

Lecture 9- Spring 2024

Villas-Boas

Lecture 9 EEP 118 Spring 2024

Please work on this notebook as your daily assignment after Lecture 10 (I will finish lecturing this material in Lecture 10 since we ran out of time in Lecture 9). See the media gallery video for Lecture 10 for the material needed to solve this notebook. I will also post this notebook solved in bcourses after lecture.

I will also post Video 4 on this notebook as a "How to EEP Series" on how to test for the equality of the proportions of yes answers in two different populations.

```
In [1]: # Load the 'pacman' package
library(pacman)
#packages to use load them now using the pacman "manager"
p_load(dplyr, readr)
#Another great feature of p_load(): if you try to load a package that is not
p_load(ggplot2)

#set scientific display off, thank you Roy
options(scipen=999)

# Loading packages
pacman::p_load(lfe, lmtest, haven, sandwich, tidyverse,psych)
# lfe for running fixed effects regression
# lmtest for displaying robust SE in output table
# haven for loading in dta files
# sandwich for producing robust Var-Cov matrix
# tidyverse for manipulating data and producing plots
# psych for using describe later on
```

```
In [2]: #-----
#1. Read in data
#-----
my_data <- read_dta("data2024.dta")
head(my_data)
```

A tibble: 6 × 9

timestamp	went2class	soccerfan	correct1	correct2	correctboth	numberCorr
<chr>	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<d
2/6/2024 10:53:19	yes	yes	1	1	1	
2/6/2024 11:01:25	yes	yes	1	0	0	
2/6/2024 11:11:24	yes	yes	1	1	1	
2/6/2024 11:11:28	yes	yes	1	0	0	
2/6/2024 11:11:52	yes	yes	1	0	0	
2/6/2024 11:11:53	yes	no	1	1	1	

```
In [3]: #describe data
describe(my_data,skew = FALSE)

# 108 total responses
```

A psych: 9 × 8

	vars	n	mean	sd	min	max	range
	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
timestamp*	1	108	53.8055556	30.9529212	1	107	106 2.
went2class*	2	108	1.9351852	0.2473466	1	2	1 0.
soccerfan*	3	108	1.5185185	0.5019864	1	2	1 0.
correct1	4	108	0.9722222	0.1651017	0	1	1 0.
correct2	5	108	0.5555556	0.4992206	0	1	1 0.
correctboth	6	108	0.5555556	0.4992206	0	1	1 0.
numberCorrect	7	108	1.5277778	0.5546462	0	2	2 0.
went2Class	8	108	0.9351852	0.2473466	0	1	1 0.
isSoccerFan	9	108	0.5185185	0.5019864	0	1	1 0.

Select coming to class subsample and how many came to class?

```
In [4]: #select coming to class subsample
my_dataClass <- filter(my_data, went2Class == 1)

#how many are there that came to class?
summarise(my_dataClass, trimmed_count = n())
```

```
#answer 101 out of 108 came to class
```

A tibble: 1 × 1

trimmed_count
<int>
101

```
In [5]: # what is the proportion of correct both for those coming to class?  
mean(my_dataClass$correctboth)
```

0.574257425742574

```
In [6]: #Lets call that phat_1  p hat for having come to class"  
  
phat_1<-mean(my_dataClass$correctboth)  
phat_1
```

0.574257425742574

```
In [7]: #Lets call that Y_1  the number oe people answering correctly among those c  
  
Y_1<-mean(my_dataClass$correctboth)*nrow(my_dataClass)  
Y_1  
  
#and N1  
N_1<-nrow(my_dataClass)
```

58

Select not coming to class subsample and how many did not come to class?

```
In [8]: #select not coming to class subsample  
my_dataNotClass <- filter(my_data, went2Class == 0)  
  
#how many are there that did not come to class?  
summarise(my_dataNotClass, trimmed_count = n())  
  
#answer 7 out of 108 did not come to class
```

A tibble: 1 × 1

trimmed_count
<int>
7

```
In [9]: #Lets call that phat_2  the proportion of people that answered correctly bc  
#come to class " p hat_2
```

```
phat_2<-mean(my_dataNotClass$correctboth)
phat_2
```

0.285714285714286

```
In [10]: #Lets call that Y_2 the number of people answering correctly among those r
Y_2<-mean(my_dataNotClass$correctboth)*nrow(my_dataNotClass)
Y_2

#and
N_2<-nrow(my_dataNotClass)
```

2

Question:

Test whether the proportion of answering both correctly for those that came to class (p_1) is statistically equal to to the one of those not coming to class (p_2) at the 10% significance against an alternative that $p_1 > p_2$

that is,

against the alternative that those coming to class have a larger proportion p_1 of answering correctly than those not coming to class p_2 .

Recall the 5 step-procedure for hypothesis testing.

let $D = p_1 - p_2$ be the difference in proportions in the population

Step 1: $D=0$ null, alternative $D>0$ one sided alternative.

