S&DS 241 Lecture 12 Mean, variance, deviation inequalities

B-H: 4.6, 10.1.3

Recall from Lecture 5: Expectation

The expectation (aka expected value or mean) of a discrete random variable X is

$$E(X) = \sum_{x \in \mathcal{X}} x p_X(x)$$

where $p_X(x) = P(X = x)$ is the PMF.

Properties of expectation

• X is a nonnegative random variable $\implies E(X) \ge 0$

Properties of expectation

- X is a nonnegative random variable $\implies E(X) \geq 0$
- Linearity of expectation: For any constants *a*, *b* and any random variables *X* and *Y* (not necessarily independent!)

$$E(aX + bY) = aE(X) + bEY$$

$$E(X) = \sum_{x \in \mathcal{X}} xp_X(x) = \text{average of values weighted by PMF}$$



- Expected value \neq typical value!
- How close a random variable is to its expectation depends on many things, e.g., variance



- Expected value \neq typical value!
- How close a random variable is to its expectation depends on many things, e.g., variance





- Expected value ≠ typical value!
- How close a random variable is to its expectation depends on many things, e.g., variance



Variance of a random variable X:

$$\operatorname{Var}(X) = E((X - E(X))^2)$$

- Other notations: $\sigma_X^2, V(X)$
- Variance = mean-squared deviation from the expectation
- Significance: measures the uncertainty of a random variable/the spread of a distribution

Standard deviation

• Variance is quadratic in nature: for example

	Units
X	ft
E(X)	ft
$\operatorname{Var}(X)$	sqft

Standard deviation

• Variance is quadratic in nature: for example

	Units
X	ft
E(X)	ft
$\operatorname{Var}(X)$	sqft

• Standard deviation of X:

$$SD(X) = \sigma_X = \sqrt{Var(X)}$$

which is the root-mean-squared deviation from expectation

Standard deviation

• Variance is quadratic in nature: for example

	Units
X	ft
E(X)	ft
$\operatorname{Var}(X)$	sqft
$\mathrm{SD}(X)$	ft

• <u>Standard deviation of X</u>:

$$\mathrm{SD}(X) = \sigma_X = \sqrt{\mathrm{Var}(X)}$$

which is the root-mean-squared deviation from expectation

Intuition

- Smaller Var(X) or $SD(X) \implies$ PMF of X is more concentrated around the mean
- Larger $\operatorname{Var}(X)$ or $\operatorname{SD}(X) \implies$ PMF of X is more spread out

Example



• Mean: E(X) = EY = 0

Example



• Mean:
$$E(X) = EY = 0$$

• Variance:

$$\operatorname{Var}(X) = E(X^2) \xrightarrow{\text{LOTUS}} \frac{1}{2} \times (-10)^2 + \frac{1}{2} \times 10^2 = 100$$
$$\operatorname{Var}(Y) = E(Y^2) \xrightarrow{\text{LOTUS}} \frac{1^2 \times 2}{8} + \frac{2^2 \times 2 + 3^2 \times 2}{16} = \frac{7}{4}$$

Example



• Mean:
$$E(X) = EY = 0$$

• Variance:

$$\operatorname{Var}(X) = E(X^2) \xrightarrow{\text{LOTUS}} \frac{1}{2} \times (-10)^2 + \frac{1}{2} \times 10^2 = 100$$
$$\operatorname{Var}(Y) = E(Y^2) \xrightarrow{\text{LOTUS}} \frac{1^2 \times 2}{8} + \frac{2^2 \times 2 + 3^2 \times 2}{16} = \frac{7}{4}$$

• Standard deviation:

$$SD(X) = 10, \quad SD(Y) \approx 1.32$$

Indeed, Y is much more concentrated than X

Alternative formula of variance

$$\operatorname{Var}(X) = E(X^2) - E(X)^2$$

Proof.

Call
$$\mu = E(X)$$
. By linearity of expectation,

$$Var(X) = E((X - \mu)^2)$$

$$= E(X^2 + \mu^2 - 2\mu X) = E(X^2) + \mu^2 - 2\mu \cdot E(X)$$

$$= E(X^2) - \mu^2$$

Alternative formula of variance

$$\operatorname{Var}(X) = E(X^2) - E(X)^2$$

Proof.

Call
$$\mu = E(X)$$
. By linearity of expectation,

$$Var(X) = E((X - \mu)^2)$$

$$= E(X^2 + \mu^2 - 2\mu X) = E(X^2) + \mu^2 - 2\mu \cdot E(X)$$

$$= E(X^2) - \mu^2$$

- First moment: E(X)
- Second moment: $E(X^2)$
- kth moment: $E(X^k)$

Alternative formula of variance

$$\operatorname{Var}(X) = E(X^2) - E(X)^2$$

Proof.

Call
$$\mu = E(X)$$
. By linearity of expectation,

$$Var(X) = E((X - \mu)^2)$$

$$= E(X^2 + \mu^2 - 2\mu X) = E(X^2) + \mu^2 - 2\mu \cdot E(X)$$

$$= E(X^2) - \mu^2$$

- First moment: E(X)
- Second moment: $E(X^2)$
- kth moment: $E(X^k)$
- Variance = second moment $-(first moment)^2$

Example: dice

Let X be the outcome of a fair die. Find Var(X).

$$E(X) = \frac{1}{6}(1+2+3+4+5+6) = \frac{7}{2}$$
$$E(X^2) = \frac{1}{6}(1^2+2^2+3^2+4^2+5^2+6^2) = \frac{91}{6}$$

Therefore

$$Var(X) = \frac{91}{6} - \left(\frac{7}{2}\right)^2 = \frac{35}{12}$$
$$SD(X) = \sqrt{\frac{35}{12}} \approx 1.71$$

Suppose X takes values in [-10, 10].

• What is the largest possible variance?

Suppose X takes values in [-10, 10].

• What is the largest possible variance? 100



Suppose X takes values in [-10, 10].

• What is the largest possible variance? 100



• What is the smallest possible variance?

Suppose X takes values in [-10, 10].

• What is the largest possible variance? 100



• What is the smallest possible variance? 0



Suppose X takes values in [-10, 10].

• What is the largest possible variance? 100



Fact: $Var(X) = 0 \Leftrightarrow X = constant (no randomness)$

Properties of mean, variance and std dev

• Shift does not change variance

$$\operatorname{Var}(X+b) = \operatorname{Var}(X)$$

• Scaling:

$$\operatorname{Var}(aX) = a^2 \operatorname{Var}(X)$$

• More generally:

$$E(aX + b) = aE(X) + b$$

Var(aX + b) = a²Var(X)
SD(aX + b) = |a|SD(X)

• For any X and Y

$$E(X+Y) = E(X) + E(Y)$$

thanks to linearity of expectation

• For any X and Y

$$E(X+Y) = E(X) + E(Y)$$

thanks to linearity of expectation

• For independent X and Y

$$\operatorname{Var}(X+Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$$

• For any X and Y

$$E(X+Y) = E(X) + E(Y)$$

thanks to linearity of expectation

• For independent X and Y

$$\operatorname{Var}(X+Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$$

This can fail without independence,

• For any X and Y

$$E(X+Y) = E(X) + E(Y)$$

thanks to linearity of expectation

• For independent X and Y

$$\operatorname{Var}(X+Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$$

This can fail without independence, e.g., Y = -X.

For independent X and Y

$$\operatorname{Var}(X+Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$$

Proof.

Since shifting does not change variance, we can assume, without loss of generality, that E(X) = E(Y) = 0.

For independent X and Y

$$\operatorname{Var}(X+Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$$

Proof.

Since shifting does not change variance, we can assume, without loss of generality, that E(X) = E(Y) = 0. Then Var(X + Y) $= E((X + Y)^2)$ $= E(X^2 + Y^2 + 2XY)$ $= E(X^2) + E(Y^2) + 2E(XY)$ linearity of expectation $= E(X^2) + E(Y^2)$ independence so E(XY) = E(X)E(Y)= Var(X) + Var(Y)

Variance of sum of independent random variables

Let X_1, \ldots, X_n be independent random variables. Then

$$\operatorname{Var}(X_1 + X_2 + \dots + X_n) = \operatorname{Var}(X_1) + \dots + \operatorname{Var}(X_n)$$

Variance of sum of iid random variables

Let X_1, \ldots, X_n be iid random variables. Then

$$\operatorname{Var}(X_1 + X_2 + \dots + X_n) = n \cdot \operatorname{Var}(X_1)$$

$$\operatorname{SD}(X_1 + X_2 + \dots + X_n) = \sqrt{n} \cdot \operatorname{SD}(X_1)$$

in contrast

$$SD(X_1 + X_1 + \dots + X_1) = n \cdot SD(X_1)$$

Note: $\sqrt{n} \ll n$ for large *n*. Why does this happen?

Variance of sum of iid random variables

Let X_1, \ldots, X_n be iid random variables. Then

$$\operatorname{Var}(X_1 + X_2 + \dots + X_n) = n \cdot \operatorname{Var}(X_1)$$

$$\operatorname{SD}(X_1 + X_2 + \dots + X_n) = \sqrt{n} \cdot \operatorname{SD}(X_1)$$

in contrast

$$SD(X_1 + X_1 + \dots + X_1) = n \cdot SD(X_1)$$

Note: $\sqrt{n} \ll n$ for large n. Why does this happen? Example: random walk $(X_i = \pm 1)$

- $X_1 + X_1 + \dots + X_1$: all steps are aligned (either all + or all -). This leads to SD on the order of n.
- X₁ + X₂ + · · · + X_n: some are + and some −. Cancellation leads to SD on the order of √n.

Statistical application: sample average

Let X_1, \ldots, X_n be iid random variables with mean μ and variance σ^2 .

• Sample average:

$$\overline{X} = \frac{1}{n}(X_1 + \ldots + X_n)$$

Then

$$E(\overline{X}) = \mu$$
, $Var(\overline{X}) = \frac{\sigma^2}{n} \ll \sigma^2$ for large n

• More data (larger sample size) \implies less uncertainty
Summary

- $\operatorname{Var}(X) = \operatorname{SD}(X)^2 = E((X E(X))^2) \ge 0$
- $Var(X) = 0 \Leftrightarrow X = constant$
- $Var(X) = E(X^2) (E(X))^2$
- $\operatorname{Var}(aX + b) = a^2 \operatorname{Var}(X)$
- $\operatorname{Var}(X+Y) \xrightarrow{\text{independence}} \operatorname{Var}(X) + \operatorname{Var}(Y)$

Variance of common distributions

Bernoulli

 $X \sim \operatorname{Bern}(p)$. Then

$$Var(X) = E(X^2) - (E(X))^2 = p - p^2 = p(1 - p)$$

Bernoulli

 $X \sim \text{Bern}(p)$. Then

$$Var(X) = E(X^{2}) - (E(X))^{2} = p - p^{2} = p(1 - p)$$



Binomial

 $X\sim {\rm Bin}(n,p).$ Then

$$\operatorname{Var}(X) = np(1-p)$$

Proof.

$$X = \underbrace{X_1 + \dots + X_n}_{\text{i.i.d. Bern}(p)} \implies \operatorname{Var}(X) = n \operatorname{Var}(X_1)$$

Binomial

 $X\sim {\rm Bin}(n,p).$ Then

$$\operatorname{Var}(X) = np(1-p)$$

Proof.

$$X = \underbrace{X_1 + \dots + X_n}_{\text{i.i.d. Bern}(p)} \implies \operatorname{Var}(X) = n\operatorname{Var}(X_1)$$

<u>Exercise</u>: Alternatively, find $E(X^2)$ using LOTUS rule and binomial PMF, then apply $E(X^2) - (E(X))^2$.

Poisson

$X \sim \mathsf{Pois}(\lambda)$. Then

$$\operatorname{Var}(X) = \lambda$$

Poisson

 $X \sim \mathsf{Pois}(\lambda)$. Then

$$\operatorname{Var}(X) = \lambda$$

Interpretation: Poisson as limiting Binomial

$$\operatorname{Bin}\left(n, \frac{\lambda}{n}\right) \xrightarrow{n \to \infty} \operatorname{Pois}(\lambda)$$

so we expect:

$$\operatorname{Var} = n \frac{\lambda}{n} \left(1 - \frac{\lambda}{n} \right) \xrightarrow{n \to \infty} \quad \lambda$$

