

Mô hình với Biến Phụ thuộc bị Giới hạn (Models with Limited Dependent Variables)

Lê Việt Phú
Trường Chính sách Công và Quản lý Fulbright

Ngày 10 tháng 1 năm 2018

Table of contents

- Khái niệm biến phụ thuộc không bị giới hạn và bị giới hạn
- Sử dụng hồi quy tuyến tính đối với biến phụ thuộc bị giới hạn
- Phương pháp tối đa hoá xác suất - MLE

Khái niệm biến phụ thuộc không bị giới hạn và bị giới hạn

- ▶ Các loại biến phụ thuộc trong mô hình hồi quy:
 - ▶ Liên tục và rời rạc: tăng trưởng GDP là liên tục, có thể có con số bất kỳ, ví dụ 6.1025%; số lần đi học muộn là rời rạc, ví dụ đi muộn 0, 1, 2 lần.
 - ▶ Không bị giới hạn và bị giới hạn: lợi nhuận của công ty là không giới hạn (lỗ thì nhận giá trị âm, lãi là dương); số nhân viên là bị giới hạn (bị chặn dưới, ít nhất 1 nhân viên trong một công ty).
 - ▶ Biến phụ thuộc định tính và định lượng: có hút thuốc lá hay không là biến định tính; hút bao nhiêu điếu thuốc một ngày là định lượng và bị giới hạn (ít nhất là một điếu).
- ▶ Hầu hết các biến số kinh tế đều bị giới hạn.
- ▶ Sử dụng hồi quy tuyến tính đối với dữ liệu bị giới hạn thì kết quả có thể bị sai lệch, hoặc khó giải thích ý nghĩa về mặt kinh tế.

Một số mô hình sử dụng biến phụ thuộc bị giới hạn (1)

- ▶ Mô hình xác suất xảy ra một sự kiện hay một biến cố nào đó. Ví dụ đối tượng vị thành niên hút thuốc, đi học đại học, phụ nữ dân tộc thiểu số tham gia lao động chính thức. Biến phụ thuộc là có hoặc không (mã hoá 1 cho câu trả lời có, 0 cho câu trả lời không). Đối với biến phụ thuộc định tính thì không có cách xếp hạng câu trả lời (có/không) như biến phụ thuộc định lượng (nhiều/ít).
- ▶ Mô hình xác suất có thể là đa lựa chọn thay vì hai lựa chọn, ví dụ anh/chị đến trường bằng phương tiện gì: ô-tô, xe máy, xe đạp, đi bộ.

Một số mô hình sử dụng biến phụ thuộc bị giới hạn (2)

- ▶ Mô hình số lần xảy ra một sự kiện nào đó. Ví dụ số lần một học viên MPP đi học muộn, số con trong một gia đình, số sản phẩm bị hỏng trong một ngày, số lần đi khám bệnh một năm. Biến phụ thuộc sẽ có giá trị 0 và số nguyên dương (1, 2, 3...).
- ▶ Mô hình mô tả xếp hạng của một sự kiện, ví dụ cảm quan của anh/chị về một môn học có thể là quá khó/khó/trung bình/tương đối dễ/quá dễ.
- ▶ Mô hình với biến phụ thuộc bị chặn trên hoặc dưới. Ví dụ thu nhập chỉ có thể là 0 hoặc dương; số tiền một người đã làm từ thiện trong một năm tối thiểu là 0 hoặc dương; số giờ làm việc trong một tuần không thể quá $24 \times 7 = 168$ giờ.

Tên gọi mô hình sử dụng biến phụ thuộc có giới hạn

- ▶ Mô hình xác suất (Logit, Probit, Multinomial Logit)
- ▶ Mô hình số lần xảy ra sự kiện (Poisson)
- ▶ Mô hình với biến phụ thuộc bị chặn (Tobit, Censored/Truncated Regression)

Điều gì xảy ra nếu sử dụng công cụ OLS cùng các giả định của mô hình CLRM vào dữ liệu có biến phụ thuộc bị giới hạn?

