

HTs and CIs for differences

Difference in means and proportions

4/16/25

Housekeeping

- Project proposals due tonight!

Recap

- Test and CI for a single mean
 - If we know σ , use standard Normal distribution
 - If we don't know σ and only have access to s , use t distribution

Difference in two proportions

Now suppose we have samples of binary outcomes from two different populations.

Difference of two proportions

Suppose we have two populations 1 and 2, and want to either estimate the value of or conduct a test for the **difference in population proportions**: $p_1 - p_2$

- We have samples of size n_1 and n_2 from each population
- Reasonable point estimate: $\hat{p}_{1,obs} - \hat{p}_{2,obs}$
- We will obtain the sampling distribution of the difference of two sample proportions
- **Now that we have two populations, conditions for CLT will look slightly different!**

Sampling dist. of difference of two proportions

- In order to use CLT approximation, we have to ensure conditions are met:
 1. **Independence (extended)**: data are independent within *and* between groups
 2. **Success-failure (extended)**: success-failure conditions holds for *both* groups
 - $n_1 p_1 \geq 10, n_1(1 - p_1) \geq 10, n_2 p_2 \geq 10, \text{ and } n_2(1 - p_2) \geq 10$
- If above hold, then:

$$\hat{p}_1 - \hat{p}_2 \sim N \left(p_1 - p_2, \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}} \right)$$

where p_1 and p_2 are the population proportions

Confidence interval for difference in proportions

If we want to obtain a $\gamma \times 100\%$ CI for $p_1 - p_2$, that means we don't know the value of $p_1 - p_2$!

- Like in the case of the CI for a single proportion, we will use our observed proportions to check success-failure
- **Success-failure condition for CI for difference in proportions:**
 - $n_1 \hat{p}_{1,obs} \geq 10$ and $n_1(1 - \hat{p}_{1,obs}) \geq 10$
 - $n_2 \hat{p}_{2,obs} \geq 10$ and $n_2(1 - \hat{p}_{2,obs}) \geq 10$
- Then our formula for the CI is the same as before:

point. est \pm critical val. $\times \widehat{SE} =$

$$(\hat{p}_{1,obs} - \hat{p}_{2,obs}) \pm z_{(1+\gamma)/2}^* \sqrt{\frac{\hat{p}_{1,obs}(1 - \hat{p}_{1,obs})}{n_1} + \frac{\hat{p}_{2,obs}(1 - \hat{p}_{2,obs})}{n_2}}$$

Diff. props CI example: offshore drilling

A survey asked 827 randomly sampled registered voters in California: Do you support or oppose drilling for oil and natural gas off the Coast of California? We have the following distribution of responses separated by whether the respondent graduated from college:

position	no	yes	total
oppose	126	180	306
support	132	154	286
total	258	334	592

Let's obtain a 95% CI via the CLT for the difference in the proportion of college and non-college Californians who support offshore drilling.

Diff. props CI example (cont.)

Let population 1 be college attendees, and population 2 be non-college attendees. We want a 95% CI for $p_1 - p_2$, where p_i is the proportion of population i who support offshore drilling.

- Obtain useful statistics

- $n_1 = 334, n_2 = 258$
- $\hat{p}_{1,obs} = \frac{154}{334} = 0.461$
- $\hat{p}_{2,obs} = \frac{132}{258} = 0.512$

- Check conditions for CLT.

- Independence (extended)? Randomly sampled
- Success-failure (extended)?
 - $n_1\hat{p}_{1,obs} = 154 \geq 10$
 - $n_1(1 - \hat{p}_{1,obs}) = 180 \geq 10$
 - $n_2\hat{p}_{2,obs} = 132 \geq 10$
 - $n_2(1 - \hat{p}_{2,obs}) = 126 \geq 10$

- Since both conditions are met, we can proceed with the CLT.

Diff. props CI example (cont.)

Collect the components of CI:

- Point estimate
 - Critical value (code)
 - SE or \widehat{SE}
- $\hat{p}_{1,obs} - \hat{p}_{2,obs} = 0.461 - 0.512 = -0.051$
 - $z_{0.975}^* = \text{qnorm}(0.975, 0, 1) \approx 1.96$
 - $\widehat{SE} = \sqrt{\frac{0.461(1-0.461)}{334} + \frac{0.512(1-0.512)}{258}} = 0.041$
- So putting it all together, our 95% CI is:
 $-0.051 \pm 1.96 \times 0.041 = (-0.131, 0.029)$
 - Interpret!

