

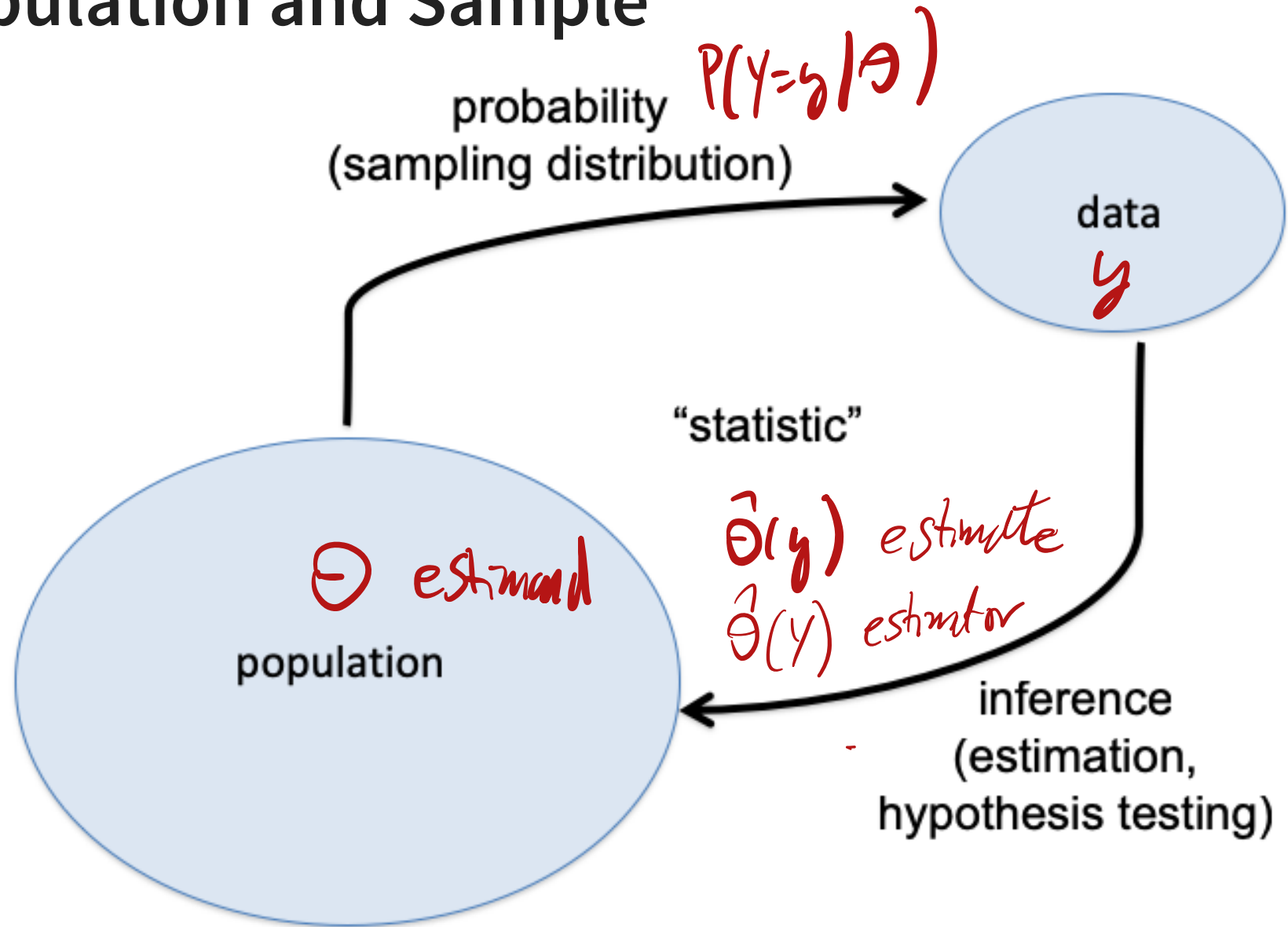
# Lecture 2: Likelihood Review

Professor Alexander Franks

# Logistics

- Read: BDA Chapters 1-2
- Sync content using link on course website:  
<https://bit.ly/3NXxk9H>
- Annotated lecture slides appear after class
- *No Midterm.*

# Population and Sample



# Independent Random Variables

- $Y_1, \dots, Y_n$  are random variables
- We say that  $Y_1, \dots, Y_n$  are *conditionally* independent given  $\theta$  if  $P(y_1, \dots, y_n \mid \theta) = \prod_i P(y_i \mid \theta)$
- Conditional independence means that  $Y_i$  gives no additional information about  $Y_j$  beyond that in knowing  $\theta$

Exchangeable :

$$P(y_1, \dots, y_n) = P(y_{\pi(1)}, \dots, y_{\pi(n)})$$

# The Likelihood Function

- The likelihood function is the probability density function of the observed data expressed as a function of the unknown parameter (conditional on observed data):
- A function of the unknown constant  $\theta$ .
- Depends on the observed data  $y = (y_1, y_2, \dots, y_n)$
- Two likelihood functions are equivalent if one is a scalar multiple of the other

$$L(\theta) \propto P(Y_1 = y_1, \dots, Y_n = y_n | \theta)$$

$$\text{If iid: } L(\theta) \propto \prod_{i=1}^n P(Y_i = y_i | \theta)$$

# Sufficient Statistics

A statistic  $s(\vec{Y})$  is sufficient for underlying parameter  $\theta$  if the conditional probability distribution of the  $Y$ , given the statistic  $s(Y)$ , does not depend on  $\theta$ .

$$P(Y_1, \dots, Y_n | S(Y_1, \dots, Y_n), \theta) = P(Y_1, \dots, Y_n | S(Y_1, \dots, Y_n))$$

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} N(\theta, 1)$$

$$\rightarrow Y_i | \bar{Y} \sim N(0, 1 - 1/n) \quad (\text{exercise})$$

Bayesian:  $P(\theta | y_1, \dots, y_n) = P(\theta | s(y_1, \dots, y_n))$

# Sufficient Statistics

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

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

*no  $\theta$*

*only  $s(y)$  &  $\theta$*

- $g$  is only a function of  $s(y)$  and  $\theta$  only
- $h$  is *not* a function of  $\theta$
- $L(\theta) \propto g(s(y), \theta)$

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

# The Likelihood Principle

- **The likelihood principle:** All information from the data that is relevant to inferences about the value of the model parameters is in the equivalence class to which the likelihood function belongs
- Two likelihood functions are equivalent if one is a scalar multiple of the other
- Frequentist testing and some design based estimators violate the likelihood principle



# Binomial vs Negative Binomial

$$Y \sim \text{Bin}(12, \theta) \quad \text{obs: } y=3$$

$$L(\theta; y=3) \propto \cancel{\binom{12}{3}} \theta^3 (1-\theta)^9$$

$$X \sim \text{NB}(3, \theta) \quad (\text{keep flipping until 3 heads})$$

$$\text{obs } x=9$$

$$L(\theta; x=9) = \cancel{\binom{11}{2}} \theta^3 (1-\theta)^9$$

$$H_0: \theta = .5$$

$$H_a: \theta < .5$$

p-values

$$\text{Bin: } \text{pbinom}(3, 12, \theta = 1/2) \\ \hookrightarrow 0.073$$

$$\text{NB: } 1 - \text{pnbinom}(8, 3, 1/2) \\ \hookrightarrow 0.033$$

# Score and Fisher Information

- The score function:  $\frac{d\ell(\theta; y)}{d\theta}$ 
  - $E\left[\frac{d\ell(\theta; Y)}{d\theta} \mid \theta\right] = 0$  (under certain regularity conditions)
- Fisher information is a measure of the amount of information a random variable carries about the parameter
  - $I(\theta) = E\left[\left(\frac{d\ell(\theta; Y)}{d\theta}\right)^2 \mid \theta\right]$  (variance of the score)
  - Equivalently:  $I(\theta) = -E_{y|\theta}\left[\frac{d^2\ell(\theta; Y)}{d^2\theta}\right]$

Observed Info:  $\frac{-d^2\ell(\theta; y)}{d^2\theta}$  <sup>obs data</sup>

# Fisher Information

$$Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$$

known

$$L(\mu) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}}$$
$$\propto e^{-\frac{(\bar{y} - \mu)^2}{2\sigma^2/n}} \quad \left. \vphantom{\prod_{i=1}^n} \right\} g(s(y), \theta)$$

$$l(\mu) = -\frac{(\bar{y} - \mu)^2}{2\sigma^2/n}$$

$$l'(\mu) = \frac{(\bar{y} - \mu)}{\sigma^2/n}, \quad l'' = \frac{-n}{\sigma^2} \Rightarrow I(\mu) = \frac{n}{\sigma^2}$$

# Data Generating Process

# Data Generating Process (DGP)

- I select 100 random students at UCSB to 10 free throw shots at the basketball court
- Assume there are two groups: experienced and inexperienced players
- Skill is identical conditional on experience level

# 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)
```

# Mixture models

*Experience*

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

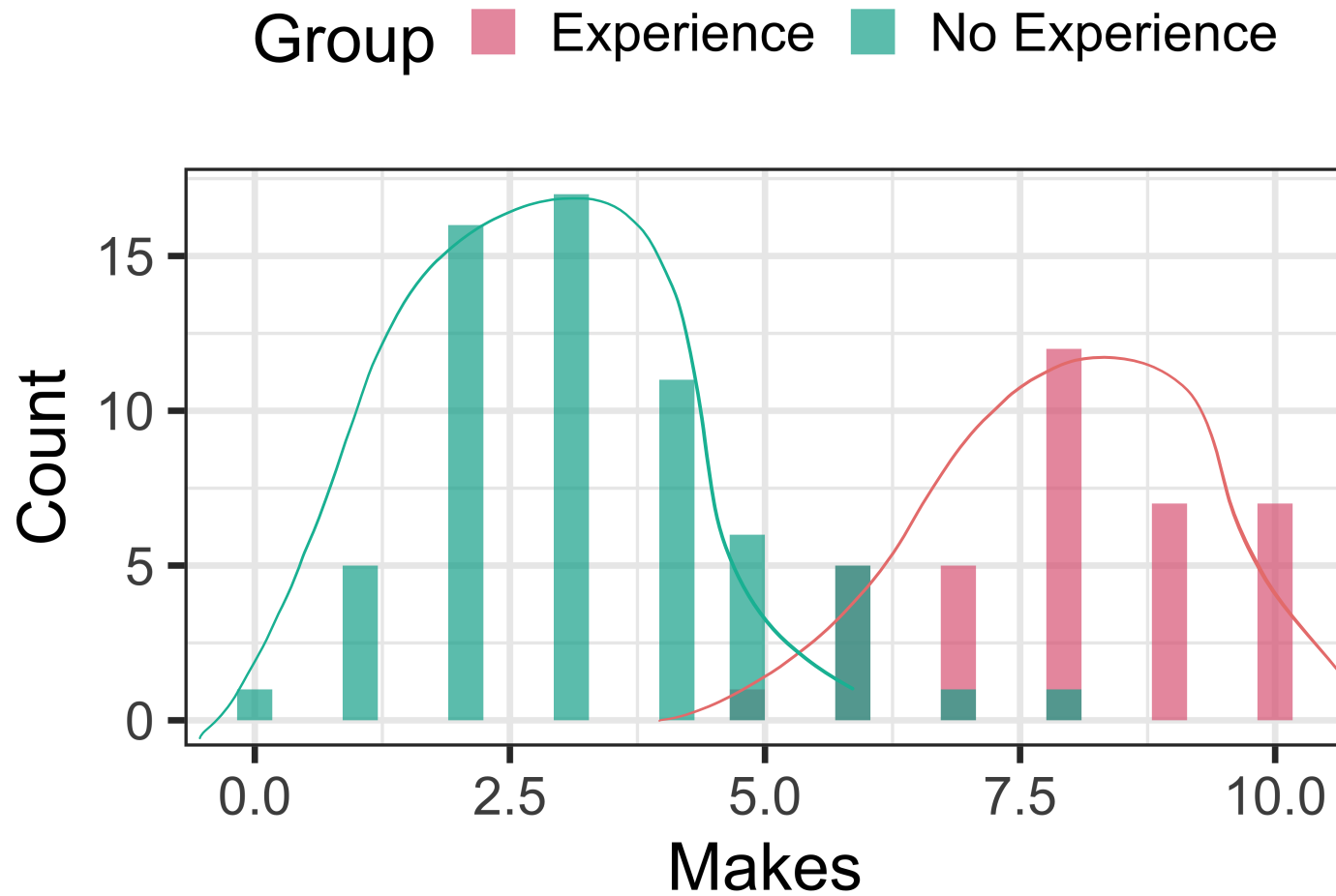
$$Z_i \sim \text{Bin}(1, \phi)$$

