

# Math 170S

## Lecture Notes Section 9.2 <sup>\*†</sup>

### Contingency tables

Instructor: Swee Hong Chan

---

**NOTE:** Materials that appear in the textbook but do not appear in the lecture notes might still be tested. Please send me an email if you find typos.

---

\*Version date: Tuesday 8<sup>th</sup> December, 2020, 22:11.

†This notes is based on Hanbaek Lyu's and Liza Rebrova's notes from the previous quarter, and I would like to thank them for their generosity. "*Nanos gigantum humeris insidentes* (I am but a dwarf standing on the shoulders of giants)".

# 1 Motivating example

Two instructors are teaching Math 170S for  $n = 50$  students, with final grade:

---

	A	B	C	D	F
Friendly instructor (FI)	8	13	16	10	3
Evil instructor (EI)	4	9	14	16	7

---

- Null hypothesis  $H_0$ : The two instructors **follow** the same grading scheme.
- Alternative hypothesis  $H_1$ : The two instructors **do not follow** the grading scheme.

Can we reject  $H_0$  with significance level  $\alpha = 0.05$ ?

## 2 Notation: Motivating examples

$p_1$  := probability to get A in FI's class;

$p_2$  := probability to get B in FI's class;

⋮

$p_5$  := probability to get F in FI's class.

$p'_1$  := probability to get A in EI's class;

$p'_2$  := probability to get B in EI's class;

⋮

$p'_5$  := probability to get F in EI's class.

These are **unknown parameters**.

$Y_1 :=$  number of people getting A in FI's class;

$Y_2 :=$  number of people getting B in FI's class;

$\vdots$

$Y_5 :=$  number of people getting F in FI's class.

$Y'_1 :=$  number of people getting A in EI's class;

$Y'_2 :=$  number of people getting B in EI's class;

$\vdots$

$Y'_5 :=$  number of people getting F in EI's class.

These are **known parameters**.

The hypothesis can then be rewritten as

- $H_0$ :  $p_i = p'_i$  for **all**  $i \in \{1, \dots, 5\}$ ;
- $H_1$ :  $p_i \neq p'_i$  for **some**  $i \in \{1, \dots, 5\}$ .

### 3 How to test the hypothesis

1. We already knew the sample mean is a good approximation for the unknown probabilities:

$$Y_1 \approx np_1; \quad Y'_1 \approx np'_1, \quad \text{and} \quad \frac{Y_1 + Y'_1}{2n} \approx \frac{p_1 + p'_1}{2}.$$

Let us write

$$\hat{p}_1 := \frac{Y_1 + Y'_1}{2n}.$$

2. Now note that, if  $p_1 \approx p'_1$ , then

$$p_1 \approx \frac{p_1 + p'_1}{2}.$$

3. Combining all these observations, if  $p_1 \approx p'_1$ :

$$Y_1 \approx np_1 \approx n \frac{p_1 + p'_1}{2} \approx n \frac{Y_1 + Y'_1}{2n} = n \hat{p}_1.$$

4. Hence we have, if  $p_1 \approx p'_1$ , then

$$Y_1 - n \hat{p}_1 \approx 0.$$

By CLT, we in fact have, if  $p_1 \approx p'_1$  :

$$\frac{(Y_1 - n \hat{p}_1)^2}{n \hat{p}_1} \text{ is small.}$$

5. On the other hand, if  $p_1$  is very far away from  $p'_1$ ,  
(e.g.,  $p_1 = 0$  and  $p'_1 = 1$ ), then

$$Y_1 = 0; \quad \hat{p}_1 = \frac{1}{2},$$

so

$$\frac{(Y_1 - n \hat{p}_1)^2}{n \hat{p}_1} = \frac{(0 - n/2)^2}{n/2} = n/2 = \text{very big.}$$

6. So we conclude that

$$\frac{(Y_1 - n \hat{p}_1)^2}{n \hat{p}_1} \text{ is small} \quad \text{if and only if} \quad p_1 \approx p'_1.$$

By the same reasoning, we have

$$\frac{(Y'_1 - n \hat{p}_1)^2}{n \hat{p}_1} \text{ is small} \quad \text{if and only if} \quad p_1 \approx p'_1.$$

7. To provide balance, we add these two tests together:

$$\frac{(Y_1 - n \hat{p}_1)^2}{n \hat{p}_1} + \frac{(Y'_1 - n \hat{p}_1)^2}{n \hat{p}_1} \text{ is small}$$

if and only if  $p_1 \approx p'_1$ .

8. By the same reasoning, for all  $i = \{1, 2, 3, 4, 5\}$ ,

$$\frac{(Y_i - n \hat{p}_i)^2}{n \hat{p}_i} + \frac{(Y'_i - n \hat{p}_i)^2}{n \hat{p}_i} \text{ is small}$$

if and only if  $p_i \approx p'_i$ .

