
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](https://www.coursera.org/learn/python-data-analysis/resources/0dhYG) (<https://www.coursera.org/learn/python-data-analysis/resources/0dhYG>) course resource.

```
In [80]: import pandas as pd
import numpy as np
from scipy.stats import ttest_ind
```

Assignment 4 - Hypothesis Testing

This assignment requires more individual learning than previous assignments - you are encouraged to check out the [pandas documentation](http://pandas.pydata.org/pandas-docs/stable/) (<http://pandas.pydata.org/pandas-docs/stable/>) to find functions or methods you might not have used yet, or ask questions on [Stack Overflow](http://stackoverflow.com/) (<http://stackoverflow.com/>) 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.
($\text{price_ratio} = \text{quarter_before_recession} / \text{recession_bottom}$)

The following data files are available for this assignment:

- From the [Zillow research data site](http://www.zillow.com/research/data/) (<http://www.zillow.com/research/data/>) there is housing data for the United States. In particular the datafile for [all homes at a city level](http://files.zillowstatic.com/research/public/City/City_Zhvi_AllHomes.csv) (http://files.zillowstatic.com/research/public/City/City_Zhvi_AllHomes.csv), `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 (https://en.wikipedia.org/wiki/List_of_college_towns#College_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 (<http://www.bea.gov/national/index.htm#gdp>) 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 [1]: # Use this dictionary to map state names to two letter acronyms
states = {'OH': 'Ohio', 'KY': 'Kentucky', 'AS': 'American Samoa', 'NV': 'Nevada', 'WY': 'Wyoming', 'NA': 'National', 'AL'
```

```
In [2]: import pandas as pd
import numpy as np
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'. '''

    temp=[]
    with open('university_towns.txt', "r") as file :
        for line in file :
            if (line.strip().endswith('[edit]')):
                state=line.split('[')[0]
                continue
            else:
                region=line.split('(')[0].strip()
                temp.append([state,region])
    df=pd.DataFrame(temp,columns=['State','RegionName'])
    return df

get_list_of_university_towns()
```

Out[2]:

	State	RegionName
0	Alabama	Auburn
1	Alabama	Florence
2	Alabama	Jacksonville
3	Alabama	Livingston
4	Alabama	Montevallo
5	Alabama	Troy
6	Alabama	Tuscaloosa

	State	RegionName
7	Alabama	Tuskegee
8	Alaska	Fairbanks
9	Arizona	Flagstaff
10	Arizona	Tempe
11	Arizona	Tucson
12	Arkansas	Arkadelphia
13	Arkansas	Conway
14	Arkansas	Fayetteville
15	Arkansas	Jonesboro
16	Arkansas	Magnolia
17	Arkansas	Monticello
18	Arkansas	Russellville
19	Arkansas	Searcy
20	California	Angwin
21	California	Arcata
22	California	Berkeley
23	California	Chico
24	California	Claremont
25	California	Cotati
26	California	Davis
27	California	Irvine
28	California	Isla Vista
29	California	University Park, Los Angeles
...
487	Virginia	Wise
488	Virginia	Chesapeake

	State	RegionName
489	Washington	Bellingham
490	Washington	Cheney
491	Washington	Ellensburg
492	Washington	Pullman
493	Washington	University District, Seattle
494	West Virginia	Athens
495	West Virginia	Buckhannon
496	West Virginia	Fairmont
497	West Virginia	Glenville
498	West Virginia	Huntington
499	West Virginia	Montgomery
500	West Virginia	Morgantown
501	West Virginia	Shepherdstown
502	West Virginia	West Liberty
503	Wisconsin	Appleton
504	Wisconsin	Eau Claire
505	Wisconsin	Green Bay
506	Wisconsin	La Crosse
507	Wisconsin	Madison
508	Wisconsin	Menomonie
509	Wisconsin	Milwaukee
510	Wisconsin	Oshkosh
511	Wisconsin	Platteville
512	Wisconsin	River Falls
513	Wisconsin	Stevens Point
514	Wisconsin	Waukesha

	State	RegionName
515	Wisconsin	Whitewater
516	Wyoming	Laramie

517 rows × 2 columns

```
In [18]: def get_recession_start():
'''Returns the year and quarter of the recession start time as a
string value in a format such as 2005q3'''
df = pd.read_excel('gdplev.xls',skiprows=7)
p=df.copy()
df=df[212:]
df=df[['Unnamed: 4','Unnamed: 5']]
df=df.rename(columns={'Unnamed: 4':'Quarter','Unnamed: 5':'GDP'})
df=df.reset_index()
pot=[]
for i in range(0,len(df)-4):
    if (df.loc[i]['GDP']>df.loc[i+1]['GDP']) & (df.loc[i+1]['GDP']>df.loc[i+2]['GDP']):
        pot.append(df.loc[i]['Quarter'])

    return pot[0]
get_recession_start()
```

Out[18]: '2008q3'

```
In [8]: def get_recession_end():  
    '''Returns the year and quarter of the recession start time as a  
    string value in a format such as 2005q3'''  
    df = pd.read_excel('gdplev.xls', skiprows=7)  
    p=df.copy()  
    df=df[212:]  
    df=df[['Unnamed: 4', 'Unnamed: 5']]  
    df=df.rename(columns={'Unnamed: 4': 'Quarter', 'Unnamed: 5': 'GDP'})  
    df=df.reset_index()  
    pot=[]  
    for i in range(0, len(df)-4):  
        if ((df.loc[i]['GDP']>df.loc[i+1]['GDP']) & (df.loc[i+1]['GDP']>df.loc[i+2]['GDP'])) & ((df.loc[i+2]['GDP']<df.lo  
            pot.append(df.loc[i+4]['Quarter'])  
  
    return pot[0]  
get_recession_end()
```

Out[8]: '2009q4'

```
In [5]: def get_recession_bottom():
        '''Returns the year and quarter of the recession bottom time as a
        string value in a format such as 2005q3'''
        df = pd.read_excel('gdplev.xls', skiprows=7)
        p = df.copy()
        df = df[212:]
        df = df[['Unnamed: 4', 'Unnamed: 5']]
        df = df.rename(columns={'Unnamed: 4': 'Quarter', 'Unnamed: 5': 'GDP'})
        df = df.reset_index()
        pot = []
        for i in range(0, len(df)-4):
            if ((df.loc[i]['GDP'] > df.loc[i+1]['GDP']) & (df.loc[i+1]['GDP'] > df.loc[i+2]['GDP'])) & ((df.loc[i+2]['GDP'] < df.loc[i+3]['GDP']) & (df.loc[i+3]['GDP'] > df.loc[i+4]['GDP'])):
                pot.append([df.loc[i]['Quarter'], df.loc[i+1]['Quarter'], df.loc[i+2]['Quarter'], df.loc[i+3]['Quarter'], df.loc[i+4]['Quarter']])
        return pot[0][2]
get_recession_bottom()
```

Out[5]: '2009q2'


```
In [6]: 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.
'''
df=pd.read_csv('City_Zhvi_AllHomes.csv')
x=df['State']
y=df['RegionName']
df2=pd.DataFrame({'State':x,'RegionName':y})
for i in range(2000,2016):
    df2[str(i)+'q1']=(df[[str(i)+'-01',str(i)+'-02',str(i)+'-03']]).mean(axis=1)
    df2[str(i)+'q2']=(df[[str(i)+'-04',str(i)+'-05',str(i)+'-06']]).mean(axis=1)
    df2[str(i)+'q3']=(df[[str(i)+'-07',str(i)+'-08',str(i)+'-09']]).mean(axis=1)
    df2[str(i)+'q4']=(df[[str(i)+'-10',str(i)+'-11',str(i)+'-12']]).mean(axis=1)
df2['2016q1']=(df[['2016-01','2016-02','2016-03']]).mean(axis=1)
df2['2016q2']=(df[['2016-04','2016-05','2016-06']]).mean(axis=1)
df2['2016q3']=(df[['2016-07','2016-08']]).mean(axis=1)

df2['State']=[states[x] for x in df2['State']]
df2=df2.set_index(["State","RegionName"])

return df2
convert_housing_data_to_quarters()
```

```
Out[6]:
```

		2000q1	2000q2	2000q3	2000q4	2001q1	2001q2	2001q3	2001q4
	State RegionName								
	New York New York	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

		2000q1	2000q2	2000q3	2000q4	2001q1	2001q2	2001q3	2001q4
State	RegionName								
California	Los Angeles	2.070667e+05	2.144667e+05	2.209667e+05	2.261667e+05	2.330000e+05	2.391000e+05	2.450667e+05	2.530333e+05
Illinois	Chicago	1.384000e+05	1.436333e+05	1.478667e+05	1.521333e+05	1.569333e+05	1.618000e+05	1.664000e+05	1.704333e+05
Pennsylvania	Philadelphia	5.300000e+04	5.363333e+04	5.413333e+04	5.470000e+04	5.533333e+04	5.553333e+04	5.626667e+04	5.753333e+04
Arizona	Phoenix	1.118333e+05	1.143667e+05	1.160000e+05	1.174000e+05	1.196000e+05	1.215667e+05	1.227000e+05	1.243000e+05
Nevada	Las Vegas	1.326000e+05	1.343667e+05	1.354000e+05	1.370000e+05	1.395333e+05	1.417333e+05	1.433667e+05	1.461333e+05
California	San Diego	2.229000e+05	2.343667e+05	2.454333e+05	2.560333e+05	2.672000e+05	2.762667e+05	2.845000e+05	2.919333e+05

In []:

```
In [81]: 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()
    headq = convert_housing_data_to_quarters()
    tempo=headq.copy()
    arr=[]
    for i in range(2000,2016):
        arr+=[str(i)+'q1',str(i)+'q2',str(i)+'q3',str(i)+'q4']
    arr+=[str(2016)+'q1',str(2016)+'q2',str(2016)+'q3']
    arr.remove(start)
    arr.remove(bottom)

    tempo=tempo.drop(arr,axis=1)
    tempo['Diff']=headq[start]-headq[bottom]
    tempo=tempo.reset_index()
    uni=pd.merge(tempo,unitowns,how='inner',left_on=['State','RegionName'],right_on=['State','RegionName'])
    uni['Status']=True
    comp=pd.merge(tempo,uni,how='outer',left_on=['State','RegionName',start,bottom,'Diff'],right_on=['State','RegionName'])
    comp['Status']=comp['Status'].fillna(False)
    uni=comp[comp['Status']==True]
    t1=uni['Diff'].dropna()
    non_uni=comp[comp['Status']==False]
    t2=non_uni['Diff'].dropna()

    t,p = ttest_ind(t1, t2)
```

```
    if p<0.01:
        different=True
    else:
        different=False

    if t1.mean()<t2.mean():
        better='university town'
    else:
        better='non-university town'

    return (different,p,better)
run_ttest()
```

Out[81]: (True, 0.0043252148535112009, 'university town')

In []: