

Cấu trúc Hàm và Lựa chọn Mô hình (Model Specifications)

Lê Việt Phú
Trường Chính sách Công và Quản lý Fulbright

5-8/01/2021

Mục đích của kiểm định

Kiểm định là một phần quan trọng của diễn giải kết quả, lựa chọn mô hình phù hợp nhất, và đảm bảo tính vững (robust/sensitivity check) của kết quả.

- ▶ Kiểm chứng kết quả và sự phù hợp với lý thuyết: Có phát hiện được tác động hay không? có tương thích với lý thuyết và các nghiên cứu đã thực hiện không?
- ▶ Lựa chọn mô hình tối ưu trong số những mô hình có thể sử dụng để diễn giải mối quan hệ giữa biến phụ thuộc và các biến giải thích. Lưu ý rằng mô hình tối ưu nên được chọn dựa trên lý thuyết thay vì kiểm định.
- ▶ Kiểm chứng và phát hiện những thuộc tính thống kê có thể có hàm ý quan trọng.

Một số kiểm định giả thuyết phổ biến đối với hồi quy đa biến

1. Giả thuyết đơn: kiểm định đối với một tham số của mô hình.
2. Kiểm định điều kiện ràng buộc đối với các tham số.
3. Giả thuyết bội: kiểm định đồng thời nhiều tham số.
4. Kiểm định cấu trúc hàm.
5. Kiểm định ước lượng khác biệt nhóm.

2. Kiểm định điều kiện ràng buộc với tham số

Ví dụ ta muốn kiểm định H_0 là tỷ suất thu nhập của đi học bằng với tỷ suất thu nhập của kinh nghiệm làm việc, $\beta_1 = \beta_2$, trong mô hình:

$$\begin{aligned} \log(\text{income}) = & \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \beta_3 \text{married} \\ & + \beta_4 \text{school} + \beta_5 \text{public} + \beta_6 \text{foreign} + \beta_7 \text{official} + u \end{aligned}$$

Trị kiểm định được tính như sau:

$$t = \frac{\hat{\beta}_1 - \hat{\beta}_2}{\text{se}(\hat{\beta}_1 - \hat{\beta}_2)}$$

Có 2 cách thực hiện trong Stata:

1. `test yoeduc = yoexper`

(1) `yoeduc - yoexper = 0`

`F(1, 7544) = 982.41`

`Prob > F = 0.0000`

2. Tạo ra biến mới `sum = yoeduc + yoexper`, ước lượng mô hình với biến `sum`, và kiểm định $\theta = 0$:

$$\begin{aligned} \log(\text{income}) = & \beta_0 + \theta \text{yoeduc} + \beta_2 \text{sum} + \beta_3 \text{married} \\ & + \beta_4 \text{school} + \beta_5 \text{public} + \beta_6 \text{foreign} + \beta_7 \text{official} + u \end{aligned}$$

Lưu ý trị kiểm định F-stat đối với một ràng buộc bằng trị kiểm định t-stat bình phương.

3. Kiểm định giả thuyết bội (multiple hypothesis test)

- ▶ Kiểm định đồng thời nhiều ràng buộc, ví dụ trong mô hình tỷ suất thu nhập ta muốn kiểm định số năm kinh nghiệm làm việc và số năm kinh nghiệm làm việc bình phương đồng thời không có tác động đến thu nhập.

$$\begin{aligned} \log(\text{income}) = & \\ & \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \beta_3 \text{yoexper}^2 + \beta_4 \text{married} \\ & + \beta_5 \text{school} + \beta_6 \text{public} + \beta_7 \text{foreign} + \beta_8 \text{official} + u \end{aligned}$$

$$H_0 : \beta_2 = 0, \beta_3 = 0$$

so với H_1 : ít nhất một trong các đẳng thức không thỏa.

- ▶ Kiểm định giả thuyết bội khác với kiểm định từng biến riêng rẽ. Có thể các biến β_2 và β_3 không có ý nghĩa thống kê nhưng không đồng thời bằng không.