$$E(X^2) = \sum_{k=0}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!}$$

$$E(X^2) = \sum_{k=1}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!}$$

$$E(X^2) = \sum_{k=1}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!}$$
$$= \sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^k}{(k-1)!}$$

$$E(X^{2}) = \sum_{k=1}^{\infty} k^{2} \frac{e^{-\lambda} \lambda^{k}}{k!}$$

= $\sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^{k}}{(k-1)!} = \sum_{k=1}^{\infty} (k-1) \frac{e^{-\lambda} \lambda^{k}}{(k-1)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^{k}}{(k-1)!}$

$$E(X^{2}) = \sum_{k=1}^{\infty} k^{2} \frac{e^{-\lambda} \lambda^{k}}{k!}$$
$$= \sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^{k}}{(k-1)!} = \sum_{k=2}^{\infty} (k-1) \frac{e^{-\lambda} \lambda^{k}}{(k-1)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^{k}}{(k-1)!}$$

$$\begin{split} E(X^2) &= \sum_{k=1}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!} \\ &= \sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^k}{(k-1)!} = \sum_{k=2}^{\infty} (k-1) \frac{e^{-\lambda} \lambda^k}{(k-1)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{k=2}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-2)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \end{split}$$

$$\begin{split} E(X^2) &= \sum_{k=1}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!} \\ &= \sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^k}{(k-1)!} = \sum_{k=2}^{\infty} (k-1) \frac{e^{-\lambda} \lambda^k}{(k-1)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{k=2}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-2)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+2}}{j!} + \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+1}}{j!} \end{split}$$

$$\begin{split} E(X^2) &= \sum_{k=1}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!} \\ &= \sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^k}{(k-1)!} = \sum_{k=2}^{\infty} (k-1) \frac{e^{-\lambda} \lambda^k}{(k-1)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{k=2}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-2)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+2}}{j!} + \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+1}}{j!} \\ &= \lambda^2 \sum_{\substack{j=0\\ =1}}^{\infty} \frac{e^{-\lambda} \lambda^j}{j!} + \lambda \sum_{\substack{j=0\\ =1}}^{\infty} \frac{e^{-\lambda} \lambda^j}{j!} \\ \end{split}$$

$$\begin{split} E(X^2) &= \sum_{k=1}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!} \\ &= \sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^k}{(k-1)!} = \sum_{k=2}^{\infty} (k-1) \frac{e^{-\lambda} \lambda^k}{(k-1)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{k=2}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-2)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+2}}{j!} + \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+1}}{j!} \\ &= \lambda^2 \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^j}{j!} + \lambda \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^j}{j!} = \lambda^2 + \lambda \end{split}$$

$$\begin{split} E(X^2) &= \sum_{k=1}^{\infty} k^2 \frac{e^{-\lambda} \lambda^k}{k!} \\ &= \sum_{k=1}^{\infty} k \frac{e^{-\lambda} \lambda^k}{(k-1)!} = \sum_{k=2}^{\infty} (k-1) \frac{e^{-\lambda} \lambda^k}{(k-1)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{k=2}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-2)!} + \sum_{k=1}^{\infty} \frac{e^{-\lambda} \lambda^k}{(k-1)!} \\ &= \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+2}}{j!} + \sum_{j=0}^{\infty} \frac{e^{-\lambda} \lambda^{j+1}}{j!} \\ &= \lambda^2 \sum_{\substack{j=0\\j=1}}^{\infty} \frac{e^{-\lambda} \lambda^j}{j!} + \lambda \sum_{\substack{j=0\\j=1}}^{\infty} \frac{e^{-\lambda} \lambda^j}{j!} = \lambda^2 + \lambda \end{split}$$

 and

$$\operatorname{Var}(X) = E(X^2) - (E(X))^2 = \lambda + \lambda^2 - \lambda^2 = \lambda$$











As λ increases, PMF of $\mathsf{Pois}(\lambda)$ shifts to the right and becomes more spread out.

