

# Probit y Logit Binomial

## Fundamentos de Econometría

Juan Palomino<sup>1</sup>

<sup>1</sup>Magister en Economía Aplicada con Mención Estudios Regionales  
juan.palominoh@pucp.pe

Departamento de Economía



- 1 Variable Dependiente Binaria
- 2 Modelo de Probabilidad Lineal
- 3 Modelo de Probabilidad No Lineal
  - Enfoque de Variable Latente
  - Probit y Logit Binomial
  - Estimación Máxima Verosimilitud
  - Identificación
  - Efectos Marginales
- 4 Bondades de Ajuste
  - $R^2_{McFadden}$
  - Criterio de Información de Akaike (AIC)
  - Criterio de Información de Bayes (BIC)

- La variable dependiente es la siguiente:

$$y_i = \begin{cases} 1 & \text{si algún evento ocurre} \\ 0 & \text{si el evento no ocurre} \end{cases}$$

- Algunos ejemplos son:

- ▶ ¿Ha migrado usted en los últimos cinco años?
- ▶ ¿Es un consumidor más probable a comprar la misma marca o intentar con una nueva?
- ▶ ¿Usted se ha enfermado en las últimas dos semanas?

- Entonces, la pregunta es: **¿Cómo estimar un modelo con variable dependiente binaria?**

El primer enfoque es aplicar MCO como si fuera la **variable dependiente continua**.

## ¿Qué pasa si estimamos por MCO... ?

- El **modelo de probabilidad lineal** es el modelo de regresión lineal para una variable dependiente binaria. El modelo estructural es:

$$y_i = x_i' \beta + \varepsilon_i$$

donde:

$$y_i = \begin{cases} 1 & \text{si algún evento ocurre} \\ 0 & \text{si el evento no ocurre} \end{cases}$$

- Cuando  $y$  es una variable aleatoria binaria, entonces:

$$E(y_i|x_i) = [1 \times \text{Pr}(y_i = 1|x_i)] + [0 \times \text{Pr}(y_i = 0|x_i)] = \text{Pr}(y_i = 1|x_i) = x_i' \beta$$

# Modelos de Probabilidad Lineal (MPL)

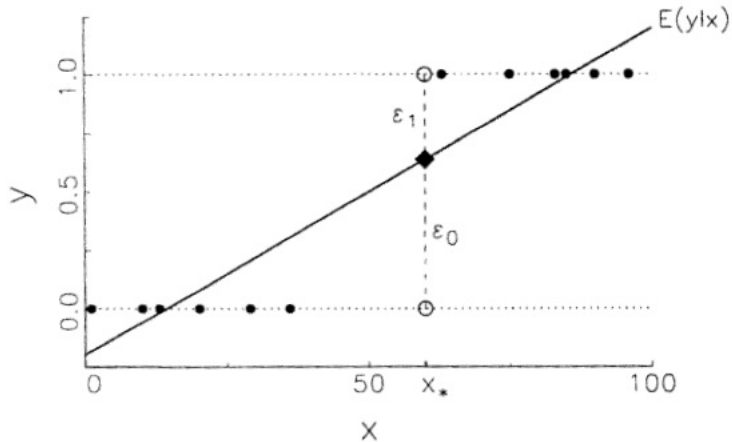


Figura: Modelo de Probabilidad Lineal

- **Heterocedasticidad:** Si una variable aleatoria binaria tiene media  $\mu$ , entonces su varianza es  $\mu(1 - \mu)$ . Luego:

$$Var(y_i|x_i) = Pr(y_i = 1|x_i)[1 - Pr(y_i = 1|x_i)] = x_i'\beta(1 - x_i'\beta)$$

lo que implica que la varianza de los errores depende de las  $x$ 's y no es constante.

- **Predictores sin sentido:** Los predictores del MLP predicen valores de  $y$  que son negativos o mayores que 1.
- **Forma Funcional:** Ya que el modelo es lineal, un incremento de una unidad en  $x_k$  resulta un cambio de  $\beta_k$  en la probabilidad de un evento. El aumento es el mismo independientemente del valor actual de  $x$ .

- Considerar las siguientes dos especificaciones:

$$y_i = \alpha + \beta x_i + \delta d_i + \varepsilon_i$$

$$y_i = F(\alpha + \beta x_i + \delta d_i)$$

donde  $d_i$  es una variable dummy.

