

Project: Investigate a Dataset

Table of Contents

- [Introduction](#)
- [Data Wrangling](#)
- [Exploratory Data Analysis](#)
- [Conclusions](#)

Introduction

Dataset Description

FREE DATA FROM WORLD BANK VIA GAPMINDER.ORG.

Gapminder has collected a lot of information about how people live their lives in different countries, tracked across the years, and on a number of different indicators. In the coal_consumption_total.csv dataset, (source: The coal consumption, total dataset in energy/oil folder on <https://www.gapminder.org/data/>) the columns are the countries and value of coal consumption each year from 1965 to 2019

In the oil_consumption_total.csv dataset (source: The oil consumption, total dataset in energy/oil folder on <https://www.gapminder.org/data/>) the columns are the countries and value of oilconsumption each year from 1965 to 2019

and in the consumption_co2_emissions_1000_tonnes.csv, (source: The consumption_co2_emissions_1000_tonnes dataset in environments/emissions/ folder on <https://www.gapminder.org/data/>) the columns are the countries and value of co2 emissions each year from 1990 to 2017

Question(s) for Analysis

- [what is the general trend in the consumption of coal between 2007 and 2017](#)
- [what is the general trend in the consumption of oil between 2007 and 2017](#)
- [what is the general trend in the emission of co2 between 2007 and 2017](#)
- [What top 10 countries have the highest coal consumption, oil consumption and emission of co2](#)
- [What top 10 countries have the lowest coal consumption, oil consumption and emission of co2](#)

Data Wrangling

Import libraries that are needed

```
In [1]: import pandas as pd
import re
import matplotlib.pyplot as plt
matplotlib inline
import seaborn as sns
```

Check relationships between consumption of coal, oil and co2 emissions

Import dataset on total consumption of coals from 1965 to 2019 in all countries

```
In [2]: coal_data = pd.read_csv('coal_consumption_total.csv')
coal_data.head()
```

```
Out[2]:
```

	country	1965	1966	1967	1968	1969	1970	1971	1972	1973	...	2010	2011	2012	2013	2014	2015	2016	2017	2018
0	United Arab Emirates	0	0	0	0	0	0	0	0	0	...	658k	445k	138M	177M	171M	1.84M	2.13M	2.41M	

1 rows × 56 columns

Import dataset on total consumption of oil from 1965 to 2019 in all countries

```
In [3]: oil_data = pd.read_csv('oil_consumption_total.csv')
oil_data.head()
```

```
Out[3]:
```

	country	1965	1966	1967	1968	1969	1970	1971	1972	1973	...	2010	2011	2012	2013	2014	2015	2016	2017	2018
0	United Arab Emirates	69.6k	74.3k	79.8k	89.1k	99k	115k	136k	179k	280k	...	31.3M	33.8M	35.2M	39M	39.3M	42.2M	44.9M	44.2M	46M

1 rows × 56 columns

Import dataset on total co2 emissions from fossil fuels consumption from 1990 to 2017 in all countries

```
In [4]: co2_data = pd.read_csv('consumption_co2_emissions_1000_tonnes.csv')
co2_data.head()
```

```
Out[4]:
```

	country	1990	1991	1992	1993	1994	1995	1996	1997	1998	...	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017
0	Albania	5620	4570	2840	2630	2290	2440	2870	2190	2510	...	6130	6370	6140	6350	6080	5870	6300	5660	5500	5650

1 rows × 29 columns

Data cleaning

Limit data to between the years 2007 to 2017

```
In [5]: cols = ['country'] + [str(x) for x in range(2007,2018)] # create list to get columns from 2009 to 2017,
#with country inclusive
Out[5]:
```

```
['country',
'2007',
'2008',
'2009',
'2010',
'2011',
'2012',
'2013',
'2014',
'2015',
'2016',
'2017']
```

To combine the data from the three indicators to a single data

```
In [6]: def clean_data(data, cols, var):
'''
docstring: clean data and rearrange data according to years https://pandas.pydata.org/docs/reference/api/pandas.melt.html
input: data - DataFrame, columns to melt
cols - string, columns to keep
var - string, name of column values
output: melted DataFrame
'''
data2 = data[cols].copy()
data2 = pd.melt(data2, id_vars=['country'], var_name='year', value_name = var)
return data2
```

```
In [7]: # merge coal_data, oil_data and co2_data after they have been melted
data3 = clean_data(data = coal_data, cols= cols, var = 'coal')\
.merge(clean_data(data = oil_data, cols= cols, var = 'oil'), on = ['country', 'year'])\
.merge(clean_data(data = co2_data, cols= cols, var = 'co2'), on = ['country', 'year'])
```

```
In [8]: data.head()
```

```
Out[8]:
```

	country	year	coal	oil	co2
0	United Arab Emirates	2007	136000	28700000	203000
1	Argentina	2007	1230M	24.2M	161k
2	Australia	2007	55900000	42100000	358000
3	Austria	2007	3.86M	12.6M	101k
4	Azerbaijan	2007	4930	4.52M	32.2k

