

# Notes on Binary Outcome Model

Max Leung

February 20, 2024

## 1 General Binary Outcome Model

$$y_i = \begin{cases} 1 & \text{with probability } p_i \\ 0 & \text{with probability } 1 - p_i \end{cases}$$

Model  $p_i$  as

$$p_i := Pr(y_i = 1 | \mathbf{x}_i) = F(\mathbf{x}_i' \boldsymbol{\beta}) = F_i$$

### 1.1 Marginal Effect

$$\begin{aligned} \frac{\partial Pr(y_i = 1 | \mathbf{x}_i)}{\partial \mathbf{x}_i} &= \frac{\partial F(\mathbf{x}_i' \boldsymbol{\beta})}{\partial \mathbf{x}_i} \\ &= \frac{\partial F(\mathbf{x}_i' \boldsymbol{\beta})}{\partial \mathbf{x}_i' \boldsymbol{\beta}} \frac{\partial \mathbf{x}_i' \boldsymbol{\beta}}{\partial \mathbf{x}_i} \\ &= F'(\mathbf{x}_i' \boldsymbol{\beta}) \boldsymbol{\beta} \end{aligned}$$

If  $F(\cdot)$  is cdf,  $F'(\cdot) > 0$ . So,  $sign(\boldsymbol{\beta})$  decide the sign

Average Marginal Effect is

$$AME := N^{-1} \sum_{i=1}^N F'(\mathbf{x}_i' \hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}}$$

or

$$AME := F'(N^{-1} \sum_{i=1}^N \mathbf{x}_i' \hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}} = F'(\bar{\mathbf{x}}' \hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}}$$

If  $\bar{\mathbf{x}}$  includes  $\overline{\ln(z)}$ , it can be replaced by  $\ln(\bar{z})$ . If  $\mathbf{x}_i$  includes  $z_i$  and  $z_i^2$ ,  $AME_i$  for  $z_i$  is  $F'(\mathbf{x}_i' \hat{\boldsymbol{\beta}})(\hat{\beta}_z + 2\hat{\beta}_z z_i)$ . Therefore,  $AME$  is  $F'(\bar{\mathbf{x}}' \hat{\boldsymbol{\beta}})(\hat{\beta}_z + 2\hat{\beta}_z \bar{z})$  where  $\bar{z}^2$  in  $\bar{\mathbf{x}}$  can be replaced by  $\bar{z}^2$ . If  $\mathbf{x}_i$  includes an indicator  $z_i$ ,  $AME$  is  $F(\bar{\mathbf{x}}'_{-z} \hat{\boldsymbol{\beta}}_{-z} + 1 \cdot \hat{\beta}_z) - F(\bar{\mathbf{x}}'_{-z} \hat{\boldsymbol{\beta}}_{-z})$ .

## 1.2 Maximum Likelihood Estimation

### 1.2.1 Probability Mass Function

As  $y_i$  is binary, it must be Bernoulli distributed. The probability mass function (pmf) of such random variable is

$$\begin{aligned} f(y_i | \mathbf{x}_i) &= Pr(y_i = 1 | \mathbf{x}_i)^{y_i} Pr(y_i = 0 | \mathbf{x}_i)^{1-y_i} \\ &= Pr(y_i = 1 | \mathbf{x}_i)^{y_i} (1 - Pr(y_i = 1 | \mathbf{x}_i))^{1-y_i} \\ &= p_i^{y_i} (1 - p_i)^{1-y_i} \\ &= F(\mathbf{x}_i' \boldsymbol{\beta})^{y_i} (1 - F(\mathbf{x}_i' \boldsymbol{\beta}))^{1-y_i} \end{aligned}$$

### 1.2.2 Log Likelihood Function

$$\begin{aligned}
\ln[L_N(\boldsymbol{\beta})] &= \ln\left[\prod_{i=1}^N f(y_i|\mathbf{x}_i)\right] && \text{assume independence} \\
&= \ln\left[\prod_{i=1}^N F(\mathbf{x}'_i\boldsymbol{\beta})^{y_i}(1 - F(\mathbf{x}'_i\boldsymbol{\beta}))^{1-y_i}\right] \\
&= \sum_{i=1}^N \ln[F(\mathbf{x}'_i\boldsymbol{\beta})^{y_i}(1 - F(\mathbf{x}'_i\boldsymbol{\beta}))^{1-y_i}] \\
&= \sum_{i=1}^N \{y_i \ln[F(\mathbf{x}'_i\boldsymbol{\beta})] + (1 - y_i) \ln[1 - F(\mathbf{x}'_i\boldsymbol{\beta})]\}
\end{aligned}$$

### 1.2.3 Gradient Vector

$$\begin{aligned}
\frac{\partial \ln[L_N(\boldsymbol{\beta})]}{\partial \boldsymbol{\beta}} &= \frac{\partial \sum_{i=1}^N \{y_i \ln[F(\mathbf{x}'_i\boldsymbol{\beta})] + (1 - y_i) \ln[1 - F(\mathbf{x}'_i\boldsymbol{\beta})]\}}{\partial \boldsymbol{\beta}} \\
&= \sum_{i=1}^N \frac{\partial \{y_i \ln[F(\mathbf{x}'_i\boldsymbol{\beta})] + (1 - y_i) \ln[1 - F(\mathbf{x}'_i\boldsymbol{\beta})]\}}{\partial \boldsymbol{\beta}} \\
&= \sum_{i=1}^N \left\{ \frac{\partial y_i \ln[F(\mathbf{x}'_i\boldsymbol{\beta})]}{\partial \boldsymbol{\beta}} + \frac{\partial (1 - y_i) \ln[1 - F(\mathbf{x}'_i\boldsymbol{\beta})]}{\partial \boldsymbol{\beta}} \right\} \\
&= \sum_{i=1}^N \left\{ y_i \frac{1}{F(\mathbf{x}'_i\boldsymbol{\beta})} F'(\mathbf{x}'_i\boldsymbol{\beta}) \mathbf{x}_i + (1 - y_i) \frac{1}{1 - F(\mathbf{x}'_i\boldsymbol{\beta})} (-1) F'(\mathbf{x}'_i\boldsymbol{\beta}) \mathbf{x}_i \right\} \\
&= \sum_{i=1}^N \left\{ \frac{y_i}{F(\mathbf{x}'_i\boldsymbol{\beta})} F'(\mathbf{x}'_i\boldsymbol{\beta}) \mathbf{x}_i - \frac{1 - y_i}{1 - F(\mathbf{x}'_i\boldsymbol{\beta})} F'(\mathbf{x}'_i\boldsymbol{\beta}) \mathbf{x}_i \right\} \\
&= \sum_{i=1}^N \left\{ \frac{y_i}{F_i} F'_i \mathbf{x}_i - \frac{1 - y_i}{1 - F_i} F'_i \mathbf{x}_i \right\} \\
&= \sum_{i=1}^N \frac{y_i F'_i \mathbf{x}_i (1 - F_i) - (1 - y_i) F'_i \mathbf{x}_i F_i}{F_i (1 - F_i)} \\
&= \sum_{i=1}^N \frac{(y_i F'_i \mathbf{x}_i - y_i F'_i \mathbf{x}_i F_i) - (F'_i \mathbf{x}_i F_i - y_i F'_i \mathbf{x}_i F_i)}{F_i (1 - F_i)} \\
&= \sum_{i=1}^N \frac{y_i F'_i \mathbf{x}_i - F'_i \mathbf{x}_i F_i}{F_i (1 - F_i)} \\
&= \sum_{i=1}^N \frac{y_i - F_i}{F_i (1 - F_i)} F'_i \mathbf{x}_i
\end{aligned}$$

### 1.2.4 First Order Condition

$$\begin{aligned}
&\sum_{i=1}^N \frac{y_i - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle})}{F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle}) (1 - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle}))} F'(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle}) \mathbf{x}_i = \mathbf{0} \\
&\sum_{i=1}^N \underbrace{\left\{ \frac{F'(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle})}{F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle}) (1 - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle}))} \right\}}_{\hat{w}_i} [y_i - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{mle})] \mathbf{x}_i = \mathbf{0}
\end{aligned}$$

There is no closed form solution for  $\hat{\boldsymbol{\beta}}_{mle}$ . We usually solve it with Newton method.  $\ln[L_N(\boldsymbol{\beta})]$  is globally concave because its Hessian matrix is negative definite. As a result, the convergence of the algorithm is fast.

### 1.2.5 Consistency of Maximum Likelihood Estimator

Correct specification of the likelihood function is a necessary condition for consistent MLE  $\hat{\beta}_{mle} \rightarrow_p \beta_0$  as  $N \rightarrow \infty$ . As  $y_i$  is binary, it must be Bernoulli distributed. However,  $p_i := Pr(y_i = 1|x_i) = F(x'_i\beta)$  can be incorrectly specified.