- El cambio discreto en  $y$  cuando  $d$  cambia de 0 a 1, manteniendo  $x$  constante como:

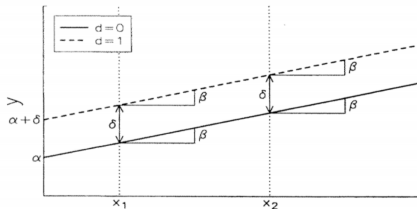
$$\frac{\Delta y}{\Delta d} = (\alpha + \beta x_i + \delta, 1) - (\alpha + \beta x_i + \delta, 0) = \delta$$

- Para nuestra segunda función el cambio discreto es...?



# Problemas con MPL

Panel A: Linear Model



Panel B: Nonlinear Model

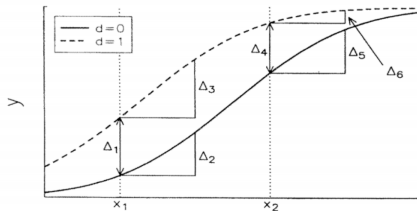


Figura: Efectos en Modelo Lineal y No Lineal

## 1 Variable Dependiente Binaria

## 2 Modelo de Probabilidad Lineal

## 3 Modelo de Probabilidad No Lineal

- Enfoque de Variable Latente
- Probit y Logit Binomial
- Estimación Máxima Verosimilitud
- Identificación
- Efectos Marginales

## 4 Bondades de Ajuste

- $R^2_{McFadden}$
- Criterio de Información de Akaike (AIC)
- Criterio de Información de Bayes (BIC)

- La variable latente  $y^*$  es asumida a ser linealmente relacionada a las variables observadas  $x$  a través del modelo estructural:

$$y_i^* = x_i' \beta + \varepsilon_i$$

- La variable latente  $y^*$  está vinculada a la variable binaria observada  $y$  por la ecuación de medida:

$$y_i = \begin{cases} 1 & \text{si } y_i^* > \tau \\ 0 & \text{si } y_i^* \leq \tau \end{cases}$$

# Enfoque de Variable Latente

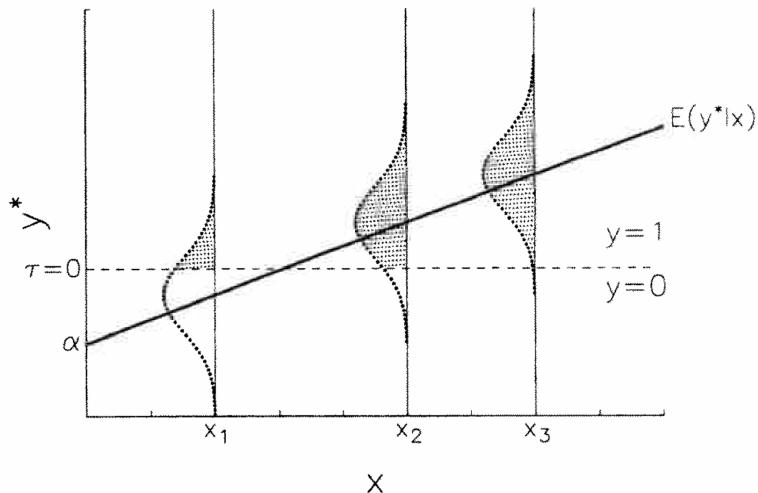


Figura: Distribución de  $y^*$  dado  $x$  en el modelo binario

$$\begin{aligned}Pr(y_i = 1|x_i) &= Pr(y_i^* > 0|x_i) \\&= Pr(x_i'\beta + \varepsilon_i > 0|x_i) \\&= Pr(\varepsilon_i > -x_i'\beta|x_i) \\&= 1 - Pr(\varepsilon_i \leq -x_i'\beta|x_i) \quad \because Pr(X > x) = 1 - Pr(X \leq x) \\&= Pr(\varepsilon_i \leq x_i'\beta|x_i) \text{ por simetría} \\&= F(x_i'\beta) \\&= \int_{-\infty}^{x_i'\beta} f(\varepsilon_i) d\varepsilon_i\end{aligned}$$

