

10 November 2020

Binary response

Data set

1. Ungrouped data: Response - y_i for $i = 1, \dots, n$
Assume $y_i \sim \text{Bernoulli}(P_i)$
2. Grouped data: Response - (n_i, r_i) for $i = 1, \dots, s$
Assume $r_i \sim \text{Binomial}(n_i, P_i)$

Problems if linear model is used

Consider the model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + e_i \quad \begin{cases} i = 1, 2, \dots, n \\ y_i = 0, 1 \end{cases}$$

If we assume the usual $E(e_i) = 0$, we have $E(y_i) = P_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$ and $\text{Var}(e_i) = P_i Q_i$

The other two problems on the above linear regression, i.e., fitting y_i on x_1, \dots, x_p

1. Estimate of y_i is imprecise for ungrouped data.
2. $\hat{P}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \dots + \hat{\beta}_p x_{ip}$ may not be within 0 and 1.

Link function

The link function provides the relationship between the linear predictor (linear combination of unknown parameter β) and the mean of the distribution function, i.e., $E(y_i) = \mu_i$. It is defined as $g(\mu)$.

For linear model,

$$\mu_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

Then, $g(\mu_i) = \mu_i$ — Identity function

For binomial data,

1.

$$P_i = \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}$$
$$\log \left[\frac{\mu_i}{1 - \mu_i} \right] = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

Then, $g(\mu_i) = \log \left[\frac{\mu_i}{1 - \mu_i} \right]$ — Logit function

2.

$$P_i = \Phi(\mathcal{X}_i^T \beta)$$
$$\Phi^{-1}(\mu_i) = \mathcal{X}_i^T \beta$$

Then, $g(\mu_i) = \Phi^{-1}(\mu_i)$ — Probit function

Weighted least squares

By minimizing

$$SS_{\text{Res}, \mathbf{V}} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}),$$

Then,

- $\tilde{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{Y}$
- $\text{Var}(\tilde{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1}$

As $\mathbf{V} = \text{Var}(\boldsymbol{\varepsilon}) = \text{diag}[\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2]$, the weighted least squares estimator of $\boldsymbol{\beta}$ can be obtained by minimizing

$$SS_{\text{Res}(\text{weighted})} = \sum_{i=1}^n w_i (y_i - \hat{y}_i)^2$$

where $w_i = 1/\sigma_i^2$.

Weighted least squares for grouped data with large n_i

We fit the model as follows:

$$\log \left[\frac{\hat{P}_i}{1 - \hat{P}_i} \right] = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + e_i \quad i = 1, 2, \dots, s$$

with $\hat{w}_i = \frac{1}{\hat{\sigma}_i^2} \approx n_i \hat{P}_i (1 - \hat{P}_i)$, where $\hat{P}_i = r_i/n_i$.