Hypothesis test for difference in proportions

Hypothesis tests for difference in proportions in this class will take the form:

$$H_0 : p_1 - p_2 = 0$$

$$H_A : p_1 - p_2 \neq 0$$

$$\text{or } <$$

$$\text{or } >$$

- For success-failure condition for difference in two proportions, we don't have null-hypothesized values for p_1 or p_2 .
- So how do we check the condition??

Pooled proportion

- Since $H_0 : p_1 = p_2$, then under the null $\hat{p}_{1,obs}$ and $\hat{p}_{2,obs}$ come from the *same* population
- So under this null, we use a special proportion called the **pooled proportion**:

$$\hat{p}_{pooled} = \frac{\text{total \# of successes from both samples}}{\text{combined sample size}}$$

- This is the best estimate of both p_1 and p_2 if $H_0 : p_1 = p_2$ is true!
- For this reason, use \hat{p}_{pooled} to verify success-failure conditions for HT for difference of proportions:
 - $n_1 \hat{p}_{pooled} \geq 10$ and $n_1(1 - \hat{p}_{pooled}) \geq 10$
 - $n_2 \hat{p}_{pooled} \geq 10$ and $n_2(1 - \hat{p}_{pooled}) \geq 10$

Hypothesis test (cont.)

3. Obtain null distribution

- If conditions satisfied, then we know the sampling distribution of $\hat{p}_1 - \hat{p}_2$
- To obtain the **null distribution** we assume $H_0 : p_1 - p_2 = 0$ is true and we use \hat{p}_{pooled} to estimate p_1 and p_2 to approximate standard error under the null:

$$\begin{aligned}\hat{p}_1 - \hat{p}_2 &\dot{\sim} N \left(p_1 - p_2, \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \right) && \text{(CLT)} \\ &\dot{\sim} N \left(0, \underbrace{\sqrt{\frac{\hat{p}_{pooled}(1-\hat{p}_{pooled})}{n_1} + \frac{\hat{p}_{pooled}(1-\hat{p}_{pooled})}{n_2}}}_{\widehat{SE}_0} \right) && (H_0)\end{aligned}$$

Hypothesis test (cont.)

Obtain test-statistic:

$$z = \frac{\text{point estimate} - \text{null value}}{\text{SE}} \approx \frac{(\hat{p}_{1,obs} - \hat{p}_{2,obs}) - 0}{\widehat{\text{SE}}_0}$$

- To obtain p-value, we want $\Pr(Z \geq z)$ and/or $\Pr(Z \leq z)$ where $Z \sim N(0, 1)$
 - Obtain using `pnorm(z, 0, 1)`

Diff. props HT example: offshore drilling (again)

Using the same data as before, let's answer the following question:

Do the data provide strong evidence at the 0.05 level that the proportion of college graduates who support off-shore drilling in California is different than that of non-college graduates?

- Let p_1 and p_2 be defined as before.
- Define hypotheses
 - $H_0 : p_1 - p_2 = 0$ and $H_A : p_1 - p_2 \neq 0$

Diff. props HT example (cont.)

Obtain pooled proportion, and use it to check conditions for CLT.

- $\hat{p}_{pooled} = \frac{154+132}{334+258} = \frac{286}{592} = 0.483$
- Conditions
 - Independence (extended): random sample
 - Success-failure (extended):
 - $n_1\hat{p}_{pooled} = 334 \times 0.483 = 161.32 \geq 10$
 - $n_1(1 - \hat{p}_{pooled}) = 334 \times (1 - 0.483) = 172.68 \geq 10$
 - $n_2\hat{p}_{pooled} = 258 \times 0.483 = 124.61 \geq 10$
 - $n_2(1 - \hat{p}_{pooled}) = 258 \times (1 - 0.483) = 133.39 \geq 10$
- Since conditions are met, we can proceed with CLT-based test!

Diff. props HT example (cont.)

- Find the null distribution for $\hat{p}_1 - \hat{p}_2$.

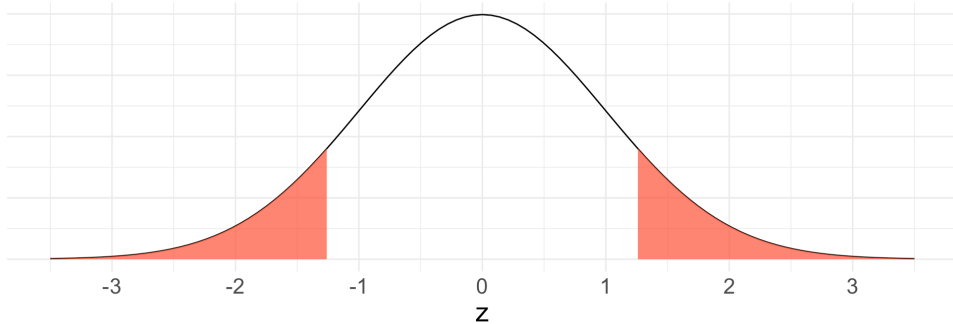
$$\hat{p}_1 - \hat{p}_2 \sim N \left(0, \sqrt{\frac{0.483(1 - 0.483)}{334} + \frac{0.483(1 - 0.483)}{258}} = 0.041 \right)$$

- Set up calculation for test statistic

$$z = \frac{(\hat{p}_{1,obs} - \hat{p}_{2,obs}) - 0}{\widehat{SE}_0} = \frac{(0.461 - 0.512) - 0}{0.041} = -1.244$$

Diff. props HT example (cont.)

Draw picture and write code for p-value



- p-value calculation:
 - $\Pr(Z \leq z) + \Pr(Z \geq -z)$
 - `2 * pnorm(-1.244, 0, 1) = 0.2134996`

Make a decision and conclusion in context.

- Since our p-value is greater than 0.05, we fail to reject H_0 . The data do not provide strong evidence of a difference between the proportions of college graduates and non-college graduates who support off-shore drilling among California voters.

Difference in two means

Difference in two means

We still have two populations, but the variable of interest is quantitative (i.e. not binary).

We are interested in learning about the difference in the means of each population.

- Let μ_1 and μ_2 represent the population means for the two populations 1 and 2
- Samples of size n_1 and n_2 from each population, respectively
- **Conditions for CLT**
 1. **Independence** (extended): need data within *and* between the two groups
 - e.g. the two data sets come from independent random samples or from a randomized experiment
 2. **Normality**: we need to check for approximate normality for *both* groups separately

CLT for difference in two sample means

If CLT conditions met, the distribution of difference in sample means is:

$$\bar{X}_1 - \bar{X}_2 \sim N \left(\mu_1 - \mu_2, \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

where n_1 and n_2 are the sample sizes.

- Remember, we often do not know σ_1 nor σ_2
- In practice, will have to estimate with s_1 and s_2 and use the t -distribution

CI for difference in two means

If the conditions hold, then our usual formula for $\gamma \times 100\%$ CI still holds:

$$\text{point estimate} \pm \text{critical value} \times \text{SE}$$

Point estimate: $\bar{x}_{1,obs} - \bar{x}_{2,obs}$

If σ_1 and σ_2 known:

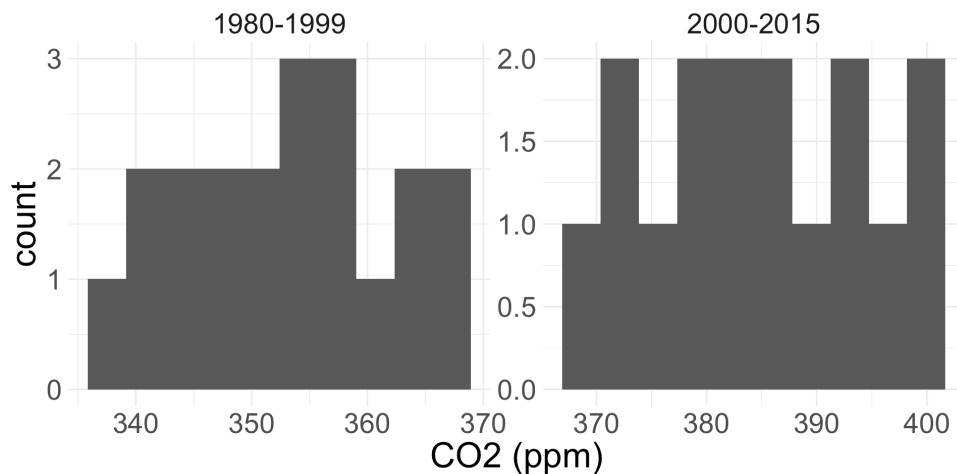
- $\text{SE} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$
- Critical value: $z_{(1+\gamma)/2}^*$
 - $(1 + \gamma)/2$ percentile of $N(0, 1)$

If σ_1 and σ_2 unknown:

- $\widehat{\text{SE}} \approx \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$
- critical value: $t_{df, (1+\gamma)/2}^*$
 - $(1 + \gamma)/2$ percentile of t_{df}
 - $df = \min\{n_1 - 1, n_2 - 1\}$

Diff. means CI example: CO2 concentrations

- The Mauna Loa Observatory in Hawaii monitors atmospheric solar, atmospheric, and meteorological parameters
- We have data on annual atmospheric CO2 concentrations from 1980-2015.
- **We will obtain a 90% confidence interval for the difference between the average atmospheric CO2 levels (ppm) from years 2000-2015 and years 1980-1999.**



group	n	xbar	s
1980-1999	20	353.12	9.0
2000-2015	16	385.02	9.9

Diff. means CI example (cont.)

Define parameters.

- Let μ_1 be the average CO2 levels from 2000-2015 and μ_2 the average CO2 levels from 1980-1999.
- Want to obtain a 90% CI for $\mu_1 - \mu_2$
 - Note: could also do $\mu_2 - \mu_1$ (interpretation just changes slightly)

Check conditions for CLT.

- Independence (extended): most likely violated because CO2 levels are probably dependent across time. BUT let's proceed with caution anyway.
- Normality: $n_1 = 16 < 30$ and $n_2 = 20 < 30$. But since histograms don't reveal outliers, Normality condition appears met.

Diff. means CI example (cont.)

Collect components for CI:

- Point estimate
- Critical value (code)
- SE or \widehat{SE}
- $\bar{x}_{1,obs} - \bar{x}_{2,obs} = 385.02 - 353.12 = 31.9$
- Since we don't know σ_1 nor σ_2 , need to use t -distribution
 - Degrees of freedom =
 $\min\{16 - 1, 20 - 1\} = 15$
 - $t_{0.95}^* = \text{qt}(0.95, \text{df} = 15) = 1.75$
- $\widehat{SE} = \sqrt{\frac{9.9^2}{16} + \frac{9^2}{20}} = 3.19$
- Put it all together:

$$\text{point est.} \pm \text{crit. val} \times \text{SE} = 31.9 \pm 1.75 \times 3.19 = (26.3175, 37.4825)$$

Diff. means CI example (cont.)

Interpret our CI of (26.3175, 37.4825) in context!

Hypothesis test for difference in means

Now suppose we're interested in testing for the difference between μ_1 and μ_2 .

- $H_0 : \mu_1 - \mu_2 = 0$ versus $H_A : \mu_1 - \mu_2 \neq 0$ (or $<$ or $>$)
- Same conditions as in CI are necessary for CLT-based inference!
 1. Independence (extended)
 2. Normality condition for **both** groups
- If CLT met, then under H_0 , the **null distribution** is

[Math Processing Error]

Test statistic for difference in means

Test-statistic is of form:

$$\frac{\text{point est.} - \text{null value}}{\text{SE}_0}$$

If σ_1, σ_2 known, our test-statistic is:

$$z = \frac{(\bar{x}_{1,obs} - \bar{x}_{2,obs}) - 0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

If σ_1, σ_2 unknown, our test-statistic is

$$t = \frac{(\bar{x}_{1,obs} - \bar{x}_{2,obs}) - 0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim t_{df}$$

$$df = \min\{n_1 - 1, n_2 - 1\}$$

Diff. means HT example: CO2

Now let's test if the mean CO2 level in 2000-2015 was greater than that mean CO2 level in 1980-1999 at the 0.05 level using CLT.

1. $H_0 : \mu_1 - \mu_2 = 0$ versus $H_A : \mu_1 > \mu_2$, where μ_1 and μ_2 were defined previously
2. Let $\alpha = 0.05$
3. Conditions for CLT are same as before (proceed with caution)

Diff. means HT example (cont.)

Obtain test-statistic and p-value.

- Find the value of the test-statistic and its distribution

$$t = \frac{(385.02 - 353.12) - 0}{3.19} = 10 \sim t_{15}$$

- Write code for p-value (optionally draw picture)
- p-value = `1 - pt(10, df = 15)` = 2.4984491×10^{-8} (tiny!)

