

Math 170S

Lecture Notes Section 8.4 ^{*†}

Tests about proportions

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: Thursday 19th November, 2020, 23:50.

†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 Setting: Bernoulli, one variable

Object: Y is a **Bernoulli** random variables with **unknown parameter** p .

Hypotheses:

- **Null Hypothesis** H_0 : p is equal to p_0 .
- **Alternative Hypothesis** H_1 : The alternative hypothesis can take one of these three forms:
 - (a) p is **strictly greater than** p_0 ;
 - (b) p is **strictly smaller than** p_0 ;
 - (c) p is **not equal to** p_0 .

Input: Random samples Y_1, \dots, Y_n for Y (which are either 0 or 1) and significance level α .

Methodology:

- Compute the critical region that depends on α (and potentially \bar{Y}); and

Output:

- Reject the null hypothesis if \bar{Y} is contained in the critical region.
- Do not reject the null hypothesis (i.e., test is inconclusive) otherwise.

2 Theorem: Bernoulli, one variable

Theorem 1. (a) For the case $p > p_0$,

$$\text{critical region} = \left[p_0 + z_\alpha \sqrt{\frac{p_0(1-p_0)}{n}}, \infty \right).$$

(b) For the case $p < p_0$,

$$\text{critical region} = \left(-\infty, p_0 - z_\alpha \sqrt{\frac{p_0(1-p_0)}{n}} \right].$$

(c) For the case $p \neq p_0$,

$$\text{critical region} = \left(-\infty, p_0 - z_{\alpha/2} \sqrt{\frac{p_0(1-p_0)}{n}} \right] \cup \left[p_0 + z_{\alpha/2} \sqrt{\frac{p_0(1-p_0)}{n}}, \infty \right).$$

3 Example: Bernoulli, one variable

Your instructor suspected that the dice used by a certain magical casino has been tampered with, so that the probability p of rolling a six with these dice is strictly higher than $1/6$.

To validate his hypothesis, he played $n = 8000$ times, and he saw that six was rolled 1375 times.

Could he reject the hypothesis that the dice is fair with significance level $\alpha = 0.05$?

4 Answer: Bernoulli, one variable

Here we have

- **Null Hypothesis:** p is equal to $\frac{1}{6}$.
- **Alternative Hypothesis:** The alternative hypothesis is $p > \frac{1}{6}$.

So we have case (a), which gives us

$$\begin{aligned} z_\alpha \sqrt{\frac{p_0(1-p_0)}{n}} &= z_{0.05} \sqrt{\frac{(1/6)(5/6)}{8000}} \\ &= (1.645) \sqrt{\frac{(1/6)(5/6)}{8000}} = 0.007. \end{aligned}$$

So the critical region is

$$\begin{aligned} \left[p_0 + z_\alpha \sqrt{\frac{p_0(1-p_0)}{n}}, \infty \right) &= \left[\frac{1}{6} + 0.007, \infty \right) \\ &= \left[0.17367, \infty \right) \end{aligned}$$

Since the sample mean $\bar{Y} = \frac{1375}{8000} = 0.171875$ is not contained in the critical region (albeit barely), the test is inconclusive.

Remark 2. Note that the textbook has an alternative formula for the critical region in Theorem 1, where

$$\sqrt{\frac{p_0(1 - p_0)}{n}} \quad \text{is replaced by} \quad \sqrt{\frac{\bar{Y}(1 - \bar{Y})}{n}}.$$

These two formulas yield approximately the same numerical result in practice.

The pros for lecture version is the critical region can be calculated *without* knowing the sample mean \bar{Y} in advance, and is consistent with Table 8.4-1 in the textbook.

The cons for lecture version is the formula is slightly inconsistent with the formula in Theorem 4 from Lecture Notes 7.3.

My conclusion is that pros outweighs cons, so we will use the lecture notes formula (unless indicated otherwise).

