

MATH 180A: INTRO TO PROBABILITY (FOR DATA SCIENCE)

www.math.ucsd.edu/~tkemp/180A

Today: § 4.3-4.4

Next: § 4.4-4.5

HW4 due **TONIGHT** by 11:59 pm

Lab 4 due next **wednesday (Nov 6)** by 11:59 pm

Example

Flip a fair coin n times. How does

$$S_{2.01\%} = 0.5001$$

$$\lim_{n \rightarrow \infty} P\left(\frac{\# \text{Heads}}{n} \geq 50.01\%\right) = 0$$

behave as $n \rightarrow \infty$?

$$\frac{1}{2} + 0.0001$$

$n = 10^6$	$\sqrt{n} = 10^3$
$2\sqrt{n}(0.0001) = 200$	
$1 - \Phi(200) = 0$	
	$\sqrt{n} = 100$
	$n = 10,000$
	$\epsilon = 0.0001$

Suppose after 10,000 flips, there are 5,001 Heads. } $1 - \Phi(0.02)$
 Should we doubt that the coin is really fair? } $\geq 40\%$
 (49%)

$\sqrt{n} = 1000$
 $\Phi(0.2)$
 $\sqrt{}$
 34%
 What if, after 1,000,000 flips, there are 500,100 Heads.
 Now how confident should we be that the coin is really fair?

$$S_n = \# \text{Heads} \sim \text{Bin}(n, \frac{1}{2})$$

$$P\left(\frac{S_n}{n} \geq \frac{1}{2} + \epsilon\right) = P\left(\frac{S_n - \frac{1}{2}n}{\sqrt{\text{Var}(S_n)}} \geq \epsilon\sqrt{n}\right) = P\left(\frac{S_n - \frac{1}{2}n}{\sqrt{n/4}} \geq 2\epsilon\sqrt{n}\right) \approx P(X \geq 2\sqrt{n}\epsilon)$$

replace $\sqrt{\text{Var}(S_n)} = \sqrt{n \cdot \frac{1}{2} \cdot \frac{1}{2}} = \frac{\sqrt{n}}{2}$

Normal \downarrow $N(0,1)$

$$= 1 - P(X < 2\sqrt{n}\epsilon) = 1 - \Phi(2\sqrt{n}\epsilon)$$

Confidence

4.3

Suppose we have a coin that is biased by some unknown amount;

$$X \sim \text{Ber}(p) \quad \text{unknown } p!$$

How can we figure out what p is?

Use the law of large numbers: $p = \lim_{n \rightarrow \infty} \frac{S_n}{n}$

We can't actually wait around for $n \rightarrow \infty$. Instead, we estimate

$$p \approx \hat{p} := \frac{S_n}{n} \quad \text{for some large } n.$$

The question is: how good an estimate is this for given n ?
Or, turning it around: how big must you take n to get an estimate of a certain accuracy?

$$|\hat{p} - p| < \varepsilon \quad (\varepsilon = 0.01)$$

$$P(|\hat{p} - p| < \varepsilon) \geq 95\%$$

" \hat{p} is within margin of error ε of p with probability 95%."

A Maximum Likelihood Estimate

want to find n large enough that (with $\hat{p} = S_n/n$)

$$P(|\hat{p} - p| < \varepsilon) \geq \text{(high probability)}$$

↑
chosen tolerance

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

$$P(|\hat{p} - p| < \varepsilon) = P\left(\frac{S_n - np}{n} < \varepsilon\right) = P\left(\frac{|S_n - np|}{\sqrt{np(1-p)}} < \frac{\varepsilon \sqrt{n}}{\sqrt{p(1-p)}}\right) \approx P(|X| < \frac{\varepsilon \sqrt{n}}{\sqrt{p(1-p)}}) \\ = \Phi\left(\frac{\varepsilon \sqrt{n}}{\sqrt{p(1-p)}}\right) - \Phi\left(-\frac{\varepsilon \sqrt{n}}{\sqrt{p(1-p)}}\right)$$

$$P(|\hat{p} - p| < \varepsilon) \approx 2\Phi\left(\frac{\varepsilon \sqrt{n}}{\sqrt{p(1-p)}}\right) - 1$$

$$0 < p < 1 \\ p(1-p) \leq \frac{1}{4} \\ \text{max @ } p = \frac{1}{2}$$

$$\frac{1}{\sqrt{p(1-p)}} \geq 2 \\ \Phi \uparrow$$

Conclusion: $P(|\hat{p} - p| < \varepsilon) \geq 2\Phi(2\varepsilon\sqrt{n}) - 1$

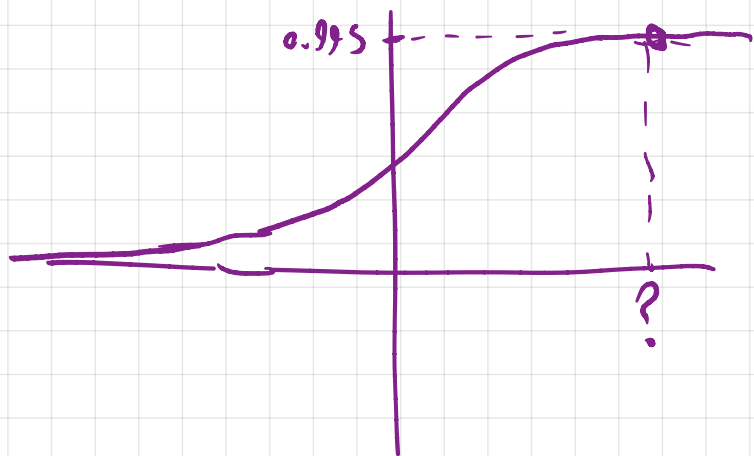
Example: (of the Beast) How many times should we flip a coin, biased an unknown amount p , so that the estimate $\hat{p} = S_n/n$ is within a tolerance of 0.05 of the true value p , with probability $\geq 99\%$?