9. We want to test all five parameters **simultaneously**. So we add all the tests up together:

$$Q := \left[ \frac{(Y_1 - n \hat{p}_1)^2}{n \hat{p}_1} + \frac{(Y_1' - n \hat{p}_1)^2}{n \hat{p}_1} \right] \\ + \left[ \frac{(Y_2 - n \hat{p}_2)^2}{n \hat{p}_2} + \frac{(Y_2' - n \hat{p}_2)^2}{n \hat{p}_2} \right] \\ + \dots + \left[ \frac{(Y_5 - n \hat{p}_5)^2}{n \hat{p}_5} + \frac{(Y_5' - n \hat{p}_5)^2}{n \hat{p}_5} \right].$$

We have

$Q$  is small if and only if  $H_0$  is true.

10. It can be shown that  $Q$  is approximately a  $\chi^2$  random variable with 4 degrees of freedom.

**Conclusion:** we reject  $H_0$  if and only if  $Q \geq \chi_\alpha^2(4)$ .

## 4 Answer: motivating examples

We have from the sample data that

$$Y_1 = 8; \quad Y_2 = 13; \quad Y_3 = 16; \quad Y_4 = 10; \quad Y_5 = 3;$$

$$Y'_1 = 4; \quad Y'_2 = 9; \quad Y'_3 = 14; \quad Y'_4 = 16; \quad Y'_5 = 7,$$

and

$$\hat{p}_1 = 0.12; \quad \hat{p}_2 = 0.22; \quad \hat{p}_3 = 0.30;$$

$$\hat{p}_4 = 0.26; \quad \hat{p}_5 = 0.10.$$

Then  $Q$  is equal to

$$\begin{aligned} Q = & \left[ \frac{((8) - (50)(0.12))^2}{(50)(0.12)} + \frac{((4) - (50)(0.12))^2}{(50)(0.12)} \right] \\ & + \left[ \frac{((13) - (50)(0.22))^2}{(50)(0.22)} + \frac{((9) - (50)(0.22))^2}{(50)(0.22)} \right] \\ & + \dots + \left[ \frac{((3) - (50)(0.10))^2}{(50)(0.10)} + \frac{((7) - (50)(0.10))^2}{(50)(0.10)} \right] = 5.18. \end{aligned}$$

On the other hand,  $\chi_\alpha^2(4)$  is equal to

$$\chi_\alpha^2(4) = \chi_{0.05}^2(4) = 9.488.$$

Since  $Q$  is smaller than  $\chi_\alpha^2(4)$ , we conclude that the test is inconclusive.

# 5 Settings: equality in distribution

**Object:**

- $X^{(1)}, X^{(2)}, \dots, X^{(h)}$  are **independent** random variables with **unknown distribution**.
- $k$  mutually exclusive, exhaustive events  $A_1, \dots, A_k$ , and we write

$p_i^{(j)} :=$  probability of the event  $A_i$  to occur for  $X^{(j)}$ ,

for  $i \in \{1, 2, \dots, k\}$  and  $j \in \{1, \dots, h\}$ .

## Hypotheses:

- **Null Hypothesis**  $H_0$ :  $X^{(1)}, X^{(2)}, \dots, X^{(h)}$  have the same distribution, i.e.,

$$p_1^{(1)} = p_1^{(2)} = \dots = p_1^{(h)}; \quad \text{and}$$

$$p_2^{(1)} = p_2^{(2)} = \dots = p_2^{(h)}; \quad \text{and}$$

$$\vdots \quad \vdots \quad \vdots$$

$$p_k^{(1)} = p_k^{(2)} = \dots = p_k^{(h)}.$$

- **Alternative Hypothesis**  $H_1$ : The null hypothesis is false.

**Input:**  $n^{(1)}$  many random samples for  $X^{(1)}$ ,  $n^{(2)}$  many random samples for  $X^{(2)}$ ,  $\dots$ ,  $n^{(h)}$  many random samples for  $X^{(h)}$ , and significance level  $\alpha$ .

**Methodology:**

- Compute  $Y_i^{(j)}$  for  $i \in \{1, 2, \dots, k\}$  and  $j \in \{1, \dots, h\}$ ,  
 $Y_i^{(j)} :=$  number of times  $A_i$  occurs in samples for  $X^{(j)}$ .
- Compute  $\hat{p}_1, \dots, \hat{p}_k$  given by

$$\hat{p}_i := \frac{Y_i^{(1)} + Y_i^{(2)} + \dots + Y_i^{(h)}}{n^{(1)} + n^{(2)} + \dots + n^{(h)}}.$$

- Compute  $Q$  given by

$$Q := \sum_{j=1}^h \sum_{i=1}^k \frac{(Y_i^{(j)} - n^{(j)} \hat{p}_i)^2}{n^{(j)} \hat{p}_i}.$$

- Reject  $H_0$  if  $Q \geq \chi_\alpha^2((h-1)(k-1))$ , and the test is inconclusive otherwise.

