

# Lecture 5: Maximum likelihood estimation

Ciaran Evans

## Recap: maximum likelihood estimation

**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)  $\theta$ . The *likelihood function* is

$$L(\theta|\mathbf{Y}) = f(\mathbf{Y}|\theta)$$

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

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

Example:  $N(\mu, \sigma^2)$

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

$$\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}$$

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

$$\sigma^2 > 0$$

$$L(\theta | Y) = (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2 \right\}$$

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

Plan: maximize wrt  $\mu$ , then maximize wrt  $\sigma^2$

$$\frac{\partial}{\partial \mu} \ell(\theta | Y) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu) \stackrel{\text{set}}{=} 0$$

$$= \sum_{i=1}^n y_i = n\mu \Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$$

For any value of  $\sigma^2$ :

$$\begin{aligned} L(\theta|\gamma) &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (\gamma_i - \mu)^2 \right\} \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (\gamma_i - \bar{\gamma})^2 \right\} \end{aligned}$$

Now we want to maximize

$$(2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (\gamma_i - \bar{\gamma})^2 \right\}$$

$$\ell^*(\sigma^2|\gamma) = \underbrace{-\frac{n}{2} \log(2\pi\sigma^2)}_{-\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2)} - \frac{1}{2\sigma^2} \sum_{i=1}^n (\gamma_i - \bar{\gamma})^2$$

$$\frac{d\ell^*}{d\sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (\gamma_i - \bar{\gamma})^2 \quad \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \boxed{\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (\gamma_i - \bar{\gamma})^2 \\ \hat{\mu} &= \bar{\gamma} \end{aligned}}$$

## Linear regression with normal errors

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

$$\mu_i = \beta_0 + \beta_1 X_{i,1} + \cdots + \beta_k X_{i,k}$$

Suppose we observe independent samples  $(X_1, Y_1), \dots, (X_n, Y_n)$ .  
Write down the likelihood function

$$L(\beta | \mathbf{X}, \mathbf{Y}) \propto \prod_{i=1}^n f(Y_i | \beta, X_i)$$

$$L(\beta | X, Y) = f(X, Y | \beta) = \prod_{i=1}^n f(x_i, y_i | \beta) \quad (\text{independence})$$

$(x_i, y_i)$  are iid from the same joint distribution

$$= \prod_{i=1}^n \underbrace{f(x_i | \beta)}_{f(x_i)} f(y_i | x_i, \beta)$$

but  $y_i | x_i$  is different depending on  $x_i$

$$= \left( \prod_{i=1}^n f(x_i) \right) \left( \prod_{i=1}^n f(y_i | x_i, \beta) \right)$$

$x_i$  doesn't depend on  $\beta$

$$\Rightarrow L(\beta | X, Y) \propto \prod_{i=1}^n f(y_i | x_i, \beta)$$

$\beta$  describes  $y_i | x_i$

$$f(y_i | x_i, \beta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - \mu_i)^2 \right\}$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - \beta^T x_i)^2 \right\}$$

$\mu_i$  can be different for each observation  
 $\mu_i = \beta_0 + \beta_1 x_{i1} + \dots$

$$L(\beta | X, Y) \propto (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta^T x_i)^2 \right\}$$

$\beta^T x_i$      $x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix}$

maximize  $L$  wrt  $\beta$ : minimize  $\sum_i (y_i - \beta^T x_i)^2$  (SSE again!)

Maximizing  $L$  wrt  $\beta$  means

$$\begin{aligned} \text{minimize} \quad \sum_{i=1}^n (\gamma_i - \beta^T x_i)^2 &= \| \gamma - X\beta \|^2 \\ &= (\gamma - X\beta)^T (\gamma - X\beta) \end{aligned}$$

$$\gamma = \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_d \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{pmatrix}$$

(design matrix)

$$\text{Let } u = \gamma - X\beta$$

$$\begin{aligned} \text{minimize} \quad & u^T u \\ \frac{\partial}{\partial \beta} \quad & u^T u \end{aligned}$$

Matrix chain rule:

$$\begin{aligned} \frac{\partial}{\partial \beta} u^T u &= \underbrace{\left( \frac{\partial u}{\partial \beta} \right)}_{\frac{\partial}{\partial \beta} (\gamma - X\beta)} \underbrace{\left( \frac{\partial}{\partial u} u^T u \right)}_{2u} \\ &= - \frac{\partial}{\partial \beta} (X\beta) = -X^T \end{aligned}$$

$$\begin{aligned} &= -2X^T (\gamma - X\beta) \stackrel{\text{set}}{=} 0 \\ \Rightarrow X^T (\gamma - X\beta) &= 0 \\ \Rightarrow X^T \gamma &= X^T X \beta \\ \Rightarrow \hat{\beta} &= (X^T X)^{-1} X^T \gamma \end{aligned}$$

## Logistic regression

$$Y_i \sim \text{Bernoulli}(p_i)$$

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_{i,1} + \cdots + \beta_k X_{i,k}$$

Suppose we observe independent samples  $(X_1, Y_1), \dots, (X_n, Y_n)$ .

Write down the likelihood function

$$L(\beta|\mathbf{X}, \mathbf{Y}) \propto \prod_{i=1}^n f(Y_i|\beta, X_i)$$