

PSTAT 115

Lecture 1: Introduction

Professor Alexander Franks

1/8/24

Class Resources

Required Textbook

- Bayes Rules: <https://www.bayesrulesbook.com/>

Course Pages

Canvas

- Class website on Canvas:
<https://www.canvas.ucsb.edu/>
- Nectir for course related questions and discussion:
[https://ucsb.nectir.io/group/pstat115-w2³₄](https://ucsb.nectir.io/group/pstat115-w23)
 - On Canvas site
- Gradescope: TBA
 - Linked on Canvas site

Grades

- 35% - ~~expect~~ 5 homeworks

- 20% - Midterm (February ~~9~~)

- 10% - Quizzes

- 5% - Participation

- 30% - Final exam (March ~~21~~)

(see catalog)

Homework

- There will be 5 homeworks (35% of your grade total)
- You will typically have 1-2 weeks to complete the homeworks
- You are allowed to work with a partner
 - Add partners name to your assignment
- Every student *must* submit their own assignment on gradescope
- Homework turned in within 24 hrs after the deadline without prior approval will receive a 10 pt deduction (out of 100)
- Homework will not be accepted more than 24 hrs late.

Homework submission format

- All code must be written to be reproducible in Rmarkdown
- All derivations can be done in any format of your choosing (latex, written by hand) but must be legible and *must be integrated into your Rmarkdown pdf.*
- All files must be zipped together and submitted to Gradescope
- Ask a TA *early* if you have problems regarding submissions.

Software and Deliverables

Software

- R ([R studio](#))

Homeworks submission format

- Electronic submission via Gradescope
- [R markdown code](#)
- Generated PDF file
- Any supplementary files (e.g. write up for math problems)

Labs and Quizzes

Section

- There will be a handful of in ~~class~~ quizzes throughout the quarter.
- You will have 5 minutes to take the quiz
- There are no makeups, but the lowest quiz grade will be dropped from your final score.
- Quizzes (10%) will be multiple choice and will test your comprehension of the basic concept.
- Participation (5%). Includes lecture attendance, section attendance, and nectir posts.

RStudio Cloud Service

- Log on to pstat115.lsit.ucsb.edu
 - Cloud based rstudio service
 - Log in with your UCSB NetID
- Use <https://bit.ly/48CD68y> to sync new material
(BOOKMARK THIS)
- Make sure you can write and compile an [R markdown](#) ([Rmd](#)) document online
- Text formatting is minimal but [syntax](#) is simple

Class Policies

- All questions should be posted on nectir, *not by email* (unless they are personal or grade-related)

Artificial intelligence

- LLMs (ChatGPT etc) are allowed BUT...
- The less you know the more likely you are to be convinced by misinformation
- It's not about getting the right answer
 - “It's the journey, not the destination”
- Ask yourself “Am I using it to avoid work? Or am I using it to help me develop an understanding?”

Markdown and mathematical formulas

The text inserted between two \$ signs will be interpreted as a Latex instruction, e.g. \$x\$

Code	Rendered math
\$x\$	x
\$\theta\$	θ
\$x_i^2\$	x_i^2
$\frac{1}{n} \sum_{i=1}^n x_i$	$\frac{1}{n} \sum_{i=1}^n x_i$

Code

```
$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2
```

Rendered math

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Rmarkdown and Latex resources

- [Introduction to RMarkdown](#)
- [Latex cheat sheet](#)
- [Introduction to Latex](#)

Other R resources

A different approach
(not "frequentist")

Bayes Theorem

What is Bayesian statistics?

Prior Info

Regression.

Sampling.

C.I.

**What is the version of
statistics you already
know?**

Hypothesis Testing

Likelihood

Frequentist

Frequentist statistics

- Associated with the *frequentist* interpretation of probability
 - For any given event, only one of two possibilities may hold: it occurs or it does not.
 - The *frequency* of an event (in repeated experiments) is the *probability* of the event

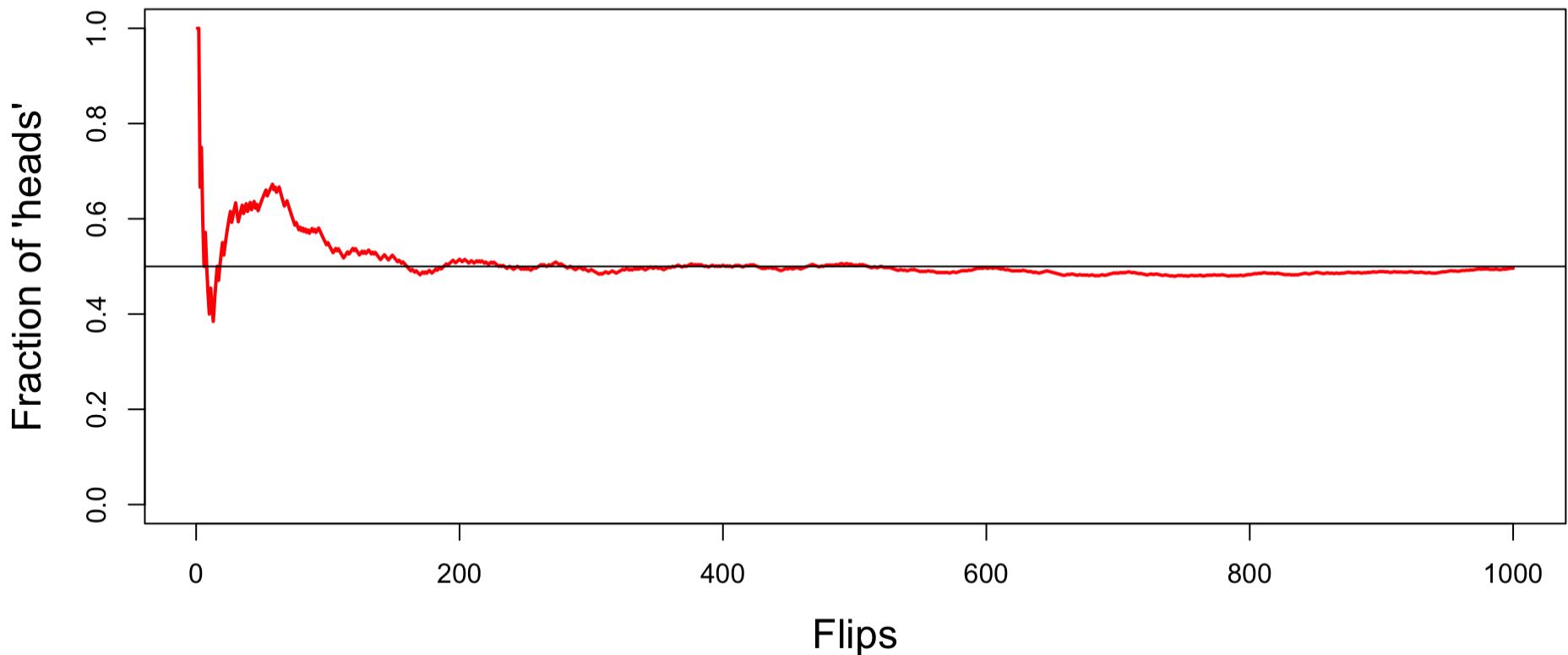
Frequentist statistics

- Associated with the *frequentist* interpretation of probability
 - For any given event, only one of two possibilities may hold: it occurs or it does not.
 - The *frequency* of an event (in repeated experiments) is the *probability* of the event

- Null Hypothesis Significant Testing (NHST) and Confidence Intervals
 - Frequentist uncertainty premised on imaginary resampling of data
 - Example: If the null model is true, and I re-run the experiment many times, how often will I reject?

Frequentist probability

The probability of a coin landing on heads is 50%

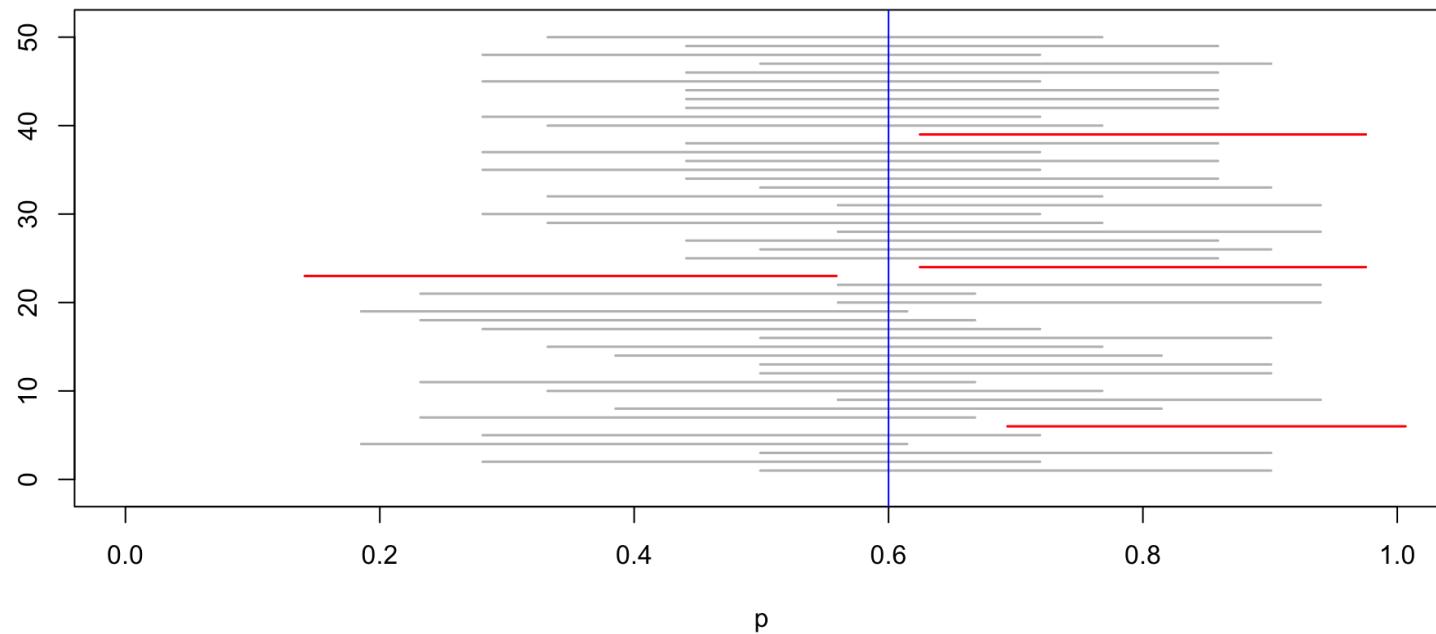


The long run fraction of heads is 50%

Confidence intervals

I have a 95% confidence interval for a parameter θ .
What does this mean?

*weight of
coin*



We expect $0.05 \times 50 = 2.5$ of the intervals to *not* cover the true parameter, $p = 0.6$, on average

Falsification



H_0 : "All swans are white"

H_a : not

Falsification



H_0 : “The Ivory-billed Woodpecker is extinct”

Falsification

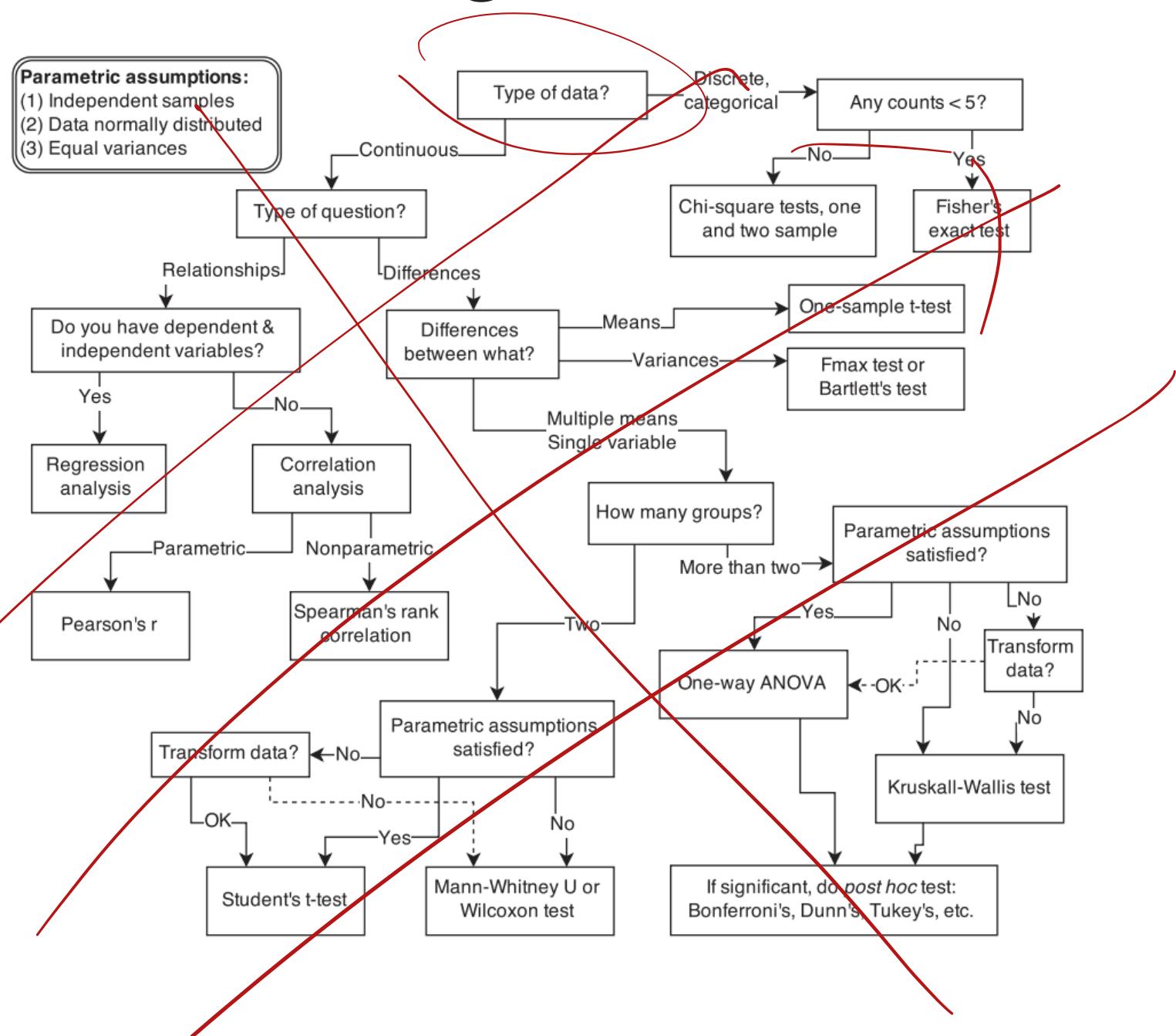


H_0 : “Black swans are rare”

Falsification

- Is an observation real or spurious?
 - Importance of measurement error
 - Natural phenomena are usually continuous in nature
- Falsification requires consensus more than logic
 - Scientific communities argue toward consensus
 - Science is messy!

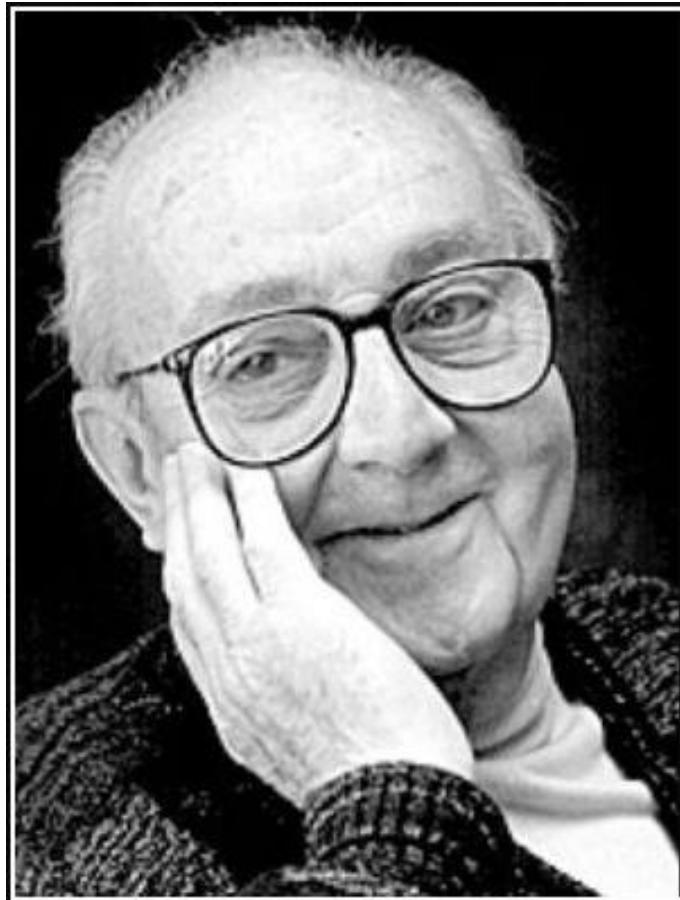
Significance Testing Flowchart



Alternative: focus on modeling!

- A statistical model represents a set of assumption about how the data was generated.
- Models can still be used to develop statistical tests.
- Can also be used to make predictions or forecasts and describe sources of variability.
- Can (and should) be continuously refined and extended!

