

# Probit y logit

Gabriel V. Montes-Rojas

# Modelos de respuesta binaria

Supongamos que la variable dependiente es una variable dummy que toma valores 0 y 1. En general, se usa la denominación de la literatura experimental éxito (success) para 1 y fracaso (failure) para 0.

Ejemplos:

- *Participación en la fuerza laboral*: Supongamos que queremos modelar el efecto de ciertas variables en la participación en la fuerza laboral o sobre el desempleo. Supongamos una variable variable  $inlf$  (in labor force) que toma 1 si el individuo trabaja y 0 si no lo hace.
- *Bancarrota*: Supongamos un modelo que estima el efecto de ciertas variables sobre la probabilidad de que una firma o un individuo se declare en bancarrota.
- *Default de la deuda soberana*: Supongamos un modelo que estima el efecto de las variables macroeconómicas (déficit, deuda/PBI) sobre la probabilidad que un país se declare en default.

# Modelo de probabilidad lineal

Sea  $y = 0, 1$  la variable dependiente. Una opción de primera instancia es usar el **modelo de probabilidad lineal** que consiste en correr una regresión MCO:

$$y = \mathbf{x}\beta + u.$$

¿Cuál es la interpretación de  $\beta$ ?

$$E[y|\mathbf{x}] = \mathbf{x}\beta = P[y = 1|\mathbf{x}].$$

Entonces  $\beta_1 = \frac{\partial P[y=1|\mathbf{x}]}{\partial x}$ . En otras palabras:  $\beta$  nos da *el efecto marginal sobre la probabilidad de tener un éxito ( $y = 1$ )*.

Nota: Se puede pensar a  $y$  como una variable Bernoulli que toma valor 1 con probabilidad  $p$  y valor 0 con probabilidad  $1 - p$ . Entonces,  $E[y] = p$ ,  $Var[y] = p(1 - p)$ .

# Modelo de probabilidad lineal

Sin embargo, este modelo tiene ciertos problemas e inconsistencias:

- 1 **Valor predicho:** El modelo no garantiza que  $0 \leq \hat{y} \leq 1$ , lo cual es un problema porque estamos hablando de una probabilidad, ya que  $\hat{y} = P[\widehat{y = 1}|\mathbf{x}]$ . Notar que sin embargo  $0 \leq \bar{y} = \hat{y}(\bar{\mathbf{x}}) = \bar{\mathbf{x}}\hat{\beta} \leq 1$ .
- 2 **Heteroscedasticidad:**  $Var(y|\mathbf{x}) = P[y = 1|\mathbf{x}] \times (1 - P[y = 1|\mathbf{x}])$ . [Ejercicio: Proponer un método para obtener errores estándar correctos usando el modelo de probabilidad lineal.]

# Modelos logit y probit

- Supongamos el siguiente modelo de variable latente:

$$y^* = \mathbf{x}\beta + e$$

donde  $e$  es una variable aleatoria continua independiente de  $\mathbf{x}$  y la distribución de  $e$  es simétrica en cero.

- Pero no observamos  $y^*$  (por eso es latente), sino

$$y = \mathbf{1}[y^* > 0] = \mathbf{1}[e > -(\mathbf{x}\beta)],$$

donde  $\mathbf{1}[\cdot]$  es la función indicador que toma valor 1 si el argumento en  $[\cdot]$  es verdadero, 0 si no lo es.

- Entonces,

$$P(y = 1|\mathbf{x}) = P(y^* > 0|\mathbf{x}) = P(e > -\mathbf{x}\beta|\mathbf{x}) = 1 - F(-\mathbf{x}\beta) = F(\mathbf{x}\beta)$$

donde  $F$  es la función de distribución acumulada (*cdf*) de  $e$ .

- Diferentes supuestos sobre  $F$  determinan distintos modelos:
  - Si  $F$  es normal, entonces tenemos el modelo **probit**;
  - Si  $F$  es logística, entonces tenemos el modelo **logit**.

