

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 [Scimago 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 [1]: 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)) # Skip the rows & footer
    energy = energy[['Unnamed: 1','Petajoules','Gigajoules','%']]
    # Set the column names
    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)

    # For converting Energy Supply to gigajoules
    energy['Energy Supply'] = energy['Energy Supply'] * 1000000

    # Rename the following list of countries
    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','2013','2014','2015']

    ScimEn = pd.read_excel(io='scimagojr-3.xlsx')
    ScimEn = ScimEn[:15] # For 15 entries
```

```

# Merge sci & energy
df = pd.merge(SciEn_m, energy, how='inner', left_on='Country', right_on='Country')

# Merge sci energy & GDP
final_df = pd.merge(df, GDP, how='inner', left_on='Country', right_on='Country')
final_df = final_df.set_index('Country')

#print(len(final_df))
return final_df

```

answer_one()

Out[1]:

	Rank	Documents	Citable documents	Citations	Self-citations	Citations per document	H index	Energy Supply	Energy Supply per Capita	% Renewable	2006	2007
Country												
China	1	127050	126767	597237	411683	4.70	138	1.271910e+11	93.0	19.754910	3.992331e+12	4.559041e+12
United States	2	966661	94747	792274	265436	8.20	230	9.083800e+10	286.0	11.570980	1.479230e+13	1.505540e+13
Japan	3	30504	30287	223024	61554	7.31	134	1.898400e+10	149.0	10.232820	5.496542e+12	5.617036e+12
United Kingdom	4	20944	20357	206091	37874	9.84	139	7.920000e+09	124.0	10.600470	2.419631e+12	2.482203e+12
Russian Federation	5	18534	18301	34266	12422	1.85	57	3.070900e+10	214.0	17.288680	1.385793e+12	1.504071e+12
Canada	6	17899	17620	215003	40930	12.01	149	1.043100e+10	296.0	61.945430	1.564469e+12	1.596740e+12
Germany	7	17027	16831	140566	27426	8.26	126	1.326100e+10	165.0	17.901530	3.332891e+12	3.441561e+12
India	8	15005	14841	128763	37209	8.58	115	3.319500e+10	26.0	14.969080	1.265894e+12	1.374865e+12
France	9	13153	12973	130632	28601	9.93	114	1.059700e+10	166.0	17.020280	2.607840e+12	2.669424e+12
South Korea	10	11983	11923	114675	22595	9.57	104	1.100700e+10	221.0	2.279353	9.410199e+11	9.924316e+11
Italy	11	10964	10794	111850	26661	10.20	106	6.530000e+09	109.0	33.667230	2.202170e+12	2.234627e+12
Spain	12	9428	9330	123336	23964	13.08	115	4.923000e+09	106.0	37.968590	1.414823e+12	1.468146e+12
Iran	13	8896	8819	57470	19125	6.46	72	9.172000e+09	119.0	5.707721	3.895523e+11	4.250646e+11
Australia	14	8831	8725	90765	15606	10.28	107	5.386000e+09	231.0	11.810810	1.021939e+12	1.060340e+12
Brazil	15	8668	8596	60702	14396	7.00	86	1.214900e+10	59.0	69.648030	1.845080e+12	1.957118e+12

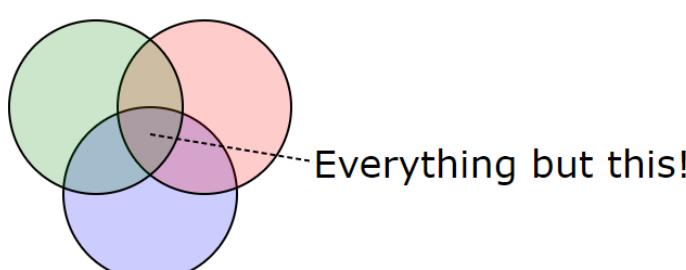
Question 2 (6.6%)

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 [2]:

```
%%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" stroke-dasharray="5,3"/>
  <text x="300" y="165" font-family="Verdana" font-size="35">Everything but this!</text>
</svg>
```



In [27]:

```

def answer_two():
    import pandas as pd
    import numpy as np

    x = pd.ExcelFile('Energy Indicators.xls')
    energy = x.parse(skiprows=17, skip_footer=(38)) # Skip the rows & footer

    energy = energy[['Unnamed: 1','Petajoules','Gigajoules','%']]