Want n large enough that

$$P(|\hat{p} - p| < 0.05) \geq 99\%$$

makes sense

We know $P(|\hat{p} - p| < 0.05) \approx 2\Phi(2(0.05)\sqrt{n}) - 1 \geq 99\%$



$$\Phi(2(0.05)\sqrt{n}) \geq 0.995$$

$$\therefore 2(0.05)\sqrt{n} \geq 2.58$$

$$\sqrt{n} \geq 25.8$$

$$n \geq 665.64$$

666

Confidence Intervals

Turning this around: if we can't control n , we would like to say **how accurate** the sample mean is as an estimate of the true mean, for a given number n of samples.

Eg. A coin (of unknown bias p) is tossed 1000 times. 450 Heads come up. Within what tolerance can we say we know the true value of p with probability $\geq 95\%$?

$$\text{Estimate } p \approx \hat{p} = \frac{S_{1000}}{1000} = 0.45$$

$$\text{Want } P(|p - \hat{p}| < \varepsilon) \geq 95\%$$

$$\text{Know: } P(|p - \hat{p}| < \varepsilon) \approx 2\Phi(2\varepsilon\sqrt{1000}) - 1 \geq 0.95$$

$$\Phi(2\varepsilon\sqrt{1000}) \geq 0.975$$

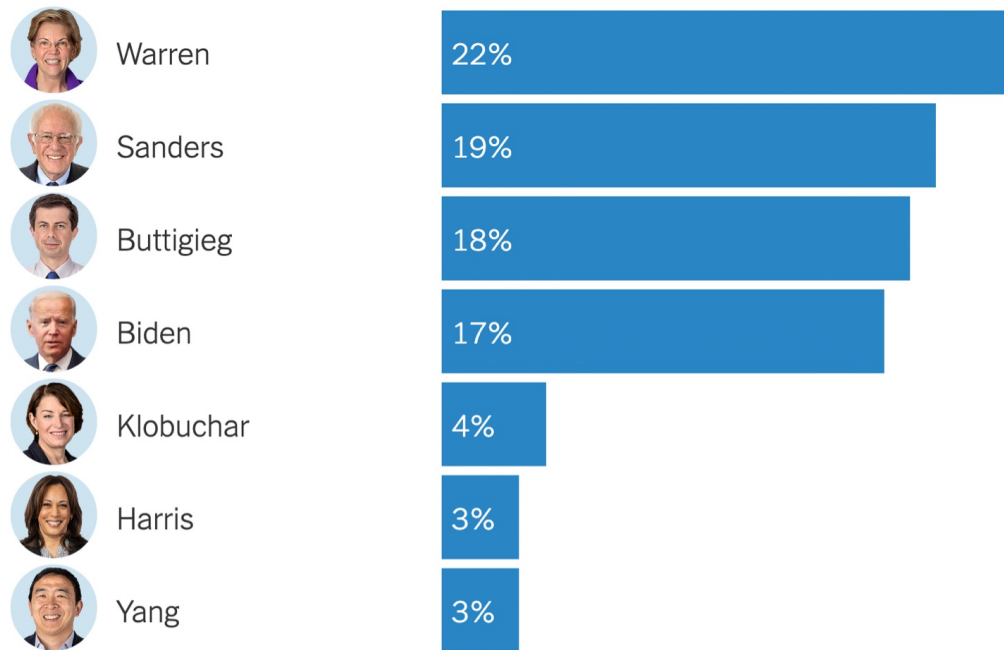
$$\left. \begin{array}{l} \text{I.e. } |p - 0.45| < 0.031 \text{ w } P \geq 95\% \\ 0.45 - 0.031 < p < 0.45 + 0.031 \end{array} \right\} \begin{array}{l} 2\varepsilon\sqrt{1000} \geq 1.96 \leadsto \varepsilon \geq \frac{1.96}{2\sqrt{1000}} \doteq 0.031 \\ p \in [0.419, 0.481] \text{ w } P \geq 95\% \\ \text{95\% confidence interval} \end{array}$$

If an experiment is repeated in many independent trials,
and the preceding (normal approximation) estimates yield

$$P(|\hat{p} - p| < \varepsilon) \geq 95\%$$

we say $[\hat{p} - \varepsilon, \hat{p} + \varepsilon]$ is the 95% confidence interval for p .

The same statement might be given as " $p = \hat{p}$ with margin of error ε
(95 times out of 100)".



Source: New York Times Upshot/Siena College poll conducted Oct. 25-30.

Poll conducted Oct 25-30
of 439 Iowa Democratic
caucusgoers.

$$P(|p - \hat{p}| < \varepsilon) \geq 2\Phi\left(\frac{2\varepsilon\sqrt{439}}{0.22}\right) - 1$$

$$(\approx) \geq 0.95$$

$$\text{i.e. } 2\varepsilon\sqrt{439} \geq 1.96$$

$$\varepsilon \geq 4.68\%$$

Margin of error: 4.7%

Poisson Approximation

4.4

$$S_n \sim \text{Bin}(n, \lambda/n) : \lim_{n \rightarrow \infty} \mathbb{P}(S_n = k) = e^{-\lambda} \frac{\lambda^k}{k!}$$

Quantitative Bound:

Theorem: If $X \sim \text{Bin}(n, p)$ and $Y \sim \text{Poisson}(np)$,
for any subset $A \subseteq \mathbb{N}$

$$|\mathbb{P}(X \in A) - \mathbb{P}(Y \in A)| \leq np^2$$

Upshot: if np^2 is small, use Poisson Approximation.
if $np(1-p)$ is big, use Normal Approximation.