

Mon

This
week

Notes 12,
Lab 06

Wed

Notes 12,
Notes 13

HW 9 due
Quiz 2 corrections due

Fri

Notes 13+
review

Next
week

Quiz 03!
HW 10 due
(completion)

asynch.
activity.
Zoom off during
class time
Lab 06 due

BREAK

After
Thanksgiving

Start
regression

Project
proposal
due
(NO HW)

12: TWO-SAMPLE INFERENCE

Larsen & Marx 9.2, 9.4

Prof Amanda Luby

Today, we're going to continue our exploration of inference for a few different settings beyond inference for the mean or proportion of a population. Specifically, we're going to derive the (approximate) sampling distributions for a *difference in means* and a *difference in proportions*. We'll see that even in simple settings where we're able to make "nice" assumptions, deriving exact test statistics quickly becomes unwieldy.

1 Inference for a difference in means

One of the most common settings for inference is comparing the means for two groups. For example, if we split a random sample of patients into a *treatment* and a *placebo* group in a clinical trial, do we obtain different amounts of improvement? We could also be interested in measuring differences between existing subgroups within a population, like those who grew up within a 50 mile radius of a superfund site compared to those who did not.

Idea: $\bar{X} \sim N(\mu_X, \sigma_X^2/n)$, $\bar{Y} \sim N(\mu_Y, \sigma_Y^2/m)$ by CLT. Interested in $\mu_X - \mu_Y \rightarrow$ Depending on what we're willing to assume about our σ_X^2, σ_Y^2 , we obtain different test statistics

1.1 Assuming $\sigma_X = \sigma_Y$

Two-sample t statistic

Let $X_1, \dots, X_n \sim N(\mu_X, \sigma^2)$ and let $Y_1, \dots, Y_m \sim N(\mu_Y, \sigma^2)$, and let all X_i 's and Y_j 's be independent. Let S_X^2 and S_Y^2 be the corresponding sample variances, and let s_p^2 be the pooled variance, where

← weighted average of $S_X^2 + S_Y^2$

$$s_p^2 = \frac{(n-1)s_X^2 + (m-1)s_Y^2}{n+m-2} = \frac{\sum(X_i - \bar{X})^2 + \sum(Y_i - \bar{Y})^2}{n+m-2}$$

Then,

$$T_{n+m-2} = \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{s_p \sqrt{1/n + 1/m}} \text{ has a } T_{n+m-2} \text{ distribution}$$

Proof: Idea: $T_v = \frac{Z}{\sqrt{V/v}}$ where $Z \sim N(0,1)$ $V \sim \chi_v^2$

① Divide top & bottom by $\sigma \sqrt{1/n + 1/m}$

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{s_p^2 / v}}$$

1

② $\bar{X} \sim N(\mu_X, \frac{\sigma^2}{n})$ $\bar{Y} \sim N(\mu_Y, \frac{\sigma^2}{m}) \rightarrow \bar{X} - \bar{Y} \sim N(\mu_X - \mu_Y, \frac{\sigma^2}{n} + \frac{\sigma^2}{m}) \rightarrow$ top is a $N(0,1)$

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{S_p^2 / \nu}} \sim N(0,1)$$

Proof (cont):

$$\textcircled{3} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^2 \sim \chi_{n-1}^2 = \frac{(n-1)s_x^2}{\sigma^2} \quad \sum_{i=1}^m \left(\frac{y_i - \bar{y}}{s} \right)^2 \sim \chi_{m-1}^2 = \frac{(m-1)s_y^2}{\sigma^2}$$

• x_i 's & y_i 's are independent, so $\frac{(n-1)s_x^2}{\sigma^2} + \frac{(m-1)s_y^2}{\sigma^2}$ are also indep.
 • Sum of χ^2 RV's, we obtain a χ^2

$$\frac{(n-1)s_x^2}{\sigma^2} + \frac{(m-1)s_y^2}{\sigma^2} \sim \chi_{n+m-2}^2$$

$$\textcircled{4} \text{ can write } S_p^2 / \sigma^2 \text{ as } \left(\frac{(n-1)s_x^2}{\sigma^2} + \frac{(m-1)s_y^2}{\sigma^2} \right) \cdot \frac{1}{n+m-2} = \frac{\chi_{n+m-2}^2}{n+m-2}$$

$$\rightarrow T_{n+m-2} = \frac{Z}{\sqrt{\chi_{n+m-2}^2 / (n+m-2)}} \text{ so } T_{n+m-2} \text{ has a } T_{n+m-2} \text{ df.}$$

Form for a $(1 - \alpha)\%$ confidence interval:

$$\text{for } \mu_X - \mu_Y: (\bar{X} - \bar{Y}) \pm t_{\alpha/2, n+m-2} \cdot s_p \sqrt{1/n + 1/m}$$

Rejection regions for α -level tests:

$$H_0: \mu_X = \mu_Y \text{ let } t = \frac{\bar{X} - \bar{Y}}{s_p \sqrt{1/n + 1/m}}$$

$$H_1: \mu_X > \mu_Y$$

$$H_1: \mu_X < \mu_Y$$

$$H_1: \mu_X \neq \mu_Y$$

$$\text{reject if } t \geq t_{\alpha, n+m-2}$$

$$t \leq -t_{\alpha, n+m-2}$$

$$t \leq -t_{\alpha/2, n+m-2}$$

$$t \geq t_{\alpha/2, n+m-2}$$

1.2 Assuming $\sigma_X \neq \sigma_Y$

\rightarrow 'Best guess': S_X, S_Y

Welch's 2-sample t statistic

$$W = \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}}$$

$$\text{Ideally, } = \frac{Z}{\sqrt{1/n + 1/m}}$$

has an approximate T_ν distribution, where

$$\nu = \frac{(\frac{S_X^2}{n} + \frac{S_Y^2}{m})^2}{\frac{1}{n-1}(\frac{S_X^2}{n})^2 + \frac{1}{m-1}(\frac{S_Y^2}{m})^2}, \text{ rounded to the nearest integer}$$

but in general
don't know
df of
 $\sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}$
:

Proof(ish):

Form for a $(1 - \alpha)\%$ confidence interval:

Rejection regions for α -level tests:

2 Inference for a difference in proportions

Suppose that m Bernoulli trials have resulted in X successes, and suppose n Bernoulli trials have resulted in Y successes; where all trials are independent. A common test is:

$$H_0 : p_x = p_y$$

$$H_1 : p_x \neq p_y$$

2.1 Deriving the GLRT

2.1.1 Approximation Using the CLT

Form for a $(1 - \alpha)\%$ confidence interval:

Rejection regions for α -level tests: