# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Fall 2022.

Lecture 7

# Quiz

- Average time spent on homework: median 15 hours, mean 23 hours.
- 8 people worked alone, 48 worked together, 48 split up questions or did some mix.

$X$ is the sum of independent random variables $X_1, \ldots, X_n$, each with mean $\mu_i$ and variance $\sigma_i$. Each $X_i$ takes on values in the range $[-5, 5]$.

Which of the following concentration bounds can you apply to show that $X$ lies close to its expectation with good probability? Check all that apply.

**Select one or more:**

- a. Markov's inquality.
- b. Chebyshev's inequality
- c. Bernstein's inequality.
- d. Chernoff bound.

Check

## Summary

### Last Class:

- Finish up exponential concentration bounds. Application to max load in hashing/load balancing.
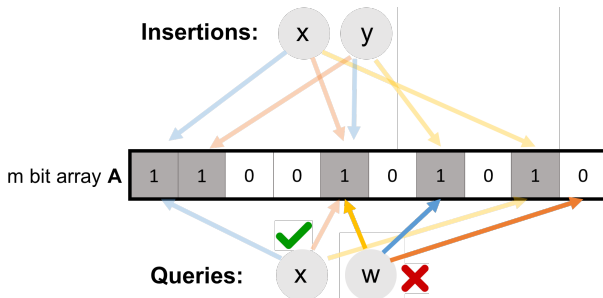- Bloom filters for storing a set with a small false positive rate.

### This Class:

- Finish up Bloom Filter Analysis.
- Start on streaming algorithms
- The distinct items problem via random hashing.

Chose $k$ independent random hash functions $\mathbf{h}_1, \ldots, \mathbf{h}_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] := 1$.
- *query(x)*: return 1 only if $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] = 1$.



No false negatives. False positives more likely with more insertions.

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0? $n \times k$ total hashes must not hit bit $i$.

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

**Step 2**: What is the probability that querying a new item $w$ gives a false positive?

$$\Pr\left(A[h_1(w)] = \ldots = A[h_k(w)] = 1\right)$$
$$= \Pr(A[h_1(w)] = 1) \times \ldots \times \Pr(A[h_k(w)] = 1)$$
$$= \left(1 - e^{-\frac{kn}{m}}\right)^k \quad \textcolor{red}{\text{Actually Incorrect!}} \text{ Dependent events.}$$

> $n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $h_1, \ldots h_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

## Correct Analysis Sketch

**Step 1**: To avoid dependence issues, condition on the event that the $A$ has $t$ zeros in it after $n$ insertions, for some $t \leq m$. For a non-inserted element $w$, after conditioning on this event we correctly have:

$$\Pr(A[\mathsf{h}_1(w)] = \ldots = A[\mathsf{h}_k(w)] = 1)$$
$$= \Pr(A[\mathsf{h}_1(w)] = 1) \times \ldots \times \Pr(A[\mathsf{h}_k(w)] = 1).$$

I.e., the events $A[\mathsf{h}_1(w)] = 1$,..., $A[\mathsf{h}_k(w)] = 1$ are independent conditioned on the number of bits set in $A$. Why?

- Conditioned on this event, for any $j$, since $\mathsf{h}_j$ is a fully random hash function, $\Pr(A[\mathsf{h}_j(w)] = 1) = 1 - \frac{t}{m}$.

- Thus conditioned on this event, the false positive rate is $\left(1 - \frac{t}{m}\right)^k$.

- It remains to show that $\frac{t}{m} \approx e^{-\frac{kn}{m}}$ with high probability. We already have that $\mathbb{E}[\frac{t}{m}] = \frac{1}{m} \sum_{i=1}^{m} \Pr(A[i] = 0) \approx e^{-\frac{kn}{m}}$.
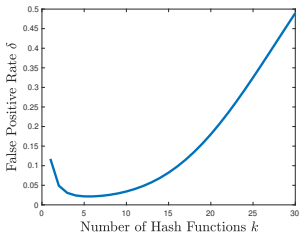
6

Need to show that the number of zeros $t$ in $A$ after $n$ insertions is bounded by $O\left(e^{-\frac{kn}{m}}\right)$ with high probability.

Can apply Theorem 2 of:
http://cglab.ca/~morin/publications/ds/bloom-submitted.pdf

**False Positive Rate:** with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$. How should we set $k$ to minimize the FPR given a fixed amount of space $m$?



- Can differentiate to show optimal number of hashes is $k = \ln 2 \cdot \frac{m}{n}$.
- Balances filling up the array vs. having enough hashes so that even when the array is pretty full, a new item is unlikely to yield a false positive.

8

# False Positive Rate

False Positive Rate: with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$.

Movies



Users

- Say we have 100 million users, each who have rated 10 movies.
- $n = 10^9 = n$ (user,movie) pairs with non-empty ratings.
- Allocate $m = 8n = 8 \times 10^9$ bits for a Bloom filter (1 GB).
- Set $k = \ln 2 \cdot \frac{m}{n} = 5.54 \approx 6$.
- False positive rate is $\approx \left(1 - e^{-k \cdot \frac{n}{m}}\right)^k \approx \frac{1}{2^k} \approx \frac{1}{2^{5.54}} = .021$.

An observation about Bloom filter space complexity:

$$\text{False Positive Rate: } \delta \approx \left(1 - e^{-\frac{kn}{m}}\right)^k.$$

For an $m$-bit bloom filter holding $n$ items, optimal number of hash functions $k$ is: $k = \ln 2 \cdot \frac{m}{n}$.

**Think Pair Share:** If we want a false positive rate $< \frac{1}{2}$ how big does $m$ need to be in comparison to $n$?

$$m = O(\log n), \ m = O(\sqrt{n}), \ m = O(n), \ m = O(n^2)?$$

If $m = \frac{n}{\ln 2}$, optimal $k = 1$, and failure rate is:

$$\delta = \left(1 - e^{-\frac{n / \ln 2}{n}}\right)^1 = \left(1 - \frac{1}{2}\right)^1 = \frac{1}{2}.$$

I.e., storing $n$ items in a bloom filter requires $O(n)$ space. So what's the point? Truly $O(n)$ bits, rather than $O(n \cdot \text{item size})$.

Questions on Bloom Filters?

## Streaming Algorithms

**Stream Processing:** Have a massive dataset $X$ with $n$ items $x_1, x_2, \ldots, x_n$ that arrive in a continuous stream. Not nearly enough space to store all the items (in a single location).

- Still want to analyze and learn from this data.
- Typically must compress the data on the fly, storing a data structure from which you can still learn useful information.
- Often the compression is randomized. E.g., bloom filters.
- Compared to traditional algorithm design, which focuses on minimizing runtime, the big question here is how much space is needed to answer queries of interest.

- **Sensor data:** images from telescopes (15 terabytes per night from the Large Synoptic Survey Telescope), readings from seismometer arrays monitoring and predicting earthquake activity, traffic cameras and travel time sensors (Smart Cities), electrical grid monitoring.



- **Internet Traffic**: 500 million Tweets per day, 5.6 billion Google searches, billions of ad-clicks and other logs from instrumented webpages, IPs routed by network switches, …

- **Datasets in Machine Learning:** When training e.g. a neural network on a large dataset (ImageNet with 14 million images), the data is typically processed in a stream due to storage

## Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream $x_1, \ldots, x_n$, estimate the number of distinct elements in the stream. E.g.,

$$1, 5, 7, 5, 2, 1 \rightarrow 4 \text{ distinct elements}$$

### Applications:

- Distinct IP addresses clicking on an ad or visiting a site.

- Distinct values in a database column (for estimating sizes of joins and group bys).

- Number of distinct search engine queries.

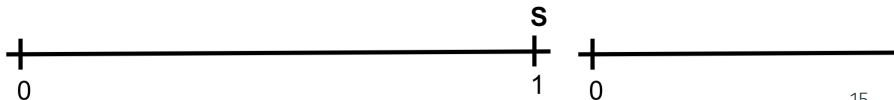- Counting distinct motifs in large DNA sequences.

Google Sawzall, Facebook Presto, Apache Drill, Twitter Algebird

## Hashing for Distinct Elements

Distinct Elements (Count-Distinct) Problem: Given a stream $x_1, \ldots, x_n$, estimate the number of distinct elements.

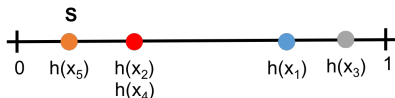Min-Hashing for Distinct Elements (variant of Flajolet-Martin):

- Let $h : U \to [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \ldots, n$
    - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$

## Hashing for Distinct Elements

Min-Hashing for Distinct Elements:

- Let $h : U \to [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \ldots, n$
  - $s := \min(s, h(x_i))$
- Return $\tilde{d} = \frac{1}{s} - 1$



- After all items are processed, $s$ is the minimum of $d$ points chosen uniformly at random on $[0, 1]$. Where $d = \#$ distinct elements.
- Intuition: The larger $d$ is, the smaller we expect $s$ to be.
- Same idea as Flajolet-Martin algorithm and HyperLogLog, except they use discrete hash functions.