    # Set the column names
    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)

    # For converting Energy Supply to gigajoules
    energy['Energy Supply'] = energy['Energy Supply'] * 1000000

    # Rename the following list of countries
    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', '2013', '2014', '2015']

ScimEn = pd.read_excel(io='scimagojr-3.xlsx')
ScimEn_m = ScimEn[:15] # For 15 entries

# Merge sci & energy
df = pd.merge(ScimEn_m, energy, how='inner', left_on='Country', right_on='Country')

# Merge sci energy & GDP
final_df = pd.merge(df, GDP, how='inner', left_on='Country', right_on='Country')
final_df = final_df.set_index('Country')

# Merge sci & energy
df2 = pd.merge(ScimEn_m, energy, how='outer', left_on='Country', right_on='Country')

# Merge sci energy & GDP
final_df2 = pd.merge(df2, GDP, how='outer', left_on='Country', right_on='Country')
final_df2 = final_df2.set_index('Country')

print(len(final_df))
print(len(final_df2))

return 156

```

answer_two()

```

15
264

```

Out[27]: 156

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 [4]: def answer_three():
    Top15 = answer_one()
    avgGDP = Top15[['2006','2007','2008','2009','2010','2011','2012','2013','2014','2015']].mean(axis=1).rename('avgGDP').sort_values(ascending=False)
    return avgGDP

answer_three()

```

Out[4]: Country

Country	avgGDP
United States	1.536434e+13
China	6.348609e+12
Japan	5.542208e+12
Germany	3.493025e+12
France	2.681725e+12
United Kingdom	2.487907e+12
Brazil	2.189794e+12
Italy	2.120175e+12
India	1.769297e+12
Canada	1.660647e+12
Russian Federation	1.565459e+12
Spain	1.418078e+12
Australia	1.164043e+12
South Korea	1.106715e+12
Iran	4.441558e+11

Name: avgGDP, dtype: float64

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 [5]: 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()

```

Out[5]: 246702696075.3999

Question 5 (6.6%)

What is the mean Energy Supply per Capita?

This function should return a single number.

```

In [6]: def answer_five():
    Top15 = answer_one()
    ans = Top15['Energy Supply per Capita'].mean()
    return ans

answer_five()

```

Out[6]: 157.5999999999999

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

Out[7]: ('Brazil', 69.64803)
```

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 [8]: def answer_seven():
    Top15 = answer_one()
    # Created col of citation ratio
    Top15['Citation Ratio'] = Top15['Self-citations'] / Top15['Citations']

    # Same as the above query
    ans = Top15[Top15['Citation Ratio'] == max(Top15['Citation Ratio'])]
    return (ans.index.tolist()[0], ans['Citation Ratio'].tolist()[0])

answer_seven()

Out[8]: ('China', 0.6893126179389422)
```

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 [9]: def answer_eight():
    Top15 = answer_one()
    Top15['Population'] = Top15['Energy Supply'] / Top15['Energy Supply per Capita']

    Top15['Population'] = Top15['Population'].sort_values(ascending=False)
    #print(Top15['Population'])

    return Top15.sort_values(by = 'Population', ascending = False).iloc[2].name

answer_eight()

Out[9]: 'United States'
```

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 [10]: def answer_nine():
    Top15 = answer_one()
    Top15['PopEst'] = Top15['Energy Supply'] / Top15['Energy Supply per Capita']

    Top15['Citable docs per Capita'] = Top15['Citable documents'] / Top15['PopEst']

    return Top15['Citable docs per Capita'].corr(Top15['Energy Supply per Capita'])

answer_nine()

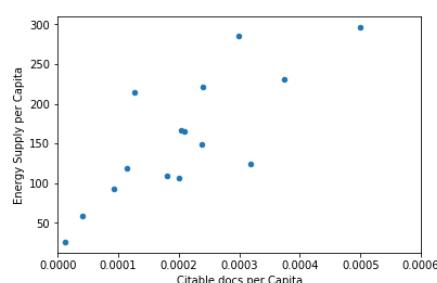
Out[10]: 0.7940010435442945
```

```
In [11]: 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])
```

```
In [12]: #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 [13]: def answer_ten():
    Top15 = answer_one()
    Top15['HighRenew'] = [1 if x >= Top15['% Renewable'].median() else 0 for x in Top15['% Renewable']]
    return Top15['HighRenew']

answer_ten()

Out[13]: Country
China           1
United States   0
Japan           0
United Kingdom  0
Russian Federation 1
Canada          1
Germany         1
India            0
France           1
South Korea     0
Italy             1
Spain            1
Iran              0
Australia        0
Brazil           1
Name: HighRenew, dtype: int64
```

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 [14]: 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()

    # Get the top continents
    Top15['Continent'] = [ContinentDict[country] for country in Top15['Country']]

    # Now set Index as Continent & Group By Population Estimate and apply the aggregate funcs
    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()
```

```
Out[14]:
```

Continent	size	sum	mean	std
Asia	5.0	2.898666e+09	5.797333e+08	6.790979e+08
Australia	1.0	2.331602e+07	2.331602e+07	NaN
Europe	6.0	4.579297e+08	7.632161e+07	3.464767e+07
North America	2.0	3.528552e+08	1.764276e+08	1.996696e+08
South America	1.0	2.059153e+08	2.059153e+08	NaN

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 [19]: 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']]

    # For bin we use pd.cut and 5 bins
    Top15['bins'] = pd.cut(Top15['% Renewable'],5)

    return Top15.groupby(['Continent','bins']).size()

answer_twelve()
```

```
Out[19]: Continent      bins
Asia          (2.212, 15.753]    4
              (15.753, 29.227]    1
Australia     (2.212, 15.753]    1
Europe         (2.212, 15.753]    1
              (15.753, 29.227]    3
              (29.227, 42.701]    2
North America (2.212, 15.753]    1
              (56.174, 69.648]    1
South America (56.174, 69.648]    1
dtype: int64
```

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 [26]: def answer_thirteen():
    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()

    tmp = list()
    Top15['PopEst'] = (Top15['Energy Supply'] / Top15['Energy Supply per Capita'])
    tmp = Top15['PopEst'].tolist()

    Top15['PopEst'] = (Top15['Energy Supply'] / Top15['Energy Supply per Capita']).apply(lambda x: "{:,}".format(x),
    tmp)
    ans = pd.Series(Top15['PopEst'])

    return ans
answer_thirteen()
```

```
Out[26]: Country
China           1,367,645,161.2903225
United States   317,615,384.61538464
Japan            127,409,395.97315437
United Kingdom  63,870,967.741935484
Russian Federation  143,500,000.0
Canada           35,239,864.86486486
Germany          80,369,696.96969697
India             1,276,730,769.2307692
France            63,837,349.39759036
South Korea       49,805,429.864253394
Italy              59,908,256.880733944
Spain              46,443,396.2264151
Iran                77,075,630.25210084
Australia         23,316,017.316017351
Brazil             205,915,254.23728815
Name: PopEst, dtype: object
```

Optional

Use the built in function `plot_optional()` to see an example visualization.

```
In [24]: 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]);

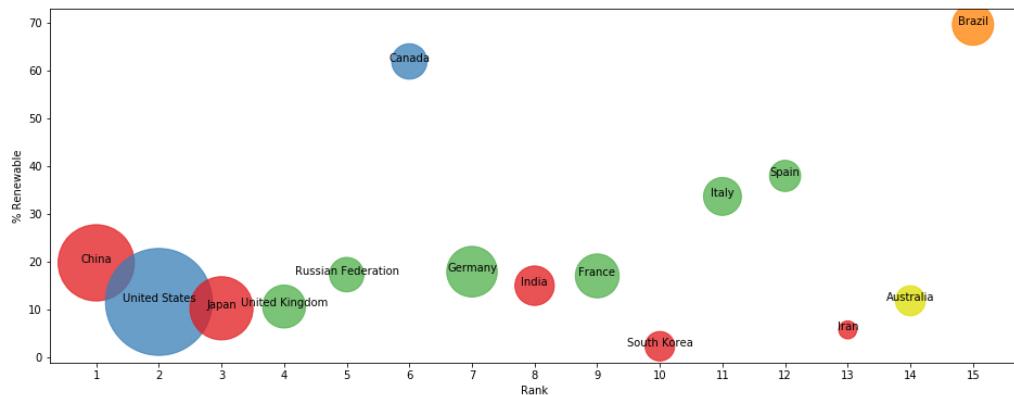
for i, txt in enumerate(Top15.index):
    ax.annotate(txt, [Top15['Rank'][i], Top15['% Renewable'][i]], ha='center')

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 [25]: `#plot_optional() # Be sure to comment out plot_optional() before submitting the assignment!`

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 []: