

# The Zipper Example<sup>1</sup>

STA2101 Fall 2019

---

<sup>1</sup>See last slide for copyright information.

# Overview

- 1 Preparation
- 2 The Example
- 3 A better model
- 4 Identifiability
- 5 Maximum Likelihood Fails

## Preparation: Indicator functions

### Conditional expectation and the Law of Total Probability

$I_A(x)$  is the *indicator function* for the set  $A$ . It is defined by

$$I_A(x) = \begin{cases} 1 & \text{for } x \in A \\ 0 & \text{for } x \notin A \end{cases}$$

Also sometimes written  $I(x \in A)$ .

Using  $E(g(X)) = \sum_x g(x)p(x)$ , we have

$$\begin{aligned} E(I_A(X)) &= \sum_x I_A(x)p(x) \\ &= \sum_{x \in A} p(x) \\ &= P\{X \in A\} \end{aligned}$$

So the expected value of an indicator is a probability.

# If $X$ is continuous

$$\begin{aligned} E(I_A(X)) &= \int_{-\infty}^{\infty} I_A(x) f(x) dx \\ &= \int_A f(x) dx \\ &= P\{X \in A\} \end{aligned}$$

# Applies to conditional probabilities too

$$\begin{aligned} E(I_A(X)|Y) &= \sum_x I_A(x)p(x|Y), \text{ or} \\ &\int_{-\infty}^{\infty} I_A(x)f(x|Y) dx \\ &= Pr\{X \in A|Y\} \end{aligned}$$

So the conditional expected value of an indicator is a *conditional* probability.

Double expectation:  $E(g(X)) = E(E[g(X)|Y])$

$E(E[I_A(X)|Y]) = E[I_A(X)] = Pr\{X \in A\}$ , so

$$\begin{aligned} Pr\{X \in A\} &= E(E[I_A(X)|Y]) \\ &= E(Pr\{X \in A|Y\}) \\ &= \int_{-\infty}^{\infty} Pr\{X \in A|Y = y\} f_Y(y) dy \\ \text{or } &\sum_y Pr\{X \in A|Y = y\} p_Y(y) \end{aligned}$$

This is known as the *Law of Total Probability*

# The Zipper Example

Members of a Senior Kindergarten class (which we shall treat as a sample) try to zip their coats within one minute. We observe whether each child succeeds.

How about a model?

$Y_1, \dots, Y_n \stackrel{i.i.d.}{\sim} B(1, \theta)$ , where  $\theta$  is the probability of success.

Do you actually like this model?

# A better model than $Y_1, \dots, Y_n \stackrel{i.i.d.}{\sim} B(1, \theta)$

- Obviously, the probability of success is not the same for each child.
- Some kids are better at it than others, and some coats have easier zippers than others.
- Some children are almost certain to succeed, and others have almost no chance.

