

# assignment2

April 19, 2021

## 1 Assignment 2

For this assignment you'll be looking at 2017 data on immunizations from the CDC. Your datafile for this assignment is in [assets/NISPUF17.csv](#). A data users guide for this, which you'll need to map the variables in the data to the questions being asked, is available at [assets/NIS-PUF17-DUG.pdf](#). **Note: you may have to go to your Jupyter tree (click on the Coursera image) and navigate to the assignment 2 assets folder to see this PDF file).**

### 1.1 Question 1

Write a function called `proportion_of_education` which returns the proportion of children in the dataset who had a mother with the education levels equal to less than high school (<12), high school (12), more than high school but not a college graduate (>12) and college degree.

*This function should return a dictionary in the form of (use the correct numbers, do not round numbers):*

```
{"less than high school":0.2,  
 "high school":0.4,  
 "more than high school but not college":0.2,  
 "college":0.2}
```

```
[2]: import pandas as pd  
  
def proportion_of_education():  
    df = pd.read_csv('assets/NISPUF17.csv', index_col=0)  
    ndf= df[df['YEAR']==2017]  
    # your code goes here  
    COUNT =ndf["EDUC1"].groupby(ndf['EDUC1']).count()  
    #print(COUNT)  
  
    lths = COUNT.iloc[0]  
    hs=COUNT.iloc[1]  
    mths= COUNT.iloc[2]  
    clg=COUNT.iloc[3]  
    total = len(ndf['EDUC1'])  
  
    result = {"less than high school":lths/total,  
             "high school":hs/total,
```

```

        "more than high school but not college":mths/total,
        "college":clg/total}
    # YOUR CODE HERE
    return result
    #raise NotImplementedError()
proportion_of_education()

```

```

[2]: {'less than high school': 0.10202002459160373,
      'high school': 0.172352011241876,
      'more than high school but not college': 0.24588090637625154,
      'college': 0.47974705779026877}

```

```

[3]: assert type(proportion_of_education())==type({}), "You must return a dictionary.
      ↪"
      assert len(proportion_of_education()) == 4, "You have not returned a dictionary_
      ↪with four items in it."
      assert "less than high school" in proportion_of_education().keys(), "You have_
      ↪not returned a dictionary with the correct keys."
      assert "high school" in proportion_of_education().keys(), "You have not_
      ↪returned a dictionary with the correct keys."
      assert "more than high school but not college" in proportion_of_education().
      ↪keys(), "You have not returned a dictionary with the correct keys."
      assert "college" in proportion_of_education().keys(), "You have not returned a_
      ↪dictionary with the correct keys."

```

## 1.2 Question 2

Let's explore the relationship between being fed breastmilk as a child and getting a seasonal influenza vaccine from a healthcare provider. Return a tuple of the average number of influenza vaccines for those children we know received breastmilk as a child and those who know did not.

*This function should return a tuple in the form (use the correct numbers):*

(2.5, 0.1)

```

[11]: def average_influenza_doses():

    # YOUR CODE HERE

    df = pd.read_csv('assets/NISPUF17.csv', index_col=0)
    ndf= df[df['YEAR']==2017]

    bf =ndf["CBF_01"]==1
    nbf =ndf["CBF_01"]==2

    abf = ndf["P_NUMFLU"].where(bf ).dropna().mean()
    anbf = ndf["P_NUMFLU"].where(nbf).dropna().mean()
    final = (abf,anbf)
    return final

```

```
raise NotImplementedError()
average_influenza_doses()
```

[11]: (1.8799187420058687, 1.5963945918878317)

[5]: `assert len(average_influenza_doses())==2, "Return two values in a tuple, the_`  
`→first for yes and the second for no."`

### 1.3 Question 3

It would be interesting to see if there is any evidence of a link between vaccine effectiveness and sex of the child. Calculate the ratio of the number of children who contracted chickenpox but were vaccinated against it (at least one varicella dose) versus those who were vaccinated but did not contract chicken pox. Return results by sex.

*This function should return a dictionary in the form of (use the correct numbers):*

```
{"male":0.2,
 "female":0.4}
```

Note: To aid in verification, the `chickenpox_by_sex()['female']` value the autograder is looking for starts with the digits 0.0077.

```
[21]: def chickenpox_by_sex():

    df = pd.read_csv('assets/NISPUF17.csv', index_col=0)
    ndf= df[df['YEAR']==2017]

    fndf = ndf[["SEX", "P_NUMVRC", "HAD_CPOX"]].dropna()

    #print(len(fndf))

    #number of male person who has cpox and vaccinated
    '''print(fndf["SEX"].unique())
    print(fndf["P_NUMVRC"].unique())
    print(fndf["HAD_CPOX"].unique())
    '''

    #number of male person who has cpox and vaccinated

    had_cpox_vac_m = fndf.where((ndf["SEX"]==1) & (ndf["P_NUMVRC"]!=0) &_
    →(ndf["HAD_CPOX"]==1))
    m_had = len(had_cpox_vac_m.dropna())

    ##number of male person who didn't have cpox and vaccinated
    not_had_cpox_vac_m = fndf.where((ndf["SEX"]==1) & (ndf["P_NUMVRC"]!=0) &_
    →(ndf["HAD_CPOX"]==2))
    m_n_had = len(not_had_cpox_vac_m.dropna())
```

```

#number of female person who has cpox and vaccinated
had_cpox_vac_m = fndf.where((ndf["SEX"]==2) & (ndf["P_NUMVRC"]!=0) &
→(ndf["HAD_CPOX"]==1))
f_had = len(had_cpox_vac_m.dropna())

#number of female person who didn't have cpox and vaccinated
not_had_cpox_vac_m = fndf.where((ndf["SEX"]==2) & (ndf["P_NUMVRC"]!=0) &
→(ndf["HAD_CPOX"]==2))
f_n_had = len(not_had_cpox_vac_m.dropna())

r_m = m_had / m_n_had

r_f = f_had / f_n_had

final = {"male": r_m , "female": r_f}
#print (len(final))
# YOUR CODE HERE
return final
raise NotImplementedError()

chickenpox_by_sex()

```

[21]: {'male': 0.009675583380762664, 'female': 0.0077918259335489565}

[22]: `assert len(chickenpox_by_sex())==2, "Return a dictionary with two items, the`  
`→first for males and the second for females."`

## 1.4 Question 4

A correlation is a statistical relationship between two variables. If we wanted to know if vaccines work, we might look at the correlation between the use of the vaccine and whether it results in prevention of the infection or disease [1]. In this question, you are to see if there is a correlation between having had the chicken pox and the number of chickenpox vaccine doses given (varicella).

Some notes on interpreting the answer. The `had_chickenpox_column` is either 1 (for yes) or 2 (for no), and the `num_chickenpox_vaccine_column` is the number of doses a child has been given of the varicella vaccine. A positive correlation (e.g.,  $\text{corr} > 0$ ) means that an increase in `had_chickenpox_column` (which means more no's) would also increase the values of `num_chickenpox_vaccine_column` (which means more doses of vaccine). If there is a negative correlation (e.g.,  $\text{corr} < 0$ ), it indicates that having had chickenpox is related to an increase in the number of vaccine doses.

Also, `pval` is the probability that we observe a correlation between `had_chickenpox_column` and `num_chickenpox_vaccine_column` which is greater than or equal to a particular value occurred by chance. A small `pval` means that the observed correlation is highly unlikely to occur by chance. In this case, `pval` should be very small (will end in  $e-18$  indicating a very small number).

[1] This isn't really the full picture, since we are not looking at when the dose was given. It's possible that children had chickenpox and then their parents went to get them the vaccine. Does

this dataset have the data we would need to investigate the timing of the dose?

```
[23]: def corr_chickenpox():
import scipy.stats as stats
import numpy as np
import pandas as pd

'''# this is just an example dataframe
df=pd.DataFrame({"had_chickenpox_column":np.random.randint(1,3,size=(100)),
                "num_chickenpox_vaccine_column":np.random.
→randint(0,6,size=(100))})

# here is some stub code to actually run the correlation
corr, pval=stats.
→pearsonr(df["had_chickenpox_column"],df["num_chickenpox_vaccine_column"])

# just return the correlation
#return corr
'''
# YOUR CODE HERE
df = pd.read_csv('assets/NISPUF17.csv',index_col=0)
ndf= df[df['YEAR']==2017]

zndf =  ndf[["P_NUMVRC", "HAD_CPOX"]].where(ndf["HAD_CPOX"]!=77).dropna()
corr, pval=stats.pearsonr(zndf["HAD_CPOX"],zndf["P_NUMVRC"])

return corr

raise NotImplementedError()

corr_chickenpox()
```

```
[23]: 0.07044873460147986
```

```
[24]: assert -1<=corr_chickenpox()<=1, "You must return a float number between -1.0_
→and 1.0."
```