To convert the format of the values e.g k is 1000, M is 10⁶ and B is 10⁹

```
In [9]: for column in data.columns[2:]:
data[column] = data[column].replace({'k':'*1e3', 'M':'*1e6', 'B':'*1e9'}, regex = True).map(pd.eval)\
data.head()
```

```
Out[9]:
```

	country	year	coal	oil	co2
0	United Arab Emirates	2007	136000	28700000	203000
1	Argentina	2007	1230000	24200000	161000
2	Australia	2007	55900000	42100000	358000
3	Austria	2007	3860000	12600000	101000
4	Azerbaijan	2007	4930	4520000	32200

To check for null values

```
In [10]: data.isna().sum()
```

```
Out[10]:
```

```
country      0
year          0
coal          0
oil           0
co2           0
dtype: int64
```

To check for duplicated values

```
In [11]: data.duplicated().sum()
```

```
Out[11]: 0
```

Check datatypes and change year column to datetime datatype

```
In [12]: data['year'] = pd.to_datetime(data['year'])
```

```
In [13]: data.dtypes
```

```
Out[13]:
```

```
country      object
year          datetime64[ns]
coal          int64
oil           int64
co2           int64
dtype: object
```

Check for invalid values, outliers

```
In [14]: data.describe()
```

```
Out[14]:
```

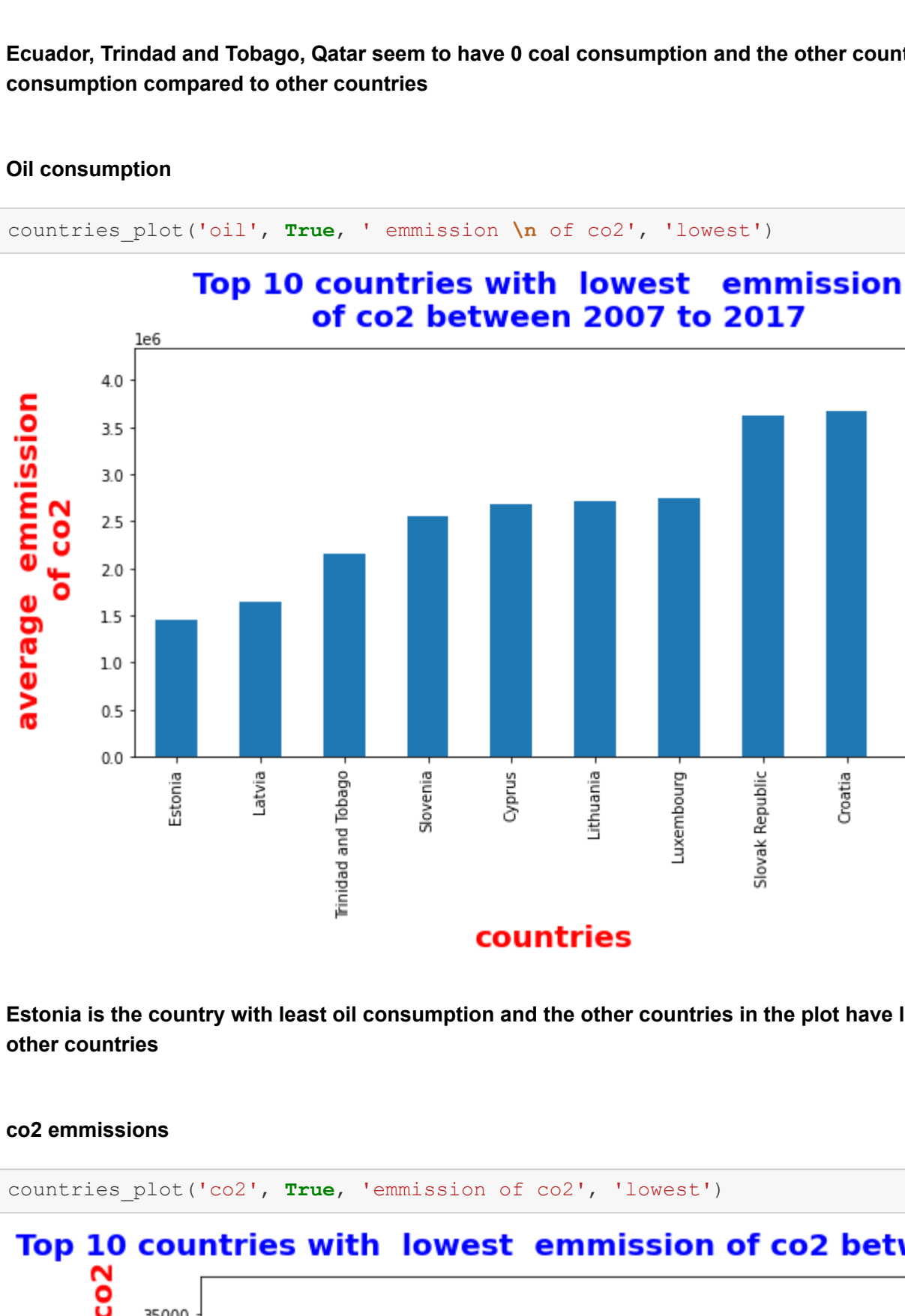
	coal	oil	co2
count	8.030000e+02	8.030000e+02	8.030000e+02
mean	4.993194e+07	5.263674e+07	4.294369e+05
std	2.022715e+08	1.126780e+08	1.133047e+06
min	0.000000e+00	1.330000e+06	7.770000e+03
25%	8.115000e+06	8.100000e+06	5.735000e+04
50%	4.280000e+06	1.590000e+07	1.140000e+05
75%	1.765000e+07	5.220000e+07	3.530000e+05
max	1.970000e+09	9.080000e+08	8.560000e+06

