

Paired Data

ST551 Lecture 20

Charlotte Wickham

2017-11-06

Review of last week's t-tests

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: Difference in population means $\mu_Y - \mu_X$

Equal variance two sample t-test

Assume $\sigma_X^2 = \sigma_Y^2$.

$$t(\delta_0) = \frac{(\bar{Y} - \bar{X}) - \delta_0}{\sqrt{s_p^2 \left(\frac{1}{n} + \frac{1}{m} \right)}}$$

Compare to $t_{(n+m-2)}$.

Welch's t-test

σ_X^2 not necessarily equal to σ_Y^2 .

$$t(\delta_0) = \frac{(\bar{Y} - \bar{X}) - \delta_0}{\sqrt{\frac{s_Y^2}{n} + \frac{s_X^2}{m}}}$$

Compare to t_v , where

$$v = \frac{(s_Y^2/n + s_X^2/m)^2}{\frac{s_Y^4}{n^2(n-1)} + \frac{s_X^4}{m^2(m-1)}}$$

In both cases

If df is the appropriate degrees of freedom for the test.

Rejection regions:

- $H_A : \mu_Y - \mu_X > 0$: Reject H_0 for $t(\delta_0) > t_{(df),1-\alpha}$
- $H_A : \mu_Y - \mu_X < 0$: Reject H_0 for $t(\delta_0) < t_{(df),\alpha}$
- $H_A : \mu_Y - \mu_X \neq 0$: Reject H_0 for $|t(\delta_0)| > t_{(df),1-\alpha/2}$

Confidence intervals:

$$\bar{Y} - \bar{X} \pm t_{(df),1-\alpha/2} \text{SE}_{\bar{Y}-\bar{X}}$$

Paired Data

Setting

Two **dependent** samples

Y_1, \dots, Y_n i.i.d from population with c.d.f F_Y , and
 X_1, \dots, X_n i.i.d from population with c.d.f F_X
 \downarrow same size as sample from F_Y

Observations come in pairs:

$$(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)$$

with joint distribution F_{YX} . Observations are somehow *matched*.

$Cov(Y_i, X_i) = \sigma_{YX}$ and $Cov(Y_i, X_j) = 0$ for all $i \neq j$.

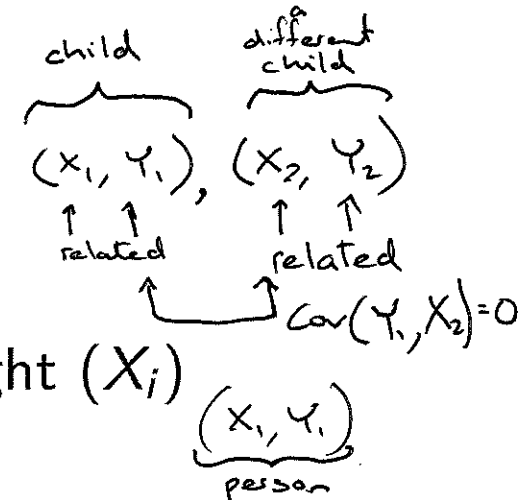
Parameter: Difference in population means $\mu_Y - \mu_X$

\hookrightarrow often
~~some~~ by
some larger
unit
• child
• person 8

Examples

We've already seen examples like this:

- **midterm** Mother's IQ (Y_i) and Father's IQ (X_i)
- **homework** Current weight (Y_i) and desired weight (X_i)



We did one sample t-tests on the differences $Y_i - X_i$. This works, but why?

Sampling Distribution of difference in sample means

Consider

$$D_i = Y_i - X_i \quad i = 1, \dots, n$$

CLT says

$$\frac{\bar{D} - E(D_i)}{\sqrt{\text{Var}(D_i)/n}} \sim N(0, 1)$$

What are \bar{D} , $E(D_i)$ and $\text{Var}(D_i)$?

What are \bar{D} , $E(D_i)$ and $Var(D_i)$?

$$\begin{aligned}\bar{D} &= \frac{1}{n} \sum_{i=1}^n D_i = \frac{1}{n} \sum_{i=1}^n (Y_i - X_i) \\ &= \frac{1}{n} \sum_{i=1}^n Y_i - \frac{1}{n} \sum_{i=1}^n X_i \\ &= \bar{Y} - \bar{X} \\ &= \text{difference in sample averages}\end{aligned}$$

Your Turn :

$$\begin{aligned}E(D_i) &= E(Y_i - X_i) \\ &\begin{array}{l} \uparrow \\ \text{one} \\ \text{difference} \end{array} = E(Y_i) - E(X_i) \\ &= \mu_Y - \mu_X \\ &= \text{difference in population means}\end{aligned}$$

linearity
of expectation

A Z-test for paired data

If σ_Y^2 , σ_X^2 , and σ_{YX} are known.

Hypothesis Test for $H_0 : \delta = \delta_0 = \underbrace{\mu_Y - \mu_X}_{\text{hypothesized difference}}$

Test Statistic:

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

Reference Distribution: Under H_0 , $Z(\delta_0) \sim N(0, 1)$

Rejection Region:

- $H_A : \delta > \delta_0$: 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 $|z(\delta_0)| > z_{1-\alpha/2}$

$$\begin{aligned}\text{Var}(D_i) &= \text{Var}(Y_i - X_i) \\ &= \text{Var}(Y_i) + \text{Var}(X_i) \\ &\quad - 2 \text{Cov}(Y_i, X_i) \\ &= \sigma_y^2 + \sigma_x^2 - 2\sigma_{yx}\end{aligned}$$

Fact:

$$\text{Var}(W - Z) =$$

$$\text{Var}(W) + \text{Var}(Z) - 2 \text{Cov}(W, Z)$$

A Z-test for paired data

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

$\sigma_{YX} > 0$
"positive correlation"
denominator is
smaller

Notice the test statistic is just like a two sample Z-test, but with a correction to $Var(\bar{Y} - \bar{X})$ for the correlation between Y_i and X_i .

What if population variances and covariances aren't known?

Plug in estimates for σ_Y^2 , σ_X^2 and σ_{YX} .

Sample covariance:

$$\sigma_{YX} = E(Y - E(Y))(X - E(X))$$

$$\hat{\sigma}_{YX} = s_{YX} = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})$$

‡

is an unbiased estimate of σ_{YX} .

Plugging in the estimates gives the estimated variance of $\bar{Y} - \bar{X}$:

estimated variance

$$\widehat{\text{Var}}(\bar{D}) = \frac{s_Y^2}{n} + \frac{s_X^2}{n} - 2 \frac{s_{YX}}{n}$$

In the unpaired case : $\sigma_{YX} = 0$

$$E(s_{YX}) = 0$$

Compare to estimated $\text{Var}(D_i)$

$$\begin{aligned} s_D^2 &= \frac{1}{n-1} \sum_{i=1}^n (D_i - \bar{D})^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left((Y_i - X_i) - (\bar{Y} - \bar{X}) \right)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left(\underbrace{(Y_i - \bar{Y})}_a - \underbrace{(X_i - \bar{X})}_b \right)^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n \left[(Y_i - \bar{Y})^2 + (X_i - \bar{X})^2 - 2(Y_i - \bar{Y})(X_i - \bar{X}) \right] \\ &= s_Y^2 + s_X^2 - 2s_{XY} \\ \frac{s_D^2}{n} &= \frac{s_Y^2}{n} + \frac{s_X^2}{n} - \frac{2s_{XY}}{n} \end{aligned}$$

sample variance of the differences

$(a-b)^2 = a^2 + b^2 - 2ab$

Paired data t-test

$$t(\delta_0) = \frac{(\bar{Y} - \bar{X}) - \delta_0}{\sqrt{\frac{s_Y^2}{n} + \frac{s_X^2}{n} - 2\frac{s_{YX}}{n}}} = \frac{\bar{D} - \delta_0}{\sqrt{\frac{s_D^2}{n}}}$$

one sample
t-statistic
 D_i

If differences are Normal, $t(\delta_0)$ has **exactly** a t-distribution with $n - 1$ degrees of freedom when the null hypothesis is true.

For paired samples:

1. Take differences $D_i = X_i - Y_i$
2. Perform a one-sample hypothesis test for the population mean difference $\mu_D = \mu_Y - \mu_X$

That is, do a one-sample t-test on the differences

This is equivalent to estimating the population covariance and appropriately adjusting the denominator of the two-sample t-test to take this covariance into account

D_i obs differences

μ_D mean difference

$Y_i - X_i$

$\mu_Y - \mu_X$

Statistical Summaries

① We estimate the mean difference between IQ for mothers and fathers of gifted children is -3 .

"mean difference"

$-\frac{1}{2}$ direction
is not obvious

② We estimate the mean IQ of mothers to be 3 points higher than the mean IQ of fathers, of gifted children.

"difference means"

② is preferred