

Lecture 1: Review and Background

Logistics

- First homework out
 - Due October 12 (Sunday)
- Use tinyurl.com/pstat115
 - Cloud based rstudio service
 - Log in with your UCSB NetID
 - Syncs all class content
- Canvas website for syllabus

Resources

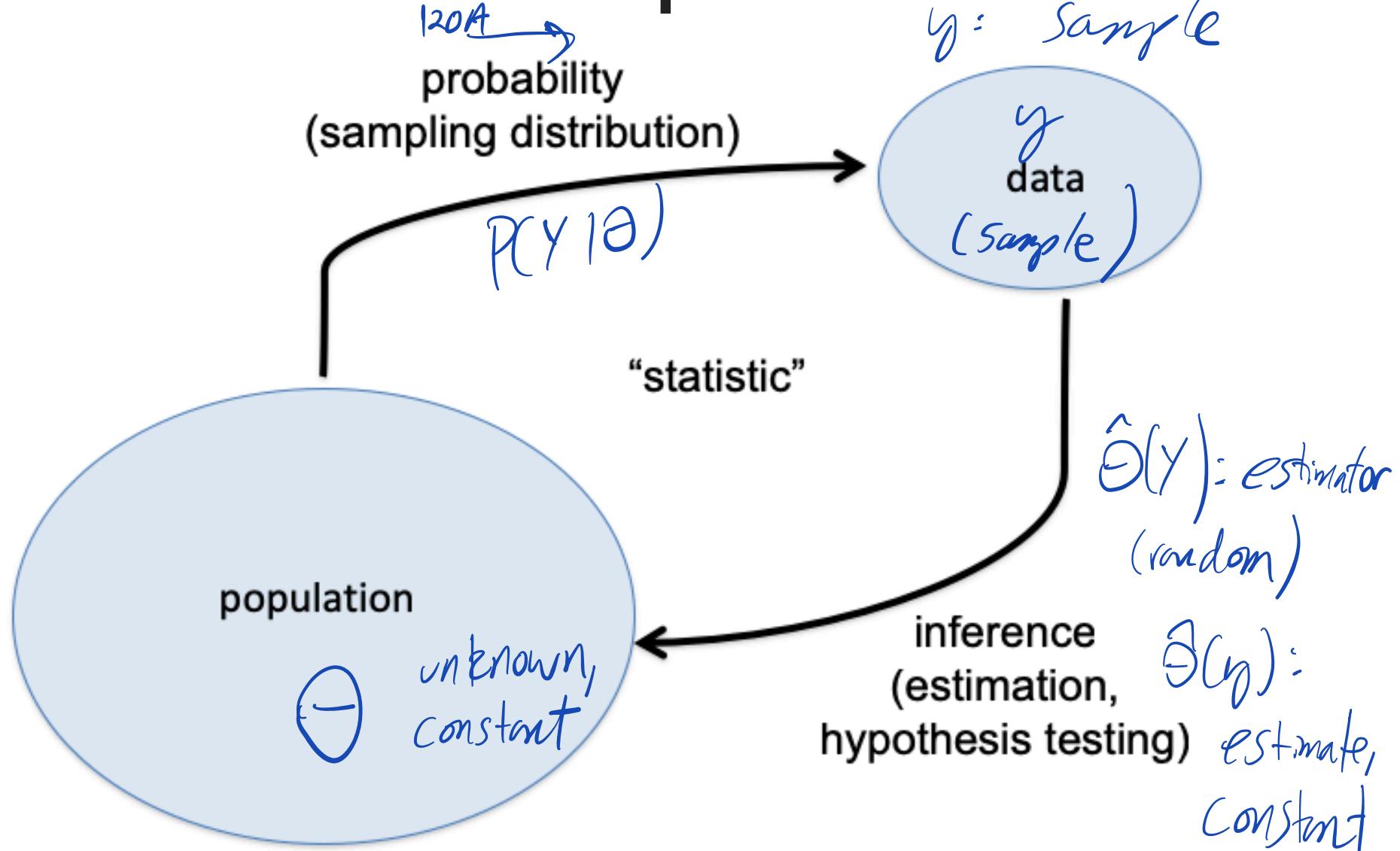
Look at the resources folder in cloud for

- A fantastic probability review sheet (120A)
- Probability density information
- Bayes Rules: Chapters 1 and 2

Homework 1

- Homework 1 out `fall25/homework/homework1.qmd`
- Do not change the name of the file or the directory
- Checking as you go
 - Leave code cells that look like `. = ottr::check("tests/q1a.R")`
 - If these cells fail you have an error in your code, fix this before proceeding

Population and Sample



Population and Sample

- The *population* is the group or set of items relevant to your question
 - Usually very large (often conceptualize a population as infinite)
- Sample: a finite subset of the population
 - How is the sampling collected (representative?)
 - Denote the sample size with n

Modeling

Population and Sample

- Our goal is (usually) to learn about the population from the sample
 - Population parameters encode relevant quantities
 - The **estimand** is the thing we want to infer and is usually a function of the population parameters

denote θ

Random variables

- A random variable, \underline{Y} has variability, can take on several different values (possibly infinitely many), and is associated with a distribution.

- The distribution determines the probability that the r.v. will take a specific value.

- Notation:

- \underline{Y} (uppercase) denotes a random variable

- y (lowercase) is a realization of that random variable and is not random

$$Y \sim \text{Bin}(n, \theta)$$

|| ||
10 0.5

$y = 5$ obs.
constant.

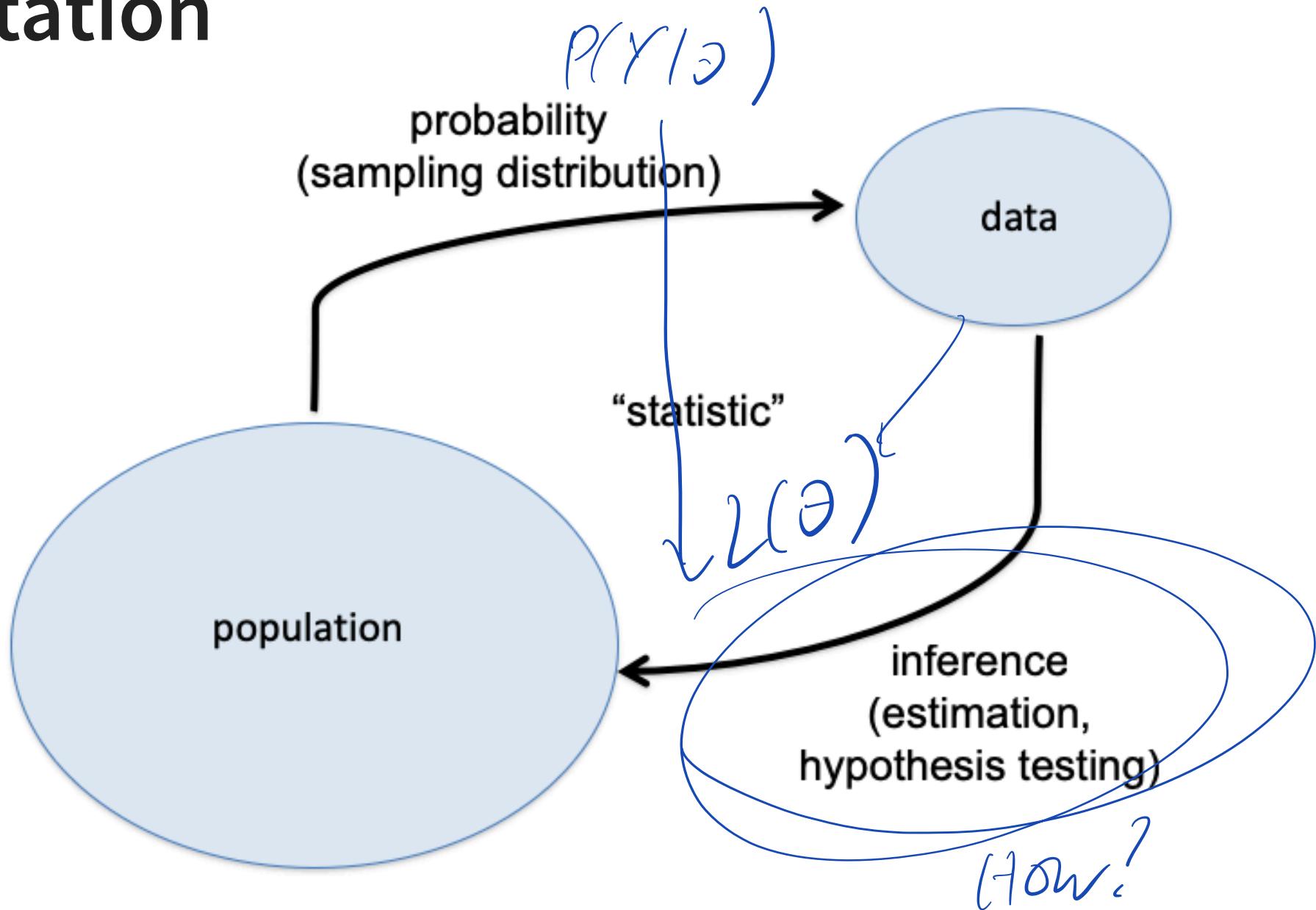
Constants

- Constants: quantities with 0 variance.
 - Constants can be *known* (e.g. observed data)
 - Constants can be *unknown* (not observed)
 - θ (Fair coin?)
"weight of the coin"
Constant in Frequentist
 - $y = 5$ heads.