### 1.2.6 Asymptotic Distribution of Maximum Likelihood Estimator

Under regularity conditions, Information Matrix equality holds and MLE is asymptotically normally distributed. Its asymptotic variance is the inverse of Fisher's Information matrix.

$$\begin{aligned}
\sqrt{N}(\hat{\beta}_{mle} - \beta_0) &\rightarrow_d N(\mathbf{0}, [\mathbf{I}(\beta_0)]^{-1}) \\
&= N(\mathbf{0}, [\sum_{i=1}^N \frac{1}{F(x'_i\beta_0)(1-F(x'_i\beta_0))} F'(x'_i\beta_0)^2 x_i x'_i]^{-1}) \\
&= N(\mathbf{0}, [\sum_{i=1}^N \frac{1}{Var(y_i|x_i)} F'(x'_i\beta_0)^2 x_i x'_i]^{-1}) \\
\mathbf{I}(\beta_0) &:= -\mathbb{E}[\frac{\partial^2 \ln L_N(\beta)}{\partial \beta \partial \beta'} |_{\beta_0} | x_i] = \mathbb{E}[\frac{\partial \ln L_N(\beta)}{\partial \beta} \cdot \frac{\partial \ln L_N(\beta)}{\partial \beta'} |_{\beta_0} | x_i] \\
&= -\mathbb{E}[\frac{\partial \sum_{i=1}^N \{ \frac{y_i}{F_i} F'_i x_i - \frac{1-y_i}{1-F_i} F'_i x_i \}}{\partial \beta'} |_{\beta_0} | x_i] \\
&= -\sum_{i=1}^N \mathbb{E}[\frac{\partial (\frac{y_i}{F_i} - \frac{1-y_i}{1-F_i}) F'_i x_i}{\partial \beta'} |_{\beta_0} | x_i] \\
&= -\sum_{i=1}^N \mathbb{E}[\frac{\partial (\frac{y_i}{F_i} - \frac{1-y_i}{1-F_i})}{\partial \beta'} F'_i x_i |_{\beta_0} + (\frac{y_i}{F_i} - \frac{1-y_i}{1-F_i}) \frac{\partial F'_i x_i}{\partial \beta'} |_{\beta_0} | x_i] \\
&= -\sum_{i=1}^N \mathbb{E}[(\frac{\partial \frac{y_i}{F_i}}{\partial \beta'} - \frac{\partial \frac{1-y_i}{1-F_i}}{\partial \beta'}) F'_i x_i + (\frac{y_i}{F_i} - \frac{1-y_i}{1-F_i}) \frac{\partial F'_i x_i}{\partial \beta'} | x_i] |_{\beta_0} \\
&= -\sum_{i=1}^N \mathbb{E}[(-y_i F_i^{-2} F'_i x'_i - (1-y_i)(1-F_i)^{-2} F'_i x'_i) F'_i x_i + (\frac{y_i}{F_i} - \frac{1-y_i}{1-F_i}) F''_i x_i x'_i | x_i] |_{\beta_0} \\
&= \sum_{i=1}^N \{ \mathbb{E}[(y_i F_i^{-2} F'_i x'_i + (1-y_i)(1-F_i)^{-2} F'_i x'_i) F'_i x_i | x_i] |_{\beta_0} + \mathbb{E}[(\frac{y_i}{F_i} - \frac{1-y_i}{1-F_i}) F''_i x_i x'_i | x_i] |_{\beta_0} \} \\
&= \sum_{i=1}^N \{ \mathbb{E}[(y_i F_i^{-2} + (1-y_i)(1-F_i)^{-2}) F_i'^2 x_i x'_i | x_i] |_{\beta_0} + \mathbb{E}[(\frac{y_i}{F_i} - \frac{1-y_i}{1-F_i}) | x_i] F''_i x_i x'_i |_{\beta_0} \} \\
&= \sum_{i=1}^N \{ (\mathbb{E}[y_i | x_i] F_i^{-2} + (1 - \mathbb{E}[y_i | x_i])(1-F_i)^{-2}) F_i'^2 x_i x'_i |_{\beta_0} + (\frac{\mathbb{E}[y_i | x_i]}{F_i} - \frac{1 - \mathbb{E}[y_i | x_i]}{1-F_i}) F''_i x_i x'_i |_{\beta_0} \} \\
&= \sum_{i=1}^N \{ (F(x'_i\beta_0) F(x'_i\beta_0)^{-2} + (1-F(x'_i\beta_0))(1-F(x'_i\beta_0))^{-2}) F'(x'_i\beta_0)^2 x_i x'_i + \\
&\quad (\frac{F(x'_i\beta_0)}{F(x'_i\beta_0)} - \frac{1-F(x'_i\beta_0)}{1-F(x'_i\beta_0)}) F''(x'_i\beta_0) x_i x'_i \} \\
&= \sum_{i=1}^N (\frac{1}{F(x'_i\beta_0)} + \frac{1}{1-F(x'_i\beta_0)}) F'(x'_i\beta_0)^2 x_i x'_i \\
&= \sum_{i=1}^N \frac{(1-F(x'_i\beta_0)) + F(x'_i\beta_0)}{F(x'_i\beta_0)(1-F(x'_i\beta_0))} F'(x'_i\beta_0)^2 x_i x'_i \\
&= \sum_{i=1}^N \frac{1}{F(x'_i\beta_0)(1-F(x'_i\beta_0))} F'(x'_i\beta_0)^2 x_i x'_i
\end{aligned}$$

As  $\mathbb{E}(y_i|x_i) = 1 \cdot Pr(y_i = 1|x_i) + 0 \cdot Pr(y_i = 0|x_i) = Pr(y_i = 1|x_i) = F(x'_i\beta_0)$   
 $Var(y_i|x_i) = \mathbb{E}[(y_i - \mathbb{E}(y_i|x_i))^2|x_i] = (1 - \mathbb{E}(y_i|x_i))^2 Pr(y_i = 1|x_i) + (0 - \mathbb{E}(y_i|x_i))^2 Pr(y_i = 0|x_i) = (1 - Pr(y_i = 1|x_i))^2 Pr(y_i = 1|x_i) + Pr(y_i = 1|x_i)^2 (1 - Pr(y_i = 1|x_i)) = (1 - Pr(y_i = 1|x_i)) Pr(y_i = 1|x_i) (1 - Pr(y_i = 1|x_i) + Pr(y_i = 1|x_i)) = (1 - Pr(y_i = 1|x_i)) Pr(y_i = 1|x_i) = F(x'_i\beta_0)(1-F(x'_i\beta_0))$

$$1|\mathbf{x}_i)) = Pr(y_i = 1|\mathbf{x}_i)(1 - Pr(y_i = 1|\mathbf{x}_i)) = F(\mathbf{x}'_i\boldsymbol{\beta}_0)(1 - F(\mathbf{x}'_i\boldsymbol{\beta}_0))$$

For binary outcome model, even not all the regularity conditions for Information Matrix equality hold, the equality still holds algebraically. i.e.  $\mathbf{A} = -\mathbf{B}$ . It can be seen

$$\begin{aligned} \mathbf{B} &= \mathbb{E}\left[\frac{\partial \ln L_N(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \cdot \frac{\partial \ln L_N(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}'} \middle| \boldsymbol{\beta}_0 | \mathbf{x}_i\right] \\ &= \mathbb{E}\left(\sum_{i=1}^N \frac{y_i - F_i}{F_i(1 - F_i)} F'_i \mathbf{x}_i \cdot \sum_{i=1}^N \frac{y_i - F_i}{F_i(1 - F_i)} F'_i \mathbf{x}'_i \middle| \boldsymbol{\beta}_0 | \mathbf{x}_i\right) \\ &= \mathbb{E}\left(\sum_{i=1}^N \frac{y_i - F_i}{F_i(1 - F_i)} \frac{y_i - F_i}{F_i(1 - F_i)} F_i'^2 \mathbf{x}_i \mathbf{x}'_i \middle| \boldsymbol{\beta}_0 | \mathbf{x}_i\right) \\ &= \sum_{i=1}^N \frac{\mathbb{E}[(y_i - F(\mathbf{x}'_i\boldsymbol{\beta}_0))^2 | \mathbf{x}_i]}{F(\mathbf{x}'_i\boldsymbol{\beta}_0)(1 - F(\mathbf{x}'_i\boldsymbol{\beta}_0))} \frac{1}{F(\mathbf{x}'_i\boldsymbol{\beta}_0)(1 - F(\mathbf{x}'_i\boldsymbol{\beta}_0))} F'(\mathbf{x}'_i\boldsymbol{\beta}_0)^2 \mathbf{x}_i \mathbf{x}'_i \\ &= \sum_{i=1}^N \frac{1}{F(\mathbf{x}'_i\boldsymbol{\beta}_0)(1 - F(\mathbf{x}'_i\boldsymbol{\beta}_0))} F'(\mathbf{x}'_i\boldsymbol{\beta}_0)^2 \mathbf{x}_i \mathbf{x}'_i \\ &= -\mathbb{E}\left[\frac{\partial^2 \ln L_N(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \middle| \boldsymbol{\beta}_0 | \mathbf{x}_i\right] \\ &:= \mathbf{I}(\boldsymbol{\beta}_0) = -\mathbf{A} \end{aligned}$$

Without satisfying all the regularity conditions, the asymptotic distribution of MLE is  $\sqrt{N}(\hat{\boldsymbol{\beta}}_{mle} - \boldsymbol{\beta}_0) \rightarrow_d N(\mathbf{0}, \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1}) = N(\mathbf{0}, -\mathbf{A}^{-1} \mathbf{A} \mathbf{A}^{-1}) = N(\mathbf{0}, -\mathbf{A}^{-1}) = N(\mathbf{0}, \mathbf{I}(\boldsymbol{\beta}_0)^{-1})$ . Thus, the asymptotic variance is still the conventional one i.e., the inverse of the Information Matrix. Therefore, we do not have to use the sandwich standard error for binary outcome model.