*Fraction Experienced*

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

*Make Probabilities*

# A Mixture Model



Note:  $z$  is observed



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

$$\propto \prod_{i=1}^n \left( \left[ \cancel{\phi} \left( \cancel{\frac{1}{y_i}} \right) \theta_1^{y_i} (1-\theta_1)^{1-y_i} \right]^{z_i} \times \right. \\ \left. \left[ (1-\phi) \left( \cancel{\frac{1}{y_i}} \right) \theta_0^{z_i} (1-\theta_0)^{1-z_i} \right]^{1-z_i} \right)$$

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

# 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)

# Mixture models

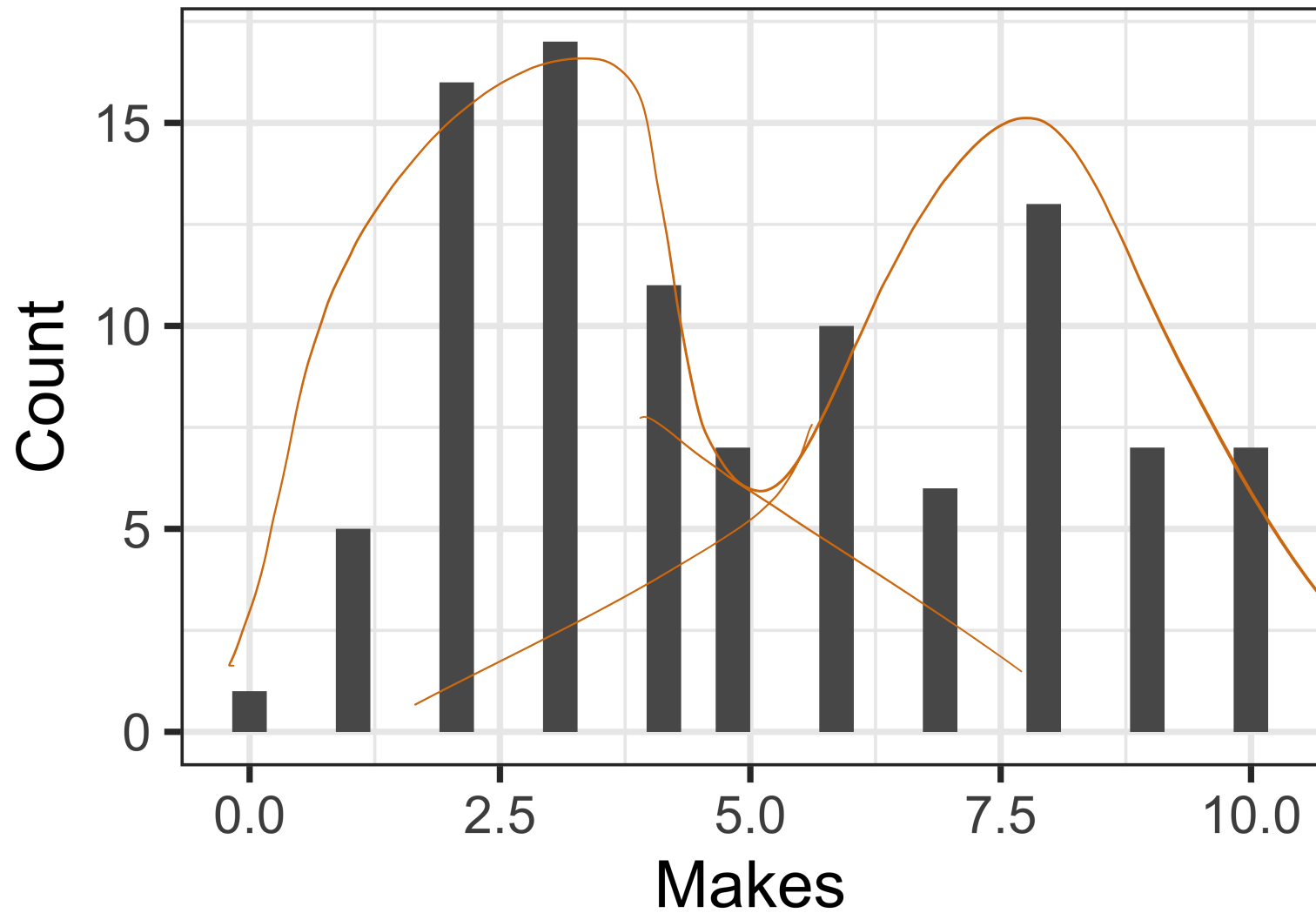
- A mixture model is a probabilistic model for representing the presence of subpopulations
- The subpopulation to which each individual belongs is not necessarily known
  - e.g. do we ask: “have you played basketball before?”
- When  $z_i$  is not observed, we sometimes refer to it as a clustering model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

# 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



$$L(\theta_1, \theta_0, \phi) \propto \prod_{i=1}^n P(y_i | \theta_0, \theta_1, \phi) \quad P(X, Y) = P(X|Y)P(Y)$$

$$\propto \prod_{i=1}^n \sum_{z=0}^1 P(y_i, z_i | \theta_1, \theta_0, \phi)$$

$$\propto \prod_{i=1}^n \left[ \sum_z P(y_i | z_i, \theta_1, \theta_0, \phi) P(z_i | \phi) \right]$$

$$\propto \prod_{i=1}^n \left[ \phi \binom{10}{y_i} \theta_1^{y_i} (1-\theta_1)^{10-y_i} + (1-\phi) \binom{10}{y_i} \theta_0^{y_i} (1-\theta_0)^{10-y_i} \right]$$

# A finite mixture model

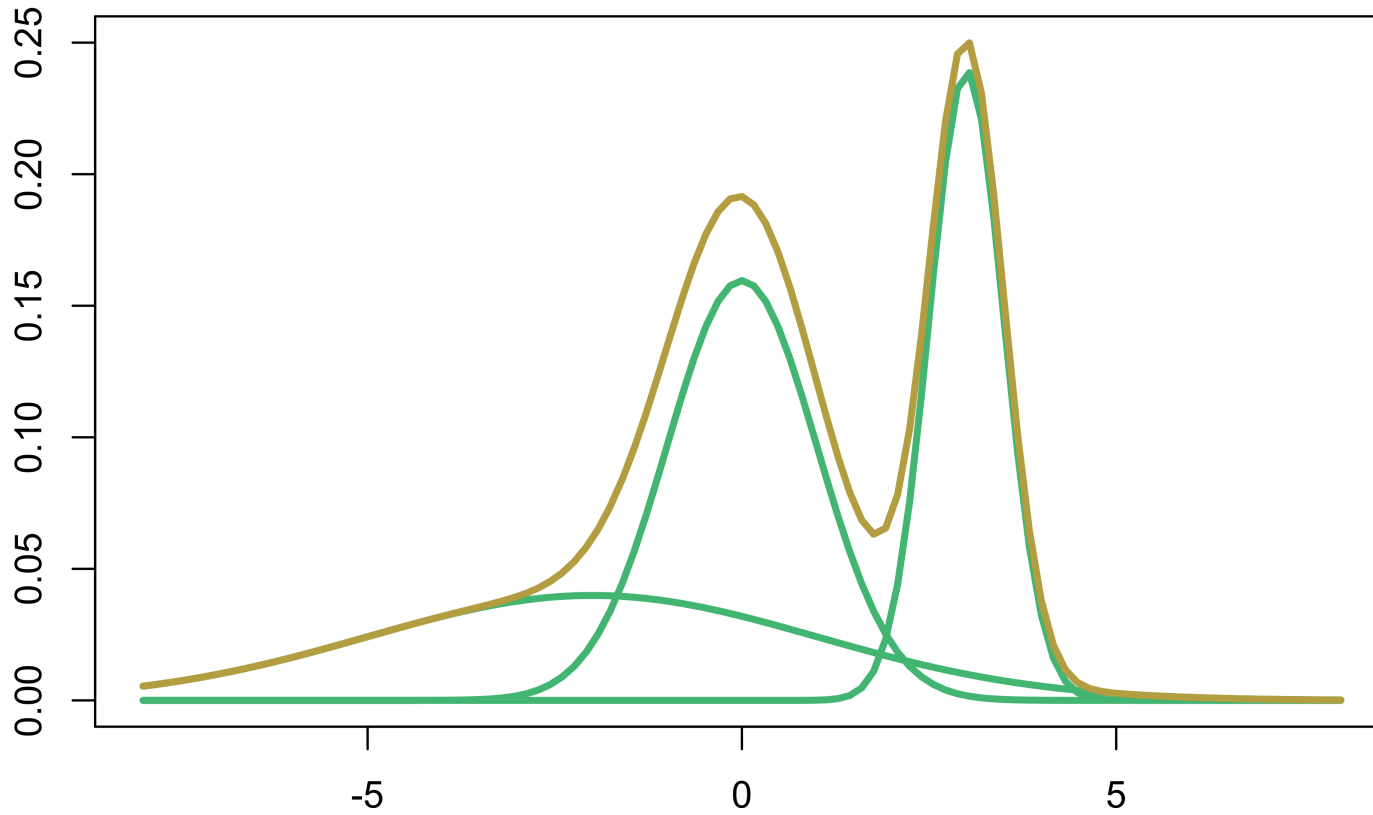
- Often crucial to understand the complete data generating process by introducing *latent* variables
- Write the *observed data likelihood* by integrating out the latent variables from the *complete data likelihood*

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

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

$$p(Y) = \sum_k^K \pi_k p_k(Y) \text{ with } \sum \pi_k = 1$$

# Finite mixture models





# Infinite Mixture Models

- Often helpful to think about infinite mixture models
- Example 1: normal observations with normally distributed mean

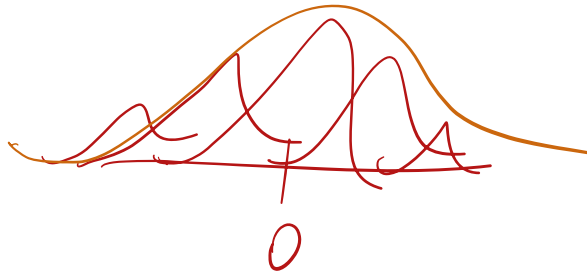
$$\mu_i \sim N(0, \tau^2)$$

$$Y_i \sim N(\mu_i, \sigma^2)$$

What is the distribution of  $Y_i$  given  $\tau^2$  and  $\sigma^2$  (integrating over  $\mu$ )?

$$P(Y_i | \tau^2, \sigma^2) = \int P(Y_i | \mu_i, \sigma^2) P(\mu_i | \tau^2) d\mu$$

- Calculus
- MGFs
- Representation



$$Y_i = \mu_i + \varepsilon \quad \begin{array}{l} \varepsilon \sim N(0, \sigma^2) \\ \mu_i \sim N(0, \tau^2) \end{array}$$

$$Y \sim N(0, \sigma^2 + \tau^2)$$

# Infinite Mixture Models

Example 2: Poisson observations with random rates

$$\lambda \sim \text{Gamma}(\alpha, \beta)$$

$$Y \sim \text{Pois}(\lambda)$$

$$\begin{aligned} P(Y|\alpha, \beta) &= \int_{\lambda} P(Y|\lambda) P(\lambda|\alpha, \beta) d\lambda \\ &= \int_{\lambda} \frac{\lambda^y e^{-\lambda}}{y!} \frac{\beta^{\alpha} \lambda^{\alpha-1} e^{-\beta\lambda}}{\Gamma(\alpha)} d\lambda \end{aligned}$$

$$= \frac{B^\alpha}{\Gamma(\alpha)\Gamma(y+1)} \int_0^\infty \lambda^{y+\alpha-1} e^{-(B+1)\lambda} d\lambda$$

kernel of Gamma  
 $\text{Gam}(\alpha+y, B+1)$

$$\frac{(B+1)^{\alpha+y}}{\Gamma(y+\alpha)}$$

$$= \frac{B^\alpha}{\Gamma(\alpha)\Gamma(y+1)} \frac{(B+1)^{\alpha+y}}{\Gamma(y+\alpha)}$$

Negative  
Binomial

# Infinite Mixture Models

- Example 3: normal observations with exponentially distributed scale

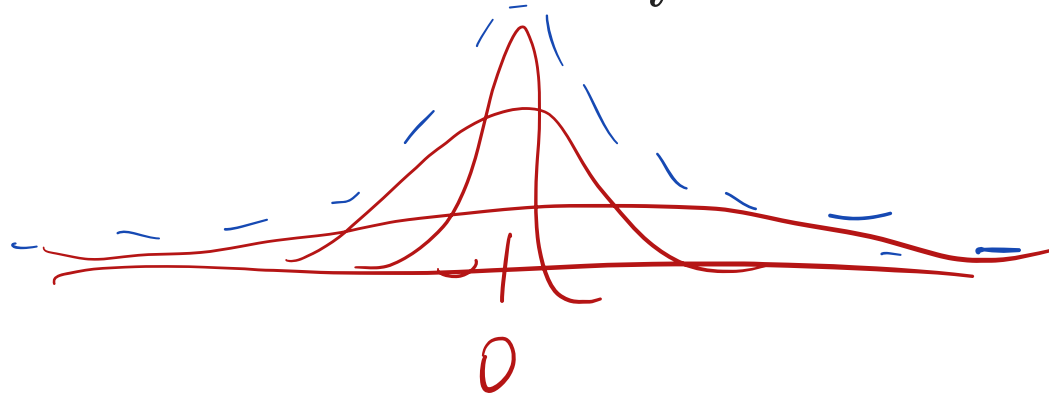
$$\sigma_i^2 \sim \text{Exponential}(1/2)$$

$$Y_i \sim N(0, \sigma_i^2)$$

*Inv-Gam( $\nu/2, \nu/2$ )*

Laplace

What is the distribution of  $Y_i$ ?



$$\left. \begin{aligned} y_i &\sim N(0, \sigma_i^2) \\ \sigma_i &= 1/|z_i| \\ z_i &\sim N(0, 1) \end{aligned} \right\}$$

$$x_i \sim N(0, 1)$$

$$z_i \sim N(0, 1)$$

$$y_i = \frac{x_i}{|z_i|}$$



Cauchy

# Summary

- Likelihood, log likelihood
- Sufficient statistics
- Fisher information
- Mixture models

# Assignments

- Read chapter 1-2 BDA3