Setup

- The *sample space* \mathcal{Y} is the set of all possible datasets we could observe. 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
- Our *sampling model* $p(y \mid \theta)$ describes our belief about what data we are likely to observe for a given value of θ .

Notation



The Likelihood Function

(density)

- The likelihood is the “probability of the observed data” expressed as a function of the unknown parameter:
- A function of the unknown constant θ .
- Depends on the observed data $y = (y_1, y_2, \dots, y_n)$

$$L(\theta) \propto P(y_1, \dots, y_n | \theta)$$

↑
obs ↗ unknown,

Likelihood is not a distribution/density.

Independent Random Variables

- Y_1, \dots, Y_n are random variables
- We say that Y_1, \dots, Y_n are *conditionally* independent given θ if
- Conditional independence means that Y_i gives no additional information about Y_j beyond that in knowing θ

$$P(Y_1, Y_2, \dots, Y_n | \theta) = \prod_{i=1}^n P(Y_i | \theta)$$

Example: A binomial model

- Assume I go to the basketball court and takes 5 free throw shots
- Model the number of made shots I make using a $\text{Bin}(5, \theta)$
 - What are the key assumptions that make these a reasonable emodel?
- θ represents my ¹true skill ¹(the fraction of shots I make)
- How can we estimate my true skill?

Likelihood:

→ (If I take ∞
many shots)

$$L(\theta) = P(Y|\theta) = \binom{5}{y} \theta^y (1-\theta)^{5-y}$$

$\cancel{\binom{5}{y}}$

No θ 's

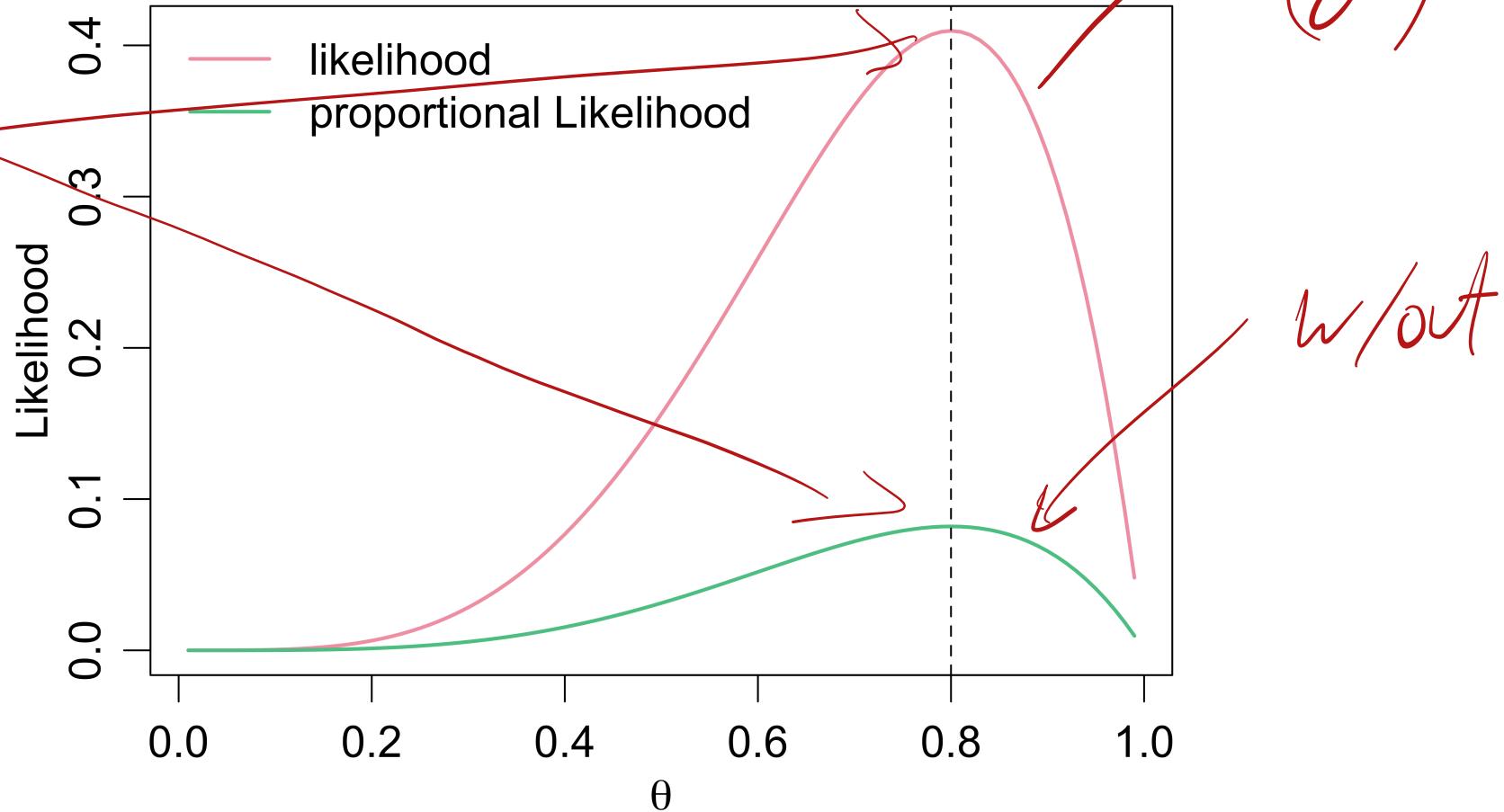
Obs $y = y$

$$L(\theta) \propto \theta^4 (1-\theta)^{5-4}$$

The binomial likelihood

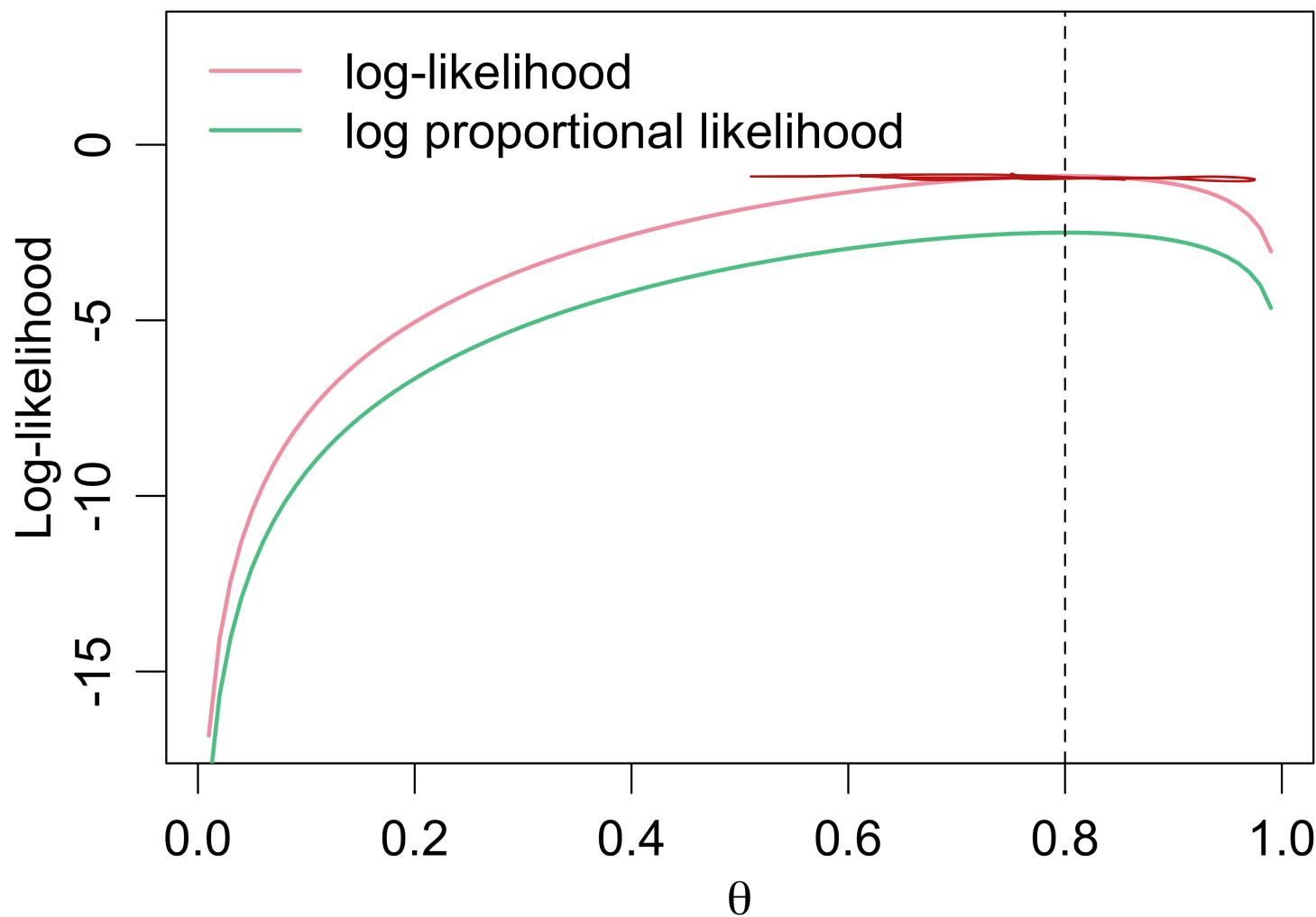
I make 4 out of 5

Max
doesn't
change.



The log-likelihood

$$l(\theta) = \log(L(\theta))$$



Maximum Likelihood Estimation

- The *maximum likelihood estimate* (MLE) is the value of θ that makes the data the most “likely”, that is, the value that maximizes $L(\theta)$
- To compute the maximum likelihood estimate:
 1. Write down the likelihood and take its log:

$$\log(L(\theta)) = \ell(\theta)$$

(log likelihood)

2. Take the derivative of $\ell(\theta)$ with respect to θ :

$$\ell'(\theta) = \frac{d\ell(\theta)}{d\theta}$$

3. Solve for $\hat{\theta}$ such that $\ell'(\theta) = 0$

Example: Binomial

$$1) L(\theta) \propto \theta^y (1-\theta)^{n-y}$$

$$1) l(\theta) = y \log(\theta) + (n-y) \log(1-\theta)$$

$$2) \frac{dl}{d\theta} = \frac{y}{\theta} - \frac{(n-y)}{(1-\theta)} = 0$$

$$\boxed{\begin{aligned} \log(a^b) &= b \log(a) \\ \log(ab) &= \\ \log(a) + \log(b) & \end{aligned}} \Rightarrow \frac{y}{\theta} = \frac{n-y}{(1-\theta)} \Rightarrow (n-y)\theta = y(1-\theta)$$
$$\hat{\theta}_{MLE} = \frac{y}{n} = \frac{4}{5} = .8$$

Example: the likelihood for independent Bernoulli's

$$\begin{aligned} p(y_1, y_2, \dots, y_n | 1, \theta) &= p(y_1, y_2, \dots, y_n | \theta) \\ &= p(y_1 | \theta) p(y_2 | \theta) \dots p(y_n | \theta) \\ &= \prod_{i=1}^n p(y_i | \theta) \\ &= \prod_{i=1}^n \binom{1}{y_i} \theta^{y_i} (1 - \theta)^{(1-y_i)} \\ &= \left[\prod_{i=1}^n \binom{1}{y_i} \right] \theta^{\sum_{i=1}^n y_i} (1 - \theta)^{n - \sum_{i=1}^n y_i} \\ &= L(\theta) \end{aligned}$$

Assume I poll the class about preferred candidate in election,

y_1, \dots, y_n , $n = 25$ people sampled.

$$y_i \sim \text{Bern}(\theta) \equiv \text{Bin}(1, \theta)$$

$$\begin{aligned} L(\theta) &= P(y_1, \dots, y_{25} | \theta) && \text{MLE:} \\ &= \prod_{i=1}^{25} P(y_i | \theta) && \xrightarrow{\sum y_i} \frac{\sum y_i}{n} \\ &= \prod_{i=1}^{25} \theta^{y_i} (1-\theta)^{1-y_i} \\ &= \theta^{\sum y_i} (1-\theta)^{\sum (1-y_i)} \end{aligned}$$

$$\text{Equiv. } \sum y_i \sim \text{Bin}(25, \theta)$$

Sufficient Statistics

- Let $L(\theta) = p(y_1, \dots, y_n | \theta)$ be the likelihood and $s(y_1, \dots, y_n)$ be a statistic *of data*,
- $s(y)$ is a sufficient statistic if we can write:

$$L(\theta) = h(y_1, \dots, y_n)g(s(y), \theta)$$

- g is only a function of $s(y)$ and θ only
 - h is *not* a function of θ
- This is known as the *factorization theorem*
- $L(\theta) \propto g(s(y), \theta)$
-
- $$(5 \choose y) \theta^y (1 - \theta)^{5-y}$$

z_1, z_2, \dots, z_s make/miss
 $\begin{cases} z=1 & \text{is make} \\ z=0 & \text{is miss} \end{cases}$

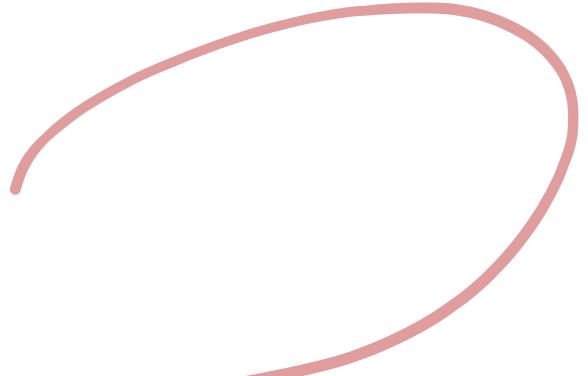
$$\begin{aligned}
 L(\theta) &= \prod_{i=1}^s \binom{1}{z_i} \theta^{z_i} (1-\theta)^{1-z_i} \\
 &= \underbrace{\left(\prod_{i=1}^s \binom{1}{z_i} \right)}_{h(z)} \theta^{\sum z_i} (1-\theta)^{\sum (1-z_i)}
 \end{aligned}$$

Sufficient Statistics

- Intuition: a sufficient statistic contains all of the information about θ
 - Many possible sufficient statistics *(not unique)*
 - Often seek a statistic of the lowest possible dimension (minimal sufficient statistic)
 - What are some sufficient statistics in the previous binomial example?

$\sum z_i$, \bar{z} are both sufficient,
 (z_1, \dots, z_5) is also sufficient.

(z_1, z_2) is not sufficient. Losing info about Q ,



Estimators and Estimates

- In classical (frequentist) statistics, θ is an unknown constant
- An **estimator** of a parameter θ is a function of the random variables, Y
 - E.g. for Binomial(1, θ): $\hat{\theta}(Y) = \frac{\sum_i Y_i}{n}$
 - An estimator is a random variable
 - Interested in properties of estimators (e.g. mean and variance)

*Capital
Y!*

$$\hat{\theta}(Y) = \frac{\sum_i Y_i}{n}$$

Estimators and Estimates

- $\hat{\theta}(y)$ as a function of realized data is called an estimate
 - Plug in observed data $y = (y_1, \dots, y_n)$ to estimate θ
 - An estimate is a non-random constant (it has 0 variability)
 - E.g. in the binomial($1, \theta$), $\hat{\theta} = \bar{y} = \frac{\sum_i y_i}{n}$ is the maximum likelihood estimate for the binomial proportion.

What are desirable properties
of estimators?

Unbiased: on average $\hat{\theta}(\gamma)$ is θ .

Small Variance

Want Both!

consistent: get θ if you have infinite data.

Simple / computational Tractable.

Low Bias & Low Variance

→ Low Error /

High Accuracy -



Bias and Variance

- Estimators are random variables. What are some r.v. properties that are desirable?

- Bias: $E[\hat{\theta}] - \theta$

- $E[\hat{\theta}] - \theta = 0$ means the estimator is unbiased

- E.g. expectation of the binomial MLE: $E[\hat{\theta}] = E\left[\frac{\sum Y_i}{n}\right] = \theta$

- $\text{Var}(\hat{\theta}) = E[(\hat{\theta} - \underline{E[\hat{\theta}]})^2]$

E.g. variance of the binomial MLE is

$$\text{Var}[\hat{\theta}] = \text{Var}\left(\frac{\sum Y_i}{n}\right) = \frac{\theta(1-\theta)}{n}$$

$$\text{Var}(a + bY) = b^2 \text{Var}(Y), \quad \text{Var}(Y_1 + Y_2) = \text{Var}(Y_1) + \text{Var}(Y_2)$$

if indep.

$$\text{Var}\left(\frac{\sum g_i}{n}\right) = \frac{1}{n^2} \sum \text{Var}(Y_i) = \frac{n\sigma^2(1-\sigma)}{n^2}$$

Bias and Variance

- Want estimators that have low bias and variance because this implies low overall error

- Mean squared error equals bias² + variance

$$MSE = B^2 + V$$

$$MSE : E[(\hat{\theta} - \theta)^2]$$

$$E[(\hat{\theta} - \theta)^2] = E[((\hat{\theta} - E[\hat{\theta}]) + (E[\hat{\theta}] - \theta))^2]$$