### 1.3 Iteratively Weighted Non-linear Least Squares

#### 1.3.1 Loss Function

$$\sum_{i=1}^N w_i (y_i - F(\mathbf{x}'_i \boldsymbol{\beta}))^2$$

#### 1.3.2 First Order Condition

$$\begin{aligned} \frac{\partial \sum_{i=1}^N w_i (y_i - F(\mathbf{x}'_i \boldsymbol{\beta}))^2}{\partial \boldsymbol{\beta}} \bigg|_{\hat{\boldsymbol{\beta}}_{wnls}} &= \mathbf{0} \\ \sum_{i=1}^N w_i \frac{\partial (y_i - F(\mathbf{x}'_i \boldsymbol{\beta}))^2}{\partial \boldsymbol{\beta}} \bigg|_{\hat{\boldsymbol{\beta}}_{wnls}} &= \mathbf{0} \\ -2 \sum_{i=1}^N w_i (y_i - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{wnls})) \frac{\partial F(\mathbf{x}'_i \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \bigg|_{\hat{\boldsymbol{\beta}}_{wnls}} &= \mathbf{0} \\ \sum_{i=1}^N w_i (y_i - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{wnls})) F'(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{wnls}) \mathbf{x}_i &= \mathbf{0} \\ \sum_{i=1}^N w_i F'(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{wnls}) [y_i - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}}_{wnls})] \mathbf{x}_i &= \mathbf{0} \end{aligned}$$

If  $w_i$  is the inverse of conditional variance of  $y_i$  i.e.,  $\frac{1}{F(\mathbf{x}'_i \hat{\boldsymbol{\beta}})[1 - F(\mathbf{x}'_i \hat{\boldsymbol{\beta}})]}$ , WNLS and ML estimation have the same FOC, except  $\hat{\boldsymbol{\beta}}$  in WNLS's weight is from NLS or last step while that in MLE is still  $\hat{\boldsymbol{\beta}}_{mle}$ .

#### 1.3.3 Algorithm

For  $t = 1 \dots$

Calculate the weight  $w_i^t$  with  $\hat{\boldsymbol{\beta}}_{wnls}^{t-1}$  from last step or  $\hat{\boldsymbol{\beta}}_{nls}^0$  from NLS if  $t = 1$ .

Get  $\hat{\boldsymbol{\beta}}_{wnls}^t = \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^N w_i^t (y_i - F(\mathbf{x}'_i \boldsymbol{\beta}))^2 = \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^N (\sqrt{w_i^t} y_i - \sqrt{w_i^t} F(\mathbf{x}'_i \boldsymbol{\beta}))^2$  with appropriate algorithm for NLS e.g., Gauss-Newton method.

Repeat above two steps until convergence of  $\hat{\beta}_{wnls}^t$ , which will be equal to  $\hat{\beta}_{mle}$ .

ML and NLS estimator are special cases of M estimator, which is a special case of extremum estimator (Amemiya, 1985). Therefore, under regularity conditions of extremum estimator, NLS estimator is consistent and asymptotic normal. However, it is inefficient. WNLS with an inverse conditional variance weight is efficient. Thus, it is not necessary to apply IWNLS.

## 1.4 GMM Estimation

### 1.4.1 Unconditional Moments

Set  $w_i = \frac{1}{F(\mathbf{x}'_i\beta_0)(1-F(\mathbf{x}'_i\beta_0))}$ ,

$$\begin{aligned}\mathbb{E}\left[\frac{\partial w_i(y_i - F(\mathbf{x}'_i\beta))}{\partial \beta} \middle| \beta_0\right] &= \mathbf{0} \\ \mathbb{E}\left[\underbrace{w_i F'(\mathbf{x}'_i\beta_0)(y_i - F(\mathbf{x}'_i\beta_0))\mathbf{x}_i}_{\mathbf{g}(y_i, \mathbf{x}_i; \beta_0)}\right] &= \mathbf{0}\end{aligned}$$

It is just-identified because the number of moments i.e.,  $\dim(\mathbf{x}_i) = \dim(\beta_0)$  i.e., the number of parameters. GMM reduced to MM estimation.

### 1.4.2 Asymptotic Distribution

$$\sqrt{N}(\hat{\beta}_{gmm} - \beta_0) \rightarrow_d N(\mathbf{0}, (\mathbf{G}'\mathbf{W}\mathbf{G})^{-1}\mathbf{G}'\mathbf{W}\mathbf{S}\mathbf{W}\mathbf{G}(\mathbf{G}'\mathbf{W}\mathbf{G})^{-1})$$

where  $\mathbf{G} = \mathbb{E}\left[\frac{\partial \mathbf{g}(y_i, \mathbf{x}_i; \beta)}{\partial \beta} \middle| \beta_0\right]$  and  $\mathbf{S} = \mathbb{E}[\mathbf{g}(y_i, \mathbf{x}_i; \beta_0)\mathbf{g}(y_i, \mathbf{x}_i; \beta_0)']$ .  $\mathbf{S}$  has this simple form because  $y_i$  is assumed to be independent and must be Bernoulli distributed i.e., i.i.d.

$$\begin{aligned}\mathbf{G} &= \mathbb{E}\left[\frac{\partial \mathbf{g}(y_i, \mathbf{x}_i; \beta)}{\partial \beta} \middle| \beta_0\right] \\ &= \mathbb{E}\left[\frac{\partial w_i F'(\mathbf{x}'_i\beta)(y_i - F(\mathbf{x}'_i\beta))\mathbf{x}_i}{\partial \beta} \middle| \beta_0\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\frac{\frac{F'(\mathbf{x}'_i\beta)}{F(\mathbf{x}'_i\beta_0)(1-F(\mathbf{x}'_i\beta_0))}(y_i - F(\mathbf{x}'_i\beta))\mathbf{x}_i}{\partial \beta} \middle| \beta_0, \mathbf{x}_i\right]\right] \\ &= \mathbb{E}\left[-\frac{(F'(\mathbf{x}'_i\beta_0))^2}{F(\mathbf{x}'_i\beta_0)(1-F(\mathbf{x}'_i\beta_0))}\mathbf{x}_i\mathbf{x}'_i\right] \\ \mathbf{S} &= \mathbb{E}[\mathbf{g}(y_i, \mathbf{x}_i; \beta_0)\mathbf{g}(y_i, \mathbf{x}_i; \beta_0)'] \\ &= \mathbb{E}[w_i F'(\mathbf{x}'_i\beta_0)(y_i - F(\mathbf{x}'_i\beta_0))\mathbf{x}_i(w_i F'(\mathbf{x}'_i\beta_0)(y_i - F(\mathbf{x}'_i\beta_0))\mathbf{x}_i)'] \\ &= \mathbb{E}[(w_i F'(\mathbf{x}'_i\beta_0))^2 \mathbb{E}[(y_i - F(\mathbf{x}'_i\beta_0))^2 | \mathbf{x}_i] \mathbf{x}_i\mathbf{x}'_i] \\ &= \mathbb{E}[(w_i F'(\mathbf{x}'_i\beta_0))^2 \text{Var}(y_i | \mathbf{x}_i) \mathbf{x}_i\mathbf{x}'_i] \\ &= \mathbb{E}\left[\left(\frac{F'(\mathbf{x}'_i\beta_0)}{F(\mathbf{x}'_i\beta_0)(1-F(\mathbf{x}'_i\beta_0))}\right)^2 F(\mathbf{x}'_i\beta_0)(1-F(\mathbf{x}'_i\beta_0))\mathbf{x}_i\mathbf{x}'_i\right] \\ &= \mathbb{E}\left[\frac{(F'(\mathbf{x}'_i\beta_0))^2}{F(\mathbf{x}'_i\beta_0)(1-F(\mathbf{x}'_i\beta_0))}\mathbf{x}_i\mathbf{x}'_i\right] = -\mathbf{G}\end{aligned}$$

