

## Logistic regression

- In many cases, we may have a binary dependent variable: That is,

$$Y_i = \begin{cases} 1 \\ 0 \end{cases}, \text{ e.g., approval vs denial, employed vs unemployed, pass vs fail...}$$

- Mortgage application : outcome is either approval or denial. This decision is made based on several factors such as income, job, debt, ... Suppose that we want to examine the relation between this mortgage approval decision and loan payment/income (P/I ratio). Let

$$Y_i = \begin{cases} 1, & \text{if denied} \\ 0, & \text{if approved} \end{cases}, \quad X_i = \text{P/I ratio},$$

- We are interested in the probability of approval given my P/I ratio,  $P(\text{denial} \mid \text{P/I ratio})$ .
- Suppose that this decision is based on latent variable  $Y^*$  (e.g., risk score) which is determined by  $X_i = \text{P/I ratio}$  and other unobserved factors, such that

$$Y^* = \beta_0 + \beta_1 X_i + \varepsilon_i$$

and

$$Y_i = \begin{cases} 1, & \text{if } Y^* \geq 0 \\ 0, & \text{if } Y^* < 0 \end{cases}.$$

- Here we assume that  $\varepsilon_i \mid X_i$  follows the logistic distribution.

$$L(\varepsilon \mid X) = \frac{1}{1 + \exp(-\varepsilon)},$$

where  $L(\cdot \mid X)$  is the conditional cdf of the logistic random variable.

- Note that

$$\begin{aligned} P(Y_i = 1 \mid X_i) &= P(Y_i^* \geq 0 \mid X_i) \\ &= P(\beta_0 + \beta_1 X_i + \varepsilon_i \geq 0 \mid X_i) \\ &= P(\varepsilon_i \geq -\beta_0 - \beta_1 X_i \mid X_i) \\ &= P(\varepsilon_i \leq \beta_0 + \beta_1 X_i \mid X_i) \\ &= L(\beta_0 + \beta_1 X_i \mid X_i) \\ &= \frac{1}{1 + \exp(-\beta_0 - \beta_1 X_i)} \end{aligned}$$

- We can estimate the coefficients using "maximum likelihood estimation" (MLE), and the estimated probability of approval given the P/I ratio is

$$\hat{P}(Y_i = 1 \mid X_i) = \frac{1}{1 + \exp(-\hat{\beta}_0 - \hat{\beta}_1 X_i)}.$$

- Example: Suppose that  $\hat{\beta}_0 = -2, \hat{\beta}_1 = 3, X = 0.4$ ,

$$P(Y = 1 \mid X = 0.4) = (-2 + 3 \times 0.4) = L(-0.8).$$

- Properties

1. It has S-shape,  $0 \leq P(Y = 1 \mid X) \leq 1$  and  $P(Y = 1 \mid X)$  increases in  $X$  when  $\beta_1 > 0$ .
2.  $\beta_1$  does not represent the change in  $P(Y = 1 \mid X)$  by increasing 1 unit increase in  $X$ . Note that the cdf has an S shape. The marginal probability is the function of  $X$ .

- Predicted probability

$$\hat{P}(Y = 1|X) = L(\hat{\beta}_0 + \hat{\beta}_1 X).$$

- Estimation: OLS does not work! We cannot use the `lm()` function!