



Program Code: J620-002-4:2020

Program Name: FRONT-END SOFTWARE DEVELOPMENT

Title : Exe12 - Confidence Intervals NHANES Exercise

Name: Chuay Xiang Ze

IC Number: 021224070255

Date : 11/7/2023

Introduction : Learning how to calculate for and using CI to solve questions.

Conclusion : Managed to complete tasks using CI.

Exercise 1: Confidence Intervals - NHANES

This exercise, we are going to practice on how to load data, clean/manipulate a dataset, and construct a confidence interval for the difference between two population proportions and means.

We will use the 2015-2016 wave of the NHANES data for our analysis.

For our population proportions, we will analyze the difference of proportion between female and male smokers. The column that specifies smoker and non-smoker is "SMQ020" in our dataset.

For our population means, we will analyze the difference of mean of body mass index within our female and male populations. The column that includes the body mass index value is "BMXBMI".

Additionally, the gender is specified in the column "RIAGENDR".

In [9]:

```
import pandas as pd
import numpy as np
import matplotlib
matplotlib.use('Agg')
import seaborn as sns
%matplotlib inline
import matplotlib.pyplot as plt
import statsmodels.api as sm
```

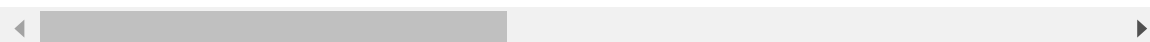
In [11]:

```
url = "./nhanes_2015_2016.csv"
da = pd.read_csv(url)
da
```

Out[11]:

	SEQN	ALQ101	ALQ110	ALQ130	SMQ020	RIAGENDR	RIDAGEYR	RIDRETH1	DMDCL
0	83732	1.0	NaN	1.0	1	1	62	3	
1	83733	1.0	NaN	6.0	1	1	53	3	
2	83734	1.0	NaN	NaN	1	1	78	3	
3	83735	2.0	1.0	1.0	2	2	56	3	
4	83736	2.0	1.0	1.0	2	2	42	4	
...
5730	93695	2.0	2.0	NaN	1	2	76	3	
5731	93696	2.0	2.0	NaN	2	1	26	3	
5732	93697	1.0	NaN	1.0	1	2	80	3	
5733	93700	NaN	NaN	NaN	1	1	35	3	
5734	93702	1.0	NaN	2.0	2	2	24	3	

5735 rows × 28 columns



Investigating and Cleaning Data

Create a new column named 'SMQ020x' and store data from column 'SMQ020' with following replacements:

- 1 to "Yes"
- 2 to "No"
- 7 to NaN
- 9 to NaN

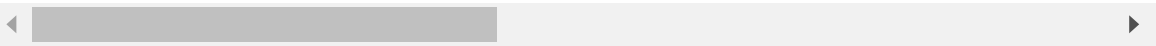
In [13]:

```
da['SMQ020x'] = da['SMQ020']
reset = {1: "Yes", 2: "No", 7: None, 9: None}
da = da.replace({"SMQ020x": reset})
da['SMQ020x']
da
```

Out[13]:

	SEQN	ALQ101	ALQ110	ALQ130	SMQ020	RIAGENDR	RIDAGEYR	RIDRETH1	DMDC
0	83732	1.0	NaN	1.0	1	1	62	3	
1	83733	1.0	NaN	6.0	1	1	53	3	
2	83734	1.0	NaN	NaN	1	1	78	3	
3	83735	2.0	1.0	1.0	2	2	56	3	
4	83736	2.0	1.0	1.0	2	2	42	4	
...	
5730	93695	2.0	2.0	NaN	1	2	76	3	
5731	93696	2.0	2.0	NaN	2	1	26	3	
5732	93697	1.0	NaN	1.0	1	2	80	3	
5733	93700	NaN	NaN	NaN	1	1	35	3	
5734	93702	1.0	NaN	2.0	2	2	24	3	

5735 rows × 29 columns



Create a new column named 'RIAGENDRx' and store data from column 'RIAGENDR' with following replacements:

- 1 to "Male"
- 2 to "Female"

In [14]:

```
da['RIAGENDRx'] = da['RIAGENDR']
reset = {1: "Male", 2: "Female"}
da = da.replace({"RIAGENDRx": reset})
da['RIAGENDRx']
```

Out[14]:

```
0      Male
1      Male
2      Male
3    Female
4    Female
...
5730  Female
5731    Male
5732  Female
5733    Male
5734  Female
Name: RIAGENDRx, Length: 5735, dtype: object
```

Drop all NAs from both `SMQ020x` & `RIAGENDRx` and store into a new dataframe named 'dx'. Plot the following crosstab using `pd.crosstab` library.

In [19]:

```
dx = da[da['SMQ020x'] != None]
pd.crosstab(dx['SMQ020x'], dx['RIAGENDRx'])
```

Out[19]:

RIAGENDRx	Female	Male
SMQ020x		
No	2066	1340
Yes	906	1413

Replace `dx['SMQ020x']` "Yes" to 1 and "No" to 0.

