# MATH 180A: Introduction to Probability

Lecture B00 (Nemish)

www.math.ucsd.edu/~ynemish/teaching/180a

Lecture C00 (Au)

www.math.ucsd.edu/~bau/f20.180a

## Today: Confidence Intervals.
## Poisson Approximation

## Next:   ASV 4.5

Video: Prof. Todd Kemp, Fall 2019

Week 6:

- Homework 5 (due Friday, November 13, 11:59 PM)

# Example

Flip a fair coin $n$ times. How does

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

behave as $n \to \infty$ ?

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

(right side annotations)

$50.01\% = 0.5001$

$n = 10^{10} \quad \sqrt{n} = 10^5$

$2\sqrt{n}(0.0001) = 200$

$1 - \Phi(200)$
$\to 0$

$\sqrt{n} = 100$
$n = 10,000$
$\varepsilon = 0.0001$

Suppose after 10,000 flips, there are 5,001 Heads.

Should we doubt that the coin is really fair?

$\left.\begin{array}{l} \\ \\ \end{array}\right\} 1 - \Phi(0.02)$
$\geq 40\%$
$(49\%)$

$\sqrt{n} = 1000$

$\Phi(0.2)$
$1 - $
$\geq 46\%$

What if, after 1,000,000 flips, there are 500,100 Heads.

Now how confident should we be that the coin is really fair?

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

$$\mathbb{P}\left(\frac{S_n}{n} \geq \frac{1}{2} + \varepsilon\right) = \mathbb{P}\left(\frac{S_n - \frac{1}{2}n}{n} \geq \varepsilon\right) = \mathbb{P}\left(\frac{S_n - \frac{1}{2}n}{\sqrt{n}/2} \geq 2\varepsilon\sqrt{n}\right) \approx \mathbb{P}(X \geq 2\sqrt{n}\varepsilon)$$

Normal. $\mathcal{N}(0,1)$

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

$= 1 - \mathbb{P}(X < 2\varepsilon\sqrt{n})$

$= 1 - \Phi(2\varepsilon\sqrt{n})$

# Confidence

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

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

How can we figure out what $p$ is?

Use the law of large numbers: $p = \lim\limits_{n \to \infty} \dfrac{S_n}{n}$

We can't actually wait around for $n \to \infty$. 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?

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

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

"$\hat{p}$ is within margin of error $\varepsilon$ of $p$, w probability 95%."

# A Maximum Likelihood Estimate

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

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

$\uparrow$ chosen tolerance

$X \sim \mathcal{N}(0, 1)$

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

$\frac{S_n}{n}$   $\sqrt{np(1-p)}$

$$= \Phi\left(\frac{\varepsilon\sqrt{n}}{\sqrt{p(1-p)}}\right) - \Phi(-\,\text{"}\,\text{"})$$

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

$0 < p < 1$

$\qquad \dfrac{1}{\sqrt{p(1-p)}} \geq 2$

$p(1-p) \leq \frac{1}{4}$

$\max @ p = \frac{1}{2}$   $\Phi \uparrow$

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

**Example:** How many times should we flip a coin, biased an unknown
**(of the Beast)** 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

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

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



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

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

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

$$n \geq 665.64$$

**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 $\geq 95\%$?

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

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

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

$$\Phi(2\varepsilon\sqrt{1000}) \geq 0.975$$
$$2\varepsilon\sqrt{1000} \geq 1.96 \rightsquigarrow \varepsilon \geq \frac{1.96}{2\sqrt{1000}} \doteq 0.031$$

[i.e. $|p-0.45| < 0.031$ w/ $\mathbb{P} \geq 95\%$

$0.45 - 0.031 < p < 0.45 + 0.031$ ]

$p \in [0.419, 0.481]$ w/ $\mathbb{P} \geq 95\%$

$\uparrow$ 95% confidence interval

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

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

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

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



Warren 22%

Sanders 19%

Buttigieg 18%

Biden 17%

Klobuchar 4%

Harris 3%

Yang 3%

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

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

$$\mathbb{P}(|p-\hat{p}| < \varepsilon) \geqslant 2\Phi(2\varepsilon\sqrt{439})-1$$
$$\underset{0.22}{} \quad (\approx) \quad \geqslant 0.95$$

ie. $2\varepsilon\sqrt{439} \geqslant 1.96$

$\varepsilon \geqslant 4.68\%$

Margin of error: 4.7%

## Theorem. Let $S_n \sim \text{Bin}(n,p)$
$X \sim \text{Poisson}(np)$
$Y \sim N(0,1)$

if $p = \frac{\lambda}{n^{0.51}}$ $(\lambda > 0)$

$np^2 = n \left(\frac{\lambda}{n^{0.51}}\right)^2 = \frac{\lambda^2 n}{n^{1.02}}$

$\to 0$
as $n \to \infty$

For any subset $A \subseteq \mathbb{N}$,

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

OTOH, for any $x \in \mathbb{R}$,

3 is not optimal

Berry -Esseen Thm.

$$\left| \mathbb{P}\left( \underbrace{\frac{S_n - np}{\sqrt{np(1-p)}} \leq x}_{\text{CDF of } \frac{S_n - np}{\sqrt{np(1-p)}}} \right) - \underbrace{\mathbb{P}(Y \leq x)}_{\Phi(x)} \right| \leq \frac{3}{\sqrt{np(1-p)}}$$

$\leftarrow$ optimal

Upshot: if $np^2$ is small, use Poisson Approximation.

if $np(1-p)$ is quite large, use Normal Approximation.

Beyond independent trials:

* The normal approximation breaks down qickly if the trials are dependent.

* The Poisson approximation holds up well under "weak dependence"

Example. A factory experiences 3 accidents per month, on average. What is the probability there will be 3 accidents this month?

$X = \#$ accidents in a given month.

well modeled by a Poisson.

$X \sim \text{Poisson}(\lambda)$

$3 = \mathbb{E}(X) = \lambda$

$$P(X = 3) = e^{-3} \frac{3^3}{3!} \doteq 22.4\%$$

$$\frac{3^3}{3!} = \frac{3^2 \cdot 3}{3 \cdot 2 \cdot 1} = \frac{3^2}{2!}$$

$$P(X = 2) = e^{-3} \frac{3^2}{2!} \doteq 22.4\%$$