Step 2: construct the test stat

Let $\hat{D} = \hat{p}_1 - \hat{p}_2$

Testing equality of proportions

The test statistic for testing the difference in two population proportions, that is, for testing the null hypothesis $H_0 : p_1 - p_2 = 0$ is:

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1 - \hat{p}) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

where:

$$\hat{p} = \frac{Y_1 + Y_2}{n_1 + n_2}$$

the proportion of "successes" in the two samples combined.

where $\hat{p}_1 = \frac{Y_1}{n_1}$ and $\hat{p}_2 = \frac{Y_2}{n_2}$ and we interpret $\hat{p} = \frac{Y}{n}$ as the overall proportion in the sample of correct answers, and \hat{p}_1 the proportion correct for the first group (coming to class), and \hat{p}_2 the proportion correct for the second group (not coming to class).

The sample estimate for the population p_1 is $\hat{p}_1 = \frac{Y_1}{n_1}$ The sample estimate for the population p_2 is $\hat{p}_2 = \frac{Y_2}{n_2}$

If under the null hypothesis $p_1 = p_2 = p$ then the sample estimate for the population p is $\hat{p} = \frac{Y}{n} = \frac{Y_1 + Y_2}{n_1 + n_2}$

Under the null of the population proportion being p , then the variance is $p(1-p)$

Recall that the sample estimate for the variance of the sample average of sample size N is $\frac{\hat{p}(1-\hat{p})}{N}$

We know that under the null hypothesis that $p_1 = p_2 = p$ then

$$\text{var}(\hat{p}_1) = \frac{\hat{p}(1-\hat{p})}{N_1}$$

and

$$\text{var}(\hat{p}_2) = \frac{\hat{p}(1-\hat{p})}{N_2}$$

So the variance(\hat{D}) is the variance $(\hat{p}_1 - \hat{p}_2) = \text{var}(\hat{p}_1) + \text{var}(\hat{p}_2)$

And the $se(\hat{D}) = \sqrt{\text{variance}(\hat{D})}$

which is equivalent to

$$se(\hat{D}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{N_1} + \frac{\hat{p}(1-\hat{p})}{N_1}} = \sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{N_1} + \frac{1}{N_1}\right)}$$

like in the formula in the Testing Equality of Proportions given above

```
In [11]: #get phat
phat<-(Y_1+Y_2)/(N_1+N_2)

phat
```

0.5555555555555555

```
In [12]: #to construct statistic
#get denominator of Z (slide 30 in Lecture 9 notes)
temp<-phat*(1-phat)*(1/N_1+1/N_2)
denom<-sqrt(temp)

denom
```

0.194211373302486

```
In [13]: #numerator is
D_hat<-phat_1-phat_2

D_hat
```

0.288543140028289

```
In [14]: #get statistic
#t= Dhat/denom

t_testStatValue<-D_hat/denom
t_testStatValue
```

1.48571700576402

you get that the value of the t is 1.48757

Step 3: given significance level $\alpha=10\%$
get the Critical T value for 106 degrees of
freedom for one sided

get one sided DF=106, 10% critical
value

go to the t table

you will see that

`t_crit<-1.2896`

```
In [15]: #or get it in R by
          tcritical<-qt(0.1,106,lower.tail=FALSE)
          tcritical

          #this gives you the t value that has 10 percent mass higher than it,
          #which is what we want in the one sided test
```

1.28958924084563

Step 4: Compare t with t critical

If the value of t from step 2 is less than the critical value from step 3 then we fall in the non rejection area,

that is, if $t < t_{critical}$ then cannot reject the null of step 1

If the t is greater than the tcritical then we fall in the rejection area, that is, if $t > t_{critical}$ then we reject the null of step 1

What do you find?

since t greater than critical value we reject the null

Step 5: Conclude.

In this case, we reject that both proportions are the same at the 10% significance level against an alternative that p_1 is greater than p_2 , that is,

against the alternative that those coming to class have a larger proportion p_1 of answering correctly than those not coming to class p_2 .

oh Yes!

Final Task, and repeat as daily assignment at some point, please

See solutions in Lecture10.R in Bcourses

Test at the 10% significance level whether you can reject the null that the average number of correct answers is the same for the sample of those coming to class as the sample of those not coming to class

Let us look at the number of correct answers by the respondent, given by the variable `my_data$numberCorrect`

How would you compute the average number of correct answers in the sample?

```
In [17]: mean(my_data$numberCorrect)
```

```
1.52777777777778
```

Solution: answer is obtained by running in the command line above the following command

```
mean(my_data$numberCorrect)
```

answer is 1.52778

Test at the 10% significance level whether you can reject the null that the average number of correct answers is the same for the sample of those coming to class as the sample of those not coming to class

Let N_c be the population mean number of correct answers for those coming to class

Let N_{nc} be the population mean number of correct answers for those not coming to class

Let $D = N_c - N_{nc}$ be the difference of the two

```
N_c <-  
mean(my_dataClass$numberCorrect)
```

```
N_nc <-  
mean(my_dataNotClass$numberCorrect)
```

