

CS 331, Fall 2024
lecture 20 (11/6)

Today: - Basic boosting
- Chebyshev's inequality
- Mean/median boosting
- Morris counter

Basic boosting (Part VII, Section 3)

Today: improving the failure probability

$$\Pr[A \text{ does ...}] \leq \delta \leftarrow$$

Basic setting: you have A , succeeds

@ outputting "good" \times w.p. $\geq p$

e.g. Find pivot

$$p = \frac{1}{2}$$

Cartesian resolution

$$p = \frac{1}{n}$$

What if p small?

If we can check whether x good...

- Run A T times
- Check all of them, return any good

We'll get good x except w.p.

$$(1-p)^T$$

Key fact: $\forall p \in (0, \frac{1}{2})$,

$$\frac{1}{4} \leq (1-p)^{\frac{1}{p}} \leq \frac{1}{e}$$

Thus, $T \geq \frac{1}{p} \log(\frac{1}{\delta})$

$$\Rightarrow (1-p)^T \leq \delta$$

What if we can't verify??

Motivation: running experiment to estimate x^*

designer algo to output x , $E[x] = x^*$

- but...
- We care about $|x - x^*|$
 - Can't check (don't know x^*)

Today: how to prove things like

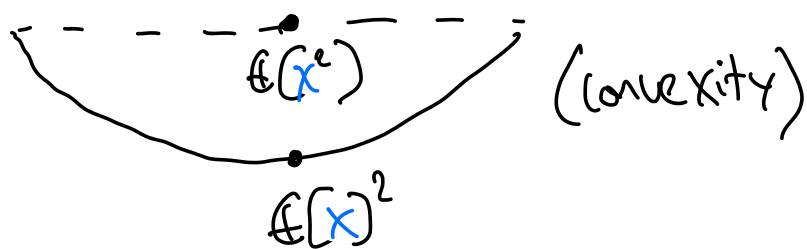
$$\Pr \left[x \notin \underbrace{[x^* - R, x^* + R]}_{\text{"confidence interval"}} \right] \leq \delta$$

- Mean boosting: improve R
- Median boosting: improve δ

Chebyshev's inequality (Part VIII, Section 3.1)

Key players:

$$\text{Var}[X] := E[X^2] - E[X]^2 \geq 0$$



$$\text{Stdev}[X] := \sqrt{\text{Var}[X]} \quad \text{"spread"}$$

Aside Variance cash course

Reminder about E : (no cheats!!!)

$$\begin{aligned} \textcircled{1} E[cX] \\ = cE[X] \end{aligned}$$

$$\begin{aligned} \textcircled{2} E[X + Y] \\ = E[X] + E[Y] \end{aligned}$$

$$\begin{aligned}
E[(X - E[X])^2] &= E[X^2] \\
&- 2E[X \cdot \underbrace{E[X]}_{\text{constant, use ①}}] + E[X]^2 \\
&= E[X^2] - E[X]^2 = \text{Var}[X]
\end{aligned}$$

Also, we have

$$\text{① } \text{Var}[cX] = c^2 \text{Var}[X]$$

$$\text{② } \text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$$

only if X, Y independent!

(Compare to E , ① & ② both always true)

Our intuition: for reasonable δ ,
pick $R \approx \text{stddev}(X)$

Markov's inequality

Let r.v. $X \geq 0$, then $\forall \delta \in (0, 1)$,

$$\Pr\left[X \geq \frac{E[X]}{\delta}\right] \leq \delta$$

e.g. $X < 10 E[X]$ except w.p. $\frac{1}{10}$

Proof: if not, let $\tau = \frac{E[X]}{\delta}$,

$$\begin{aligned} E[X] &= \Pr[X \geq \tau] \cdot \tau + \Pr[X < \tau] \cdot 0 \\ &> \delta \cdot \tau > E[X] \end{aligned}$$

Chebyshev's inequality

Apply Markov's with $X \leftarrow (X - E[X])^2$

$$\Pr \left[(X - E[X])^2 \geq \frac{\text{Var}[X]}{\delta} \right] \leq \delta$$

$$\Leftrightarrow |X - E[X]| \geq \frac{\text{stddev}[X]}{\sqrt{\delta}}$$

Thus we have shown except w.p. δ ,

$$X \in \left[E[X] - \frac{\text{stddev}[X]}{\sqrt{\delta}}, E[X] + \frac{\text{stddev}[X]}{\sqrt{\delta}} \right]$$

e.g. $\delta = \frac{1}{4}$, confidence interval

$$\left[E[X] - \underbrace{2\text{stddev}[X]}_{\text{radius } R}, E[X] + 2\text{stddev}[X] \right]$$

Mean / median boosting (Part VIII, Section 3.2)

Idea (: Improve R (mean boosting)

Recall that $R \propto \text{stdv}(X)$

How to halve R ? Decrease Var !

Basic fact:

Let X_1, X_2, \dots, X_k independent copies of X ,

$$\bar{X} = \frac{1}{k} \sum_{i \in [k]} X_i$$

$$\text{Var}(\bar{X}) = \frac{1}{k^2} \text{Var}\left(\sum_{i \in [k]} X_i\right)$$

$$= \frac{k}{k^2} \text{Var}(X) = \frac{1}{k} \text{Var}(X)$$

Takeaway: if $k \approx \frac{1}{\epsilon^2}$,

then our confidence gets ϵ^x better

$$\text{Var}(\bar{X}) = \frac{1}{k} \text{Var}(X), \quad \sqrt{\frac{1}{k}} \times \text{for stdev.}$$

Idea 2: Improve δ (median boosting)

Suppose r.v. X has

$$\Pr(X \in I) \geq \frac{3}{4} \quad \text{for interval } I \subset \mathbb{R}$$

Claim: let $\hat{X} = \text{median}(X_1, X_2, \dots, X_k)$

$$k \geq 12 \log\left(\frac{1}{\delta}\right)$$

Then, $\hat{X} \in I$ except w.p. δ

How to prove?

Step 1: it's enough to show $\geq \frac{k}{2}$ copies $\in I$

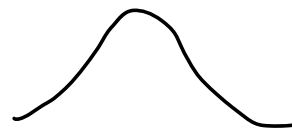
e.g. if $\hat{X} \notin I$ to the right, clearly $< \frac{k}{2} \in I$



Our goal: we have coin that H w.p. $\geq \frac{3}{4}$
toss k coins, $\geq \frac{k}{2}$ are H.

Step 2: "binomial concentration"

Suppose $i \geq \frac{k}{2}$ coins miss I



$$\begin{aligned} \Pr(\geq \frac{k}{2} \text{ are } T) &= \sum_{i=\frac{k}{2}}^k \binom{k}{i} (1-p)^i p^{k-i} \\ &\leq \sum_{i=\frac{k}{2}}^k \binom{k}{i} \left(\frac{1}{4}\right)^i \left(\frac{3}{4}\right)^{k-i} \end{aligned}$$

Observe if $i \leftarrow i+1$, $\binom{k}{i} \left(\frac{1}{4}\right)^i \left(\frac{3}{4}\right)^{k-i}$
 decreases decreases $3x$

Thus geom sequence. Enough to get first term

$$\binom{k}{k/2} \left(\frac{1}{4}\right)^{k/2} \left(\frac{3}{4}\right)^{k/2} = O(\delta)$$

$$\leq 2^k \cdot \left(\frac{3}{16}\right)^{k/2} = \left(\frac{3}{4}\right)^{k/2}$$

Good enough if $k = O\left(\log\left(\frac{1}{\delta}\right)\right)!$

- 3-stage plan:
- Basic interval via $\sqrt{\text{Var}(X)}$
 $R = 2 \text{stdev}(X)$, $\delta = 1/4$
 - Mean boost $\times \frac{4}{\epsilon^2}$
 $R = \epsilon \text{stdev}(X)$, $\delta = 1/4$
 - Median boost $\times 12 \log\left(\frac{1}{\delta}\right)$
 $R = \epsilon \text{stdev}(X)$, $\delta = \delta$

Morris Counter (Part VIII, Section 3.3)

Goal: build data structure to store counter $i \in \mathbb{O}$

- Count(): $i \leftarrow i+1$
- Report(): return i

How much space needed to Count() n times?

Trivial: $O(\log(n))$ space (store i)

Today: $O(\log \log(n))$ space (randomized)

Motivation:

- Count many large things
- Constant factors help a lot!
- web crawler, search engine, etc.
- Morris: spellchecker, needed 26^3 trigram counters

Algo: $X \leftarrow 0,$

Count(): w.p. $2^{-X}, X \leftarrow X+1$

Report(): return $2^X - 1$

It works ????

Intuition: $X \approx \log(i+1)$

Increase once in each bucket

$i \in (3, 4), [5, 8), [9, 16), \dots$

Let $X_k =$ value of X after $i = k$

Claim: $E[2^{X_k}] = k+1$ $\forall k \in (n)$

$$\text{Var}[2^{X_k}] \leq 2k^2$$

3-step plan: gives estimate in

$$[(1-\varepsilon)n, (1+\varepsilon)n] \text{ w.p. } \geq 1-\delta$$

$$(NY22): \text{ space } O\left(\log \log \frac{n}{\delta} + \log \frac{1}{\varepsilon}\right)$$

Proof of \mathbb{E} : true when $k=0$. Induct!

$$\mathbb{E}\left(2^{X_{k+1}}\right) = 2^j + 1$$

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

$$= \sum_{j=0}^{\infty} \left(\Pr[X_k = j] 2^j + \Pr[X_k = j]\right)$$

$$= \mathbb{E}\left(2^{X_k}\right) + 1. \quad \square \quad (\text{Var: see notes})$$