All models are wrong

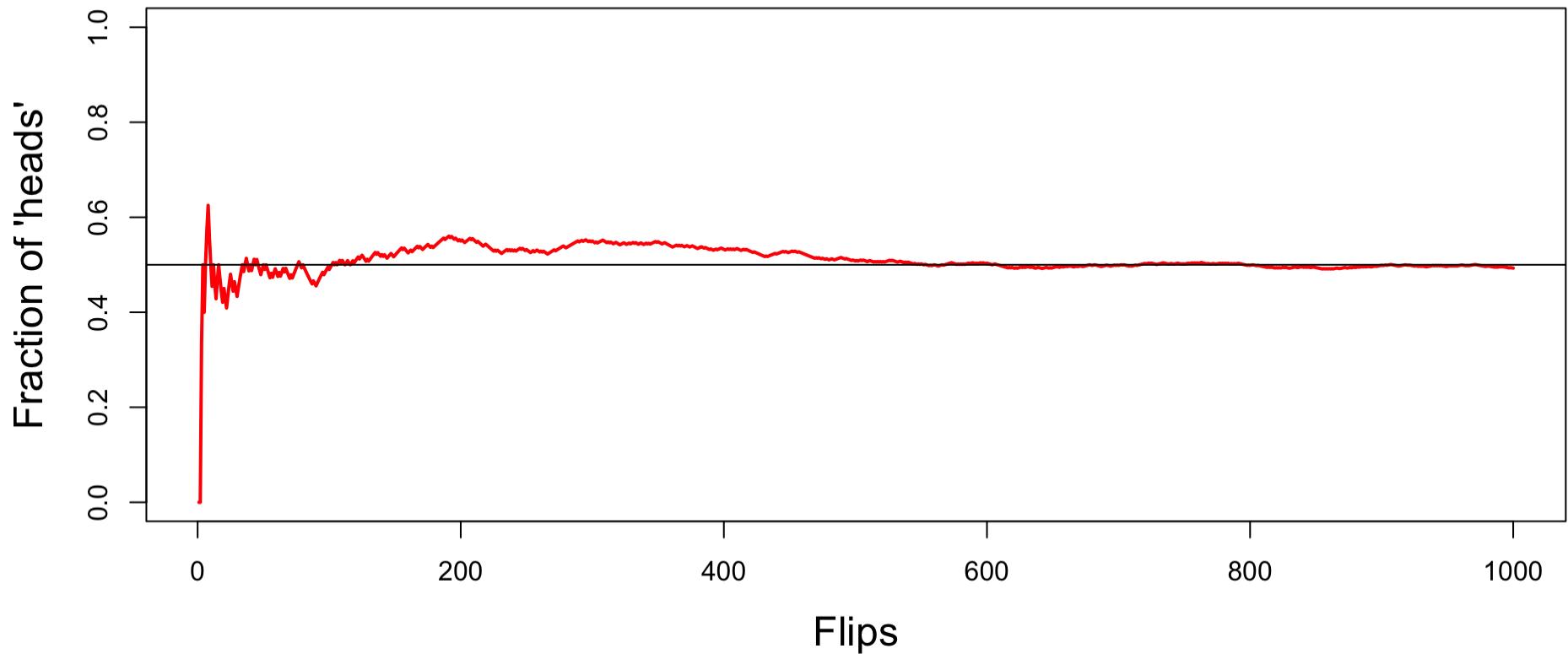


All models are wrong, but some are useful.

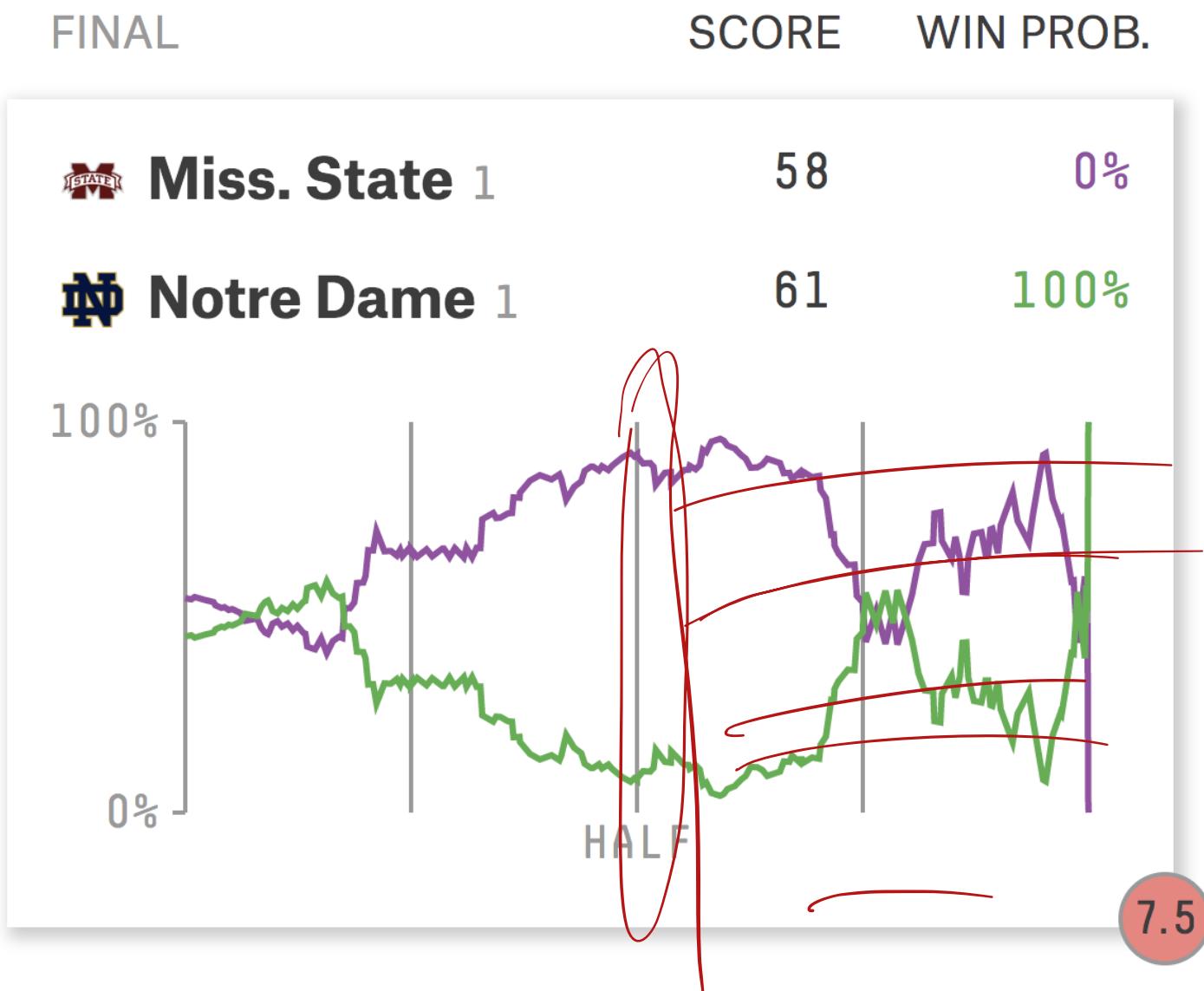
— *George E. P. Box* —

https://en.wikipedia.org/wiki/All_models_are_wrong

Frequentist probability

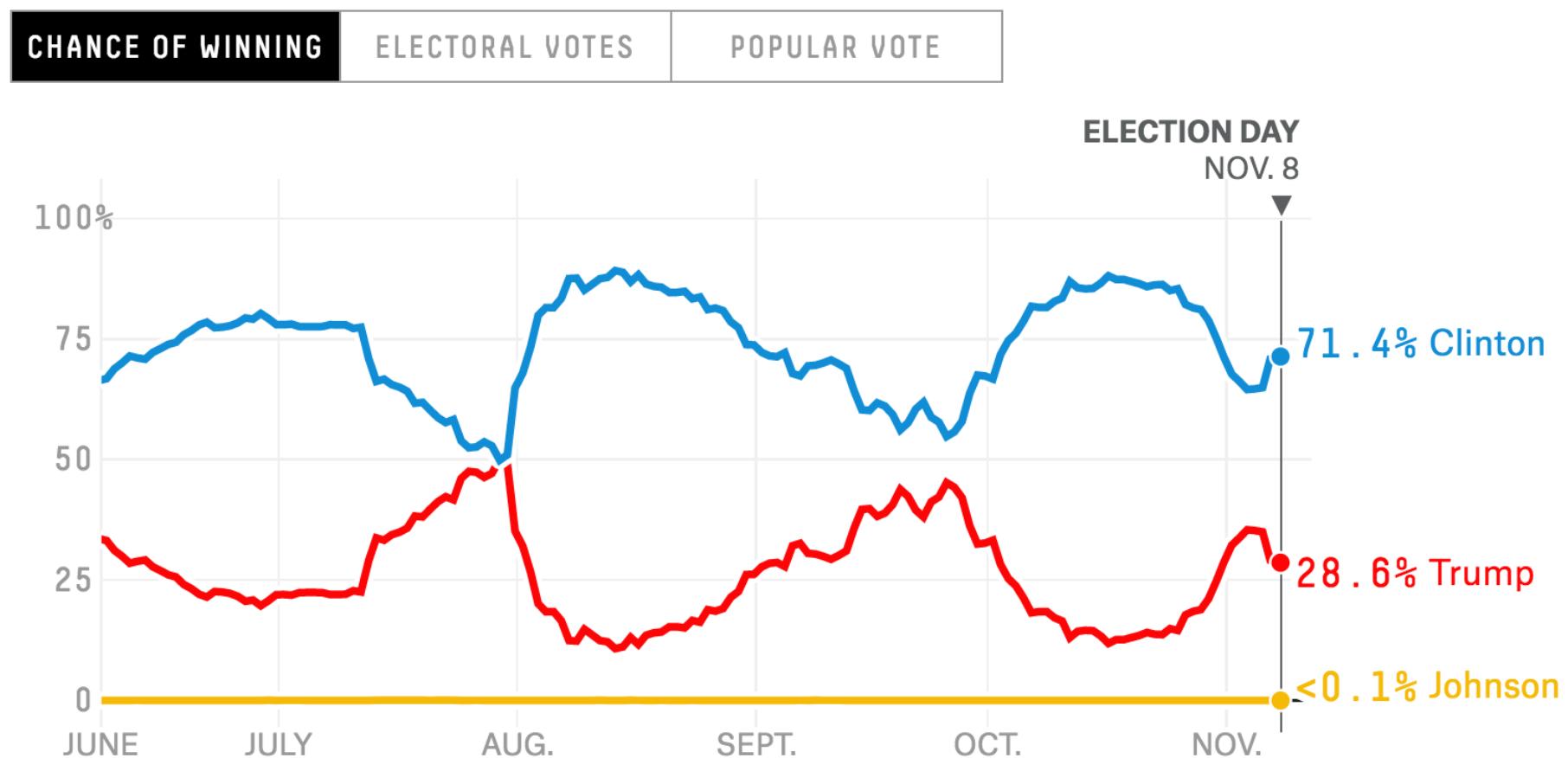


Win probability



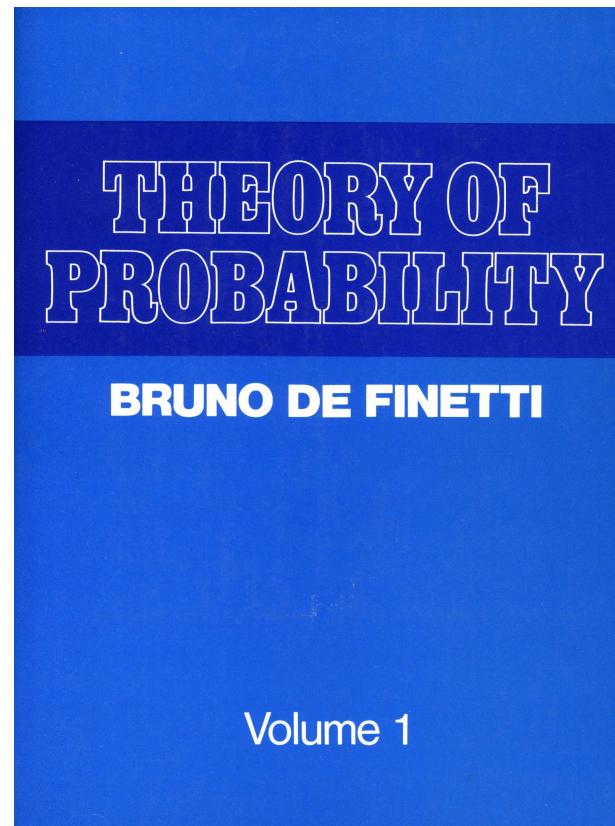
source: fivethirtyeight.com

Win probability



source: fivethirtyeight.com

Bayesian probability



Bruno de Finetti began his book on probability with:
“PROBABILITY DOES NOT EXIST”

Bayesian probability

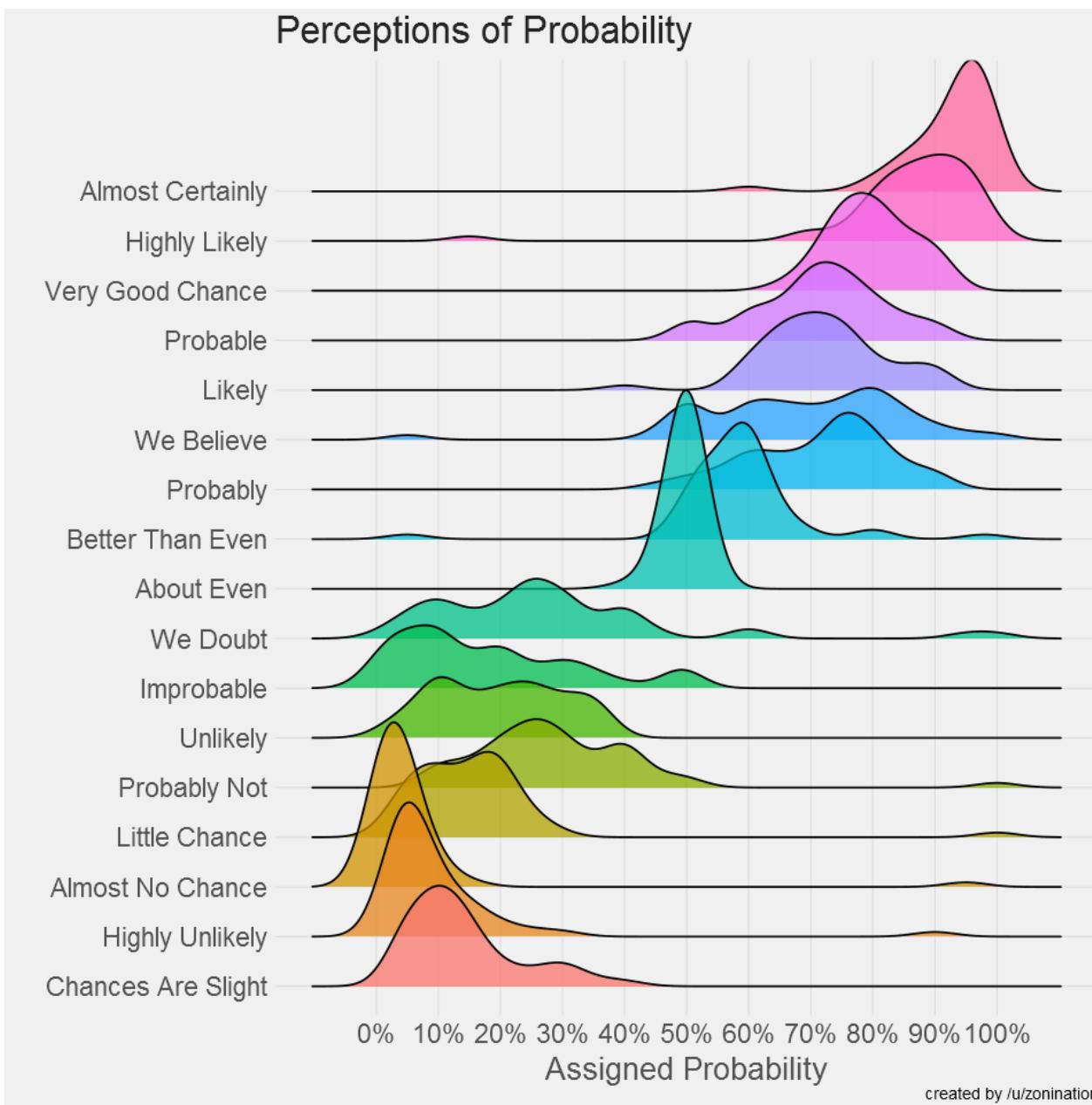
- de Finetti is arguing that probability is about *belief*
 - Probability doesn't exist in an *objective* sense
 - “The coin is fair” means *I believe* that its equally likely to be heads or tails.
 - “Hillary Clinton has a 71% chance to win” reflects a belief, since the election happens only once
- Rarely, if ever, get *true* replications to estimate frequentist probabilities
- Bayesian idea: focus statistical practice around belief about parameters

Bayesian probability

“The terms *certain* and *probable* describe the various degrees of rational belief about a proposition which different amounts of knowledge authorise us to entertain. All propositions are true or false, but the knowledge we have of them depends on our circumstances

— John M Keynes

Perceptions of Probability



Why Bayesian statistics?

- Classical Stats toolbox too rigid
 - What if assumptions of test fail?
- Bayesian Statistics gives us a procedure for building our tests/tools given a model.
 - Design, Build, Refine...
- Philosophy.
- Powerful tools + Computer Simulation

Setup

- The sample space \mathcal{Y} is the set of all possible datasets.
 - Y is a random variable with support in \mathcal{Y}
 - We observe one dataset y from which we hope to learn about the world.
- The parameter space Θ is the set of all possible parameter values θ
- θ encodes the population characteristics that we want to learn about!

weight of the coin.

Prob (120A)

Sample
 $P(Y|\theta)$

Population
 θ

data

y

$\hat{\theta}(y)$
(estimate)

Stats (120B)

$\hat{\theta}(Y)$ (estimators)
are random.

Three steps of Bayesian data analysis

1. Construct a plausible probability model governed by parameters θ $P(\gamma | \theta)$

- This includes specifying your belief about θ before seeing data (*the prior*) $P(\theta)$

2. Condition on the observed data and compute *the posterior distribution for θ*

3. Evaluate the model fit, revise and extend. Then repeat.

Bayesian Inference in a Nutshell

1. The *prior distribution* $p(\theta)$ describes our belief about the true population characteristics, for each value of $\theta \in \Theta$.
2. Our *sampling model* $p(y | \theta)$ describes our belief about what data we are likely to observe if θ is true.
3. Once we actually observe data, y , we update our beliefs about θ by computing the *posterior distribution* $p(\theta | y)$. We do this with Bayes' rule!



Key difference: θ is random!

Bayes' Rule

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

$P(B \mid A)P(A)$
 $\sim P(B)$

$P(B \mid A)$ $P(A)$

- $P(A \mid B)$ is the conditional probability of A given B
- $P(B \mid A)$ is the conditional probability of B given A
- $P(A)$ and $P(B)$ are called the marginal probability of A and B (unconditional)

Bayes' Rule for Bayesian Statistics

$$P(\theta \mid y) = \frac{P(y \mid \theta)P(\theta)}{\cancel{P(y)}}$$

- $P(\theta \mid y)$ is the posterior distribution
- $P(y \mid \theta)$ is the likelihood
- $P(\theta)$ is the prior distribution
- $P(y) = \int_{\Theta} p(y \mid \tilde{\theta})p(\tilde{\theta})d\tilde{\theta}$ is the model evidence

