

# Lecture 1: Likelihood Review

Professor Alexander Franks

# Logistics

- Read: BDA Chapters 1-2

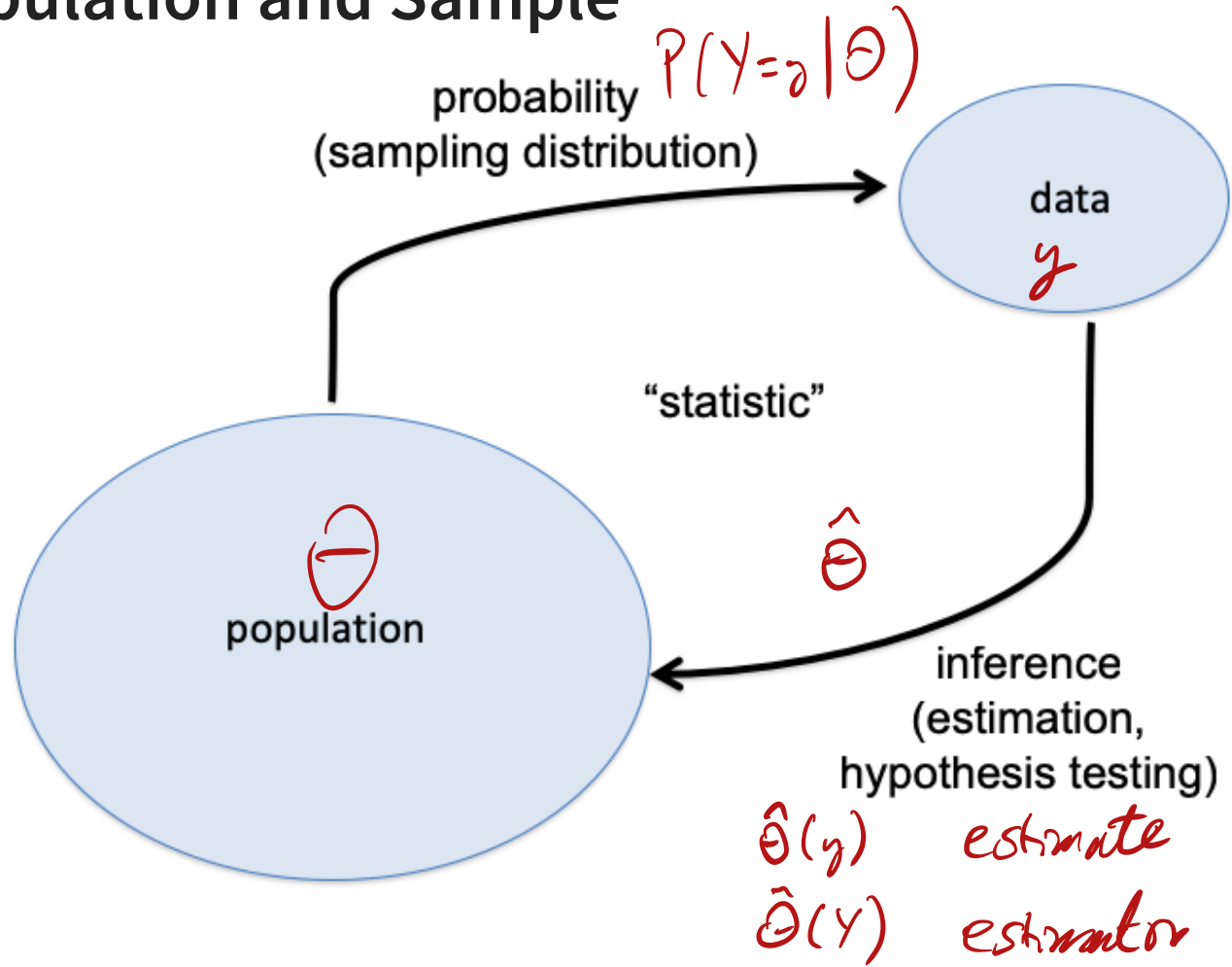
- Nectar

- Jupyterhub + sync Content.

# Logistics

- Use this link to pull all course content into your environment
- Link is on the course website. Will be used to sync all assignments.

