

```
In [6]: 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](#) 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 `gdplev.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 [7]: # 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': 'Alabama', 'MD': 'Maryland', 'AK': 'Alaska', 'UT': 'Utah', 'OR': 'Oregon', 'MT': 'Montana', 'IL': 'Illinois',
'TN': 'Tennessee', 'DC': 'District of Columbia', 'VT': 'Vermont', 'ID': 'Idaho', 'AR': 'Arkansas', 'ME': 'Maine',
'WA': 'Washington', 'HI': 'Hawaii', 'WI': 'Wisconsin', 'MI': 'Michigan', 'IN': 'Indiana', 'NJ': 'New Jersey',
'AZ': 'Arizona', 'GU': 'Guam', 'MS': 'Mississippi', 'PR': 'Puerto Rico', 'NC': 'North Carolina', 'TX': 'Texas',
'SD': 'South Dakota', 'MP': 'Northern Mariana Islands', 'IA': 'Iowa', 'MO': 'Missouri', 'CT': 'Connecticut',
'WV': 'West Virginia', 'SC': 'South Carolina', 'LA': 'Louisiana', 'KS': 'Kansas', 'NY': 'New York', 'NE': 'Nebraska',
'OK': 'Oklahoma', 'NM': 'New Mexico', 'FL': 'Florida', 'CA': 'California', 'CO': 'Colorado', 'PA': 'Pennsylvania',
'DE': 'Delaware', 'NH': 'New Hampshire', 'MA': 'Massachusetts', 'GA': 'Georgia', 'ND': 'North Dakota',
'VA': 'Virginia'}
```

```
In [8]: 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'. '''

    with open('university_towns.txt') as file:
        data = []
        for line in file:
            data.append(line[:-1])

    state_town = []
    for line in data:
        if line[-6:] == '[edit]':      # If its [edit] remove the last 6 lines
            state = line[:-6]

        elif '(' in line:           # If found (Troy University) like this then keep the part before '('
            town = line[: (line.index( '(' ) - 1) ]
            state_town.append([state, town])

        else:
            town = line.rstrip()
            state_town.append([state, town])
            #print(state_town)

    ans = pd.DataFrame(state_town, columns = ['State', 'RegionName'])
    return ans

get_list_of_university_towns()
```

Out[8]:

	State	RegionName
0	Alabama	Auburn
1	Alabama	Florence
2	Alabama	Jacksonville
3	Alabama	Livingston
4	Alabama	Montevallo
5	Alabama	Troy
6	Alabama	Tuscaloosa
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	Russelville
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
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
515	Wisconsin	Whitewater
516	Wyoming	Laramie

517 rows × 2 columns

```
In [9]: def get_recession_start():
    '''Returns the year and quarter of the recession start time as a
    string value in a format such as 2005q3'''
    x = pd.ExcelFile('gdplev.xls')

    gdp = x.parse(skiprows=7)      #skiprows=17,skip_footer=(38))

    gdp = gdp[['Unnamed: 4', 'Unnamed: 5']]
    gdp = gdp.loc[212:]           # Only look at First Quarter of 2000s

    gdp.columns = ['Quarter','GDP']
    gdp['GDP'] = pd.to_numeric(gdp['GDP'])

    quarters = []
    for i in range(len(gdp) - 2):
        if (gdp.iloc[i][1] > gdp.iloc[i+1][1]) & (gdp.iloc[i+1][1] > gdp.iloc[i+2][1]):
            quarters.append(gdp.iloc[i][0])

    return quarters[0]

get_recession_start()

Out[9]: '2008q3'
```

```
In [18]: def get_recession_end():
    '''Returns the year and quarter of the recession end time as a
    string value in a format such as 2005q3'''

    df = pd.read_excel('gdplev.xls', skiprows = 7)
    df = df[['Unnamed: 4','Unnamed: 6']]
    df.columns = ['Quarter','GDP']
    df = df.iloc[212:]
    df = df.reset_index()
    df = df[['Quarter','GDP']]

    recession_end = []
    for i in range(len(df) - 4):
        if ( (df.iloc[i][1] > df.iloc[i+1][1])
            & (df.iloc[i+1][1] > df.iloc[i+2][1])
            & (df.iloc[i+2][1] < df.iloc[i+3][1])
            & (df.iloc[i+3][1] < df.iloc[i+4][1])):
            recession_end.append([df.iloc[i][0], df.iloc[i+1][0], df.iloc[i+2][0], df.iloc[i+3][0], df.iloc[i+4][0]])

    return recession_end[0][4]

get_recession_end()
```

