

## ⇒ Generalized Linear Models for Binary Data:

Many categorical response variables have only two categories. For example, whether you take public transport today (Yes, No) - Denote a binary response variable by  $Y$  and its two possible outcomes by 1 (success) and 0 (failure).

The distribution of  $Y$  is specified by probabilities  $P(Y=1) = \pi$  of success and  $P(Y=0) = (1-\pi)$  of failure. For each binary observation its mean is  $E(Y) = \pi$ .

The following section introduces GLMs for binary responses. Although GLMs can have multiple explanatory variables, for simplicity we introduce them using single  $x$ . The value of  $\pi$  can vary as the value of  $x$  changes, and we replace  $\pi$  by  $\pi(x)$  when we want to describe its dependence on that value.

### ① Linear Probability Model:

In ordinary regression,  $\mu = E(Y)$  is a linear function of  $x$ . For a binary response, an analogous model is

$$\pi(x) = \alpha + \beta x$$

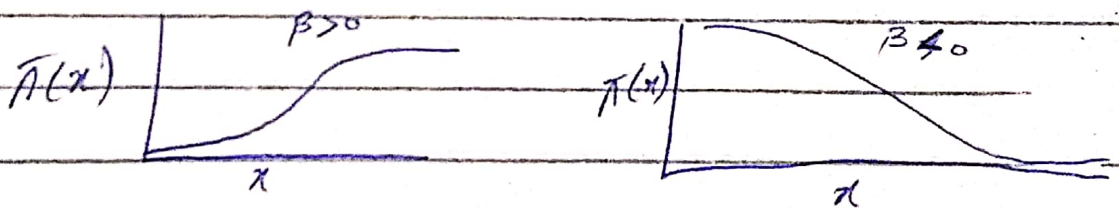
This is called a linear probability model because the probability of success changes linearly in  $x$ . The parameter  $\beta$  represents the change in the probability per unit change in  $x$ . This model is a GLM with binomial random component and identity link function.



This model is simple, but unfortunately it has a structural defect. Probability fall between 0 and 1, whereas linear functions take value over the entire real line. This model predict,  $\pi(x) < 0$  and  $\pi(x) > 1$  for sufficiently large or small  $x$  values. The model can fit adequate over a restricted range of  $x$  values.

(a) Logistic Regression Model:

Relationship between  $\pi(x)$  and  $x$  are usually non linear rather than linear. (Prob. of buying a new car increase when <sup>increase in annual incomes</sup>  $x$  increase. In practice,  $\pi(x)$  often either increases continuously or decreases continuously as  $x$  increase. The S-shaped curves are often realistic shapes for the relationship



The most important mathematical function with this shape has formula:

$$\pi(x) = \frac{\exp(\alpha + \beta x)}{1 + \exp(\alpha + \beta x)} = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

This is called the logistic regression function.



The corresponding logistic regression model form is

$$\log\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x$$

$$\frac{\pi(x)}{1-\pi(x)} = \frac{e^{\alpha+\beta x}}{1+e^{\alpha+\beta x}}$$
$$= \frac{e^{\alpha+\beta x}}{1+e^{\alpha+\beta x} - e^{\alpha+\beta x}}$$
$$= \frac{e^{\alpha+\beta x}}{1+e^{\alpha+\beta x} - e^{\alpha+\beta x}}$$

$$\frac{\pi(x)}{1-\pi(x)} = e^{\alpha+\beta x}$$

$$\log\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x \log(e)$$

$$\log\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x \Rightarrow \text{logit}(\pi(x)) = \alpha + \beta x$$

$\searrow$  logit function of  $\pi(x)$

This is a GLM with binomial random component and logit link function. The logistic regression models are often called logit models.

Whereas  $\pi$  is restricted to the 0-1 range, the logit can be any real number. The real numbers are also the potential range for linear predictors (such as  $\alpha + \beta x$ ) that form the systematic component of a GLM, so this model does not have the structural problem that the linear



probability model has.

The parameter  $\beta$  determines the rate of increase or decrease of the curve.

when  $\beta > 0$ ,  $\pi(x)$  increase as  $x$  increase.

when  $\beta < 0$ ,  $\pi(x)$  decrease as  $x$  increase.

The magnitude of  $\beta$  determines how fast the curve increases or decreases.

when  $\beta = 0$ , the curve flattens to a horizontal straight line.

### ③ Probit Regression Model:

Another model that has S-shaped curves is called probit model. The link function for the model, called the probit link, transform probabilities to  $z$ -scores from the standard normal distribution. The probit model has expression

$$\text{Probit}[\pi(x)] = \alpha + \beta x$$

The probit link function applied to  $\pi(x)$  gives the standard normal  $z$ -score at which the left-tail probability equals  $\pi(x)$ . For instance  $\text{probit}(0.05) = -1.645$  because 5% of the standard normal distribution falls below  $-1.645$ .

like wise

$$\text{Probit}(0.50) = 0$$

$$\text{probit}(0.95) = 1.645$$

$$\text{probit}(0.975) = 1.96$$