

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Spring 2020.

Lecture 4

Last Class:

- **Union bound:** $\Pr(A \cup B) \leq \Pr(A) + \Pr(B)$ for *any* events A, B . Application to bounding the maximum server load when using randomized routing.
- **From Markov's inequality to exponential tail bounds.**
- Bernstein, Chernoff bounds and the Central Limit Theorem.

This Time:

- **Bloom Filters:** Random hashing to maintain a large sets in very small space.
- **Distinct Elements:** Estimating the number of unique items in a data stream via hashing. Prelude to audio fingerprinting, document search, etc.

A GENERAL COMMENT

We have covered **a lot** in the first three classes.

The proofs in class are meant to illustrate techniques that can be used to tackle many algorithmic and data related problems. You **do not** need to have these proofs or their conclusions memorized.

- Know and be comfortable applying: independence, linearity of expectation, linearity of variance, union bound, Markov's inequality, Chebyshev's inequality, basic probability calculations.
- Able to use techniques like: breaking random variables into sums of indicator variables and analyzing expected collisions (used in CATCHA analysis, two-level hashing analysis, *and* variance calculations for randomized load balancing.)
- Know definitions of 2-universal and pairwise independent hash functions and why they are useful.
- Able to apply exponential tail bounds (not have them memorized.)
- Understand law of large numbers and central limit theorem at a high level.

Bernstein Inequality (Single Variable): Consider independent random variables X_1, \dots, X_n random variable X_1 in $[-M, M]$. Let $\mu = \mathbb{E}[\sum_{i=1}^n X_i X_1]$ and $\sigma^2 = \text{Var}[\sum_{i=1}^n X_i] = \sum_{i=1}^n \text{Var}[X_i] \text{Var}[X_1]$. For any $t \geq 0$:

$$\Pr \left(\left| \sum_{i=1}^n X_i X_1 - \mu \right| \geq t \right) \leq 2 \exp \left(-\frac{t^2}{2\sigma^2 + \frac{4}{3}Mt} \right).$$

Why can't this just applied when $n = 1$?

Observation 1: Only interesting to apply when $t \leq 2M$. Why?

Observation 2: When $t \leq 2M$,

$$2 \exp \left(-\frac{t^2}{2\sigma^2 + \frac{4}{3}Mt} \right) \geq 2 \exp \left(-\frac{4M^2}{\frac{8}{3}M^2} \right) \geq 2 \exp(-3/2) \approx .44.$$

Very weak bound.

Where does sample size come in?

Bernstein Inequality: Consider independent random variables X_1, \dots, X_n in $[-M, M]$. Let $\mu = \mathbb{E}[\sum_{i=1}^n X_i]$ and $\sigma^2 = \text{Var}[\sum_{i=1}^n X_i] = \sum_{i=1}^n \text{Var}[X_i]$. For any $t \geq 0$:

$$\Pr\left(\left|\sum_{i=1}^n X_i - \mu\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2\sigma^2 + \frac{4}{3}Mt}\right).$$

Bernstein Inequality: Consider independent random variables $\frac{1}{n}X_1, \dots, \frac{1}{n}X_n$ in $[-\frac{M}{n}, \frac{M}{n}]$. Let $\mu = \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n X_i\right]$ and $\sigma^2 = \text{Var}\left[\sum_{i=1}^n \frac{1}{n}X_i\right] = \frac{1}{n^2}\sum_{i=1}^n \text{Var}[X_i]$. For any $t \geq 0$:

$$\Pr\left(\left|\frac{1}{n}\sum_{i=1}^n X_i - \mu\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2\sigma^2 + \frac{4}{3}Mt}\right)$$

We can then rework the inequality letting $\bar{M} = -M/n$ and $\bar{\sigma}^2 = \frac{\sigma^2}{n} = \frac{1}{n} \text{Var}[\sum_{i=1}^n \mathbf{X}_i]$. Giving:

Bernstein Inequality: (Sample Mean) Consider independent random variables $\mathbf{X}_1, \dots, \mathbf{X}_n$ in $[-\bar{M}, \bar{M}]$. Let $\mu = \mathbb{E}[\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i]$ and $\bar{\sigma}^2 = \frac{1}{n} \text{Var}[\sum_{i=1}^n \mathbf{X}_i] = \frac{1}{n} \sum_{i=1}^n \text{Var}[\mathbf{X}_i]$. For any $t \geq 0$:

$$\Pr \left(\left| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i - \mu \right| \geq t \right) \leq 2 \exp \left(-\frac{t^2 n}{2\bar{\sigma}^2 + \frac{4}{3}\bar{M}t} \right).$$

Bound is $< \delta$ when $n \geq 2 \log(1/\delta) \cdot \left(\frac{\bar{\sigma}^2 + \bar{M}t}{t^2} \right)$

Want to store a set S of items from a massive universe of possible items (e.g., images, text documents, IP addresses).

Goal: support $insert(x)$ to add x to the set and $query(x)$ to check if x is in the set. Both in $O(1)$ time.

- Allow small probability $\delta > 0$ of false positives. I.e., for any x ,

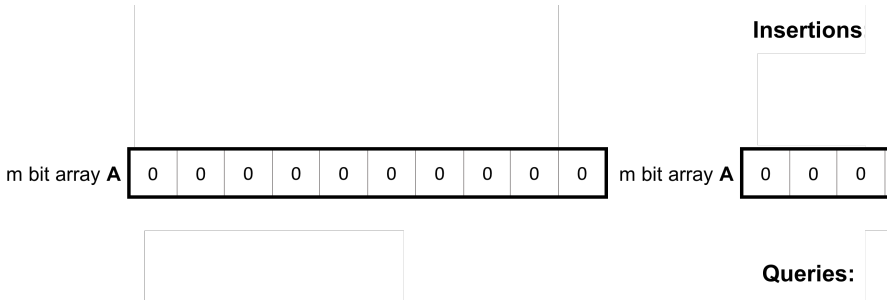
$$\Pr(query(x) = 1 \text{ and } x \notin S) \leq \delta.$$

Solution: Bloom filters (repeated random hashing).

BLOOM FILTERS

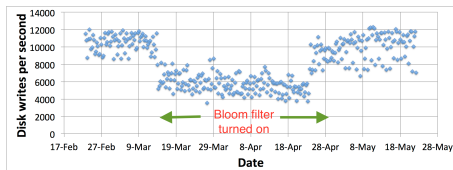
Chose k random hash functions h_1, \dots, h_k mapping the universe of elements $U \rightarrow [m]$.

- Maintain an array A containing m bits, all initially 0.
- *insert*(x): set all bits $A[h_1(x)] = \dots = A[h_k(x)] := 1$.
- *query*(x): return 1 only if $A[h_1(x)] = \dots = A[h_k(x)] = 1$.



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

Akamai (Boston-based company serving 15 – 30% of all web traffic) applies bloom filters to prevent caching of ‘one-hit-wonders’ – pages only visited once fill over 75% of cache.



- When url x comes in, if $query(x) = 1$, cache the page at x . If not, run $insert(x)$ so that if it comes in again, it will be cached.
- **False positive:** A new url (possible one-hit-wonder) is cached. If the bloom filter has a false positive rate of $\delta = .05$, the number of cached one-hit-wonders will be reduced by at least 95%.

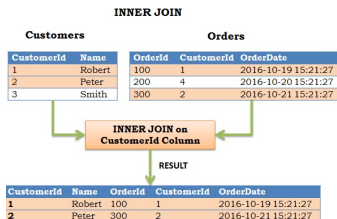
Distributed database systems, including Google Bigtable, Apache HBase, Apache Cassandra, and PostgreSQL use bloom filters to prevent expensive lookups of non-existent data.

Movies

	5			1	4				
		3						5	
Users					4				
		5							5
	1			2					

- When a new rating is inserted for $(user_x, movie_y)$, add $(user_x, movie_y)$ to a bloom filter.
- Before reading $(user_x, movie_y)$ (possibly requiring an out of memory access), check the bloom filter, which is stored in memory.
- **False positive:** A read is made to a possibly empty cell. A $\delta = .05$ false positive rate gives a 95% reduction in these empty reads.

Bloom filters are used by Oracle and other database companies to speed up database *joins*.



- Matches up a key in column **A** of one table to a key in column **B** of another, and merges corresponding information.
- A bloom filter can be used to quickly eliminate entries that appear in **A** but not in **B**.
- A false positive rate of δ means that a $1 - \delta$ fraction of these entries can be eliminated in the initial bloom filter check.

- **Recommendation systems** (Netflix, Youtube, Tinder, etc.) use bloom filters to prevent showing users the same recommendations twice.
- **Spam/Fraud Detection:**
 - Bit.ly and Google Chrome use bloom filters to quickly check if a url maps to a flagged site and prevent a user from following it.
 - Can be used to detect repeat clicks on the same ad from a single IP-address, which may be the result of fraud.
- **Digital Currency:** Some Bitcoin clients use bloom filters to quickly pare down the full transaction log to transactions involving bitcoin addresses that are relevant to them (SPV: simplified payment verification).

For a bloom filter with m bits and k hash functions, the insertion and query time is $O(k)$. How does the false positive rate δ 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 .

$$\begin{aligned}
 \Pr(A[i] = 0) &= \Pr(\mathbf{h}_1(x_1) \neq i \cap \dots \cap \mathbf{h}_k(x_k) \neq i \\
 &\quad \cap \mathbf{h}_1(x_2) \neq i \dots \cap \mathbf{h}_k(x_2) \neq i \cap \dots) \\
 &= \underbrace{\Pr(\mathbf{h}_1(x_1) \neq i) \times \dots \times \Pr(\mathbf{h}_k(x_1) \neq i) \times \Pr(\mathbf{h}_1(x_2) \neq i) \dots}_{k \cdot n \text{ events each occurring with probability } 1-1/m} \\
 &= \left(1 - \frac{1}{m}\right)^{kn}
 \end{aligned}$$

How does the false positive rate δ 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?

