

Maximum likelihood estimation

Recap: ways of fitting linear regression models

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

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

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$

Possible methods

- Minimize SSE
- Projection (\Leftrightarrow to minimizing SSE)
- Maximize likelihood

$$Y_i \sim N(\beta_0 + \beta_1 X_{i1} + \cdots + \beta_n X_{in}, \sigma^2_\varepsilon)$$

$$L(\beta_0, \dots, \beta_n, \sigma^2_\varepsilon \mid (x_1, y_1), \dots, (x_n, y_n))$$

$$\begin{aligned} &= \prod_{i=1}^n f(Y_i \mid \beta_0, \dots, \beta_n, \sigma^2_\varepsilon, x_i) \\ &= \prod_{i=1}^n \underbrace{\frac{1}{\sqrt{2\pi\sigma^2_\varepsilon}} \exp\left\{-\frac{1}{2\sigma^2_\varepsilon} (Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_n X_{in})^2\right\}} \\ &= (2\pi)^{-\frac{n}{2}} \sigma^{-n} \exp\left\{-\frac{1}{2\sigma^2_\varepsilon} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_n X_{in})^2\right\} \end{aligned}$$

SSE!

maximizing likelihood \Leftrightarrow minimize SSE
 (for normal data)

Summary: three ways of fitting linear regression models

- + Minimize SSE, via derivatives of

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

- + Minimize $\|Y - \hat{Y}\|$ (equivalent to minimizing SSE)
- + Maximize likelihood (for *normal* data, equivalent to minimizing SSE) \hat{Y} could be appropriate for logistic regression, but need to use Bernoulli instead of Normal

linear regression:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{i,1} + \cdots + \hat{\beta}_k X_{i,k}$$

$$\Rightarrow Y_i - \hat{Y}_i \text{ makes sense}$$

not appropriate for logistic regression b/c we don't have the same idea of a residual

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?

Idea: pick parameter value for which the observed data has the highest "probability"

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 (closer to 0.7)

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

$$P(\text{data} \mid p = 0.7) = 0.7^3 0.3^2 = 0.031$$

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

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

function of \mathbf{Y} , given θ

- $L(\theta|\mathbf{Y})$: condition on the observed data. Want to know how "probability" (joint distribution) 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 $\gamma_1, \dots, \gamma_n$ $\stackrel{iid}{\sim}$ Bernoulli(p) $f(y|p) = p^y (1-p)^{1-y}$

$$\begin{aligned} L(p|y) &= \prod_{i=1}^n f(\gamma_i|p) = \prod_{i=1}^n p^{\gamma_i} (1-p)^{1-\gamma_i} \\ &= p^{\sum_{i=1}^n \gamma_i} (1-p)^{n - \sum_{i=1}^n \gamma_i} \end{aligned}$$

Example : $y = (1, 1, 0, 0, 1)$

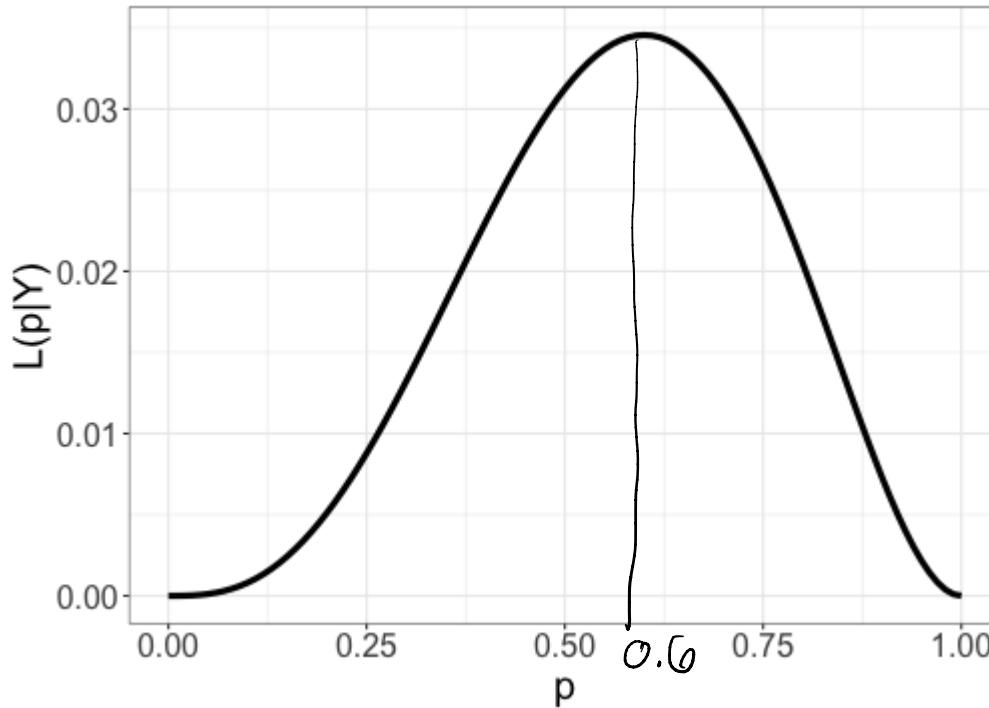
$$L(p|y) = p^3 (1-p)^2$$

Example: Bernoulli data

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

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

$$L(p|\mathbf{Y}) = p^3(1-p)^2$$



Maximum likelihood estimator

Definition: Let $\mathbf{Y} = (Y_1, \dots, Y_n)$ be a sample of n observations.

The *maximum likelihood estimator* (MLE) is

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

$\operatorname{argmax}_{\theta}$ means "value of θ that maximizes..."¹¹

Example: Bernoulli(p)

$\gamma_1, \dots, \gamma_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$

$$L(p|\gamma) = p^{\sum_i \gamma_i} (1-p)^{n - \sum_i \gamma_i}$$

Now maximize!

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

$$\ell(p|\gamma) = \log L(p|\gamma) = (\sum_i \gamma_i) \log p + (n - \sum_i \gamma_i) \log(1-p)$$

- ② Differentiate wrt parameter of interest:

$$\frac{\partial}{\partial p} \ell(p|\gamma) = \frac{\sum_i \gamma_i}{p} - \frac{(n - \sum_i \gamma_i)}{1-p}$$

③ set $\frac{\partial}{\partial p} \ell(p|Y) = 0$ & solve

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

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

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

$$\Rightarrow \frac{1}{p} = \frac{n}{\sum_i Y_i} \Rightarrow p = \frac{\sum_i Y_i}{n} \quad (\text{sample proportion})$$

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

④

Check max. or min:

$$\left. \frac{\partial^2}{\partial p^2} \ell(p|Y) \right|_{p=\bar{Y}} = \left. -\frac{\sum_i Y_i}{p^2} - \frac{(n - \sum_i Y_i)}{(1-p)^2} \right|_{p=\bar{Y}} < 0$$

\Rightarrow maximum!

⑤ Check boundary points: $p=0$ and $p=1$

$$l(p|Y) = \begin{cases} n \log(1-p) & p=0 \\ n \log p & p=1 \end{cases}$$

If either $p=0$ or $p=1$ has highest likelihood,

$$p = \frac{\sum_i Y_i}{n} \text{ anyway}$$

So $\hat{p} = \frac{\sum_i Y_i}{n} = \bar{Y}$

Example: $N(\theta, 1)$

$\theta \in (-\infty, \infty)$

$y_1, \dots, y_n \stackrel{iid}{\sim} N(\theta, 1)$

$$f(y|\theta) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-\theta)^2}$$

$$L(\theta|y) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(y_i - \theta)^2\right\}$$

$$= (2\pi)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2} \sum_i (y_i - \theta)^2\right\}$$

$$\Rightarrow \ell(\theta|y) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n (y_i - \theta)^2$$

$$\frac{\partial}{\partial \theta} \ell(\theta|y) = -\frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \theta} ((y_i - \theta)^2)$$

$$= -\frac{1}{2} \sum_{i=1}^n 2(y_i - \theta)(-1) = \sum_{i=1}^n (y_i - \theta)$$

$$\Rightarrow \sum_{i=1}^n y_i = \sum_{i=1}^n \overset{\text{set}}{=} \theta = n\theta \Rightarrow \theta = \frac{\sum y_i}{n} = \bar{y}$$

Example: $Uniform(0, \theta)$