

# Data Analysis using Energy, Power and GDP Dataset

1

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, for 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 [Sciago 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]: import pandas as pd
import numpy as np
import re

energy = pd.read_excel('E:\Protofolio\Data Analysis using Energy, Power and GDP Dataset\Energy Indicator s.xls',skiprows=18, usecols = [2,3,4,5],header=None, skip_footer=38,na_values=...)
energy.columns = ['Country', 'Energy Supply', 'Energy Supply per Capita', '% Renewable']
energy = energy.replace(r'\s+', '', inplace=True)
energy['Energy Supply'] = energy['Energy Supply']/1000000
energy['Country'] = energy['Country'].str.split("(")[0].str[0]
energy['Country'] = energy['Country'].str.replace('(','+', '')
energy.Country = energy.Country.replace(("United States of America": "United States","United Kingdom of Great Britain and Northern Ireland": "United Kingdom","China, Hong Kong Special Administrative Region" : "Hong Kong","Republic of Korea": "South Korea"))

GDP = pd.read_csv('E:\Protofolio\Data Analysis using Energy, Power and GDP Dataset/world_bank.csv', header=0)
GDP['Country Name'] = GDP['Country Name'].replace(("Korea, Rep.": "South Korea", "Iran, Islamic Rep.": "Iran","Hong Kong SAR, China": "Hong Kong"))
GDP=GDP.rename(columns={"Country Name": "Country", })

ScimEn = pd.read_excel('E:\Protofolio\Data Analysis using Energy, Power and GDP Dataset/scimagojr-3.xls x')

df=pd.merge(GDP,energy, on='Country', how='inner')
Country=pd.merge(df,ScimEn, on='Country', how='inner')
columns_to_keep = ['Rank', 'Country', 'Documents', 'Citable documents', 'Citations', 'Self-citations', 'Citati ons per document', 'H index', 'Energy Supply', 'Energy Supply per Capita', '% Renewable', '2006', '2007', '2008', '2009', '2010', '2011', '2012', '2013', '2014', '2015']
Country = Country[columns_to_keep]

Country=Country.sort_values(by=['Rank'])
Country=Country.set_index('Country')
Country15=Country.head(15)
def answer_one():
    return Country15

answer_one()
```

```
Out[1]:
```

	Rank	Documents	Citable documents	Citations	Self-citations	Citations per document	H index	Energy Supply	Energy Supply per Capita	% Renewable	2006
Country											
China	1	127050	126767	597237	411683	4.70	138	1.271910e+11	93.0	19.754910	3.992331e+12 4.5598
United States	2	96661	94747	792274	265436	8.20	230	9.083800e+10	286.0	11.570980	1.478230e+13 1.5051
Japan	3	30504	30287	223024	61554	7.31	134	1.898400e+10	149.0	10.232820	5.496542e+12 5.6171
United Kingdom	4	20944	20357	206091	37874	9.84	139	7.920000e+09	124.0	10.604700	2.419631e+12 2.4821
Russian Federation	5	18534	18301	34266	12422	1.85	57	3.070900e+10	214.0	17.288680	1.385793e+12 1.5041
Canada	6	17899	17620	215003	40930	12.01	149	1.043100e+10	296.0	61.945430	1.564469e+12 1.5961
Germany	7	17027	16831	140566	27426	8.26	126	1.326100e+10	165.0	17.901530	3.332891e+12 3.4411
India	8	15005	14841	128763	37209	8.58	115	3.195000e+10	26.0	14.969080	1.265894e+12 1.3741
France	9	13153	12973	130632	28801	9.93	114	1.059700e+10	166.0	17.020280	2.607840e+12 2.6669
South Korea	10	11983	11923	114675	22595	9.57	104	1.100700e+10	221.0	2.279353	9.410199e+11 9.924
Italy	11	10964	10794	111850	26661	10.20	106	6.530000e+09	109.0	33.667230	2.202170e+12 2.2341
Spain	12	9428	9330	123336	23064	13.08	115	4.923000e+09	106.0	37.968590	1.414823e+12 1.468
Iran	13	8896	8819	57470	19125	6.46	72	9.172000e+09	119.0	5.707721	3.895523e+11 4.260
Australia	14	8831	8725	90765	15606	10.28	107	5.386000e+09	231.0	11.810810	1.021939e+12 1.0601
Brazil	15	8668	8596	60702	14396	7.00	86	1.214900e+10	59.0	68.648030	1.845080e+12 1.957

2

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="100" 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-dasharr
ay="5,3"/>
<text x="300" y="165" font-family="Verdana" font-size="35">Everything but this!</text>
</svg>
```