Xem xét mô hình:

$$SMOKING_i = \beta_0 + \beta_1 * PRICE_i + u_i \quad (1)$$

trong đó $SMOKING_i$ là biến định tính cho hành vi hút thuốc lá của trẻ vị thành niên, nhận giá trị 1 nếu có hút thuốc và 0 nếu không. Biến giải thích là giá bán lẻ.

$$SMOKING_i = \begin{cases} 1 & \text{for smoker} \\ 0 & \text{for non-smoker} \end{cases}$$

- ▶ Trong mô hình thông thường, β_1 là thay đổi của biến phụ thuộc $SMOKING$ nếu biến giải thích $PRICE$ tăng một đơn vị.
- ▶ Đối với biến phụ thuộc nhị phân, $SMOKING_i$ chỉ nhận giá trị 0 hoặc 1, ý nghĩa của β_1 là gì?

Mô hình xác suất tuyến tính - Linear Probability Model (LPM)

- ▶ Với giả thiết kỳ vọng của biến dư bằng 0, $E[u|PRICE] = 0$:

$$E[SMOKING|PRICE] = \beta_0 + \beta_1 * PRICE \quad (2)$$

- ▶ Đồng thời:

$$\begin{aligned} E[SMOKING] &= \\ &= 1 * P(SMOKING = 1) + 0 * P(SMOKING = 0) \\ &= P(SMOKING = 1) \end{aligned}$$

$$\Rightarrow P(SMOKING = 1|PRICE) = \beta_0 + \beta_1 * PRICE$$

- ▶ Điều này có nghĩa là xác suất quan sát được một vị thành niên hút thuốc là mô hình tuyến tính của biến giải thích $PRICE$. Ví dụ $\beta = -0.1$, nếu giá bán tăng 1 đơn vị thì xác suất vị thành niên hút thuốc sẽ giảm 10%.

Những vấn đề của mô hình xác suất tuyến tính

- ▶ Nếu $\beta_1 = -0.1$ thì tăng giá bán thêm 20 đơn vị có làm cho xác suất hút thuốc giảm về 0 hay thậm chí âm không?
- ▶ Tác động biên của giá bán là cố định có hợp lý không? Ví dụ nếu giá thuốc lá tăng từ 10.000đ lên 20.000đ/bao có khác so với tăng từ 100.000đ lên 110.000đ/bao không?
- ▶ Giả định về phương sai không đổi trong mô hình CLRM, $Var(u_i) = \sigma^2$, bị vi phạm.¹

$$Var(u_i|X_i) = P_i * (1 - P_i) , \text{ với}$$

$$P_i = \beta_0 + \beta_1 * PRICE_i$$

$\Rightarrow Var(u_i|PRICE_i) \in PRICE_i$, hay nói cách khác, phương sai của sai số thay đổi.

¹Biến phụ thuộc Y_i phân phối Bernoulli với xác suất $P_i = \beta_0 + \beta_1 * X_i$ nên u_i cũng phân phối Bernoulli với xác suất $P_{u_i} = 1 - \beta_0 - \beta_1 * X_i$. Phương sai của phân phối Bernoulli là $Var(u_i) = P_{u_i} * (1 - P_{u_i})$.

Phương pháp xác suất tối đa - Maximum Likelihood Estimation (MLE)

- ▶ Khắc phục các nhược điểm đã nêu trên, để (a) ước lượng xác suất luôn nằm trong khoảng $[0,1]$ với mọi giá trị của biến giải thích PRICE, và (b) tác động biên của biến giải thích không cố định, chúng ta cần cách tiếp cận mới không sử dụng phương pháp OLS.
- ▶ Giả định xác suất của việc hút thuốc được xác định bởi hàm phân phối xác suất tích lũy $G(\cdot)$:

$$P(SMOKING_i = 1 | PRICE) = G(\beta_0 + \beta_1 * PRICE_i) \quad (3)$$

Với hàm $G(\beta_0 + \beta_1 * PRICE_i)$ nhận giá trị nằm trong khoảng $[0,1]$ với mọi giá trị của biến giải thích PRICE.

