

## Logistic Regression

Recall: linear regression is used when response variable is quantitative, and ideally when error distribution is normal and linearity is a good assumption. However, in practice, other types of response arise. Two examples:

- ① response is binary/categorical/qualitative (e.g. presence/absence of a disease)
- ② response occurs as counts/integer (e.g. arrivals in a queue)

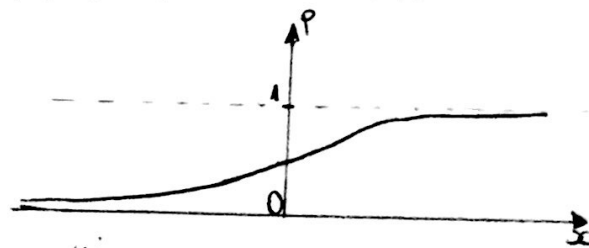
Today we're going to study models of type ①, which can be modeled by Logistic Regression, and response is binary/binomial.

- Simple Logistic Regression  $Y \in \{0, 1\} \Leftrightarrow$  binary/Bernoulli
- Multinomial Logistic Regression  $Y \in \{0, 1, 2, \dots\} \Leftrightarrow$  Binomial

Why Logistic Regression?

① When you must have  $0 \leq E[Y|X] \leq 1$

② s-shaped curve arises empirically



$$\text{logit (logistic) function} = \ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x \quad \text{or} \quad \ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_i \quad (\text{one covariate})$$

$$\text{or} \quad \ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \sum_j \beta_j x_{ij} = \mathbf{x}_i^T \boldsymbol{\beta} \quad (j \text{ covariates})$$

\* log odds = logit function =  $\ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x$

\* odds =  $\frac{p}{1-p} = e^{\beta_0 + \beta_1 x} = e^{\text{logit function}}$

\* odds ratio (A vs B) =  $\frac{e^{\beta_0 + \beta_1 A}}{e^{\beta_0 + \beta_1 B}} = e^{\beta_1(A-B)}$

Interpretation: for any value of  $x$ , increasing the covariate  $x$  by one unit increases the log odds by  $\beta_1$ .

$$\Rightarrow \frac{p}{1-p} = e^{\beta_0 + \beta_1 x} \Rightarrow p = e^{\beta_0 + \beta_1 x} - p e^{\beta_0 + \beta_1 x}$$

$$\Leftrightarrow p = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

where the last equality comes from:  $\frac{e^{\beta_0 + \beta_1 x}}{e^{\beta_0 + \beta_1 x} + 1} \cdot \frac{(1)}{(1/e^{\beta_0 + \beta_1 x} + 1)} = \frac{1}{1 + e^{-\beta_0 - \beta_1 x}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$

Model Estimation using MLE - Response is Binary  $\therefore y_i \in \{0, 1\}$   
(Simple Logistic Regression)

① For  $p$  ( $\hat{p}_{MLE}$ )  $y_i | x_i, x_i \sim \text{Bernoulli}(p)$

$$\text{Likelihood: } L(y_i; p) = \prod_{i=1}^n p^{y_i} (1-p)^{1-y_i}$$

$$\text{loglikelihood: } l(y_i; p) = \sum_{i=1}^n \{y_i \log p + (1-y_i) \log(1-p)\}$$

$\arg\max_p l(y_i; p) = ? \Rightarrow$  Take 1<sup>st</sup> derivatives wrt  $p$  and set to 0

$$\frac{\partial l(y_i; p)}{\partial p} = \sum_{i=1}^n \frac{y_i}{p} - \frac{1-y_i}{1-p} = 0 \Leftrightarrow \frac{\sum y_i}{p} = \frac{n - \sum y_i}{1-p}$$

$$\Leftrightarrow \sum y_i - p \sum y_i = pn - p \sum y_i$$

$$\boxed{\hat{p}_{MLE} = \frac{\sum y_i}{n} = \bar{y}}$$

② for  $\beta$  ( $\hat{\beta}_{MLE} \rightarrow$  no closed form  $\rightarrow$  solve numerically)

$$\text{loglikelihood } l(y_i; \beta_0, \beta) = \sum_{i=1}^n \left\{ y_i \log \left( \frac{1}{1 + e^{-(\beta_0 + x_i^T \beta)}} \right) + (1-y_i) \log \left( 1 - \frac{1}{1 + e^{-(\beta_0 + x_i^T \beta)}} \right) \right\}$$

$$\text{notice: } h(x_i; \beta_0, \beta) = \frac{1}{1 + e^{-(\beta_0 + x_i^T \beta)}}$$

$$\Rightarrow \log h(x_i; \beta_0, \beta) = \log 1 - \log(1 + e^{-(\beta_0 + x_i^T \beta)}) = -\log(1 + e^{-(\beta_0 + x_i^T \beta)})$$

$$\text{loglikelihood } l(y_i; \beta_0, \beta) = \sum_{i=1}^n \left\{ -y_i \log(1 + e^{-(\beta_0 + x_i^T \beta)}) + (1-y_i) \log \left( \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} \right) \right\}$$

Now, to apply Newton's Method (to get approx.  $\hat{\beta}$ ) we need to derive gradients and Hessian.

\* Getting Gradients:

$$\begin{aligned}
 \nabla l(\beta_0) &= \sum_{i=1}^n \left\{ y_i \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} + (1 - y_i) \frac{1}{\frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}}} \frac{(-e^{-(\beta_0 + x_i^T \beta)} (1 + e^{-(\beta_0 + x_i^T \beta)}) - (-e^{-(\beta_0 + x_i^T \beta)}) e^{-(\beta_0 + x_i^T \beta)})}{(1 + e^{-(\beta_0 + x_i^T \beta)})^2} \right\} \\
 &= \sum_{i=1}^n \left\{ y_i \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} + (1 - y_i) \frac{1}{e^{-(\beta_0 + x_i^T \beta)}} \left( -e^{-(\beta_0 + x_i^T \beta)} + \frac{(e^{-(\beta_0 + x_i^T \beta)})^2}{1 + e^{-(\beta_0 + x_i^T \beta)}} \right) \right\} \\
 &= \sum_{i=1}^n \left\{ y_i \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} + (1 - y_i) \left( -1 + \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} \right) \right\} \\
 &= \sum_{i=1}^n \left\{ -1 + \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} + y_i \right\} = \sum_{i=1}^n \left\{ y_i - \frac{1}{(1 + e^{-(\beta_0 + x_i^T \beta)})} \right\} = \sum_{i=1}^n \{ y_i - h(x_i, \beta_0, \beta) \}
 \end{aligned}$$

$$\begin{aligned}
 \nabla l(\beta) &= \sum_{i=1}^n \left\{ y_i \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} x_i + (1 - y_i) \frac{1}{\frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}}} \frac{(-e^{-(\beta_0 + x_i^T \beta)} (1 + e^{-(\beta_0 + x_i^T \beta)}) - (-e^{-(\beta_0 + x_i^T \beta)}) e^{-(\beta_0 + x_i^T \beta)})}{(1 + e^{-(\beta_0 + x_i^T \beta)})^2} x_i \right\} \\
 &= \sum_{i=1}^n \left\{ y_i \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} x_i + (1 - y_i) \left( -x_i + \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} x_i \right) \right\} \\
 &= \sum_{i=1}^n \left\{ -x_i + \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} x_i + y_i x_i \right\} = \sum_{i=1}^n \left\{ x_i \left( \frac{e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} - 1 + y_i \right) \right\} \\
 &= \sum_{i=1}^n \left\{ x_i \left( -\frac{1}{1 + e^{-(\beta_0 + x_i^T \beta)}} + y_i \right) \right\} = \sum_{i=1}^n \left\{ x_i (y_i - h(x_i, \beta_0, \beta)) \right\}
 \end{aligned}$$

Hence,  $\nabla l(\beta_0, \beta) = \sum_{i=1}^n \begin{bmatrix} 1 \\ x_i \end{bmatrix} (y_i - h(x_i, \beta_0, \beta))$

$$\text{where } h(x_i, \beta_0, \beta) = \frac{1}{1 + e^{-(\beta_0 + x_i^T \beta)}} = \frac{e^{\beta_0 + x_i^T \beta}}{1 + e^{\beta_0 + x_i^T \beta}} = p$$

