

Assignment 4

May 14, 2019

You are currently looking at **version 1.1** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the [Jupyter Notebook FAQ](#) course resource.

```
In [ ]: import pandas as pd
import numpy as np
from scipy.stats import ttest_ind
```

1 Assignment 4 - Hypothesis Testing

This assignment requires more individual learning than previous assignments - you are encouraged to check out the [pandas documentation](#) to find functions or methods you might not have used yet, or ask questions on [Stack Overflow](#) and tag them as pandas and python related. And of course, the discussion forums are open for interaction with your peers and the course staff.

Definitions: * A *quarter* is a specific three month period, Q1 is January through March, Q2 is April through June, Q3 is July through September, Q4 is October through December. * A *recession* is defined as starting with two consecutive quarters of GDP decline, and ending with two consecutive quarters of GDP growth. * A *recession bottom* is the quarter within a recession which had the lowest GDP. * A *university town* is a city which has a high percentage of university students compared to the total population of the city.

Hypothesis: University towns have their mean housing prices less effected by recessions. Run a t-test to compare the ratio of the mean price of houses in university towns the quarter before the recession starts compared to the recession bottom. (price_ratio=quarter_before_recession/recession_bottom)

The following data files are available for this assignment: * From the [Zillow research data site](#) there is housing data for the United States. In particular the datafile for [all homes at a city level](#), City_Zhvi_AllHomes.csv, has median home sale prices at a fine grained level. * From the Wikipedia page on college towns is a list of [university towns in the United States](#) which has been copy and pasted into the file university_towns.txt. * From Bureau of Economic Analysis, US Department of Commerce, the [GDP over time](#) of the United States in current dollars (use the chained value in 2009 dollars), in quarterly intervals, in the file gdp1ev.xls. For this assignment, only look at GDP data from the first quarter of 2000 onward.

Each function in this assignment below is worth 10%, with the exception of run_ttest(), which is worth 50%.

```

In [2]: # Use this dictionary to map state names to two letter acronyms
        states = {'OH': 'Ohio', 'KY': 'Kentucky', 'AS': 'American Samoa', 'NV': 'Nevada', 'WY':

In [4]: def get_list_of_university_towns():
        '''Returns a DataFrame of towns and the states they are in from the
        university_towns.txt list. The format of the DataFrame should be:
        DataFrame( [ ["Michigan", "Ann Arbor"], ["Michigan", "Yipsilanti" ]],
        columns=["State", "RegionName" ] )

        The following cleaning needs to be done:

        1. For "State", removing characters from "[" to the end.
        2. For "RegionName", when applicable, removing every character from " (" to the end.
        3. Depending on how you read the data, you may need to remove newline character '\n'

        data = []
        state = None
        state_towns = []
        with open('university_towns.txt') as file:
            for line in file:
                thisLine = line[:-1]
                if thisLine[-6:] == '[edit]':
                    state = thisLine[:-6]
                    continue
                if '(' in line:
                    town = thisLine[:thisLine.index('(')-1]
                    state_towns.append([state,town])
                else:
                    town = thisLine
                    state_towns.append([state,town])
                data.append(thisLine)
        df = pd.DataFrame(state_towns,columns = ['State','RegionName'])
        return df

In [6]: def get_recession_start():
        '''Returns the year and quarter of the recession start time as a
        string value in a format such as 2005q3'''
        gdplev = pd.ExcelFile('gdplev.xls')
        gdplev = gdplev.parse("Sheet1", skiprows=219)
        gdplev = gdplev[['1999q4', 9926.1]]
        gdplev.columns = ['Quarter','GDP']
        for i in range(2, len(gdplev)):
            if (gdplev.iloc[i-2][1] > gdplev.iloc[i-1][1]) and (gdplev.iloc[i-1][1] > gdplev

In [7]: def get_recession_end():
        '''Returns the year and quarter of the recession end time as a
        string value in a format such as 2005q3'''

```

```

gdplev = pd.ExcelFile('gdplev.xls')
gdplev = gdplev.parse("Sheet1", skiprows=219)
gdplev = gdplev[['1999q4', 9926.1]]
gdplev.columns = ['Quarter', 'GDP']
start = get_recession_start()
start_index = gdplev[gdplev['Quarter'] == start].index.tolist()[0]
gdplev=gdplev.iloc[start_index:]
for i in range(2, len(gdplev)):
    if (gdplev.iloc[i-2][1] < gdplev.iloc[i-1][1]) and (gdplev.iloc[i-1][1] < gdplev
        return gdplev.iloc[i][0]

