

Expectation is not everything....

Which game would you prefer?

- ① With probability $\frac{1}{2}$ win \$1, with probability $\frac{1}{2}$ pay \$1.
- ② With probability $\frac{1}{2}$ win \$100,000, with probability $\frac{1}{2}$ pay \$100,000.
- ③ With probability $\frac{1}{1,000,000}$ win \$1,000,000, with probability $\frac{1}{2}$ pay \$5, else \$0.

$$(1) E(X) = 1 \cdot \frac{1}{2} + (-1) \cdot \frac{1}{2} = 0$$

$$(2) E(X) = 10^5 \cdot \frac{1}{2} + (-10^5) \cdot \frac{1}{2} = 0$$

$$(3) E(X) = 10^6 \cdot \frac{1}{10^6} + (-5) \cdot \frac{1}{2} + 0 \cdot \left(1 - \frac{1}{2} - \frac{1}{10^6}\right) = 1 - \frac{5}{2}$$

Which Job Would You Prefer?

- A job that pays \$1000 a week. — $E(X) = 1000$
- A job that pays \$1 a week plus a bonus of \$1,000,000 with probability $\frac{1}{1000}$. $\rightarrow E(X) = 1 + 10^3$

Bounding Deviation from Expectation

Theorem

[Markov Inequality] For any non-negative random variable

$$\Pr(X \geq a) \leq \frac{E[X]}{a}.$$

Proof.

$$E[X] = \sum i \Pr(X = i) \geq a \sum_{i \geq a} \Pr(X = i) = a \Pr(X \geq a).$$



Example: What is the probability of getting more than $\frac{3N}{4}$ heads in N coin flips? $\leq \frac{N/2}{3N/4} \leq \frac{2}{3}$.

Variance

Definition

The **variance** of a random variable X is

$$\text{Var}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2] = \underbrace{\mathbf{E}[X^2]}_{\text{hard to compute}} - (\mathbf{E}[X])^2.$$

→ Easy to compute

↳ hard to compute

Definition

The **standard deviation** of a random variable X is

$$\sigma(X) = \sqrt{\text{Var}[X]}.$$

Example: Let X be a 0-1 random variable with $Pr(X = 0) = Pr(X = 1) = 1/2$.

$$[E[X] = 1/2.]^*$$

$$Var[X] = \frac{1}{2}(1 - \frac{1}{2})^2 + \frac{1}{2}(0 - \frac{1}{2})^2 = \frac{1}{4}.$$

$$\hookrightarrow = E((X - E(X))^2)^*$$

Chebyshev's Inequality

Theorem 1

For **any** random variable

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

Proof.

$$\Pr(|X - E[X]| \geq a) = \Pr((X - E[X])^2 \geq a^2)$$

By Markov inequality

$$\begin{aligned} \Pr((X - E[X])^2 \geq a^2) &\leq \frac{E[(X - E[X])^2]}{a^2} \\ &= \frac{\text{Var}[X]}{a^2} \end{aligned}$$

GENERAL USEFUL CONCENTRATION BOUNDS

Theorem 2

For **any** random variable

$$\sigma[X] = \sqrt{\text{VAR}[X]}$$
$$\Pr(|X - E[X]| \geq a\sigma[X]) \leq \frac{1}{a^2}.$$

Theorem 3

For **any** random variable

$$\Pr(|X - E[X]| \geq \epsilon E[X]) \leq \frac{\text{Var}[X]}{\epsilon^2 (E[X])^2}.$$

Theorem 4

If X and Y are independent random variable

$$E[XY] = E[X] \cdot E[Y],$$

Proof.

$$\begin{aligned} E[XY] &= \sum_i \sum_j i \cdot j \Pr((X = i) \cap (Y = j)) = \\ &= \sum_i \sum_j ij \Pr(X = i) \cdot \Pr(Y = j) = \\ &= \left(\sum_i i \Pr(X = i) \right) \left(\sum_j j \Pr(Y = j) \right). \end{aligned}$$



