

Chương Trình Giảng Dạy Kinh tế Fulbright

Học kỳ Thu năm 2011

Các Phương Pháp Phân Tích Định Lượng

Gợi ý giải bài tập 5

KIỂM ĐỊNH THỐNG KÊ

Ngày Phát: Thứ Hai, 31/10/2011

Ngày Nộp: 8:20 sáng, Thứ Hai, 07/11/2011

Bản in nộp tại Hộp nộp bài tập trong phòng Lab

Bản điện tử gửi lên <http://intranet.fetp.edu.vn:81>

Bài 1 (20 điểm)

Trong năm 1960 một cuộc điều tra dân số đã chỉ ra rằng độ tuổi mà người đàn ông Mỹ lập gia đình lần đầu có trung bình là 23.3. Tuy nhiên, vào năm 201X, có một đánh giá được chấp thuận khá rộng rãi là nam thanh niên ở Mỹ ngày nay lập gia đình trễ hơn. Chúng ta muốn kiểm định xem độ tuổi trung bình của lần lập gia đình đầu tiên của họ có tăng lên sau hơn 50 năm hay không?

- a) Hãy lập ra các giả thuyết hợp lý (H_0 và H_a)?

$H_0: \mu = 23.3$ (độ tuổi trung bình của lần đầu tiên lập gia đình của đàn ông Mỹ là 23.3 tuổi)

$H_a: \mu > 23.3$ (độ tuổi trung bình của lần đầu tiên lập gia đình của đàn ông Mỹ lớn hơn 23.3 tuổi)

- b) Lấy mẫu ngẫu nhiên 120 đàn ông ở Mỹ mới lập gia đình lần đầu tiên năm ngoái. Kết quả thống kê mô tả từ mẫu cho thấy độ tuổi trung bình của họ khi lập gia đình lần đầu tiên là 26.2 tuổi và độ lệch chuẩn là 4.5 năm. Thực hiện việc kiểm định giả thuyết và tìm p_{value} ?

Trị thống kê kiểm định cho cỡ mẫu lớn

$$z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} = \frac{26.2 - 23.3}{4.5 / \sqrt{120}}$$

$$z = 2.9 / 0.41 = 7.06$$

Khi mức ý nghĩa là 0.05, giá trị tới hạn 1 phía $z_{0.05} = 1.645$

$$p_{\text{value}} = P(z > 7.06) = 8.35302E-13 \approx 0.000$$

Do $z = 7.06 > z_{0.05} = 1.645$ (hay $P_{\text{value}} < 0.05$) nên bác bỏ H_0 , chấp nhận H_a ở độ tin cậy 95%.

- c) Giải thích ý nghĩa của p_{value} trong trường hợp này?

Xác suất mắc sai lầm loại I khi bác bỏ H_0 là 0.000 ; Do xác suất này nhỏ hơn xác suất tối đa mà chúng ta đặt ra (0.05) nên có thể bác bỏ giả thuyết H_0 .

d) Các anh chị hãy cho biết kết luận của mình?

Kết luận đưa ra là: tuổi trung bình của lần lập gia đình đầu tiên của người dân hiện nay (năm 201X) đã tăng lên so với năm 1960.

Bài 2 (20 điểm)

Một sinh viên điều tra những ảnh hưởng tiềm tàng của chất caffeine đối với việc học thi. Có 40 sinh viên tình nguyện tham gia quá trình điều tra này và được chia làm 2 nhóm. Đầu tiên, mỗi nhóm sẽ làm một bài kiểm tra về trí nhớ. Sau đó, một nhóm sẽ được cho uống hai ly nước ngọt có caffeine và nhóm còn lại sẽ uống hai ly không có caffeine. Ba mươi phút sau mỗi nhóm sẽ thi một bài kiểm tra khác có độ khó tương đương với bài kiểm tra đầu tiên. So sánh kết quả thi trong 2 lần cho thấy đối với nhóm 20 sinh viên uống nước ngọt có chất caffeine thì điểm trung bình giảm xuống 0.925 điểm với độ lệch chuẩn là 2.955 điểm. Đối với nhóm uống nước ngọt không có caffeine thì điểm trung bình tăng lên 1.552 điểm với độ lệch chuẩn là 2.441 điểm. Giả sử khác biệt về điểm thi tuân theo phân phối chuẩn.

a) Điểm của nhóm sinh viên uống nước ngọt có chất caffeine có thay đổi một cách có ý nghĩa không? Kiểm định các giả thuyết và cho biết kết luận của các bạn?