# 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  $\theta$



*Exchangeability*



# 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

$$y_i \stackrel{iid}{\sim} P(Y|\theta)$$

$$L(\theta) \propto \prod_{i=1}^n P(y_i = y_i | \theta)$$

# Sufficient Statistics

vector  $(Y_1, \dots, Y_n)$

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

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

$$\text{E.g. } Y_1, \dots, Y_n \mid \bar{Y} \sim N(0, 1 - 1/n) \quad (\text{exercise})$$

$$\text{Bayesian: } P(\theta \mid y_1, \dots, y_n) = P(\theta \mid s(y_1, \dots, y_n))$$

# 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:

$$L(\theta) = \underbrace{h(y_1, \dots, y_n)}_{\text{no } \theta} \underbrace{g(s(y), \theta)}_{\text{has } \theta \text{ and sufficient stat.}}$$

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



# 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

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

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

$$X \sim \text{NB}(3, \theta) \quad \text{Obs: } x = 9$$

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

$$H_0: \theta = 1/2$$

$$\text{Bin: } p_{\text{binom}}(3, 12, \theta=1/2) \\ \approx 0.073$$

$$H_a: \theta < 1/2$$

$$\text{NB: } 1 - p_{\text{NBinom}}(8, 3, \theta=1/2)$$

# 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:  $I(\theta) = -E\left[\frac{d^2\ell(\theta; Y)}{d^2\theta}\right]$

Observed Info:  $-\frac{d^2\ell(\theta; y)}{d^2\theta}$

$L$ : likelihood  
 $\ell$ : log-likelihood

0.033

$$y_1, \dots, y_n \sim N(\mu, \sigma^2) \text{ known}$$

$$L(\mu) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}}$$

$$\propto \left( e^{-\frac{(\bar{y} - \mu)^2}{2\sigma^2/n}} \right)$$

$$l(\mu) = -\frac{(\bar{y} - \mu)^2}{2\sigma^2/n} \quad \leftarrow g(y, \theta)$$

$$I(\mu) = \frac{n}{\sigma^2}$$

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

$$l'' = -\frac{n}{\sigma^2}$$

# Fisher Information

# 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 w/ experience.*

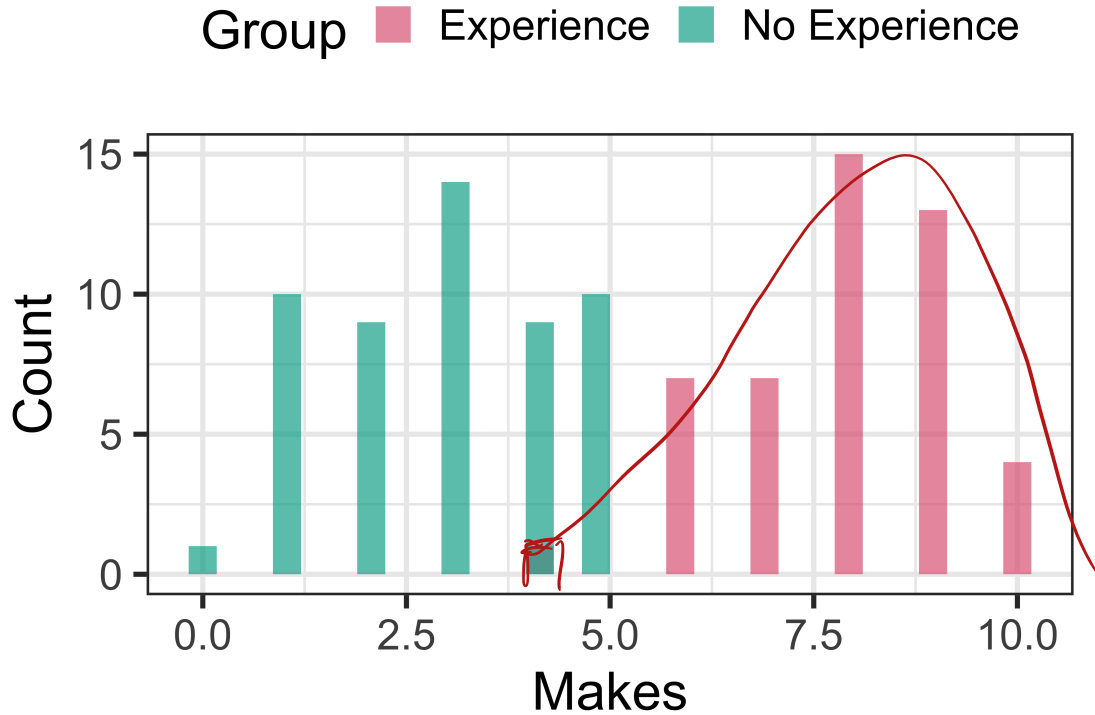
$$Y_i \sim \begin{cases} \text{Bin}(10, \theta_0) & \text{if } Z_i = 0 \\ \text{Bin}(10, \theta_1) & \text{if } Z_i = 1 \end{cases}$$

*Makes (out of 10)*

*Assume  $\theta_1 > \theta_0$*

*Make probabilities.*

# A Mixture Model



Note:  $z$  is observed

$$\prod_{i=1}^n P(Y_i, Z_i | \theta_1, \theta_0, \phi) = \prod_{i=1}^n [P(Y_i | Z_i, \theta_0, \theta_1) P(Z_i | \phi)]$$

$$\prod_{i=1}^n \left[ \phi \left( \cancel{\binom{10}{y_i}} \theta_1^{y_i} (1-\theta_1)^{10-y_i} \right)^{Z_i} \times \right. \\ \left. \left[ (1-\phi) \left( \cancel{\binom{10}{y_i}} \theta_0^{y_i} (1-\theta_0)^{10-y_i} \right)^{1-Z_i} \right] \right]$$

simplify



# 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

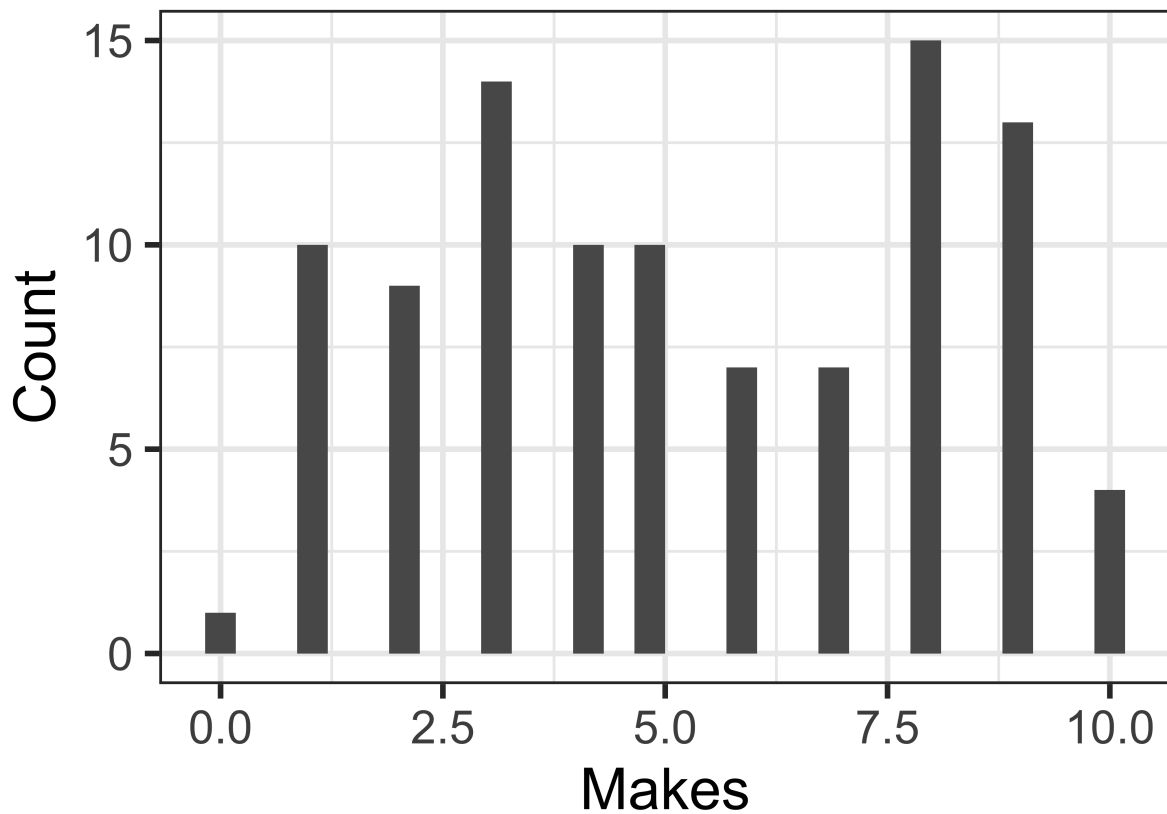
- A mixture model is a probabilistic model for representing the presence of subpopulations
- The subpopulation to which each individual belongs is not necessarily known
  - e.g. do we ask: “have you played basketball before?”
- When  $z_i$  is not observed, we sometimes refer to it as a clustering model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



# 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*

$$\begin{aligned} p(Y | \theta) &= \sum_z p(Y, Z = z | \theta) \\ &= \sum_z \underbrace{p(Y | Z = z, \theta)}_{\text{Mixture Weights}} p(Z = z | \theta) \end{aligned}$$