- ▶ Hàm phân phối xác suất $G(\cdot)$ thường không biết trước, và phải dựa vào giả định hoặc các lý thuyết kinh tế.

Các hàm phân phối xác suất thông dụng

- ▶ Nếu $G(\cdot)$ có phân phối tích lũy Logistic, khi đó ta có hồi quy "Logit":

$$G(z) = \frac{e^z}{1 + e^z}$$

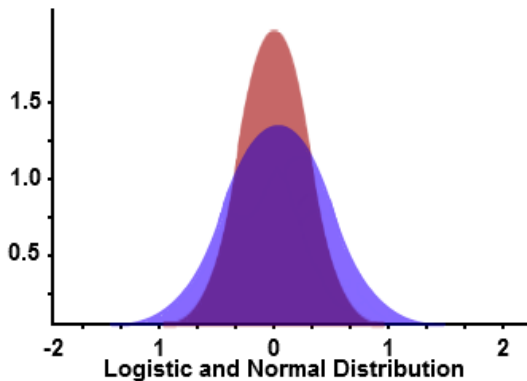
với hàm mật độ phân phối Logistic $g(z) = G'(z) = \frac{e^z}{(1+e^z)^2}$

- ▶ Nếu $G(\cdot)$ có phân phối tích lũy chuẩn \Rightarrow hồi quy Probit:

$$G(z) = \Phi(z) = \int_{-\infty}^z \phi(x) dx$$

với hàm mật độ phân phối chuẩn $\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$

Đồ thị Hàm Mật độ Phân phối Logit (Tím) và Chuẩn (Cam)



Hàm Logistic có mức độ phân tán cao hơn so với phân phối chuẩn.

Ước lượng mô hình hồi quy Logit và Probit (1)

- ▶ Khác với phương pháp bình phương phần dư tối thiểu OLS, mô hình hồi quy dựa trên hàm phân phối xác suất như Logit hay Probit dùng phương pháp xác suất tối đa (Maximum Likelihood Estimation-MLE).
- ▶ Hàm mục tiêu của phương pháp OLS là tối thiểu tổng bình phương phần dư của biến phụ thuộc, còn hàm mục tiêu của phương pháp MLE là tối đa xác suất quan sát được mẫu với thuộc tính cho trước.

Ước lượng mô hình hồi quy Logit và Probit (2)

- ▶ Xác suất quan sát được vị thành niên i có hút thuốc hay không có thể viết như sau:

$$P(SMOKING_i | PRICE_i) = [G(.)]^{SMOKING_i} \times [1 - G(.)]^{1 - SMOKING_i} \quad (4)$$

- ▶ Nếu $SMOKING_i = 1$ thì $P(SMOKING_i | PRICE_i) = G(.)$
 - ▶ Nếu $SMOKING_i = 0$ thì $P(SMOKING_i | PRICE_i) = 1 - G(.)$
- ▶ $G(.)$ là hàm đơn điệu (do $G(.)$ là hàm phân phối xác suất tích lũy, $G(.)$ chỉ tăng hoặc giảm theo biến giải thích), có thể đơn giản hoá bằng cách chuyển đổi từ hàm tích (4) sang hàm logarithm :

$$\ell_i = \ln[P(.)] = SMOKING_i \times \ln[G(.)] + [1 - SMOKING_i] \times \ln[1 - G(.)] \quad (5)$$

Ước lượng mô hình hồi quy Logit và Probit (3)

- ▶ Nếu mẫu dữ liệu có N thành viên thì hàm xác suất tổng thể được tính bằng cách lấy tổng của xác suất của các quan sát:

$$\mathbf{L} = \sum_{i=1}^N \ell_i \quad (6)$$

và việc ước lượng theo phương pháp MLE được thực hiện bằng cách tối đa hoá tổng xác suất \mathbf{L} .

$$\text{Max } \mathbf{L}_{\beta_{MLE}} = \sum_i \left\{ S_i * \ln[G(.)] + [1 - S_i] * \ln[1 - G(.)] \right\} \quad (7)$$

với S_i là biến phụ thuộc $SMOKING_i$, và $G(.)$ là hàm phân phối xác suất tích lũy $G(\beta_0 + \beta_1 * PRICE_i)$.