Como siempre asumimos que  $E(\varepsilon | X) = 0$

# Enfoque de Variable Latente

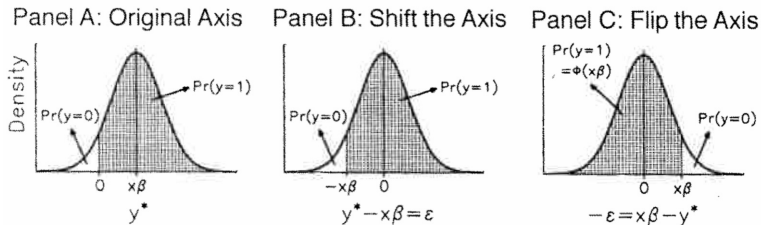


Figura: Computar  $Pr(y = 1 | x)$  en el modelo binario

# Enfoque de Variable Latente

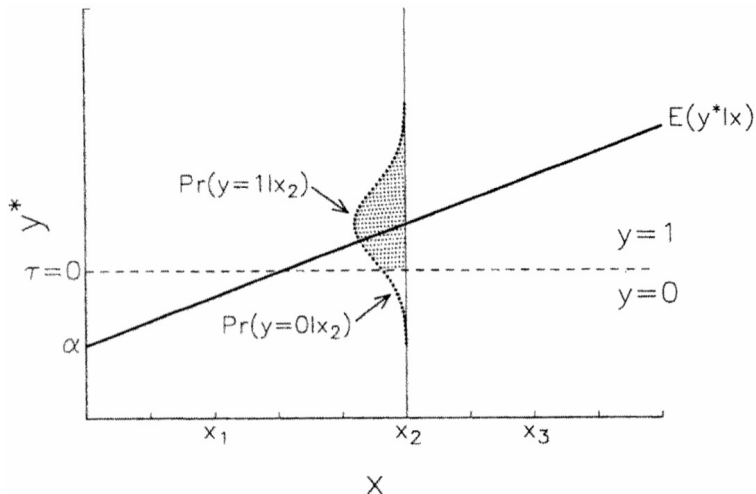


Figura: Probabilidad de valores observados en el modelo binario

## 1 Variable Dependiente Binaria

## 2 Modelo de Probabilidad Lineal

## 3 Modelo de Probabilidad No Lineal

- Enfoque de Variable Latente
- **Probit y Logit Binomial**
- Estimación Máxima Verosimilitud
- Identificación
- Efectos Marginales

## 4 Bondades de Ajuste

- $R^2_{McFadden}$
- Criterio de Información de Akaike (AIC)
- Criterio de Información de Bayes (BIC)



- En el modelo probit, el término de error  $\varepsilon$  se distribuye como normal con  $E(\varepsilon|X) = 0$  y  $Var(\varepsilon|X) = 1$ .
- Entonces, la pdf es:

$$\phi(\varepsilon) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\varepsilon^2}{2}\right)$$

- y la función de distribución acumulada (cdf) es:

$$\Phi(\varepsilon) = \int_{-\infty}^{\varepsilon} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt$$

- En el modelo logístico, los errores son asumidos a tener una distribución logística standard con media 0 y varianza  $\pi^2/3$ .
- Se elige esta variación inusual porque da como resultado una ecuación particularmente simple para el pdf

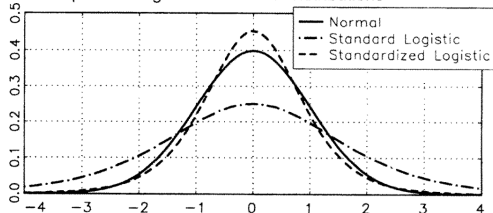
$$\lambda(\varepsilon) = \frac{\exp(\varepsilon)}{[1 + \exp(\varepsilon)]^2}$$