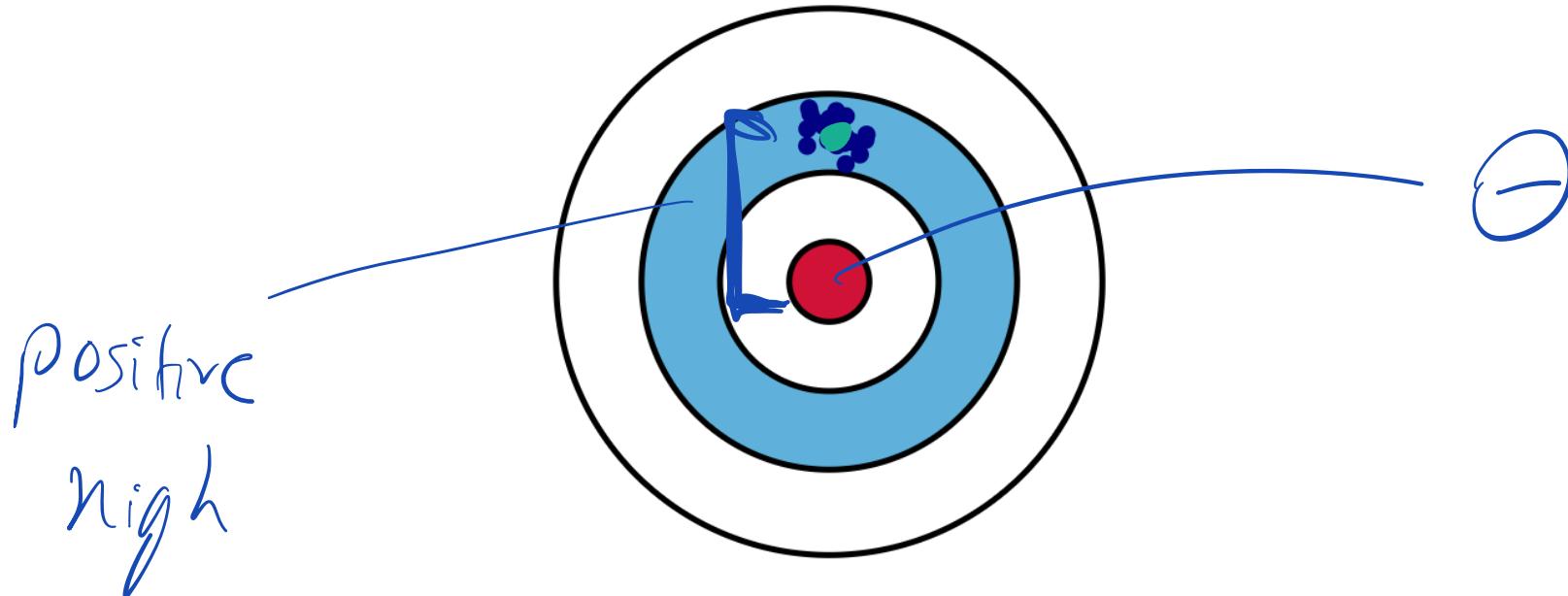
$$= E[(\hat{\theta} - E[\hat{\theta}])^2] + E[(E[\hat{\theta}] - \theta)^2] + E[2(\hat{\theta} - E[\hat{\theta}])(E[\hat{\theta}] - \theta)]$$

Varian (L) + Bias² + const

Bias / Variance Tradeoff

Bias

The average difference between the prediction and the response

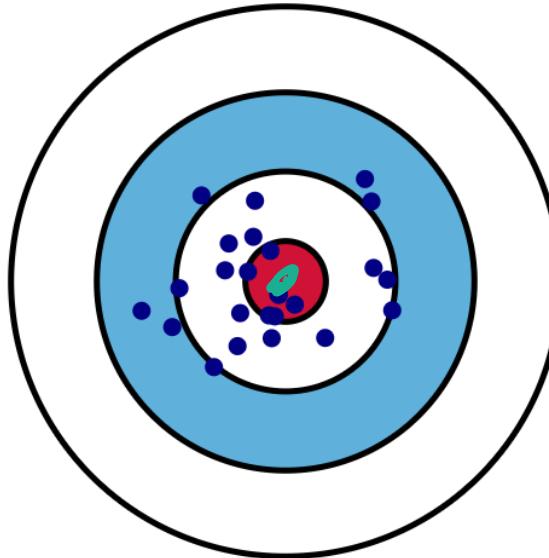


Statistical definition of bias:

$$E[\hat{\theta} - \theta]$$

Variance

How variable is the prediction about its mean?

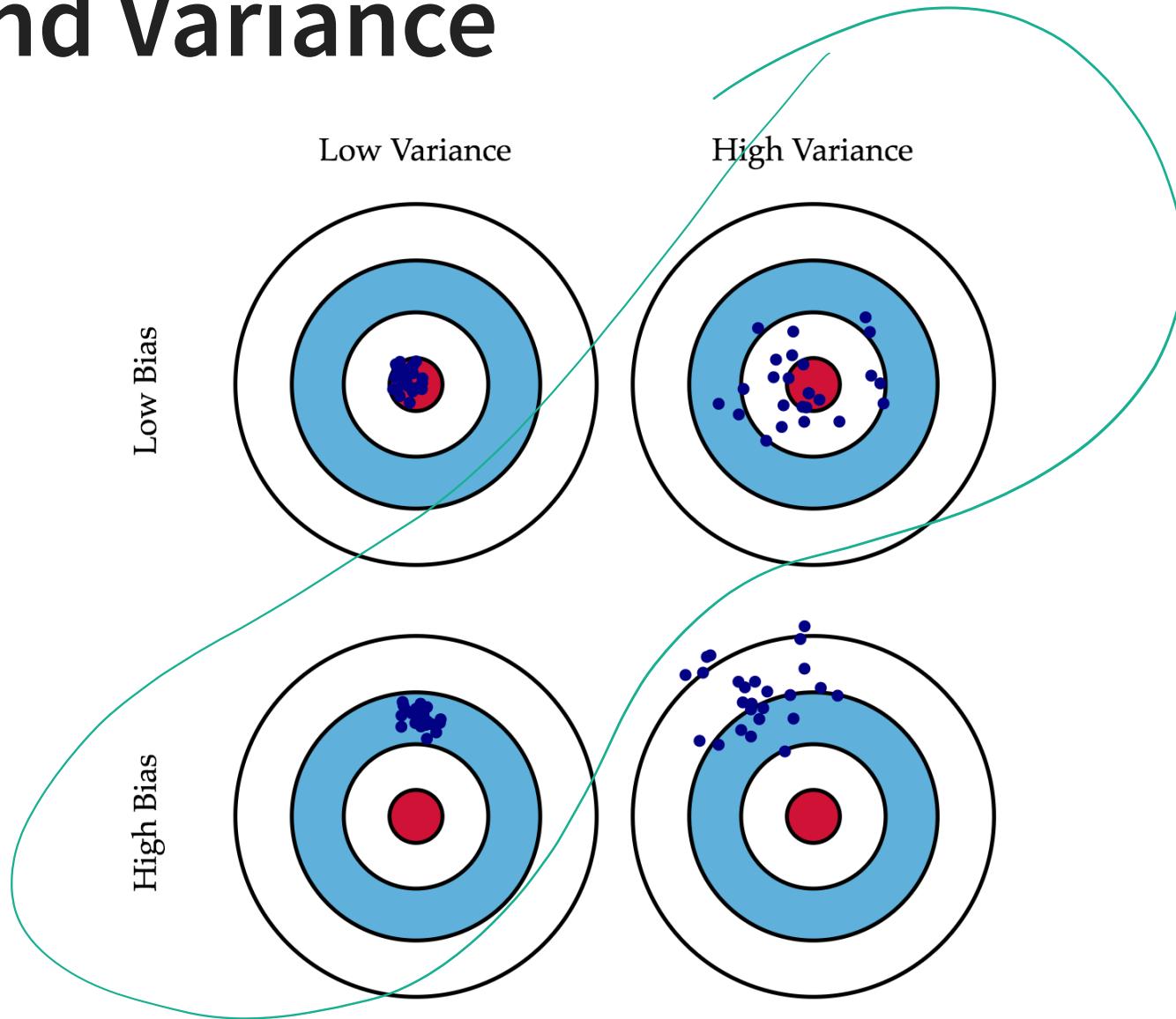


Low Bias /
High Variance.

Statistical definition of variance:

$$E[\hat{\theta} - E[\hat{\theta}]]^2$$

Bias and Variance



Maximum Likelihood Estimators

Under relatively weak conditions:

- The MLE is *consistent*. It converges to the true value as the sample size goes to infinity.
 - Need bias and variance to go to 0 as sample size increases
- The MLE is *asymptotically optimal*. For “large” sample sizes it has the lowest variance.
- *Equivariance*: if $\hat{\theta}$ is the MLE for θ then $g(\hat{\theta})$ is the MLE for $g(\theta)$

Confidence Interval Simulations

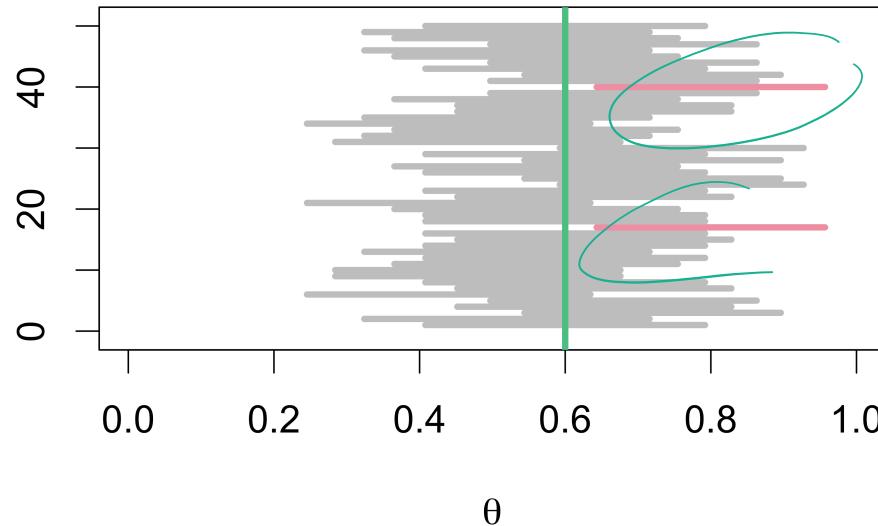
Let's do 50 hypothetical replications to illustrate confidence intervals

```
1 for i in 1 to 50:  
2   - Draw Y_i from Bin(25, 0.6)  
3   - Compute and plot the 95% confidence interval
```

- Will have 50 confidence intervals based on 50 simulated datasets.
- A 95% interval means that on average 95% of these 50 intervals will cover the true value

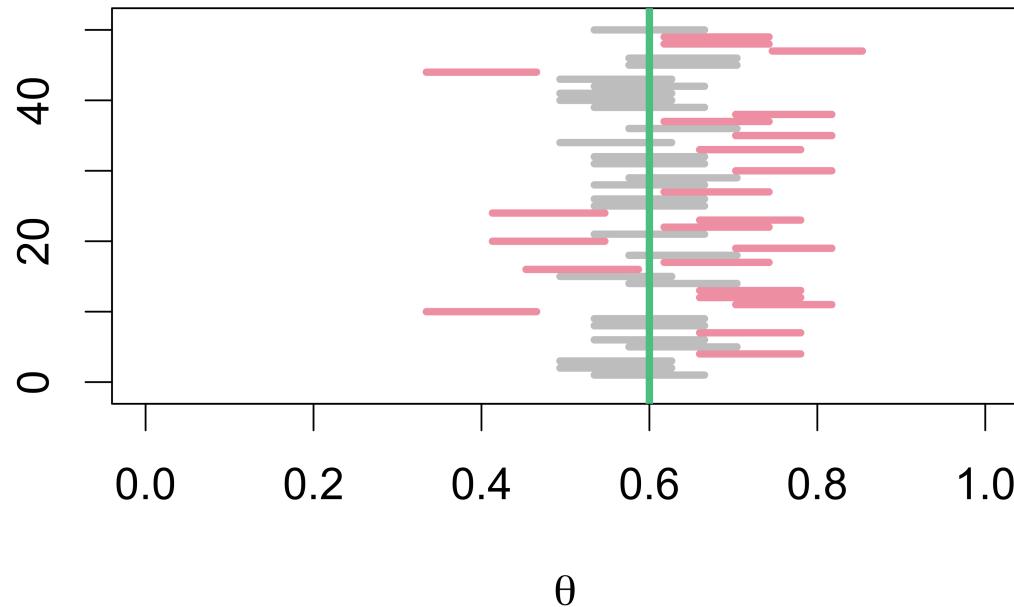
95% Confidence Intervals

In truth, 60% of the population will vote for “candidate 1”



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

50% Confidence Intervals



We expect $0.50 \times 50 = 25$ of the intervals to *not* cover the true parameter, 0.6

$$\Pr(\text{Low}(Y_1, \dots, Y_n) \leq \theta \leq \text{Up}(Y_1, \dots, Y_n)) = 0.95$$

(lower endpt)

Random



(upper endpt)

Random

constant.

$$\Pr(L(y_1, \dots, y_n) \leq \theta \leq U(y_1, \dots, y_n)) = 0.95$$

1

.5

3

Data Generating Process (DGP)

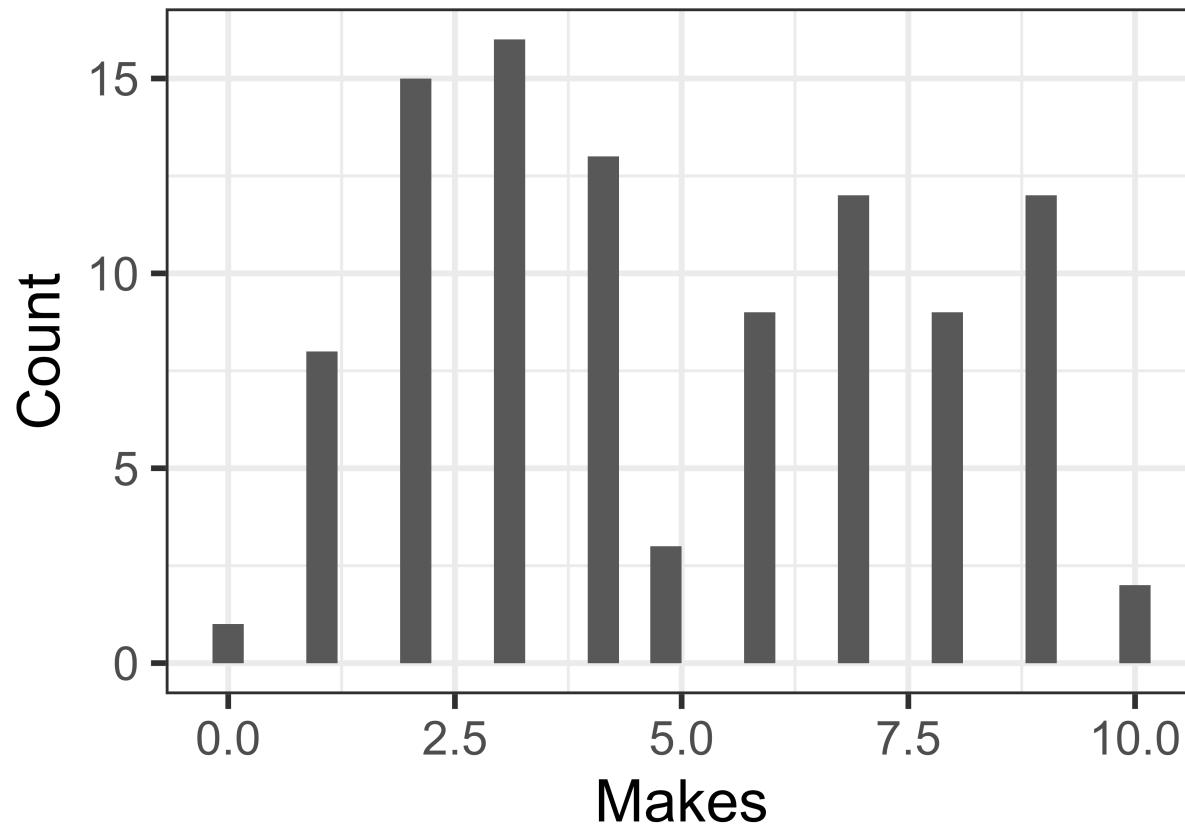
- DGP: a statistical model for how the observed data might have been generated
- Often write the DGP using pseudo-code:

```
1 for (i in 1:N)
2   - Generate y_i from a Normal(0, 1)
3 return y = (y_1, ... y_N)
```

- The DGP should tell a story about how the data came to be
- Can translate the DGP into a statistical model

Data Generating Process (DGP)

Assume everybody in this class goes to a basketball court and takes 10 free throw shots:





Data Generating Process (DGP)

Tell a plausible story: some students play basketball and some don't. Before you take your shots we record whether or not you have played before.

```
1 assume theta_1 > theta_0
2 for (i in 1:100)
3   - Generate z_i from Bin(1, phi)
4   - p_i = theta_0 if z_i=0
5   - p_i = theta_1 if z_i=1
6   - Generate y_i from a Binom(10, p_i)
7 return y = (y_1, ... y_100) and z = (z_1, ..., z_100)
```

ϕ is fraction of people who played.

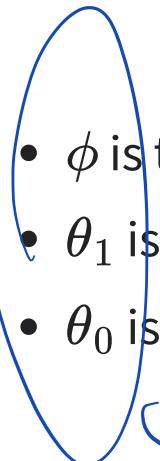
θ_0 = Prob(make) / no experience
 θ_1 = Prob(make) / experience,
Is this a reasonable model?

Mixture Models

$$\underline{Z_i} = \begin{cases} 0 & \text{if the } i^{\text{th}} \text{ if student doesn't play basketball} \\ 1 & \text{if the } i^{\text{th}} \text{ if student does play basketball} \end{cases}$$

R.V.s:

$$Z_1, \dots, Z_{100}$$



$$Z_i \sim \text{Bin}(1, \phi)$$

$$Y_i \sim \begin{cases} \text{Bin}(10, \theta_0) & \text{if } Z_i = 0 \\ \text{Bin}(10, \theta_1) & \text{if } Z_i = 1 \end{cases}$$

$$\theta_0 < \theta_1$$

- ϕ is the fraction of students that have experience playing basketball
- θ_1 is the probability of making a shot for an experienced player
- θ_0 is the probability of making a shot for an inexperienced player

estimands

Table of relevant quantities

- Can be a fixed constant (no variability) or a random variable (has variability)
- Can be observed (known) or unobserved (unknown)
- Helpful for to keep track of all of the relevant quantities

	Obs	Unknown
$\text{Var} > 0$	y_1, \dots, y_{100} z_1, \dots, z_{100}	
$\text{Var} = 0$ (constants)	$n,$ y_1, \dots, y_{100} z_1, \dots, z_{100}	$\phi, \theta_0,$ θ_1

Bayes

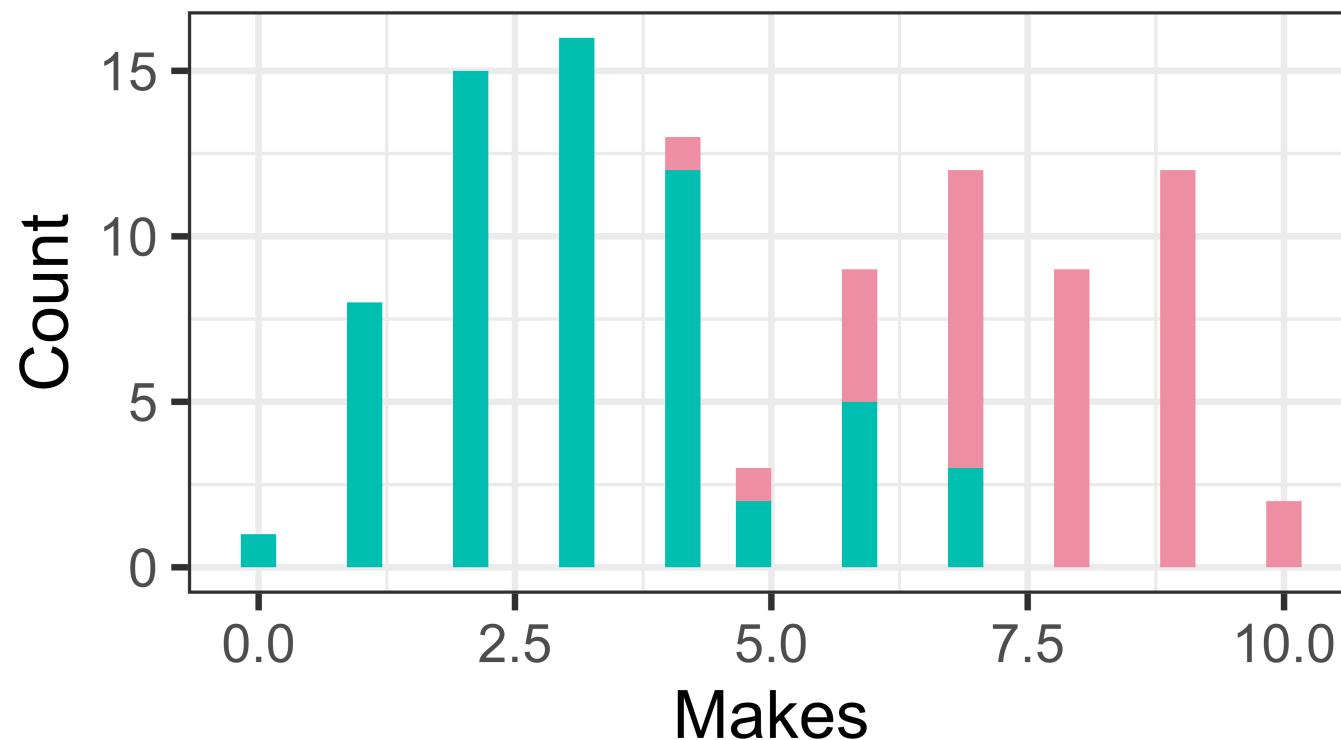
Mixture models

- A mixture model is a probabilistic model for representing the presence of subpopulations
- The subpopoluation to which each individual belongs is not necessarily known
 - e.g. do we ask: “have you played basketball before?”
- When z_i is not observed, we sometimes refer to it as a clustering model
 - *unsupervised learning*

A Mixture Model

(Supervised Learning)

Group ■ Experience ■ No Experience



Note: z is observed

Mixture Model Likelihood

$$\theta^a \theta^b = \theta^{a+b}$$

Z is observed

$$L(\phi, \theta_0, \theta_1) \propto P(y_1, \dots, y_{100}, z_1, \dots, z_{100} | \phi, \theta_0, \theta_1)$$

$$= \prod_{i=1}^{100} P(y_i = y_i | z_i, \theta_0, \theta_1) P(z_i = z_i | \phi)$$

$$= \prod_{i=1}^{100} \left[\frac{\phi^{z_i}}{y_i!} \theta_1^{y_i} (1-\theta_1)^{10-y_i} \right]^{z_i} \left[\frac{(1-\phi)^{1-z_i}}{y_i!} \theta_0^{y_i} (1-\theta_0)^{10-y_i} \right]^{1-z_i} \times$$

$$\cancel{\propto \theta_1^{\sum z_i y_i} (1-\theta_1)^{\sum z_i (10-y_i)} \theta_0^{\sum (1-z_i) y_i} (1-\theta_0)^{\sum (1-z_i)(10-y_i)} \phi^{\sum z_i} (1-\phi)^{\sum (1-z_i)}}$$

Sufficient statistics When Z_i is observed

Together, the following quantities are sufficient for $(\theta_0, \theta_1, \phi)$

- $\sum y_i z_i$ (total number of shots made by experienced players)
- $\sum y_i(1 - z_i)$ (total number of shots made by inexperienced players)
- $\sum z_i$ (total number experienced players)

$$\hat{\theta}_{1, MLE} = \frac{\sum y_i z_i}{10 \sum z_i}, \quad \hat{\theta}_{0, MLE} = \frac{\sum y_i(1 - z_i)}{10 \sum (1 - z_i)}$$

$$\phi_{MLE} = \frac{\sum z_i}{100}$$

Data Generating Process (DGP)

```
1 for (i in 1:100)
2   - Generate z_i from Bin(1, phi)
3   - p_i = theta_1 if z_i=1
4   - p_i = theta_0 if z_i=0
5   - Generate y_i from a Binom(10, p_i)
6 return y = (y_1, ... y_100)
```

This time we don't record who has experience with basketball.