**Alternative Model:**  $Y_1, \dots, Y_n$  are independent random variables, with  $Y_i \sim B(1, \theta_i)$ .

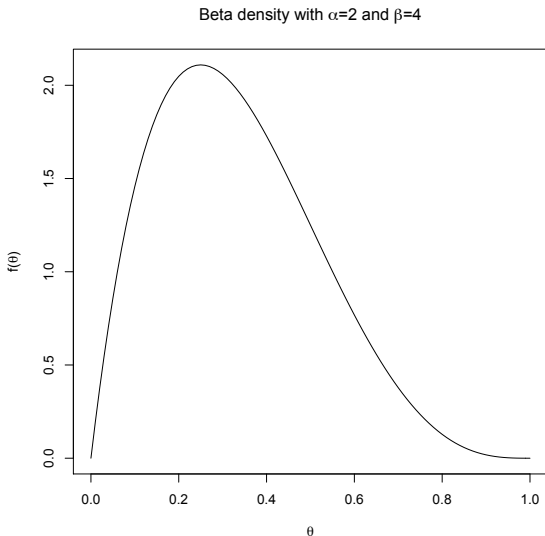


$Y_1, \dots, Y_n$  independent  $B(1, \theta_i)$ 

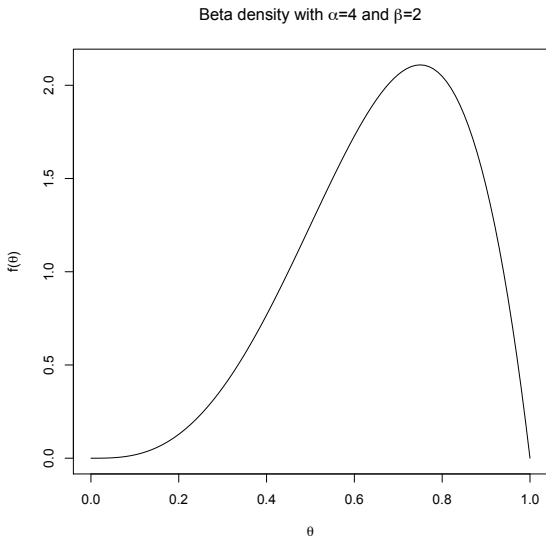
- This is a two-stage sampling model.
- First, sample from a population of children in which each child has a personal probability of success,  $\theta_i$ .
- Then for child  $i$ , use  $\theta_i$  to randomly generate success or failure.
- Note that  $\theta_1, \dots, \theta_n$  are random variables with some probability distribution.
- This distribution is supported on  $[0, 1]$
- How about a beta distribution?

$$f(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

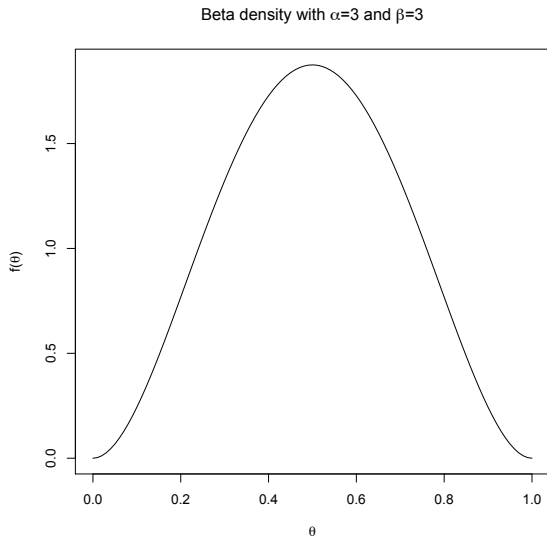
# Beta density is flexible



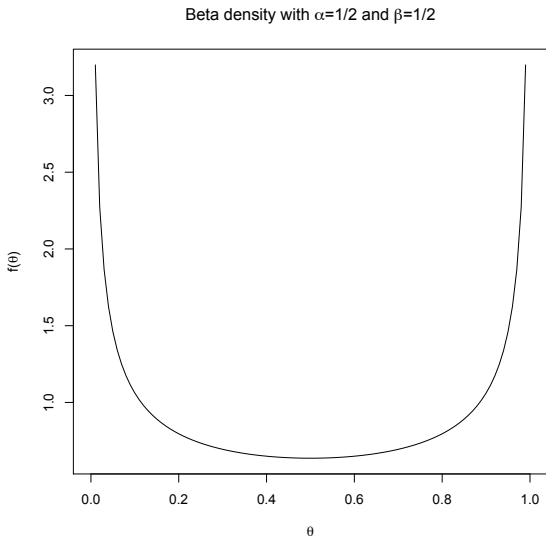
# Beta density is flexible



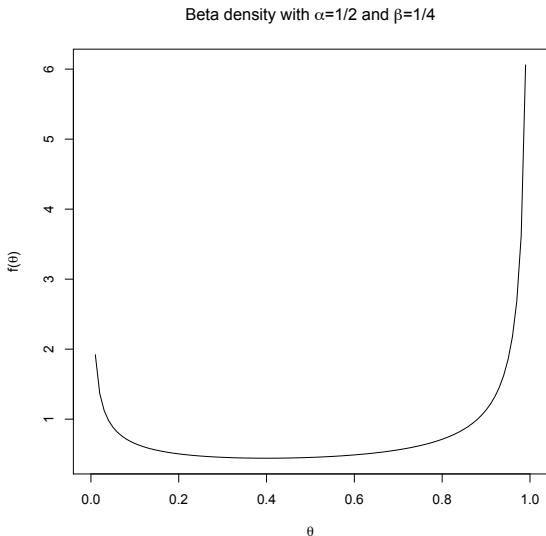
# Beta density is flexible



# Beta density is flexible



# Beta density is flexible



# Law of total probability

## Double expectation

$$\begin{aligned} P(Y_i = 1) &= \int_0^1 P(Y_i = 1|\theta_i) f(\theta_i) d\theta_i \\ &= \int_0^1 \theta_i f(\theta_i) d\theta_i \\ &= \int_0^1 \theta_i \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta_i^{\alpha-1} (1 - \theta_i)^{\beta-1} d\theta_i \\ &= \frac{\alpha}{\alpha + \beta} \end{aligned}$$

# Distribution of the observable data

Corresponds to the likelihood

$$P(\mathbf{Y} = \mathbf{y} | \alpha, \beta) = \prod_{i=1}^n \left( \frac{\alpha}{\alpha + \beta} \right)^{y_i} \left( 1 - \frac{\alpha}{\alpha + \beta} \right)^{1 - y_i}$$

- Distribution of the observable data depends on the parameters  $\alpha$  and  $\beta$  only through  $\frac{\alpha}{\alpha + \beta}$ .
- Infinitely many  $(\alpha, \beta)$  pairs yield the same distribution of the data.
- How could you use the data to decide which one is right?



# Parameter Identifiability

## The general idea

- The parameters of the Zipper Model are not *identifiable*.
- The model parameters cannot be recovered from the distribution of the sample data.
- And all you can ever learn from sample data is the distribution from which it comes.
- So there will be problems using the sample data for estimation and inference about the parameters.
- This is true *even if the model is completely correct*.

# Definitions

Connected to parameter identifiability

- A *Statistical Model* is a set of assertions that partly specify the probability distribution of the observable data.
- Suppose a statistical model implies  $\mathbf{D} \sim P_{\theta}, \theta \in \Theta$ . If no two points in  $\Theta$  yield the same probability distribution, then the parameter  $\theta$  is said to be *identifiable*.
- That is, identifiability means that  $\theta_1 \neq \theta_2$  implies  $P_{\theta_1} \neq P_{\theta_2}$ .
- On the other hand, if there exist distinct  $\theta_1$  and  $\theta_2$  in  $\Theta$  with  $P_{\theta_1} = P_{\theta_2}$ , the parameter  $\theta$  is *not identifiable*.

## An equivalent definition

Equivalent to  $\theta_1 \neq \theta_2 \Rightarrow P_{\theta_1} \neq P_{\theta_2}$

- The probability distribution is always a function of the parameter vector.
- If that function is one-to-one, the parameter vector is identifiable, because then  $\theta_1 \neq \theta_2$  yielding the same distribution could not happen.
- That is, if the parameter vector can somehow be recovered from the distribution of the data, it is identifiable.

# Theorem

If the parameter vector is not identifiable, consistent estimation for all points in the parameter space is impossible.



- Suppose  $\theta_1 \neq \theta_2$  but  $P_{\theta_1} = P_{\theta_2}$
- $T_n = T_n(D_1, \dots, D_n) \xrightarrow{P} \theta$  for all  $\theta \in \Theta$ .
- Distribution of  $T_n$  is identical for  $\theta_1$  and  $\theta_2$ .