## \* Getting Hessian Matrix

$$\begin{aligned} \nabla^2 l(\beta_0) &= \frac{\partial^2 l(\beta_0, \beta)}{\partial \beta^2} = \sum_{i=1}^n \left\{ - \left( \frac{e^{-(\beta_0 + x_i^T \beta)}}{(1 + e^{-(\beta_0 + x_i^T \beta)})^2} \right) \right\} = \sum_{i=1}^n \left\{ \overbrace{\frac{1}{1 + e^{-(\beta_0 + x_i^T \beta)}}}^h \overbrace{\left( \frac{-e^{-(\beta_0 + x_i^T \beta)}}{1 + e^{-(\beta_0 + x_i^T \beta)}} \right)}^{(h-1)} \right\} \\ &= \sum_{i=1}^n \left\{ h(x_i; \beta_0, \beta) (h(x_i; \beta_0, \beta) - 1) \right\} \end{aligned}$$

$$\nabla^2 l(\beta) = \frac{\partial^2 l(\beta_0, \beta)}{\partial \beta \partial \beta^T} = \sum_{i=1}^n \left\{ -x_i \left( \frac{x_i^T (e^{-(\beta_0 + x_i^T \beta)})}{(1 + e^{-(\beta_0 + x_i^T \beta)})^2} \right) \right\} = \sum_{i=1}^n \left\{ x_i x_i^T h(x_i; \beta_0, \beta) (h(x_i; \beta_0, \beta) - 1) \right\}$$

$$\nabla^2 l(\beta_0, \beta) = \frac{\partial^2 l(\beta_0, \beta)}{\partial \beta_0 \partial \beta^T} = \sum_{i=1}^n \left\{ - \left( \frac{x_i^T e^{-(\beta_0 + x_i^T \beta)}}{(1 + e^{-(\beta_0 + x_i^T \beta)})^2} \right) \right\} = \sum_{i=1}^n \left\{ x_i^T h(x_i; \beta_0, \beta) (h(x_i; \beta_0, \beta) - 1) \right\}$$

$$\text{and } \frac{\partial^2 l(\beta_0, \beta)}{\partial \beta \partial \beta_0} = \sum_{i=1}^n \left\{ x_i h(x_i; \beta_0, \beta) (h(x_i; \beta_0, \beta) - 1) \right\}$$

Hence, the Hessian:

$$\nabla^2(\beta_0, \beta) = \sum_{i=1}^n h(x_i; \beta_0, \beta) (h(x_i; \beta_0, \beta) - 1) \begin{bmatrix} 1_{(K+1)} & x_{i(K+1)}^T \\ x_{i(K+1)} & x_i x_i^T \end{bmatrix}_{(K+1)(K+1)}$$

$$\text{or equivalently } \nabla^2(\beta_0, \beta) = \sum_{i=1}^n \begin{bmatrix} 1 \\ x_i \end{bmatrix} \begin{bmatrix} 1 \\ x_i^T \end{bmatrix}^T h(x_i; \beta_0, \beta) (h(x_i; \beta_0, \beta) - 1)$$

## \* Applying Newton's Method

$$\hat{\beta}^{\text{new}} = \hat{\beta}^{\text{old}} - [\text{Hessian}]^{-1} \cdot \text{gradient}$$

$$\begin{bmatrix} \beta_0 \\ \beta \end{bmatrix}^{(t+1)} = \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix}^{(t)} - \underset{\substack{\downarrow \\ \text{step}}}{\eta} \left[ \nabla^2 l(\beta_0^{(t)}, \beta^{(t)}) \right]^{-1} \cdot \nabla l(\beta_0^{(t)}, \beta^{(t)})$$

define  $\begin{bmatrix} \beta_0 \\ \beta \end{bmatrix}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  and choose a step (e.g. 0.01, 1, etc.) to get convergence

→ Another common notation:

if we write the loglikelihood as:

$$l(y|\beta) = -y^T \log(1 + e^{-X^T \beta}) + (1^T - y^T) \log\left(\frac{e^{-X^T \beta}}{1 + e^{-X^T \beta}}\right)$$

Then  $\nabla l(y|\beta) = X^T (y - p)$  where  $p$  is the vector  $\begin{pmatrix} h(x_1|\beta) \\ h(x_2|\beta) \\ \vdots \\ h(x_n|\beta) \end{pmatrix}$

And  $\nabla^2 l(y|\beta) = -X^T W X$  where  $W$  is the matrix of weights :  $W = [P(x_i)/(1 - P(x_i))]_{n \times n}$

Therefore, we rewrite Newton's Method as:

$$\boxed{\hat{\beta}^{\text{new}} = \hat{\beta}^{\text{old}} + (X^T W X)^{-1} X^T (y - p)}$$

\* Standard error  $se(\hat{\beta}) = (X^T W X)^{-1}$