3

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 [21]: def answer_three():
    Top15 = answer_one()
    Top15=Top15.reset_index()

    Top15['avgGDP']=Top15['2006']+Top15['2007']+Top15['2008']+Top15['2009']+Top15['2010']+Top15['2011']
    +Top15['2012']+Top15['2013']+Top15['2014']+Top15['2015']
    Top15['avgGDP']=Top15['avgGDP']/10
    Top15=Top15.sort_values(by=['avgGDP'],ascending=False)
    Top15=Top15.set_index('Country')

    avgGDP=Top15['avgGDP']
    return avgGDP

answer_three()
```

```
Out[21]:
```

Country	
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	1.160443e+12
South Korea	1.106715e+12
Iran	9.405429e+11
Name: avgGDP, dtype: float64	

4

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 [22]: def answer_four():
    Top15 = answer_one()
    a = Top15.reset_index()
    GDP_15=a.loc[3,'2015']
    GDP_06=a.loc[3,'2006']
    GDP_06-GDP_15-GDP_06
    return GDP_06

answer_four()
```

```
Out[22]: 246702696075.3999
```

5

What is the mean `Energy Supply per Capita`?

This function should return a single number.

```
In [23]: def answer_five():
    Top15 = answer_one()
    return Top15['Energy Supply per Capita'].mean(axis=0)

answer_five()
```

```
Out[23]: 157.6
```

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 [24]: def answer_six():
    Top15 = answer_one()
    b=Top15.sort_values(by=['% Renewable'],ascending=False)
    b=b.head(1).reset_index()
    b=b[['Country','% Renewable']]
    return b.apply(tuple, axis=1)

answer_six()
```

```
Out[24]: 0 (Brazil, 69.64803)
dtype: object
```

7

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 [25]: def answer_seven():
    Top15 = answer_one()
    Top15['Ratio_Citation'] = Top15['Self-citations']/Top15['Citations']
    q7=Top15.sort_values(by=['Ratio_Citation'],ascending=False)
    q7=q7.head(1).reset_index()
    q7=q7[['Country','Ratio_Citation']]
    return q7.apply(tuple, axis=1)

answer_seven()
```

```
<ipython-input-25-f3ee39a6338b>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['Ratio_Citation'] = Top15['Self-citations']/Top15['Citations']

Out[25]: 0 (China, 0.6893126179389422)
dtype: object
```

8

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 [26]: def answer_eight():
    Top15 = answer_one()
    Top15['PopEst'] = Top15['Energy Supply']/Top15['Energy Supply per Capita']
    q8=Top15.sort_values(by=['PopEst'],ascending=False).reset_index()
    return q8.loc[2,'Country']

answer_eight()
```

```
<ipython-input-26-f49b8a0ee6eb>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['PopEst'] = Top15['Energy Supply']/Top15['Energy Supply per Capita']

Out[26]: 'United States'
```

9

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 [27]: 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()
```

```
<ipython-input-27-54f28dbc36c1>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['PopEst'] = Top15['Energy Supply']/Top15['Energy Supply per Capita']
Top15['Citable docs per Capita'] = Top15['Citable documents']/Top15['PopEst']

Out[27]: 0.7940010435442942
```

```
In [11]: def plot9():
    import matplotlib as plt
    import 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()
```

```
<ipython-input-11-f3f4960fca0b>:6: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['PopEst'] = Top15['Energy Supply'] / Top15['Energy Supply per Capita']
<ipython-input-11-f3f4960fca0b>:7: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['Citable docs per Capita'] = Top15['Citable documents'] / Top15['PopEst']

Out[27]: 0.7940010435442942
```

10

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 [28]: def answer_ten():
    Top15 = answer_one()
    median=Top15['% Renewable'].median(axis=0)
    Top15['HighRenew'] = np.where(Top15['% Renewable']>=median, 1, 0)
    return Top15.loc[:, 'HighRenew']

answer_ten()
```

```
<ipython-input-28-721d1ed8e397>:4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['HighRenew'] = np.where(Top15['% Renewable']>=median, 1, 0)

Out[28]:
```

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: int32	

11

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 [29]: 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'}

Country15['Continent']=Country15.index.map(continentDict)
Country15['PopEst'] = Country15['Energy Supply']/Country15['Energy Supply per Capita']
a = Country15.groupby('Continent').PopEst.agg(['sum', 'mean', 'std'])

Country15.reset_index(inplace=True)

b = Country15.reset.groupby('Continent').Country.agg(['size'])

Continent=pd.merge(b,a, on='Continent', how='outer')
def answer_eleven():
    return Continent

answer_eleven()
```

```
<ipython-input-29-836d1ff6b4ba>:17: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Country15['Continent']=Country15.index.map(continentDict)
<ipython-input-29-836d1ff6b4ba>:18: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['PopEst'] = Country15['Energy Supply']/Country15['Energy Supply per Capita']

Out[29]:
```

Continent	size	sum	mean	std
Asia	5	2.898606e+09	5.797333e+08	6.790979e+08
Australia	1	2.331602e+07	2.331602e+07	NaN
Europe	6	4.579297e+08	7.632161e+07	3.464767e+07
North America	2	3.528552e+08	1.764276e+08	1.996696e+08
South America	1	2.059153e+08	2.059153e+08	NaN

12

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 Multindex of `Continent`, then the bins for % Renewable. Do not include groups with no countries.

```
In [30]: def answer_twelve():
    Top15 = answer_one()
    return "ANSWER"

answer_twelve()
```

```
Out[30]: 'ANSWER'
```

13

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 [31]: def answer_thirteen():
    Top15 = answer_one()
    PopEst=Country15['PopEst'].map('{:,}'.format)

    return PopEst

answer_thirteen()
```

```
<ipython-input-31-9a9a956ad053>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexin
g.html#returning-a-view-versus-a-copy
Top15['PopEst'] = Top15['Energy Supply'] / Top15['Energy Supply per Capita']

Out[31]:
```

Country	
China	1,267,645,161.2903225
United States	317,615,384.61538464
Japan	127,409,395.97315437
United Kingdom	63,870,967.741935484
Russian Federation	143,500,000.0
Canada	35,239,664.86486486
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	25,316,017.316017315
Brazil	205,915,254.23726815
Name: PopEst, dtype: object	

14

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

```
In [17]: def plot_optional():
    import matplotlib as plt
    import matplotlib inline

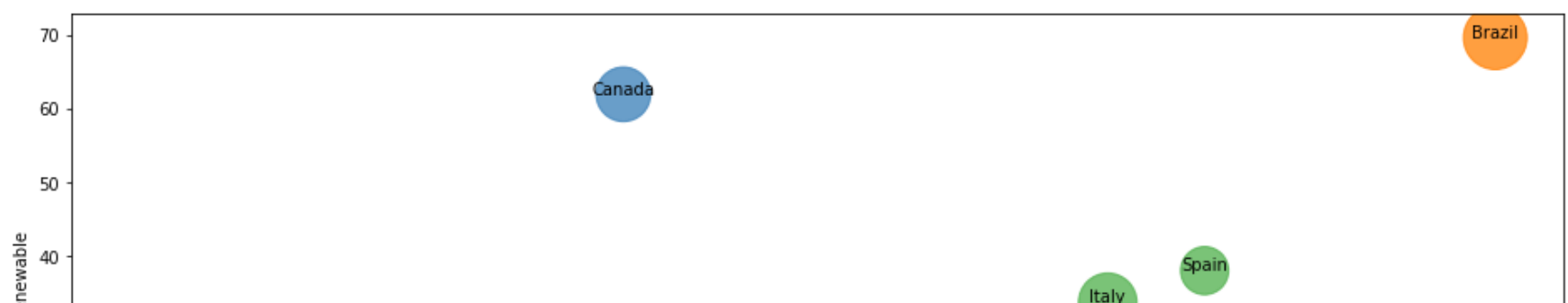
    Top15 = answer_one()
    ax = Top15.plot(x='Rank', y='% Renewable', kind='scatter',
                    c=[ '#4169E1', '#4377ED', '#4DAF4A', '#4DAF4A', '#377EB8', '#4DAF4A', '#E41A1C',
                        '#4DAF4A', '#E41A1C', '#4DAF4A', '#4DAF4A', '#E41A1C', '#D62728', '#F781BF' ],
                    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 countrie s\
2014 GDP, and the color corresponds to the continent.")
```

```
In [32]: plot_optional()
```

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



```
In [1]:
```