Geometric

 $X\sim \mathsf{Geom}(p).$ Then $\mathrm{Var}(X)=\frac{1-p}{p^2}$

Exercise (see B-H Example 4.6.4)

Deviation inequalities

Motivations

In statistics, we want to answer questions like the following:

1 Toss a fair coin 100 times. Getting at least 75 heads is probably unlikely. How unlikely is it?

In statistics, we want to answer questions like the following:

- Toss a fair coin 100 times. Getting at least 75 heads is probably unlikely. How unlikely is it?
- 2 Toss a coin 100 times. Turns out there are 75 heads, which indicates the coin is likely biased. How biased is it and how confident are we about our estimate?

In statistics, we want to answer questions like the following:

- Toss a fair coin 100 times. Getting at least 75 heads is probably unlikely. How unlikely is it?
- 2 Toss a coin 100 times. Turns out there are 75 heads, which indicates the coin is likely biased. How biased is it and how confident are we about our estimate?

Useful tools: deviation inequalities (tail bounds)

Preparation





then for any random variable X

 $E(f(X)) \leq E(g(X))$

Markov's inequality

Let X be a non-negative random variable and a > 0. Then

$$P(X \ge a) \le \frac{E(X)}{a}$$



Markov's inequality

Let X be a non-negative random variable and a > 0. Then

$$P(X \ge a) \le \frac{E(X)}{a}$$

- <u>Intuition</u>: if the average salary is \$100, then the fraction of people earning \$300+ cannot exceed 1/3.
- <u>Lesson</u>: A non-negative random variable cannot far exceed its expectation with high probability.

Application: coin

Question

Toss a fair coin 100 times. How unlikely is it to get at least 75 heads?

Application: coin

Question

Toss a fair coin 100 times. How unlikely is it to get at least 75 heads?

Let $X \sim Bin(100, 1/2)$. Then E(X) = 50

Markov's inequality
$$\implies P(X \ge 75) \le \frac{E(X)}{75} = \frac{2}{3}$$

Application: coin

Question

Toss a fair coin 100 times. How unlikely is it to get at least 75 heads?

Let $X \sim \text{Bin}(100, 1/2)$. Then E(X) = 50

Markov's inequality
$$\implies P(X \ge 75) \le \frac{E(X)}{75} = \frac{2}{3}$$

- Useful estimate, but probably too conservative...
- Can we do better (using more information, e.g., variance)?

Chebyshev's inequality

Let X be an arbitrary random variable and d > 0. Then

$$P(|X - E(X)| \ge d) \le \frac{\operatorname{Var}(X)}{d^2}$$

Chebyshev's inequality

Let X be an arbitrary random variable and d > 0. Then

$$P(|X - E(X)| \ge d) \le \frac{\operatorname{Var}(X)}{d^2}$$

Proof.

Let $\mu = E(X)$. LHS = $P((X - \mu)^2 \ge d^2)$. Apply Markov's inequality to the nonnegative random variable $(X - \mu)^2$. Or, see picture below:
Chebyshev's inequality

Let X be an arbitrary random variable and d > 0. Then

$$P(|X - E(X)| \ge d) \le \frac{\operatorname{Var}(X)}{d^2}$$

Proof.

Let $\mu = E(X)$. LHS = $P((X - \mu)^2 \ge d^2)$. Apply Markov's inequality to the nonnegative random variable $(X - \mu)^2$. Or, see picture below:



Question

Toss a fair coin 100 times. How unlikely is it to get at least 75 heads?

Question

Toss a fair coin 100 times. How unlikely is it to get at least 75 heads?

Let $X \sim \text{Bin}(100, 1/2)$. Then E(X) = 50 and Var(X) = 25

$$\mathsf{Chebyshev} \implies P(X \ge 75) \le P(|X - 50| \ge 25) \le \frac{25}{25^2} = 4\%$$

Question

Toss a fair coin 100 times. How unlikely is it to get at least 75 heads?

Let $X \sim \text{Bin}(100, 1/2)$. Then E(X) = 50 and Var(X) = 25

Chebyshev
$$\implies P(X \ge 75) \le P(|X - 50| \ge 25) \le \frac{25}{25^2} = 4\%$$

- Much better than Markov (we used both mean and variance)
- Actually value: 2.8×10^{-7}

Question

Toss a fair coin 100 times. How unlikely is it to get at least 75 heads?