Ước lượng mô hình hồi quy Logit và Probit (4)

- ▶ Để tìm tham số β_0 và β_1 nhằm tối đa giá trị \mathbf{L} , sử dụng điều kiện tối ưu bậc nhất (first-order condition). Ví dụ với β_1 , sử dụng quy tắc chuỗi (chain-rule) khi lấy đạo hàm bậc nhất:
 - ▶ $G'(\beta_0 + \beta_1 * X_i) = g(.) * X_i$
 - ▶ $\frac{\partial \ln[G(.)]}{\partial \beta_1} = \frac{1}{G(.)} * g(.) * X_i$
- ▶ Do đó, điều kiện bậc nhất với β_1 :

$$\frac{\partial \mathbf{L}}{\partial \beta_1} = \sum_i \left\{ \frac{S_i}{G(.)} * g(.) * X_i - \frac{1 - S_i}{1 - G(.)} * g(.) * X_i \right\} = 0 \quad (8)$$

Ước lượng mô hình hồi quy Logit và Probit (5)

- ▶ Ví dụ đối với hồi quy Logit, $G(z) = \frac{e^z}{1+e^z}$ và $g(z) = \frac{e^z}{(1+e^z)^2}$. Sau khi biến đổi, điều kiện bậc nhất đối với β_1 là:

$$\frac{\partial \mathbf{L}}{\partial \beta_1} = \sum_i S_i * X_i - \sum_i \frac{e^{\beta_0 + \beta_1 * X_i}}{1 + e^{\beta_0 + \beta_1 * X_i}} * X_i = 0 \quad (9)$$

và áp dụng đối với β_0 :

$$\frac{\partial \mathbf{L}}{\partial \beta_0} = \sum_i S_i - \sum_i \frac{e^{\beta_0 + \beta_1 * X_i}}{1 + e^{\beta_0 + \beta_1 * X_i}} = 0 \quad (10)$$

- ▶ Trong phương pháp MLE, do tính phi tuyến của điều kiện bậc nhất (9) và (10) nên không có công thức cụ thể để tính $\hat{\beta}_0$ và $\hat{\beta}_1$ như phương pháp OLS.
- ▶ Việc ước lượng β_0 và β_1 phải sử dụng các phần mềm chuyên dụng.
- ▶ Với hàm Probit thì phương pháp ước lượng cũng tương tự.

Giải thích ý nghĩa của mô hình Logit và Probit (1)

- ▶ Từ giả định xác suất của hành vi hút thuốc (3):

$$P(SMOKING_i = 1 | PRICE) = G(\beta_0 + \beta_1 * PRICE_i) \quad (11)$$

Với những thay đổi nhỏ của giá bán lẻ $PRICE$ thì tác động biên lên xác suất hút thuốc có thể được tính như sau:

$$\frac{\partial P(SMOKING)}{\partial PRICE} = g(\beta_0 + \beta_1 * PRICE_i) * \beta_1 \quad (12)$$

với $g(\beta_0 + \beta_1 * PRICE_i)$ là hàm mật độ phân phối xác suất.

- ▶ Trong phương pháp MLE, tác động biên của giá lên hành vi hút thuốc thay đổi tùy thuộc vào giá trị của hàm mật độ $g(\cdot)$ tại giá bán gốc, khác với tác động biên cố định trong phương pháp hồi quy xác suất tuyến tính LPM!

Giải thích ý nghĩa của mô hình Logit và Probit (2)

- ▶ Thông thường chúng ta tính tác động biên tại mức giá trung bình, tại các tứ phân vị, giá trị tối đa/tối thiểu.
- ▶ Nếu biến giải thích là biến rời rạc (ví dụ có thêm biến giới tính hay số con trong gia đình trong hồi quy Logit đa biến) thì không áp dụng được công thức (12). Khi đó, tác động của giới tính đến hành vi hút thuốc có thể ước lượng trực tiếp từ công thức (11):

$$\begin{aligned}\Delta P &= P(\text{SMOKING}|\text{MALE}) - P(\text{SMOKING}|\text{FEMALE}) \quad (13) \\ &= G(\beta_0 + \beta_1 * \text{PRICE} + D) - G(\beta_0 + \beta_1 * \text{PRICE})\end{aligned}$$

với D là biến giả đại diện cho giới tính.

So sánh giữa LPM, Logit và Probit

Sử dụng bộ dữ liệu mô phỏng SMOKE.dta

Regression Results

	LPM b/se	Logit b/se	Probit b/se
main			
sex	0.0050 (0.0526)	0.0214 (0.2227)	0.0132 (0.1377)
price	-0.0028 (0.0036)	-0.0116 (0.0152)	-0.0072 (0.0094)
Constant	0.5461* (0.2270)	0.2082 (0.9530)	0.1263 (0.5910)
N	807	807	807

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Lưu ý trị kiểm định của mô hình LPM là t-test, của mô hình Logit hoặc Probit là z-test.

Phương trình hồi quy sau ước lượng

- ▶ LMP:

$$\widehat{SMOKE} = .5461 + .0050 * sex - .0028 * price$$

Với mô hình Logit và Probit, phương trình hồi quy được viết dưới dạng log của tỷ lệ thành công (odds ratio):

- ▶ Logit:

$$\log\left(\frac{\widehat{SMOKE}}{1 - \widehat{SMOKE}}\right) = .2082 + .0214 * sex - .0116 * price$$

- ▶ Probit:

$$\log\left(\frac{\widehat{SMOKE}}{1 - \widehat{SMOKE}}\right) = .1263 + .0132 * sex - .0072 * price$$

Diễn giải ý nghĩa các tham số của LPM, Logit và Probit

- ▶ Với mô hình LPM, nam có xác suất hút thuốc cao hơn nữ là 0.5%. Tác động biên của tăng giá thuốc lá 1 cent/bao, xác suất hút sẽ giảm 0.28%.
- ▶ Tác động biên là hằng số, không phụ thuộc vào giá gốc.

Diễn giải ý nghĩa các tham số mô hình Logit và Probit (1)

Với mô hình Logit và Probit, cần tính giá trị hàm mật độ tại các mốc tham chiếu cho trước. Ví dụ đối với quan sát là nam ($sex = 1$), tại mức giá trung bình ($price = 60.03$), tác động biên của tăng giá lên xác suất hút thuốc là:

- ▶ Logit:

$$\begin{aligned}\frac{\partial P(SMOKE)}{\partial price} &= g(\cdot) * \beta_{price} = \frac{e^z}{(1 + e^z)^2} * \beta_{price} \\ &= \frac{e^{(.0214 - .0116 * 60.03 + .2082)}}{(1 + e^{(.0214 - .0116 * 60.03 + .2082)})^2} * (-.0116) \\ &= -.0027451\end{aligned}$$

⇒ tăng giá thuốc lá 1 cent/bao từ mức giá trung bình làm giảm xác suất hút thuốc là 0.27% với đối tượng là nam.

- ▶ Nếu mức giá gốc lần lượt là 44 cent/bao và 70 cent/bao. Tác động biên là bao nhiêu?

Diễn giải ý nghĩa các tham số mô hình Logit và Probit (2)

Với mô hình Probit:

$$\begin{aligned}\frac{\partial P(SMOKE)}{\partial price} &= g(\cdot) * \beta_{price} = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} * \beta_{price} \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{(0.0132 - .0072 * 60.03 + .1263)^2}{2}} * (-.0072) \\ &= -.0027423\end{aligned}$$

⇒ tăng giá thuốc lá 1 cent/bao từ mức giá trung bình làm giảm xác suất hút thuốc là .27% với đối tượng là nam.

- ▶ Nếu mức giá gốc lần lượt là 44 cent/bao và 70 cent/bao. Tác động biên là bao nhiêu?

Diễn giải ý nghĩa các tham số mô hình Logit và Probit (3)

- ▶ Khác biệt về xác suất hút thuốc giữa nhóm nam và nữ như thế nào, tại mức giá trung bình?

$$\Delta P = P(SMOKING|MALE) - P(SMOKING|FEMALE)$$

- ▶ Hàm phân phối tích lũy Logit là $G(z) = \frac{e^z}{1+e^z} \Rightarrow \Delta P = .0050433 \approx 0.5\%$.

- ▶ Hàm phân phối chuẩn (Probit) là $G(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \Rightarrow \Delta P = .0050428 \approx 0.5\%$.

Khả năng dự báo của mô hình xác suất (1)

- ▶ Khả năng dự báo của mô hình: thể hiện xác suất mô hình dự đoán đúng thực tế, bao gồm cả dự báo đúng hành vi hút thuốc và không hút thuốc.
- ▶ Một dự báo được coi là đúng nếu xác suất hút thuốc ước lượng được > 0.5 đối với người có hút thuốc, và xác suất không hút thuốc ước lượng được < 0.5 đối với người không hút thuốc.

Ma trận dự báo

		Dự báo		Tổng số
		Có	Không	
Thực tế	Có	X1 [Đúng]	X2 [Sai]	X1 + X2
	Không	X3 [Sai]	X4 [Đúng]	X3 + X4
Tổng số		X1 + X3	X2 + X4	$\sum X_i$

- ▶ Khả năng dự báo của mô hình $\gamma = \frac{X1+X4}{X1+X2+X3+X4}$

Khả năng dự báo của mô hình xác suất (2)

Ma trận dự báo

	Thực tế		Tổng số
	Có	Không	
Dự báo			
Có	0	0	0
Không	310	497	807
Tổng số	310	497	807

- ▶ $\gamma = \frac{(0+497)}{807} = 61.59\%$
- ▶ Do dữ liệu tự mô phỏng dẫn đến mô hình này dự đoán sai hoàn toàn đối với những người hút thuốc!

Khả năng dự báo của mô hình xác suất (3)

Có thể làm đơn giản hơn bằng lệnh:

```
. estat classification
```

```
Logistic model for SMOKE
```

Classified	True		Total
	D	~D	
+	0	0	0
-	310	497	807
Total	310	497	807

Kiểm định hồi quy Logit (1)

- ▶ Đối với kiểm định đơn biến, sử dụng z-test.
- ▶ Đối với kiểm định đa biến, sử dụng kiểm định Likelihood Ratio (LR). Ví dụ kiểm định k tham số ước lượng đồng thời không có ý nghĩa thống kê:

$H_0 : \beta_1 = \dots = \beta_k = 0$ với $H_1 : \text{Ít nhất một } \beta_j \neq 0$

Kiểm định hồi quy Logit (2)

Cách thực hiện kiểm định LR:

- ▶ Ước lượng hai mô hình riêng biệt: mô hình không giới hạn (unrestricted, u) với đầy đủ biến giải thích, và mô hình giới hạn (restricted, r) không có biến giải thích X_1, \dots, X_k .
- ▶ Tính trị kiểm định $LR = 2 * (\mathbf{L}_u - \mathbf{L}_r)$, với \mathbf{L}_u và \mathbf{L}_r là giá trị log-likelihood từ công thức (7) và tương ứng với mô hình không giới hạn và mô hình giới hạn.
- ▶ LR có phân phối χ_k^2 với số bậc tự do k .
- ▶ Bác bỏ giả thuyết $H_0 \Rightarrow$ ít nhất một trong các tham số kiểm định $\beta_j \neq 0$.

Thực hành trên Stata với bộ dữ liệu SMOKE.dta.