A Mixture Model

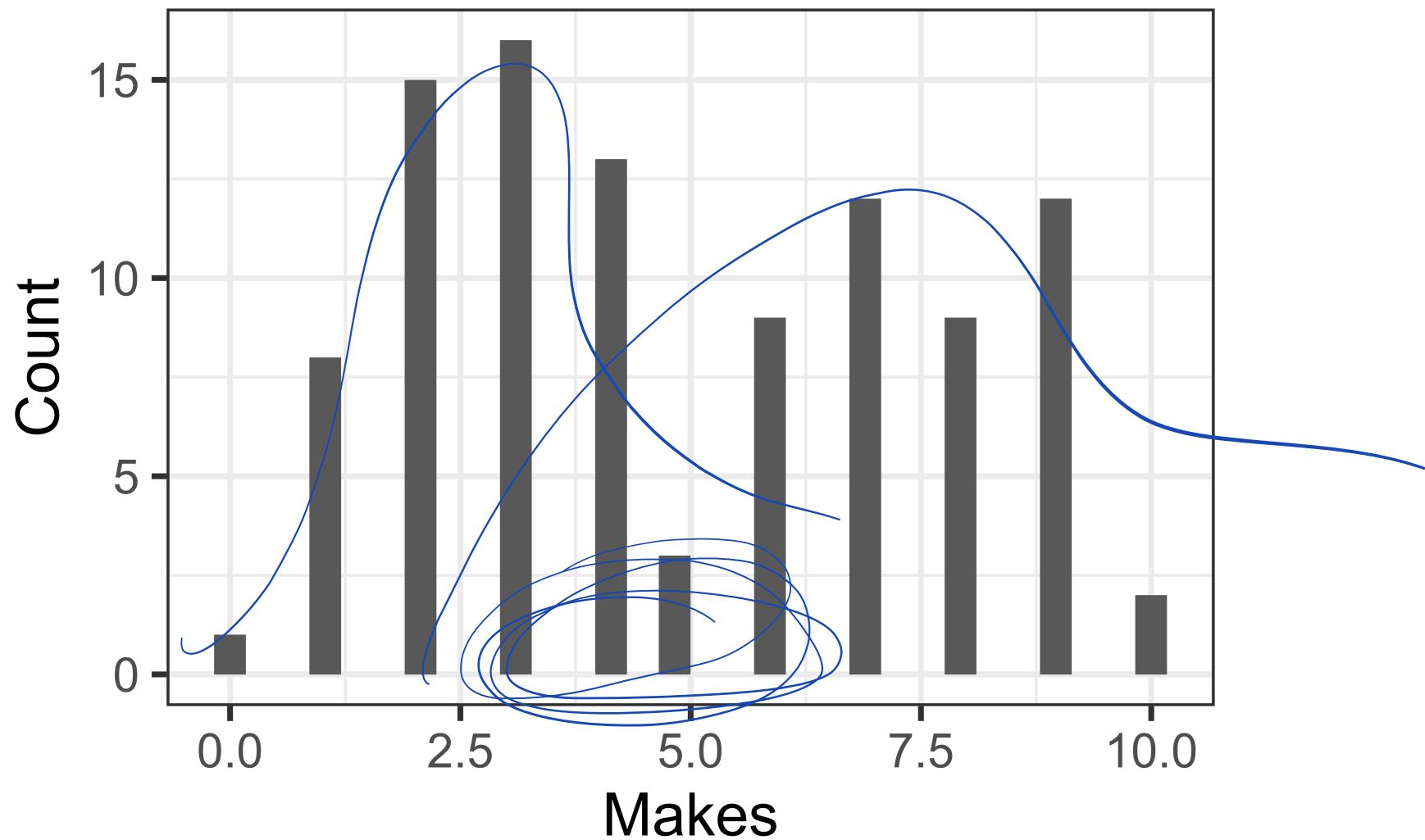


Table of Relevant Quantities

	Obs	Unknown
$\text{Var} > 0$	y_1, \dots, y_{100}	z_1, \dots, z_{100}
$\text{Var} = 0$ (constants)	$n,$ y_1, \dots, y_{100}	$\phi_1, \theta_0,$ θ_1 z_1, \dots, z_{100}

A finite mixture model

- Even if we don't observe Z , it's often useful to introduce it as a *latent* variable
- Write the observed data likelihood by integrating out the latent variables from the complete data likelihood

$$\begin{aligned} p(Y | \theta) &= \sum_z p(Y, Z = z | \theta) \\ &= \sum_z p(Y | Z = z, \theta) p(Z = z | \theta) \end{aligned}$$

In general we can write a K component mixture model as:

$$p(Y) = \sum_k^K \pi_k p_k(Y)$$

weight

$$\text{with } \sum \pi_k = 1$$

Mixture Model Likelihood

$$L(\Theta_0, \Theta_1, \phi) = \prod_{i=1}^{100} P(y_i | \Theta_0, \Theta_1, \phi)$$

z unobserved

$$= \prod_{i=1}^{100} \left[\sum_{z_i \geq 0} P(y_i, z_i | \Theta_0, \Theta_1, \phi) \right]$$

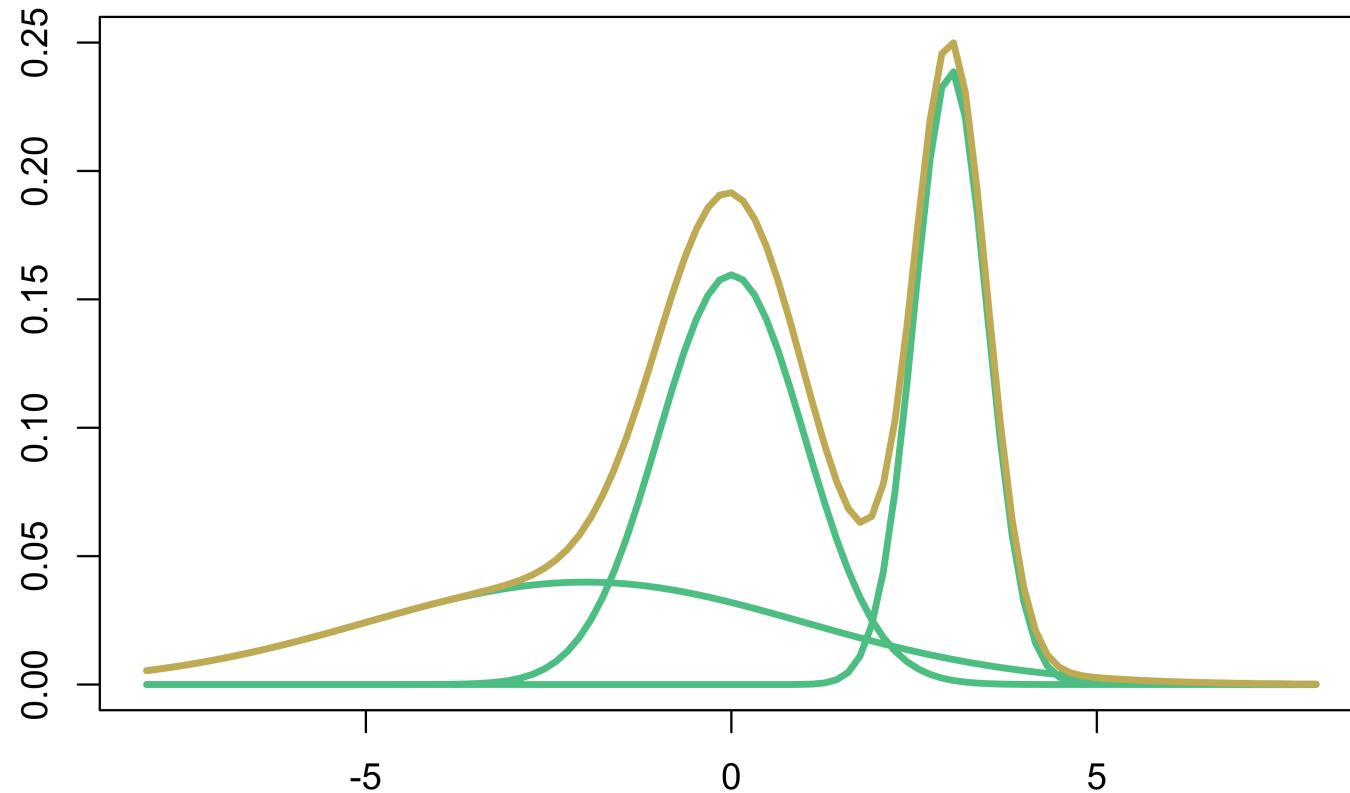
$$= \prod_{i=0}^{100} \left[(1-\phi) \binom{10}{y_i} \Theta_0^{y_i} (1-\Theta_0)^{10-y_i} + \phi \binom{10}{y_i} \Theta_1^{y_i} (1-\Theta_1)^{10-y_i} \right]$$

Observed Data Likelihood

Complete Data Likelihood

Sufficient stats are
full data (no compression
possible!)

Finite Mixture models



Infinite Mixture Models

- In the previous example the latent variable had finitely many outcomes
- Latent variables can have infinitely many outcomes in which case we have any infinite mixture
- Example:

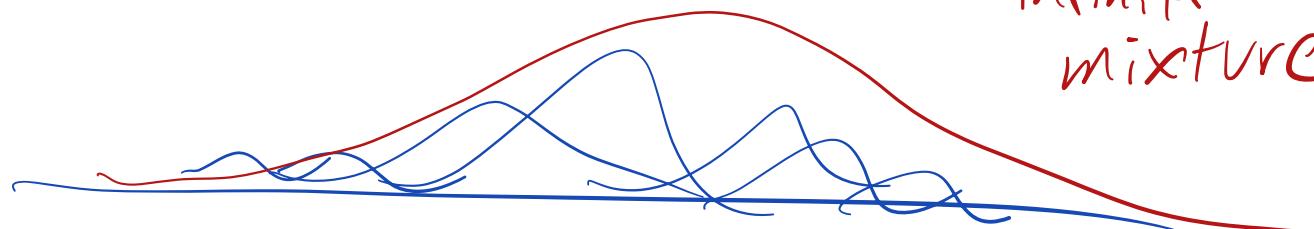
$$\begin{aligned} \mu &\sim N(0, \tau^2) \\ Y &\sim N(\mu, \sigma^2) \end{aligned}$$

$\epsilon_x \sim N(0, \sigma^2)$
 $\epsilon_\mu \sim N(0, \tau^2)$

$$Y = \mu + \epsilon_x$$
$$\mu = \epsilon_\mu$$
$$Y = \epsilon_\mu + \epsilon_x$$
$$\underline{p(Y | \sigma^2, \tau^2)} = \int p(Y, \mu | \sigma^2, \tau^2) d\mu$$

What is the *marginal* distribution of Y ?

infinite
mixture.



Bayesian Inference

- In frequentist inference, θ is treated as a fixed unknown constant
- In Bayesian inference, θ is treated as a random variable
- Need to specify a model for the joint distribution

$$p(y, \theta) = p(y | \theta)p(\theta)$$

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!

Bayes' Rule

$$P(A | B) = \frac{P(B | A)PAB}{P(B)}$$

- $P(A | B)$ is the conditional probability of A given B
- $P(B | A)$ 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

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

- $P(\theta \mid y)$ 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

$$\begin{aligned} P(\theta \mid y) &= \frac{P(y \mid \theta)P(\theta)}{P(y)} \\ &\propto P(y \mid \theta)P(\theta) \end{aligned}$$

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

The posterior is proportional to the likelihood times the prior!

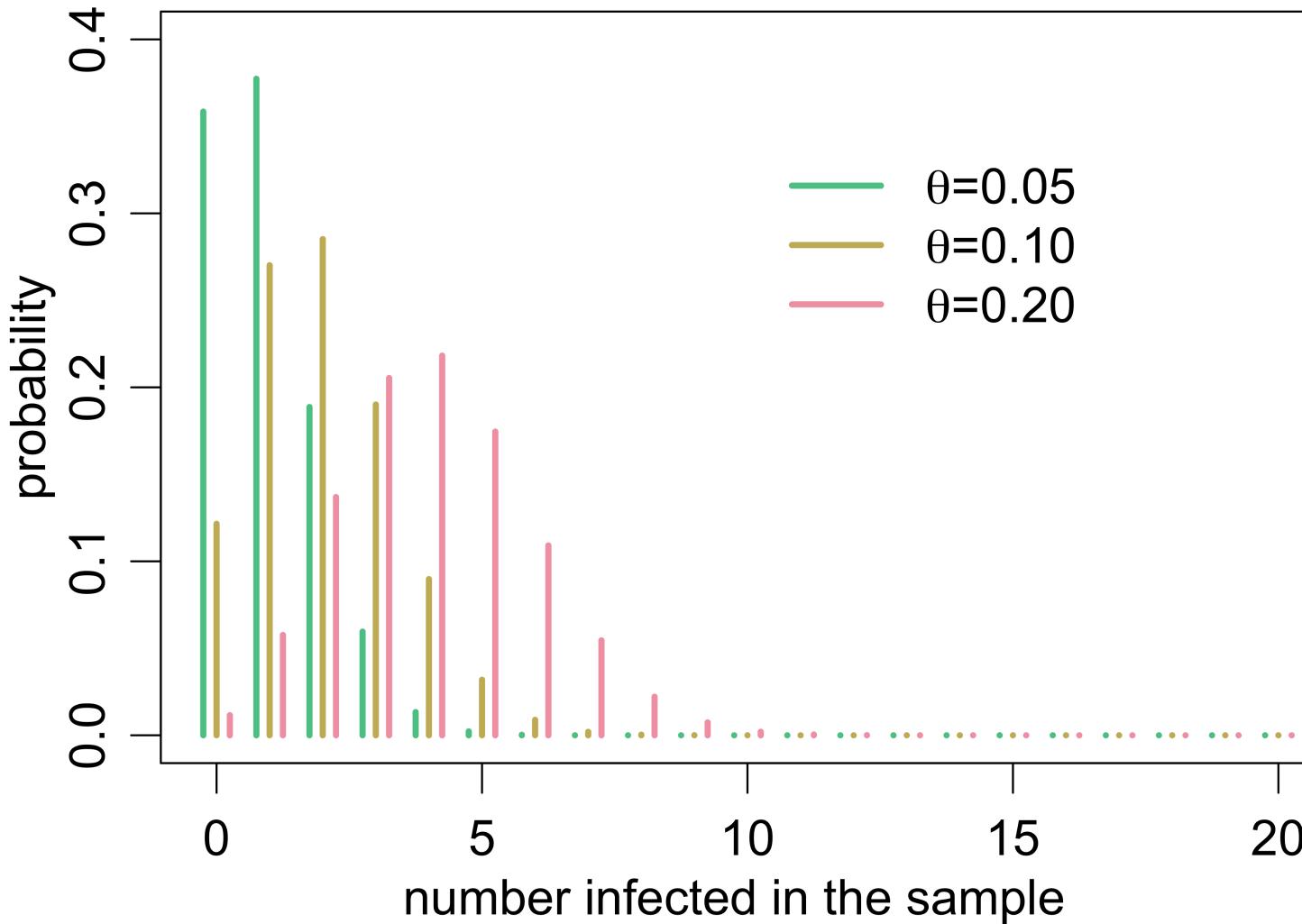
Example: Estimating COVID Infection Rates

- 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: $Y \sim \text{Binom}(20, \theta)$

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

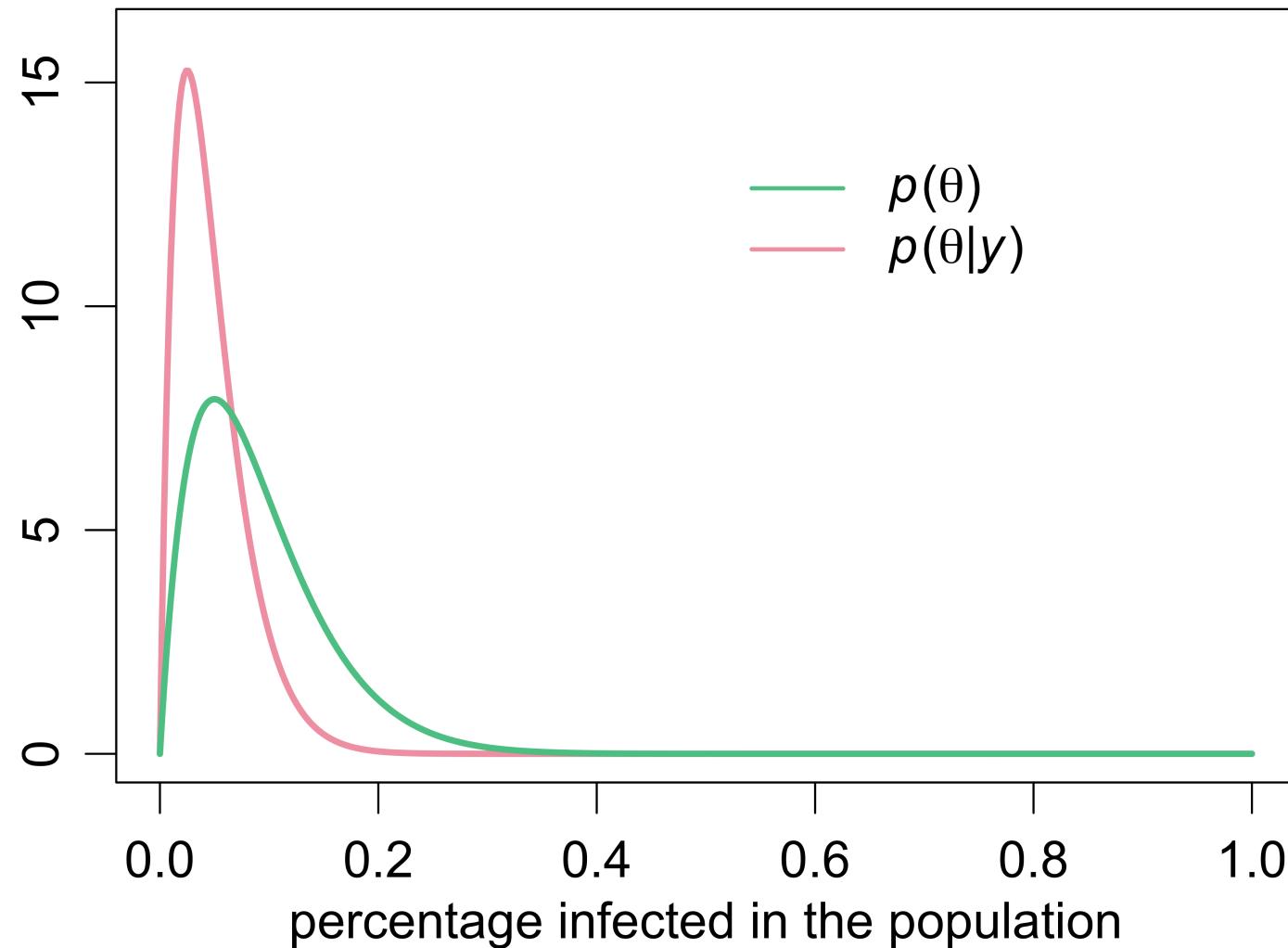


Table of Relevant Quantities

Summary

- Likelihood, log likelihood in MLE
- Confidence intervals (how they are defined in frequentist inference)
- Sufficient statistics
- Mixture models

Summary

- In frequentist inference, unknown parameters treated as constants
 - Estimators are random (due to sampling variability)
 - Asks: “how would my results change if I repeated the experiment?”

Summary

- In Bayesian inference, unknown parameters are random variables.
 - Need to specify a prior distribution for θ (not easy)
 - Asks: “what do I *believe* are plausible values for the unknown parameters?”
 - Who cares what might have happened, focus on what *did* happen!

Assignments

- Read chapters 1 and 2 of BR
- Homework 1 due 10/12