]^2 = 2^{-j} 2^{-2j} < 2^{-j} = E[W_x]$

$$E[W_x] = 2^{-j}$$
 and $Var[W_x] < E[W_x]$

Geometric distribution over counter bits

• Z_j = number of stream items x s.t. $t(h(x)) \ge j$ = $\sum_{x \in U \cap \Sigma} W_x$

•
$$E[Z_j] = \sum_{x \in U \cap \Sigma} E[W_x] = \sum_{x \in U \cap \Sigma} 2^{-j} = \frac{D}{2^j}$$

- Due to pairwise independence of W_x and W_y , $Var[W_x + W_y] = Var[W_x] + Var[W_y]$
- $Var[Z_j] = \sum_{x \in U \cap \Sigma} Var[W_x] < \sum_{x \in U \cap \Sigma} E[W_x] = E[Z_j]$

$$E[Z_j] = rac{D}{2^j}$$
 and $Var[Z_j] < E[Z_j]$

•
$$R = \max j$$
 such that $Z_j > 0$

Probability of overestimating

Let c > 2. $Pr\{2^R > cD\} = ?$

By Markov's inequality (Z_j takes only non-negative values):

$$Pr\{Z_j \ge 1\} \le \frac{E[Z_j]}{1} = \frac{D}{2^j}$$
 (1)

$$\begin{array}{ll} 2^R > cD & \Rightarrow \exists j \text{ such that } C[j] = 1 \& 2^j > cD \\ & \Rightarrow C[j] = 1 \& j > \log_2(cD) \\ & \Rightarrow Z_{\log_2(c\,D)} \ge 1 \end{array}$$

Thus:

$$\Pr\{2^{R} > cD\} \le \Pr\{Z_{\log_{2}(cD)} \ge 1\} \le_{(1)} \frac{D}{2^{\log_{2}(cD)}} = \frac{1}{c}$$

Probability of underestimating

Let c > 2. $Pr\{2^R < \frac{D}{c}\} = ?$

By Chebyshev inequality (Z_j takes only non-negative values):

$$Pr\{Z_{j} = 0\} = Pr\{|Z_{j} - E[Z_{j}]| \ge E[Z_{j}]\}$$

$$\le \frac{Var[Z_{j}]}{E[Z_{j}]^{2}} < \frac{1}{E[Z_{j}]} = \frac{2^{j}}{D}$$
(2)

$$2^{R} < \frac{D}{c} \quad \Rightarrow C[p] = 0 \quad \forall p \ge \log_{2}(D/c)$$
$$\Rightarrow Z_{\log_{2}(D/c)} = 0$$

Thus:

$$\Pr\left\{2^{R} < \frac{D}{c}\right\} \le \Pr\{Z_{\log_{2}(D/c)} = 0\} \le_{(2)} \frac{2^{\log_{2}(D/c)}}{D} = \frac{1}{c}$$

Distinct items: summing up

Let D be the exact number of distinct values and let 2^R be the output of the probabilistic counter.

For any c > 2, the probability that 2^R is not between D/c and cD is at most 2/c.

Frequency moments

Stream $\Sigma = \langle x_1, x_2, ..., x_n \rangle$ of tokens drawn from universe U $f_i = |\{j : x_j = i\}|$



Useful statistical information:

- F₀ = distinct items
- $F_1 = \text{stream length}$
- $F_2 = \text{Gini's index (skew of the data)}$
- F_{∞} related to maximum frequency element, i.e., $\max_{i \in U} f_i$

AMS sketch for F_2

Fundamental technique introduced by Alon, Matias, and Szegedy

AMS sketches = randomized linear projections

Define a random variable Z such that $E[Z^2] = F_2$:

- select at random a hash function $\xi: U \rightsquigarrow \{-1, +1\}$ from a family of 4-wise independent hash functions
- $Z = \sum_{u \in U} f_u \xi(u)$

random linear projection (inner product) of frequency vector $\langle f_1, f_2, ..., f_u \rangle$ with random vector $\{-1, +1\}^u$

• Z incrementally updated upon arrival of x_t by adding $\xi(x_t)$

AMS sketch: expectation

$$Z = \sum_{u \in U} f_u \xi(u)$$

 $\xi : U \rightsquigarrow \{-1, +1\}$ 4-wise independent
 $E[\xi(i)] = (-1)\frac{1}{2} + (1)\frac{1}{2} = 0$

$$E[Z^2] = E\left[\left(\sum_{i \in U} f_i \xi(i)\right)^2\right]$$

 $= E\left[\sum_{i \in U} f_i^2 (\xi(i))^2 + 2\sum_{i \neq j \in U} f_i f_j \xi(i)\xi(j)\right]$
 $= \sum_{i \in U} f_i^2 E\left[(\xi(i))^2\right] + 2\sum_{i \neq j \in U} f_i f_j E\left[\xi(i)\xi(j)\right]$
 $= \sum_{i \in U} f_i^2 = F_2$

since $(\xi(i))^2 = 1$ and by pair-wise independence $E[\xi(i)\xi(j)] = E[\xi(i)]E[\xi(j)] = 0 \cdot 0 = 0$

Median of the averages

Still need small variance and good confidence:

- Compute μ random variables Y₁, ..., Y_μ and output their median Y as the estimator for F₂
- Each Y_i is the average of α independent, identically distributed random variables X_{ij} computed as random linear projections

Averaging X_{ij} implies each Y_i has small variance Computing Y as the median of the Y_i allows it to boost confidence using Chernoff bounds

F_2 : summing up

For every $\lambda, \delta > 0$, there exists a randomized algorithm that computes a number Y that deviates from F_2 by more than λF_2 with probability at most δ .

The algorithm uses only

$$O\left(\frac{\log(1/\delta)}{\lambda^2}(\log u + \log n)\right)$$

memory bits and performs one pass over the data.

Similar results for frequency moments F_k , with k > 2

Mining graphs

68 / 99

Models for graph streams

- G = (V, E) graph with |V| = n nodes and |E| = m edges, possibly weighted
- Observe edges of G in a stream, one by one
- What order do we see the edges in?
 - Arbitrary (adversarial) order
 - Incidence streams: all edges incident to one vertex appear sequentially (easier, stronger bounds)



- How many passes over the data can we take (one or many?)
- How much space?

Counting triangles

- Finding frequent graph patterns and dense subgraphs are basic tools in the analysis of the structure of large networks (e.g., social networks, Web graph)
- Exact triangle counting reduces to matrix multiplication: unfeasible even for networks of medium size
- Resort to random sampling
- We'll present an algorithm for the arbitrary order model

A 3-pass algorithm

Algorithm SampleTriangle

1st pass. Count number of edges *m* in the stream

2nd pass. Sample an edge e = (a, b) uniformly from *E* and a node *v* uniformly from $V \setminus \{a, b\}$

3rd pass. If $(a, v) \in E$ and $(b, v) \in E$ then $\beta = 1$, else $\beta = 0$



A useful property



 T_i = triples with *i* edges, $0 \le i \le 3$

$$E[\beta] = \frac{3|T_3|}{m \cdot (n-2)} = \frac{3|T_3|}{|T_1| + 2|T_2| + 3|T_3|}$$

m · (*n* − 2) ways to select an edge (*a*, *b*) and a node *v* ≠ *a*, *b i*|*T_i*| ways to select a triple with *i* edges, *i* > 0

The complete 3-pass algorithm

• Start s parallel instances of algorithm SampleTriangle, where

$$s \geq \frac{3}{\varepsilon^2} \frac{|T_1| + 2|T_2| + 3|T_3|}{|T_3|} \ln\left(\frac{2}{\delta}\right)$$

• Each instance returns a value β_i

• Return
$$\widetilde{T}_3 = \left(\frac{1}{s}\sum_{i=1}^s \beta_i\right) \frac{m \cdot (n-2)}{3}$$
 as an estimation for T_3

$$E[\widetilde{T_3}] = |T_3|$$
 because $E[\beta_i] = \frac{3|T_3|}{m \cdot (n-2)}$
OK, but how far from the mean?
Chernoff bounds

X₁, X₂, ... X_n independent Bernoulli trials: X_i indicator random variable, Pr{X_i = 1} = p, X_i all independent

•
$$X = \sum_{i=1}^{n} X_i$$

•
$$E[X] = \mu = n p$$

Lower tail bound

For any
$$arepsilon\in(0,1]$$
 $Pr\{X<(1-arepsilon)\mu\}< e^{-rac{\muarepsilon^2}{2}}$

Upper tail bound

For any
$$\varepsilon \in (0,1]$$
 $\Pr\{X > (1+\varepsilon)\mu\} < e^{-rac{\mu\varepsilon^{-}}{3}}$

2

In triangle counting,
$$X = \sum_{i=1}^{s} \beta_i$$
 and $p = \frac{3|\mathcal{T}_3|}{|\mathcal{T}_1|+2|\mathcal{T}_2|+3|\mathcal{T}_3|}$

 $\Pr\{X < (1-\varepsilon) ps \mid\mid X > (1+\varepsilon) ps\} \quad < e^{-\frac{ps\varepsilon^2}{2}} + e^{-\frac{ps\varepsilon^2}{3}}$

In triangle counting,
$$X = \sum_{i=1}^{s} \beta_i$$
 and $p = \frac{3|T_3|}{|T_1|+2|T_2|+3|T_3|}$

$$\begin{aligned} \Pr\{X < (1-\varepsilon)ps \mid| X > (1+\varepsilon)ps\} &< e^{-\frac{ps\varepsilon^2}{2}} + e^{-\frac{ps\varepsilon^2}{3}} \\ &\leq 2e^{-\frac{sp\varepsilon^2}{3}} \end{aligned}$$

In triangle counting,
$$X = \sum_{i=1}^{s} \beta_i$$
 and $p = \frac{3|T_3|}{|T_1|+2|T_2|+3|T_3|}$

$$Pr\{X < (1-\varepsilon)ps \mid | X > (1+\varepsilon)ps\} < e^{-\frac{ps\varepsilon^2}{2}} + e^{-\frac{ps\varepsilon^2}{3}} \le 2e^{-\frac{sp\varepsilon^2}{3}} \le \delta$$

as long as
$$s \ge \frac{3}{\varepsilon^2} \frac{|I_1| + 2|I_2| + 3|I_3|}{|T_3|} \ln \left(\frac{2}{\delta}\right)$$

In triangle counting,
$$X = \sum_{i=1}^{s} \beta_i$$
 and $p = \frac{3|T_3|}{|T_1|+2|T_2|+3|T_3|}$

$$\begin{aligned} \Pr\{X < (1-\varepsilon)ps \mid\mid X > (1+\varepsilon)ps\} &< e^{-\frac{ps\varepsilon^2}{2}} + e^{-\frac{ps\varepsilon^2}{3}} \\ &\leq 2e^{-\frac{sp\varepsilon^2}{3}} \leq \delta \end{aligned}$$

as long as $s \geq \frac{3}{\varepsilon^2} \, \frac{|\mathcal{T}_1| + 2|\mathcal{T}_2| + 3|\mathcal{T}_3|}{|\mathcal{T}_3|} \, \ln \left(\frac{2}{\delta} \right)$

$$X < (1-\varepsilon)ps \Leftrightarrow \underbrace{\left(\frac{1}{s}\sum_{i=1}^{s}\beta_{i}\right)\frac{m\cdot(n-2)}{3}}_{\widetilde{T_{3}}} < (1-\varepsilon)\underbrace{p\frac{m(n-2)}{3}}_{T_{3}}$$