5 Setting: Bernoulli, two variables

Object: Y_1 and Y_2 are **independent Bernoulli** random variables with **unknown parameter** p_1 and p_2 .

Hypotheses:

- **Null Hypothesis** H_0 : p_1 is equal to p_2 .
- **Alternative Hypothesis** H_1 : The alternative hypothesis can take one of these three forms:
 - (a) p_1 is **strictly greater than** p_2 ;
 - (b) p_1 is **strictly smaller than** p_2 ;
 - (c) p_1 is **not equal to** p_2 .

Input: Significance level α , n_1 many random samples for Y_1 , and n_2 many random samples for Y_2 .

Methodology:

- Compute the critical region that depends on α and the given random samples.

Output:

- Reject the hypothesis if $\bar{Y}_1 - \bar{Y}_2$ is contained in the critical region.
- Do not reject the hypothesis (i.e., test is inconclusive) otherwise.

6 Theorem: Bernoulli, two variables

Theorem 3. (a) For the case $p_1 > p_2$, the critical region is

$$\left[z_\alpha \sqrt{\left(\frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(1 - \frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}, \infty \right).$$

(b) For the case $p_1 < p_2$, the critical region is

$$\left(-\infty, -z_\alpha \sqrt{\left(\frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(1 - \frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \right).$$

(c) For the case $p_1 \neq p_2$, the critical region is

$$\left(-\infty, -z_{\alpha/2} \sqrt{\left(\frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(1 - \frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \right) \cup \left[z_{\alpha/2} \sqrt{\left(\frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(1 - \frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2} \right) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}, \infty \right).$$

7 Example: Bernoulli, two variables

There is a superstition among Dungeons and Dragons players that one needs to “pre-roll” the dice to get bad rolls out of the way.

The friendly instructor performed an experiment with one pre-rolled die and one new die.

- For the pre-rolled die, six was rolled 137 out of 800 observations.
- For the new die, that six was rolled 99 out of 600 observations.

Can he reject that the pre-rolling ritual is a mere superstition with significance level $\alpha = 0.05$?

8 Answer: Bernoulli, two variables

We are in case (c), so we have

$$\begin{aligned}\bar{Y}_1 &= \frac{137}{800}; & \bar{Y}_2 &= \frac{99}{600}; \\ z_{\alpha/2} &\sqrt{\left(\frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2}\right) \left(1 - \frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2}\right) \left(\frac{1}{n_1} + \frac{1}{n_2}\right)} \\ &= (1.96) \sqrt{\left(\frac{\frac{137}{800} + \frac{99}{600}}{800 + 600}\right) \left(1 - \frac{\frac{137}{800} + \frac{99}{600}}{800 + 600}\right) \left(\frac{1}{800} + \frac{1}{600}\right)} \\ &= 0.0016.\end{aligned}$$

So the critical region is

$$\left(-\infty, -0.0016\right] \cup \left[0.0016, \infty\right).$$

Since $\bar{Y}_1 - \bar{Y}_2 = \frac{137}{800} - \frac{99}{600} = 0.00625$ is contained in the critical region, we reject the null hypothesis.

Remark 4. Note that the textbook has an alternative formula for the critical region in Theorem 3, where

$$\sqrt{\left(\frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2}\right) \left(1 - \frac{\bar{Y}_1 + \bar{Y}_2}{n_1 + n_2}\right) \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$$

is replaced by

$$\sqrt{\frac{\bar{Y}_1(1 - \bar{Y}_1)}{n_1} + \frac{\bar{Y}_2(1 - \bar{Y}_2)}{n_2}}.$$

These two formulas yield approximately the same numerical result in practice.

The lecture version is chosen to be consistent with Table 8.4-2.

This choice is slightly inconsistent with the formula in Theorem 6 from Lecture Notes 7.3.

My conclusion is that being consistent with the textbook will reduce logistical problems in the long run; hence the decision.