- y una ecuación más simple para el cdf:

$$\Lambda(\varepsilon) = \frac{\exp(\varepsilon)}{1 + \exp(\varepsilon)}$$

# Distribución Logística y Normal

Panel A: pdf's for logistic and normal distributions



Panel B: cdf's for logistic and normal distributions

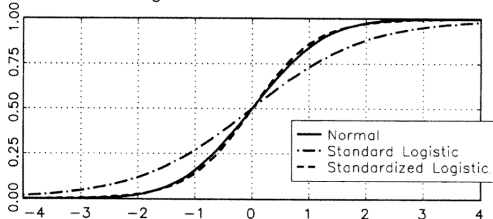


Figura: Distribución Logística y Normal

# Modelos de Variable Dependiente Binaria

- Si los  $\varepsilon_i$ 's son independientemente y normalmente distribuido,  $\varepsilon_i \sim N(0, \sigma^2)$ . Entonces:

$$\begin{aligned} Pr(y_i = 1 | x_i) &= Pr\left(\frac{\varepsilon_i}{\sigma} > -\frac{x_i' \beta}{\sigma} \mid x_i\right) \\ &= 1 - \Phi\left(-\frac{x_i' \beta}{\sigma}\right) \\ &= \Phi\left(\frac{x_i' \beta}{\sigma}\right) \text{ por simetría de distribución normal standard} \end{aligned}$$

- La probabilidad depende de  $\beta$  y  $\sigma$ , pero solo la fracción  $\beta/\sigma$  es identificada, pero no los parámetros  $\beta$  y  $\sigma$
- Por ejemplo, si  $\beta$  y  $\sigma$  son multiplicados por una constante  $c$ , entonces la probabilidad permanece sin cambios.
- Tipicamente, sea  $\sigma = 1$  como normalización

## 1 Variable Dependiente Binaria

## 2 Modelo de Probabilidad Lineal

## 3 Modelo de Probabilidad No Lineal

- Enfoque de Variable Latente
- Probit y Logit Binomial
- **Estimación Máxima Verosimilitud**
- Identificación
- Efectos Marginales

## 4 Bondades de Ajuste

- $R^2_{McFadden}$
- Criterio de Información de Akaike (AIC)
- Criterio de Información de Bayes (BIC)

# Estimación Máxima Verosimilitud

- El resultado es distribuido Bernoulli, la distribución binomial con solo una prueba. Una notación compacta muy conveniente para la densidad de  $y_i$ , o más formalmente su función de probabilidad, es:

$$f(y_i|x_i) = P_i^{y_i}(1 - P_i)^{1-y_i}, \quad y_i = 1, 0$$

donde  $P_i = F(x_i'\beta)$ . Este produce probabilidades  $P_i$  y  $(1 - P_i)$  ya que  $f(1) = P^1(1 - P)^0 = P$  y  $f(0) = P^0(1 - P)^1 = 1 - P$ . Suponiendo que cada probabilidad es independiente de la otra, la función de probabilidad conjunta (o función de verosimilitud) es:

$$Pr(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n | X) = \prod_{i=1}^N f(y_i|x_i)$$

- La función de verosimilitud para una muestra de  $n$  observaciones pueden ser escritos como:

$$L(\underbrace{\beta}_{data} | y, X) = \prod_{i=1}^N [F(x_i'\beta)]^{y_i} [1 - F(x_i'\beta)]^{1-y_i}$$

# Función Log-Verosimilitud

- Tomando logs, obtenemos la función de Log-Likelihood, que debe ser maximizado:

$$\begin{aligned}\log L(\beta|data) &= \log\left(\prod_{i=1}^N [F(x_i'\beta)]^{y_i} [1 - F(x_i'\beta)]^{1-y_i}\right) \\ &= \sum_{i=1}^n \{\log([F(x_i'\beta)]^{y_i}) + \log([1 - F(x_i'\beta)]^{1-y_i})\} \\ &= \sum_{i=1}^n \{y_i \log([F(x_i'\beta)]) + (1 - y_i) \log([1 - F(x_i'\beta)])\}\end{aligned}$$