EDA

average consumption of coal by all countries each year between 2007 and 2017

```
In [15]: def trend_plot(column, name):
'''
docstring: plot trend in average consumption of coal, oil and co2 emissions by all countries each year
between 2007 and 2017
input: column - string, columns to plot
name - string, name of plot
output: plot
'''
plt.figure(figsize = (10,5))
data.groupby('year').mean()[column].plot()
plt.title(f'Trend in {name} from 2007 to 2017', fontdict = {'fontsize': 14,
'verticalalignment': 'center_baseline', 'color': 'red', 'pad = 15})
plt.xlabel('years', fontdict = {'fontsize': 14,
'verticalalignment': 'center_baseline', 'color': 'red'})
plt.ylabel(f'{average (name)}', fontdict = {'fontsize': 14,
'verticalalignment': 'center_baseline', 'color': 'red'})
```

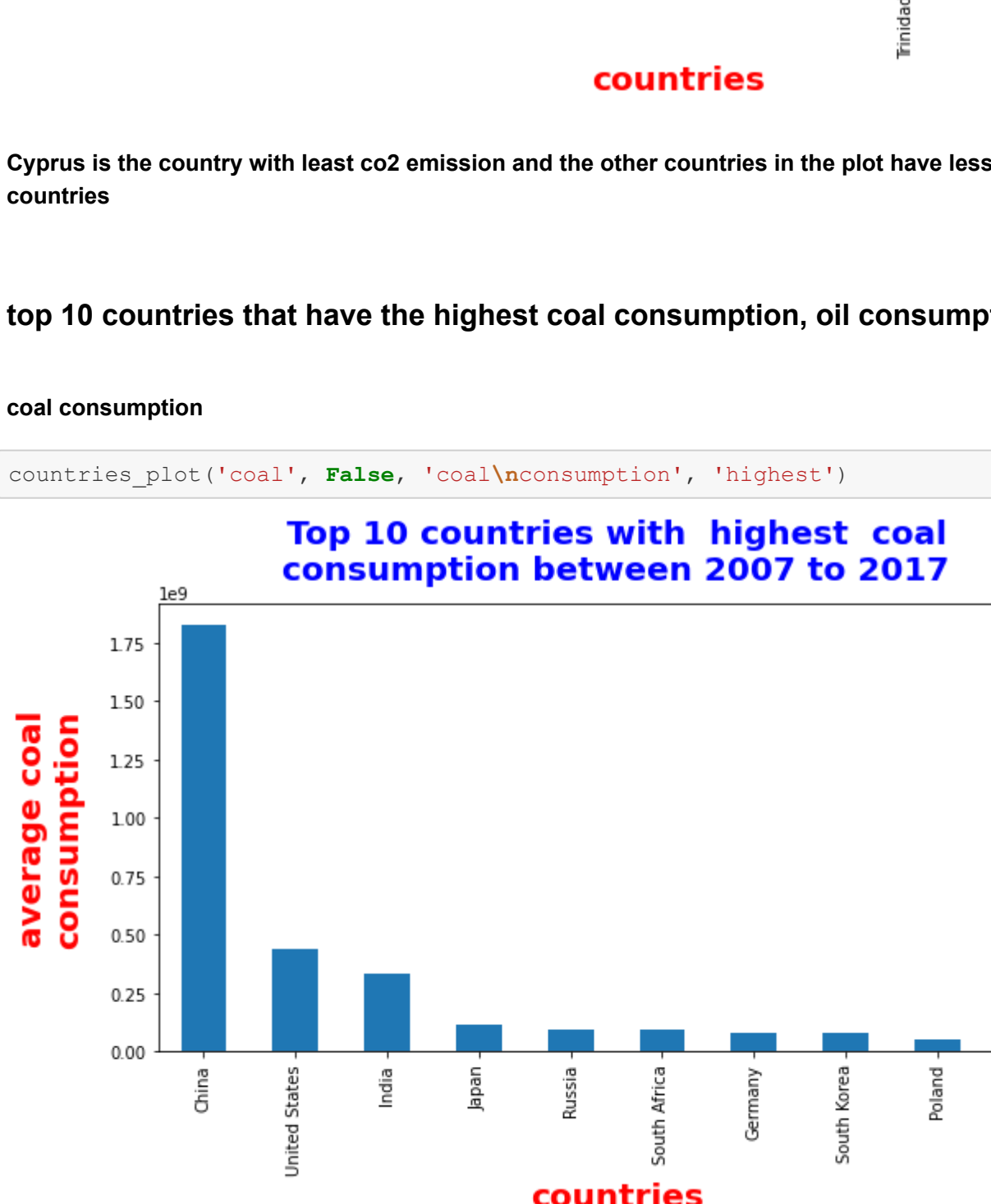
```
In [16]: trend_plot('coal', 'consumption of coal')
```