Step 1 specify the null and the alternative

Answer : $H_0: N_c = N_{nc}$ or $N_c - N_{nc} = 0$ or $D = 0$

$H_a: D$ not equal to zero

step 2: Construct the t stat

Things you need to run in the line above

what is the average number of correct answers for those coming to class
 \bar{N}_{bar_c}

```
Nhat_c<-mean(my_dataClass$numberCorrect)
```

answer is 1.574257

what is the average number of correct answers for those coming to class
 \bar{N}_{bar_nc}

```
Nhat_nc<-mean(my_dataNotClass$numberCorrect)
```

it is 0.8571429

then get the variances and std errors for average Number correct for those coming to class $\text{Var_}\bar{N}_{bar_c}$ and $\text{se_}\bar{N}_{bar_c}$

and for average number correct for those not coming to class $\text{Var_}\bar{N}_{bar_nc}$ and $\text{se_}\bar{N}_{bar_nc}$

get variance and std error for number of correct answers for those coming to Class

```
var_Nhat_c<-var(my_dataClass$numberCorrect)/nrow(my_dataClass)
```

```
se_Nhat_c<-sqrt(var_Nhat_c)
```

get variance and std error for number of correct answers for those not coming to class

```
var_Nhat_nc<-var(my_dataNotClass$numberCorrect)/nrow(my_dataNotClass)
```

```
se_Nhat_nc<-sqrt(var_Nhat_nc)
```

now we can do hypothesis testing using the five steps

step 1

$H_0: D=0$

$H_a: D \text{ not equal to } 0$

where D =population Mean of Number correct for those coming to class minus

population Mean of number correct of those not coming to class

step 2: construct the t statistic

$$t = (\hat{D} - D_{\text{under null}}) / \text{se}(\hat{D})$$

$D_{\text{under null}}=0$

$$\hat{D} = \hat{N}_{\text{hat}_c} - \hat{N}_{\text{hat}_{nc}}$$

$se(\hat{D}) = \text{square root of } [\text{var}(\hat{N}_c) + \text{var}(\hat{N}_{nc})]$

step 2: construct the t statistic of daily assignment (DA)

$t_{DA} = (\hat{D} - \text{Dunder null}) / se(\hat{D})$

where Dunder null=0

$\hat{D}_{DA} = \hat{N}_c - \hat{N}_{nc}$

```
DhatDA<-Nhat_c-Nhat_nc
```

$se(\hat{D}) = \text{square root of } [\text{var}(\hat{N}_c) + \text{var}(\hat{N}_{nc})]$

```
se_DhatDA<-sqrt(var_Nhat_c+var_Nhat_nc)
```

t of daily assignment t_{DA}

```
t_DA <- DhatDA / se_DhatDA
```

```
t_DA
```

answer is 2.086796

do it in the code cell below

```
In [20]: Nhat_c<-mean(my_dataClass$numberCorrect)

#answer is 1.574257

#what is the average number of correct answers for those coming to class Nba

Nhat_nc<-mean(my_dataNotClass$numberCorrect)
```

```

#it is 0.8571429

#then get the variances and std errors for average Number correct for those

#get variance and std error for number of correct answers for those coming t
var_Nhat_c<-var(my_dataClass$numberCorrect)/nrow(my_dataClass)

se_Nhat_c<-sqrt(var_Nhat_c)

#get variance and std error for number of correct answers for those not comi
var_Nhat_nc<-var(my_dataNotClass$numberCorrect)/nrow(my_dataNotClass)

se_Nhat_nc<-sqrt(var_Nhat_nc)
#step 2: construct the t statistic of daily assignment (DA)

#t_DA = (Dhat-Dunder null) / se(Dhat)

#where Dunder null=0

#DhatDA=Nhat_c - Nhat_nc

DhatDA<-Nhat_c-Nhat_nc

#se(Dhat)=square root of [ var(Nhat_c)+var(Nhat_nc) ]

se_DhatDA<-sqrt(var_Nhat_c+var_Nhat_nc)

#t of daily assignment t_DA

t_DA <- DhatDA / se_DhatDA

t_DA

```

2.08679568747402

Step 3: given significance level 10% what is the critical value we know its 1.659 from above.

what about for significance level 5%?

find it in the cell below

answer is obtained by the commands typed (first to get, then to show in output)

```
t_criticalDA<- qt(0.025, 106, lower.tail=FALSE)
```

```
t_criticalDA
```

```
In [21]: t_criticalDA<- qt(0.025, 106, lower.tail=FALSE)
t_criticalDA
```

```
1.9825972617655
```

Step 4: compare the absolute value of your tDA with the absolute value of the tcritical value. what do you conclude in Step 5?

Answer step 4 : the absolute value of the t stat is 2.08 is greater than the absolute value of the critical t of 1.9825972617655

so we land in the rejection area of this two-sided test

Answer step 5: so we reject the null that the mean of number of correct answers for those coming to class is the same as the mean of the number of correct answers for those not coming to class, against the alternative of the means being different at the 5 percent significance level.

based on this test, reject that it does not matter coming to class in terms of the Null that the number of correct answers in the population would have the same mean regardless of coming or not coming to class !

So, do come to class :-)

You can't write the entire code now to solve the task in the cell below

```
In [ ]:
```

the end

```
In [ ]:
```