## Truco útil

Si la distribución es simétrica, como la normal y logística, entonces

$1 - F(x_i'\beta) = F(-x_i'\beta)$ . Sea  $q_i = 2y_i - 1$ . Entonces:

$$\log L(\beta|data) = \sum_{i=1}^n \log F(q_i x_i' \beta)$$

- Las condiciones de primer orden:

$$\begin{aligned}\underbrace{\frac{\partial \log L(\beta | \text{data})}{\partial \beta}}_{(K \times 1)} &= \frac{\partial}{\partial \beta} \sum_{i=1}^N \{y_i \log[F(x'_i \beta)] + (1 - y_i) \log[1 - F(x'_i \beta)]\} \\&= \sum_{i=1}^n \left\{ y_i \frac{f(x'_i \beta) x_i}{F(x'_i \beta)} + (1 - y_i) \frac{-f(x'_i \beta) x_i}{1 - F(x'_i \beta)} \right\} \frac{\partial F(x'_i \beta)}{\partial \beta} = f(x'_i \beta) x_i \\&= \sum_{i=1}^n \left\{ \frac{y_i}{F(x'_i \beta)} - \frac{(1 - y_i)}{1 - F(x'_i \beta)} \right\} f(x'_i \beta) x_i \\&= \sum_{i=1}^n \left\{ \frac{y_i(1 - F(x'_i \beta)) - (1 - y_i)F(x'_i \beta)}{F(x'_i \beta)[1 - F(x'_i \beta)]} \right\} f(x'_i \beta) x_i \\&= \sum_{i=1}^n \left\{ \frac{y_i - F(x'_i \beta)}{F(x'_i \beta)[1 - F(x'_i \beta)]} \right\} f(x'_i \beta) x_i\end{aligned}$$



- Asimismo, el estimador MV  $\hat{\beta}_{MV}$  es dado por la solución de:

$$\underbrace{\sum_{i=1}^n \left\{ \frac{y_i - F(x_i' \beta)}{F(x_i' \beta) [1 - F(x_i' \beta)]} \right\}}_{(1 \times 1)} \underbrace{f(x_i' \beta)}_{(1 \times 1)} \underbrace{x_i}_{(K \times 1)} = \underbrace{0}_{K \times 1}$$

Esta ecuación no tiene una solución analítica de  $\hat{\beta}_{MV}$  y tenemos que resolverlo por métodos computacionales.

- Usando nuestro truco previo, tenemos:

$$\frac{\partial \log L(\beta)}{\partial \beta} = \sum_{i=1}^n \left[ \frac{q_i f(q_i x_i' \beta)}{F(x_i' \beta)} \right] x_i = \sum_{i=1}^n \lambda_i x_i = 0$$

- La Hesiana es:

$$H = \frac{\partial^2 \log L}{\partial \beta \partial \beta'} = \sum_{i=1}^N \left[ \frac{q_i^2 f'(q_i x_i' \beta)}{F(q_i x_i' \beta)} - \lambda_i^2 \right] x_i x_i'$$

- Nuevamente, tenga en cuenta que  $F(\cdot)$  depende de la distribución del término de error. Si asumimos un modelo probit, entonces:

$$F(\varepsilon) = \Phi(\varepsilon) = \int_{-\infty}^{\varepsilon} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right) dt$$

$$f(\varepsilon) = \phi(\varepsilon) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\varepsilon^2}{2}\right)$$

$$f'(\varepsilon) = -\varepsilon \phi(\varepsilon)$$

- y para el modelo logit, tenemos:

$$F(\varepsilon) = \Lambda(\varepsilon)$$

$$f(\varepsilon) = \Lambda(\varepsilon)[1 - \Lambda(\varepsilon)]$$

$$f'(\varepsilon) = \Lambda(\varepsilon)[1 - \Lambda(\varepsilon)][1 - 2\Lambda(\varepsilon)]$$

- El hessiano siempre es definido negativo, por lo que la probabilidad logarítmica es globalmente cóncava.
- El método de Newton generalmente convergerá al máximo de la probabilidad de registro en solo unas pocas iteraciones, a menos que los datos estén especialmente mal condicionados.

## 1 Variable Dependiente Binaria

## 2 Modelo de Probabilidad Lineal

## 3 Modelo de Probabilidad No Lineal

- Enfoque de Variable Latente
- Probit y Logit Binomial
- Estimación Máxima Verosimilitud
- **Identificación**
- Efectos Marginales

## 4 Bondades de Ajuste

- $R^2_{McFadden}$
- Criterio de Información de Akaike (AIC)
- Criterio de Información de Bayes (BIC)

- Asumir que el término de error tiene media diferente de cero  $\theta$ . Entonces:

$$\begin{aligned}Pr(y_i | x_i) &= Pr(\varepsilon_i \leq \alpha + x_i' \beta) \\&= Pr(\varepsilon_i - \theta \leq (\alpha - \theta) + x_i' \beta) \\&= Pr(\varepsilon_i^* \leq \alpha^* + x_i' + x_i' \beta)\end{aligned}$$

- Mientras el modelo de elección binaria contenga un término constante, no hay pérdida de generalidad, suponiendo que la media del término aleatorio sea cero.
- Una media distinta de cero desaparecería en el término constante de la función de utilidad.

