Lecture 1

TTIC 41000: Algorithms for Massive Data Toyota Technological Institute at Chicago Spring 2021

Instructor: Sepideh Mahabadi

Welcome to Algorithms for Massive Data!

- Lectures: Monday and Wednesday 4:10-5:30pm
- Office Hours: Monday 6-7pm
- Zoom link (for both): <u>Here</u>

Make sure to check

- <u>Course page</u>
- Canvas page

Grading

- Two PSets (Each 25%)
- Final Project (40%)
- Class Participation (10%)

Problem Sets

- Free to discuss the problems
- Each person should completely understand and write their own solutions
- Write the name of the collaborators

Projects

- Can be done individually or in groups of size two
- Goal: Motivate you to write a paper in this area
- Each project includes a 10-20 minutes presentation in class, and a 5-10 pages report
- Two types of projects:
 - Summary paper
 - Research Project
- Important dates
 - April 30th: proposal
 - May 30th: first draft
 - June 10th: final draft

Introduction

Large Data Restrictions and Models

Cannot even read the whole data

Sub-linear Time Algorithms



Cannot afford storing all the data

Streaming Algorithms



Single Processor does not suffice



≻E.g. Map-Reduce algorithms

Large Data Tasks

- □ Searching through the data
- (Approximate Nearest Neighbor)



- Extracting the most relevant data (efficient summary)
- Sketching, Composable Coresets)





Large Data Techniques

□ Sampling, Importance Sampling

□ Dimensionality Reduction



This Course

Models

Tasks

Techniques

D ...

Streaming Model of Computation

Problems:

- Distinct Elements
- Morris Counter

Streaming Model

- Huge data set (does not fit into the main memory)
- Only sequential access to the data
 - One pass
 - Few passes (the data is stored somewhere else)
- Use little memory
 - Sublinear in input parameters
 - Sublinear in the input size
- Solve the problem (approximately)



Streaming Model

- Huge data set (does not fit into the main memory)
- Only sequential access to the data
 - One pass
 - Few passes (the data is stored somewhere else)
- Use little memory
 - Sublinear in input parameters
 - Sublinear in the input size
- Solve the problem (approximately)

Parameters of Interest:

- 1. Memory usage
- 2. Number of passes
- 3. Approximation Factor
- 1. (Sometimes) query/update time



Restrictions

Approximation

- Randomization
 - Analyze expectation
 - Concentration: But also this happens with a good probability Prob[being approximately correct] is high

Homework: take a look at the inequalities.

- Markov
- Chebychev
- Chernoff/Hoeffding

Probability Recap

Markov's Inequality

- Argue about concentration using the expected value
- Applies to positive random variables
- If X is a non-negative random variable, then

$$\forall \lambda > 0$$
, $\Pr(X > \lambda \cdot \mathbb{E}[X]) \leq \frac{1}{\lambda}$

The probability of exceeding expectation by more than a factor of λ is at most $\frac{1}{\lambda}$

Chebyshev's Inequality

- Argue about concentration using the variance
- Applies to r.v. with finite expected value and non-zero variance

$$\forall \lambda > 0, \quad \Pr(|\mathbf{X} - \mathbb{E}[X]| > \lambda \cdot \sigma) \le \frac{1}{\lambda^2}$$

The probability of deviating from expectation by more than $\lambda \cdot \sigma$ is at most $\frac{1}{\lambda^2}$

Or,

$$\forall \lambda > 0$$
, $\Pr(|X - \mathbb{E}[X]| > \lambda) \le \frac{\sigma^2}{\lambda^2}$

Chebyshev's Inequality

- Argue about concentration using the variance
- Applies to r.v. with finite expected value and non-zero variance

$$\forall \lambda > 0, \quad \Pr(|\mathbf{X} - \mathbb{E}[X]| > \lambda \cdot \sigma) \le \frac{1}{\lambda^2}$$

The probability of deviating from expectation by more than $\lambda \cdot \sigma$ is at most $\frac{1}{\lambda^2}$

Or,

$$\forall \lambda > 0$$
, $\Pr(|X - \mathbb{E}[X]| > \lambda) \le \frac{\mathbb{E}[(X - \mathbb{E}[X])^2]}{\lambda^2}$

More generally,

 $\forall p \ge 1, \forall \lambda > 0, \qquad \Pr(|X - \mathbb{E}[X]| > \lambda) \le \frac{\mathbb{E}[(X - \mathbb{E}[X])^p]}{\lambda^p}$