## Why don't we hear more about identifiability?

- Consistent estimation indirectly proves identifiability.
- Because without identifiability, consistent estimation would be impossible.
- Any *function* of the parameter vector that can be estimated consistently is identifiable.

# Maximum likelihood fails for the Zipper Example

It has to fail.

$$L(\alpha, \beta) = \left(\frac{\alpha}{\alpha + \beta}\right)^{\sum_{i=1}^n y_i} \left(1 - \frac{\alpha}{\alpha + \beta}\right)^{n - \sum_{i=1}^n y_i}$$

$$\ell(\alpha, \beta) = \log \left( \left(\frac{\alpha}{\alpha + \beta}\right)^{\sum_{i=1}^n y_i} \left(1 - \frac{\alpha}{\alpha + \beta}\right)^{n - \sum_{i=1}^n y_i} \right)$$

Partially differentiate with respect to  $\alpha$  and  $\beta$ , set to zero, and solve.

## Two equations in two unknowns

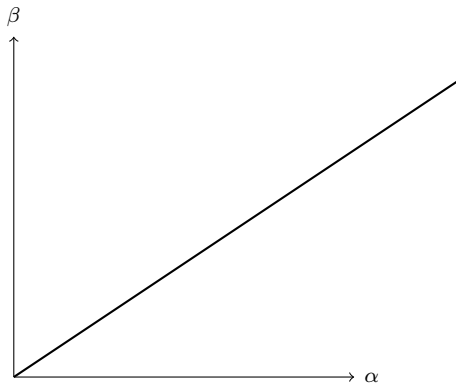
$$\begin{aligned}\frac{\partial \ell}{\partial \alpha} \stackrel{\text{set}}{=} 0 &\Rightarrow \frac{\alpha}{\alpha + \beta} = \bar{y} \\ \frac{\partial \ell}{\partial \beta} \stackrel{\text{set}}{=} 0 &\Rightarrow \frac{\alpha}{\alpha + \beta} = \bar{y}\end{aligned}$$

Any pair  $(\alpha, \beta)$  with  $\frac{\alpha}{\alpha + \beta} = \bar{y}$  will maximize the likelihood.

The MLE is not unique.

# What is happening geometrically?

$$\frac{\alpha}{\alpha + \beta} = \bar{y} \Leftrightarrow \beta = \left( \frac{1 - \bar{y}}{\bar{y}} \right) \alpha$$





# Look what has happened to us.

- We made an honest attempt to come up with a better model.
- And it *was* a better model.
- But the result was disaster.

## There is some good news.

Remember from earlier that by the Law of Total Probability,

$$P(Y_i = 1) = \int_0^1 \theta_i f(\theta_i) d\theta_i = E(\Theta_i)$$

- Even when the probability distribution of the (random) probability of success is completely unknown,
- We can estimate its expected value (call it  $\mu$ ) consistently with  $\bar{Y}_n$ .
- So that *function* of the unknown probability distribution is identifiable.
- And often that's all we care about anyway, say for comparing group means.
- So the usual procedures, based on a model nobody can believe (Bernoulli), are actually informative about a much more realistic model whose parameter vector is not fully identifiable.
- We don't often get this lucky.

## One more question about the parametric version

What would it take to estimate  $\alpha$  and  $\beta$  successfully?

- Get the children to try zipping their coats twice, say on two consecutive days.
- Assume their ability does not change, and conditionally on their ability, the two tries are independent.
- That will do it.
  
- This kind of thing often happens. When the parameters of a reasonable model are not identifiable, it is often possible to design a different way of collecting data so that the parameters *can* be identified.

## Moral of the story

- If you think up a better model for standard kinds of data, the parameters of the model may not be identifiable. You need to check.
- The problem is not with the model. It's with the data.
- The solution is better *research design*.

## Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistics, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The  $\text{\LaTeX}$  source code is available from the course website:  
<http://www.utstat.toronto.edu/~brunner/oldclass/2101f19>