In triangle counting,
$$X = \sum_{i=1}^{s} \beta_i$$
 and $p = \frac{3|T_3|}{|T_1|+2|T_2|+3|T_3|}$

$$\begin{aligned} \Pr\{X < (1-\varepsilon)ps \mid| X > (1+\varepsilon)ps\} &< e^{-\frac{ps\varepsilon^2}{2}} + e^{-\frac{ps\varepsilon^2}{3}} \\ &\leq 2e^{-\frac{sp\varepsilon^2}{3}} \leq \delta \end{aligned}$$

as long as $s \geq \frac{3}{\varepsilon^2} \frac{|T_1| + 2|T_2| + 3|T_3|}{|T_3|} \ln\left(\frac{2}{\delta}\right)$

$$X < (1 - \varepsilon)ps \Leftrightarrow \underbrace{\left(\frac{1}{s}\sum_{i=1}^{s}\beta_{i}\right)\frac{m \cdot (n-2)}{3}}_{\widetilde{T_{3}}} < (1 - \varepsilon)\underbrace{p\frac{m(n-2)}{3}}_{T_{3}}$$

Similarly, $X > (1 + \varepsilon)ps \leftrightarrow \widetilde{T} > (1 + \varepsilon)T$

Similarly $X > (1 + \varepsilon)ps \Leftrightarrow \widetilde{T_3} > (1 + \varepsilon)T_3$

- when edge (a, b) and node v sampled, hash missing edges
 (a, v) and (b, v) to a set M
- in the third pass, lookup each edge (x, y) in M, and mark it if present
- Itriangles determined in a postprocessing step

- when edge (a, b) and node v sampled, hash missing edges
 (a, v) and (b, v) to a set M
- in the third pass, lookup each edge (x, y) in M, and mark it if present
- **③** triangles determined in a postprocessing step
- 1-pass: exploit reservoir sampling

- when edge (a, b) and node v sampled, hash missing edges
 (a, v) and (b, v) to a set M
- in the third pass, lookup each edge (x, y) in M, and mark it if present
- **③** triangles determined in a postprocessing step
- 1-pass: exploit reservoir sampling
- $\bullet\,$ Other minors and cliques of size $\alpha\,$

- when edge (a, b) and node v sampled, hash missing edges
 (a, v) and (b, v) to a set M
- in the third pass, lookup each edge (x, y) in M, and mark it if present
- **③** triangles determined in a postprocessing step
- 1-pass: exploit reservoir sampling
- $\bullet\,$ Other minors and cliques of size $\alpha\,$
- Better space bounds for incidence streams

Semi-streaming model

For many graph problems space × passes = Ω(n), even using randomization and approximation

 \Rightarrow Cannot achieve O(1) passes and polylog working space

• Semi-streaming model: polylog space requirement is relaxed

working memory size *O*(*n* polylog *n*) for input graph with *n* nodes

enough space to store nodes, not enough for edges

• Problems solvable in semi-streaming: spanners, matching, diameter estimation...

Maximum weight matching

- Edge weighted, undirected graph G(V, E, w)
- No two edges in a matching have a common endpoint



Optimization problem: find a maximum weight matching M^*

1-pass semi-streaming algorithm with approximation ratio 1/6:

 $w(M) \geq \frac{w(M^*)}{6}$

where M returned matching

Semi-streaming algorithm

Data structure: matching *M* maintained in main memory

Query algorithm: return M

Update algorithm: upon arrival of edge e, consider set $C \subseteq M$ of conflicting edges (edges in M that share an endpoint with e)

- if w(e) > 2w(C), replace C with $\{e\}$ in M
- if $w(e) \leq 2w(C)$, ignore e

$$\begin{split} \Sigma &= \langle \; (c,f,2) \; (b,e,10) \; (h,i,4) \; (e,f,30) \; (h,f,50) \\ & (e,g,40) \; (d,e,62) \; (a,d,120) \; (d,g,130) \; \rangle \end{split}$$



Every edge $e \in M$ is root of a replacement tree T_e

$$\begin{split} \Sigma &= \langle \; (c,f,2) \; (b,e,10) \; (h,i,4) \; (e,f,30) \; (h,f,50) \\ & (e,g,40) \; (d,e,62) \; (a,d,120) \; (d,g,130) \; \rangle \end{split}$$



Every edge $e \in M$ is root of a replacement tree T_e $R(e) = \text{nodes in } T_e \text{ except for root } e$

$$\begin{split} \Sigma &= \langle \; (c,f,2) \; (b,e,10) \; (h,i,4) \; (e,f,30) \; (h,f,50) \\ & (e,g,40) \; (d,e,62) \; (a,d,120) \; (d,g,130) \; \rangle \end{split}$$



Every edge $e \in M$ is root of a replacement tree T_e $R(e) = \text{nodes in } T_e \text{ except for root } e$

80 / 99

$$\begin{split} \Sigma &= \langle \; (c,f,2) \; (b,e,10) \; (h,i,4) \; (e,f,30) \; (h,f,50) \\ & (e,g,40) \; (d,e,62) \; (a,d,120) \; (d,g,130) \; \rangle \end{split}$$



Every edge $e \in M$ is root of a replacement tree T_e $R(e) = \text{nodes in } T_e \text{ except for root } e$

Replacement edges have small weight

 $w(R(e)) \leq w(e)$

Replacement edges have small weight

 $w(R(e)) \leq w(e)$



Replacement edges have small weight

 $w(R(e)) \leq w(e)$



A charging scheme for M^*

M^{*} maximum weight matching

H = history edges part of the matching at some point

Charge weight of M^* to H. For each $o \in M^*$:

• $o \in H$: charge w(o) to o itself

② *o* ∉ *H*:

• $C = \text{edges conflicting with } o \text{ it was examined for insertion:} w(o) \le 2 w(C)$, since o was not inserted

• If
$$C = \{e\}$$
: charge $w(o) \leq 2w(e)$ to e

•
$$\frac{w(o)w(e_1)}{w(e_1) + w(e_2)} \le 2 w(e_1)$$
 to e_1
• $\frac{w(o)w(e'')}{w(e') + w(e'')} \le 2 w(e_2)$ to e_2

(a) Charge of $o \in M^*$ to any edge $e \in H \leq 2 w(e)$

Initial charging

$$\Sigma = \langle (c, f, 2) (b, e, 10) (h, i, 4) (e, f, 30) (h, f, 50) \\ (e, g, 40) (d, e, 62) (a, d, 120) (d, g, 130) \rangle$$
$$M^* = \{ (a, d), (e, g), (h, f) \}$$



(b) Any edge of H charged by at most two edges of M^* , one per endpoint.

Charging redistribution

If $o \in M^*$ charges $e \in H$, e replaced by $e' \in H$, e' and o incident, transfer charge of o from e to e'.



(a) Charge of $o \leq 2 w(e) \leq 2 w(e')$

(b) Any edge of H charged by at most two edges of M^* , one per endpoint (redistribution preserves incidence)

(c) Each edge $e \in H \setminus M$ charged by at most one edge in M^*

Charging redistribution

If $o \in M^*$ charges $e \in H$, e replaced by $e' \in H$, e' and o incident, transfer charge of o from e to e'.



(a) Charge of $o \leq 2 w(e) \leq 2 w(e')$

(b) Any edge of H charged by at most two edges of M^* , one per endpoint (redistribution preserves incidence)

(c) Each edge $e \in H \setminus M$ charged by at most one edge in M^*

Charging redistribution

If $o \in M^*$ charges $e \in H$, e replaced by $e' \in H$, e' and o incident, transfer charge of o from e to e'.



(a) Charge of $o \leq 2 w(e) \leq 2 w(e')$

(b) Any edge of H charged by at most two edges of M^* , one per endpoint (redistribution preserves incidence)

(c) Each edge $e \in H \setminus M$ charged by at most one edge in M^*

Analysis: summing up

- Charge of $o \in M^*$ to any edge $e \in H \leq 2 w(e)$
- Edges in $H \setminus M$ charged by at most one edge in M^*
- Edges in M charged by at most two edges in M^*

$$w(M^*) \leq \sum_{x \in H \setminus M} 2w(x) + \sum_{e \in M} 4w(e)$$

Analysis: summing up

- Charge of $o \in M^*$ to any edge $e \in H \leq 2 w(e)$
- Edges in $H \setminus M$ charged by at most one edge in M^*
- Edges in M charged by at most two edges in M^*

$$w(M^*) \leq \sum_{x \in H \setminus M} 2w(x) + \sum_{e \in M} 4w(e)$$

Since
$$H \setminus M = \bigcup_{e \in M} R(e)$$
:
 $w(M^*) \le \sum_{x \in H \setminus M} 2w(x) + \sum_{e \in M} 4w(e) = \sum_{e \in M} 2w(R(e)) + \sum_{e \in M} 4w(e)$

Analysis: summing up

- Charge of $o \in M^*$ to any edge $e \in H \leq 2 w(e)$
- Edges in $H \setminus M$ charged by at most one edge in M^*
- Edges in M charged by at most two edges in M^*

$$w(M^*) \leq \sum_{x \in H \setminus M} 2w(x) + \sum_{e \in M} 4w(e)$$

Since
$$H \setminus M = \bigcup_{e \in M} R(e)$$
:
 $w(M^*) \le \sum_{x \in H \setminus M} 2w(x) + \sum_{e \in M} 4w(e) = \sum_{e \in M} 2w(R(e)) + \sum_{e \in M} 4w(e)$

Since replacement edges have small weight $w(R(e)) \le w(e)$:

$$w(M^*) \leq \sum_{e \in M} 6w(e) = 6w(M)$$

Lower bounds

Communication complexity

Important technique for proving streaming lower bounds: reducing communication complexity problems to streaming problems

Lower bounds known in communication complexity yield streaming lower bounds

Example related to triangle counting:

To determine whether $T_3 > 0$, we need $\Omega(n^2)$ space, even using a randomized algorithm

 $T_3 =$ number of triangles

Communication complexity

2-player set-disjointness

Alice has $n \times n$ matrix ABob has $n \times n$ matrix B





Alice and Bob wish to determine if $A \cap B \neq \emptyset$ $A \cap B \neq \emptyset \iff \exists i, j : A[i, j] = 1 \text{ and } B[i, j] = 1$

By a communication complexity lower bound, this requires $\Omega(n^2)$ bits even for protocols that are correct with probability 3/4

Intro Puzzles Sampling Sketches Graphs Lower bounds

Communication complexity

Is $T_3 > 0$? Graph construction

Alice has $n \times n$ matrix ABob has $n \times n$ matrix B

 $A \cap B \neq \emptyset$?





Build graph G = (V, E) as follows:

•
$$V = \{u_1, u_2, \dots, u_n\} \cup \{v_1, v_2, \dots, v_n\} \cup \{w_1, w_2, \dots, w_n\}$$

• $E = \{(u_i, v_i) : i \in [1, n]\} \cup \{(u_i, w_j) : A[i, j] = 1\} \cup \{(v_i, w_j) : B[i, j] = 1\}$

Triangles can only have the form $\langle u_i, v_i, w_j \rangle$ *G* contains a triangle $\Leftrightarrow \exists j : A[i,j] = 1$ and B[i,j] = 1

The reduction

 $\mathcal{A}=\textit{s}\text{-bit}$ streaming algorithm that determines whether $\mathcal{T}_3>0$

Use \mathcal{A} to solve set disjointness as follows:

- Alice creates a stream with blue and red edges, and runs the algorithm on the stream
- ② Then she sends *s* bits (her memory content) to Bob
- Bob runs the algorithm, starting from Alice memory content, on the remaining yellow edges
- G He finally communicates 1 bit (the result) to Alice

Communication: s + 1 bits

 $\Rightarrow s = \Omega(n^2)$

Conclusions

More streaming algorithms...

Many others fundamentals have been studied, not covered here

- Different stream data types:
 - geometric data (location streams)
 - permutations
 - graphs and hypergraphs
- Different streaming models:
 - time-conscious models: sliding windows, exponential decay
 - non adversarial models: random order streams, skewed streams
- Different streaming scenarios:
 - distributed computations
 - sensor network computations

Directions: time-conscious models



Which is more popular between Star Wars - Episode IV (1977) and Mission Impossible -Ghost Protocol (2011)?



Are N tickets sold in each of the last 20 years better than N tickets sold in the last week?

Recent past in some cases more important than distant past \Rightarrow windowed streaming:

- fixed size window
- decaying window: influence of items on the result decreases exponentially
Directions: graphs

Rich graph structure in Web data: conversations, friendships, video, images...

Billions of dollar industry applications rely on analyzing Web info

Graph problems are very challenging:

- More dense graph problems in semi-streaming (so far, matching, spanners, shortest paths and diameter)
- Space/passes tradeoffs: reduce or annotate the stream, taking multiple passes on less and less elements
- Look at graphs as matrices: can we compute fundamental properties such as eigenvalues?
- Many natural graph questions are "hard" in standard models: more realistic and tractable models?

Directions: distributed streams

Data progressively seen from distributed sources, a central monitor (coordinator) needs to estimate some quantity

Goal: minimize total number of bits communicated by the distributed streams to the coordinator



- Can we continuously track a (global) query over streams while bounding the communication with the coordinator?
- Can we design stream summary data structures that can be combined to summarize the union of streams?

Directions: beyond adversarial order

In practice, not all frequency distributions are worst case

Can we prove stronger algorithmic results for:

- Skewed data (e.g., "Zipfian" distribution)
- Small-world scale-free models for graphs
- Random order streams
- Semi-random streams: can we develop algorithms whose performance degrades smoothly as the stream ordering becomes "less-random"?

Results in these lectures: references

- Reservoir sampling. J. S. Vitter. Random Sampling with a Reservoir, ACM Transactions on Mathematical Software, 11(1), 37-57, 1985
- *Heavy hitters.* G. S. Manku & R. Motwani. Approximate Frequency Counts over Data Streams. VLDB 2002
- Distinct items. P. Flajolet, G. N. Martin. Probabilistic Counting Algorithms for Data Base Applications. J. Comput. Syst. Sci. 1985
- Frequency moments. N. Alon, Y. Matias and M. Szegedy. The Space Complexity of Approximating the Frequency Moments. J. Comput. Syst. Sci. 1999
- Triangle counting. L. Buriol, G. Frahling, S. Leonardi, A. Marchetti-Spaccamela, & C. Sohler. Counting Triangles in Data Streams. **PODS 2006**
- Weighted matching. J. Feigenbaum, S. Kannan, A. McGregor, S. Suri, J. Zhang. On graph problems in a semi-streaming model. Theor. Comput. Sci. 2005

Online resources

Too many papers to be comprehensive... Some surveys and interesting pointers:

- Data streams: algorithms and applications, S. Muthukrishnan http://www.cs.rutgers.edu/~muthu/
- Sketch techniques for massive data, G. Cormode Continuous distributed monitoring: a short survey, G. Cormode http://dimacs.rutgers.edu/~graham/
- Algorithms for data streams, C. Demetrescu & I. Finocchi twiki.di.uniroma1.it/pub/Ing_algo/WebHome/DFchapter08.pdf
- Andrew McGregor's crash course and blog http://polylogblog.wordpress.com/2010/09/08/some-slides/
- IITK Workshop on Algorithms for Processing Massive Data Sets, IIT-Kanpur, India, 2009 http://www2.cse.iitk.ac.in/~fsttcs/2009/wapmds/
- Open problems in data streams, property testing, and related topics, Indyk et al., 2011 (the Bertinoro and Kanpur lists) http://polylogblog.wordpress.com/category/open-problems/

Thanks!