From 2009 there has been exponential growth in coal consumption and a decline from 2014

average consumption of oil between 2007 and 2017

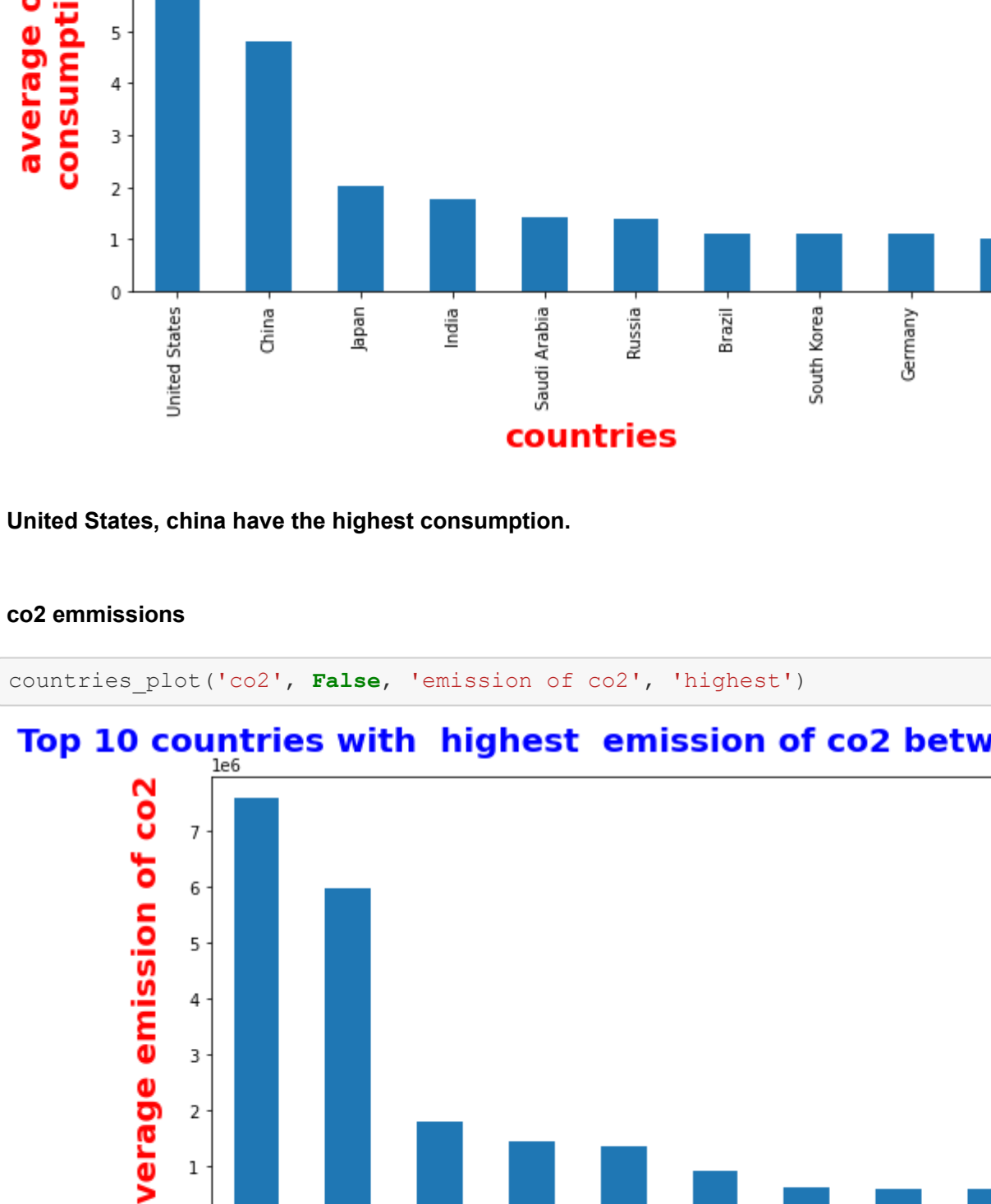
```
In [17]: trend_plot('oil', 'consumption of oil')
```



From 2009 there has been exponential growth in oil consumption

average co2 emissions from fossil fuels consumption between 2007 and 2017

```
In [18]: trend_plot('co2', 'co2 emissions')
```



From 2009 there has been exponential growth in co2 emission

average coal consumption, oil consumption and emission of co2 per country

```
In [19]: def countries_plot(column, sort, name, level):
'''
docstring: plot top 10 countries that have the lowest and highest coal consumption,
oil consumption and emission of co2
input: column - string, columns to plot
sort - bool, if it should sort in ascending order
name - string, name of plot
level - string, lowest or highest
output: plot
'''
plt.figure(figsize = (10,5))
data.groupby('country').mean()[column].sort_values(ascending = sort).head(10).plot(kind = 'bar')
plt.xlabel('countries', fontdict = {'fontsize': 12, 'fontweight': 'bold', 'color': 'red'})
plt.ylabel(f'{average (name)}', fontdict = {'fontsize': 12, 'fontweight': 'bold', 'color': 'red'})
plt.title(f'Top 10 countries with {level} {name} between 2007 to 2017',
fontdict = {'fontsize': 12, 'fontweight': 'bold', 'color': 'blue', 'pad = 15})
```

top 10 countries that have the lowest coal consumption, oil consumption and emission of co2

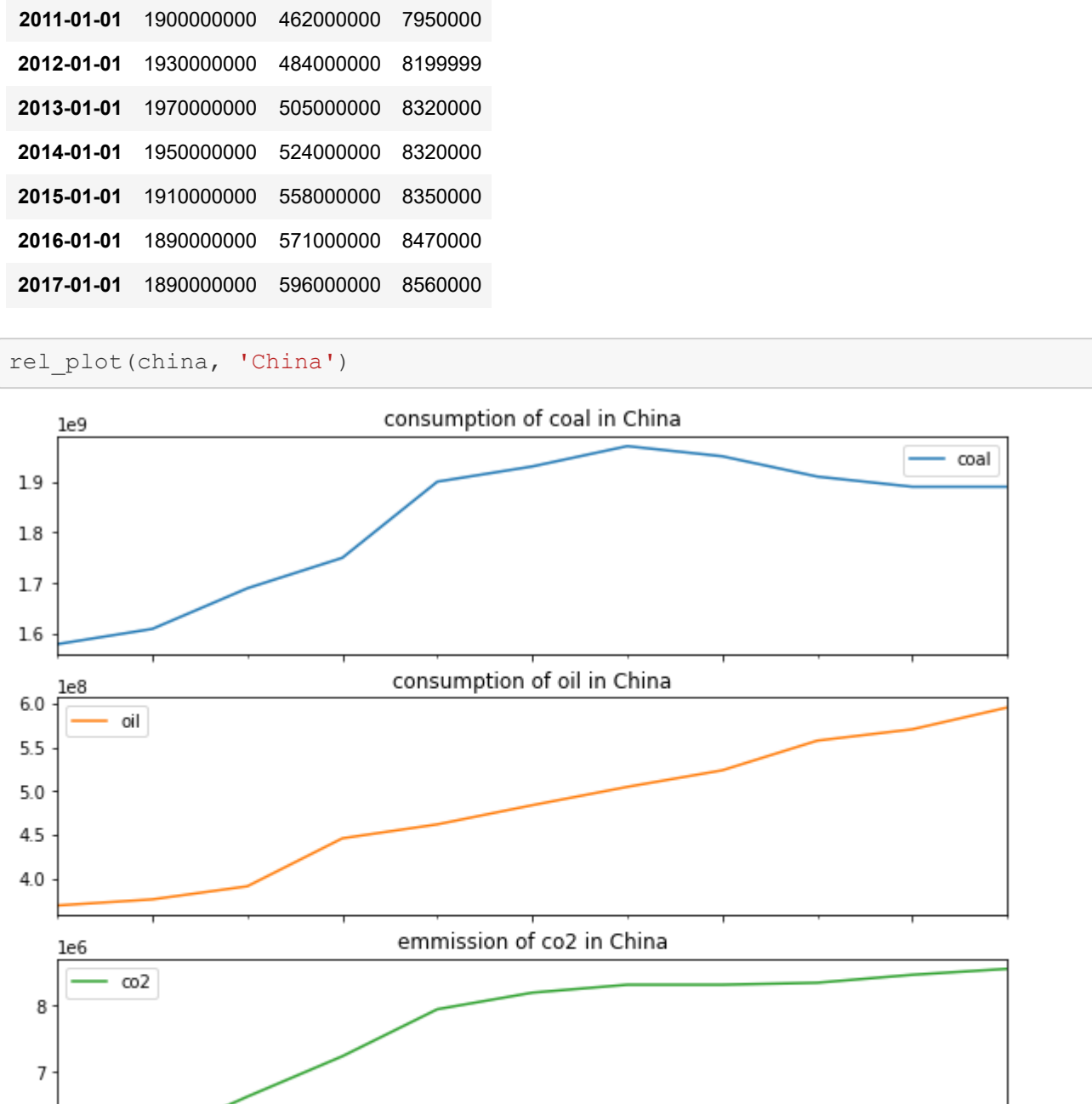
```
In [20]: countries_plot('coal', True, 'coal\consumption', 'lowest')
```



Ecuador, Trinidad and Tobago, Qatar seem to have 0 coal consumption and the other countries in the plot have lesser coal consumption compared to other countries