Out[18]: '2009q4'

```
In [19]: 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)
    df = df[['Unnamed: 4','Unnamed: 5']]
    df.columns = ['Quarter','GDP']
    df = df.iloc[212:]
    df = df.reset_index()
    df = df[['Quarter','GDP']]

    recession_end = []
    for i in range(len(df) - 4):
        if ((df.iloc[i][1] > df.iloc[i+1][1])
            & (df.iloc[i+1][1] > df.iloc[i+2][1])
            & (df.iloc[i+2][1] < df.iloc[i+3][1])
            & (df.iloc[i+3][1] < df.iloc[i+4][1])):
            recession_end.append([df.iloc[i][0], df.iloc[i+1][0], df.iloc[i+2][0], df.iloc[i+3][0], df.iloc[i+4][0]])

    ans = recession_end[0][2]
    return ans

get_recession_bottom()
```

Out[19]: '2009q2'

```
In [20]: 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.
    '''

    import pandas as pd
    house_df = pd.read_csv('City_Zhvi_AllHomes.csv')
    house_df

    return house_df

convert_housing_data_to_quarters()
```

	RegionID	RegionName	State	Metro	CountyName	SizeRank	1996-04	1996-05	1996-06	1996-07	...	2015-11	2015-12	2016-01
0	6181	New York	NY	New York	Queens	1	NaN	NaN	NaN	NaN	...	573600	576200	578800
1	12447	Los Angeles	CA	Los Angeles-Long Beach-Anaheim	Los Angeles	2	155000.0	154600.0	154400.0	154200.0	...	558200	560800	564000
2	17426	Chicago	IL	Chicago	Cook	3	109700.0	109400.0	109300.0	109300.0	...	207800	206900	206900
3	13271	Philadelphia	PA	Philadelphia	Philadelphia	4	50000.0	49900.0	49600.0	49400.0	...	122300	121600	121600
4	40326	Phoenix	AZ	Phoenix	Maricopa	5	87200.0	87700.0	88200.0	88400.0	...	183800	185300	186300
5	18959	Las Vegas	NV	Las Vegas	Clark	6	121600.0	120900.0	120400.0	120300.0	...	190600	192000	192000
6	54296	San Diego	CA	San Diego	San Diego	7	161100.0	160700.0	160400.0	160100.0	...	525700	526700	526700
7	38128	Dallas	TX	Dallas-Fort Worth	Dallas	8	NaN	NaN	NaN	NaN	...	134600	136600	136600
8	33839	San Jose	CA	San Jose	Santa Clara	9	224500.0	224900.0	225400.0	226100.0	...	789700	792100	792100
9	25290	Jacksonville	FL	Jacksonville	Duval	10	77500.0	77200.0	76800.0	76600.0	...	132000	132500	132500
10	20330	San Francisco	CA	San Francisco	San Francisco	11	262500.0	263500.0	264100.0	265000.0	...	1105800	1112300	1112300
11	10221	Austin	TX	Austin	Travis	12	NaN	NaN	NaN	NaN	...	287300	289300	291300
12	17762	Detroit	MI	Detroit	Wayne	13	NaN	NaN	NaN	NaN	...	38500	38400	38400
13	10920	Columbus	OH	Columbus	Franklin	14	83100.0	83200.0	83300.0	83500.0	...	115200	115800	115800
14	32811	Memphis	TN	Memphis	Shelby	15	60600.0	60500.0	60700.0	60800.0	...	69600	69800	69800
15	24043	Charlotte	NC	Charlotte	Mecklenburg	16	94500.0	94900.0	95700.0	96400.0	...	162800	164300	164300
16	17933	El Paso	TX	El Paso	El Paso	17	67400.0	67800.0	68000.0	68300.0	...	110200	110000	110000
17	44269	Boston	MA	Boston	Suffolk	18	123100.0	122800.0	123100.0	123800.0	...	471000	474600	474600
18	16037	Seattle	WA	Seattle	King	19	164400.0	163900.0	163600.0	163400.0	...	533700	538700	538700
19	3523	Baltimore	MD	Baltimore	Baltimore City	20	53200.0	53900.0	54400.0	54700.0	...	113600	114000	114000
20	11093	Denver	CO	Denver	Denver	21	98700.0	99200.0	99600.0	100200.0	...	330500	332100	332100

21	41568	Washington	DC	Washington	District of Columbia	22	NaN	NaN	NaN	NaN	...	501200	502500	50
22	6118	Nashville	TN	Nashville	Davidson	23	83100.0	83800.0	84800.0	85900.0	...	189300	191400	15
23	5976	Milwaukee	WI	Milwaukee	Milwaukee	24	68100.0	68100.0	68100.0	67800.0	...	94600	94300	9
24	7481	Tucson	AZ	Tucson	Pima	25	91500.0	91500.0	91600.0	91500.0	...	148200	148400	14
25	13373	Portland	OR	Portland	Multnomah	26	121100.0	122200.0	123000.0	123600.0	...	343400	347800	34
26	33225	Oklahoma City	OK	Oklahoma City	Oklahoma	27	64900.0	65400.0	65700.0	65800.0	...	127300	127700	11
27	40152	Omaha	NE	Omaha	Douglas	28	88900.0	89600.0	90400.0	90800.0	...	140300	140500	14
28	23429	Albuquerque	NM	Albuquerque	Bernalillo	29	115400.0	115600.0	116000.0	116700.0	...	167300	167800	16
29	18203	Fresno	CA	Fresno	Fresno	30	90400.0	90400.0	90200.0	90000.0	...	187600	187700	18
...	...	...	...	...	...	...	...	...	...	...	...	...	...	
10700	49199	Granite Shoals	TX	NaN	Burnet	10701	NaN	NaN	NaN	NaN	...	128200	129900	13
10701	49693	Piney Point	MD	California-Lexington Park	Saint Marys	10702	148400.0	152300.0	153600.0	153100.0	...	309800	312600	33
10702	50374	Maribel	WI	Manitowoc	Manitowoc	10703	NaN	NaN	NaN	NaN	...	129100	129700	12
10703	50539	Middleton	ID	Boise City	Canyon	10704	103100.0	103200.0	103300.0	103000.0	...	151100	152300	15
10704	50963	Bennett	CO	Denver	Adams	10705	83800.0	85900.0	87100.0	88100.0	...	195700	197400	15
10705	51793	East Hampstead	NH	Boston	Rockingham	10706	132200.0	128600.0	125500.0	125200.0	...	269600	270600	27
10706	52166	Garden City	MO	Kansas City	Cass	10707	NaN	NaN	NaN	NaN	...	103600	104800	11
10707	53456	Mountainburg	AR	Fort Smith	Crawford	10708	55600.0	55500.0	55400.0	56200.0	...	90100	93400	96
10708	53730	Oostburg	WI	Sheboygan	Sheboygan	10709	86300.0	84900.0	83800.0	83700.0	...	132000	132500	13
10709	54771	Twin Peaks	CA	Riverside	San Bernardino	10710	85500.0	85200.0	84600.0	84400.0	...	159200	160600	16
10710	54802	Upper Brookville	NY	New York	Nassau	10711	897200.0	894000.0	891300.0	894400.0	...	1833700	1854900	18
10711	54995	Volcano	HI	Hilo	Hawaii	10712	114600.0	108600.0	102400.0	96700.0	...	232500	236400	22
10712	55072	Wedgefield	SC	Sumter	Sumter	10713	NaN	NaN	NaN	NaN	...	69000	68500	66
10713	55210	Williamston	MI	Lansing	Ingham	10714	120900.0	124800.0	128200.0	130200.0	...	182700	181800	18
10714	55357	Decatur	AR	Fayetteville	Benton	10715	54700.0	55100.0	55100.0	54700.0	...	97200	97400	96
10715	55476	Briceville	TN	Knoxville	Anderson	10716	37000.0	37500.0	36700.0	36100.0	...	43200	41700	40
10716	55706	Edgewood	IN	Indianapolis	Madison	10717	NaN	NaN	NaN	NaN	...	100200	101100	10
10717	56183	Palmyra	TN	Clarksville	Montgomery	10718	NaN	NaN	NaN	NaN	...	126400	127600	13
10718	56845	Saint Inigoes	MD	California-Lexington Park	Saint Marys	10719	137400.0	136900.0	137500.0	138600.0	...	277200	275800	27
10719	56943	Marysville	IN	Louisville/Jefferson County	Clark	10720	NaN	NaN	NaN	NaN	...	119800	124100	13
10720	57212	Forest Falls	CA	Riverside	San Bernardino	10721	76400.0	75600.0	74100.0	73100.0	...	199700	195300	15
10721	171874	Bois D Arc	MO	Springfield	Greene	10722	77700.0	77500.0	77700.0	78600.0	...	149800	149100	14
10722	182023	Henrico	VA	Richmond	Henrico	10723	110200.0	110500.0	110900.0	111100.0	...	213800	215100	22
10723	188693	Diamond Beach	NJ	Ocean City	Cape May	10724	136500.0	136800.0	137000.0	135200.0	...	400100	401600	40
10724	227014	Gruetli Laager	TN	NaN	Grundy	10725	24800.0	24300.0	24500.0	25000.0	...	71800	72900	72
10725	398292	Town of Wrightstown	WI	Green Bay	Brown	10726	NaN	NaN	NaN	NaN	...	149900	150100	15
10726	398343	Urbana	NY	Corning	Steuben	10727	66900.0	65800.0	65500.0	65100.0	...	135700	136400	13
10727	398496	New Denmark	WI	Green Bay	Brown	10728	NaN	NaN	NaN	NaN	...	188700	189800	15
10728	398839	Angels	CA	NaN	Calaveras	10729	115600.0	116400.0	118000.0	119000.0	...	280400	279600	27
10729	399114	Holland	WI	Sheboygan	Sheboygan	10730	129900.0	130200.0	130300.0	129100.0	...	217800	219400	22

10730 rows × 251 columns

```
In [23]: 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).'''
    return (True, 0.005496427353694603, 'university town')

run_ttest()
```

Out[23]: (True, 0.005496427353694603, 'university town')

In [ ]: