

Chernoff bounds: Let X_1, \dots, X_n be indep. r.v.'s

where $0 \leq X_i \leq 1$. Let $X = \sum_{i=1}^n X_i$ & $\mu = E[X]$.

For $0 \leq \delta \leq 1$,

$$\Pr(X \geq (1+\delta)\mu) \leq e^{-\delta^2 \mu / 3}$$

$$\Pr(X \leq (1-\delta)\mu) \leq e^{-\delta^2 \mu / 2}$$

Warm-up example: Median estimate

Let $S = [X_1, \dots, X_m]$

for simplicity assume X_i 's are distinct
& let $\text{rank}(y) = |\{x \in S : x \leq y\}|$

Goal: Find ϵ -approx. median(S),

more precisely, given $\epsilon > 0$, find y where:

$$\frac{m}{2} - \epsilon m = \frac{m}{2} (1 - 2\epsilon) < \text{rank}(y) < \frac{m}{2} (1 + 2\epsilon) = \frac{m}{2} + \epsilon m$$

Alg.: Choose t random elements from S .
 Let $R = [r_1, \dots, r_t]$
 Return median(R).

Claim: if $t \geq \frac{2}{\epsilon^2} \log(\frac{2}{\delta})$, then the alg. returns an ϵ -approx. median with prob. $\geq 1 - \delta$.
 This is an (ϵ, δ) -approximation of the median.

Proof: Let $S_L = \{x \in S : \text{rank}(x) \leq \frac{m}{2} - \epsilon m\}$
 $S_M = \{x \in S : \frac{m}{2} - \epsilon m \leq \text{rank}(x) < \frac{m}{2} + \epsilon m\}$
 ~~$S_U = \{x \in S : \text{rank}(x) \geq \frac{m}{2} + \epsilon m\}$~~
 $S_U = \{x \in S : \text{rank}(x) \geq \frac{m}{2} + \epsilon m\}$

If $< \frac{1}{2}$ elements of S_L & $< \frac{1}{2}$ of S_U
 are in ~~S_M~~ R then the median of R is in S_M ,
 and it's an ϵ -approx median.

Let $X_i = \begin{cases} 1 & \text{if } r_i \in S_L \\ 0 & \text{o/w} \end{cases}$

$$\& X = \sum_{i=1}^m X_i,$$

$$\text{Then, } E[X] = t \left(\frac{\frac{m}{2} - \epsilon m}{m} \right) = \frac{t}{2} - \epsilon t$$

$$\Pr\left(X \geq \frac{t}{2}\right) \leq \Pr\left(X \geq E[X](1 + 2\epsilon)\right)$$

$$\leq e^{-4\epsilon^2 \left(\frac{t}{2} - \epsilon t\right)^2 / 3}$$

$$\leq e^{-\frac{4\epsilon^2}{3} t} \leq \frac{\delta}{2} \text{ for } t \geq \frac{3}{4\epsilon^2} \log\left(\frac{2}{\delta}\right).$$

& $\Pr\left(\geq \frac{t}{2} \text{ elements from } S_U \text{ in } R\right) \leq \frac{\delta}{2}$ by analogous argument.

Hence, $\Pr(\text{alg. outputs } \epsilon\text{-approx. median}) \geq 1 - \delta.$ \square

(4)

Streaming: we get one-by-one m elements

$$x_1, x_2, \dots, x_m$$

where $x_i \in \{1, 2, \dots, n\}$.

Think of m as huge so can't store entire stream.

Let f_i be the frequency of number i in ^{the} stream.

$$\& f = (f_1, \dots, f_n).$$

Reservoir Sampling:

We don't know m beforehand, but we want to choose an element s uniformly at random from X_1, \dots, X_m .

Alg.:

Set $s = X_1$

for $t > 1$:

- upon seeing t^{th} element X_t ,
with prob. $\frac{1}{t}$ set $s = X_t$

What's Prob. that $s = X_i$ for some time $t \geq i$?

$$\begin{aligned}
 \Pr(s = X_i) &= \frac{1}{i} \times \left(1 - \frac{1}{i+1}\right) \times \left(1 - \frac{1}{i+2}\right) \times \dots \times \left(1 - \frac{1}{t}\right) \\
 &= \frac{1}{i} \times \left(\frac{i}{i+1}\right) \times \left(\frac{i+1}{i+2}\right) \times \dots \times \left(\frac{t-1}{t}\right) \\
 &= \frac{1}{t}
 \end{aligned}$$

— Takes $O(\log n)$ bits of space to get s
& $O(k \log n)$ bits to get k samples.

AMS = [Alon, Matias, Szegedy '99]

Have function g where $g(0) = 0$.

Goal: Compute $\sum_{i=1}^n g(f_i)$

Unbiased estimator: r.v. X where $E[X] = \sum_{i=1}^n g(f_i)$

Alg.: First choose random ~~sample~~ ^{index} $J \in \{1, \dots, m\}$ → do this using reservoir sampling scheme
& compute ~~r~~

$$r = |\{j \geq J : x_j = x_J\}|$$

= # of occurrences of x_J after time J .

Then, output $X = m \cdot (g(r) - g(r-1))$

Claim: $E[X] = \sum_{i=1}^n g(f_i)$

Proof: $E[X] = \sum_i \Pr(x_J = i) E[X | x_J = i] = \sum_i \frac{f_i}{m} \sum_{r=1}^{f_i} \frac{m(g(r) - g(r-1))}{f_i}$

f_i elements = i
 ~~f_i~~ elements of value i afterwards

How to ensure the output is (ϵ, δ) -approx?
idea: Take t estimators & average?

Example: Frequency moments

For integer $k \geq 1$, let $F_k = \sum_{i=1}^n f_i^k$

$$\text{So } g(r) = r^k$$

Run AMS alg. & we have $X = m(r^k - (r-1)^k)$

$$\text{We know: } E[X] = F_k$$

Idea: Do t indep't. trials to get X_1, \dots, X_t
& output $\frac{1}{t} \sum_{i=1}^t X_i$

Hope to use Chernoff to show it's $(1 \pm \delta)$ approx.

But: X_i can be huge.

So show $\text{Var}(X)$ is small & apply Chebyshev's.

Lemma: $\text{Var}(X) \leq k n^{1-\frac{1}{k}} F_k^2$

$$\text{For } \delta = \frac{3 \text{Var}(X)}{\epsilon^2 E[X]^2} = \frac{3 k n^{1-\frac{1}{k}} F_k^2}{\epsilon^2 F_k^2} = 3 k n^{1-\frac{1}{k}} \epsilon^{-2},$$

let $Y = \frac{1}{t} \sum_{i=1}^t X_i = \text{mean of } t \text{ trials.}$

$$E[Y] = E[X_i] = F_k$$

$$\text{Var}(Y) = \frac{1}{l^2} \sum_{i=1}^l \text{Var}(X_i) = \frac{\text{Var}(X)}{l}$$

By Chebyshev's ineq.,

$$\Pr(|Y - E[Y]| \geq \epsilon E[Y])$$

$$= \Pr(|Y - F_k| \geq \epsilon F_k)$$

$$\leq \frac{\text{Var}(Y)}{(\epsilon F_k)^2} \leq \frac{\text{Var}(X)}{l \epsilon^2 F_k^2} \leq \frac{kn^{1-\frac{1}{2}} F_k^2}{3kn^{1-\frac{1}{2}} \epsilon^2 F_k^2} = \frac{1}{3}$$

So with prob. $\geq \frac{2}{3}$ it's an ϵ -approx. of F_k

How to boost this prob. to $1 - \delta$?

Repeat t times & take the median.

Now we do this $t = \frac{1}{\epsilon} \log(1/\delta)$ times, call them Y_1, \dots, Y_t

Let $Z_j = \begin{cases} 1 & \text{if } |Y_j - F_k| \leq \epsilon F_k \text{ so it's an } \epsilon\text{-approx} \\ 0 & \text{o/w} \end{cases}$

Then, for $Z = \sum Z_j$, $E[Z] \geq \frac{2}{3}t$

If $Z \geq \frac{t}{2}$ then the median is an ϵ -approx.

$$\Pr(Z < \frac{t}{2}) \leq \Pr(Z \leq E[Z] - \frac{t}{6}) \leq e^{-\frac{2t}{3} \frac{1}{36\epsilon^2}} \leq \delta \text{ for big enough } t.$$

Proof of Lemma:

$$\begin{aligned} \text{Var}(x) &\leq E[x^2] \\ &= \sum_{j=1}^n \frac{f_j}{m} \sum_{r=1}^{f_j} \frac{1}{f_j} m^2 (r^k - (r-1)^k)^2 \\ &= m \sum_j \sum_r (r^k - (r-1)^k)^2 \end{aligned}$$

see claim on next page →

$$\leq m \sum_j \sum_r k r^{k-1} (r^k - (r-1)^{k-1})$$

$$\leq m \sum_j k f_j^{k-1} \sum_r (r^k - (r-1)^{k-1})$$

$$= m \sum_j k f_j^{k-1} f_j^k$$

$$= k F_1 F_{2k-1}$$

$$\leq k n^{1-\frac{1}{k}} F_k^2$$

See claim 2 on Page 11

Claim: $r^k - (r-1)^k \leq kr^{k-1}$

Mean-value theorem:

for function f which is continuous on $[a, b]$
& differentiable on (a, b)

then $\exists c \in (a, b)$ where:

$$f'(c) = \frac{f(b) - f(a)}{b - a}$$

Thus, $\exists l \in [r-1, r]$ where: for $f(z) = z^k$,

$$f'(l) = \frac{f(r) - f(r-1)}{1}$$

$$kl^{k-1} = r^k - (r-1)^k$$

Therefore,

$$r^k - (r-1)^k = kl^{k-1} \leq kr^{k-1}$$

□

$$F_1 F_{2k-1} = \left(\sum_i f_i \right) \left(\sum_i f_i^{2k-1} \right)$$

$$\leq \left(\sum_i f_i \right) \left(f_*^{k-1} \sum_i f_i^k \right) \quad \text{where } f_* = \max_i f_i$$

$$\leq \left(\sum_i f_i \right) \left(f_*^k \right)^{k-1/2} \left(\sum_i f_i^k \right)$$

$$\leq \left(\sum_i f_i \right) \left(\sum_i f_i^k \right)^{k-1/2} \left(\sum_i f_i^k \right)$$

$$\stackrel{?}{\leq} n^{1-1/2} \left(\sum_i f_i^k \right)^{1/2} \left(\sum_i f_i^k \right)^{k-1/2} \left(\sum_i f_i^k \right)$$

$$= n^{1-1/2} F_k^2$$

$$\frac{1}{n} \sum_i x_i \leq \left(\frac{1}{n} \sum_i x_i^k \right)^{1/k}$$

AM \leq k^{th} power mean

Power-mean inequality.