

Lecture 3: Maximum likelihood estimation

Recap: ways of fitting a *linear* regression model

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \cdots + \beta_k X_{i,k} + \varepsilon_i \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma_\varepsilon^2)$$

Suppose we observe data $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$, where $X_i = (1, X_{i,1}, \dots, X_{i,k})^T$.

How do we fit this linear regression model? That is, how do we estimate

$$\beta = (\beta_0, \beta_1, \dots, \beta_k)^T$$

Minimize sum of squared residuals

Maximum likelihood estimation (assume distribution for Y)

Projection

Minimize SSE

$$SSE = \sum_{i=1}^n (\tau_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots - \beta_n x_{in})^2$$

Find $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_n)$ minimize SSE

$$\frac{\partial SSE}{\partial \beta_0}$$

set

0

n+1 equations

n+1 unknowns

$$\frac{\partial SSE}{\partial \beta_1}$$

set

0

$\hat{\beta}$ vector that solves

this system

$$\frac{\partial SSE}{\partial \beta_n}$$

set

0

Projection

$$\gamma = \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_n \end{pmatrix} \in \mathbb{R}$$

$$X = \begin{pmatrix} 1 & x_{11} & \cdots \\ 1 & x_{21} & \cdots \\ \vdots & \vdots & \ddots \\ 1 & x_{n1} & \cdots \end{pmatrix} \in \mathbb{R}^{n \times n+1}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_n \end{pmatrix}$$

$$\hat{\gamma} = X \hat{\beta} \quad (\hat{\gamma} \in \text{colspace}(X))$$

want $\hat{\gamma}$ "close" to γ

$$\text{minimize } \|\gamma - \hat{\gamma}\| = \sqrt{\text{SSE}}$$

(Euclidean norm)

\Leftrightarrow minimize SSE

$$\gamma_i \sim N(\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}, \sigma^2_\varepsilon)$$

$$e^a e^b = e^{a+b}$$

$$L(\beta_0, \dots, \beta_k, \sigma^2_\varepsilon \mid (x_1, \gamma_1), \dots, (x_n, \gamma_n))$$

(independence)

$$\begin{aligned}
 &= \prod_{i=1}^n f(\gamma_i \mid \beta_0, \dots, \beta_k, \sigma^2_\varepsilon, x_i) \\
 &\quad \underbrace{\qquad\qquad\qquad}_{\sqrt{2\pi\sigma^2_\varepsilon}} \exp \left\{ -\frac{1}{2\sigma^2_\varepsilon} (\gamma_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2 \right\} \\
 &= (2\pi)^{-\frac{n}{2}} \sigma^{-n} \exp \left\{ -\frac{1}{2\sigma^2_\varepsilon} \sum_{i=1}^n (\gamma_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2 \right\} \\
 &\quad \underbrace{\qquad\qquad\qquad}_{SSE}
 \end{aligned}$$

maximizing likelihood \Leftrightarrow minimizing SSE
 (for normal data)

Summary: three ways of fitting linear regression models

- Minimize SSE, via derivatives of

$$\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i,1} - \cdots - \beta_k X_{i,k})^2$$

SSE: might not work for
binary data (what
should the residual
be??)

- Minimize $\|Y - \hat{Y}\|$ (equivalent to minimizing SSE)

- Maximize likelihood (for *normal* data, equivalent to
minimizing SSE)

← might work for logistic, but
w/ different distribution
(Bernoulli?)

Which of these three methods, if any, is appropriate for
fitting a logistic regression model? Do any changes need to
be made for the logistic regression setting?

Step back: likelihoods and estimation

Let $Y \sim \text{Bernoulli}(p)$ be a Bernoulli random variable, with $p \in [0, 1]$. We observe 5 independent samples from this distribution:

$$Y_1 = 1, Y_2 = 1, Y_3 = 0, Y_4 = 0, Y_5 = 1$$

The true value of p is unknown, so two friends propose different guesses for the value of p : 0.3 and 0.7. Which do you think is a “better” guess?

Sample proportion: 0.6 (observed to 0.7)

$$\begin{aligned} p(\text{data} \mid p=0.3) &= (0.3)(0.3)(1-0.3)(1-0.3)(0.3) \\ &= 0.3^3 (0.7)^2 = 0.013 \end{aligned}$$
$$p(\text{data} \mid p=0.7) = 0.7^3 0.3^2 = 0.031$$

intuition: choose value of p which makes data more “likely”

Likelihood

Definition: Let $\mathbf{Y} = (Y_1, \dots, Y_n)$ be a sample of n observations, and let $f(\mathbf{y}|\theta)$ denote the joint pdf or pmf of \mathbf{Y} , with parameter(s) θ . The *likelihood function* is

$$\underbrace{L(\theta|\mathbf{Y})}_{\text{function of } \theta, \text{ given } \mathbf{Y}} = f(\mathbf{Y}|\theta) \leftarrow \begin{array}{l} \text{"probability"} \\ \text{of the observed data,} \\ \text{if } \theta \text{ is the} \\ \text{parameter} \end{array}$$

$L(\theta|\mathbf{Y})$: condition on the observed data. Want to know how "probability" (joint density / mass function) of \mathbf{Y} changes as a function of θ

Since $f(\mathbf{y}|\theta) \geq 0$, $L(\theta|\mathbf{Y}) \geq 0 \quad \forall \theta$

