

Paired Data

ST551 Lecture 20

Charlotte Wickham

2017-11-06

Review of last week's t-tests

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:

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

Paired Data

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

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 $\mu_Y - \mu_X$

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

$$D_i = Y_i - X_i$$

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 $Var(D_i)$?

A Z-test for paired data

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

Hypothesis Test for $H_0 : \delta = \delta_0$

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

A Z-test for paired data

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

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:

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

is an unbiased estimate of σ_{YX} .

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

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

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

$$s_D^2 =$$

Paired data t-test

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

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

Summary

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