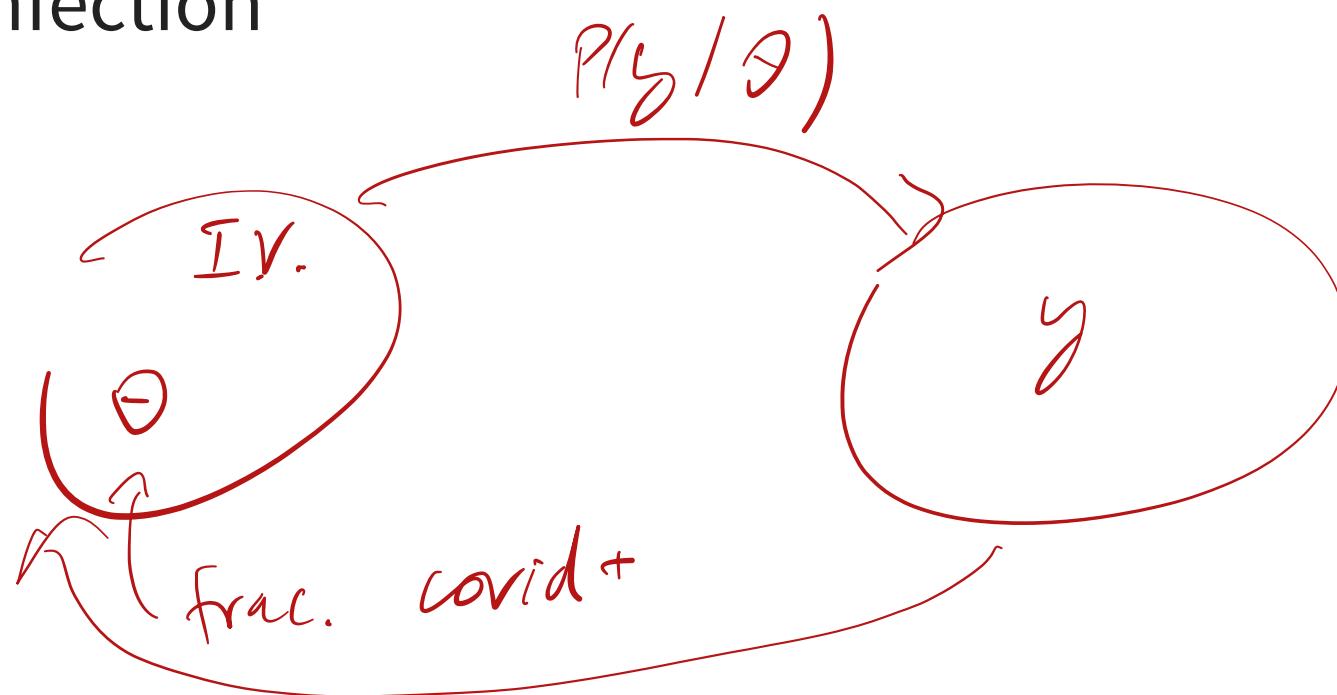
Bayes' Rule for Bayesian Statistics

$$P(\theta \mid y) = \frac{P(y \mid \theta)P(\theta)}{P(y)} \propto P(y \mid \theta)P(\theta)$$

- Start with a subjective belief (prior)
- Update it with evidence from data (likelihood)
- Summarize what you learn (posterior)

Example: Estimating COVID Infection Rates

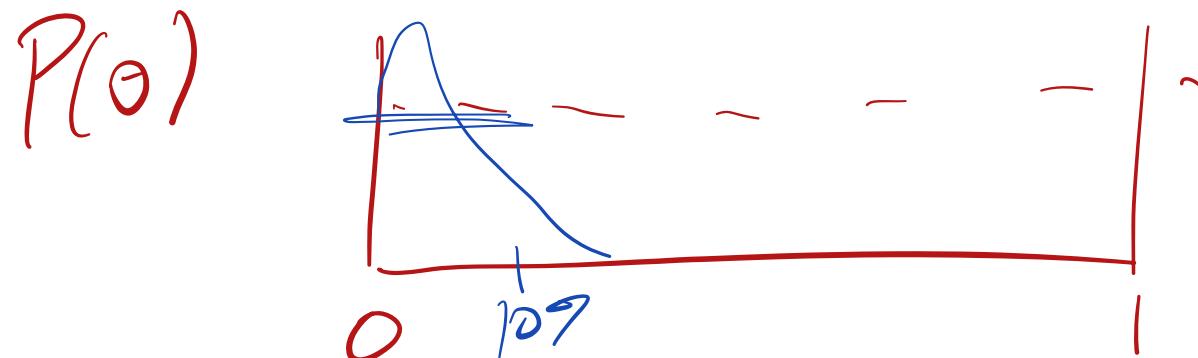
- We need to estimate the prevalence of a COVID in Isla Vista
- Get a small random sample of 20 individuals to check for infection



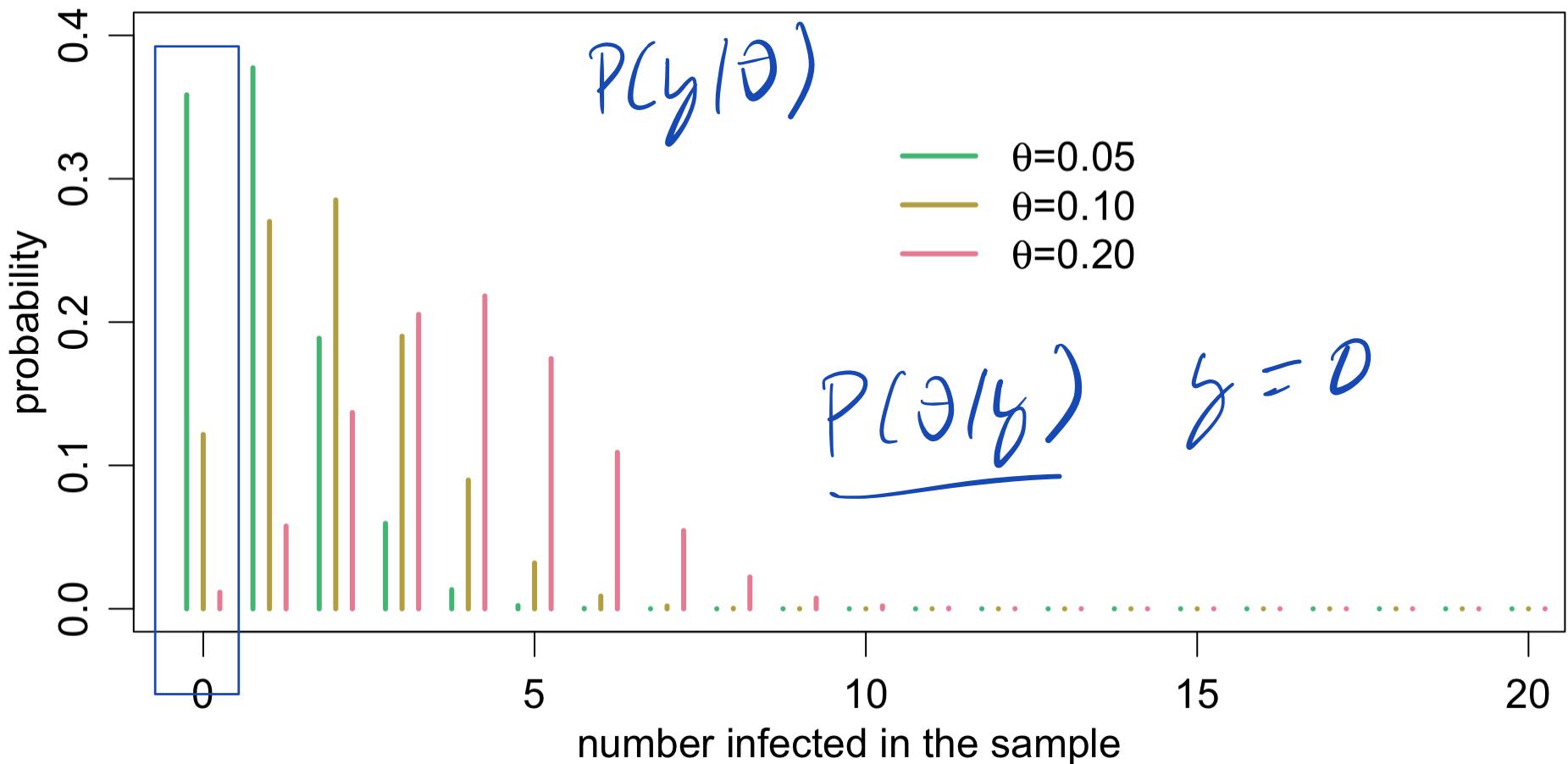
Example: Estimating Infection Rates

- θ represents the population fraction of infected
- Y is a random variable reflecting the number of infected in the sample
- $\Theta = [0, 1]$ $\mathcal{Y} = \{0, 1, \dots, 20\}$
- Sampling model: $Y \sim \text{Binom}(20, \theta)$

positive



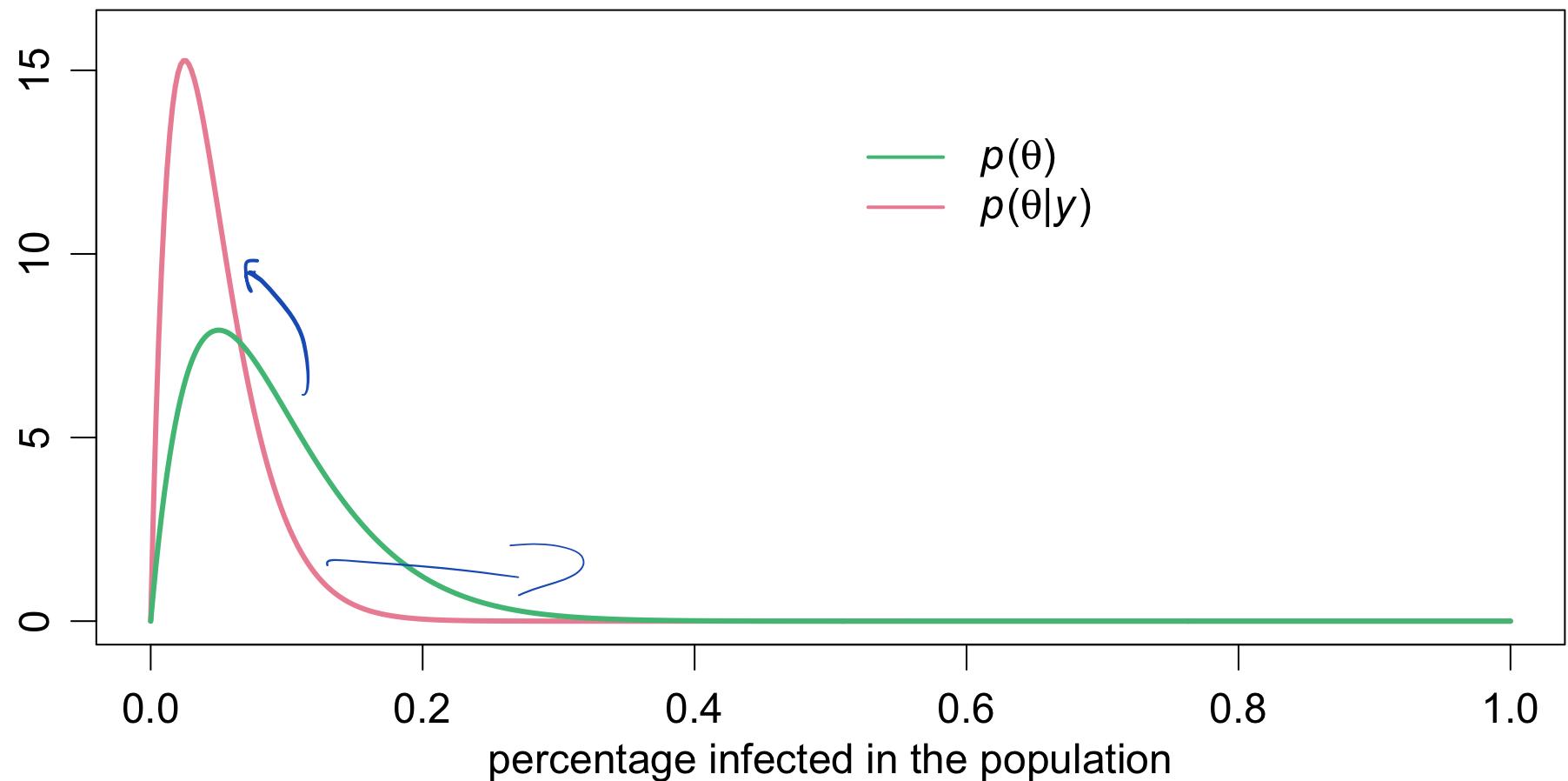
Example: Estimating Infection Rates



Example: Estimating Infection Rates

- Assume *a priori* that the population rate is low
 - The infection rate in comparable cities ranges from about 0.05 to 0.20
- Assume we observe $Y = 0$ infected in our sample
- What is our estimate of the true population fraction of infected individuals?

Example: Estimating Infection Rates



Tentative syllabus

- One parameter models (binomial, poisson, and normal)
- Monte Carlo methods (i.e. simulation-based inference)
- Markov chain Monte Carlo (MCMC)
- Hierarchical modeling

Assignment

- Check Nectir
- Lab starts Wednesday (2, 3, 4, and 6pm)
- Start reviewing probability cheat sheet!
- Read chapters 1 and 2 of Bayes Rules