Theorem §

If X and Y are *independent* random variable

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

Proof.

$$\begin{aligned}\text{Var}[X + Y] &= E[(X + Y - E[X] - E[Y])^2] = \\ E[(X - E[X])^2 + (Y - E[Y])^2 + 2(X - E[X])(Y - E[Y])] &= \\ \text{Var}[X] + \text{Var}[Y] + 2E[X - E[X]]E[Y - E[Y]]\end{aligned}$$

Since the random variables $X - E[X]$ and $Y - E[Y]$ are independent.

But $E[X - E[X]] = E[X] - E[X] = 0$.



Back to Coin Flips

Assume again that we flip N coins. Let X be the number of heads.

$X_i = 1$ if the i -th flip was a head else $X_i = 0$.

$E[X_i] = 1/2$. $\text{Var}[X_i] = 1/4$.

$$\Pr(X \geq 3N/4) \leq \Pr(|X - E[X]| \geq N/4) =$$

$$\Pr(|X - E[X]| \geq E[X]/2) \leq \frac{\text{Var}[X]}{(E[X])^2(1/4)} =$$

THM 3

THM 5

$$\frac{N/4}{(N^2/4)(1/4)} = \boxed{4/N} \quad !!$$

A significantly better bound than $3/4$.

Bernoulli Trial

Let X be a 0-1 random variable such that

$$Pr(X = 1) = p, \quad Pr(X = 0) = 1 - p.$$

$$E[X] = 1 \cdot p + 0 \cdot (1 - p) = p.$$

$$Var[X] = p(1 - p)^2 + (1 - p)(0 - p)^2 = p(1 - p)[1 - p + p] =$$

$$p(1 - p).$$

A Binomial Random variable

Consider a sequence of n independent Bernoulli trials X_1, \dots, X_n .

Let

$$X = \sum_{i=1}^n X_i.$$

X has a **Binomial** distribution $X \sim B(n, p)$.

$$Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}.$$

$$E[X] = np.$$

$$Var[X] = np(1 - p).$$

Algorithm for Computing the Median

The **median** of a set X of n distinct elements is the $\lceil \frac{n}{2} \rceil$ largest element in the set.

If $n = 2k + 1$, the median element is the $k + 1$ -th element in the sorted order.

Easily computed through sorting in $O(n \log n)$ time. There exists a complicated $O(n)$ deterministic algorithm.)

↳ not suitable for most

Randomized Median Algorithm

S

Input: A set of $n = 2k + 1$ elements from a totally ordered universe.

Output: The $k + 1$ -th largest element in the set.

$$|R| = n^{3/4}$$

- 1 Pick a (multi)-set R of $s = n^{3/4}$ elements in S , chosen independently and uniformly at random with replacement. Sort the set R . }*

- 2 Let d be the $(\frac{1}{2}n^{3/4} - \sqrt{n})$ th smallest element in the sorted set R .

- 3 Let u be the $(\frac{1}{2}n^{3/4} + \sqrt{n})$ th smallest element in the sorted set R .

- 4 By comparing every element in S to d and u compute the set

$C = \{x \in S : d \leq x \leq u\}$, and the numbers

$l_d = |\{x \in S : x < d\}|$ and $l_u = |\{x \in S : x > u\}|$.

linear time

linear time

- 5 If $l_d > n/2$ or $l_u > n/2$ then FAIL.

- 6 If $|C| \leq 4n^{3/4}$ then sort the set C , otherwise FAIL.

- 7 Output the $(\lfloor \frac{n}{2} \rfloor - l_d + 1)$ st element in the sorted order of C .

$\hookrightarrow C$ "starts" from e_d

CONSTRUCTION of R from S' ($|S|=n$)

- $R = \emptyset$

- For $i = 1$ to $n^{3/4}$ do

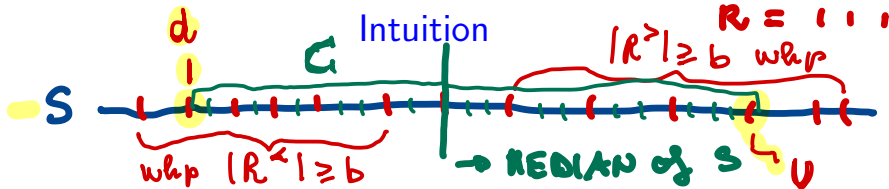
pick $u \in S'$ and do $\left[R = R \cup \{u\} \right]$
Keep u in S

- STOP

FACT. $\forall u \in S: \Pr[u \in R] = 1 - \left(1 - \frac{1}{n}\right)^{n^{3/4}} \quad (1)$

for suff. large n , $(1) \sim 1 - e^{-\frac{n^{3/4}}{n}} \sim \frac{1}{n^{1/4}}$

$\Pr[u \in R] \sim \frac{1}{n^{1/4}} \quad \forall u \in S'$



- We can sort sets of size $< n/\log n$ in linear time.
- The sample of R elements are spaced “more or less” evenly among the elements of X .

• W.h.p. more than $\frac{1}{2}n^{3/4} - \sqrt{n}$ samples are smaller than the median. $\Rightarrow |R^<| \geq b$

• W.h.p. more than $\frac{1}{2}n^{3/4} - \sqrt{n}$ samples are larger than the median. $|R^>| \geq b$

• W.h.p. the median is in the set C , and $|C| < n/\log n$.

Put all elements x of S s.t. $d \leq x \leq u$ into C

What might turn WRONG? BAD EVENTS

$$Y_1 \equiv \sum_{i=1}^{n^{3/4}} Y_1^i \quad Y_1^i = 1 \text{ iff } x_i < \text{median in the sample}$$

Let Y_1 be the number of samples below the median. in R

Let Y_2 be the number of samples above the median. in R

The algorithm fails to compute the median in $O(n)$ time iff at least one of the following three events occurs: BAD EVENTS

- 1 $E_1: Y_1 < \frac{1}{2}n^{3/4} - \sqrt{n}$. $\rightarrow R$ is too shifted on the Right
d is too large
- 2 $E_2: Y_2 < \frac{1}{2}n^{3/4} - \sqrt{n}$. $\rightarrow R$ is ... symmetric
- 3 $E_3: |C| > n/\log n$. $\rightarrow R$ is not well distributed

What is the probability that the three random variables Y_1, Y_2 and $|C|$ are all within the required ranges?

The sample space in execution of this algorithm is the set of all possible choices of $n^{3/4}$ elements from n , with repetitions. (The sample space has $n^{n^{3/4}}$ points.)

Each point in the sample space defines values for Y_1 , Y_2 and $|C|$. Computing the probabilities directly is too complicated, instead we use bounds on deviation from the expectation.

$$y_i^i = \begin{cases} 1 & \text{i-th} < \text{median} \\ 0 & \text{otherwise} \end{cases}$$

$$Y_1 = \sum_{i=1}^n y_i^i$$

Y_1 be the number of samples below the median.

What is the probability that $Y_1 < \frac{1}{2}n^{3/4} - \sqrt{n}$

Viewing Y_1 as the sum of $n^{3/4}$ independent 0-1 random variable, each with expectation $1/2$ and variance $1/4$ we prove (not counting the median itself):

$$E[Y_1] = \frac{1}{2}n^{3/4}.$$

$$\begin{aligned} \text{Var}[Y_1] &= \frac{1}{4}n^{3/4}. \quad \text{VAR}[Y_1] = \\ &= \sum \text{VAR}[Y_1^i] = \\ & \quad n^{3/4} \cdot \frac{1}{4} \end{aligned}$$

THH1

Applying Chebyshev Inequality we get:

$$Pr(E_1 : Y_1 < \frac{1}{2}n^{3/4} - \sqrt{n}) \leq Pr(|Y_1 - E[Y_1]| > \sqrt{n}) \leq \frac{a}{n}$$

$$\frac{Var[Y_1]}{n} = \frac{n^{3/4}/4}{n} = \frac{1}{4}n^{-1/4}. \quad \text{w hyp}$$

\searrow
 a^2

Similarly

$$Pr(E_2 : Y_2 < \frac{1}{2}n^{3/4} - \sqrt{n}) \leq \frac{1}{4}n^{-1/4}.$$

$$Pr(E_1 \cup E_2) \leq \frac{2}{4}n^{-1/4}.$$

$$C \equiv \{x \in S \mid d \leq x \leq u\} \quad d \equiv \left(\frac{1}{2}n^{3/4} - \sqrt{n}\right) \text{ elem in } R$$

$$U \equiv \dots + \sqrt{n} \dots$$

Recall: $E_3 : |C| > n/\log n$.

Lemma

$$\Pr(E_3) \leq \frac{1}{2}n^{-1/4}.$$

Define the following two events:

- ① $\mathcal{E}_{3,1}$: at least $2n^{3/4}$ elements of C are greater than the median;
- ② $\mathcal{E}_{3,2}$: at least $2n^{3/4}$ elements of C are smaller than the median.

If $|C| > 4n^{3/4}$, then at least one of the above two events occurs.

↳ this is deterministically true by def. of MEDIAN.

$U \equiv$ is the $(\frac{1}{2}n^{3/4} + \sqrt{n})$ -th element in sorted R
 * by def of $C \Rightarrow U$ is the largest element of C

We bound $\mathcal{E}_{3,1}$: at least $2n^{3/4}$ elements of C are greater than the median;

At least $2n^{3/4}$ elements of C above the median $\Rightarrow *$

u is at least the $\frac{1}{2}n + 2n^{3/4}$ largest in $S \Rightarrow$

R had at least $\frac{1}{2}n^{3/4} - \sqrt{n}$ samples among the $\frac{1}{2}n - 2n^{3/4}$ largest $**$
 elements in S .

Let X be the number of samples among the $\frac{1}{2}n - 2n^{3/4}$ largest $**$
 elements in S . Let $X = \sum_{i=1}^{n^{3/4}} X_i$ where

$$X_i = \begin{cases} 1 & \text{the } i\text{-th sample in } \frac{1}{2}n - 2n^{3/4} \\ & \text{largest elements in } S \\ 0 & \text{otherwise.} \end{cases}$$

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

$$\frac{\frac{1}{2}n - 2n^{3/4}}{n} = \rho_n$$

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

$$E[X] = \frac{1}{2}n^{3/4} - 2\sqrt{n}$$

$$\text{Var}[X] \leq \frac{1}{4}n^{3/4}$$

Applying Chebyshev's Inequality yields

$$\begin{aligned} \Pr(\mathcal{E}_{3,1}) &= \Pr(X \geq \frac{1}{2}n^{3/4} - \sqrt{n}) \\ &\leq \Pr(|X - E[X]| \geq \sqrt{n}) \\ &\leq \frac{\text{Var}[X]}{n} = \frac{\frac{n^{3/4}}{4}}{n} = \frac{1}{4}n^{-1/4}. \end{aligned}$$

Similarly,

$$\Pr(\mathcal{E}_{3,2}) \leq \frac{1}{4}n^{-\frac{1}{4}},$$

and

$$\Pr(E_3) \leq \Pr(\mathcal{E}_{3,1}) + \Pr(\mathcal{E}_{3,2}) \leq \frac{1}{2}n^{-\frac{1}{4}}.$$

The probability that the algorithm succeeds is

$$\geq 1 - (\Pr(E_1) + \Pr(E_2) + \Pr(E_3)) \geq 1 - \frac{1}{n^{1/4}}.$$

The Geometric Distribution

- How many times we need to perform a trial with probability p for success till we get the first success?
- How many times do we need to roll a dice until we get the first 6?

Definition

A geometric random variable X with parameter p is given by the following probability distribution on $n = 1, 2, \dots$

$$\Pr(X = n) = (1 - p)^{n-1} p.$$

Memoryless Distribution

Lemma

For a geometric random variable with parameter p and $n > 0$,

$$\Pr(X = n + k \mid X > k) = \Pr(X = n).$$

Proof.

$$\begin{aligned}\Pr(X = n + k \mid X > k) &= \frac{\Pr((X = n + k) \cap (X > k))}{\Pr(X > k)} \\&= \frac{\Pr(X = n + k)}{\Pr(X > k)} = \frac{(1 - p)^{n+k-1}p}{\sum_{i=k}^{\infty} (1 - p)^i p} \\&= \frac{(1 - p)^{n+k-1}p}{(1 - p)^k} = (1 - p)^{n-1}p = \Pr(X = n).\end{aligned}$$



Expectation

- Let X be a geometric random variable with parameter p .
- Let $Y = 1$ if the first trial is a success, $Y = 0$ otherwise.
-

$$\begin{aligned}\mathbf{E}[X] &= \Pr(Y = 0)\mathbf{E}[X \mid Y = 0] + \Pr(Y = 1)\mathbf{E}[X \mid Y = 1] \\ &= (1 - p)\mathbf{E}[X \mid Y = 0] + p\mathbf{E}[X \mid Y = 1].\end{aligned}$$

- If $Y = 0$ let Z be the number of trials after the first one.
- $\mathbf{E}[X] = (1 - p)\mathbf{E}[Z + 1] + p \cdot 1 = (1 - p)\mathbf{E}[Z] + 1$
- But $\mathbf{E}[Z] = \mathbf{E}[X]$, giving $\mathbf{E}[X] = 1/p$.

Example: Coupon Collector's Problem

- We place balls independently and uniformly at random in n boxes.
- Let X be the number of balls placed until all boxes are not empty.
- What is $E[X]$?

- Let X_i = number of balls placed when there were exactly $i - 1$ non-empty boxes.
- $X = \sum_{i=1}^n X_i$.
- X_i is a geometric random variable with parameter $p_i = 1 - \frac{i-1}{n}$.
-

$$\mathbf{E}[X_i] = \frac{1}{p_i} = \frac{n}{n - i + 1}.$$

$$\begin{aligned} \mathbf{E}[X] &= E \left[\sum_{i=1}^n X_i \right] = \sum_{i=1}^n \mathbf{E}[X_i] \\ &= \sum_{i=1}^n \frac{n}{n - i + 1} = n \sum_{i=1}^n \frac{1}{i} = n \ln n + \Theta(n). \end{aligned}$$

Variance of a Geometric Random Variable

- We use

$$\text{Var}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[X^2] - (\mathbf{E}[X])^2.$$

- To compute $\mathbf{E}[X^2]$, let $Y = 1$ if the first trial is a success, $Y = 0$ otherwise.

-

$$\begin{aligned}\mathbf{E}[X^2] &= \Pr(Y = 0)\mathbf{E}[X^2 \mid Y = 0] + \Pr(Y = 1)\mathbf{E}[X^2 \mid Y = 1] \\ &= (1 - p)\mathbf{E}[X^2 \mid Y = 0] + p\mathbf{E}[X^2 \mid Y = 1].\end{aligned}$$

- If $Y = 0$ let Z be the number of trials after the first one.

-

$$\begin{aligned}\mathbf{E}[X^2] &= (1 - p)\mathbf{E}[(Z + 1)^2] + p \cdot 1 \\ &= (1 - p)\mathbf{E}[Z^2] + 2(1 - p)\mathbf{E}[Z] + 1,\end{aligned}$$

- $\mathbf{E}[Z] = 1/p$ and $\mathbf{E}[Z^2] = \mathbf{E}[X^2]$.



$$\begin{aligned}\mathbf{E}[X^2] &= (1-p)\mathbf{E}[(Z+1)^2] + p \cdot 1 \\ &= (1-p)\mathbf{E}[Z^2] + 2(1-p)\mathbf{E}[Z] + 1,\end{aligned}$$



$$\mathbf{E}[X^2] = (1-p)\mathbf{E}[X^2] + 2(1-p)/p + 1 = (1-p)\mathbf{E}[X^2] + (2-p)/p,$$

- $\mathbf{E}[X^2] = (2-p)/p^2$.

$$\begin{aligned} \text{Var}[X] &= \mathbf{E}[X^2] - \mathbf{E}[X]^2 \\ &= \frac{2-p}{p^2} - \frac{1}{p^2} \\ &= \frac{1-p}{p^2}. \end{aligned}$$

Back to the Coupon Collector's Problem

- We place balls independently and uniformly at random in n boxes.
- Let X be the number of balls placed until all boxes are not empty.
- $E[X] = nH_n = n \ln n + \Theta(n)$
- What is $\Pr(X \geq 2E[X])$?
- Applying Markov's inequality

$$\Pr(X \geq 2nH_n) \leq \frac{1}{2}.$$

- Can we do better?

- Let X_i = number of balls placed when there were exactly $i - 1$ non-empty boxes.
- $X = \sum_{i=1}^n X_i$.
- X_i is a geometric random variable with parameter $p_i = 1 - \frac{i-1}{n}$.
- $Var[X_i] \leq \frac{1}{p^2} \leq (\frac{n}{n-i+1})^2$.
-

$$Var[X] = \sum_{i=1}^n Var[X_i] \leq \sum_{i=1}^n \left(\frac{n}{n-i+1} \right)^2 = n^2 \sum_{i=1}^n \left(\frac{1}{i} \right)^2 \leq \frac{\pi^2 n^2}{6}.$$

- By Chebyshev's inequality

$$\Pr(|X - nH_n| \geq nH_n) \leq \frac{n^2 \pi^2 / 6}{(nH_n)^2} = \frac{\pi^2}{6(H_n)^2} = O\left(\frac{1}{\ln^2 n}\right).$$

Direct Bound

- The probability of not obtaining the i -th coupon after $n \ln n + cn$ steps:

$$\left(1 - \frac{1}{n}\right)^{n(\ln n + c)} < e^{-(\ln n + c)} = \frac{1}{e^c n}.$$

- By a union bound, the probability that some coupon has not been collected after $n \ln n + cn$ step is e^{-c} .
- The probability that all coupons are not collected after $2n \ln n$ steps is at most $1/n$.

The Advantage of Multiple Samples

Theorem

For any constant a ,

$$\text{Var}[aX] = a^2 \text{Var}[X].$$

Proof.

$$\begin{aligned}\text{Var}[aX] &= E[(aX - E[aX])^2] = E[a^2(X - E[X])^2] \\ &= a^2 E[(X - E[X])^2] = a^2 \text{Var}[X].\end{aligned}$$



Theorem

Let X_1, \dots, X_n be n independent, identically distributed random variable. Let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$.

$$\text{Var}[\bar{X}] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n X_i\right] = \frac{1}{n} \text{Var}[X_i].$$

The (Weak) Law of Large Numbers

Theorem

Let X_1, \dots, X_n be independent, identically distributed, random variables. Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. For any constant $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \Pr(|\bar{X}_n - \mathbf{E}[X]| \leq \epsilon) = 1.$$

Proof.

$\text{Var}[\bar{X}_n] = \frac{1}{n} \text{Var}[X_i]$. Applying Chebyshev's bound

$$\Pr(|\bar{X}_n - \mathbf{E}[X]| > \epsilon) \leq \frac{\text{Var}[X_i]}{n\epsilon^2}.$$



[Can be proven even when $\text{Var}[X_i]$ is not bounded.]