

# Báo cáo bài tập 7

1612174 - Phùng Tiến Hào - [tienhaophung@gmail.com](mailto:tienhaophung@gmail.com)

12/05/2019

# Contents

<b>1</b>	<b>Chi-square Goodness of Fit Test for single catogorical variable</b>	<b>1</b>
1.1	DecisionMale (Yes/No) . . . . .	1
1.2	RaceF (Caucasian, Asian,..., Other) . . . . .	2
<b>2</b>	<b>Chi-Square Test of Independence (significant association) for two catogorical variables</b>	<b>3</b>

**Dữ liệu khảo sát:** SpeedDating trong package Lock5withR

Load package và thêm các thư viện cần thiết trước khi đi vào xử lý:

```
1 require(Lock5withR) # Load package
2 library(Lock5withR)
3 library(mosaic)
4
5 # head(SpeedDating)
6 attach(SpeedDating) # Avoid dollar sign before each variables name
7
```

# 1 Chi-square Goodness of Fit Test for single catogorical variable

## 1.1 DecisionMale (Yes/No)

Giả sử, ta cần khảo sát tỉ lệ nam phản hồi (Yes/No) cho quần thể (population) là toàn bộ học sinh nam của trường Columbia. Từ tổng thể, ta thu thập được một mẫu dữ liệu ngẫu nhiên (random sample) gồm 276 quan sát trong đó có 146 phản hồi "Yes" và 130 phản hồi "No". Dựa vào mẫu dữ liệu này, ta kiểm định nghi vấn "tỉ lệ phản hồi Yes và phản hồi No bằng nhau" với mức ý nghĩa (significance level) 5%.

Tỉ lệ kỳ vọng (Expected proportion) là 0.5 cho các phản hồi "Yes" và "No".

Ta thực hiện kiểm định giả thuyết:

$$\begin{cases} H_0 : p_{Yes} = p_{No} = 0.5 \\ H_1 : p_{Yes} \neq p_{No} \end{cases}$$

với mức ý nghĩa  $\alpha = 0.5\%$

Lưu ý: Để dùng Chi-square để kiểm định thì các ô quan sát phải ít nhất 5.

```
1 > # Load data
2 > data <- DecisionMale
3 > # Significant level
4 > alpha <- 0.05
5 > # Frequency table
6 > t <- table(data); t
7 data
8 No Yes
9 130 146
10 > # Chisq test
11 > res <- chisq.test(t); res
12
13 Chi-squared test for given probabilities
14
15 data: t
16 X-squared = 0.92754, df = 1, p-value = 0.3355
17
18 > # If p.value < alpha, we ignore H0
19 > (res$p.value < alpha)
20 [1] FALSE
21 > # Expected values
22 > res$expected
```

No Yes  
138 138

Vì  $p\text{-value} > \alpha$  nên ta không bác bỏ  $H_0$ . Như vậy, với mức ý nghĩa 5%, ta không có đủ căn cứ để bác bỏ "tỉ lệ phản hồi "Yes" bằng "No".

### Nhận xét:

- Ta thấy rằng expected values và observed values chênh lệch không nhiều, do đó không đóng góp nhiều vào  $\chi^2$ .

## 1.2 RaceF (Caucasian, Asian,..., Other)

Giả sử, ta cần khảo sát tỉ lệ dân tộc nữ (Caucasian, Asian,..., Other) cho quần thể (population) là toàn bộ học sinh nữ của trường Columbia. Từ tổng thể, ta thu thập được một mẫu dữ liệu ngẫu nhiên (random sample) gồm 276 quan sát trong đó 6 dân tộc: 4 rỗng, 70 Asians, 15 Blacks, 148 Caucasians, 23 Latinos và 16 Others.

Dựa vào mẫu dữ liệu này, ta kiểm định nghi vấn "Tỉ lệ các dân tộc phân bố không đều nhau" với mức ý nghĩa (significance level) 5%.