## 1 Variable Dependiente Binaria

## 2 Modelo de Probabilidad Lineal

## 3 Modelo de Probabilidad No Lineal

- Enfoque de Variable Latente
- Probit y Logit Binomial
- Estimación Máxima Verosimilitud
- Identificación
- Efectos Marginales

## 4 Bondades de Ajuste

- $R^2_{McFadden}$
- Criterio de Información de Akaike (AIC)
- Criterio de Información de Bayes (BIC)

# Efectos Marginales

- Recordar que:

$$E(y_i|x_i) = [1 \times \text{Pr}(y_i = 1|x_i)] + [0 \times \text{Pr}(y_i = 0|x_i)] = \text{Pr}(y_i = 1|x_i) = F(x_i'\beta)$$

El efecto marginal es dado por:

$$\underbrace{\frac{\partial E(y_i|x_i)}{\partial x_i}}_{(K \times 1)} = \left[ \frac{dF(x_i'\beta)}{d(x_i'\beta)} \right] \beta = \underbrace{f(x_i'\beta)}_{(1 \times 1)} \underbrace{\beta}_{(K \times 1)}$$

donde  $f(\cdot)$  es la función de densidad de probabilidad.

- Entonces:

$$\text{Probit} \Rightarrow \frac{\partial E(y_i|x_i)}{\partial x_i} = \phi(x_i'\beta_i)\beta$$

$$\text{Logit} \Rightarrow \frac{\partial E(y_i|x_i)}{\partial x_i} = \Lambda(x_i'\beta_i)[1 - \Lambda(x_i'\beta_i)]\beta$$

- Algunos aspectos son importantes para comentar:
  - ▶ Efectos marginales pueden variar con los valores de  $x$
  - ▶ Practicas comunes:
    - ★ Calcular efectos marginales en el promedio de las variables
    - ★ Calcular efectos marginales en valores específicos
    - ★ Evaluar efectos marginales en cada observación y usar el promedio muestral de los efectos marginales individuales.



- El efecto marginal apropiado para una variable independiente binaria, es decir,  $d$ , sería:

$$ME = [Pr(y_i = 1 | \bar{x}_{(d)}, d_i = 1)] - [Pr(y_i = 1 | \bar{x}_{(d)}, d_i = 0)]$$

donde  $\bar{x}_{(d)}$  denota los promedios de todas las otras variables en el modelo.

- Es común reportar elasticidades de probabilidades, en vez de derivadas. Estos son computados como:

$$\begin{aligned}\epsilon_{i,k} &= \frac{\partial \ln \Pr(y_i = 1|x)}{\partial \ln x_{i,k}} \\ &= \frac{\partial \ln \Pr(y_i = 1|x)}{\partial x_{i,k}} \frac{x_{i,k}}{\Pr(y_i = 1|x)}\end{aligned}$$

- Como es una relación de cambios porcentuales, la elasticidad no es útil para las variables dummies.

- Se han realizado varias sugerencias sobre cómo evaluar la calidad general de un modelo de respuesta binaria.
  - ▶ Primer enfoque: imitar la medida  $R^2$
  - ▶ Evaluar el rendimiento predictivo del modelo
- $R^2$  no son directamente aplicables en modelos no lineales, como los modelos de respuesta binaria, ya que no tenemos un resultado de descomposición de la varianza adecuado
- Se sugiere la medida del pseudo  $R^2$

- 1 Variable Dependiente Binaria
- 2 Modelo de Probabilidad Lineal
- 3 Modelo de Probabilidad No Lineal
  - Enfoque de Variable Latente
  - Probit y Logit Binomial
  - Estimación Máxima Verosimilitud
  - Identificación
  - Efectos Marginales
- 4 Bondades de Ajuste
  - $R^2_{McFadden}$
  - Criterio de Información de Akaike (AIC)
  - Criterio de Información de Bayes (BIC)

