

Two sample inference

ST551 Lecture 18

Charlotte Wickham

2017-11-01

Finish last time's slides

Two sample inference

Two sample setting

Setting: two **independent** samples

Y_1, \dots, Y_n i.i.d from population with c.d.f F_Y , and
 X_1, \dots, X_m i.i.d from population with c.d.f F_X

Parameter: now focus on some comparison between the two populations F_Y and F_X

Alternative view

Setting: two **independent** samples

one big dataset, measure two variables

$(Y_1, G_1), (Y_2, G_2), \dots, (Y_n, G_n), (Y_{n+1}, G_{n+1}), \dots, (Y_{n+m}, G_{n+m})$

$Y_1, \dots, Y_n, X_1, \dots, X_m$

where G is a binary *grouping* variable which indicates which population the observation came from:

$$G_i = \begin{cases} 0, & \text{observation from } Y \\ 1, & \text{observation from } X \end{cases}$$

Two views are equivalent

Depending on sampling scheme one view may seem more natural:

- I sample 40 OSU graduate students and 20 OSU undergraduate students: *fixed by design*

- Y_i = graduate student time to complete 1 mile run, $i = 1, \dots, \underline{40}$
- X_i = undergraduate student time to complete 1 mile run, $i = 1, \dots, \underline{20}$

- I sample 60 OSU students and record:

- Y_i = time to complete 1 mile run, $i = 1, \dots, 60$
- G_i = student's level (0 = graduate, 1 = undergraduate), $i = 1, \dots, 60$

$$\frac{1}{60} \sum_{i=1}^{60} 1(G_i = 0) \leftarrow \text{estimate proportion of OSU students that are graduate}$$

In second view, if we condition on the counts in each group, inference is the same as first view.

Two sample inference for difference in population means

To compare population means: $\mu_Y = E(Y_i)$, $\mu_X = E(X_i)$, we might look at their difference:

Parameters
of interest

$$\delta = \mu_Y - \mu_X$$

\mathcal{I}_δ

(In alternative view: equivalent to

$$\delta = E(Y_i | G_i = 0) - E(Y_i | G_i = 1))$$

- Estimate for δ ?
- Test for $H_0 : \delta = \delta_0$?
- Confidence interval for δ ?

Difference in sample means

It seems reasonable to use:

$$\hat{\delta} = \bar{Y} - \bar{X}$$

as a good starting point for inference on $\delta = \mu_X - \mu_Y$.

Complete worksheet (Charlotte will provide)

Properties of difference in sample means

ST551 Lecture 18

Charlotte Wickham

2017-11-01

Let Y_1, \dots, Y_n be an i.i.d sample of size n from a population with mean μ_Y and variance σ_Y^2 , and let X_1, \dots, X_m be an i.i.d sample of size m from a population with mean μ_X and variance σ_X^2 .

The samples are drawn independently of each other.

Q1 Using the Central Limit Theorem, find the approximate distribution of \bar{Y} and \bar{X} for large sample sizes.

$$\bar{Y} \sim N\left(\mu_Y, \frac{\sigma_Y^2}{n}\right)$$

$$\bar{X} \sim N\left(\mu_X, \frac{\sigma_X^2}{m}\right)$$

Q2 \bar{Y} and \bar{X} are independent. Justify this fact.

X_i to be independent of Y_j for all i, j

Q3 Derive the distribution of $\bar{Y} - \bar{X}$. A useful fact (from earlier in the quarter) is provided in the box below.

If $X \sim N(\mu_X, \sigma_X^2)$ and $Y \sim N(\mu_Y, \sigma_Y^2)$, independent of X .
Then,

$$Z = X + Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

$$\bar{Y} - \bar{X} \sim N\left(\mu_Y - \mu_X, \frac{\sigma_Y^2}{n} + \frac{\sigma_X^2}{m}\right)$$

$X - Y$
 $X + (-Y)$
 $X + (-1)Y$

Q4 Using your result from Q3, suggest a test statistic for testing the null hypothesis $H_0: \mu_Y - \mu_X = \delta_0$ that would have a (approximately) standard Normal distribution, when the null hypothesis is true.

$$\hat{\delta} = \frac{\bar{Y} - \bar{X} - \delta_0}{\sqrt{\frac{\sigma_Y^2}{n} + \frac{\sigma_X^2}{m}}} \sim N(0, 1)$$

Leads to two sample Z-test and intervals

Assume known population variances: $\text{Var}(Y_i) = \sigma_Y^2$ $\text{Var}(X_i) = \sigma_X^2$.

$$Z(\delta_0) = \frac{(\bar{Y} - \bar{X}) - \delta_0}{\sqrt{\sigma_Y^2/n + \sigma_X^2/m}}$$

Reference Distribution: If null hypothesis $H_0 : \delta = \delta_0$ is true, then

$$Z(\delta_0) \sim N(0, 1)$$

Rejection Regions:

- $H_A : \delta > \delta_0$, reject H_0 for $Z(\delta_0) > z_{1-\alpha}$
- $H_A : \delta < \delta_0$, reject H_0 for $Z(\delta_0) < z_\alpha$
- $H_A : \delta \neq \delta_0$, reject H_0 for $|Z(\delta_0)| > z_{1-\alpha/2}$

Leads to two sample Z-test and intervals

$(1 - \alpha)100\%$ Confidence interval for $\delta = \mu_Y - \mu_X$

$$(\bar{Y} - \bar{X}) \pm z_{1-\alpha/2} \sqrt{\frac{\sigma_Y^2}{n} + \frac{\sigma_X^2}{m}}$$

Next time...

What if population variances aren't known?