Special case: Y_1, \dots, Y_n iid w/ pdf or pmf f

$$L(\theta|\mathbf{Y}) = \prod_{i=1}^n f(Y_i|\theta)$$

Example: Bernoulli data

Let $\gamma_1, \dots, \gamma_n$ ~ Bernoulli (p)

$$f(y|p) = p^y (1-p)^{1-y}$$

$$P(\gamma=1) = p^1 (1-p)^{1-1} = p$$

$$P(\gamma=0) = p^0 (1-p)^{1-0} = 1-p$$

$$\begin{aligned} L(p | \gamma_1, \dots, \gamma_n) &= \prod_{i=1}^n f(\gamma_i | p) \\ &= \prod_{i=1}^n p^{\gamma_i} (1-p)^{1-\gamma_i} \\ &= \frac{\sum \gamma_i}{p} (1-p)^{n - \sum \gamma_i} \end{aligned}$$

Ex: $\gamma = (1, 1, 0, 0, 1)$

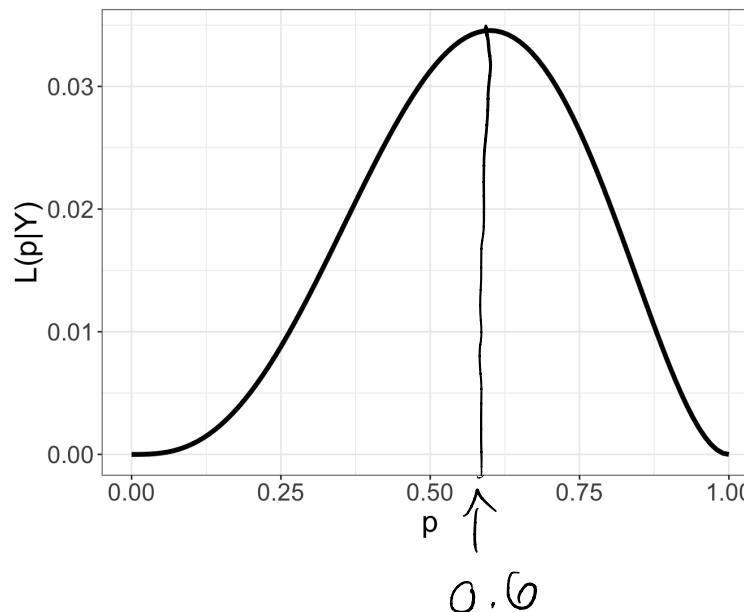
$$L(p | \gamma) = p^3 (1-p)^2$$

Example: Bernoulli data

$Y_1, \dots, Y_5 \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$, with observed data

$$Y_1 = 1, Y_2 = 1, Y_3 = 0, Y_4 = 0, Y_5 = 1$$

$$L(p|Y) = p^3(1-p)^2$$



Maximum likelihood estimator

Definition: Let $\mathbf{Y} = (Y_1, \dots, Y_n)$ be a sample of n observations. The *maximum likelihood estimator* (MLE) is

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta | \mathbf{Y})$$

argmax θ means "value of θ that maximizes ..."

Example: Bernoulli(p)

$y_1, \dots, y_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$

$$L(p|y) = p^{\sum_i y_i} (1-p)^{n - \sum_i y_i}$$

Now maximize to estimate p

- ① Take log to make life easier (log is monotone increasing function, so if \hat{p} maximizes $\log L(p|y)$, then \hat{p} maximizes $L(p|y)$)

$$\ell(p|y) = \log L(p|y) = (\sum_i y_i) \log p + (n - \sum_i y_i) \log(1-p)$$

- ② Differentiate wrt parameter of interest:

$$\frac{\partial}{\partial p} \ell(p|y) = \frac{\sum_i y_i}{p} - \frac{(n - \sum_i y_i)}{1-p} \stackrel{\text{set}}{=} 0$$

$$\frac{\partial}{\partial p} \ell(p|Y) = \frac{\sum_i Y_i}{p} - \frac{(n - \sum_i Y_i)}{1-p} \stackrel{\text{set}}{=} 0$$

$$\frac{\sum_i Y_i}{p} = \frac{(n - \sum_i Y_i)}{1-p} \Rightarrow \frac{1-p}{p} = \frac{n - \sum_i Y_i}{\sum_i Y_i}$$

$$\Rightarrow \frac{1}{p} = \frac{n}{\sum_i Y_i} \Rightarrow p = \frac{1}{n} \sum_i Y_i$$

(sample proportion!)

So $\ell(p|Y)$ has a max or min at $p = \frac{1}{n} \sum_i Y_i$

check second derivative:

$$\frac{\partial^2}{\partial p^2} \ell(p|Y) \Big|_{p=\frac{1}{n} \sum_i Y_i} = -\frac{\sum_i Y_i}{p^2} - \frac{(n - \sum_i Y_i)}{(1-p)^2}$$

$\Rightarrow \hat{p} = \frac{1}{n} \sum_i Y_i$ maximizes $\ell(p|Y)$

check $p = 0$ and $p = 1$

if $p = 0$: $L(p|Y) = \begin{cases} 0 & \text{if any } Y_i = 1 \\ 1 & \text{if all } Y_i = 0 \end{cases}$

$p = 1$: $L(p|Y) = \begin{cases} 0 & \text{if any } Y_i = 0 \\ 1 & \text{if all } Y_i = 1 \end{cases}$

if all $Y_i = 1 \Rightarrow \hat{p} = 1 = \frac{1}{n} \sum_i Y_i$

if all $Y_i = 0 \Rightarrow \hat{p} = 0 = \frac{1}{n} \sum_i Y_i$