# Modelo probit

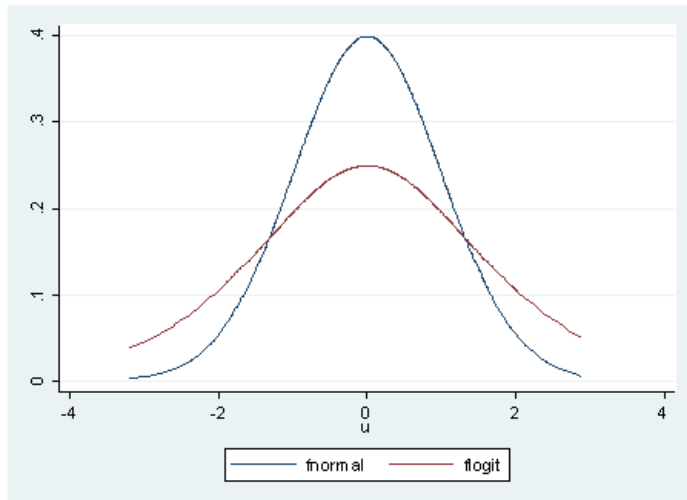
- Supongamos  $e \sim N(0, \sigma_e^2)$ , modelo **probit**.
- En este caso  $F(z) = P[e \leq z] = \int_{-\infty}^z \phi(v) dv = \Phi(z)$  donde  $\phi$  es la función de densidad normal,  $\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$ , y  $\Phi$  la *cdf* normal.
- $P[y = 1|\mathbf{x}] = P[e > -(\mathbf{x}\beta)] = 1 - F(-\mathbf{x}\beta) = F(\mathbf{x}\beta) = \Phi(\mathbf{x}\beta)$
- $P[y = 0|\mathbf{x}] = P[e \leq -(\mathbf{x}\beta)] = F(-\mathbf{x}\beta) = 1 - F(\mathbf{x}\beta) = 1 - \Phi(\mathbf{x}\beta)$

# Modelo logit

- Supongamos  $e \sim \text{logística}$ , modelo **logit**.
- En este caso:  $F(z) = P[e \leq z] = \frac{\exp(z)}{1+\exp(z)} = \Lambda(z)$ , donde  $\Lambda$  es la *cdf* de una variable aleatoria logística. Notar que la *pdf* es  $\lambda(z) = \frac{\exp(z)}{(1+\exp(z))^2}$ .
- $P[y = 1|\mathbf{x}] = P[e > -(\mathbf{x}\boldsymbol{\beta})] = 1 - F(-\mathbf{x}\boldsymbol{\beta}) = F(\mathbf{x}\boldsymbol{\beta}) = \Lambda(\mathbf{x}\boldsymbol{\beta})$
- $P[y = 0|\mathbf{x}] = P[e \leq -(\mathbf{x}\boldsymbol{\beta})] = F(-\mathbf{x}\boldsymbol{\beta}) = 1 - F(\mathbf{x}\boldsymbol{\beta}) = 1 - \Lambda(\mathbf{x}\boldsymbol{\beta})$