Oil consumption

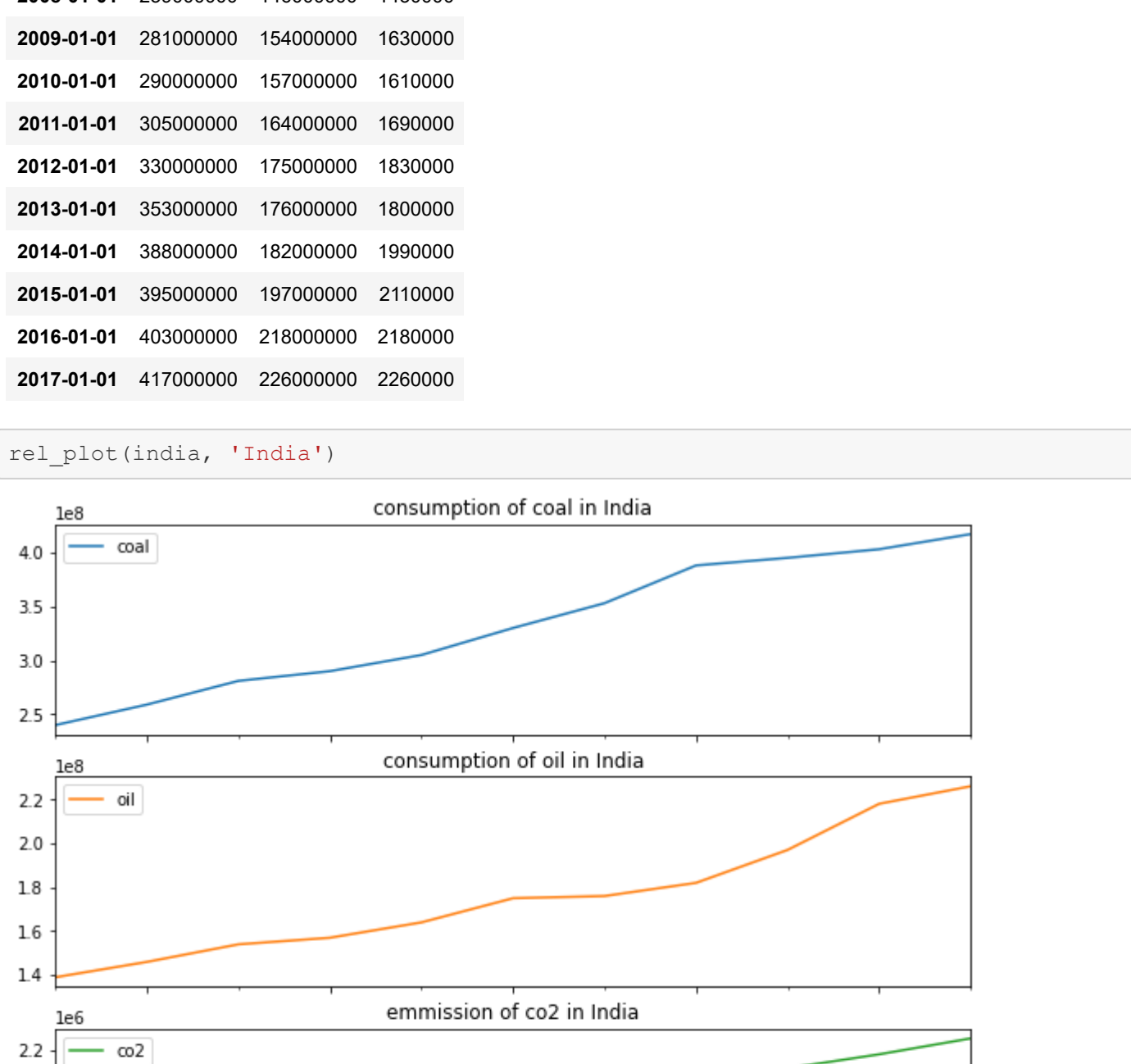
```
In [21]: countries_plot('oil', True, 'oil\consumption', 'lowest')
```



Estonia is the country with least oil consumption and the other countries in the plot have lesser oil consumption compared to other countries

co2 emissions

```
In [22]: countries_plot('co2', True, 'emission of co2', 'lowest')
```



Cyprus is the country with least co2 emission and the other countries in the plot have lesser co2 emission compared to other countries

top 10 countries that have the highest coal consumption, oil consumption and emission of co2

```
In [23]: countries_plot('coal', False, 'coal\consumption', 'highest')
```



United States, china have the highest consumption.

Oil consumption

```
In [24]: countries_plot('oil', False, 'oil\consumption', 'highest')
```


United States, china have the highest consumption.

co2 emissions

```
In [25]: countries_plot('co2', False, 'emission of co2', 'highest')
```


There seems to be a positive correlation between the three indicators in United States

Data on China

```
In [30]: china = count(df = data, countryname = "country == 'China'")
china
```

```
Out[30]:
```

	coal	oil	co2
year			
2007-01-01	15800000000	3690000000	5570000
2008-01-01	16100000000	3760000000	5990000
2009-01-01	16900000000	3910000000	6630000
2010-01-01	17500000000	4480000000	7240000
2011-01-01	19000000000	4620000000	7950000
2012-01-01	19300000000	4840000000	8199999
2013-01-01	19700000000	5050000000	8320000
2014-01-01	19500000000	5240000000	8350000
2015-01-01	19100000000	5580000000	8350000
2016-01-01	18900000000	5710000000	8470000
2017-01-01	18900000000	5960000000	8560000

```
In [31]: rel_plot(china, 'China')
```



