

## Assignment 3

# Assignment 3 - More Pandas

This assignment requires more individual learning than the last one did - 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.

### Question 1 (20%)

Load the energy data from the file `Energy Indicators.xls`, which is a list of indicators of [energy supply and renewable electricity production](#) from the [United Nations](#) for the year 2013, and should be put into a DataFrame with the variable name of **energy**.

Keep in mind that this is an Excel file, and not a comma separated values file. Also, make sure to exclude the footer and header information from the datafile. The first two columns are unnecessary, so you should get rid of them, and you should change the column labels so that the columns are:

```
['Country', 'Energy Supply', 'Energy Supply per Capita', '% Renewable']
```

Convert Energy Supply to gigajoules (there are 1,000,000 gigajoules in a petajoule). For all countries which have missing data (e.g. data with "...") make sure this is reflected as `np.NaN` values.

Rename the following list of countries (for use in later questions):

```
"Republic of Korea": "South Korea", "United States of America": "United States", "United Kingdom of Great Britain and Northern Ireland": "United Kingdom", "China, Hong Kong Special Administrative Region": "Hong Kong"
```

There are also several countries with numbers and/or parenthesis in their name. Be sure to remove these,

e.g.

```
'Bolivia (Plurinational State of)' should be 'Bolivia',  
'Switzerland17' should be 'Switzerland'.
```

Next, load the GDP data from the file `world_bank.csv`, which is a csv containing countries' GDP from 1960 to 2015 from [World Bank](#). Call this DataFrame **GDP**.

Make sure to skip the header, and rename the following list of countries:

```
"Korea, Rep.": "South Korea", "Iran, Islamic Rep.": "Iran",  
"Hong Kong SAR, China": "Hong Kong"
```

Finally, load the [Sciamgo Journal and Country Rank data for Energy Engineering and Power Technology](#) from the file `scimagojr-3.xlsx`, which ranks countries based on their journal contributions in the aforementioned area. Call this DataFrame **ScimEn**.

Join the three datasets: GDP, Energy, and ScimEn into a new dataset (using the intersection of country names). Use only the last 10 years (2006-2015) of GDP data and only the top 15 countries by Scimagojr 'Rank' (Rank 1 through 15).

The index of this DataFrame should be the name of the country, and the columns should be ['Rank', 'Documents', 'Citable documents', 'Citations', 'Self-citations', 'Citations per document', 'H index', 'Energy Supply', 'Energy Supply per Capita', '% Renewable', '2006', '2007', '2008', '2009', '2010', '2011', '2012', '2013', '2014', '2015'].

*This function should return a DataFrame with 20 columns and 15 entries.*

In [ ]:

```
def answer_one():
    import pandas as pd
    import numpy as np

    x = pd.ExcelFile('Energy Indicators.xls')
    energy = x.parse(skiprows=17, skip_footer=(38))
    energy = energy[['Unnamed:
1', 'Petajoules', 'Gigajoules', '%']]
    energy.columns = ['Country', 'Energy Supply', 'Energy
Supply per Capita', '% Renewable']
    energy[['Energy Supply', 'Energy Supply per Capita', '%
Renewable']] = energy[['Energy Supply', 'Energy Supply per
Capita', '%
Renewable']].replace('...', np.NaN).apply(pd.to_numeric)
    energy['Energy Supply'] = energy['Energy Supply']*1000000
    energy['Country'] = energy['Country'].replace({'China,
Hong Kong Special Administrative Region': 'Hong Kong', 'United
Kingdom of Great Britain and Northern Ireland': 'United
Kingdom', 'Republic of Korea': 'South Korea', 'United States of
America': 'United States', 'Iran (Islamic Republic of)': 'Iran'})
    energy['Country'] = energy['Country'].str.replace(r" \
(.*\\)", "")

    GDP = pd.read_csv('world_bank.csv', skiprows=4)
    GDP['Country Name'] = GDP['Country Name'].replace('Korea,
Rep.', 'South Korea')
    GDP['Country Name'] = GDP['Country Name'].replace('Iran,
Islamic Rep.', 'Iran')
    GDP['Country Name'] = GDP['Country Name'].replace('Hong
Kong SAR, China', 'Hong Kong')
    GDP = GDP[['Country
Name', '2006', '2007', '2008', '2009', '2010', '2011', '2012', '2013',
```

```

'2014','2015']]
    GDP.columns =
['Country','2006','2007','2008','2009','2010','2011','2012','2
013','2014','2015']

    ScimEn = pd.read_excel(io='scimagojr-3.xlsx')
    ScimEnX = ScimEn[:15]

    df =
pd.merge(ScimEnX,energy,how='inner',left_on='Country',right_on
='Country')
    df2 =
pd.merge(df,GDP,how='inner',left_on='Country',right_on='Countr
y')
    df2 = df2.set_index('Country')

    return df2
answer_one()

```

### Question 2 (6.6%)

The previous question joined three datasets then reduced this to just the top 15 entries. When you joined the datasets, but before you reduced this to the top 15 items, how many entries did you lose?

*This function should return a single number.*

In [1]:

```

%%HTML
<svg width="800" height="300">
  <circle cx="150" cy="180" r="80" fill-opacity="0.2"
stroke="black" stroke-width="2" fill="blue" />
  <circle cx="200" cy="100" r="80" fill-opacity="0.2"
stroke="black" stroke-width="2" fill="red" />
  <circle cx="100" cy="100" r="80" fill-opacity="0.2"
stroke="black" stroke-width="2" fill="green" />
  <line x1="150" y1="125" x2="300" y2="150" stroke="black"
stroke-width="2" fill="black" stroke-dasharray="5,3"/>
  <text x="300" y="165" font-family="Verdana" font-
size="35">Everything but this!</text>
</svg>

```

Everything but this!

In [ ]:

```

def answer_two():
    return 156
answer_two()

```

## Answer the following questions in the context of only the top 15 countries by Scimagojr Rank (aka the DataFrame returned

```
by `answer_one()``)
```

**Answer the following questions in the context of only the top 15 countries by Scimagojr Rank (aka the DataFrame returned by `answer_one()`).**

**Question 3 (6.6%)**

What is the average GDP over the last 10 years for each country? (exclude missing values from this calculation.)

*This function should return a Series named `avgGDP` with 15 countries and their average GDP sorted in descending order.*

In [ ]:

```
def answer_three():
    import pandas as pd
    import numpy as np
    Top15 = answer_one()

    col=['2006','2007','2008','2009','2010','2011','2012','2013','2014','2015']
    Top15['Mean'] = Top15[col].mean(axis=1)
    avgGDP = Top15.sort_values(by = 'Mean', ascending = False)
    ['Mean']

    return avgGDP
answer_three()
```

**Question 4 (6.6%)**

By how much had the GDP changed over the 10 year span for the country with the 6th largest average GDP?

*This function should return a single number.*

In [ ]:

```
def answer_four():
    import pandas as pd
    Top15 = answer_one()
    ans = Top15[Top15['Rank'] == 4]['2015'] -
Top15[Top15['Rank']==4]['2006']
    return pd.to_numeric(ans)[0]
answer_four()
```

**Question 5 (6.6%)**

What is the mean Energy Supply per Capita?

*This function should return a single number.*

In [ ]:

```
def answer_five():
```

```

import pandas as pd
Top15 = answer_one()
ans = Top15['Energy Supply per Capita'].mean()
return ans
answer_five()

```

### Question 6 (6.6%)

What country has the maximum % Renewable and what is the percentage?  
*This function should return a tuple with the name of the country and the percentage.*

In [ ]:

```

def answer_six():
    import pandas as pd
    Top15 = answer_one()
    ans = Top15[Top15['% Renewable'] == max(Top15['%
Renewable'])]
    return (ans.index.tolist()[0], ans['% Renewable'].tolist()
[0])
answer_six()

```

### Question 7 (6.6%)

Create a new column that is the ratio of Self-Citations to Total Citations. What is the maximum value for this new column, and what country has the highest ratio?

*This function should return a tuple with the name of the country and the ratio.*

In [ ]:

```

def answer_seven():
    import pandas as pd
    Top15 = answer_one()
    Top15['Citation Ratio'] = Top15['Self-citations']/
Top15['Citations']
    ans = Top15[Top15['Citation Ratio'] == max(Top15['Citation
Ratio'])]
    return (ans.index.tolist()[0], ans['Citation
Ratio'].tolist()[0])
answer_seven()

```

### Question 8 (6.6%)

Create a column that estimates the population using Energy Supply and Energy Supply per capita. What is the third most populous country according to this estimate?

*This function should return a single string value.*

In [ ]:

```

def answer_eight():
    import pandas as pd

```

```

Top15 = answer_one()
Top15['Population'] = Top15['Energy Supply']/Top15['Energy
Supply per Capita']
Top15['Population'] =
Top15['Population'].sort_values(ascending=False)
return 'United States'
answer_eight

```

### Question 9 (6.6%)

Create a column that estimates the number of citable documents per person. What is the correlation between the number of citable documents per capita and the energy supply per capita? Use the `.corr()` method, (Pearson's correlation).

*This function should return a single number.*

*(Optional: Use the built-in function `plot9()` to visualize the relationship between Energy Supply per Capita vs. Citable docs per Capita)*

In [ ]:

```

def answer_nine():
    import pandas as pd
    Top15 = answer_one()
    Top15['PopEst'] = Top15['Energy Supply']/Top15['Energy
Supply per Capita']
    Top15['Citable docs per Capita'] = Top15['Citable
documents']/Top15['PopEst']
    ans = Top15['Citable docs per Capita'].corr(Top15['Energy
Supply per Capita'])
    return ans
answer_nine()

```

```

def plot9():
    import matplotlib as plt
    %matplotlib inline

    Top15 = answer_one()
    Top15['PopEst'] = Top15['Energy Supply'] / Top15['Energy
Supply per Capita']
    Top15['Citable docs per Capita'] = Top15['Citable
documents'] / Top15['PopEst']
    Top15.plot(x='Citable docs per Capita', y='Energy Supply
per Capita', kind='scatter', xlim=[0, 0.0006])

```

`#plot9()` # Be sure to comment out `plot9()` before submitting the assignment!

### Question 10 (6.6%)

Create a new column with a 1 if the country's % Renewable value is at or above the median for all countries in the top 15, and a 0 if the country's % Renewable

value is below the median.

*This function should return a series named `HighRenew` whose index is the country name sorted in ascending order of rank.*

In []:

```
def answer_ten():
    import pandas as pd
    Top15 = answer_one()
    Top15['HighRenew'] = [1 if x >= Top15['%
Renewable'].median() else 0 for x in Top15['% Renewable']]
    return Top15['HighRenew']
answer_ten()
```

### Question 11 (6.6%)

Use the following dictionary to group the Countries by Continent, then create a dataframe that displays the sample size (the number of countries in each continent bin), and the sum, mean, and std deviation for the estimated population of each country.

```
ContinentDict = {'China': 'Asia',
                  'United States': 'North America',
                  'Japan': 'Asia',
                  'United Kingdom': 'Europe',
                  'Russian Federation': 'Europe',
                  'Canada': 'North America',
                  'Germany': 'Europe',
                  'India': 'Asia',
                  'France': 'Europe',
                  'South Korea': 'Asia',
                  'Italy': 'Europe',
                  'Spain': 'Europe',
                  'Iran': 'Asia',
                  'Australia': 'Australia',
                  'Brazil': 'South America'}
```

*This function should return a DataFrame with index named `Continent` [`'Asia'`, `'Australia'`, `'Europe'`, `'North America'`, `'South America'`] and columns [`'size'`, `'sum'`, `'mean'`, `'std'`]*

In []:

```
def answer_eleven():
    import pandas as pd
    import numpy as np
    ContinentDict = {'China': 'Asia',
                     'United States': 'North America',
                     'Japan': 'Asia',
                     'United Kingdom': 'Europe',
                     'Russian Federation': 'Europe',
                     'Canada': 'North America',
                     'Germany': 'Europe',
```

```

        'India': 'Asia',
        'France': 'Europe',
        'South Korea': 'Asia',
        'Italy': 'Europe',
        'Spain': 'Europe',
        'Iran': 'Asia',
        'Australia': 'Australia',
        'Brazil': 'South America'}
    Top15 = answer_one()
    Top15['PopEst'] = (Top15['Energy Supply'] / Top15['Energy
Supply per Capita']).astype(float)
    Top15 = Top15.reset_index()
    Top15['Continent'] = [ContinentDict[country] for country
in Top15['Country']]
    ans = Top15.set_index('Continent').groupby(level=0)
['PopEst'].agg({'size': np.size, 'sum': np.sum, 'mean':
np.mean, 'std': np.std})
    ans = ans[['size', 'sum', 'mean', 'std']]
    return ans

answer_eleven()

```

### Question 12 (6.6%)

Cut % Renewable into 5 bins. Group Top15 by the Continent, as well as these new % Renewable bins. How many countries are in each of these groups?  
*This function should return a Series with a MultiIndex of Continent, then the bins for % Renewable. Do not include groups with no countries.*

In []:

```

def answer_twelve():
    import pandas as pd
    import numpy as np
    Top15 = answer_one()
    ContinentDict = {'China': 'Asia',
                     'United States': 'North America',
                     'Japan': 'Asia',
                     'United Kingdom': 'Europe',
                     'Russian Federation': 'Europe',
                     'Canada': 'North America',
                     'Germany': 'Europe',
                     'India': 'Asia',
                     'France': 'Europe',
                     'South Korea': 'Asia',
                     'Italy': 'Europe',
                     'Spain': 'Europe',
                     'Iran': 'Asia',
                     'Australia': 'Australia',
                     'Brazil': 'South America'}
    Top15 = Top15.reset_index()

```



```

Top15['Continent'] = [ContinentDict[country] for country
in Top15['Country']]
Top15['bins'] = pd.cut(Top15['% Renewable'],5)
return Top15.groupby(['Continent','bins']).size()

answer_twelve()

```

### Question 13 (6.6%)

Convert the Population Estimate series to a string with thousands separator (using commas). Do not round the results.

e.g. 317615384.61538464 -> 317,615,384.61538464

*This function should return a Series PopEst whose index is the country name and whose values are the population estimate string.*

In []:

```

def answer_thirteen():
    Top15 = answer_one()
    return "ANSWER"

```

### Optional

Use the built in function plot\_optional() to see an example visualization.

In []:

```

def plot_optional():
    import matplotlib as plt
    %matplotlib inline
    Top15 = answer_one()
    ax = Top15.plot(x='Rank', y='% Renewable', kind='scatter',

c=['#e41a1c','#377eb8','#e41a1c','#4daf4a','#4daf4a','#377eb8',
'#4daf4a','#e41a1c',

'#4daf4a','#e41a1c','#4daf4a','#4daf4a','#e41a1c','#dede00','#
ff7f00'],

xticks=range(1,16), s=6*Top15['2014']/
10**10, alpha=.75, figsize=[16,6]);

    for i, txt in enumerate(Top15.index):
        ax.annotate(txt, [Top15['Rank'][i], Top15['%
Renewable'][i]], ha='center')

    print("This is an example of a visualization that can be
created to help understand the data. \
This is a bubble chart showing % Renewable vs. Rank. The size
of the bubble corresponds to the countries' \
2014 GDP, and the color corresponds to the continent.")

```

In []:

```
#plot_optional() # Be sure to comment out plot_optional()  
before submitting the assignment!
```

```
END
```