Chernoff Bound

• Suppose X_1, \dots, X_n are independent r.v. with $X_i \in [0, 1]$. Let $X = \sum_i X_i$ and $\mu = \mathbb{E}[X]$. Then

$$\begin{aligned} \forall \lambda > 0, \ \Pr(X > (1 + \lambda)\mu) < e^{-\frac{\lambda^2 \mu}{3}} & (upper \ tail) \\ \text{and,} \\ \forall \lambda > 0, \ \Pr(X < (1 - \lambda)\mu) < e^{-\frac{\lambda^2 \mu}{2}} & (lower \ tail) \end{aligned}$$

• Strong guarantee but requires the independence of r.v.

- Input: a stream of *n* numbers from [*m*]
- Goal: estimate the number of distinct elements in the stream



- Input: a stream of *n* numbers from [*m*]
- Goal: estimate the number of distinct elements in the stream

3, 6, 9, 9, 3, 4, 5, 4, 4, 5, 4, ...



- Input: a stream of *n* numbers from [*m*]
- Goal: estimate the number of distinct elements in the stream

$$m = 10$$

{3,6,9,4,5}

Sol = 5

- Input: a stream of *n* numbers from [*m*]
- Goal: estimate the number of distinct elements in the stream
- Trivial solution : min(m, n[log m])
 - Keep a counter for each of the *m* elements, space: *m* bits.
 - Keep everyone in the stream, space: $n \log m$
- We want to do much better

- Input: a stream of *n* numbers from [*m*]
- Goal: estimate the number of distinct elements in the stream
- Approximate (up to a factor of $1 + \epsilon$)
- Randomized (with probability $1-\delta$)

Vector Notation

- Receive updates to coordinates of a vector $x \in \mathbb{R}^m$
 - Coordinate *i* is denoted as *x_i*
 - Upon receiving i, we increment x_i by one
- Estimate the number of non-zero coordinates.
- L_0 -norm estimation
- *L*_p-norm estimation: Next Lecture!

$$x = \begin{bmatrix} 0, 0, 2, 4, 2, 1, 0, 0, 2, 0 \end{bmatrix}$$

$$x = \begin{bmatrix} 0 \\ m = 10 \end{bmatrix}$$

$$3, 6, 9, 9, 3, 4, 5, 4, 4, 5, 4, ...$$

- Relax the conditions: allow approximation and randomization:
- Find an estimate \widetilde{DE} of #distinct elements DE, such that

 $\Pr[\left|\widetilde{DE} - DE\right| > \epsilon DE] < \delta$

- Simpler version: Decision problem where given D, identify
 - YES case: $DE \ge D(1 + \epsilon)$ • NO case: $DE < D(1 - \epsilon)$
- Solve threshold variant, and run for $D = \{1, 1 + \epsilon, \dots, (1 + \epsilon)^i, \dots, n\}$
- > Total space will be multiplied by $\log_{1+\epsilon} n \approx \frac{\log n}{\epsilon}$
- Probability of failure multiplied by the same factor

- Relax the conditions: allow approximation and randomization:
- Find an estimate \widetilde{DE} of #distinct elements DE, such that

 $\Pr[\left|\widetilde{\mathrm{DE}} - \mathrm{DE}\right| > \epsilon \mathrm{DE}] < \delta$

- Simpler version: Decision problem where given D, identify
 - YES case: $DE \ge D(1 + \epsilon)$ • NO case: $DE < D(1 - \epsilon)$
- Solve threshold variant, and run for $D = \{1, 1 + \epsilon, \dots, (1 + \epsilon)^{\iota}, \dots, n\}$
- > Total space will be multiplied by $\log_{1+\epsilon} n \approx \frac{\log n}{\epsilon}$
- Probability of failure multiplied by the same factor

Basic Algorithm

YES case: $DE \ge D(1 + \epsilon)$ NO case: $DE < D(1 - \epsilon)$

• Sample each of *m* coordinate w.p. $\frac{1}{D}$

Basic Algorithm

- Sample each of m coordinate w.p. $\frac{1}{p}$
- If all sampled coordinates are 0, return NO
- Otherwise, return YES

Example.

•
$$D = 10$$

• $S = \{8, 10\}$

• zero



 $DE \ge D(1 + \epsilon)$

 $DE < D(1 - \epsilon)$

- D = 3• $S = \{1, 7, 9, 10\}$
 - **non-zero (** $x_9 = 2$ **)**



3, 6, 9, 9, 3, 4, 5, 4, 4, 5, 4, ...

YES case:

NO case:







Decision Variant

- Sample each of m coordinate w.p. $\frac{1}{p}$
- If all sampled coordinates are 0, return NO
- Otherwise, return YES

•
$$\Pr[NO] = \left(1 - \frac{1}{D}\right)^{DE} \approx e^{-\frac{DE}{D}}$$

- For small enough ϵ
 - If $DE > (1 + \epsilon)D$, then $e^{-(1+\epsilon)} < \frac{1}{e} \frac{\epsilon}{3}$

• If
$$DE < (1 - \epsilon)D$$
, then $e^{-(1-\epsilon)} > \frac{1}{e} + \frac{\epsilon}{3}$

• Repeat $k = O(\frac{\log \frac{1}{\delta}}{\epsilon^2})$ times • S_1, \dots, S_k

• Define
$$C \triangleq \left| \left\{ S_j \right| \sum_{i \in S_j} x_i = 0 \right\} \right|$$

• Compare *C* with *k*/*e*

Decision Variant

Boost the success probability using repetition and Chernoff

- Repeat $k = O\left(\frac{\log\left(\frac{1}{\delta}\right)}{\epsilon^2}\right)$ times and let *C* be #times we get 0
- Using Chernoff bound, with probability at least (1δ)
 - 1. If $DE > (1 + \epsilon)D$ then $C < \frac{k}{e}$ 2. If $DE < (1 - \epsilon)D$ then $C > \frac{k}{e}$

Case 1.

 $\forall \lambda > 0, \Pr(X > (1 + \lambda)\mu) < e^{-\frac{\lambda^2 \mu}{3}}$

•
$$\mathbb{E}[C] < \frac{k}{e} \left(1 - \frac{\epsilon e}{3}\right)$$
 and set $\lambda = \frac{\epsilon e}{3}$
• $\Pr\left(C > \frac{k}{e}\right) < \Pr\left(C > \mathbb{E}[C]\left(1 + \frac{\epsilon e}{3}\right)\right) < e^{-\frac{(\epsilon e/3)^2 \cdot O\left((\log\frac{1}{\delta})/\epsilon^2\right)}{3}} = \delta$

Decision Variant

Boost the success probability using repetition and Chernoff

- Repeat $k = O\left(\frac{\log\left(\frac{1}{\delta}\right)}{c^2}\right)$ times and let *C* be #times we get 0
- Using Chernoff bound, with probability at least (1δ)

 - 1. If $DE > (1 + \epsilon)D$ then $C < \frac{k}{e}$ 2. If $DE < (1 \epsilon)D$ then $C > \frac{k}{e}$
- Space usage of the algorithm: $O\left(\frac{\log\left(\frac{1}{\delta}\right)}{\epsilon^2}\right)$ numbers from [n]
- Probability of error δ



• Space usage of the algorithm: $O\left(\frac{\log n}{\epsilon} \cdot \frac{\log(\frac{1}{\delta})}{\epsilon^2}\right)$ numbers from [n](Can handle deletions too) • Probability of error $\delta \cdot \frac{\log n}{c}$ • Number of passes: 1 By **1**, DE > $(1 + \epsilon)^{j-1} \cdot (1 - \epsilon)$ By **2**, DE < $(1 + \epsilon)^j \cdot (1 + \epsilon)$ • Approximation $(1 + \epsilon)$ i = j - 1 $i = j + 1 \mid ...$ i = 1i = j $i = \log n$ i = 0YES YES YES NO NO NO

- Space usage of the algorithm: $O\left(\frac{\log n}{\epsilon} \cdot \frac{\log(\frac{1}{\delta})}{\epsilon^2}\right)$ numbers from [n] (Can handle deletions too)
- Probability of error $\delta \cdot \frac{\log n}{\epsilon}$
- Number of passes: 1



Some Comments

1. Deletions

- The algorithm can handle deletion as long as values of coordinates remain non-negative
- Dynamic or turnstile streaming
- 2. Implementing the sample set S
 - Pick a hash function $h: [m] \rightarrow [D]$ and $S = \{i | h(i) = 1\}$
 - *h* can be implemented using pseudo-random generators

State of the art

Uses space:
$$\Theta(\epsilon^{-2}\log(\frac{1}{\delta}) + \log n)$$

- Upper bound:
 - [Jaroslaw Blasiok'20]. Optimal streaming and tracking distinct elements with high probability
- Lower bound:
 - [Noga Alon, Yossi Matias, and Mario Szegedy'99]. The space complexity of approximating the frequency moments
 - [T. S. Jayram and David P. Woodruff'13]. Optimal bounds for Johnson-Lindenstrauss transforms and streaming problems with subconstant error

Morris Counter

Counting Problem

• Space efficient algorithm that monitors a sequence of events, then at any *given query time t*, returns an *estimate of the number of events* seen so far (i.e. up to time t)

update(): increment the number of events by 1
query(): output (an estimate of) the number of events

- Trivial Bound?
 - If #events is *n*, log *n* bits suffice to maintain #events exactly
 - Impossible to solve this problem exactly using o(log n)

Approximately Counting

- In most practical cases, it suffices to provide an estimate of #events
- **Goal:** Given error parameters (ϵ, δ) , return an estimate \tilde{n} of #events n, s.t.

$$\Pr[|\tilde{n} - n| > \epsilon n] < \delta$$

Morris Algorithm

1. Initialize $X \leftarrow 0$

2. Per each *update()*, increment X with probability $\frac{1}{2^X}$ \Box For *query()*, output $\tilde{n} = 2^X - 1$

Intuition: X is attempting to store a value that is roughly $\approx \log n$ The amount of space required to store X is $O(\log \log n)$

Analysis of Morris's Algorithm

Claim 1. Let X_n denote X after n updates. Then, $\mathbb{E}[2^{X_n}] = n + 1$. *Proof by induction on n.*

Base case n = 0, X = 0 trivially holds. Consider the inductive step.

$$\mathbb{E}[2^{X_{n+1}}] = \sum_{j=0}^{\infty} \Pr\left(X_n = j\right) \cdot \mathbb{E}[\left(2^{X_{n+1}} | X_n = j\right)]$$

Analysis of Morris's Algorithm

Claim 1. Let X_n denote X after n updates. Then, $\mathbb{E}[2^{X_n}] = n + 1$. *Proof by induction on n.*

Base case n = 0, X = 0 trivially holds. Consider the inductive step.

$$\mathbb{E}[2^{X_{n+1}}] = \sum_{j=0}^{\infty} \Pr\left(X_n = j\right) \cdot \mathbb{E}[\left(2^{X_{n+1}} | X_n = j\right)]$$
$$= \sum_{j=0}^{\infty} \Pr\left(X_n = j\right) \cdot \left(2^j \cdot \left(1 - \frac{1}{2^j}\right) + 2^{j+1} \cdot \frac{1}{2^j}\right) \xrightarrow{\text{Probability of}}_{\text{incrementing } X}$$

Analysis of Morris's Algorithm

Claim 1. Let X_n denote X after n updates. Then, $\mathbb{E}[2^{X_n}] = n + 1$. *Proof by induction on n.*

Base case n = 0, X = 0 trivially holds. Consider the inductive step.

$$\mathbb{E}[2^{X_{n+1}}] = \sum_{j=0}^{\infty} \Pr(X_n = j) \cdot \mathbb{E}[(2^{X_{n+1}} | X_n = j)]$$

= $\sum_{j=0}^{\infty} \Pr(X_n = j) \cdot (2^j \cdot \left(1 - \frac{1}{2^j}\right) + 2^{j+1} \cdot \frac{1}{2^j})$
= $\sum_{j=0}^{\infty} \Pr(X_n = j) \cdot 2^j + \sum_{j=0}^{\infty} \Pr(X_n = j)$
= $\mathbb{E}[2^{X_n}] + 1$
= $n + 1$

• We just showed that \tilde{n} is an unbiased estimator of n

• i.e., $\tilde{n} \stackrel{\text{\tiny def}}{=} \mathbb{E}[2^X] = n+1$

• We just showed that \tilde{n} is an unbiased estimator of n

• i.e., $\tilde{n} \stackrel{\text{\tiny def}}{=} \mathbb{E}[2^X] = n+1$

• By bounding the variance of \tilde{n} and an application of Chebyshev's, we can achieve the desired accuracy guarantee.

- We just showed that \tilde{n} is an unbiased estimator of n
 - i.e., $\tilde{n} \stackrel{\text{\tiny def}}{=} \mathbb{E}[2^X] = n + 1$
- By bounding the variance of \tilde{n} and an application of Chebyshev's, we can achieve the desired accuracy guarantee.
 - To recall, goal is to show $\Pr[|\tilde{n} n| > \epsilon n] < \delta$

$$\Pr[|\tilde{n} - n| > \epsilon n] < \frac{1}{\epsilon^2 n^2} \cdot \mathbb{E}[(\tilde{n} - n)^2] = \frac{1}{\epsilon^2 n^2} \cdot \mathbb{E}[(2^X - 1 - n)^2]$$

Claim 2. $\mathbb{E}[2^{2X_n}] = \frac{3}{2}n^2 + \frac{3}{2}n + 1.$ Proof. Similar to Claim 1 (exercise)

- We just showed that \tilde{n} is an unbiased estimator of n
 - i.e., $\tilde{n} \stackrel{\text{\tiny def}}{=} \mathbb{E}[2^X] = n + 1$

С

Ρ

- By bounding the variance of \tilde{n} and an application of Chebyshev's, we can achieve the desired accuracy guarantee.
 - To recall, goal is to show $\Pr[|\tilde{n} n| > \epsilon n] < \delta$

$$\Pr[|\tilde{n} - n| > \epsilon n] < \frac{1}{\epsilon^2 n^2} \cdot \mathbb{E}[(\tilde{n} - n)^2] = \frac{1}{\epsilon^2 n^2} \cdot \mathbb{E}[(2^X - 1 - n)^2]$$

$$|\text{aim 2. } \mathbb{E}[2^{2X_n}] = \frac{3}{2}n^2 + \frac{3}{2}n + 1.$$

$$\Pr[|\tilde{n} - n| > \epsilon n] < \frac{1}{\epsilon^2 n^2} \cdot \frac{n^2}{2} = \frac{1}{2\epsilon^2}$$

$$\Pr[|\tilde{n} - n| > \epsilon n] < \frac{1}{\epsilon^2 n^2} \cdot \frac{n^2}{2} = \frac{1}{2\epsilon^2}$$

$$\Pr[|\tilde{n} - n| > \epsilon n] < \frac{1}{\epsilon^2 n^2} \cdot \frac{n^2}{2} = \frac{1}{2\epsilon^2}$$

• Not very meaningful! RHS is better than ½ only when $\epsilon > 1$ (for which we can instead always return 0 !)

$$\Pr[|\tilde{n} - n| > \epsilon n] < \frac{1}{\epsilon^2 n^2} \cdot \frac{n^2}{2} = \frac{1}{2\epsilon^2}$$

- Not very meaningful! RHS is better than ½ only when $\epsilon > 1$ (for which we can instead always return 0 !)
- How to decrease the failure probability?

• Morris+

Average of *s* Morris estimators. Variance is multiplied by $(\frac{1}{s})$. Setting $s = \Theta(\frac{1}{\epsilon^2 \delta})$ suffices to get failure probability δ

• Morris+

Average of *s* Morris estimators. Variance is multiplied by $(\frac{1}{s})$. Setting $s = \Theta(\frac{1}{\epsilon^2 \delta})$ suffices to get failure probability δ

• Morris++

Median of t Morris+ estimators. Setting $s = \frac{1}{\epsilon^2}$, each Morris+ estimator succeeds w.p. at least $\frac{2}{3}$. By Chernoff and setting $t = \Theta(\log \frac{1}{\delta})$, the failure probability becomes at most δ

• Morris+

Average of *s* Morris estimators. Variance is multiplied by $(\frac{1}{s})$. Setting $s = \Theta(\frac{1}{\epsilon^2 \delta})$ suffices to get failure probability δ

• Morris++

Median of t Morris+ estimators.

Setting $s = \frac{1}{\epsilon^2}$, each Morris+ estimator succeeds w.p. at least $\frac{2}{3}$. By Chernoff and setting $t = \Theta(\log \frac{1}{\delta})$, the failure probability becomes at most δ **Total Space of Morris++**: $\Theta(\frac{1}{\epsilon^2} \cdot \log \frac{1}{\delta} \cdot \log \log n)$ w.p. at least $1 - \delta$