# Probit vs. Logit

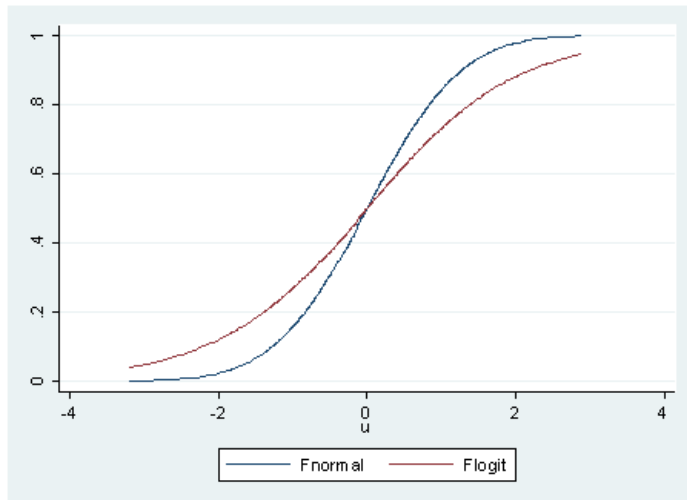
## Funciones de densidad





# Probit vs. Logit

## Funciones de distribución



# Modelos logit y probit

## ¿Cómo interpretar los coeficientes?

- Si  $\mathbf{x}$  es una variable continua,

$$\frac{\partial P[y = 1 | \mathbf{x}]}{\partial \mathbf{x}} = \frac{\partial F(\mathbf{x}\beta)}{\partial \mathbf{x}} = f(\mathbf{x}\beta)\beta.$$

Entonces  $\beta \neq \frac{\partial P[y=1|\mathbf{x}]}{\partial \mathbf{x}}$  y no se pueden interpretar los coeficientes de los modelos probit y logit directamente.

- Necesitamos  $f(z) = \frac{dF}{dz}(z)$ , la densidad asumida de  $e$  en el modelo de variable latente. Notar que cada individuo  $i$  tiene un valor marginal potencialmente diferente:  
 $\frac{\partial P[y=1|\mathbf{x}_i]}{\partial \mathbf{x}} = f(\mathbf{x}_i\beta)$ .
- Sí se puede interpretar el signo:  $sign(\beta) = sign\left(\frac{\partial P[y=1|\mathbf{x}]}{\partial \mathbf{x}}\right)$ .
- Pero, ¿qué valor de  $\mathbf{x}$  usar para computar  $f(\mathbf{x}\beta)$ ? En general se usa  $\mathbf{x} = \bar{\mathbf{x}}$ .
- También podemos computar el efecto marginal promedio como  $E_{\mathbf{x}}[f(\mathbf{x}\beta)\beta]$ .
- Si  $x_K$  es una variable dummy, entonces se mide el efecto de un cambio de 0 a 1:

$$F(\beta_1 + \beta_2\bar{x}_2 + \dots + \beta_{K-1}\bar{x}_{K-1} + \beta_K) - F(\beta_1 + \beta_2\bar{x}_2 + \dots + \beta_{K-1}\bar{x}_{K-1}).$$

## Comparación de modelos

- Modelo de probabilidad lineal:  $F(x) = x$
- Modelo probit:  $F(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du = \Phi(x)$
- Modelo logit:  $F(x) = \frac{u^x}{1+u^x} = \Lambda(x)$
- $\hat{\beta}_\Lambda \simeq 1.6\hat{\beta}_\Phi$ .
- $\hat{\beta}_{pl} \simeq 0.4\hat{\beta}_\Phi$  excepto para la constante.
- $\hat{\beta}_{pl} \simeq 0.4\hat{\beta}_\Phi + 0.5$  para la constante.

# Modelos logit

- Para los modelos logit tenemos una interpretación de los coeficientes.
- Notar que  $P[y = 1|\mathbf{x}]/P[y = 0|\mathbf{x}] = e^{\mathbf{x}\beta}$ . Esto se conoce como el ratio de odds (odds ratio) y mide cuan probable es un evento en proporción a cuan probable es otro.
- Aplicando logs,  $\log(P[y = 1|\mathbf{x}]/P[y = 0|\mathbf{x}]) = \mathbf{x}\beta$ , entonces

$$\frac{\partial \log(P[y = 1|\mathbf{x}]/P[y = 0|\mathbf{x}])}{\partial x_j} = \beta_j.$$

Los coeficientes se interpretan como el efecto porcentual sobre el ratio de odds (semi-elasticidad).

# Una introducción a máxima verosimilitud

- La variable dependiente muestral es  $\{y_i\}_{i=1}^N$ , que son 0 o 1 para cada observación.
- Si observamos un 1,  $y_i = 1$ , ¿cuál es la probabilidad asociada a tener este valor en particular en dicha observación  $i$ ? Usando el modelo de variable latente,  $y_i^* = \mathbf{x}\beta_i + e_i > 0$ , y como  $e$  se asumió logit/probit  $P[y_i = 1|\mathbf{x}_i] = F(\mathbf{x}_i\beta)$ .
- Si observamos un 0,  $y_i = 0$ , ¿cuál es la probabilidad asociada a tener este valor en particular en dicha observación  $i$ ? Usando el modelo de variable latente,  $y_i^* = \mathbf{x}\beta_i + e_i > 0$ , y como  $e$  se asumió logit/probit  $P[y_i = 0|\mathbf{x}_i] = 1 - F(\mathbf{x}_i\beta)$ .
- En general tenemos (en forma condensada),

$$P[y|\mathbf{x}] = [F(\mathbf{x}\beta)]^y [1 - F(\mathbf{x}\beta)]^{1-y}.$$

# Una introducción a máxima verosimilitud

- Si tomamos TODAS las observaciones juntas,  $\{y_i\}_{i=1}^N$ , asumiendo independencia entre las observaciones,

$$P[y_1, y_2, \dots, y_N | \mathbf{x}] = \prod_{i=1}^N P[y_i | \mathbf{x}_i] = \prod_{i=1}^N [F(\mathbf{x}_i \boldsymbol{\beta})]^{y_i} [1 - F(\mathbf{x}_i \boldsymbol{\beta})]^{1-y_i}$$

Esta es la función de verosimilitud (likelihood function).

*Dos eventos A y B son independientes si  $P[A \& B] = P[A] \times P[B]$ .*

- En general es más fácil trabajar con el log de la función de verosimilitud (log-likelihood function).

$$\mathcal{L}(\boldsymbol{\beta}) = \sum_{i=1}^N \ell_i(\boldsymbol{\beta})$$

$$\ell_i(\boldsymbol{\beta}) = \log P[y_i | \mathbf{x}_i] = y_i \log F(\mathbf{x}_i \boldsymbol{\beta}) + (1 - y_i) \log [1 - F(\mathbf{x}_i \boldsymbol{\beta})]$$

- Entonces, el estimador de máxima verosimilitud (maximum likelihood estimator, MLE) es  $\hat{\boldsymbol{\beta}}$  que maximiza  $\mathcal{L}(\boldsymbol{\beta})$ . En otras palabras, para cada valor posible de  $\boldsymbol{\beta}$ ,  $\mathcal{L}(\hat{\boldsymbol{\beta}}) \geq \mathcal{L}(\boldsymbol{\beta})$ .

# Una introducción a máxima verosimilitud

- La función score de este modelo es

$$\mathbf{s}_i(\boldsymbol{\beta}) = \frac{f(\mathbf{x}_i\boldsymbol{\beta})\mathbf{x}'_i[y_i - F(\mathbf{x}_i\boldsymbol{\beta})]}{F(\mathbf{x}_i\boldsymbol{\beta})[1 - F(\mathbf{x}_i\boldsymbol{\beta})]}$$

- Ejercicio: Probar que  $E[\mathbf{s}_i(\boldsymbol{\beta}_0)|\mathbf{x}_i] = 0$ , donde  $\boldsymbol{\beta}_0$  es el valor verdadero de  $\boldsymbol{\beta}$ .
- El hessiano es

$$\mathbf{H}_i(\boldsymbol{\beta}) = \left\{ \begin{aligned} &\frac{f'(\mathbf{x}_i\boldsymbol{\beta})[y_i - F(\mathbf{x}_i\boldsymbol{\beta})]}{F(\mathbf{x}_i\boldsymbol{\beta})[1 - F(\mathbf{x}_i\boldsymbol{\beta})]} - \frac{f^2(\mathbf{x}_i\boldsymbol{\beta})}{F(\mathbf{x}_i\boldsymbol{\beta})[1 - F(\mathbf{x}_i\boldsymbol{\beta})]} \\ &- \frac{f(\mathbf{x}_i\boldsymbol{\beta})[y_i - F(\mathbf{x}_i\boldsymbol{\beta})]}{(F(\mathbf{x}_i\boldsymbol{\beta})[1 - F(\mathbf{x}_i\boldsymbol{\beta})])^2} (f(\mathbf{x}_i\boldsymbol{\beta}) - 2F(\mathbf{x}_i\boldsymbol{\beta})f(\mathbf{x}_i\boldsymbol{\beta})) \end{aligned} \right\} (\mathbf{x}'_i\mathbf{x}_i).$$

