

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  $\sigma_X^2, \sigma_Y^2$ , 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  $\sigma \sqrt{1/n + 1/m}$

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{s_p^2 / v_2}}$$

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):

$$\textcircled{3} \sum_i \left( \frac{x_i - \bar{x}}{\sigma} \right)^2 \sim \chi^2_{n-1} = \frac{(n-1)S_x^2}{\sigma^2} \quad \sum_i \left( \frac{y_i - \bar{y}}{\sigma} \right)^2 \sim \chi^2_{m-1} = \frac{(m-1)S_y^2}{\sigma^2}$$

•  $X_i$ '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  $\chi^2$  RV's, we obtain a  $\chi^2$

$$\frac{(n-1)S_x^2}{\sigma^2} + \frac{(m-1)S_y^2}{\sigma^2} \sim \chi^2_{n+m-2}$$

$$\textcircled{4} \text{ can write } S_p^2 / \sigma^2 \text{ as } \left( \frac{(n-1)S_x^2}{\sigma^2} + \frac{(m-1)S_y^2}{\sigma^2} \right) \cdot \frac{1}{n+m-2} = \frac{\chi^2_{n+m-2}}{n+m-2}$$

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

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  $\alpha$ -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

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

$$\text{Ideally, } = \frac{Z}{\sqrt{\nu/\nu}}$$

has an approximate  $T_\nu$  distribution, where

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

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

Divide top & bottom of W  
by  $\sqrt{\sigma_X^2/n + \sigma_Y^2/m}$

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

$N(0,1)$  by same reasoning

$$\text{Ideally, } \frac{S_X^2/n + S_Y^2/m}{\sigma_X^2/n + \sigma_Y^2/m} = \frac{\nu}{\nu}$$

↑  
not generally true

$$\frac{S_x^2/n + S_y^2/m}{\sigma_x^2/n + \sigma_y^2/m} = \frac{V}{v} \quad \text{where } v \sim \chi_v^2$$

$$\text{LHS} \rightarrow \sum \left( \frac{x_i - \bar{x}}{\sigma_x} \right)^2 + \sum \left( \frac{y_i - \bar{y}}{\sigma_y} \right)^2$$

sums of squared  $N(0,1) \rightarrow \chi^2$

Proof(ish):

$$\text{Rearrange: } \frac{S_x^2}{n} + \frac{S_y^2}{m} = \left( \frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m} \right) \cdot \frac{V}{v}$$

By equating means and variances of both sides,

$$v = \frac{\left( \sigma_x^2/n + \sigma_y^2/m \right)^2}{\sigma_x^4/n^2(n-1) + \sigma_y^4/m^2(m-1)} \quad \theta = \frac{\sigma_x^2}{\sigma_y^2} \quad \text{and divide by } \sigma_y^4$$

$$= \frac{\left( \frac{1}{n} \cdot \frac{\sigma_x^2}{\sigma_y^2} + \frac{1}{m} \right)^2}{\frac{1}{n^2(n-1)} \left( \frac{\sigma_x^2}{\sigma_y^2} \right)^2 + \frac{1}{m^2(m-1)}} = \frac{\left( \frac{1}{n} \theta + \frac{1}{m} \right)^2}{\frac{1}{n^2(n-1)} \theta^2 + \frac{1}{m^2(m-1)}} \quad \text{multiply by } n^2$$

$$= \frac{\left( \theta + \frac{n}{m} \right)^2}{\frac{1}{n-1} \theta^2 + \frac{1}{m-1} \left( \frac{n}{m} \right)^2}$$

← but don't know  $\theta$ , plug in  $\hat{\theta} = \frac{S_x^2}{S_y^2}$

$$= \frac{\left( \frac{S_x^2}{S_y^2} + \frac{n}{m} \right)^2}{\frac{1}{n-1} \left( \frac{S_x^2}{S_y^2} \right)^2 + \frac{1}{m-1} \left( \frac{n}{m} \right)^2}$$

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

$$(\bar{x} - \bar{y}) \pm t_{\alpha/2, v} \cdot \sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{m}}$$

← based on  $v$  test statistic we just derived

Rejection regions for  $\alpha$ -level tests:

## 2 Inference for a difference in proportions

$$\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right)$$

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

$$H_0 : p_x = p_y$$

$$X \sim \text{Binom}(n, p_x)$$

$$H_1 : p_x \neq p_y$$

$$Y \sim \text{Binom}(m, p_y)$$

$$H_0: P_X = P_Y$$

$$H_1: P_X \neq P_Y$$

## 2.1 Deriving the GLRT

$$\Omega_0 = \{(P_X, P_Y) : 0 \leq P_X = P_Y \leq 1\}$$

$$\Omega_1 = \{(P_X, P_Y) : 0 \leq P_X \leq 1, 0 \leq P_Y \leq 1, P_X \neq P_Y\}$$

Since  $X$  &  $Y$  independent,

$$L(P_X, P_Y) : \prod_{i,j} f_X(X_i; P_X) f_Y(Y_j; P_Y) = P_X^x (1-P_X)^{n-x} P_Y^y (1-P_Y)^{m-y}$$

Under  $H_0$ :

$$P_X = P_Y = P_0$$

$$L = P_0^{x+y} (1-P_0)^{n+m-x-y}$$

$$\ln L = (x+y) \ln(P_0) + (n+m-x-y) \ln(1-P_0)$$

$$\frac{\partial \ln L}{\partial P_0} = \frac{(x+y)}{P_0} + \frac{(n+m-x-y)}{1-P_0} \cdot -1 = 0$$

$$\frac{(x+y)}{P_0} = \frac{(n+m-x-y)}{1-P_0}$$

$$\frac{1-P_0}{P_0} = \frac{n+m-(x+y)}{(x+y)}$$

$$\frac{1-P_0}{P_0} = \frac{(n+m)}{(x+y)} - \frac{(x+y)}{(x+y)}$$

$$\frac{1}{P_0} = \frac{n+m}{x+y}$$

$$\hat{P}_0 = \frac{x+y}{n+m}$$

← pooled success proportion

$$\lambda = \frac{\max_{\theta \in \Omega_0} L(\theta)}{\max_{\theta \in \Omega_1} L(\theta)}$$

under  $H_1$ : need to maximize  $P_X$  &  $P_Y$  separately

→ simplifies to the 1 sample MLE's

$$\hat{P}_X = \frac{x}{n} \quad \hat{P}_Y = \frac{y}{m}$$

Back to GLRT:

$$\lambda = \frac{\left(\frac{x+y}{n+m}\right)^{x+y} \left(1 - \frac{x+y}{n+m}\right)^{n+m-x-y}}{\left(\frac{x}{n}\right)^x \left(1 - \frac{x}{n}\right)^{n-x} \left(\frac{y}{m}\right)^y \left(1 - \frac{y}{m}\right)^{m-y}}$$

Intuitive, but ugly. RV's  $X$  &  $Y$  in numerator, denominator, exponents, etc.

hard to derive  $\lambda$  directly → approximate

### 2.1.1 Approximation Using the CLT

$$\text{Under } H_0, \quad \frac{X}{n} - \frac{Y}{m} \sim N\left(0, \frac{P(1-P)}{n} + \frac{P(1-P)}{m}\right)$$

plug in MLE (under  $H_0$ ) for  $p$  to obtain

$$Z = \frac{X/n - Y/m}{\sqrt{\frac{P_P(1-P_P)}{n} + \frac{P_P(1-P_P)}{m}}} \sim N(0,1)$$

$$P_P = \frac{x+y}{n+m}$$

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

$$\left(\frac{X}{n} - \frac{Y}{m}\right) \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{P}_X(1-\hat{P}_X)}{n} + \frac{\hat{P}_Y(1-\hat{P}_Y)}{m}} \rightarrow \text{Probably better, does not assume } P_X = P_Y$$

Rejection regions for  $\alpha$ -level tests:

$$\cdot \sqrt{\frac{\hat{P}_P(1-\hat{P}_P)}{n} + \frac{\hat{P}_P(1-\hat{P}_P)}{m}}$$

$$\cdot \sqrt{\frac{1}{4n} + \frac{1}{4m}}$$

→ most conservative