Let $X \sim \text{Bin}(100, 1/2)$. Then E(X) = 50 and Var(X) = 25

Chebyshev
$$\implies P(X \ge 75) \le P(|X - 50| \ge 25) \le \frac{25}{25^2} = 4\%$$

- Much better than Markov (we used both mean and variance)
- Actually value: 2.8×10^{-7}
- Central limit theorem (Lec 16) gives accurate estimate: $2.9 imes 10^{-7}$

Equivalent formulation of Chebyshev:

$$P(|X - E(X)| \ge C \cdot \sigma_X) \le \frac{1}{C^2}$$

Question

Toss a coin 100 times. Turns out there are 75 heads. How biased is the coin and how confident are we about our estimate?

Question

Toss a coin 100 times. Turns out there are 75 heads. How biased is the coin and how confident are we about our estimate?

- Let $X \sim \text{Bin}(n, p)$, where n = 100 and the bias p is unknown.
- Observe X = 75.

Question

Toss a coin 100 times. Turns out there are 75 heads. How biased is the coin and how confident are we about our estimate?

- Let $X \sim Bin(n, p)$, where n = 100 and the bias p is unknown.
- Observe X = 75.
- Empirical frequency of heads is

$$\hat{p} = \frac{X}{n} = 0.75,$$

a reasonable estimate of p.

Question

Toss a coin 100 times. Turns out there are 75 heads. How biased is the coin and how confident are we about our estimate?

- Let $X \sim Bin(n, p)$, where n = 100 and the bias p is unknown.
- Observe X = 75.
- Empirical frequency of heads is

$$\hat{p} = \frac{X}{n} = 0.75,$$

a reasonable estimate of p.

More refined estimate: confidence interval. If

$$P(p \in [\hat{p} - \epsilon, \hat{p} + \epsilon]) \ge c$$
, for any p

we say $[\hat{p}-\epsilon,\hat{p}+\epsilon]$ is a confidence interval with confidence level c , e.g. c=95%.

- Conflicting goals: higher confidence and narrower interval.
 - \blacktriangleright [0,1] is a CI of confidence level 100%, but not very useful

- Conflicting goals: higher confidence and narrower interval.
 - \blacktriangleright [0,1] is a CI of confidence level 100%, but not very useful
 - Need a tool to assess the chance of being off: let's use Chebyshev inequality

- Conflicting goals: higher confidence and narrower interval.
 - \blacktriangleright [0,1] is a CI of confidence level 100%, but not very useful
 - Need a tool to assess the chance of being off: let's use Chebyshev inequality

•
$$\operatorname{Var}(\hat{p}) = \operatorname{Var}(X/n) = \frac{1}{n^2} \operatorname{Var}(X) = \frac{p(1-p)}{n}$$

Thus

$$P(p \notin [\hat{p} - \epsilon, \hat{p} + \epsilon]) = P(|\hat{p} - p| > \epsilon) \le \frac{\operatorname{Var}(\hat{p})}{\epsilon^2} = \frac{p(1 - p)}{n\epsilon^2} \le \frac{1}{4n\epsilon^2}$$

- Conflicting goals: higher confidence and narrower interval.
 - $\blacktriangleright~[0,1]$ is a CI of confidence level 100%, but not very useful
 - Need a tool to assess the chance of being off: let's use Chebyshev inequality

•
$$\operatorname{Var}(\hat{p}) = \operatorname{Var}(X/n) = \frac{1}{n^2} \operatorname{Var}(X) = \frac{p(1-p)}{n}$$

Thus

$$P(p \notin [\hat{p} - \epsilon, \hat{p} + \epsilon]) = P(|\hat{p} - p| > \epsilon) \le \frac{\operatorname{Var}(\hat{p})}{\epsilon^2} = \frac{p(1 - p)}{n\epsilon^2} \le \frac{1}{4n\epsilon^2}$$

- Since n = 100, if we take $\epsilon = 0.1$, then $\frac{1}{4n\epsilon^2} = 0.25$.
- Upon observing X = 75, $[0.75 \pm 0.1]$ is a Cl of level 75%

• Want confidence level 99%. Then

$$P(p \notin [\hat{p} - \epsilon, \hat{p} + \epsilon]) \le \frac{1}{4n\epsilon^2} \le 1\% \xrightarrow{n=100} \epsilon \ge 0.5$$

therefore the CI is $[\hat{p}\pm 0.5]$ — too wide to be very useful.

• Want confidence level 99%. Then

$$P(p \notin [\hat{p} - \epsilon, \hat{p} + \epsilon]) \le \frac{1}{4n\epsilon^2} \le 1\% \xrightarrow{n=100} \epsilon \ge 0.5$$

therefore the CI is $[\hat{p}\pm0.5]$ — too wide to be very useful.

• How to get a CI $[\hat{p} \pm 0.1]$ of level 99%?

• Want confidence level 99%. Then

$$P(p \notin [\hat{p} - \epsilon, \hat{p} + \epsilon]) \le \frac{1}{4n\epsilon^2} \le 1\% \xrightarrow{n=100} \epsilon \ge 0.5$$

therefore the CI is $[\hat{p}\pm 0.5]$ — too wide to be very useful.

• How to get a CI $[\hat{p} \pm 0.1]$ of level 99%? Increase the sample size!

$$\frac{1}{4n\epsilon^2} \le 1\% \xrightarrow{\epsilon=0.1} n \ge 2500$$

• Want confidence level 99%. Then

$$P(p \notin [\hat{p} - \epsilon, \hat{p} + \epsilon]) \le \frac{1}{4n\epsilon^2} \le 1\% \xrightarrow{n=100} \epsilon \ge 0.5$$

therefore the CI is $[\hat{p}\pm0.5]$ — too wide to be very useful.

• How to get a CI $[\hat{p} \pm 0.1]$ of level 99%? Increase the sample size!

$$\frac{1}{4n\epsilon^2} \le 1\% \xrightarrow{\epsilon=0.1} n \ge 2500$$

(Better tools than Chebyshev $\implies n \ge 169$)

• Lesson: More data leads to more accurate estimate. Here $\operatorname{Var}(\hat{p}) \propto \frac{1}{n}$.





- Poll *n* people out of a population <u>uniformly</u> at random <u>with</u> <u>replacements</u>. Let *X* be the number of surveyed people who would vote for Biden.
- Then $X \sim Bin(n, p)$, where p = fraction of people in the entire population voting for Biden.
- Suppose we want $[\hat{p} \pm 0.01]$ to have confidence 95%

$$\frac{1}{4n\epsilon^2} \le 5\% \xrightarrow{\epsilon=0.01} \boxed{n \ge 50000}$$

- Poll *n* people out of a population <u>uniformly</u> at random <u>with</u> <u>replacements</u>. Let *X* be the number of surveyed people who would vote for Biden.
- Then $X \sim Bin(n, p)$, where p = fraction of people in the entire population voting for Biden.
- Suppose we want $[\hat{p} \pm 0.01]$ to have confidence 95%

$$\frac{1}{4n\epsilon^2} \le 5\% \xrightarrow{\epsilon=0.01} n \ge 50000$$

Surprise

The number of people needed to poll to reach a desired accuracy and confidence does NOT depend on the population size!!!

- Poll *n* people out of a population <u>uniformly</u> at random <u>with</u> <u>replacements</u>. Let *X* be the number of surveyed people who would vote for Biden.
- Then $X \sim Bin(n, p)$, where p = fraction of people in the entire population voting for Biden.
- Suppose we want $[\hat{p} \pm 0.01]$ to have confidence 95%

$$\frac{1}{4n\epsilon^2} \le 5\% \xrightarrow{\epsilon=0.01} \boxed{n \ge 50000}$$

Surprise

The number of people needed to poll to reach a desired accuracy and confidence does NOT depend on the population size!!!

So what's the catch?

- Poll *n* people out of a population <u>uniformly</u> at random <u>with</u> <u>replacements</u>. Let *X* be the number of surveyed people who would vote for Biden.
- Then $X \sim Bin(n, p)$, where p = fraction of people in the entire population voting for Biden.
- Suppose we want $[\hat{p} \pm 0.01]$ to have confidence 95%

$$\frac{1}{4n\epsilon^2} \le 5\% \xrightarrow{\epsilon=0.01} \boxed{n \ge 50000}$$

Surprise

The number of people needed to poll to reach a desired accuracy and confidence does NOT depend on the population size!!!

So what's the catch?

True story

Polling before phone and internet...



By mailing out millions of postcards to readers and simply counting the returns, *The Literary Digest* correctly predicted four presidential elections in a roll (1920,1924,1928,1932)

1936: Landon vs Roosevelt

 Literary Digest predicted Landon would win by 57% based on a sample of size 2.3 million

The Literary Digest NEW YORK OCTOBER 31, 1936

Topics of the day

LANDON, 1,293,669; ROOSEVELT, 972,897

Final Returns in The Digest's Poll of Ten Million Voters

Well, the great battle of the ballots in the lean National Committee purchased Turn Poll of ten million voters, scattered Larmaxe Donort? And all types and vari-throughout the forty-right States of the division ("Larma the done purchased vari-minim of non of the wordt muscle of times").

returned and let the people of the Nation draw their conclusions as to our accuracy So far, we have been right in every Pol Will we be right in the current Poll? That as Mrs. Roosevelt said concerning the Presi

1936: Landon vs Roosevelt

 Literary Digest predicted Landon would win by 57% based on a sample of size 2.3 million



Well, the great battle of the ballots in the lean National Committee purchased Tam Foll of ten million voters, scattered LTMTRARY DROMP?" And all types and vari-throughout the forty-right States of the etics, including: "Have the dress purchased

"We never make any claims before ele tion but we respectfully refer you to

George Gallup's company conducted a survey of much smaller size 50000 and predicted FDR would win by 56%



As it turns out...

• FDR won by a landslide (61%)

- FDR won by a landslide (61%)
- Literary Digest had their final issue in 1938

- FDR won by a landslide (61%)
- Literary Digest had their final issue in 1938
- Gallup went on to make a career in polling

- FDR won by a landslide (61%)
- Literary Digest had their final issue in 1938
- Gallup went on to make a career in polling
- "The Poll That Changed Polling"

- FDR won by a landslide (61%)
- Literary Digest had their final issue in 1938
- Gallup went on to make a career in polling
- "The Poll That Changed Polling"

- FDR won by a landslide (61%)
- Literary Digest had their final issue in 1938
- Gallup went on to make a career in polling
- "The Poll That Changed Polling"

So why was Literary Digest so wrong?

- FDR won by a landslide (61%)
- Literary Digest had their final issue in 1938
- Gallup went on to make a career in polling
- "The Poll That Changed Polling"

So why was Literary Digest so wrong? Selection bias!

Why was Literary Digest so wrong?

- Unbeknownst to the magazine, their subscribers tend to be more affluent Americans who favor Republicans.
- Gallup noticed this bias, and polled a much smaller but demographically representative of the population

Why was Literary Digest so wrong?

- Unbeknownst to the magazine, their subscribers tend to be more affluent Americans who favor Republicans.
- Gallup noticed this bias, and polled a much smaller but demographically representative of the population
- To drive his point home, he even predicted the result of the Literary Digest poll to within about 1%

Why was Literary Digest so wrong?

- Unbeknownst to the magazine, their subscribers tend to be more affluent Americans who favor Republicans.
- Gallup noticed this bias, and polled a much smaller but demographically representative of the population
- To drive his point home, he even predicted the result of the Literary Digest poll to within about 1%
- There is lots of science behind opinion polling and survey sampling
 - Selection bias
 - Non-response bias
 - Predict turnout
 - ▶ ...
Contemporary story of selection bias



23 Retweets 1 Quote Tweet 30 Likes

Contemporary story of selection bias



 Breaking911
 ~

 WHO WON THE FIRST
 PRESIDENTIAL DEBATE?

 President Donald Trump
 79%

 Former VP Joe Biden
 21%

 76,192 votes · 6 days 23 hours left
 10:42 PM · 9/29/20 · Twitter Web App

 10:42 PM · 9/29/20 · Twitter Web App
 3,579 Retweets
 236 Quote Tweets
 2,861 Likes

Tweet

Echo chamber effect in social media

A bit preview

How to obtain better estimate than Chebyshev:



A bit preview

How to obtain better estimate than Chebyshev:



A bit preview

How to obtain better estimate than Chebyshev: continuous approximation!