In [29]:

```
reset = {"Yes": 1, "No": 0}
dx = dx.replace({"SMQ020x": reset})
dx = dx[['SMQ020x', 'RIAGENDRx']]
dx
```

Out[29]:

	SMQ020x	RIAGENDRx
0	1.0	Male
1	1.0	Male
2	1.0	Male
3	0.0	Female
4	0.0	Female
...
5730	1.0	Female
5731	0.0	Male
5732	1.0	Female
5733	1.0	Male
5734	0.0	Female

5735 rows × 2 columns

Calculate the 'mean' and 'size' and store into a new dataframe called dz

In [30]:

```
dz = dx[['SMQ020x', 'RIAGENDRx']].groupby('RIAGENDRx').agg(['mean', 'size'])
dz
```

Out[30]:

	SMQ020x	
	mean	size
RIAGENDRx		
Female	0.304845	2976
Male	0.513258	2759

Constructing Confidence Intervals

Now that we have the population proportions of male and female smokers, we can begin to calculate confidence intervals. From lecture, we know that the equation is as follows:

$$\text{Best Estimate} \pm \text{Margin of Error}$$

Where the *Best Estimate* is the **observed population proportion or mean** from the sample and the *Margin of Error* is the **t-multiplier**.

The equation to create a 95% confidence interval can also be shown as:

$$\text{Population Proportion or Mean} \pm (t - \text{multiplier} * \text{Standard Error})$$

The Standard Error is calculated differently for population proportion and mean:

$$\text{Standard Error for Population Proportion} = \sqrt{\frac{\text{Population Proportion} * (1 - \text{Population Proportion})}{\text{Number Of Observations}}}$$

$$\text{Standard Error for Mean} = \frac{\text{Standard Deviation}}{\sqrt{\text{Number Of Observations}}}$$

Lastly, the standard error for difference of population proportions and means is:

$$\text{Standard Error for Difference of Two Population Proportions Or Means} = \sqrt{SE_{\text{Proportion 1}}^2 + SE_{\text{Proportion 2}}^2}$$

Difference of Two Population Proportions

Calculate the standard error for female

In [31]:

```
import math
p = dz['SMQ020x']['mean']['Female']
n = dz['SMQ020x']['size']['Female']
sef = math.sqrt(p*(1-p)/n)
sef
```

Out[31]:

0.008438475404634192

Calculate the standard error for male

In [38]:

```
p = dz['SMQ020x']['mean']['Male']
n = dz['SMQ020x']['size']['Male']
sem = math.sqrt((p*(1-p))/n)
sem
```

Out[38]:

0.009515714829783395

Calculate the difference between these two Standard Errors

In [39]:

```
dse = math.sqrt(sef**2+sem**2)
dse
```

Out[39]:

0.012718360581316123

Calculate the confidence Interval

In [40]:

```
ci = dz['SMQ020x']['mean']['Male'] - dz['SMQ020x']['mean']['Female']
x = ci-2*dse
y = ci+2*dse
print(x,y)
print(dse)
```

```
0.1829763204770033 0.23384976280226777
0.012718360581316123
```

Difference of Two Population Means

Now we look into the differences between 2 population means

In [41]:

```
da["BMXBMI"].head()
```

Out[41]:

```
0    27.8
1    30.8
2    28.8
3    42.4
4    20.3
Name: BMXBMI, dtype: float64
```

In [44]:

```
x = da[['BMXBMI', 'RIAGENDRx']].groupby('RIAGENDRx').agg(['mean', 'std', 'size'])
x
```

Out[44]:

	BMXBMI		
	mean	std	size
RIAGENDRx			
Female	29.939946	7.753319	2976
Male	28.778072	6.252568	2759

Calculate the Standard Error for Mean for both female and male

In [45]:

```
s = x['BMXBMI']['std']['Female']
n = x['BMXBMI']['size']['Female']
semf = s/math.sqrt(n)
s = x['BMXBMI']['std']['Male']
n = x['BMXBMI']['size']['Male']
semm = s/math.sqrt(n)
print(semf,semm)
```

0.14212522940758335 0.11903715722332033

Calculate the difference between 2 Standard Error for Mean

In [46]:

```
dsem = math.sqrt(semf**2+ semm**2)
dsem
```

Out[46]:

0.18538992862064455

The difference between two means for male and female

In [47]:

```
dbm = math.sqrt(x['BMXBMI']['mean']['Female']**2 + x['BMXBMI']['mean']['Male']**2)
dbm
```

Out[47]:

41.52803607359479

Calculate the confidence interval between two population means

In [48]:

```
cipm = x['BMXBMI']['mean']['Female'] - x['BMXBMI']['mean']['Male']
y = cipm-2*dsem
z = cipm+2*dsem
print(y,z)
```

0.7910936830856763 1.5326533975682544

In []: