

Chebyshev's Inequality

Let X be a random variable

For any $a > 0$,

$$P(|X - E[X]| \geq a) \leq \frac{\text{Var}[X]}{a^2} = \frac{E[X^2] - E[X]^2}{a^2}$$

probability that X deviates from the expected value of X

Proof:

$$P(|X - E[X]| \geq a) \leq P((X - E[X])^2 \geq a^2) \stackrel{\text{Markov}}{\leq} \frac{E[(X - E[X])^2]}{a^2}$$

$$\begin{aligned} E[(X - E[X])^2] &= E[X^2 - 2 \cdot X \cdot E[X] + E[X]^2] \stackrel{\text{linearity of expectation}}{=} E[X^2] - 2 \cdot E[X] \cdot E[X] + E[E[X]^2] \\ &= E[X^2] - 2E[X]^2 + E[X]^2 = E[X^2] - E[X]^2 \\ &= \text{Variance of } X \end{aligned}$$

$E[2 \cdot X \cdot E[X]] = 2 \cdot E[X] \cdot E[X]$
 $E[E[X]^2] = E[X]^2$ (number)

$E[3X] = 3E[X]$

$$P(|X - E[X]| \geq a) \leq \frac{E[X^2] - E[X]^2}{a^2}$$

$$P(|X - E[X]| \geq a) \leq \frac{E[X^2] - E[X]^2}{a^2} = \frac{\text{Var}(X)}{a^2}$$

Chebyshev's Inequality.

This is ~~the~~ ^{very} useful definition when we have independence.

$X_1, X_2, X_3, \dots, X_N$ are **iid RV** - independent identically distributed Random Variables.

Define $X = X_1 + X_2 + \dots + X_N$

$$\text{Var}(X) = \text{Var}\left(\sum_{i=1}^N X_i\right) = \sum_{i=1}^N \text{Var}(X_i) \quad \text{when you have independence!}$$

Sort of like linearity of expectation, except now you need independence.

- Is Chebyshev's inequality is a stronger bound than Markov's?
- If you remember how we proved Chebyshev's inequality - we just applied Markov.

ex. Coin flips. Flip a fair coin N times.

$$X = \# \text{ of Heads} = \sum_{i=1}^N X_i, \quad X_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ flip is Heads} \\ 0, & \text{otherwise} \end{cases}$$

With linearity of expectation we showed that $E[X] = N/2$

With Markov's inequality we showed: $P(X \geq \frac{3}{4}N) \leq \frac{N/2}{(3/4)N} = \frac{2}{3}$

Let's see what Chebyshev's can offer: ← positive random variables, sum of

$$\text{Var}(X_i) = E[X_i^2] - E[X_i]^2 = \left(\frac{1}{2}\right) - \frac{1}{4} = \frac{1}{4}$$

because X_i can be 1 or 0, so X_i^2 can be $1^2=1$ or $0^2=0$.

$$E[X_i^2] = 1 \cdot \frac{1}{2} + 0 \cdot \frac{1}{2} = \frac{1}{2}$$

$$\text{Var}(X) = \sum_{i=1}^N \text{Var}(X_i) = \frac{N}{4} \quad \text{since each coin flip is independent.}$$

If A then B
then: $P_r(A) \leq P_r(B)$

$$P(X \geq \frac{3}{4}N) = P(X \geq \underbrace{E[X]}_{N/2} + \frac{1}{4}N) = P(X - E[X] \geq \frac{N}{4}) \leq$$

if we take absolute value of left-hand side then the probability gets bigger.

$$\leq P(|X - E[X]| \geq \frac{N}{4}) \leq \frac{\text{Var}(X)}{(\frac{N}{4})^2} = \frac{\frac{N}{4}}{(\frac{N}{4})^2} = \frac{4}{N}$$

Apply Chebyshev's Inequality

← That is much smaller, than we can show with Markov.

So, Chebyshev's Inequality is much stronger!!!

Randomized Median Finding using Random Sampling.

Median of N numbers is a number with rank $N/2$.

half the numbers are smaller than the Median, and half bigger.

Quick Select ^{← randomized} solves this in $O(N)$ expected. (with a worst case $O(N^2)$)
 (like Median of Medians) $\leftarrow T(n) \leq 10 \cdot c \cdot N$

There are deterministic algorithms to solve this in $O(N)$ time. (but they are very complicated)

• Input: Set S of N elements.

• Output: Median
There is nothing special about $N^{3/4}$, we just want things to cancel nicely. We want N in some power smaller than 1.

$O(N^{3/4})$ 1) Pick $N^{3/4}$ elements at random from S with replacement. Call this set R . (I pick any element with probability $1/N$), but don't remove it from S
[We want R to sort of look like S , to represent S .]
but also to be small.

$O(N)$ 2) Sort R in $O(N^{3/4} \cdot \log N)$ time which is $o(N)$
 $\log N^c = c \cdot \log N$
const

$O(1)$ 3) Let $d =$ element of rank $\frac{N^{3/4}}{2} - \sqrt{N}$ in R .
We want true median lie in between. We want u and d to be close together - to have C small, But also far apart enough so the will contain the median with high probability

$O(1)$ 4) Let $u =$ element of rank $\frac{N^{3/4}}{2} + \sqrt{N}$ in R .

5) Let $C = \{x \in S \mid d \leq x \leq u\}$

$l_d = |\{x \in S \mid x < d\}|$

$l_u = |\{x \in S \mid x > u\}|$

\rightleftarrows cardinality of a set

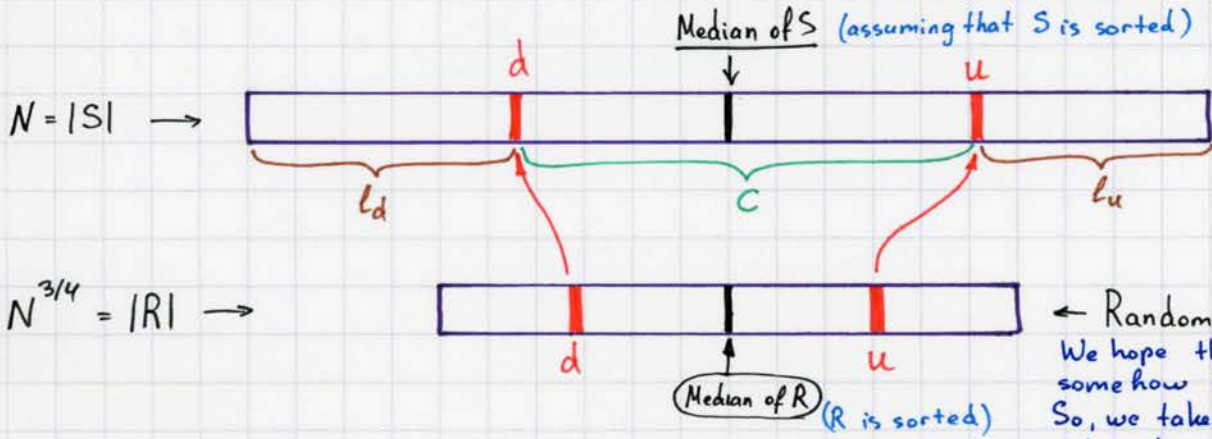
$O(n)$
 every other step is sublinear

6) If $l_d > \frac{N}{2}$ or $l_u > \frac{N}{2}$ FAIL (because the Median is not in C.)

7) If $|C| > 4N^{3/4}$ then FAIL

8) Sort C and output element of rank: $\frac{N}{2} - l_d + 1$

Compare all the items of S to d and u and discard them. Sort elements that are left: C and output the median, because we know how many elements before d and after u.



← Random sample of S
We hope that median of R is somehow close to median of S. So, we take an interval around u of R and we hope that it will give an interval around median of S.

If algorithm doesn't fail it runs linear time.

- What are the bad things that make algorithm fail?

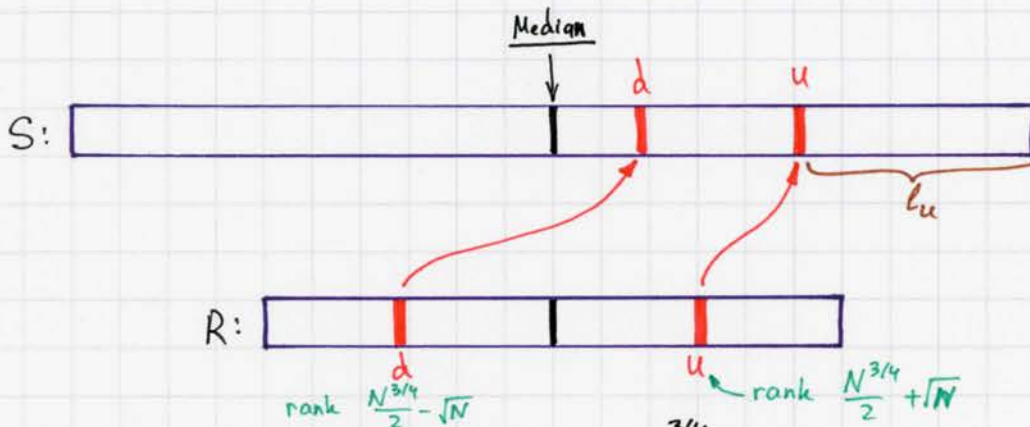
Bad events: $E_1: l_d > N/2$

$E_2: l_u > N/2$

$E_3: |C| > 4N^{3/4}$

E_1

• If $l_d > \frac{N}{2}$ then Median of S is less than d:



How much distance is expected in S between two consecutive elements in R? $\rightarrow N^{1/4}$

We give a really big window of $2 \cdot \sqrt{N}$
The window in S is expected to be $2 \cdot N^{1/2} \cdot N^{1/4} = 2N^{3/4}$

$$y_1 = |\{r \in R \mid r \leq \text{Median of } S\}| < \frac{N^{3/4}}{2} - \sqrt{N}$$

$$P(l_d > \frac{N}{2}) \leq P(y_1 < \frac{N^{3/4}}{2} - \sqrt{N})$$

If this is true then this is true.

We want to compute $P(Y_2 < \frac{N^{3/4}}{2} - \sqrt{N})$. Let's say we want to apply Chebyshev. Random Sampling

Compute $E[Y_2]$ and $\text{Var}(Y_2) = E[Y_2^2] - E[Y_2]^2$.

What is the operation I used to build R ?
We define an indicator random variable that looks at each of those selections.

$X_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ sample placed in } R \text{ is } \leq \text{Median} \\ 0, & \text{otherwise} \end{cases}$ [or 1, if i^{th} sample falls in Y_2]

$$Y_2 = \sum_{i=1}^{N^{3/4}} X_i$$

$E[X_i] = \frac{1}{2}$ (half elements is smaller than Med. and half (bigger))

$$E[Y_2] = E\left(\sum_{i=1}^{N^{3/4}} X_i\right) = \sum_{i=1}^{N^{3/4}} E[X_i] = \sum_{i=1}^{N^{3/4}} \frac{1}{2} = \frac{N^{3/4}}{2}$$

$$\text{Var}(X_i) = E[X_i^2] - E[X_i]^2 = \frac{1}{2} - \frac{1}{4} = \frac{1}{4}$$

That's why we did random sampling with replacement.

For indicator R.V.s $E[X_i^2] = E[X_i]$

because all variables are independent by construction

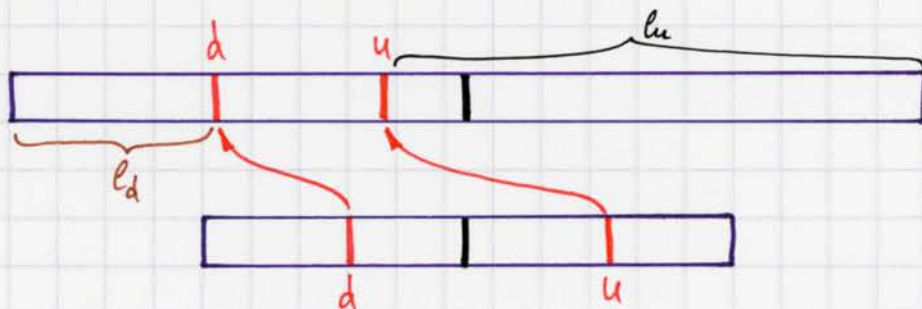
$$\text{Var}(Y_2) = \sum_{i=1}^{N^{3/4}} \text{Var}(X_i) = \sum_{i=1}^{N^{3/4}} \frac{1}{4} = \frac{N^{3/4}}{4}$$

$$P\left(Y_2 < \underbrace{\frac{N^{3/4}}{2}}_{E[Y_2]} - \sqrt{N}\right) \leq P(|Y_2 - E[Y_2]| > \sqrt{N}) \stackrel{\text{Chebyshev}}{\leq} \frac{\text{Var}(Y_2)}{(\sqrt{N})^2} \leq \frac{N^{3/4} \cdot \frac{1}{4}}{N} = \frac{1}{4N^{1/4}}$$

small

E_2

• If $u > \frac{N}{2}$ then Median of S is bigger than u :



$$Y_2 = |\{r \in R \mid r \geq \text{Median of } S\}|$$

In the exact way as with E_1 case we can show that

$$P\left(Y_2 < \frac{N^{3/4}}{2} - \sqrt{N}\right) \leq \frac{1}{4N^{1/4}}$$

$N^{3/4} - \left(\frac{N^{3/4}}{2} + \sqrt{N}\right)$

(Use exactly the same indicator variables and so on...)

$E_3: |C| > 4N^{3/4}$ ← is a bad event

We break this event into two events:

E_1 : at least $2N^{3/4}$ elements in C are \geq Median.

E_2 : at least $2N^{3/4}$ elements in C are \leq Median.

$$P(|C| > 4N^{3/4}) \leq P(E_1) + P(E_2)$$

$P(E_1)$: What is a rank of u in R ? By definition, it is $\frac{N^{3/4}}{2} + \sqrt{N}$ (that's how we picked u).

How many elements in R are bigger ^{than} or equal to u ? $N^{3/4} - (\frac{N^{3/4}}{2} + \sqrt{N}) = \frac{N^{3/4}}{2} - \sqrt{N}$.

What is the rank of u in S ? - At least $\frac{N}{2} + 2N^{3/4}$ because we are in case E_1 .

What is the size of l_u ? - at most $N - (\frac{N}{2} + 2N^{3/4}) = \frac{N}{2} - 2N^{3/4}$

Among $\frac{N}{2} - 2N^{3/4}$ elements in l_u we picked $\frac{N^{3/4}}{2} - \sqrt{N}$ of them and put them in R .

$X_i = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ sample is among the } \frac{N}{2} - 2N^{3/4} \text{ elements.} \\ 0 & \text{otherwise.} \end{cases}$

$$X = \sum_{i=1}^{N^{3/4}} X_i; \quad E[X] = E\left[\sum_{i=1}^{N^{3/4}} X_i\right] = \sum_{i=1}^{N^{3/4}} E[X_i] = \sum_{i=1}^{N^{3/4}} P(X_i = 1) =$$

All choices

$$\frac{\frac{N}{2} - 2N^{3/4}}{N} \leftarrow \text{Choices for } X_i \text{ to be 1} = \frac{1}{2} - \frac{2}{N^{1/4}}$$

What is $P(X_i = 1)$? Each element is equally likely;

$$\sum_{i=1}^{N^{3/4}} P(X_i = 1) = N^{3/4} \cdot \left(\frac{1}{2} - \frac{2}{N^{1/4}}\right) = \frac{N^{3/4}}{2} - 2\sqrt{N} = E[X]$$

$$P(\mathcal{E}_1) = P(X \geq \frac{N^{3/4}}{2} - \sqrt{N}) = P(X - E[X] \geq \frac{N^{3/4}}{2} - \sqrt{N} - \frac{N^{3/4}}{2} + 2\sqrt{N}) =$$

$$= P(X - E[X] \geq \sqrt{N}) \stackrel{\text{Chebyshev?}}{\leq} P(|X - E[X]| \geq \sqrt{N}) \leq \frac{\text{Var}(X)}{(\sqrt{N})^2} \leq \frac{N^{3/4}}{4N} =$$

$$\text{HW: } \text{Var}(X) = \sum_{i=1}^{N^{3/4}} \text{Var}(X_i) \leq \frac{N^{3/4}}{4}$$

Very small probability that \mathcal{E}_1 happens

$$\text{So, } P(|C| > 4N^{3/4}) \leq P(\mathcal{E}_1) + P(\mathcal{E}_2) \leq \frac{1}{4N^{1/4}} + \frac{1}{4N^{1/4}} = \frac{1}{2N^{1/4}}$$

$$P(l_u > \frac{N}{2}) + P(l_d > \frac{N}{2}) + P(|C| > 4N^{3/4}) \leq \frac{1}{4N^{1/4}} + \frac{1}{4N^{1/4}} + \frac{1}{2N^{1/4}} \leq \frac{1}{N^{1/4}}$$

$$Pr(X_i = 1) = \frac{1}{2} - \frac{2}{N^{1/4}}$$

$$\text{Var}(X_i) = E(X_i^2) - E(X_i)^2 = E(X_i) - E(X_i)^2 = \frac{1}{2} - \frac{2}{N^{1/4}} - \left(\frac{1}{2} - \frac{2}{N^{1/4}}\right)^2 =$$

$$= \frac{1}{2} - \frac{2}{N^{1/4}} - \frac{1}{4} + \frac{2 \cdot 2}{N^{1/4}} - \frac{4}{N^{1/2}} = \frac{1}{4} - \frac{4}{N^{1/2}}$$

$$\text{Var}(X) = \text{Var}\left(\sum_{i=1}^{N^{3/4}} X_i\right) = \sum_{i=1}^{N^{3/4}} \text{Var}(X_i) = N^{3/4} \cdot \left(\frac{1}{4} - \frac{4}{N^{1/2}}\right) = \frac{N^{3/4}}{4} - 4N^{1/4} \leq \frac{N^{3/4}}{4}$$