

# COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

---

Cameron Musco

University of Massachusetts Amherst. Fall 2021.

Lecture 7

- Problem Set 1 due tomorrow at **11:59pm**.
- My office hours are this evening at 5pm.

## Last Class:

- Bloom filters for storing a set with a small false positive rate.
- Space usage of  $O(n)$  bits vs.  $O(n \cdot \text{item size})$  for hash tables.

## This Class:

- Start on streaming algorithms
- The distinct items problem via random hashing.
- Distinct elements in practice: Flajolet-Martin and HyperLogLog.

**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

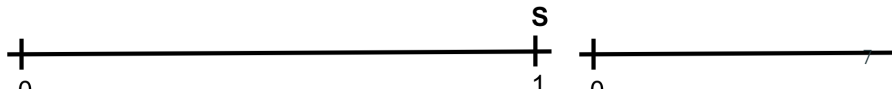
**Think Pair Share:** Discuss ways you might solve this problem without storing the full list of items seen.



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

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

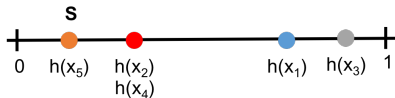
- Let  $h : U \rightarrow [0, 1]$  be a random hash function (with a real valued output)
- $s := 1$
- For  $i = 1, \dots, n$ 
  - $s := \min(s, h(x_i))$
- Return  $\tilde{d} = \frac{1}{s} - 1$





## Min-Hashing for Distinct Elements:

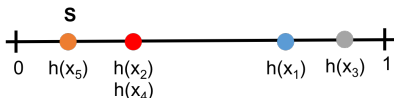
- Let  $h : U \rightarrow [0, 1]$  be a random hash function (with a real valued output)
- $s := 1$
- For  $i = 1, \dots, n$ 
  - $s := \min(s, h(x_i))$
- Return  $\tilde{d} = \frac{1}{s} - 1$



- 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.

## PERFORMANCE IN EXPECTATION

$s$  is the minimum of  $d$  points chosen uniformly at random on  $[0, 1]$ .  
Where  $d = \#$  distinct elements.



$$\mathbb{E}[s] = \frac{1}{d+1} \text{ (using } \mathbb{E}(s) = \int_0^\infty \Pr(s > x)dx \text{ + calculus)}$$

- So estimate of  $\hat{d} = \frac{1}{s} - 1$  output by the algorithm is correct if  $s$  exactly equals its expectation. Does this mean  $\mathbb{E}[\hat{d}] = d$ ? No, but:
- **Approximation is robust:** if  $|s - \mathbb{E}[s]| \leq \epsilon \cdot \mathbb{E}[s]$  for any  $\epsilon \in (0, 1/2)$  and a small constant  $c \leq 4$ :

$$(1 - c\epsilon)d \leq \hat{d} \leq (1 + c\epsilon)d$$

So question is how well  $\mathbf{s}$  concentrates around its mean.

$$\mathbb{E}[\mathbf{s}] = \frac{1}{d+1} \text{ and } \text{Var}[\mathbf{s}] \leq \frac{1}{(d+1)^2} \text{ (also via calculus).}$$

**Chebyshev's Inequality:**

$$\Pr [|\mathbf{s} - \mathbb{E}[\mathbf{s}]| \geq \epsilon \mathbb{E}[\mathbf{s}]] \leq \frac{\text{Var}[\mathbf{s}]}{(\epsilon \mathbb{E}[\mathbf{s}])^2} = \frac{1}{\epsilon^2}.$$

Bound is vacuous for any  $\epsilon < 1$ . **How can we improve accuracy?**

$\mathbf{s}$ : minimum of  $d$  distinct hashes chosen randomly over  $[0, 1]$ , computed by hashing algorithm.  $\hat{\mathbf{d}} = \frac{1}{\mathbf{s}} - 1$ : estimate of # distinct elements  $d$ .

Leverage the law of large numbers: improve accuracy via repeated independent trials.

## Hashing for Distinct Elements (Improved):

- Let  $h : U \rightarrow [0, 1]$  be a random hash function  
Let  $h_1, h_2, \dots, h_k : U \rightarrow [0, 1]$  be random hash functions
- $s := 1$
- $s_1, s_2, \dots, s_k := 1$
- For  $i = 1, \dots, n$ 
  - $s := \min(s, h(x_i))$
  - For  $j=1, \dots, k$ ,  $s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^k s_j$
- Return  $\hat{d} = \frac{1}{s} - 1$



$\mathbf{s} = \frac{1}{k} \sum_{j=1}^k \mathbf{s}_j$ . Have already shown that for  $j = 1, \dots, k$ :

$$\mathbb{E}[\mathbf{s}_j] = \frac{1}{d+1} \implies \mathbb{E}[\mathbf{s}] = \frac{1}{d+1} \text{ (linearity of expectation)}$$

$$\text{Var}[\mathbf{s}_j] \leq \frac{1}{(d+1)^2} \implies \text{Var}[\mathbf{s}] \leq \frac{1}{k \cdot (d+1)^2} \text{ (linearity of variance)}$$

**Chebyshev Inequality:**

$$\Pr[|\mathbf{s} - \mathbb{E}[\mathbf{s}]| \geq \epsilon \mathbb{E}[\mathbf{s}]] = \Pr\left[|d - \hat{d}| \geq 4\epsilon \cdot d\right] \leq \frac{\text{Var}[\mathbf{s}]}{(\epsilon \mathbb{E}[\mathbf{s}])^2} = \frac{\mathbb{E}[\mathbf{s}]^2/k}{\epsilon^2 \mathbb{E}[\mathbf{s}]^2} = \frac{1}{k \cdot \epsilon^2} = \frac{\epsilon^2}{k}$$

How should we set  $k$  if we want  $4\epsilon \cdot d$  error with probability  $\geq 1 - \delta$ ?

$$k = \frac{1}{\epsilon^2 \cdot \delta}.$$

$\mathbf{s}_j$ : minimum of  $d$  distinct hashes chosen randomly over  $[0, 1]$ .  $\mathbf{s} = \frac{1}{k} \sum_{j=1}^k \mathbf{s}_j$ .  
 $\hat{d} = \frac{1}{\mathbf{s}} - 1$ : estimate of # distinct elements  $d$ .

## Hashing for Distinct Elements:

- Let  $h_1, h_2, \dots, h_k : U \rightarrow [0, 1]$  be random hash functions
- $s_1, s_2, \dots, s_k := 1$
- For  $i = 1, \dots, n$ 
  - For  $j=1, \dots, k$ ,  $s_j := \min(s_j, h_j(x_i))$
- $s := \frac{1}{k} \sum_{j=1}^k s_j$
- Return  $\hat{d} = \frac{1}{s} - 1$



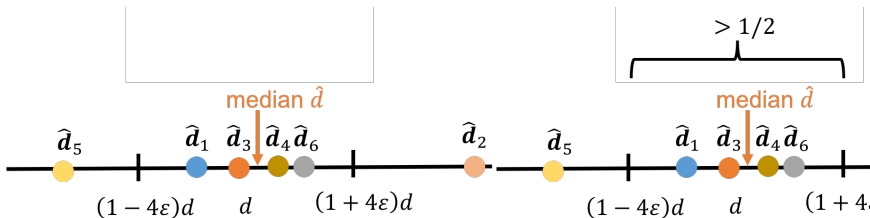
- Setting  $k = \frac{1}{\epsilon^2 \cdot \delta}$ , algorithm returns  $\hat{d}$  with  $|d - \hat{d}| \leq 4\epsilon \cdot d$  with probability at least  $1 - \delta$ .
- Space complexity is  $k = \frac{1}{\epsilon^2 \cdot \delta}$  real numbers  $s_1, \dots, s_k$ .
- $\delta = 5\%$  failure rate gives a factor 20 overhead in space complexity.

## IMPROVED FAILURE RATE

How can we improve our dependence on the failure rate  $\delta$ ?

**The median trick:** Run  $t = O(\log 1/\delta)$  trials each with failure probability  $\delta' = 1/5$  – each using  $k = \frac{1}{\delta'\epsilon^2} = \frac{5}{\epsilon^2}$  hash functions.

- Letting  $\hat{d}_1, \dots, \hat{d}_t$  be the outcomes of the  $t$  trials, return  $\hat{d} = \text{median}(\hat{d}_1, \dots, \hat{d}_t)$ .



- If  $> 1/2 > 2/3$  of trials fall in  $[(1-4\epsilon)d, (1+4\epsilon)d]$ , then the median will.
- Have  $< 1/2 < 1/3$  of trials on both the left and right.

## THE MEDIAN TRICK

- $\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t$  are the outcomes of the  $t$  trials, each falling in  $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$  with probability at least  $4/5$ .
- $\hat{\mathbf{d}} = \text{median}(\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_t)$ .

What is the probability that the median  $\hat{\mathbf{d}}$  falls in  $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ ?

- Let  $\mathbf{X}$  be the # of trials falling in  $[(1 - 4\epsilon)d, (1 + 4\epsilon)d]$ .  $\mathbb{E}[\mathbf{X}] = \frac{4}{5} \cdot t$ .

$$\Pr(\hat{\mathbf{d}} \notin [(1 - 4\epsilon)d, (1 + 4\epsilon)d]) \leq \Pr\left(\mathbf{X} < \frac{2}{3} \cdot \frac{5}{6} \cdot \mathbb{E}[\mathbf{X}]\right) \leq \Pr\left(|\mathbf{X} - \mathbb{E}[\mathbf{X}]| \geq \frac{1}{6} \mathbb{E}[\mathbf{X}]\right)$$

Apply Chernoff bound:

$$\Pr\left(|\mathbf{X} - \mathbb{E}[\mathbf{X}]| \geq \frac{1}{6} \mathbb{E}[\mathbf{X}]\right) \leq 2 \exp\left(-\frac{\frac{1^2}{6} \cdot \frac{4}{5} t}{2 + 1/6}\right) = O(e^{-ct}).$$

- Setting  $t = O(\log(1/\delta))$  gives failure probability  $e^{-\log(1/\delta)} = \delta$ .



**Upshot:** The median of  $t = O(\log(1/\delta))$  independent runs of the hashing algorithm for distinct elements returns  $\hat{d} \in [(1 - 4\epsilon)d, (1 + 4\epsilon)d]$  with probability at least  $1 - \delta$ .

**Total Space Complexity:**  $t$  trials, each using  $k = \frac{1}{\epsilon^2 \delta'}$  hash functions, for  $\delta' = 1/5$ . Space is  $\frac{5t}{\epsilon^2} = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$  real numbers (the minimum value of each hash function).

No dependence on the number of distinct elements  $d$  or the number of items in the stream  $n$ ! Both of these numbers are typically very large.

**A note on the median:** The median is often used as a robust alternative to the mean, when there are outliers (e.g., heavy tailed distributions, corrupted data).