Because of just-identification, GMM reduces to MM estimation,  $\mathbf{G}$  is a square matrix and thus invertible.

$$\begin{aligned}\sqrt{N}(\hat{\beta}_{gmm} - \beta_0) &\rightarrow_d N(\mathbf{0}, (\mathbf{G}'\mathbf{W}\mathbf{G})^{-1}\mathbf{G}'\mathbf{W}\mathbf{S}\mathbf{W}\mathbf{G}(\mathbf{G}'\mathbf{W}\mathbf{G})^{-1}) \\ &= N(\mathbf{0}, \mathbf{G}^{-1}\mathbf{W}^{-1}\mathbf{G}'^{-1}\mathbf{G}'\mathbf{W}\mathbf{S}\mathbf{W}\mathbf{G}\mathbf{G}^{-1}\mathbf{W}^{-1}\mathbf{G}'^{-1}) \\ &= N(\mathbf{0}, \mathbf{G}^{-1}\mathbf{S}\mathbf{G}'^{-1}) \\ &= N(\mathbf{0}, [\mathbf{I}(\beta_0)]^{-1})\end{aligned}$$

Because

$$\mathbf{G}^{-1}\mathbf{S}\mathbf{G}'^{-1} = -\mathbf{G}'^{-1} = \{\mathbb{E}\left[\frac{(F'(\mathbf{x}'_i\beta_0))^2}{F(\mathbf{x}'_i\beta_0)(1-F(\mathbf{x}'_i\beta_0))}\mathbf{x}_i\mathbf{x}'_i\right]\}^{-1} = [\mathbf{I}(\beta_0)]^{-1}$$

Thus, GMM/MM and ML estimator has the same asymptotic distribution.

## 1.5 Special Case: Logit Model

If  $F(.) = \Lambda(.)$  i.e., Logistic function (the c.d.f. of standard Logistic random variable),

$$p_i := Pr(y_i = 1 | \mathbf{x}_i) = \Lambda(\mathbf{x}'_i \boldsymbol{\beta}) = \frac{e^{\mathbf{x}'_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}'_i \boldsymbol{\beta}}} = \frac{1}{1 + e^{-\mathbf{x}'_i \boldsymbol{\beta}}}$$

$$\ln \frac{p_i}{1 - p_i} = \Lambda^{-1}(p_i) = \mathbf{x}'_i \boldsymbol{\beta}$$

$\Lambda^{-1}(.)$  is called Logit function

### 1.5.1 Marginal Effect

$$\begin{aligned} \Lambda'(z) &= \frac{d(1 + e^{-z})^{-1}}{dz} \\ &= -(1 + e^{-z})^{-2} e^{-z} (-1) \\ &= (1 + e^{-z})^{-1} (1 + e^{-z})^{-1} e^{-z} \\ &= \Lambda(z) \frac{e^{-z}}{1 + e^{-z}} \\ &= \Lambda(z) \frac{1}{1 + e^z} \\ &= \Lambda(z) \frac{1 + e^z - e^z}{1 + e^z} \\ &= \Lambda(z)(1 - \Lambda(z)) \end{aligned}$$

$$\Lambda'(\mathbf{x}'_i \boldsymbol{\beta}) \boldsymbol{\beta} = \Lambda(\mathbf{x}'_i \boldsymbol{\beta})(1 - \Lambda(\mathbf{x}'_i \boldsymbol{\beta})) \boldsymbol{\beta} \leq 0.25 \boldsymbol{\beta}$$

As

$$\begin{aligned} \left. \frac{dz(1 - z)}{dz} \right|_{z^*} &= 1 - 2z^* = 0 \\ z^* &= 1/2 \\ z^*(1 - z^*) &= 1/2 \cdot 1/2 = 1/4 = 0.25 \end{aligned}$$

$\Lambda(\mathbf{x}'_i \boldsymbol{\beta})(1 - \Lambda(\mathbf{x}'_i \boldsymbol{\beta})) = 0.25$  when  $\Lambda(\mathbf{x}'_i \boldsymbol{\beta}) = 1/2$  which happens when  $\mathbf{x}'_i \boldsymbol{\beta} = 0$  as p.d.f. of standard Logistic random variable is symmetric at 0 (c.d.f. = 0.5 at 0).

### 1.5.2 Odds Ratio

$$\ln \frac{p_i}{1 - p_i} = \Lambda^{-1}(p_i) = \mathbf{x}'_i \boldsymbol{\beta}$$

If  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are different in  $x_k$  by 1 unit, then

$$\begin{aligned} \mathbf{x}'_i \boldsymbol{\beta} - \mathbf{x}'_j \boldsymbol{\beta} &= \ln \frac{p_i}{1 - p_i} - \ln \frac{p_j}{1 - p_j} \\ \mathbf{x}'_j \boldsymbol{\beta} + 1 \cdot \beta_k - \mathbf{x}'_j \boldsymbol{\beta} &= \beta_k = \ln \frac{p_i/(1 - p_i)}{p_j/(1 - p_j)} \\ \exp(\beta_k) &= \frac{p_i/(1 - p_i)}{p_j/(1 - p_j)} := OR \end{aligned}$$

Odds Ratio (OR) can be interpreted as

$$\begin{aligned} \frac{p_j}{1 - p_j} &= \exp(\mathbf{x}'_j \boldsymbol{\beta}) \\ \exp(\mathbf{x}'_i \boldsymbol{\beta}) &= \exp(\mathbf{x}'_j \boldsymbol{\beta} + 1 \cdot \beta_k) = \exp(\mathbf{x}'_j \boldsymbol{\beta}) \exp(\beta_k) \\ &= \frac{p_j}{1 - p_j} \exp(\beta_k) \end{aligned}$$

So, an unit increase in  $x_k$  means the odds  $\frac{p_j}{1 - p_j}$  is multiplied by the Odds Ratio (OR)  $\exp(\beta_k)$ .

### 1.5.3 First Order Condition

$$\begin{aligned}
\sum_{i=1}^N \frac{y_i - \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})}{\Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})(1 - \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}}))} \Lambda'(\mathbf{x}'_i \hat{\boldsymbol{\beta}}) \mathbf{x}_i &= \mathbf{0} \\
\sum_{i=1}^N \frac{y_i - \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})}{\Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})(1 - \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}}))} \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})(1 - \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})) \mathbf{x}_i &= \mathbf{0} \\
\sum_{i=1}^N (y_i - \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})) \mathbf{x}_i &= \mathbf{0} \\
\sum_{i=1}^N (y_i - \mathbb{E}(y_i | \mathbf{x}_i)) \mathbf{x}_i &= \mathbf{0}
\end{aligned}$$

This is similar to the first order condition of OLS estimation of linear model. Moreover, if intercept is included in  $\mathbf{x}_i$ .

$$\begin{aligned}
\sum_{i=1}^N (y_i - \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}})) \cdot 1 &= 0 && \text{"residual" sum to 0} \\
N^{-1} \sum_{i=1}^N \Lambda(\mathbf{x}'_i \hat{\boldsymbol{\beta}}) &= \bar{y}
\end{aligned}$$

Interesting result,  $\bar{y}$  is the percentage of one in the sample, which is equal to the average predicted probability of Logit Model.

## 1.6 Special Case: Probit Model

If  $F(\cdot) = \Phi(\cdot)$  i.e., the c.d.f. of standard Normal random variable,

$$p_i := \Pr(y_i = 1 | \mathbf{x}_i) = \Phi(\mathbf{x}'_i \boldsymbol{\beta}) = \int_{-\infty}^{\mathbf{x}'_i \boldsymbol{\beta}} \phi(z) dz$$

$$\Phi^{-1}(p_i) = \mathbf{x}'_i \boldsymbol{\beta}$$

$\Phi^{-1}(\cdot)$  is called Probit function, no closed form

### 1.6.1 Marginal Effect

$$\begin{aligned}
\Phi'(z) &= \frac{d \int_{-\infty}^z \phi(a) da}{dz} \\
&= \phi(z) \\
&= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z^2\right)
\end{aligned}$$

by First Fundamental Theorem of Calculus

$$\begin{aligned}
\Phi'(\mathbf{x}'_i \boldsymbol{\beta}) \boldsymbol{\beta} &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (\mathbf{x}'_i \boldsymbol{\beta})^2\right) \boldsymbol{\beta} \\
&\leq \frac{1}{\sqrt{2\pi}} \cdot 1 \cdot \boldsymbol{\beta} \\
&\approx 0.4 \boldsymbol{\beta}
\end{aligned}$$

as  $0 < \exp(z) \leq 1$  if  $z \leq 0$

### 1.6.2 First Order Condition

$$\sum_{i=1}^N \underbrace{\left\{ \frac{\Phi'(\mathbf{x}'_i \hat{\boldsymbol{\beta}})}{\Phi(\mathbf{x}'_i \hat{\boldsymbol{\beta}})(1 - \Phi(\mathbf{x}'_i \hat{\boldsymbol{\beta}}))} \right\}}_{\hat{w}_i} [y_i - \Phi(\mathbf{x}'_i \hat{\boldsymbol{\beta}})] \mathbf{x}_i = \mathbf{0}$$

## 1.7 Special Case: Robit Model

if  $F(\cdot) = F_{t,\nu}(\cdot)$  i.e., the c.d.f. of standard student's t random variable,

$$p_i := \Pr(y_i = 1 | \mathbf{x}_i) = F_{t,\nu}(\mathbf{x}_i' \boldsymbol{\beta}) = \int_{-\infty}^{\mathbf{x}_i' \boldsymbol{\beta}} t(z) dz$$

$$F_{t,\nu}^{-1}(p_i) = \mathbf{x}_i' \boldsymbol{\beta}$$

$F_{t,\nu}^{-1}(\cdot)$  is called Robit function, no closed form

As standard student's t random variable converges to standard normal random variable when the degree of freedom  $\nu \rightarrow \infty$ , Robit model is a generalization of Probit model. The extra flexibility provided by  $\nu$  seems to handle outliers well. Except some special cases,  $F_{t,\nu}(\cdot)$  does not have a closed form. The additional parameter  $\nu$  can be estimated by ML method with  $\boldsymbol{\beta}$  simultaneously, or a grid of  $\nu$  can be pre-specified.

## 1.8 Special Case: Linear Probability Model (LPM)

If  $F(\cdot)$  is an identity function,

$$p_i := \Pr(y_i = 1 | \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}$$

However, it is not likely that  $F(\cdot)$  is an identity function because the resulting predicted probability can be larger than 1 or smaller than 0. MLE's first order condition is

$$\sum_{i=1}^N \left\{ \frac{1}{\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle} (1 - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})} \right\} [y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle}] \mathbf{x}_i = \mathbf{0}$$

There is default heteroskedasticity problem. As shown before,  $\text{Var}(y_i | \mathbf{x}_i) = p_i(1 - p_i) = \mathbf{x}_i' \boldsymbol{\beta} (1 - \mathbf{x}_i' \boldsymbol{\beta})$  which depend on  $i$ . Heteroskedasticity can be solved by using GLS estimation i.e., WLS estimation with an inverse conditional variance weight. First, use  $\hat{\boldsymbol{\beta}}_{ols}$  to get the weight  $[\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols} (1 - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols})]^{-1}$  for  $\forall i$ . Second, get

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{wls} &= \underset{\boldsymbol{\beta}}{\text{argmin}} \sum_{i=1}^N \left\{ \frac{1}{\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols} (1 - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols})} \right\} [y_i - \mathbf{x}_i' \boldsymbol{\beta}]^2 \\ &= \underset{\boldsymbol{\beta}}{\text{argmin}} \sum_{i=1}^N \left\{ \frac{1}{\sqrt{\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols} (1 - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols})}} \right\}^2 [y_i - \mathbf{x}_i' \boldsymbol{\beta}]^2 \\ &= \underset{\boldsymbol{\beta}}{\text{argmin}} \sum_{i=1}^N \left[ \frac{y_i}{\sqrt{\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols} (1 - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols})}} - \frac{\mathbf{x}_i'}{\sqrt{\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols} (1 - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols})}} \boldsymbol{\beta} \right]^2 \\ &= \underset{\boldsymbol{\beta}}{\text{argmin}} \sum_{i=1}^N [\tilde{y}_i - \tilde{\mathbf{x}}_i' \boldsymbol{\beta}]^2 \end{aligned}$$

Therefore,  $\hat{\boldsymbol{\beta}}_{wls} = (\sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i')^{-1} \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{y}_i$ . Its first order condition is

$$\sum_{i=1}^N \left\{ \frac{1}{\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols} (1 - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{ols})} \right\} [y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{wls}] \mathbf{x}_i = \mathbf{0}$$

If  $\hat{\boldsymbol{\beta}}_{wls}$  is then used in estimating the weight and loop the above process until convergence, the resulting IWLS estimate is equal to ML estimate. It is a special case of IWNLS = ML estimation derived in the previous section for general model.  $\mathbf{x}_i' \boldsymbol{\beta} \rightarrow 0$  or  $\rightarrow 1$ ,  $\hat{w}_i$  is large and lead to numerical instability during optimization.

If  $F(\cdot)$  is truly an identity function, OLS estimator is still unbiased and consistent. The default heteroskedasticity can be tackled by using robust standard error e.g., Eicker-Huber-White standard error. However, it is still inefficient.

Some scholars e.g., Angrist advocate using LPM e.g., if your target is to estimate Average Treatment Effect (ATE) and you have only one indicator regressor, the OLS slope estimate is estimated ATE under independent assumption, which can be justified by random assignment. This result holds even  $y$  is not binary.  $\beta = \frac{\text{Cov}(Y, D)}{\text{Var}(D)} = \mathbb{E}(Y | D = 1) - \mathbb{E}(Y | D = 0)$ , second equality uses double expectation to derive.

$$\begin{aligned} \mathbb{E}(Y | D = 1) - \mathbb{E}(Y | D = 0) &= \mathbb{E}(Y_1 | D = 1) - \mathbb{E}(Y_0 | D = 0) \\ &= \mathbb{E}(Y_1) - \mathbb{E}(Y_0) \\ &= \mathbb{E}(Y_1 - Y_0) = ATE \end{aligned}$$

$Y_1, Y_0 \perp D$  under random assignment



Stoker's average derivatives provide the average marginal effect without specifying  $F(\cdot)$ . Define  $\mathbb{E}[y|\mathbf{x}] = m(\mathbf{x})$ ,

$$\begin{aligned}
AME &:= \mathbb{E}\left[\frac{\partial m(\mathbf{x})}{\partial \mathbf{x}}\right] = -\mathbb{E}\left[m(\mathbf{x}) \frac{\partial \ln(f(\mathbf{x}))}{\partial \mathbf{x}}\right] && \text{by generalized Information matrix equality} \\
&= -\mathbb{E}\left[\mathbb{E}[y|\mathbf{x}] \frac{f'(\mathbf{x})}{f(\mathbf{x})}\right] \\
&= -\mathbb{E}\left[\mathbb{E}\left[y \frac{f'(\mathbf{x})}{f(\mathbf{x})} \middle| \mathbf{x}\right]\right] \\
&= -\mathbb{E}\left[y \frac{f'(\mathbf{x})}{f(\mathbf{x})}\right]
\end{aligned}$$

Thus,

$$\widehat{AME} = -N^{-1} \sum_{i=1}^N y_i \frac{\widehat{f}'(\mathbf{x}_i)}{\widehat{f}(\mathbf{x}_i)}$$

If  $\mathbf{x}$  follow multivariate normal,

$$\begin{aligned}
AME &= -\mathbb{E}[y \cdot -(\mathbf{x} - \boldsymbol{\mu}_x) \boldsymbol{\Sigma}_x^{-1}] \\
&= \mathbb{E}[y(\mathbf{x} - \boldsymbol{\mu}_x) \{\mathbb{E}[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)']\}^{-1}] \\
&= \{\mathbb{E}[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)']\}^{-1} \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu}_x)y]
\end{aligned}$$

Thus, it is an OLS estimator of regressing  $y$  on  $\mathbf{x} - \boldsymbol{\mu}_x$ . This result holds no matter  $y$  is binary or not. If your target is AME with multivariate normal  $\mathbf{x}$ , you can simply regress  $y$  on  $\mathbf{x} - \boldsymbol{\mu}_x$ , the estimated coefficients is estimated AME. Even  $\mathbf{x}$  is not multivariate normal,  $AME$  can be non-parametrically estimated using  $-N^{-1} \sum_{i=1}^N y_i \frac{\widehat{f}'(\mathbf{x}_i)}{\widehat{f}(\mathbf{x}_i)}$ .

## 1.9 Model Evaluation: Pseudo- $R^2$

McFadden (1974) suggests

$$\begin{aligned}
R_{Binary}^2 &= 1 - \frac{\ln L_{fit}}{\ln L_0} \\
&= 1 - \frac{\sum_{i=1}^N \{y_i \ln \widehat{p}_i + (1 - y_i) \ln(1 - \widehat{p}_i)\}}{\sum_{i=1}^N \{y_i \ln \bar{y} + (1 - y_i) \ln(1 - \bar{y})\}} \\
&= 1 - \frac{\sum_{i=1}^N \{y_i \ln \widehat{p}_i + (1 - y_i) \ln(1 - \widehat{p}_i)\}}{(\sum_{i=1}^N y_i) \ln \bar{y} + (N - \sum_{i=1}^N y_i) \ln(1 - \bar{y})} \\
&= 1 - \frac{\sum_{i=1}^N \{y_i \ln \widehat{p}_i + (1 - y_i) \ln(1 - \widehat{p}_i)\}}{N \bar{y} \ln \bar{y} + N(1 - \bar{y}) \ln(1 - \bar{y})} \\
&= 1 - \frac{\sum_{i=1}^N \{y_i \ln \widehat{p}_i + (1 - y_i) \ln(1 - \widehat{p}_i)\}}{N(\bar{y} \ln \bar{y} + (1 - \bar{y}) \ln(1 - \bar{y}))}
\end{aligned}$$

## 1.10 Other Model Evaluation Methods

In-sample accuracy, out-of-sample accuracy, cross validation accuracy, confusion matrix, ROC, etc.

## 1.11 The Motivation of the choice of $F(\cdot)$

It can be motivated by Latent Variable Models and Generalized Linear Model (GLM) discussed below.

## 2 Latent Variable Models

### 2.1 Index Function Model

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + u_i \quad y_i^* \text{ is unobservable}$$

But we can observe  $y_i$ ,

$$\begin{aligned} y_i &= 1(y_i^* > 0) \\ \mathbb{E}(y_i | \mathbf{x}_i) &= 1 \cdot Pr(y_i = 1 | \mathbf{x}_i) + 0 \cdot Pr(y_i = 0 | \mathbf{x}_i) \\ &= Pr(y_i = 1 | \mathbf{x}_i) \\ &= Pr(y_i^* > 0 | \mathbf{x}_i) \\ &= Pr(\mathbf{x}_i' \boldsymbol{\beta} + u_i > 0 | \mathbf{x}_i) \\ &= Pr(u_i > -\mathbf{x}_i' \boldsymbol{\beta} | \mathbf{x}_i) \\ &= Pr(u_i \leq \mathbf{x}_i' \boldsymbol{\beta} | \mathbf{x}_i) \quad \text{if } u_i \text{ is symmetric at 0} \\ &= F_u(\mathbf{x}_i' \boldsymbol{\beta}) \end{aligned}$$

If  $u_i$  follows standard Logistic distribution, the model is Logit model. If  $u_i$  follows standard Normal distribution, it is Probit model. If  $u_i$  follows standard student's t distribution, it is Robit model.

### 2.2 Identification of parameters

$$y_i = 1 \implies y_i^* > 0 \implies \mathbf{x}_i' \boldsymbol{\beta} + u_i > 0$$

However, for any constant  $c > 0$ ,

$$\mathbf{x}_i' \boldsymbol{\beta} + u_i > 0 \iff \mathbf{x}_i' c\boldsymbol{\beta} + cu_i > 0$$

So,  $\boldsymbol{\beta}$  is not identified with  $y_i = 1(y_i^* > 0)$ . Thus, we restrict  $Var(u_i | \mathbf{x}_i)$  to identify  $\boldsymbol{\beta}$ .

If  $u_i$  follows standard Logistic distribution,  $Var(u_i | \mathbf{x}_i) = \pi^2/3$ . If  $u_i$  follows standard Normal distribution,  $Var(u_i | \mathbf{x}_i) = 1$ .

### 2.3 Additive Random Utility Model (ARUM)

$y = 0$  means choosing option 0. Utility obtained from this is  $U_0$ ;  $y = 1$  means choosing option 1. Utility obtained from this is  $U_1$ .

$$\begin{aligned} U_0 &= V_0 + \varepsilon_0 \\ U_1 &= V_1 + \varepsilon_1 \end{aligned} \quad V_0 \text{ is deterministic component of utility}$$

$$\begin{aligned} y &= 1(U_1 > U_0) \\ \mathbb{E}(y | \mathbf{x}) &= 1 \cdot Pr(y = 1 | \mathbf{x}) + 0 \cdot Pr(y = 0 | \mathbf{x}) \\ &= Pr(y = 1 | \mathbf{x}) \\ &= Pr(U_1 > U_0 | \mathbf{x}) \\ &= Pr(V_1 + \varepsilon_1 > V_0 + \varepsilon_0 | \mathbf{x}) \\ &= Pr(V_1 - V_0 > \varepsilon_0 - \varepsilon_1 | \mathbf{x}) \\ &= Pr(\varepsilon_0 - \varepsilon_1 < V_1 - V_0 | \mathbf{x}) \\ &= F_{\varepsilon_0 - \varepsilon_1}(V_1 - V_0) \end{aligned}$$

#### 2.3.1 Special Case: Logit Model

If  $\varepsilon_0$  and  $\varepsilon_1$  are independent and both follows Type 1 Extreme Value distribution (log Weibull distribution). It can be shown  $\varepsilon_0 - \varepsilon_1$  follows standard Logistic distribution i.e.,  $F_{\varepsilon_0 - \varepsilon_1}(\cdot) = \Lambda(\cdot)$ .

It can also be shown by direct integration,

$$\begin{aligned}
Pr(y = 1|\mathbf{x}) &= Pr(\varepsilon_0 - \varepsilon_1 < V_1 - V_0|\mathbf{x}) \\
&= Pr(\varepsilon_0 < \varepsilon_1 + V_1 - V_0|\mathbf{x}) \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\varepsilon_1 + V_1 - V_0} f(\varepsilon_0, \varepsilon_1) \partial\varepsilon_0 \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\varepsilon_1 + V_1 - V_0} f_{\varepsilon_0}(\varepsilon_0) f_{\varepsilon_1}(\varepsilon_1) \partial\varepsilon_0 \partial\varepsilon_1 && \text{as independence} \\
&= \int_{-\infty}^{\infty} f_{\varepsilon_1}(\varepsilon_1) \left[ \int_{-\infty}^{\varepsilon_1 + V_1 - V_0} f_{\varepsilon_0}(\varepsilon_0) \partial\varepsilon_0 \right] \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} f_{\varepsilon_1}(\varepsilon_1) \left[ \int_{-\infty}^{\varepsilon_1 + V_1 - V_0} e^{-\varepsilon_0} \exp(-e^{-\varepsilon_0}) \partial\varepsilon_0 \right] \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} f_{\varepsilon_1}(\varepsilon_1) \exp(-e^{-\varepsilon_0}) \Big|_{-\infty}^{\varepsilon_1 + V_1 - V_0} \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} f_{\varepsilon_1}(\varepsilon_1) [\exp(-e^{-(\varepsilon_1 + V_1 - V_0)}) - \exp(-e^{-\infty})] \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} f_{\varepsilon_1}(\varepsilon_1) \exp(-e^{-(\varepsilon_1 + V_1 - V_0)}) \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} e^{-\varepsilon_1} \exp(-e^{-\varepsilon_1}) \exp(-e^{-(\varepsilon_1 + V_1 - V_0)}) \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} e^{-\varepsilon_1} \exp(-e^{-\varepsilon_1} - e^{-(\varepsilon_1 + V_1 - V_0)}) \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} e^{-\varepsilon_1} \exp(-e^{-\varepsilon_1} - e^{-\varepsilon_1} e^{-(V_1 - V_0)}) \partial\varepsilon_1 \\
&= \int_{-\infty}^{\infty} e^{-\varepsilon_1} \exp(-e^{-\varepsilon_1} (1 + e^{-(V_1 - V_0)})) \partial\varepsilon_1 \\
&= 1/(1 + e^{-(V_1 - V_0)}) && \text{as } \int_{-\infty}^{\infty} a e^{-\varepsilon} \exp(-a e^{-\varepsilon}) d\varepsilon = 1 \\
&= \Lambda(V_1 - V_0)
\end{aligned}$$

If  $V_1 - V_0 = \mathbf{x}'\boldsymbol{\beta}$ ,  $Pr(y = 1|\mathbf{x}) = \Lambda(\mathbf{x}'\boldsymbol{\beta})$ . It is Logit model.

### 2.3.2 Special Case: Probit Model

If  $\varepsilon_0$  and  $\varepsilon_1$  are multivariate (bivariate here) standard normally distributed, any linear combination of  $\varepsilon_0$  and  $\varepsilon_1$  also follow standard normal. So,  $\varepsilon_0 - \varepsilon_1$  follows univariate standard normal. i.e.,  $F_{\varepsilon_0 - \varepsilon_1}(\cdot) = \Phi(\cdot)$ .

### 3 Lagrange Multiplier Test, that is LM Test

that is easy to go through, no reason to take a break.

$$\mathbf{S} = \begin{pmatrix} s_1(\hat{\boldsymbol{\beta}})' \\ \vdots \\ s_N(\hat{\boldsymbol{\beta}})' \end{pmatrix} \quad \text{where } s_i(\hat{\boldsymbol{\beta}}) = \frac{\partial \ln l_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\hat{\boldsymbol{\beta}}}$$

$$\alpha_{LM} = N \frac{\mathbf{1}' \mathbf{S} (\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}' \mathbf{1}}{\mathbf{1}' \mathbf{1}} = N \cdot R^2 \rightarrow_d \chi^2$$

$\mathbf{S}(\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}'$  is the Hat matrix. Thus,  $\mathbf{S}(\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}' \mathbf{1}$  is an OLS predicted value of Outer-Product-of-the-Gradient (OPG) regression in which  $\mathbf{1}$  is dependent variable and  $\mathbf{S}$  is regressors. The Sum of Squares Regression i.e., Explained Sum of Squares is

$$\begin{aligned} (\mathbf{S}(\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}' \mathbf{1})' \mathbf{S}(\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}' \mathbf{1} &= \mathbf{1}' \mathbf{S}(\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}' \mathbf{S}(\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}' \mathbf{1} \\ &= \mathbf{1}' \mathbf{S}(\mathbf{S}' \mathbf{S})^{-1} \mathbf{S}' \mathbf{1} \end{aligned}$$

$\mathbf{1}' \mathbf{1}$  is Sum of Squares Total. Thus,  $\alpha_{LM}$  is  $N$  times  $R^2$ .

#### 3.1 Normality Test

If the density function of error term  $u_i$  in latent variable model is within the Pearson family, its c.d.f. is

$$Pr(u_i \leq t) = \Phi(t + \gamma_1 t^2 + \gamma_2 t^3)$$

where  $\gamma_1$  controls the third moment i.e., skewness while  $\gamma_2$  controls the fourth moment i.e., excess kurtosis. Clearly,  $u_i$  is standard normally distributed if both  $\gamma_1$  and  $\gamma_2$  are zero. This is the zero hypothesis.

Extended probit model becomes

$$Pr(y_i = 1 | \mathbf{x}_i) = \Phi(\mathbf{x}_i' \boldsymbol{\beta} + \gamma_1 (\mathbf{x}_i' \boldsymbol{\beta})^2 + \gamma_2 (\mathbf{x}_i' \boldsymbol{\beta})^3)$$

$$s_i \left( \begin{pmatrix} \hat{\boldsymbol{\beta}}_{mle} \\ \hat{\gamma}_{1,mle} \\ \hat{\gamma}_{2,mle} \end{pmatrix} \right) = \begin{pmatrix} \frac{\partial \ln l_i(\boldsymbol{\beta}, \gamma_1, \gamma_2)}{\partial \boldsymbol{\beta}} \Big|_{\hat{\boldsymbol{\beta}}_{mle}} \\ \frac{\partial \ln l_i(\boldsymbol{\beta}, \gamma_1, \gamma_2)}{\partial \gamma_1} \Big|_{\hat{\gamma}_{1,mle}} \\ \frac{\partial \ln l_i(\boldsymbol{\beta}, \gamma_1, \gamma_2)}{\partial \gamma_2} \Big|_{\hat{\gamma}_{2,mle}} \end{pmatrix} = \begin{pmatrix} \hat{\varepsilon}_i^G \mathbf{x}_i \\ \hat{\varepsilon}_i^G (\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})^2 \\ \hat{\varepsilon}_i^G (\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})^3 \end{pmatrix}$$

where  $\hat{\varepsilon}_i^G := \frac{y_i - \Phi(\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})}{\Phi(\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})(1 - \Phi(\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle}))} \Phi'(\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})$  is the generalized residual. We know that matrix  $\mathbf{S}$  has variables  $\hat{\varepsilon}_i^G \mathbf{x}_i$ ,  $\hat{\varepsilon}_i^G (\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})^2$  and  $\hat{\varepsilon}_i^G (\mathbf{x}_i' \hat{\boldsymbol{\beta}}_{mle})^3$ . We can then regress 1 on these variables. The  $R^2$  of such model times  $N$  is  $\alpha_{LM}$ , which is asymptotically Chi-squared distributed with two degree of freedom under zero hypothesis. If  $\alpha_{LM}$  is large enough, we reject zero hypothesis i.e., reject normality of  $u_i$  and  $F(\cdot)$  should not be  $\Phi(\cdot)$ .

## 4 Generalized Linear Model

GLM is originated from Statistics, and usually not covered in Econometrics textbook.

### 4.1 Exponential Family of Distribution

#### 4.1.1 Probability Mass Function

$$f(y|\theta, \eta) = a(y) \exp\left[\frac{y \cdot \theta - b(\theta)}{\eta} + c(y, \eta)\right]$$

#### 4.1.2 Moments

$$\mathbb{E}(y) = b'(\theta)$$

$$Var(y) = b''(\theta)\eta$$

If  $\theta = \mathbf{x}'_i \boldsymbol{\beta}$ , we have  $\mathbb{E}(y) = b'(\mathbf{x}'_i \boldsymbol{\beta})$ . Thus,  $b'(\cdot)$  is the canonical mean function and  $b'^{-1}(\cdot)$  is the canonical link function.

#### 4.1.3 Special Case: Binomial Distribution

$$\begin{aligned} f(y|p) &= \binom{n}{p} p^y (1-p)^{n-y} \\ &= \exp[\ln(\binom{n}{p} p^y (1-p)^{n-y})] \\ &= \exp[\ln \binom{n}{p} + \ln(p^y (1-p)^{n-y})] \\ &= \exp[\ln \binom{n}{p} + y \cdot \ln(p) + (n-y) \ln(1-p)] \\ &= \exp[\ln \binom{n}{p} + y \cdot \ln(p) + n \cdot \ln(1-p) - y \cdot \ln(1-p)] \\ &= \exp[\ln \binom{n}{p} + y(\ln(p) - \ln(1-p)) + n \cdot \ln(1-p)] \\ &= \exp[y \ln \frac{p}{1-p} - (-n \cdot \ln(1-p)) + \ln \binom{n}{p}] \end{aligned}$$

Thus,  $\eta = 1$ ,  $a(y) = 1$ ,  $\theta = \ln \frac{p}{1-p}$ ,  $b(\theta) = -n \cdot \ln(1-p)$ ,  $c(y, \eta) = \ln \binom{n}{p}$

$$\theta = \ln \frac{p}{1-p} = \text{logit}(p) \iff \text{logistic}(\theta) = p$$

$$\begin{aligned} b(\theta) &= -n \cdot \ln(1-p) \\ &= -n \cdot \ln(1 - \text{logistic}(\theta)) \\ &= -n \cdot \ln\left(1 - \frac{1}{1 + e^{-\theta}}\right) \\ &= -n \cdot \ln\left(\frac{e^{-\theta}}{1 + e^{-\theta}}\right) \\ &= -n \cdot \ln\left(\frac{1}{1 + e^{\theta}}\right) \\ &= n \cdot \ln(1 + e^{\theta}) \end{aligned}$$

$$b'(\theta) = \frac{d \cdot n \cdot \ln(1 + e^{\theta})}{d\theta} = n \frac{1}{1 + e^{\theta}} e^{\theta} = n \cdot \text{logistic}(\theta)$$

Thus,  $\mathbb{E}(y) = b'(\theta) = n \cdot \text{logistic}(\theta)$ . Bernoulli distribution is a special case of Binomial distribution with  $n = 1$ . Therefore, binary outcome model, i.e.,  $y$  is Bernoulli distributed, is  $\mathbb{E}(y) = \text{logistic}(\mathbf{x}' \boldsymbol{\beta})$  by setting  $n = 1$  and  $\theta = \mathbf{x}' \boldsymbol{\beta}$ . Logistic function is the canonical mean function for binary outcome model.  $F(\cdot)$  is thus selected to be logistic function.

## 5 Conditional Logistic Regression

Conditional MLE is applied in non-linear panel model in Econometrics.

$$Pr(y_i = 1) = \text{logistic}(\gamma + \mathbf{x}'_i \boldsymbol{\beta}) = \frac{e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}}{1 + e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}}$$

### 5.1 Full Likelihood Function

$$\begin{aligned} L(\mathbf{y}; \boldsymbol{\beta}, \gamma) &= \prod_{i=1}^N \frac{e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta}) y_i}}{1 + e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}} && \text{assume independence of } y_i \\ &= \frac{\prod_{i=1}^N e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta}) y_i}}{\prod_{i=1}^N [1 + e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}]} \\ &= \frac{e^{\sum_{i=1}^N (\gamma + \mathbf{x}'_i \boldsymbol{\beta}) y_i}}{\prod_{i=1}^N [1 + e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}]} \\ &= \frac{e^{\gamma \sum_{i=1}^N y_i + \sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) y_i}}{\prod_{i=1}^N [1 + e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}]} \end{aligned}$$

### 5.2 Conditional Likelihood Function

Define the set  $G_t = \{\mathbf{y} \in \mathbb{R}^{\dim(\mathbf{y})}; \mathbf{y}' \mathbf{1} = t\}$

$$\begin{aligned} L(\mathbf{y} | \mathbf{y}' \mathbf{1} = t; \boldsymbol{\beta}, \gamma) &= \frac{Pr(\mathbf{y} \cap \mathbf{y}' \mathbf{1} = t)}{Pr(\mathbf{y}' \mathbf{1} = t)} && \text{conditional probability} \\ &= \frac{Pr(\mathbf{y} \cap \mathbf{y}' \mathbf{1} = t)}{Pr(\cup_{\mathbf{z} \in \mathbb{R}^{\dim(\mathbf{z})}} \{\mathbf{z} \cap \mathbf{z}' \mathbf{1} = t\})} && \text{total probability} \\ &= \frac{Pr(\mathbf{y} \cap \mathbf{y}' \mathbf{1} = t)}{\sum_{\mathbf{z} \in \mathbb{R}^{\dim(\mathbf{z})}} Pr(\mathbf{z} \cap \mathbf{z}' \mathbf{1} = t)} && \text{no intersection of sets} \\ &= \frac{Pr(\mathbf{y} \cap \mathbf{y}' \mathbf{1} = t)}{\sum_{\mathbf{z} \in G_t} Pr(\mathbf{z} \cap \mathbf{z}' \mathbf{1} = t)} \\ &= \frac{\frac{e^{\gamma t + \sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) y_i}}{\prod_{i=1}^N [1 + e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}]} }{\sum_{\mathbf{z} \in G_t} \frac{e^{\gamma t + \sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) z_i}}{\prod_{i=1}^N [1 + e^{(\gamma + \mathbf{x}'_i \boldsymbol{\beta})}]} } \\ &= \frac{e^{\gamma t + \sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) y_i}}{\sum_{\mathbf{z} \in G_t} e^{\gamma t + \sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) z_i}} \\ &= \frac{e^{\gamma t} e^{\sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) y_i}}{e^{\gamma t} \sum_{\mathbf{z} \in G_t} e^{\sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) z_i}} \\ &= \frac{e^{\sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) y_i}}{\sum_{\mathbf{z} \in G_t} e^{\sum_{i=1}^N (\mathbf{x}'_i \boldsymbol{\beta}) z_i}} \end{aligned}$$

Parameter  $\gamma$  is gone.

## 6 Berkson's Minimum Chi-square Estimator

$$p_t := Pr(y_i = 1 | \mathbf{x}_i = \mathbf{x}_t) = F(\mathbf{x}'_t \boldsymbol{\beta})$$

$$F^{-1}(p_t) = \mathbf{x}'_t \boldsymbol{\beta}$$

$p_t$  can be estimated by  $\bar{p}_t = \bar{y}_t = N_t^{-1} \sum_{i=1}^{N_t} y_{it}$

$$F^{-1}(\bar{p}_t) - F^{-1}(\bar{p}_t) + F^{-1}(p_t) = \mathbf{x}'_t \boldsymbol{\beta}$$

$$F^{-1}(\bar{p}_t) = \mathbf{x}'_t \boldsymbol{\beta} + \underbrace{F^{-1}(\bar{p}_t) - F^{-1}(p_t)}_{v_t}$$

As  $Var(v_t | \mathbf{x}_t)$  depends on  $t$ , there is heteroskedasticity. GLS (WLS here) can be used to estimate  $\boldsymbol{\beta}$  efficiently.

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \sum_{t=1}^T Var(v_t | \mathbf{x}_t)^{-1} (F^{-1}(\bar{p}_t) - \mathbf{x}'_t \boldsymbol{\beta})^2$$

## 7 Semi-parametric Estimation

### 7.1 Maximum Score Estimation (Manski, 1975, 1985)

We predict  $y_i = 1$  if  $\mathbf{x}'_i \boldsymbol{\beta} > 0$ . We predict  $y_i = 0$  if  $\mathbf{x}'_i \boldsymbol{\beta} \leq 0$ . The Score function, which counts the number of correct observations, is defined as

$$\begin{aligned}
 S_N(\boldsymbol{\beta}) &:= \sum_{i=1}^N \{y_i 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + (1 - y_i) 1(\mathbf{x}'_i \boldsymbol{\beta} \leq 0)\} \\
 &= \sum_{i=1}^N \{y_i 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + (1 - y_i)(1 - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0))\} \\
 &= \sum_{i=1}^N \{y_i 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + 1 - y_i - (1 - y_i) 1(\mathbf{x}'_i \boldsymbol{\beta} > 0)\} \\
 &= \sum_{i=1}^N \{y_i 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + 1 - y_i - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + y_i 1(\mathbf{x}'_i \boldsymbol{\beta} > 0)\} \\
 &= \sum_{i=1}^N \{(2y_i - 1) 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + 1 - y_i\} \\
 &= \sum_{i=1}^N (2y_i - 1) 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + N - \sum_{i=1}^N y_i
 \end{aligned}$$

The Maximum Score Estimator is

$$\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} \sum_{i=1}^N (2y_i - 1) 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) \quad \text{can not use differentiation}$$

MSE can be regarded as Least Absolute Deviation (LAD) Estimator. It can be seen

$$\begin{aligned}
 Q_N(\boldsymbol{\beta}) &= N - S_N(\boldsymbol{\beta}) \\
 &= \sum_{i=1}^N (1 - \{(2y_i - 1) 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) + 1 - y_i\}) \\
 &= \sum_{i=1}^N \{y_i - (2y_i - 1) 1(\mathbf{x}'_i \boldsymbol{\beta} > 0)\} \\
 &= \sum_{i=1}^N \begin{cases} 1 - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) & \text{if } y_i = 1 \\ 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) & \text{if } y_i = 0 \end{cases} \\
 &= \sum_{i=1}^N \begin{cases} y_i - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) & \text{if } y_i = 1 \\ -(y_i - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0)) & \text{if } y_i = 0 \end{cases} \\
 &= \sum_{i=1}^N \begin{cases} |y_i - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0)| & \text{if } y_i = 1 \text{ as } 1 - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) \geq 0 \\ |y_i - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0)| & \text{if } y_i = 0 \text{ as } 0 - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0) \leq 0 \end{cases} \\
 &= \sum_{i=1}^N |y_i - 1(\mathbf{x}'_i \boldsymbol{\beta} > 0)| \\
 &= \sum_{i=1}^N |y_i - \text{Median}(y_i | \mathbf{x}_i)|
 \end{aligned}$$



The last line since

$$\begin{aligned}
\text{Median}(y_i|\mathbf{x}_i) &= \text{Median}(1(y_i^* > 0)|\mathbf{x}_i) \\
&= 1(\text{Median}(y_i^*|\mathbf{x}_i) > 0) \\
&= 1(\text{Median}(\mathbf{x}_i'\boldsymbol{\beta} + u_i|\mathbf{x}_i) > 0) \\
&= 1(\underbrace{\text{Median}(\mathbf{x}_i'\boldsymbol{\beta}|\mathbf{x}_i) + \text{Median}(u_i|\mathbf{x}_i)}_0 > 0) \quad \text{assume 0} \\
&= 1(\mathbf{x}_i'\boldsymbol{\beta} > 0)
\end{aligned}$$

Assume  $\text{Median}(u_i|\mathbf{x}_i) = 0$ ,  $\hat{\boldsymbol{\beta}}$  is consistent but  $N^{1/3}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$  converges in distribution to non-normal distribution.

## 7.2 Smooth Maximum Score Estimation (Horowitz, 1992)

$$\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} \sum_{i=1}^N (2y_i - 1)K(\mathbf{x}_i'\boldsymbol{\beta}/h_N)$$

## 7.3 Semi-parametric MLE (Klein & Spady, 1993)

$$\sum_{i=1}^N \underbrace{\left\{ \frac{\hat{F}'(\mathbf{x}_i'\hat{\boldsymbol{\beta}})}{\hat{F}(\mathbf{x}_i'\hat{\boldsymbol{\beta}})(1 - \hat{F}(\mathbf{x}_i'\hat{\boldsymbol{\beta}}))} \right\}}_{\hat{w}_i} [y_i - \hat{F}(\mathbf{x}_i'\hat{\boldsymbol{\beta}})]\mathbf{x}_i = \mathbf{0}$$

Initialize  $\boldsymbol{\beta}^{(1)}$ , estimate  $F^{(1)}$  by kernel estimation.

Given  $F^{(1)}$  and  $\boldsymbol{\beta}^{(1)}$ , estimate  $\boldsymbol{\beta}^{(2)}$  by gradient descent method i.e.,  $\boldsymbol{\beta}^{(2)} = \boldsymbol{\beta}^{(1)} + \mathbf{A}_N \frac{\partial Q_N(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \big|_{\boldsymbol{\beta}^{(1)}}$

Given  $\boldsymbol{\beta}^{(2)}$ , estimate  $F^{(2)}$ . Repeat until convergence of  $\boldsymbol{\beta}$

## 8 References

- Cameron, A. C., & Trivedi, P. K. (2005). *Microeconometrics: Methods and Applications*.
- Davidson, R., & MacKinnon, J. G. (1993). *Estimation and Inference in Econometrics*.
- Ruud, P. A. (2000). *An Introduction to Classical Econometric Theory*.