- Sea
  - ▶  $\log L(\hat{\beta}_r)$  el valor de la función de log-verosimilitud maximizada en el modelo de solo constante.
  - ▶  $\log L(\hat{\beta}_u)$  : es el valor de log-verosimilitud maximizada en el modelo completo.
- Notar que el valor de las funciones de log-verosimilitud es siempre negativo, entonces:

$$\log L(\hat{\beta}_u) \geq \log L(\hat{\beta}_r) \Rightarrow \left| \log L(\hat{\beta}_u) \right| \leq \left| \log L(\hat{\beta}_r) \right|$$

- De tal manera que:

$$0 \leq 1 - \frac{\log L(\hat{\beta}_u)}{\log L(\hat{\beta}_r)} = R_{McFadden}^2 \leq 1$$

- El  $R^2$  McFadden podría ser cero si el modelo completo no tiene poder explicativo.

- Otra medida es el McKelvey y Zavoina (1975). Se base en el modelo lineal latente  $y_i^* = x_i' \beta + \varepsilon_i$ . En particular, si  $\hat{y}_i^* = x_i' \hat{\beta}$ , entonces escribimos:

$$R_{MZ}^2 = \frac{SSE^*}{SSR^* + SSE^*} = \frac{\sum_{i=1} (\hat{y}_i^* - \bar{\hat{y}}^*)^2}{n\sigma^2 + \sum_{i=1} (\hat{y}_i^* - \bar{\hat{y}}^*)^2}$$

- donde  $SSE^*$  denota la suma de cuadrados explicada, y  $SSR^*$  denota la suma de residuos al cuadrado del modelo latente.
  - ▶ Para probit  $\sigma^2 = 1$
  - ▶ Para logit  $\sigma^2 = \pi^2/3$

- 1 Variable Dependiente Binaria
- 2 Modelo de Probabilidad Lineal
- 3 Modelo de Probabilidad No Lineal
  - Enfoque de Variable Latente
  - Probit y Logit Binomial
  - Estimación Máxima Verosimilitud
  - Identificación
  - Efectos Marginales
- 4 Bondades de Ajuste
  - $R^2_{McFadden}$
  - Criterio de Información de Akaike (AIC)
  - Criterio de Información de Bayes (BIC)

## Criterio de Información de Akaike (AIC)

El criterio de Información de Akaike es definido como:

$$AIC = \frac{-2\log\hat{L} + 2K}{n}$$

donde  $\log\hat{L}$  es la probabilidad del modelo y  $K$  es la cantidad de parámetros en el modelo.

- Valores grandes de  $\log\hat{L}$  indica un mejor ajuste.
- $-2\log\hat{L}$  oscila entre 0 y  $+\infty$  con valores más pequeños que indican un mejor ajuste.
- A medida que  $K$  aumenta,  $-2\log\hat{L}$  se convierte más pequeño ya que **más parámetros hacen que lo que se observa sea más probable**.
- $2K$  es agregado como un penal.
- Todo lo demás es igual, los valores más pequeños sugieren un modelo de mejor ajuste.
- Usar para comparar modelos en diferentes muestras o para comparar modelos no anidados.



- 1 Variable Dependiente Binaria
- 2 Modelo de Probabilidad Lineal
- 3 Modelo de Probabilidad No Lineal
  - Enfoque de Variable Latente
  - Probit y Logit Binomial
  - Estimación Máxima Verosimilitud
  - Identificación
  - Efectos Marginales
- 4 Bondades de Ajuste
  - $R^2_{McFadden}$
  - Criterio de Información de Akaike (AIC)
  - Criterio de Información de Bayes (BIC)

## Criterio de Información de Bayes (BIC)

El criterio de Información BIC es definido como:

$$BIC = \frac{-2\log\hat{L} + K\log(n)}{n}$$

donde  $\log\hat{L}$  es la probabilidad del modelo y  $K$  es la número de parámetros en el modelo y  $n$  es el número de individuos.

- Es posible aumentar la probabilidad agregando parámetros, pero hacerlo puede resultar en un ajuste excesivo. Tanto BIC como AIC resuelven este problema al introducir un término de penalización por la cantidad de parámetros en el modelo; el término de penalización es mayor en el BIC que el AIC.