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

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

$$\begin{aligned} \Pr(A[\mathbf{h}_1(w)] = \dots = A[\mathbf{h}_k(w)] = 1) \\ &= \Pr(A[\mathbf{h}_1(w)] = 1) \times \dots \times \Pr(A[\mathbf{h}_k(w)] = 1) \\ &= \left(1 - e^{-\frac{kn}{m}}\right)^k \quad \text{Actually Incorrect! Dependent events.} \end{aligned}$$

n : total number items in filter, m : number of bits in filter, k : number of random hash functions, $\mathbf{h}_1, \dots, \mathbf{h}_k$: hash functions, A : bit array, δ : false positive rate.

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:

$$\begin{aligned} \Pr(A[\mathbf{h}_1(w)] = \dots = A[\mathbf{h}_k(w)] = 1) \\ = \Pr(A[\mathbf{h}_1(w)] = 1) \times \dots \times \Pr(A[\mathbf{h}_k(w)] = 1). \end{aligned}$$

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

- Conditioned on this event, for any j , since \mathbf{h}_j is a fully random hash function, $\Pr(A[\mathbf{h}_j(w)] = 1) = \frac{t}{m}$.
- Thus conditioned on this event, the false positive rate is $(1 - \frac{t}{m})^k$.
- It remains to show that $\frac{t}{m} \approx e^{-\frac{kn}{m}}$ with high probability. We already have that $\mathbb{E}[\frac{t}{m}] = \frac{1}{m} \sum_{i=1}^m \Pr(A[i] = 0) \approx e^{-\frac{kn}{m}}$.

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

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

	5			1	4						
		3						5			
Users					4						
			5							5	
	1			2							

- We have 100 million users and 10,000 movies. On average each user has rated only 10 movies so of these 10^{12} possible (user,movie) pairs, only $10 * 100,000,000 = 10^9 = n$ (user,movie) pairs have non-empty entries in our table.
- We allocate $m = 8n = 8 \times 10^9$ bits for a Bloom filter (1 GB). **How should we set k to minimize the number of false positives?**

What if we wanted to maintain a set with possible false negatives but no false positives?

Turns out that this is extremely difficult.

Questions on Bloom Filters?

Stream Processing: Have a massive dataset X with n items x_1, x_2, \dots, 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.

SOME EXAMPLES

- **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 limitations

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

1, 5, 7, 5, 2, 1 \rightarrow 4 distinct elements

Applications:

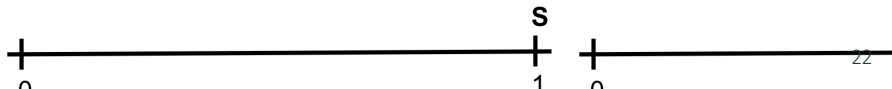
- Distinct IP addresses clicking on an ad or visiting a site.
- Distinct values in a database column (for estimating sizes of joins and group bys).
- Number of distinct search engine queries.
- Counting distinct motifs in large DNA sequences.

Google Sawzall, Facebook Presto, Apache Drill, Twitter Algebird

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

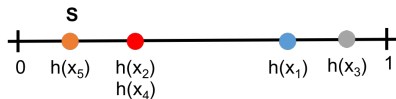
Hashing for Distinct Elements (variant of Flajolet-Martin):

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



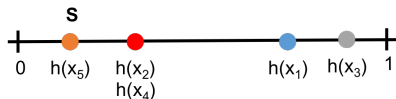
Hashing for Distinct Elements:

- Let $h : U \rightarrow [0, 1]$ be a random hash function (with a real valued output)
- $s := 1$
- For $i = 1, \dots, n$
 - $s := \min(s, h(x_i))$
- 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.

s computed by hashing algorithm is the minimum of d values chosen randomly in $[0, 1]$.



$$\mathbb{E}[s] = \frac{1}{d+1} \text{ (Interesting to prove to yourself.)}$$

- So estimate of $\tilde{d} = \frac{1}{s} - 1$ is correct if s exactly equals its expectation.
- If $|s - \mathbb{E}[s]| \leq \epsilon \cdot \mathbb{E}[s]$ for any $\epsilon \in (0, 1/2)$ can show:

$$(1 - 2\epsilon)d \leq \tilde{d} \leq (1 + 4\epsilon)d.$$

Next Time:

- Complete the analysis of hashing algorithm for distinct elements.
- Use a min-of-hashes technique for a different problem: estimating the similarity between two bit strings.
- A key idea behind audio fingerprint search (Shazam), document search (plagiarism and copyright violation detection), recommendation systems, etc.

Questions?