

# Lecture 2

Barna Saha

AT&T-Labs Research

September 12, 2013

# Outline

Concentration Inequalities Revisited

Universal Family of Hash Functions

Counting Distinct Items

Analysis of Algorithm from Lecture 0

AMS Algorithm for Counting Distinct Element

# First and Second Moment Bounds

- ▶ **Markov Inequality** For any positive random variable  $X$  and  $t > 0$

$$\Pr[X > t] \leq \frac{E[X]}{t}$$

- ▶ **Chebyshev Inequality** For any random variable  $X$  and  $t > 0$

$$\Pr[|X - E[X]| > t] \leq \frac{\text{Var}[X]}{t^2}$$

# The Chernoff Bound

- ▶ Let  $X_1, X_2, \dots, X_n$  be  $n$  independent Bernoulli random variables with  $\Pr(X_i = 1) = p_i$ . Let  $X = \sum X_i$ . Hence,

$$E[X] = E\left[\sum X_i\right] = \sum E[X_i] = \sum \Pr(X_i = 1) = \sum p_i = \mu \text{ (say).}$$

Then the Chernoff Bound says for any  $\epsilon > 0$

$$\Pr(X > (1 + \epsilon)\mu) \leq \left(\frac{e^\epsilon}{(1 + \epsilon)^\epsilon}\right)^\mu \text{ and}$$

$$\Pr(X < (1 - \epsilon)\mu) \leq \left(\frac{e^{-\epsilon}}{(1 - \epsilon)^{1-\epsilon}}\right)^\mu$$

When  $0 < \epsilon < 1$  the above expression can be further simplified to

$$\Pr(X > (1 + \epsilon)\mu) \leq e^{\frac{-\mu\epsilon^2}{3}} \text{ and}$$

$$\Pr(X < (1 - \epsilon)\mu) \leq e^{\frac{-\mu\epsilon^2}{2}}$$

Hence

$$\Pr(|X - \mu| > \epsilon\mu) \leq 2e^{\frac{-\mu\epsilon^2}{3}}$$

## Universal Hash Family

A family of hash functions  $\mathcal{H} = \{h \mid h : [N] \rightarrow [M]\}$  is called a pairwise independent family of hash functions if for all  $i \neq j \in [N]$  and any  $k, l \in [M]$

$$\Pr_{h \leftarrow \mathcal{H}} [h(i) = k \wedge h(j) = l] = \frac{1}{M^2} \text{strongly universal hash family} \quad (1)$$

Hash functions are uniform over  $[M]$ ,

$$\Pr_{h \leftarrow \mathcal{H}} [h(i) = k] = \frac{1}{M} \quad (2)$$

$$\Pr_{h \leftarrow \mathcal{H}} [h(i) = h(j)] = \frac{1}{M} \text{weakly universal hash family} \quad (3)$$

- ▶ Construction Let  $p$  be a prime. For any  $a, b \in \mathbb{Z}_p = \{0, 1, 2, \dots, p-1\}$ , define  $h_{a,b} : \mathbb{Z}_p \rightarrow \mathbb{Z}_p$  by  $h_{a,b}(x) = ax + b \pmod{p}$ . Then the collection of functions  $\mathcal{H} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p\}$  is a pairwise independent hash family.

# Counting Distinct Items

---

**Algorithm 1**  $[\mathbf{a}, \epsilon, \delta]$ 

---

$$\epsilon' = \epsilon/2$$

**for**  $t = 1, \lceil(1 + \epsilon')\rceil, \lceil(1 + \epsilon')^2\rceil, \dots, \lceil(1 + \epsilon')^{\log_{1+\epsilon'} n}\rceil$  **do**

$$\delta' = \frac{\epsilon' \delta}{\log n} \text{ \{Run in parallel\}}$$

$b_t = \text{ESTIMATE}(\mathbf{a}, t, \epsilon', \delta')$   $\{b_t \text{ is a boolean variable YES/NO}\}$

**end for**

**return** the smallest value of  $t$  such that  $b_{t-1} = \text{YES}$  and  $b_t = \text{NO}$ ,  
if no such  $t$  exists, return  $n$

---

## Counting Distinct Items

---

**Algorithm 2** [ESTIMATE( $\mathbf{a}$ ,  $t$ ,  $\epsilon'$ ,  $\delta'$ )]

---

$count \leftarrow 0$

**for**  $i = 1$  to  $\frac{c}{\epsilon'^2} \log \frac{1}{\delta'}$  **do**

    Select a hash function  $h_i$  uniformly and randomly from a fully-independent hash family  $\mathcal{H}$  {run in parallel}

$b_t^i \leftarrow \text{NO}$

**repeat**

        Consider the current element in the stream  $\mathbf{a}$ , say  $a_l = (j, \nu)$

**if**  $h_i(j) == 1$  **then**

$b_t^i \leftarrow \text{YES}$ , BREAK

**end if**

**until**  $\mathbf{a}$  is exhausted

**if**  $b_t^i == \text{NO}$  **then**

$count = count + 1$

**end if**

**end for**

---

## Counting Distinct Items

---

**Algorithm 3** [ESTIMATE( $\mathbf{a}$ ,  $t$ ,  $\epsilon'$ ,  $\delta'$ )]continued

---

```
if  $count \geq \frac{1}{e} \frac{c}{\epsilon'^2} \log \frac{1}{\delta'}$  then
    return NO
else
    return YES
end if
```

---

- ▶ **Space Complexity:**  $O(\frac{1}{\epsilon^3} \log n (\log \frac{1}{\delta} + \log \log n + \log \frac{1}{\epsilon}))$
- ▶ **Time Complexity:**  $O(\frac{1}{\epsilon^3} \log n (\log \frac{1}{\delta} + \log \log n + \log \frac{1}{\epsilon}))$



# Counting Distinct Items

## ▶ Lemma

Consider the  $i$ th round of  $ESTIMATE(\mathbf{a}, t, \epsilon', \delta')$  for any  $i \in [\frac{C}{\epsilon^2} \log \frac{1}{\delta'}]$

- ▶ If  $DE > (1 + \epsilon')t$  then  $\Pr[b_t^i == NO] \leq \frac{1}{e} - \frac{\epsilon}{2e}$ .
- ▶ If  $DE < (1 - \epsilon')t$  then  $\Pr[b_t^i == NO] \geq \frac{1}{e} + \frac{\epsilon}{2e}$ .

## ▶ Lemma

- ▶ If  $DE > (1 + \epsilon')t$  then  $\Pr[b_t == NO] \leq \frac{\delta'}{2}$ .
- ▶ If  $DE < (1 - \epsilon')t$  then  $\Pr[b_t == YES] \leq \frac{\delta'}{2}$ .

## ▶ Lemma

- ▶ If  $|DE - t| > \epsilon't$  then  $\Pr[ERROR] \leq \delta'$ .

# Counting Distinct Items

- ▶ Lemma

*For all  $t$  such that  $|DE - t| > \epsilon' t$   $\Pr[ERROR] \leq \delta$ .*

- ▶ Theorem

*Algorithm 1 returns an estimate of  $DE$  within  $(1 \pm \epsilon)$  with probability  $\geq (1 - \delta)$ .*

# AMS Sketch for Counting Distinct Element

- ▶ Uses pair-wise independent hash function
- ▶ Improved space and time complexity
- ▶ Worse approximation

---

## Algorithm 4 AMS Counting Distinct Items

---

Initialize

$z \leftarrow 0$

End Initialize

Process( $a_l = (j, \nu)$ )

**if**  $\text{zeros}(h(j)) > z$  **then**

$z \leftarrow \text{zeros}(h(j))$

**end if**

End Process

Estimate

**return**  $2^{z+\frac{1}{2}}$

End Estimate

---

# AMS Sketch for Counting Distinct Element

- ▶ Define  $X_j^r = 1$  if  $\text{zeros}(h(j)) \geq r$  and 0 otherwise. Define  $Y^r = \sum_j X_j^r$ .

## ▶ Lemma

- ▶  $E[X_j^r] = \frac{1}{2^r}$
- ▶  $E[Y^r] = \frac{DE}{2^r}$
- ▶  $\text{Var}[Y^r] \leq \frac{1}{2^r}$

## ▶ Lemma

- ▶ Consider the largest level  $a$  such that  $2^{a+\frac{1}{2}} < \frac{DE}{3}$ .  
 $\Pr[z \leq a] < \frac{\sqrt{2}}{3}$ .
- ▶ Consider the smallest level  $b$  such that  $2^{b+\frac{1}{2}} > 3DE$ .  
 $\Pr[z \geq b] < \frac{\sqrt{2}}{3}$ .
- ▶  $\Pr\left[\frac{DE}{3} < 2^{z+\frac{1}{2}} < 3DE\right] \geq 1 - \frac{2\sqrt{2}}{3}$ .

# AMS Sketch for Counting Distinct Element

- ▶ **Boosting the confidence** Median Trick.

Keep  $C \log \frac{1}{\delta}$  copies and return the median estimate

- ▶ **Theorem**

*There exists a randomized algorithm that returns an estimate of  $DE$  satisfying  $\Pr\left[\frac{DE}{3} < 2^{z+\frac{1}{2}} < 3DE\right] \geq 1 - \delta$  using space  $O(\log \frac{1}{\delta} \log n)$*