Mô hình gốc (còn gọi là mô hình không bị ràng buộc - **unrestricted model**) là:

$$\begin{aligned} \log(\text{income}) = & \\ & \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \beta_3 \text{yoexper}^2 + \beta_4 \text{married} \\ & + \beta_5 \text{school} + \beta_6 \text{public} + \beta_7 \text{foreign} + \beta_8 \text{official} + u \end{aligned}$$

Mô hình bị ràng buộc (**restricted model**) theo giả thuyết là:

$$\begin{aligned} \log(\text{income}) = & \beta_0 + \beta_1 \text{yoeduc} + \beta_4 \text{married} \\ & + \beta_5 \text{school} + \beta_6 \text{public} + \beta_7 \text{foreign} + \beta_8 \text{official} + u \end{aligned}$$

- ▶ Để kiểm định giả thuyết bội ta dựa vào tổng bình phương của phần dư SSR.
- ▶ Mô hình càng nhiều biến thì SSR càng nhỏ.
- ▶ Sự khác biệt giữa SSR của mô hình bị ràng buộc (SSR_R) và mô hình không bị ràng buộc (SSR_U) có thể dùng để kiểm định của việc thiếu biến trong mô hình.
- ▶ Trị kiểm định có phân phối $F_{q, n-k-1}$, với q là số ràng buộc của mô hình bị ràng buộc:

$$F = \frac{(SSR_R - SSR_U)/q}{SSR_U/(n - k - 1)}$$

- ▶ Kiểm định F còn gọi là kiểm định Wald.

Ví dụ với mô hình tỷ suất thu nhập

$$H_0 : \beta_2 = 0, \beta_3 = 0 \Rightarrow q = 2, n - k - 1 = 7543$$

```
. reg lnincome yoeduc yoexper yoexpersq married publicSchool public foreign official
```

Source	SS	df	MS	Number of obs	=	7,552
Model	1753.70541	8	219.213176	F(8, 7543)	=	409.20
Residual	4040.86526	7,543	.535710627	Prob > F	=	0.0000
				R-squared	=	0.3026
				Adj R-squared	=	0.3019
Total	5794.57067	7,551	.767391162	Root MSE	=	.73192

lnincome	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
yoeduc	.0926075	.0027428	33.76	0.000	.0872309 .0979841
yoexper	.061687	.0025081	24.60	0.000	.0567705 .0666035
yoexpersq	-.0012002	.0000488	-24.58	0.000	-.0012959 -.0011044
married	.0352395	.0221221	1.59	0.111	-.0081259 .078605
publicSchool	-.1145887	.0423549	-2.71	0.007	-.1976161 -.0315613
public	-.1042541	.0329488	-3.16	0.002	-.1688429 -.0396652
foreign	.4499482	.0363715	12.37	0.000	.37865 .5212464
official	.2705426	.0359373	7.53	0.000	.2000956 .3409897
_cons	8.493551	.0474837	178.87	0.000	8.40047 8.586633

```
. test yoexper = yoexpersq = 0
```

```
( 1) yoexper - yoexpersq = 0
```

```
( 2) yoexper = 0
```

```
F( 2, 7543) = 310.83  
Prob > F = 0.0000
```

4. Dùng kiểm định bội để xác định cấu trúc hàm

- ▶ R^2 , R_{adj}^2 đã được sử dụng để lựa chọn biến số và cấu trúc hàm số.
- ▶ F-test cũng có thể sử dụng để kiểm định cấu trúc hàm số trong các mô hình lồng ghép/mô hình gộp (nested models). Ví dụ mô hình (1) được lồng ghép trong mô hình (2):

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u \quad (1)$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + u \quad (2)$$

- Kiểm định $H_0 : \beta_3 = \beta_4 = 0$ để biết liệu hai mô hình trên là tương đương hay không. Nếu bác bỏ H_0 thì mô hình (1) được lồng ghép trong mô hình (2).

- ▶ Nếu dùng kiểm định F để kiểm định tất cả các tham số trong mô hình \Rightarrow Phát hiện ý nghĩa thống kê của mô hình tổng quát (overall significance of the regression).
 - Trong mô hình tỷ suất thu nhập, trị kiểm định $F_{8,7543} = 409.02$, p-value = 0.000.

5. Kiểm định khác biệt giữa các nhóm trong cùng một mô hình - Chow test

Chúng ta muốn kiểm định liệu mô hình tỷ suất thu nhập của việc đi học giống nhau giữa nhóm nam và nữ.

$$\log(\text{income}) = \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \beta_3 \text{yoexpersq} + \beta_4 \text{married} \\ + \beta_5 \text{school} + \beta_6 \text{public} + \beta_7 \text{foreign} + \beta_8 \text{official} + u$$

- ▶ Chúng ta đã ước lượng mô hình trên cho nhóm nam và nữ riêng biệt và quan sát thấy tỷ suất thu nhập của việc đi học với nhóm nữ cao hơn nhóm nam.
- ▶ Câu hỏi: Sự khác biệt có ý nghĩa thống kê hay không?
 - Nếu chỉ một tham số trong tất cả các tham số có sự khác biệt thì các nhóm là không tương đồng.

Trị kiểm định của Chow-test $F_{k+1, n-2(k+1)}$ được tính như sau:

$$F = \frac{[SSR_p - (SSR_1 + SSR_2)] / (k + 1)}{[SSR_1 + SSR_2] / (n - 2(k + 1))}$$

trong đó

- ▶ Giả thuyết H_0 : Tất cả các tham số ước lượng của mô hình nam và nữ là giống nhau.
- ▶ k là số biến giải thích trong mô hình (+1 do thêm tham số tung độ gốc)
- ▶ SSR_p , SSR_1 , SSR_2 là tổng bình phương phần dư của hồi quy gộp toàn bộ dữ liệu, của nhóm nam, và nhóm nữ.

Ví dụ với mô hình tỷ suất thu nhập

$$F = \frac{[4040.8653 - (2234.8287 + 1649.6582)]/(8 + 1)}{[2234.8287 + 1649.6582]/(7552 - 2(8 + 1))} = 33.699694$$

- ▶ Giá trị cực trị của $F_{k+1, n-2(k+1)}$ tại mức tin cậy 99% là $F(9, 7534, .99) = 2.4096768 \Rightarrow$ Bác bỏ H_0 .

Thực hiện Chow-test bằng kiểm định bội

$$\begin{aligned} \log(\text{income}) = & \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \beta_3 \text{yoexpersq} + \beta_4 \text{married} \\ & + \beta_5 \text{school} + \beta_6 \text{public} + \beta_7 \text{foreign} + \beta_8 \text{official} \\ & + \mathbf{male} * [\beta_0 + \beta_1 \text{yoeduc} + \dots + \beta_8 \text{official}] + \mathbf{u} \end{aligned}$$

- ▶ Tạo biến tương tác giữa biến giới tính (male) với các biến giải thích.
- ▶ Ước lượng mô hình gộp bao gồm $2(k+1) = 18$ biến giải thích.
- ▶ Nếu không có sự khác biệt giữa các nhóm nam và nữ thì các tham số ứng với các biến tương tác sẽ đồng thời bằng không.
- ▶ Dùng kiểm định bội đối với mô hình ràng buộc (nhóm nam và nữ giống nhau) và không ràng buộc (nam và nữ khác nhau).
- ▶ Đối chiếu với cách kiểm định dựa trên SSR phía trên.

Thực hiện Chow-test với một số biến giải thích

Ví dụ chúng ta chỉ muốn kiểm định tỷ suất thu nhập của việc đi học giữa hai nhóm nam và nữ.

- ▶ Tạo tương tác giữa biến *male* với biến *yoeduc*,
 $dyoeduc = male * yoeduc$.
- ▶ Đưa 2 biến *male* và *dyoeduc* vào mô hình và ước lượng.
- ▶ Kiểm định $H_0 : male = dyoeduc = 0$.
- ▶ Nếu H_0 bị bác bỏ nghĩa là $male \neq 0$ (tung độ gốc khác nhau) hoặc $dyoeduc \neq 0$ (hệ số góc khác nhau), hoặc cả hai.

Hồi quy với Biến Định tính

(Regression with Qualitative Variables)

Biến định tính là gì

- ▶ Là biến mô tả trạng thái (nam/nữ, đi làm/đi học, làm nông/công chức)
- ▶ Có thể là biến nhị phân (có/không) hoặc biến nhóm (categorical variable - có nhiều hơn 2 trạng thái giá trị, ví dụ phương tiện đi lại là ô tô/xe máy/xe đạp/đi bộ)
- ▶ Đa số trường hợp các biến định tính không thể xếp được thứ bậc (ví dụ làm việc trong khu vực nhà nước/tư nhân/nước ngoài).
- ▶ Một số trường hợp biến định tính có thể xếp được thứ bậc, ví dụ bằng cấp cao nhất có được là gì, từ không có bằng cấp, bằng tiểu học, THCS, THPT, cao đẳng, đại học, thạc sĩ, tiến sĩ.

- ▶ Không nhầm lẫn với biến số đếm rời rạc, ví dụ biến số con cái trong gia đình không phải là biến định tính.
- ▶ Thống kê mô tả biến định tính khác với biến định lượng.
 - Cần xác định nhóm tham chiếu (baseline/reference group) và nhóm được tham chiếu. Ví dụ với biến giới tính thì có thể đặt nhóm tham chiếu là nữ và nhóm được tham chiếu là nam.
 - Giá trị trung bình hoặc tỷ lệ diễn giải xác suất xảy ra một sự kiện.
 - Giá trị lớn nhất và nhỏ nhất không có ý nghĩa kinh tế.
 - Sai số chuẩn liên quan đến xác suất quan sát được sự kiện.
 - Hệ số tương quan mẫu (correlation coefficient) không có ý nghĩa.
 - Thường dùng biến định tính để phân tách và so sánh giữa các nhóm, ví dụ nhóm nam và nữ, đô thị và nông thôn...

Xử lý biến định tính

Sử dụng lại bộ dữ liệu VHLSS 2010.

- ▶ Cần hiểu cách mã hóa biến trong bảng dữ liệu.
- ▶ Có thể gộp biến nhóm thành biến nhị phân.
- ▶ Có thể tách biến nhóm thành nhiều biến nhị phân.
- ▶ Bẫy biến giả (dummy trap): Một biến định tính có n giá trị thì có thể tách ra tối đa là $n - 1$ biến giả. Nếu tách làm n biến giả đưa vào mô hình sẽ có hiện tượng đa cộng tuyến hoàn hảo.

Hồi quy với biến định tính

Ước lượng mô hình tỷ suất thu nhập của đi học với các biến định tính là có gia đình, học trường công, làm nhà nước, làm nước ngoài, là công chức:

$$\log(\text{income}) = \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \beta_3 \text{yoexpersq} + \beta_4 \text{married} \\ + \beta_5 \text{school} + \beta_6 \text{public} + \beta_7 \text{foreign} + \beta_8 \text{official} + u$$

Biến định tính trong mô hình hồi quy còn được gọi là biến giả (dummy variable).

Giải thích ý nghĩa của biến định tính

```
. reg lnincome yoeduc yoexper yoexpersq married publicSchool public foreign official
```

Source	SS	df	MS	Number of obs	=	7,552
Model	1753.70541	8	219.213176	F(8, 7543)	=	409.20
Residual	4040.86526	7,543	.535710627	Prob > F	=	0.0000
				R-squared	=	0.3026
				Adj R-squared	=	0.3019
Total	5794.57067	7,551	.767391162	Root MSE	=	.73192

lnincome	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
yoeduc	.0926075	.0027428	33.76	0.000	.0872309	.0979841
yoexper	.061687	.0025081	24.60	0.000	.0567705	.0666035
yoexpersq	-.0012002	.0000488	-24.58	0.000	-.0012959	-.0011044
married	.0352395	.0221221	1.59	0.111	-.0081259	.078605
publicSchool	-.1145887	.0423549	-2.71	0.007	-.1976161	-.0315613
public	-.1042541	.0329488	-3.16	0.002	-.1688429	-.0396652
foreign	.4499482	.0363715	12.37	0.000	.37865	.5212464
official	.2705426	.0359373	7.53	0.000	.2000956	.3409897
_cons	8.493551	.0474837	178.87	0.000	8.40047	8.586633

Diễn giải ý nghĩa của tham số ước lượng đối với biến định tính

- ▶ Nếu biến phụ thuộc là **thu nhập** thì tham số ước lượng là tác động tăng thêm của nhóm được tham chiếu so với nhóm tham chiếu.
- ▶ Nếu biến phụ thuộc là **log của thu nhập** thì diễn giải tham số ước lượng tùy thuộc vào biến giải thích là biến liên tục hay biến rời rạc.
 - Với **biến liên tục**, ví dụ số năm đi học *yoeduc*, hệ số ước lượng là % tăng thêm của thu nhập. Ví dụ 1 năm đi học làm tăng thu nhập 9.26%.

- ▶ Với **biến rời rạc**, ví dụ các biến định tính, hoặc nếu có biến số con trong gia đình, thì:
 - Nếu β nhỏ, β có thể coi là phần trăm tăng thêm của biến phụ thuộc.
 - Công thức tính chính xác đối với tác động của biến rời rạc lên biến phụ thuộc **log(Y)** là:

$$\frac{Y_1 - Y_0}{Y_0} = e^\beta - 1$$

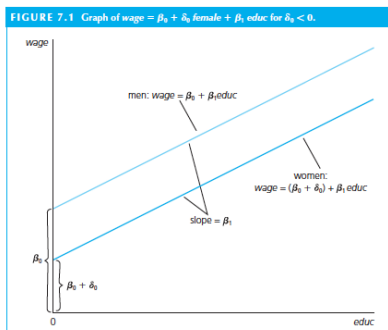
- ▶ Trong ví dụ trên:
 - Làm việc trong khu vực nước ngoài thu nhập cao hơn khu vực tư là: $2.718^{.45} - 1 = .5682$ hay 56.82% (chứ không phải là 45%).
 - Làm việc trong khu vực nhà nước thu nhập thấp hơn khu vực tư là: $2.718^{-.1043} - 1 = -.099$ hay 9.9%.
 - Nếu coi *yoeduc* là biến rời rạc thì với mỗi năm học tăng thêm thu nhập là $2.718^{.0926} - 1 = .097$ hay 9.7%.

Tung độ gốc trong mô hình hồi quy

Với biến giới tính *male* trong mô hình:

$$\log(\text{income}) = \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \dots + \sigma_0 \text{male} + u$$

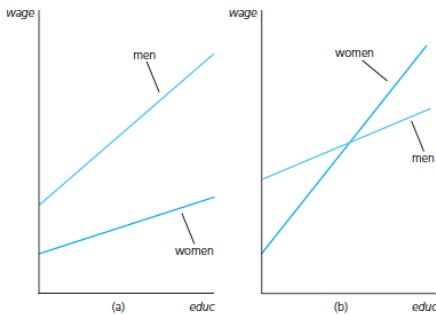
- ▶ Tung độ gốc là β_0 với nhóm nữ, và $\beta_0 + \sigma_0$ với nhóm nam
- ▶ Hệ số góc là β_1 giống nhau với cả hai nhóm (đường hồi quy song song)
- ▶ Nếu $\sigma_0 = 0$ thì hai đường hồi quy trùng nhau



Tung độ gốc và hệ số góc trong mô hình hồi quy với biến tương tác

$$\log(\text{income}) = \beta_0 + \beta_1 \text{yoeduc} + \beta_2 \text{yoexper} + \dots + \sigma_0 \text{male} + \sigma_1 \text{male} * \text{yoeduc} + u$$

- ▶ Tung độ gốc là β_0 với nhóm nữ, và $\beta_0 + \sigma_0$ với nhóm nam
- ▶ Hệ số góc là β_1 với nhóm nữ, và $\beta_1 + \sigma_1$ với nhóm nam.
- ▶ Hai đường hồi quy chỉ trùng nhau khi σ_0 và σ_1 đồng thời bằng 0.



Kiểm định khác biệt theo nhóm

- ▶ Tung độ gốc khác nhau \Rightarrow t-test nếu $\sigma_0 = 0$
- ▶ Tung độ gốc và hệ số góc khác nhau \Rightarrow F-test nếu $\sigma_0 = \sigma_1 = 0$
- ▶ Để kiểm định tất cả các tham số của hai nhóm khác nhau \Rightarrow Chow test

Ôn tập các loại kiểm định

- ▶ Kiểm định đơn: $H_0 : \sigma_0 = 0$

$$t_{\hat{\sigma}_0} \sim t_{n-k-1}$$

- ▶ Kiểm định bội: $H_0 : \sigma_0 = \sigma_1 = 0$

$$F = \frac{(SSR_R - SSR_U)/q}{SSR_U/(n-k-1)} \sim F_{q, n-k-1}$$

- ▶ Kiểm định khác biệt nhóm (tất cả các tham số):

$$H_0 : \sigma_0 = \sigma_1 = \dots = \sigma_k = 0$$

$$F = \frac{[SSR_p - (SSR_1 + SSR_2)]/(k+1)}{[SSR_1 + SSR_2]/(n-2(k+1))} \sim F_{k+1, n-2(k+1)}$$