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

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

$p(Y|\theta)$        $p_k(Y|\theta)$

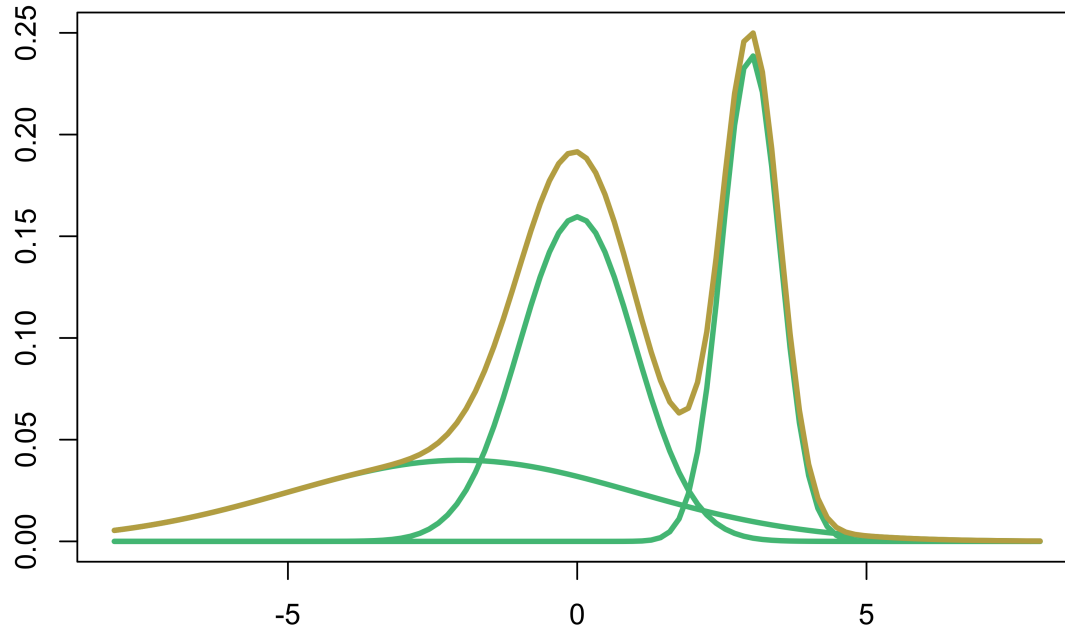


$$L(\theta_1, \theta_0, \phi) \propto \prod_{i=1}^n P(y_i | \theta_0, \theta_1, \phi)$$

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

$$\propto \prod_{i=1}^{100} \left[ \phi \cancel{\binom{10}{y_i}} \theta_1^{y_i} (1-\theta_1)^{10-y_i} + (1-\phi) \cancel{\binom{10}{y_i}} \theta_0^{y_i} (1-\theta_0)^{10-y_i} \right]$$

# Finite mixture models



# Infinite Mixture Models

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

Calculus, MOFs

$$P(Y|\tau^2, \sigma^2) = \int_{-\infty}^{\infty} P(Y|\mu, \sigma^2) P(\mu|\tau^2) d\mu$$

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

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



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

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

$$Y_i = \mu_i + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

# Infinite Mixture Models

$$\theta = \frac{1}{\beta} \rightarrow \text{Var} = 2\theta^2$$

Example 2: Poisson observations with random rates

$$\lambda \sim \text{Gamma}(\alpha, \beta) \rightarrow E[\lambda] = \frac{\alpha}{\beta}$$
$$Y \sim \text{Pois}(\lambda) \quad \text{Var}(\lambda) = \frac{\alpha}{\beta^2}$$

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

$$= \frac{\beta^\alpha}{\Gamma(\alpha)\Gamma(y+1)} \int \underbrace{\lambda^{y+\alpha-1} e^{-(\beta+1)\lambda}}_{\text{kernel of the gamma density.}} d\lambda$$

Gamma( $\alpha+y$ ,  $\beta+1$ )

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

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

# Infinite Mixture Models

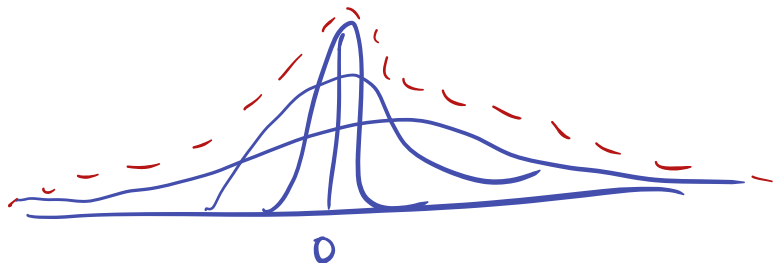
- Example 3: normal observations with exponentially distributed scale

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

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

Inv-Gauss

What is the distribution of  $Y_i$  given  $u$ ?



$$\frac{Z_1}{|Z_2|}$$

# Summary

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

# Assignments

- Read chapter 1-2 BDA3