Gọi  $p_i$  là tỉ lệ dân tộc nữ (Caucasian, Asian,..., Other) trong trường  $\hat{p}_i$  là tỉ lệ dân tộc nữ trong mẫu dữ liệu. Với  $i = 1, 2, \dots, 6$

Ta thực hiện kiểm định giả thuyết:

$$\begin{cases} H_0 : p_i = 0.1666667 \\ H_1 : p_i \neq 0.1666667 \end{cases}$$

với mức ý nghĩa  $\alpha = 0.5\%$

Lưu ý: Để dùng Chi-square để kiểm định thì các ô quan sát phải ít nhất 5.

```
> # Load data
> data <- RaceF
> # Significant level
> alpha <- 0.05
>
> # Frequency table
> t <- table(data); t
data
Asian Black Caucasian Latino Other
4 70 15 148 23 16
> t.prob <- prop.table(t); t.prob
data
Asian Black Caucasian Latino Other
0.01449275 0.25362319 0.05434783 0.53623188 0.08333333 0.05797101
>
> # Chisq test
> res <- chisq.test(t, p = rep(1/6, 6)); res

Chi-squared test for given probabilities

data: t
X-squared = 329, df = 5, p-value < 2.2e-16
```

```

24 > # If p.value < alpha, we ignore H0
25 > (res$p.value < alpha)
26 [1] TRUE
27 > # Expected values
28 > res$expected
29 Asian   Black Caucasian   Latino   Other
30 46      46      46      46      46
31

```

Vì  $p - value < \alpha$  nên ta bác bỏ  $H_0$  và chấp nhận  $H_1$ . Như vậy, với mức ý nghĩa 5%, ta chấp nhận "tỉ lệ các dân tộc nữ phân bố không đều".

### Nhận xét:

- Ta thấy rằng giữa observed values và expected values chênh lệch khá nhiều. Đặc biệt là Caucasian, Asian và Null đóng góp nhiều vào  $\chi^2$ .
- Vì thế nên khả năng để giả thuyết xảy ra là rất thấp, cụ thể ta có  $p - value < 2.2e - 16$  rất bé.

## 2 Chi-Square Test of Independence (significant association) for two categorical variables

Chọn 2 biến định tính: DecisionMale (Yes/No) và RaceF (Asian, Black, Caucasian, Latino, Other)

Khảo sát 2 biến định tính DecisionMale và RaceF

```

1  # 2 biến định tính
2  tab1 = table(DecisionMale, RaceF)
3  # Them margin
4  addmargins(tab1)
5  >
6  RaceF
7  DecisionMale   Asian Black Caucasian Latino Other Sum
8  No    2    32    7    72    7    10 130
9  Yes   2    38    8    76   16    6 146
10 Sum   4    70   15   148   23   16 276
11
12 # 2-way table
13 # Tỉ lệ chung tộc nữ (Asian, Black, ...) nhan phan hoi
14 prop.table(tab1, margin = 1)
15 >
16 RaceF
17 DecisionMale   Asian   Black Caucasian   Latino   Other
18 No  0.01538462 0.24615385 0.05384615 0.55384615 0.05384615 0.07692308
19 Yes 0.01369863 0.26027397 0.05479452 0.52054795 0.10958904 0.04109589
20
21 barplot(tab1, legend = TRUE)
22

```

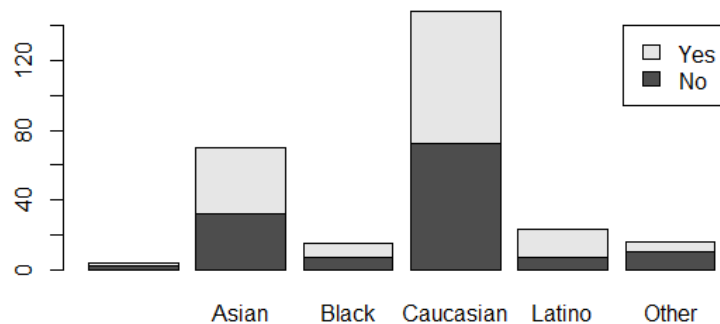


Figure 1: Segmented barchart của DecisionMale và RaceF

Giả sử, ta cần khảo sát sự liên kết giữa dân tộc nữ và sự phản hồi (Yes/No) của nam cho quần thể (population). Từ tổng thể, ta thu thập được một mẫu dữ liệu ngẫu nhiên (random sample) gồm 276 quan sát.

Dựa vào mẫu dữ liệu này, ta kiểm định nghi vấn "Giữa RaceF và DecisionMale có mối liên kết với nhau" với mức ý nghĩa (significance level) 5%.

Ta thực hiện kiểm định giả thuyết:

$$\begin{cases} H_0 : & \text{RaceF và DecisionMale độc lập nhau} \\ H_1 : & \text{RaceF và DecisionMale có sự liên kết với nhau} \end{cases}$$

với mức ý nghĩa  $\alpha = 0.5\%$

### Trực quan hóa bảng tần suất (2-way frequency table)

```

1 > # Load data
2 > data <- data.frame(DecisionMale, RaceF)
3 > # Significant level
4 > alpha <- 0.05
5 >
6 > # Frequency table
7 > t <- table(data); t
8 RaceF
9 DecisionMale Asian Black Caucasian Latino Other
10 No 2 32 7 72 7 10
11 Yes 2 38 8 76 16 6
12 > # P(RaceF|DecisionMale)
13 > t.prob <- prop.table(t, margin = 1); t.prob
14 RaceF
15 DecisionMale Asian Black Caucasian Latino Other
16 No 0.01538462 0.24615385 0.05384615 0.55384615 0.05384615 0.07692308
17 Yes 0.01369863 0.26027397 0.05479452 0.52054795 0.10958904 0.04109589
18
19 > library("graphics")
20 > # shade: color graph
21 > # las = 1: horizontal labels
22 > mosaicplot(t(t), shade = TRUE, las = 1, main = "data")

```

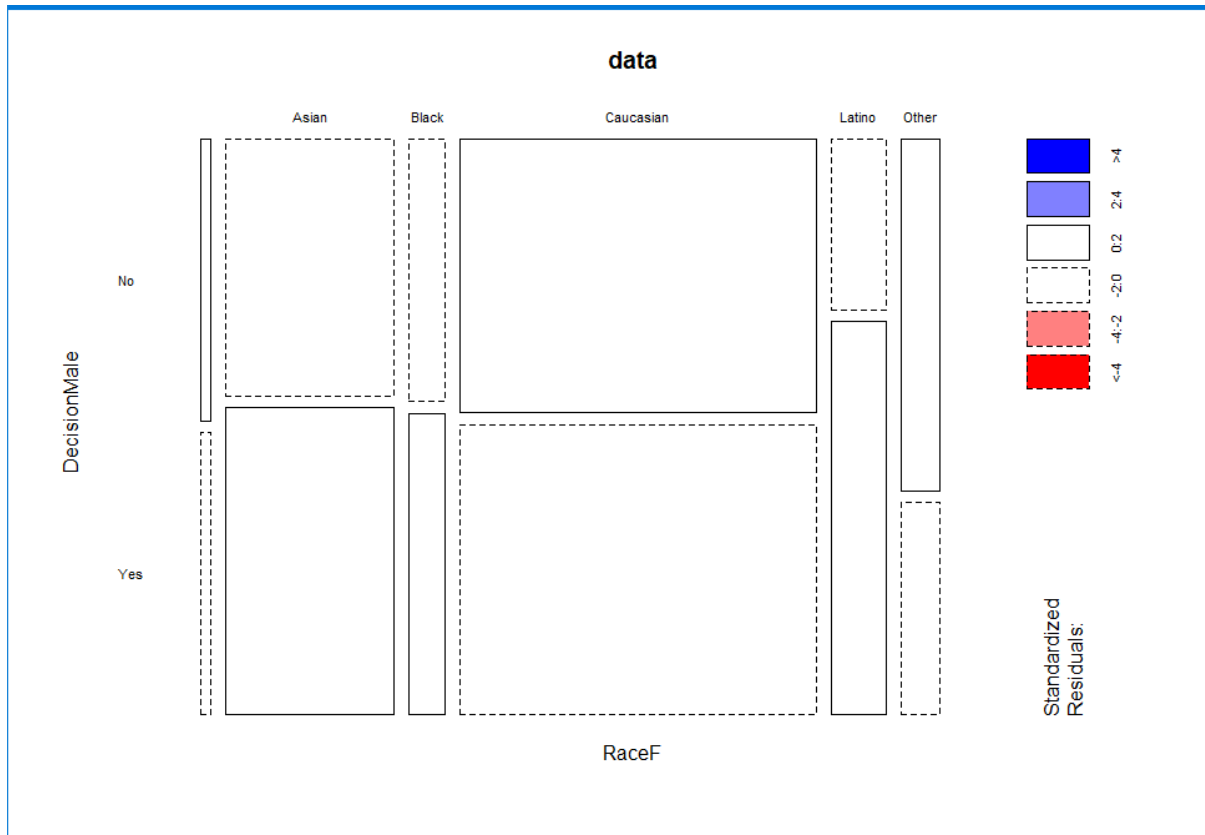


Figure 2: Màu đỏ biểu thị observed values bé hơn expected values  
 Màu xanh biểu thị observed values lớn hơn expected values  
 (với điều kiện dữ liệu phải là ngẫu nhiên)

### Nhận xét

- Nhìn vào mosaicplot thì ta thấy rằng bảng dữ liệu của chúng ta thấy rằng sự chênh lệch giữa observed values và expected values rất bé. Một phần nào cho ta thấy được giữa RaceF và DecisionMale hầu như không có liên kết.
- Các ô (cell) màu trắng nét liền biểu thị độ lệch dương và màu trắng nét đứt biểu thị độ lệch âm nhưng ta thấy rằng các độ lệch này rất bé.