$$-E[\mathbf{H}_i(\boldsymbol{\beta})|\mathbf{x}_i] = \mathbf{A}_i(\boldsymbol{\beta}) = \frac{f^2(\mathbf{x}_i\boldsymbol{\beta})(\mathbf{x}'_i\mathbf{x}_i)}{F(\mathbf{x}_i\boldsymbol{\beta})[1 - F(\mathbf{x}_i\boldsymbol{\beta})]}$$

# Una introducción a máxima verosimilitud

- Por normalidad asintótica de máxima verosimilitud

$$\sqrt{N}(\hat{\beta} - \beta_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{A}_0^{-1}),$$

donde  $\mathbf{A}_0 = -E[\mathbf{H}(\beta_0)]$ .

- Con los elementos de la diapositiva anterior tenemos todo lo que necesitamos para armar la matriz de varianzas y covarianzas.
- Probar que  $-\mathbf{A}_0 = -E[\mathbf{H}(\beta_0)] = \mathbf{B}_0 = E[\mathbf{s}(\beta_0)\mathbf{s}(\beta_0)']$ . ¿Por qué?
- ¿Cómo estimar la varianza en forma consistente?

$$AVar(\hat{\beta}) = \left( \sum_{i=1}^N \frac{f^2(\mathbf{x}_i; \hat{\beta})(\mathbf{x}'_i \mathbf{x}_i)}{F(\mathbf{x}_i; \hat{\beta})[1 - F(\mathbf{x}_i; \hat{\beta})]} \right)^{-1}$$



## Pseudo R cuadrado

- En los modelos de regresión, el R-cuadrado tiene una interpretación intuitiva y simple.
- Eso no ocurre por fuera de los modelos OLS...
- Una primera aproximación puede ser obtener el resultado predecido como  $\hat{y}_i = \max\{j = 0, 1 : p_j(\mathbf{x}_i, \hat{\boldsymbol{\beta}})\}$ , es decir la probabilidad que sea mayor a 0.5.  
Entonces,  $\text{pseudo} - R^2 = \frac{\sum_{i=1}^N 1[\hat{y}_i = y_i]}{N}$ .
- McFadden (1974) propone usar el log likelihood ratio index:

$$LRI = 1 - \frac{\mathcal{L}(\hat{\boldsymbol{\beta}})}{\mathcal{L}(\boldsymbol{\beta} = 0)}$$

En general LRI no se acerca a 1 cuando se incrementan la cantidad de variables. Y si es 1, en realidad es un problema ya que dice que hay una variable que predice perfectamente el resultado, ej.  $x > 0$  entonces  $y = 1$ .

# STATA

- Es útil seguir el manual de STATA (v.13):  
<http://www.stata.com/manuals13/rprobit.pdf> para probit,  
<http://www.stata.com/manuals13/rlogit.pdf> para logit.
- `probit` y `x1 x2`
- `logit` y `x1 x2`
- Dado que los coeficientes no se pueden interpretar directamente (sólo el signo), necesitamos computar los efectos marginales:
  - `dprobit` y `x1 x2`
  - `logit` y `x1 x2`
  - `margins, dydx(*)` (calculado como el promedio de las densidades)
  - `margins, dydx(*) atmeans` (calculado con la densidad en el promedio de las  $X$ )

# Ejemplos

- Considere los ejemplos de <http://fmwww.bc.edu/gstat/examples/wooldridge/wooldridge17.html>
- Descripción de la base de datos de Mroz:  
<http://fmwww.bc.edu/ec-p/data/wooldridge/mroz.des>
- Ejemplo de programación en STATA:  
<https://stats.idre.ucla.edu/stata/code/imple-linear-and-nonlinear-models-using-statas-ml-command/>