

Estimating parameters

- Reminder: Hw 1 due Friday at 12pm (noon!)
- Study session:

TA: Jady Li

Location: Kirby 120

Time: Thursdays 7-8pm

Goal

Logistic regression model:

Random Component *Systematic*

$$Y_i \sim \text{Bernoulli}(\pi_i) \quad \log\left(\frac{\pi_i}{1 - \pi_i}\right) = \beta_0 + \beta_1 X_i$$

Given data $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$, how do I estimate β_0 and β_1 ?

- Start with: ignoring explanatory variable
focus on response

Motivating example

$$\approx \gamma \quad \approx 0$$

Y_i = result of flipping a coin (Heads or Tails)

Is Y_i a random variable? Yes!

- two possible outcomes
- we don't know what will happen when we flip the coin

\Rightarrow Bernoulli distribution for Y_i

Motivating example

Y_i = result of flipping a coin (Heads or Tails)

Let's make a model:

- + Step 1: Distribution of the response

$$Y_i \sim Bernoulli(\pi)$$

- + Step 2: Construct a model for the parameters

$$\pi = ??$$

Motivating example

Y_i = result of flipping a coin (Heads or Tails)

Suppose your friend estimates that the probability of heads is 0.9

- + $Y_i \sim Bernoulli(\pi)$
- + $\hat{\pi} = 0.9$

How can we assess whether this estimate $\hat{\pi}$ is reasonable?

Flip the coin (many times, ideally)

Motivating example

Suppose we flip the coin 5 times, and observe

$$y_1, \dots, y_5 = T, T, T, T, H$$

What is the probability of (i.e., how *likely* is) getting this string of flips if $\pi = 0.9$? Discuss with your neighbor for a minute, then we will discuss as a class.

$$\begin{aligned} P(T, T, T, T, H) &= (0.1)(0.1)(0.1)(0.1)(0.9) && \text{multiplication rule} \\ &= (0.1)^4 (0.9) && \text{for independent events!} \\ &= 0.00009 \end{aligned}$$

Likelihood

Definition: The *likelihood* $L(\text{Model}) = P(\text{Data}|\text{Model})$ of a model is the probability of the observed data, given that we assume a certain model and certain values for the parameters that define that model.

- + Model: $Y_i \sim \text{Bernoulli}(\pi)$, and $\hat{\pi} = 0.9$
- + Data: $y_1, \dots, y_5 = T, T, T, T, H$ ← "given" , or "conditional"
- + Likelihood: $L(\hat{\pi}) = P(y_1, \dots, y_5 | \pi = 0.9) = 0.00009$

"given $\pi = 0.9$, the probability of
 T, T, T, T, H = 0.00009"

Class Activity, Part I

https://sta214-s23.github.io/class_activities/ca_lecture_4.html

Class Activity

$$\begin{aligned}L(0.2) &= P(y_1, \dots, y_5 | \pi = 0.2) \\&= (0.2)(0.8)(0.8)(0.2)(0.8) = 0.020\end{aligned}$$

$$\begin{aligned}L(0.3) &= P(y_1, \dots, y_5 | \pi = 0.3) \\&= (0.3)(0.7)(0.7)(0.3)(0.7) = 0.031\end{aligned}$$

Which value, 0.2 or $\underline{0.3}$, seems more reasonable?



higher likelihood!

\Rightarrow Data has a higher probability if $\pi = 0.3$
than if $\pi = 0.2$

Class Activity

Which value of $\hat{\pi}$ in the table would you pick?

$\hat{\pi} = 0.4$ has the highest likelihood

(of the values considered in the table)

Maximum likelihood

Maximum likelihood principle: estimate the parameters to be the values that maximize the likelihood

(i.e., value that makes the observed data most probable)

$\hat{\pi}$	Likelihood
0.30	0.031
0.35	0.033
0.40	0.036
0.45	0.033

Maximum likelihood estimate: $\hat{\pi} = 0.4$

Maximum likelihood

Maximum likelihood principle: estimate the parameters to be the values that maximize the likelihood

Steps for maximum likelihood estimation:

- + *Likelihood:* For each potential value of the parameter, compute the likelihood of the observed data
- + *Maximize:* Find the parameter value that gives the largest likelihood

Maximum likelihood for logistic regression

$$\log\left(\frac{\pi_i}{1 - \pi_i}\right) = \beta_0 + \beta_1 \text{Size}_i \quad \pi_i = \frac{\exp\{\beta_0 + \beta_1 \text{Size}_i\}}{1 + \exp\{\beta_0 + \beta_1 \text{Size}_i\}}$$

Observed data:

$$\hat{\pi}_i = \frac{\exp\{-2 + 0.5(6)\}}{1 + \exp\{-2 + 0.5(6)\}} = \frac{e}{1+e}$$

Tumor cancerous	Yes	No	No	Yes	No
Size of tumor (cm)	6	1	0.5	4	1.2

Suppose $\beta_0 = -2$, $\beta_1 = 0.5$. How would I compute the likelihood?

- 1) Plug in $\beta_0 = -2$, $\beta_1 = 0.5$, and tumor size to get $\hat{\pi}_i$
- 2) Multiply probabilities together to get likelihood.

Class Activity, Part II

https://sta214-s23.github.io/class_activities/ca_lecture_4.html

Class Activity

$$\hat{\pi}_i = \frac{\exp\{-2 + 0.5 \text{Size}_i\}}{1 + \exp\{-2 + 0.5 \text{Size}_i\}}$$

Tumor cancerous	Yes	No	No	Yes	No
Size of tumor (cm)	6	1	0.5	4	1.2
$\hat{\pi}_i$	0.73	0.18	0.15	0.5	0.2

Likelihood = $(0.73)(1 - 0.18)(1 - 0.15)(0.5)(1 - 0.2)$
 ≈ 0.2

Maximum likelihood for logistic regression

Likelihood:

- + For estimates $\hat{\beta}_0$ and $\hat{\beta}_1$, $\hat{\pi}_i = \frac{\exp\{\hat{\beta}_0 + \hat{\beta}_1 X_i\}}{1 + \exp\{\hat{\beta}_0 + \hat{\beta}_1 X_i\}}$
- + $L(\hat{\beta}_0, \hat{\beta}_1) = P(Y_1, \dots, Y_n | \hat{\beta}_0, \hat{\beta}_1)$

Maximize:

- + Choose $\hat{\beta}_0, \hat{\beta}_1$ to maximize $L(\hat{\beta}_0, \hat{\beta}_1)$

- 1) Do it all by hand
- 2) Do it on a computer
- 3) Calculus? (next time)

So far, we only considered a few values for β_0 and β_1 . How should we check other values, to make sure our estimates actually maximize likelihood?

Computing likelihood in R

Observed data: T, T, T, T, H

- + We are going to consider several different potential values for $\hat{\pi}$:

$$0, 0.1, 0.2, 0.3, \dots, 0.9, 1$$

- + For each potential value, we will compute the likelihood:

$$L(\hat{\pi}) = (1 - \hat{\pi})^4(\hat{\pi})$$

- + We then see which value has the highest likelihood.
- + Is this all possible values? No, but let's start here.

R code

Step 1: Specify potential values of $\hat{\pi}$

```
# List the values for  $\hat{\pi}$ 
pi_hat <- seq(from = 0, to = 1, by = 0.1)
```

"sequence" → 0 ← 0.1 0.2 0.3 ... 0.9 → 1
0.1 0.1

To get values 0 0.05 0.1 0.15 etc.
change by = 0.05

R code

```
# List the values for pi hat  
pi_hat <- seq(from = 0, to = 1, by = 0.1)
```

```
# Create a space to store the likelihoods  
likelihood <- rep(0, length(pi_hat))
```

"repeat" ↙ O O O ... O (repeated 11 times)

will fill in likelihood as we do calculations

R code

```
# List the values for pi hat
pi_hat <- seq(from = 0, to = 1, by = 0.1)

# Create a space to store the likelihoods
likelihood <- rep(0,length(pi_hat))

# Compute and store the likelihoods
for( i in 1:length(pi_hat) ){
  likelihood[i] <- pi_hat[i]*(1-pi_hat[i])^4
}
```

R code

```
# List the values for pi hat
pi_hat <- seq(from = 0, to = 1, by = 0.1)

# Create a space to store the likelihoods
likelihood <- rep(0,length(pi_hat))

# Compute and store the likelihoods
for( i in 1:length(pi_hat) ){
  likelihood[i] <- pi_hat[i]*(1-pi_hat[i])^4
}
```

Run this code in your R console. Which value of $\hat{\pi}$ gives the highest likelihood?

Results

pi_hat	likelihood
0.0	0.00000
0.1	0.06561
0.2	0.08192
0.3	0.07203
0.4	0.05184
0.5	0.03125
0.6	0.01536
0.7	0.00567
0.8	0.00128
0.9	0.00009
1.0	0.00000