Đây là bài toán kiểm định sự khác biệt cặp giữa hai trung bình tổng thể

$H_0: \mu_d = 0$ (không có sự khác biệt về điểm trung bình của nhóm sinh viên khi uống nước ngọt có chất caffeine so với khi không uống)

$H_a: \mu_d \neq 0$ (có sự khác biệt về điểm trung bình của nhóm sinh viên trước và sau khi uống nước ngọt có chất caffeine)

Trị thống kê kiểm định cho cỡ mẫu nhỏ:

$$t = \frac{\bar{d} - 0}{s_d \sqrt{n}} = \frac{-0.925}{2.955 / \sqrt{20}} = -1.399$$

$|t| = 1.399 < 2.093 \rightarrow$ Không bác bỏ H_0

b) Điểm của nhóm sinh viên uống nước ngọt không có chất caffeine có thay đổi một cách có ý nghĩa không? Kiểm định các giả thuyết và cho biết kết luận của các bạn?

Giả thuyết thống kê:

$H_0: \mu_d = 0$ (không có sự khác biệt về điểm trung bình của nhóm sinh viên trước và sau khi uống nước ngọt không có chất caffeine)

$H_a: \mu_d \neq 0$ (có sự khác biệt về điểm trung bình của nhóm sinh viên trước và sau khi uống nước ngọt không có chất caffeine)

Trị thống kê kiểm định cho cỡ mẫu nhỏ:

$$t = \frac{\bar{d} - 0}{s_d \sqrt{n}} = \frac{1.552}{2.441 / \sqrt{20}} = 2.843$$

$|t|=2.843 > 2.093 \rightarrow$ bác bỏ H_0 , Chấp nhận H_a . Nói cách khác, điểm trung bình của nhóm sinh viên uống nước ngọt không có chất caffeine thay đổi đáng kể so với lúc trước khi uống nước ngọt.

Bài 3 (20 điểm)

Một Trung tâm Quốc gia về Thống kê Giáo dục thường xuyên giám sát các khía cạnh khác nhau của giáo dục phổ thông. Các số liệu thống kê của họ trong năm 2005 và 2010 được xem là cơ sở để đánh giá những thay đổi trong giáo dục phổ thông. Trong năm 2005 có 34% học sinh đã không nghỉ học một ngày nào. Một cuộc điều tra trong năm 2010 cho thấy khi khảo sát 9000 học sinh phổ thông thì tỷ lệ này giảm xuống còn 30%. Các nhà quản lý giáo dục rất quan tâm đến vấn đề này và cho rằng tỷ lệ học sinh tham gia lớp học đã có sự thay đổi.

a) Hãy lập ra các giả thuyết hợp lý (H_0 và H_a)?

$H_0: p_{2010} = 0.34$ (tỷ lệ học sinh tham gia lớp học không có sự thay đổi)

$H_a: p_{2010} \neq 0.34$ (tỷ lệ học sinh tham gia lớp học có sự thay đổi)

b) Thực hiện việc kiểm định giả thuyết và tìm giá trị p_{value} ?

Trị thống kê kiểm định cho cỡ mẫu lớn:

$$z = \frac{\hat{p} - p_0}{\sigma_{\hat{p}}} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0 q_0}{n}}} = \frac{0.3 - 0.34}{\sqrt{\frac{0.34 * 0.66}{9000}}}$$

$$\rightarrow z = -8.01$$

$$p_{value} = 2 * P(z > 8.01) \approx 0.000$$

Do $|z|=8.01 > 1.96$ (hay $P\text{-value} < 0.05$) nên bác bỏ H_0 , chấp nhận H_a ở độ tin cậy 95%.

c) Các bạn có kết luận gì về tình trạng tham gia lớp học của các học sinh?

Tình trạng tham gia lớp học của học sinh đã thay đổi một cách đáng kể so với năm 2005.

d) Bạn có nghĩ rằng sự khác biệt giữa tỷ lệ học sinh tham gia lớp học trong hai năm 2010 và 2005 là có ý nghĩa không? Giải thích?