Với mỗi ô thì ta có thể tính được expected value tương ứng:

$$e = \frac{row.sum * col.sum}{grand.total}$$

Chi-square statistic tính như sau:

$$\chi^2 = \sum \frac{(o - e)^2}{e}$$

<sup>1</sup>mosaicplot() là hàm built-in của R package **graphics**

với  $\begin{cases} o: \text{observed value} \\ e: \text{expected value} \end{cases}$

## Tính Chi-square statistic trong R

```

1 > # Chisq test
2 > res <- chisq.test(t); res
3 Warning message:
4 In chisq.test(t) : Chi-squared approximation may be incorrect
5
6 Pearson Chi-squared test
7
8 data: t
9 X-squared = 4.2977, df = 5, p-value = 0.5074
10
11 > # If p.value < alpha, we ignore H0
12 > (res$p.value < alpha)
13 [1] FALSE
14
15 > # Observed values
16 > res$observed
17 RaceF
18 DecisionMale Asian Black Caucasian Latino Other
19 No 2 32 7 72 7 10
20 Yes 2 38 8 76 16 6
21 > # Expected values
22 > round(res$expected, 2)
23 RaceF
24 DecisionMale Asian Black Caucasian Latino Other
25 No 1.88 32.97 7.07 69.71 10.83 7.54
26 Yes 2.12 37.03 7.93 78.29 12.17 8.46
27

```

Vì  $p - value > \alpha$  nên ta không bác bỏ  $H_0$ . Như vậy, với mức ý nghĩa 5%, ta không có đủ căn cứ để bác bỏ "RaceF và DecisionMale độc lập hay nói cách khác cả hai không có sự liên kết".

Nhìn vào observed values table và expected values table, ta thấy được rằng chênh lệch ở đây rất ít.

Để biết rõ về bản chất của sự phụ thuộc giữa 2 biến RaceF và DecisionMale, ta sẽ tiếp tục tính dư lượng chuẩn hóa (Standardized residuals hoặc Pearson residuals) cho từng ô để biết được ô nào đóng góp nhiều vào Chi-square  $\chi^2$ :

$$r = \frac{o - e}{\sqrt{e}}$$

Pearson residuals được lấy từ kết quả của `chisq.test()`:

```

1 > # Pearson Residuals: Do lech giua observed values and expected values
2 > round(res$residuals, 3)
3 RaceF
4 DecisionMale Asian Black Caucasian Latino Other
5 No 0.084 -0.169 -0.025 0.274 -1.165 0.897
6 Yes -0.080 0.160 0.023 -0.259 1.099 -0.847
7
8 > # Visualize Pearson residuals
9 > library(corrplot)
10 corrplot 0.84 loaded

```



```

11 Warning message:
12 package ‘‘corrplot was built under R version 3.5.3
13 > corrplot(res$residuals, is.cor = FALSE)
14

```

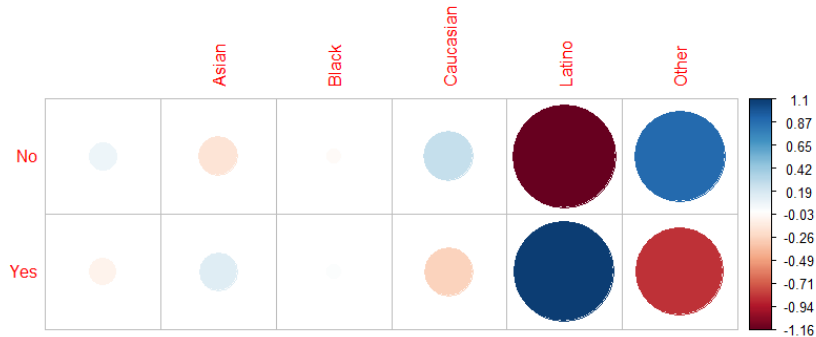


Figure 3: Correlation plot

### Chú thích:

- Kích thước của hình tròn là tỉ lệ thuận với mức độ đóng góp của ô đó.
- Hệ số tương quan ở đây khác với của 2 biến định lượng là không có dao động trong miền  $[-1, 1]$

### Nhận xét:

- Màu đỏ biểu thị mối liên kết âm (negative association). Như ta thấy, hàng No và cột Latino; hàng Yes và các cột Other, Caucasian có mối liên kết âm (negative association) - nghĩa là khi cái này tăng thì cái kia giảm. Ngụ ý là sự không thích (repulsion), ta có thể thấy dân tộc Latino là nhận nhiều phản hồi No nhất (tức là bị từ chối) và ngược lại, dân tộc Other và Caucasian lại nhận liên kết âm với phản hồi Yes.
- Ngược lại, màu xanh biểu thị mối liên kết dương (positive association) - nghĩa là cả hai đều cùng tăng. Như ta thấy, hàng Yes và cột Latino; hàng No và cột Other và Caucasian có liên kết dương. Ngụ ý là sự thu hút (attraction). Điều đó cho thấy người Latino có xu hướng nhận được phản hồi Yes cao. Tương tự, hàng No, cột Other và Caucasian có liên kết dương mạnh, có thể hiểu đơn giản là khi số lượng dân tộc Other và Caucasian tăng thì khả năng họ nhận được phản No (bị từ chối) cũng tăng theo.

Bây giờ, để biết được mức độ phần trăm đóng góp của các ô cho Chi-square  $\chi^2$ , ta tính theo công thức sau:

$$contrib = \frac{r^2}{\chi^2}$$

```

1 > # Contribution (Percentage %) of given cell to total chi-square
2 > contrib <- 100*res$residuals^2 / res$statistic
3 > round(contrib, 3)
4 RaceF
5 DecisionMale    Asian Black Caucasian Latino  Other
6 No  0.166 0.665 0.014  1.750 31.561 18.742
7 Yes 0.148 0.592 0.012  1.558 28.102 16.688

```

```

8
9 > # Visualiza contribution
10 > corrplot(contrib, is.cor = FALSE)
11

```

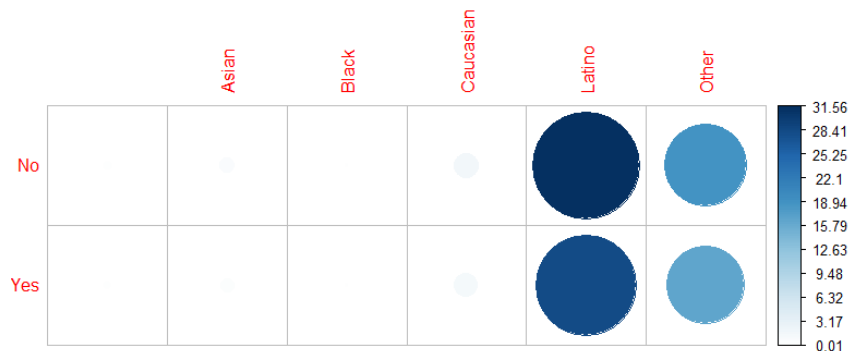


Figure 4: Contribution in percentage (%) plot

Sự đóng góp tương đối của mỗi ô vào tổng Chi bình phương cho thấy một số dấu hiệu về bản chất của sự phụ thuộc giữa các hàng và cột của bảng tần suất.

Từ ảnh trên, ta kết luận được:

- Hàng "No" có liên kết mạnh với cột Latino và Other
- Hàng "Yes" cũng có liên kết mạnh với Latino và Other
- Từ bảng trên tính bằng R, ta thấy các ô đóng góp nhiều cho Chi-square là No/Latino (31.561%), No/Other (18.742%), Yes/Latino (28.102%) và Yes/Other (16.688%).
- Tổng cộng 4 ô trên đóng góp tới tận 95.048% vào tổng Chi bình phương và vì vậy chúng chiếm phần lớn sự khác biệt giữa các giá trị kì vọng và giá trị quan sát.

## References

- [1] Randall Pruim and Lana Park. Lock5WithR. Chapter 7: Chi-Squared Tests for Categorical Variables. PDF.
- [2] R Users Guide. Chapter 7: Chi-Squared Tests for Categorical Variables. PDF.
- [3] Chi-Square Test of Independence in R. (n.d.). Retrieved from <http://www.sthda.com/english/wiki/chi-square-test-of-independence-in-r>
- [4] Chi-square Goodness of Fit Test in R. (n.d.). Retrieved from <http://www.sthda.com/english/wiki/chi-square-goodness-of-fit-test-in-r>