There seems to be positive correlation between the three indicators in China

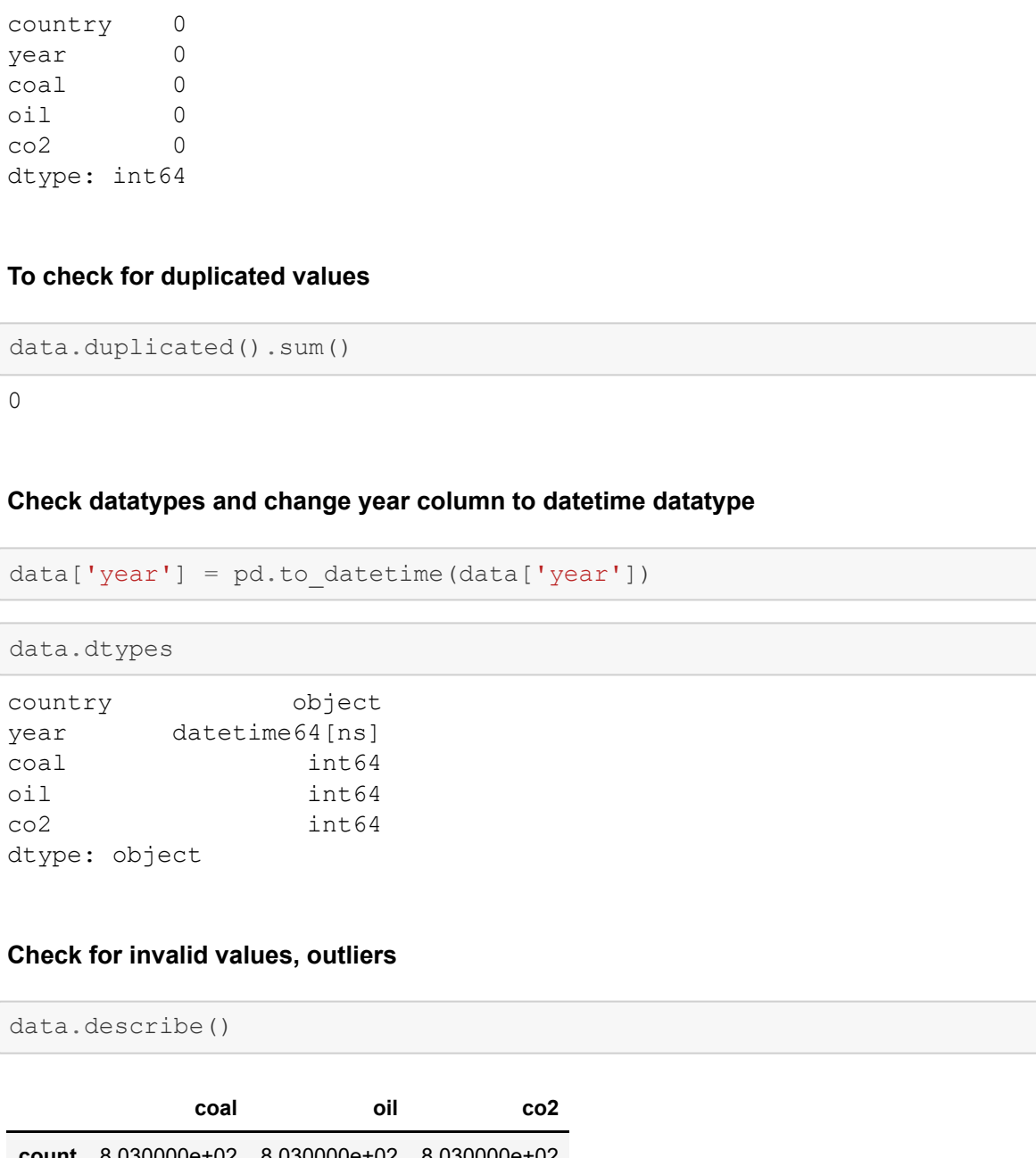
Data on India

```
In [32]: india = count(df = data, countryname = "country == 'India'")
india
```

```
Out[32]:
```

	coal	oil	co2
year			
2007-01-01	2400000000	1390000000	1300000
2008-01-01	2590000000	1460000000	1430000
2009-01-01	2810000000	1540000000	1630000
2010-01-01	2900000000	1570000000	1610000
2011-01-01	3050000000	1640000000	1690000
2012-01-01	3300000000	1760000000	1830000
2013-01-01	3530000000	1780000000	1800000
2014-01-01	3880000000	1820000000	1990000
2015-01-01	3950000000	1970000000	2110000
2016-01-01	4030000000	2180000000	2180000
2017-01-01	4170000000	2260000000	2260000

```
In [33]: rel_plot(india, 'India')
```