Sự khác biệt giữa tỷ lệ học sinh tham gia lớp học trong hai năm 2010 và 2005 là có ý nghĩa về mặt thống kê. Nếu như không có vấn đề tiêu cực, hay “bệnh thành tích trong giáo dục” thì việc giảm tỷ lệ này là đáng lo ngại cho những người quản lý giáo dục.

Bài 4 (20 điểm).

Kết quả điều tra mức sống hộ gia đình Việt Nam năm 2008 có thông tin việc làm của trẻ em từ 6 đến dưới 15 tuổi. Giả sử mẫu được chọn theo phương pháp ngẫu nhiên. Trong số 2991 em nam có 251 em đang làm việc; và có 222 em nữ đang làm việc trong 2907 em nữ được điều tra. Liệu có

đủ bằng chứng để kết luận rằng tỷ lệ trẻ em nam đang làm việc khác với tỷ lệ trẻ em nữ đang làm việc hay không?

Tỷ lệ lao động ở trẻ em:

$$p^{\text{nam}} = 251/2991 = 0.084; q^{\text{nam}} = 1 - 0.084 = 0.916$$

$$p^{\text{nữ}} = 222/2907 = 0.076; q^{\text{nữ}} = 1 - 0.076 = 0.924$$

Giả thuyết thống kê:

$H_0: p_{\text{nam}} - p_{\text{nữ}} = 0$ (tỷ lệ trẻ em nam và nữ đang làm việc không khác nhau)

$H_a: p_{\text{nam}} - p_{\text{nữ}} \neq 0$ (tỷ lệ trẻ em nam đang làm việc khác với tỷ lệ trẻ em nữ đang làm việc)

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - D_0}{\sigma_{(\hat{p}_1 - \hat{p}_2)}} = \frac{(\hat{p}_{\text{nam}} - \hat{p}_{\text{nữ}}) - D_0}{\sqrt{\frac{\hat{p}_{\text{nam}}\hat{q}_{\text{nam}}}{n_1} + \frac{\hat{p}_{\text{nữ}}\hat{q}_{\text{nữ}}}{n_2}}} = \frac{0.084 - 0.076}{\sqrt{\frac{0.084 * 0.916}{2991} + \frac{0.076 * 0.924}{2907}}}$$

$z = 1.068 \rightarrow |z| < 1.96 \rightarrow$ Không bác bỏ giả thuyết H_0 với mức ý nghĩa 5%

\rightarrow Chưa đủ cơ sở để bác bỏ giả thuyết H_0 , hay sự khác biệt về tỷ lệ lao động trẻ em giữa nam và nữ là không đáng kể (không có ý nghĩa thống kê).

Bài 5 (20 điểm) Bạn hãy sử dụng phần mềm Stata và dữ liệu VHLSS2008, tính toán các kết quả cần thiết để trả lời các câu hỏi sau (để đơn giản, bạn chưa cần tính đến yếu tố trọng số). Dưới đây là một số gợi ý, kết quả của bạn có thể khác một chút, tùy mục tiêu, giới hạn, giả định của bạn.

a. Có người cho rằng: chỉ tiêu trung bình cho việc học tập trong một năm của học sinh bậc trung học phổ thông là từ 2 triệu đồng trở lên (tạm thời sử dụng biến m2ac13k để đo lường chỉ tiêu cho việc học tập; chỉ xét những người đang đi học hoặc đang nghỉ hè tại thời điểm điều tra). Theo bạn, phát biểu này là đúng hay sai?

```
. ttest m2ac13k=2000 if (m2ac5<=2) & ( m2ac8==3)
```

```
One-sample t test
-----+-----
Variable |      Obs      Mean   Std. Err.   Std. Dev.   [95% Conf. Interval]
-----+-----
m2ac13k |    1693   1633.879   43.97914   1809.569   1547.62   1720.138
-----+-----
      mean = mean(m2ac13k)
Ho: mean = 2000
      t = -8.3249
      degrees of freedom = 1692
      Ha: mean < 2000
      Ha: mean != 2000
      Ha: mean > 2000
Pr(T < t) = 0.0000
Pr(|T| > |t|) = 0.0000
Pr(T > t) = 1.0000
```

Ta thấy, kiểm định giả thuyết 1 phía $H_0: \text{Trung bình} \geq 2000$ bị bác bỏ do P-value < 0.05 . Phát biểu của tác giả trên là không đúng.

b. Có sự khác biệt về chỉ tiêu trung bình cho việc học tập của học sinh trung học phổ thông giữa trường công lập và ngoài công lập (bao gồm bán công, dân lập, tư thục, và hệ khác) hay không? Nếu có, hãy ước lượng sự khác biệt này?

```
. gen conglap=m2ac9==1
```

. tab conglap

conglap	Freq.	Percent	Cum.
0	26,194	74.51	74.51
1	8,960	25.49	100.00
Total	35,154	100.00	

. tab conglap if (m2ac5<=2) & (m2ac8==3)

conglap	Freq.	Percent	Cum.
0	272	16.07	16.07
1	1,421	83.93	100.00
Total	1,693	100.00	

. display 188+66+14+4
 272

. robvar m2ac13k if (m2ac5<=2) & (m2ac8==3), by(conglap)

Summary of 13k.Tæng sè (a+b+...+i)			
conglap	Mean	Std. Dev.	Freq.
0	2578.1691	3019.0531	272
1	1453.1281	1399.5677	1421
Total	1633.8789	1809.5693	1693

W0 = 42.793877 df(1, 1691) Pr > F = 0.00000000

W50 = 25.466884 df(1, 1691) Pr > F = 0.00000050

W10 = 26.857094 df(1, 1691) Pr > F = 0.00000025

Do P-value của thống kê W0 trong kiểm định về sự bằng nhau của phương sai =0.000. Nên phương sai giữa 2 nhóm là khác nhau.

. ttest m2ac13k if (m2ac5<=2) & (m2ac8==3), by(conglap) unequal

Two-sample t test with unequal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
0	272	2578.169	183.057	3019.053	2217.775	2938.564
1	1421	1453.128	37.1276	1399.568	1380.297	1525.959
combined	1693	1633.879	43.97914	1809.569	1547.62	1720.138
diff		1125.041	186.7841		757.4358	1492.646

diff = mean(0) - mean(1) t = 6.0232
 Ho: diff = 0 Satterthwaite's degrees of freedom = 293.659

Ha: diff < 0 Pr(T < t) = 1.0000
 Ha: diff != 0 Pr(|T| > |t|) = 0.0000
 Ha: diff > 0 Pr(T > t) = 0.0000

Chi tiêu cho giáo dục ở bậc học trung học phổ thông có sự khác biệt giữa những trường công lập và ngoài công lập. Chi tiêu của nhóm ngoài công lập cao hơn so với công lập. Ở độ tin cậy 95%, sự khác biệt này là từ 757.4 ngàn đồng đến 1492.6 ngàn đồng.

c. Có sự khác biệt giữa trung bình chi cho việc học thêm các môn học thuộc chương trình (biến m2ac13h) và trung bình chi học phí (biến m2ac13a) của học sinh trung học phổ thông ở khu vực thành thị hay không? Sự khác biệt này như thế nào?

. ttest m2ac13h =m2ac13a if m2ac5<=2 & m2ac8==3 & urban08==1 & m2ac13h!=-2 & m2ac13a!=-2

không. Nếu các giả định bị vi phạm bạn hãy sử dụng một kiểm định phi tham số phù hợp để thay thế.

Ở câu b:

Nếu các giả định của Kiểm định T cho 2 mẫu độc lập không thỏa mãn, bạn hãy sử dụng kiểm định Mann-Whitney U - cũng cần giải quyết kiểm định tăng-hàng Wilcoxon (Hamilton, 2006). Nhưng trình bày sau đây bạn cần chú ý đến kiểm định Wilcoxon: biến của kiểm định cả thang đo thứ bậc, giả định phân phối chuẩn không thỏa mãn, vì mẫu nhỏ, cần chú ý vấn đề do các giá trị bất thường gây nên

Ở câu c:

Kiểm định phi tham số thay thế: signrank (kiểm định dấu vụ hướng)