## 6 Example: equality in distribution

A survey was conducted, asking for the education level and the media preference for news sources:

---

	Newspaper	Television	Radio
Grade school	45	22	6
High School	94	115	30
College	49	52	13

---

Let  $X^{(1)}$  be the (random) media preference for grade schoolers,  $X^{(2)}$  the (random) media preference for high schoolers, and  $X^{(3)}$  be the (random) media preference for college schoolers.

Can we reject the hypothesis  $X^{(1)} = X^{(2)} = X^{(3)}$  with significance level  $\alpha = 0.05$ ?

# 7 Answer: equality in distribution

From the sample data, we have

$$n^{(1)} = 45 + 22 + 6 = 73;$$

$$n^{(2)} = 94 + 115 + 30 = 239;$$

$$n^{(3)} = 49 + 52 + 13 = 114,$$

and

$$Y_1^{(1)} = 45; \quad Y_2^{(1)} = 22; \quad Y_3^{(1)} = 6;$$

$$Y_1^{(2)} = 94; \quad Y_2^{(2)} = 115; \quad Y_3^{(2)} = 30;$$

$$Y_1^{(3)} = 49; \quad Y_2^{(3)} = 52; \quad Y_3^{(3)} = 13.$$

Note that here  $h = k = 3$ .

So  $\hat{p}_i$ 's are given by

$$\begin{aligned}\hat{p}_1 &:= \frac{\text{first column}}{\text{total samples}} = \frac{45 + 94 + 49}{73 + 239 + 114} = \frac{188}{426}; \\ \hat{p}_2 &:= \frac{\text{second column}}{\text{total samples}} = \frac{22 + 115 + 52}{73 + 239 + 114} = \frac{189}{426}; \\ \hat{p}_3 &:= \frac{\text{third column}}{\text{total samples}} = \frac{6 + 30 + 13}{73 + 239 + 114} = \frac{49}{426}.\end{aligned}$$

Then  $Q$  is equal to

$$\begin{aligned}Q &= \sum_{j=1}^h \sum_{i=1}^k \frac{(Y_i^{(j)} - n^{(j)} \hat{p}_i)^2}{n^{(j)} \hat{p}_i} \\ &= \frac{((45) - (73)(\frac{188}{426}))^2}{(73)(\frac{188}{426})} + \frac{((22) - (73)(\frac{189}{426}))^2}{(73)(\frac{189}{426})} + \frac{((6) - (73)(\frac{49}{426}))^2}{(73)(\frac{49}{426})} \\ &+ \frac{((94) - (239)(\frac{188}{426}))^2}{(239)(\frac{188}{426})} + \frac{((115) - (239)(\frac{189}{426}))^2}{(239)(\frac{189}{426})} + \frac{((30) - (239)(\frac{49}{426}))^2}{(239)(\frac{49}{426})} \\ &+ \frac{((49) - (114)(\frac{188}{426}))^2}{(114)(\frac{188}{426})} + \frac{((52) - (114)(\frac{189}{426}))^2}{(114)(\frac{189}{426})} + \frac{((13) - (114)(\frac{49}{426}))^2}{(114)(\frac{49}{426})} \\ &= \frac{27503239}{3026142} + \frac{235591}{105399} + \frac{354817}{4725756} = \frac{939839042381}{82450264932} \approx 11.40.\end{aligned}$$

On the other hand,  $\chi_\alpha^2((h-1)(k-1))$  is equal to

$$\chi_\alpha^2((h-1)(k-1)) = \chi_{0.05}^2(4) = 9.488.$$

Since  $Q$  is greater than  $\chi_\alpha^2(4)$ , we reject the null hypothesis.

**Remark 1.** The textbook uses different notations. The  $Y_i^{(j)}$  here is written  $Y_{ij}$  in the textbook,  $n^{(j)}$  here is written  $n_j$  in the textbook, and  $p_i^{(j)}$  is written  $p_{ij}$  in the textbook.

## 8 Example: Contingency tables

Four hundred UCLA undergraduate students are classified according to their college and their gender:

---

	Bsns	Engnrg	Lib. Arts	Nursing	Phrmcy
Male	21	16	145	2	6
Female	14	4	175	13	4

---

Test at  $\alpha = 0.01$  whether gender and choice of college are independent.

# 9 Setting: Contingency tables

## Object:

- $X$  is an **unknown** random variables.
- Two different attributes:
  - $k$  mutually exclusive, exhaustive events  $A_1, \dots, A_k$ ;
  - $h$  mutually exclusive, exhaustive events  $B_1, \dots, B_h$ ;

## Hypotheses:

- **Null Hypothesis  $H_0$** : The two attributes **are independent**, i.e., for  $i \in \{1, \dots, k\}, j \in \{1, \dots, h\}$ :

$$P[A_i \cap B_j] = P[A_i]P[B_j].$$

- **Alternative Hypothesis  $H_1$** : The two attributes **are not independent**.

**Input:**  $n$  many random samples for  $X$ , and significance level  $\alpha$ .

**Methodology:**

- Compute  $Y(A_i, B_j)$  for  $i \in \{1, 2, \dots, k\}$  and  $j \in \{1, \dots, h\}$  by

$Y(A_i, B_j) :=$  number of times  $A_i$  and  $B_j$  occurs in samples.

- Compute  $a_1, \dots, a_k$  given by

$$a_i := \frac{1}{n}(\text{number of times } A_i \text{ occurs in the samples}).$$

- Compute  $b_1, \dots, b_h$  given by

$$b_j := \frac{1}{n}(\text{number of times } B_j \text{ occurs in the samples}).$$

- Compute  $Q$  given by

$$Q := \sum_{j=1}^h \sum_{i=1}^k \frac{(Y(A_i, B_j) - na_i b_j)^2}{na_i b_j}.$$

- Reject  $H_0$  if  $Q \geq \chi_\alpha^2((h-1)(k-1))$ , and the test is inconclusive otherwise.

## 10 Answer: Contingency tables

Let  $A_1, \dots, A_5$  be the event that a student is in the college of business, engineering, liberal arts, nursing, pharmacy, respectively.

Let  $B_1$  be the event that the student is male, and let  $B_2$  is the event that the student is female.

We have from the sample data that

$$a_1 = \frac{35}{400}; \quad a_2 = \frac{20}{400}; \quad a_3 = \frac{320}{400}; \quad a_4 = \frac{15}{400}; \quad a_5 = \frac{10}{400};$$

and

$$b_1 = \frac{190}{400}; \quad b_2 = \frac{210}{400}.$$

So  $Q$  is equal to

$$\begin{aligned}
Q &= \frac{((21) - (400)(\frac{35}{400})(\frac{190}{400}))^2}{(400)(\frac{35}{400})(\frac{190}{400})} + \frac{((16) - (400)(\frac{20}{400})(\frac{190}{400}))^2}{(400)(\frac{20}{400})(\frac{190}{400})} \\
&+ \frac{((145) - (400)(\frac{320}{400})(\frac{190}{400}))^2}{(400)(\frac{320}{400})(\frac{190}{400})} + \frac{((2) - (400)(\frac{15}{400})(\frac{190}{400}))^2}{(400)(\frac{15}{400})(\frac{190}{400})} \\
&+ \frac{((6) - (400)(\frac{10}{400})(\frac{190}{400}))^2}{(400)(\frac{10}{400})(\frac{190}{400})} + \frac{((14) - (400)(\frac{35}{400})(\frac{210}{400}))^2}{(400)(\frac{35}{400})(\frac{210}{400})} \\
&+ \frac{((4) - (400)(\frac{20}{400})(\frac{210}{400}))^2}{(400)(\frac{20}{400})(\frac{210}{400})} + \frac{((175) - (400)(\frac{320}{400})(\frac{210}{400}))^2}{(400)(\frac{320}{400})(\frac{210}{400})} \\
&+ \frac{((13) - (400)(\frac{15}{400})(\frac{210}{400}))^2}{(400)(\frac{15}{400})(\frac{210}{400})} + \frac{((4) - (400)(\frac{10}{400})(\frac{210}{400}))^2}{(400)(\frac{10}{400})(\frac{210}{400})} \\
&= 18.93.
\end{aligned}$$

On the other hand,  $\chi_\alpha^2((h-1)(k-1))$  is equal to

$$\chi_\alpha^2((h-1)(k-1)) = \chi_{0.01}^2(4) = 13.28.$$

Since  $Q$  is greater than  $\chi_\alpha^2(4)$ , we reject the null hypothesis.

**Remark 2.** This test can be extended to test more than two attributes. Check the textbook for exercises.

**Remark 3.** The textbook uses different notations. The  $Y(A_i, B_j)$  here is written  $Y_{ij}$  in the textbook,  $a_i$  here is written  $Y_{i.}/n$  in the textbook, and  $b_j$  is written  $Y_{.j}/n$  in the textbook.