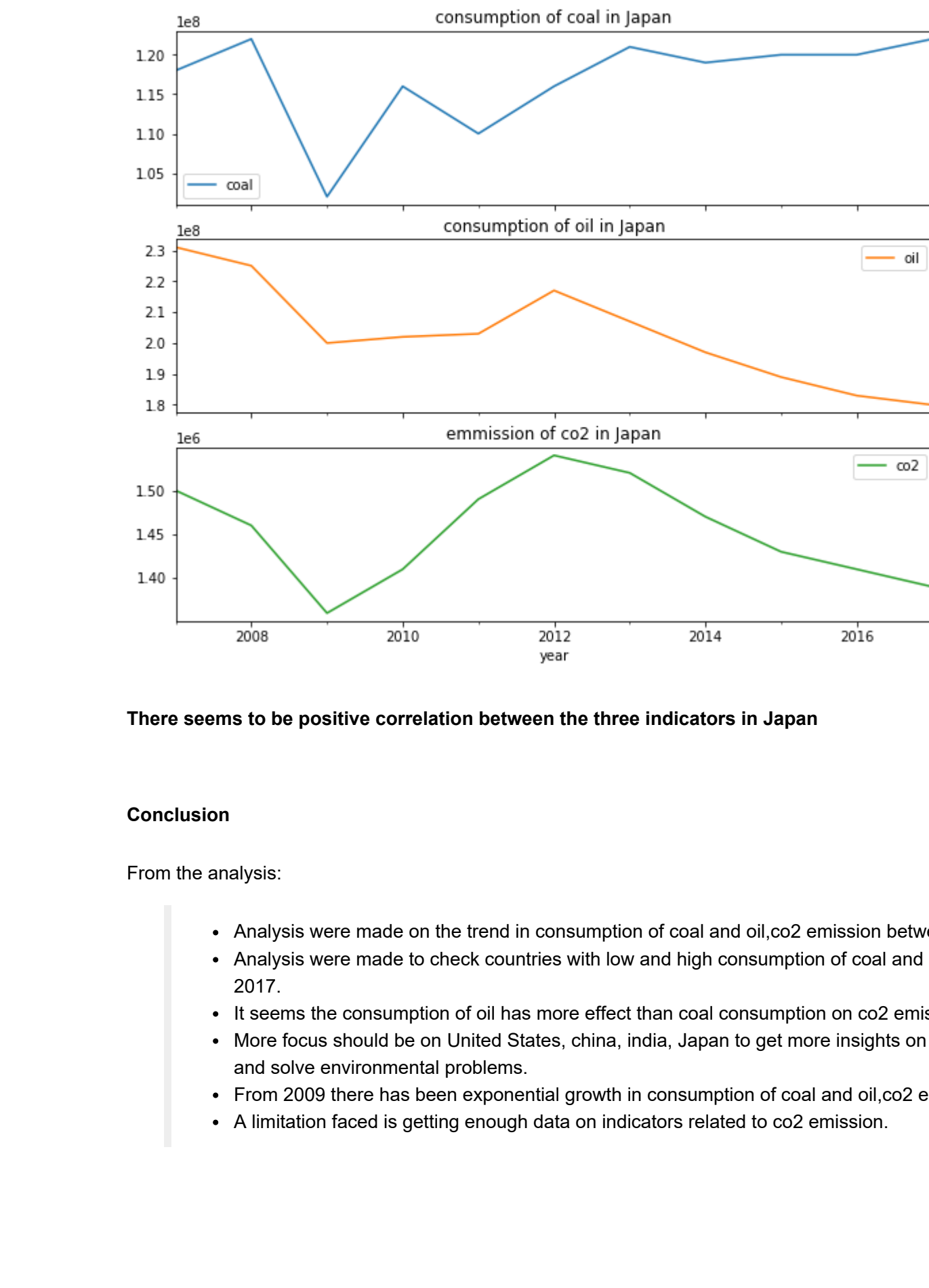
There seems to be positive correlation between the three indicators in India

Data on Japan

```
In [34]: japan = count(df = data, countryname = "country == 'Japan'")
japan
```

```
Out[34]:
```

	coal	oil	co2
year			
2007-01-01	1180000000	2310000000	1500000
2008-01-01	1220000000	2250000000	1460000
2009-01-01	1020000000	2000000000	1300000
2010-01-01	1160000000	2020000000	1410000
2011-01-01	1100000000	2030000000	1490000
2012-01-01	1160000000	2070000000	1540000
2013-01-01	1190000000	1970000000	1470000
2014-01-01	1200000000	1890000000	1430000
2015-01-01	1200000000	1830000000	1410000
2016-01-01	1220000000	1800000000	1390000
2017-01-01	1220000000	1800000000	1390000



There seems to be positive correlation between the three indicators in Japan

Conclusion

From the analysis:

- Analysis were made on the trend in consumption of coal and oil,co2 emission between 2007 and 2017.
- Analysis were made to check countries with low and high consumption of coal and oil,co2 emission between 2007 and 2017.
- It seems the consumption of oil has more effect than coal consumption on co2 emission.
- More focus should be on United States, china, india, Japan to get more insights on why they have high co2 emission and solve environmental problems.
- From 2009 there has been exponential growth in consumption of coal and oil,co2 emission.
- A limitation faced is getting enough data on indicators related to co2 emission.