# Lecture 3: Maximum likelihood estimation

Ciaran Evans

## Logistics

- HW 1 due Friday on Canvas
- Office hours:
    - Wednesday 2-3 pm
    - Thursday 9:30-10:30 am

- Bowling on Friday!

# Motivation: fitting a *linear* regression model

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

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

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

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

Minimize sum of squared errors (aka residual sum of squares)

choose $\beta = (\beta_0, \beta_1, ..., \beta_k)^T$ to minimize

$$SSE = \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_{i,1} - \beta_2 X_{i,2} - \cdots - \beta_k X_{i,k})^2$$

$$\frac{\partial SSE}{\partial \beta_0} \overset{set}{=} 0$$

$$\vdots$$

$$\frac{\partial SSE}{\partial \beta_k} \overset{set}{=} 0$$

$k+1$ unknowns

$k+1$ equations

estimate $\hat{\beta} = \begin{pmatrix} \hat{\beta_0} \\ \hat{\beta_1} \\ \vdots \\ \hat{\beta_k} \end{pmatrix}$ solves this system
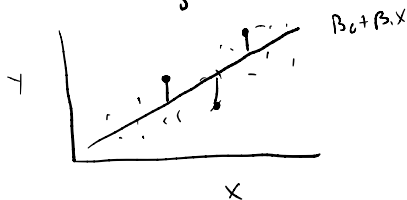
# Fitting a *logistic* regression model?

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

**Question:** Should we minimize a similar sum of squares for a *logistic* regression model?  No
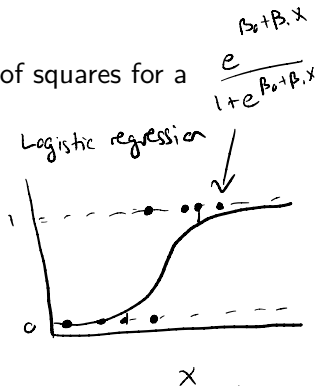
$\beta_0 + \beta_1 X$

$\dfrac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$

Linear regression

$Y$

$\beta_0 + \beta_1 X$

$X$

Logistic regression

1

0

$X$

Linear regression:

$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$

error terms

$Y_i \sim Bernoulli(P_i)$

$\log\left(\dfrac{P_i}{1 - P_i}\right) = \beta_0 + \beta_1 X_i$

not same idea of "residual"

# Motivation: likelihoods and estimation

Let $Y \sim Bernoulli(p)$ 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)

$P(\text{data} \mid p = 0.3) = (0.3)^3 (0.7)^2 = 0.013$

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

Intuition: choose value of $p$ which makes data more "likely"

# Likelihood

**Definition:** Let $\mathbf{Y} = (Y_1, ..., Y_n)$ 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

$$\underbrace{L(\theta|\mathbf{Y})}_{\substack{\text{function of } \theta, \\ \text{given observed} \\ \text{data } \mathbf{Y}}} = f(\mathbf{Y}|\theta) \underset{\substack{\text{"probability"} \\ \text{of the observed data,} \\ \underline{\text{if}} \quad \theta \quad \text{is true} \\ \text{true parameter}}}{\longleftarrow}$$

$L(\theta|Y)$ : condition on observed data, and we want to know how the joint density/ mass function of $Y$ changes as a function of $\theta$

$L(\theta|Y) \geq 0$    since    $f(y|\theta) \geq 0$

<u>Special case</u>:   $Y_1, ..., Y_n$ iid    $L(\theta|Y) = \prod_{i=1}^{n} f(Y_i|\theta)$

# Example: Bernoulli data

Let $Y_1, ..., Y_n \overset{iid}{\sim}$ Bernoulli$(p)$

$P(Y=y)$

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

(single observation)

$y \in \{0, 1\}$

$$L(p | Y_1, ..., Y_n) = \prod_{i=1}^{n} f(Y_i | p)$$

$$= \prod_{i=1}^{n} p^{Y_i} (1-p)^{1-Y_i}$$

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

$\underline{Ex}$: $Y = (1, 1, 0, 0, 1)$
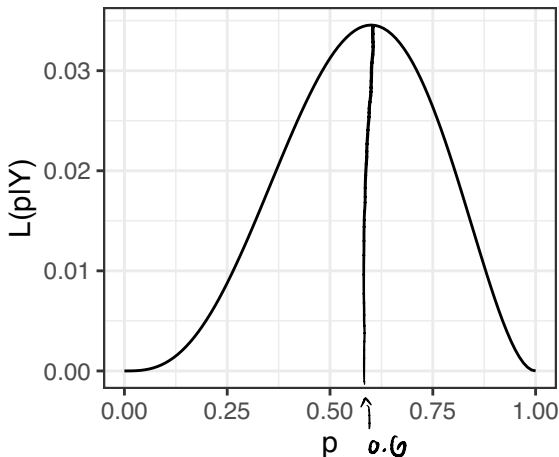
$$L(p | Y) = p^3 (1-p)^2$$

# Example: Bernoulli data

$Y_1, ..., Y_5 \overset{iid}{\sim} Bernoulli(p)$, 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 $\mathbf{Y} = (Y_1, ..., Y_n)$ be a sample of $n$ observations. The *maximum likelihood estimator* (MLE) is

$$\widehat{\theta} = \text{argmax}_\theta \; L(\theta | \mathbf{Y})$$

$\text{argmax}_\theta$ means "value of $\theta$ that maximizes..."

# Example: *Bernoulli(p)*

$Y_1, \ldots, Y_n \overset{iid}{\sim} \text{Bernoulli}(p)$

$$L(p \mid Y) = p^{\sum_i Y_i} (1-p)^{n - \sum_i Y_i}$$

maximize to estimate $p$:

① Take log to make life easier
   - log is monotone, increasing, so if $\hat{p}$ maximize
      $$L(p \mid Y) \iff \log L(p \mid Y)$$

$$\ell(p \mid Y) = \log L(p \mid Y) = \left( \sum_i Y_i \right) \log p + (n - \sum_i Y_i) \log(1-p)$$

② Differentiate wrt parameter of interest:

$$\frac{\partial}{\partial p} \ell(p \mid Y) = \frac{\sum_i Y_i}{p} - \frac{(n - \sum_i Y_i)}{1-p} \overset{set}{=} 0$$

$$\frac{\Sigma_i Y_i}{P} - \frac{(n - \Sigma_i Y_i)}{1-P} \overset{\text{set}}{=} 0$$

$$\frac{\Sigma_i Y_i}{P} = \frac{(n - \Sigma_i Y_i)}{1-P} \qquad \Rightarrow \qquad P = \frac{1}{n} \Sigma_i Y_i$$

(sample proportion!)

Check second derivative:

$$\frac{d^2}{dp^2} \ell(p|Y) \Big|_{P = \frac{1}{n} \Sigma_i Y_i} = -\frac{\Sigma_i Y_i}{P^2} - \frac{(n - \Sigma_i Y_i)}{1-P^2} \Big|_{P = \bar{Y}}$$

$$< 0$$

$$\Rightarrow \quad \hat{p} = \frac{1}{n} \Sigma_i Y_i = \bar{Y} \qquad \text{maximizes } \ell(p|Y)$$

Sidebar : $\sum_i (Y_i - \beta_0)^2$    is minimized by

$$\beta_0 = \bar{Y}$$

i.e.    $\sum_i (Y_i - \bar{Y})^2 \leq \sum_i (Y_i - \beta_0)^2$

$$\forall \beta_0$$