```

```

In [8]: def get_recession_bottom():
    '''Returns the year and quarter of the recession bottom time as a
    string value in a format such as 2005q3'''

    gdplev = pd.ExcelFile('gdplev.xls')
    gdplev = gdplev.parse("Sheet1", skiprows=219)
    gdplev = gdplev[['1999q4', 9926.1]]
    gdplev.columns = ['Quarter', 'GDP']
    start = get_recession_start()
    start_index = gdplev[gdplev['Quarter'] == start].index.tolist()[0]
    end = get_recession_end()
    end_index = gdplev[gdplev['Quarter'] == end].index.tolist()[0]
    gdplev=gdplev.iloc[start_index:end_index+1]
    bottom = gdplev['GDP'].min()
    bottom_index = gdplev[gdplev['GDP'] == bottom].index.tolist()[0]-start_index
    return gdplev.iloc[bottom_index]['Quarter']

```

```

In [9]: def convert_housing_data_to_quarters():
    '''Converts the housing data to quarters and returns it as mean
    values in a dataframe. This dataframe should be a dataframe with
    columns for 2000q1 through 2016q3, and should have a multi-index
    in the shape of ["State", "RegionName"].

    Note: Quarters are defined in the assignment description, they are
    not arbitrary three month periods.

    The resulting dataframe should have 67 columns, and 10,730 rows.
    '''

    hdata = pd.read_csv('City_Zhvi_AllHomes.csv')
    hdata = hdata.drop(hdata.columns[[0]+list(range(3,51))],axis=1)
    hdata2 = pd.DataFrame(hdata[['State', 'RegionName']])
    for year in range(2000, 2016):
        hdata2[str(year)+'q1'] = hdata[[str(year)+'-01', str(year)+'-02', str(year)+'-03']
        hdata2[str(year)+'q2'] = hdata[[str(year)+'-04', str(year)+'-05', str(year)+'-06']
        hdata2[str(year)+'q3'] = hdata[[str(year)+'-07', str(year)+'-08', str(year)+'-09']
        hdata2[str(year)+'q4'] = hdata[[str(year)+'-10', str(year)+'-11', str(year)+'-12']

```

```

year = 2016
hdata2[str(year)+'q1'] = hdata[[str(year)+'-01',str(year)+'-02',str(year)+'-03']].me
hdata2[str(year)+'q2'] = hdata[[str(year)+'-04',str(year)+'-05',str(year)+'-06']].me
hdata2[str(year)+'q3'] = hdata[[str(year)+'-07',str(year)+'-08']].mean(axis=1)
hdata2 = hdata2.replace({'State':states})
hdata2 = hdata2.set_index(['State','RegionName'])
return hdata2

```

```
In [10]: def run_ttest():
```

```

    '''First creates new data showing the decline or growth of housing prices
    between the recession start and the recession bottom. Then runs a ttest
    comparing the university town values to the non-university towns values,
    return whether the alternative hypothesis (that the two groups are the same)
    is true or not as well as the p-value of the confidence.

```

```

    Return the tuple (different, p, better) where different=True if the t-test is
    True at a p<0.01 (we reject the null hypothesis), or different=False if
    otherwise (we cannot reject the null hypothesis). The variable p should
    be equal to the exact p value returned from scipy.stats.ttest_ind(). The
    value for better should be either "university town" or "non-university town"
    depending on which has a lower mean price ratio (which is equivalent to a
    reduced market loss).'''

```

```

    unitowns = get_list_of_university_towns()
    bottom = get_recession_bottom()
    start = get_recession_start()
    hdata = convert_housing_data_to_quarters()
    bstart = hdata.columns[hdata.columns.get_loc(start) -1]

    hdata['ratio'] = hdata[bottom] - hdata[bstart]
    hdata = hdata[[bottom,bstart,'ratio']]
    hdata = hdata.reset_index()
    unitowns_hdata = pd.merge(hdata,unitowns,how='inner',on=['State','RegionName'])
    unitowns_hdata['uni'] = True
    hdata2 = pd.merge(hdata,unitowns_hdata,how='outer',on=['State','RegionName',bottom,
    hdata2['uni'] = hdata2['uni'].fillna(False)

```

```

    ut = hdata2[hdata2['uni'] == True]
    nut = hdata2[hdata2['uni'] == False]

```

```

    t,p = ttest_ind(ut['ratio'].dropna(),nut['ratio'].dropna())

```

```

    different = True if p < 0.01 else False

```

```

    better = "non-university town" if ut['ratio'].mean() < nut['ratio'].mean() else "un

```

```

    return different, p, better

```

```
In [ ]:
```