

# Data Processing Project

Import useful python libraries such as Numpy for numerical computation, Pandas for data analysis, Matplotlib & Seaborn for Data Visualization

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

## 1. understand the data

we have First understand the data that means how much rows, columns present in it. Then we have to understand maximum, minimum value of data, standard deviation, mean, mode.

```
In [2]: co2 = pd.read_csv('D:\machine learning project\co2_emission.csv')
```

Out[2]:

	Entity	Code	Year	Annual CO <sub>2</sub> emissions (tonnes)
0	Alghanistan	AFG	1949	14656.00
1	Alghanistan	AFG	1950	84272.00
2	Alghanistan	AFG	1951	91600.00
3	Alghanistan	AFG	1952	91600.00
4	Alghanistan	AFG	1953	106256.00
5	Alghanistan	AFG	1954	106256.00
6	Alghanistan	AFG	1955	136960.00
7	Alghanistan	AFG	1956	132000.00
8	Alghanistan	AFG	1957	290320.00
9	Alghanistan	AFG	1958	329760.00
10	Alghanistan	AFG	1959	384571.42
11	Alghanistan	AFG	1960	413883.42
12	Alghanistan	AFG	1961	480797.70
13	Alghanistan	AFG	1962	68854.27
14	Alghanistan	AFG	1963	70073.98
15	Alghanistan	AFG	1964	83850.83
16	Alghanistan	AFG	1965	100901.53
17	Alghanistan	AFG	1966	109118.82
18	Alghanistan	AFG	1967	126186.11
19	Alghanistan	AFG	1968	122339.69
20	Alghanistan	AFG	1969	941231.98
21	Alghanistan	AFG	1970	167039.69
22	Alghanistan	AFG	1971	169393.41
23	Alghanistan	AFG	1972	153116.69
24	Alghanistan	AFG	1973	163724.40
25	Alghanistan	AFG	1974	191567.68
26	Alghanistan	AFG	1975	212455.68
27	Alghanistan	AFG	1976	198425.38
28	Alghanistan	AFG	1977	256964.67
29	Alghanistan	AFG	1978	215543.66

20823	Zimbabwe	ZWE	1988	16029402.93
20824	Zimbabwe	ZWE	1989	16113953.54
20825	Zimbabwe	ZWE	1990	15562924.66
20826	Zimbabwe	ZWE	1991	1585834.05
20827	Zimbabwe	ZWE	1992	1651743.67
20828	Zimbabwe	ZWE	1993	1625910.52
20829	Zimbabwe	ZWE	1994	1768446.69
20830	Zimbabwe	ZWE	1995	1502945.61
20831	Zimbabwe	ZWE	1996	1405950.82
20832	Zimbabwe	ZWE	1997	1429421.06
20833	Zimbabwe	ZWE	1998	1413953.04
20834	Zimbabwe	ZWE	1999	1572748.13
20835	Zimbabwe	ZWE	2000	136181.23
20836	Zimbabwe	ZWE	2001	1250195.38
20837	Zimbabwe	ZWE	2002	1169700.10
20838	Zimbabwe	ZWE	2003	1060340.13
20839	Zimbabwe	ZWE	2004	942548.13
20840	Zimbabwe	ZWE	2005	1069642.16
20841	Zimbabwe	ZWE	2006	1058646.23
20842	Zimbabwe	ZWE	2007	953041.17
20843	Zimbabwe	ZWE	2008	771874.78
20844	Zimbabwe	ZWE	2009	9513908.33
20845	Zimbabwe	ZWE	2010	768101.98
20846	Zimbabwe	ZWE	2011	943980.98
20847	Zimbabwe	ZWE	2012	761460.24
20848	Zimbabwe	ZWE	2013	1152629.29
20849	Zimbabwe	ZWE	2014	1166348.41
20850	Zimbabwe	ZWE	2015	1059760.94
20851	Zimbabwe	ZWE	2016	962646.88
20852	Zimbabwe	ZWE	2017	1059717.47

20853 rows x 4 columns

```
In [3]: co2.shape
```

Out[3]: (20853, 4)

```
In [4]: co2.head()
```

Out[4]:

	Entity	Code	Year	Annual CO <sub>2</sub> emissions (tonnes)
0	Alghanistan	AFG	1949	14656.0
1	Alghanistan	AFG	1950	84272.0
2	Alghanistan	AFG	1951	91600.0
3	Alghanistan	AFG	1952	91600.0
4	Alghanistan	AFG	1953	106256.0

```
In [5]: co2.tail()
```

Out[5]:

	Entity	Code	Year	Annual CO <sub>2</sub> emissions (tonnes)
20848	Zimbabwe	ZWE	2013	1152629.29
20849	Zimbabwe	ZWE	2014	1166348.41
20850	Zimbabwe	ZWE	2015	1059760.94
20851	Zimbabwe	ZWE	2016	962646.88
20852	Zimbabwe	ZWE	2017	1059717.47

```
In [6]: co2.info()
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 20853 entries, 0 to 20852  
Data columns (total 4 columns):  
Entity 20853 non-null object  
Code 18646 non-null object  
Year 20853 non-null int64  
Annual CO<sub>2</sub> emissions (tonnes) 20853 non-null float64  
dtypes: float64(1), int64(1), object(2)  
memory usage: 651.7 KB

```
In [7]: co2.describe()
```

Out[7]:

	Year	Annual CO <sub>2</sub> emissions (tonnes)
count	20853.000000	2.795272e+04
mean	1993.339424	1.936517e+08
min	1751.000000	1.345153e+09
max	1978.000000	4.255223e+08
25%	1932.000000	3.187680e+05
50%	1971.000000	3.828880e+06
75%	1996.000000	3.706896e+07
max	1971.000000	3.615320e+10

```
In [8]: co2.columns
```

Out[8]: Index(['Entity', 'Code', 'Year', 'Annual CO<sub>2</sub> emissions (tonnes)'], dtype='object')

```
In [9]: co2.nunique()
```

Out[9]:

Entity	233
Code	222
Year	267
Annual CO <sub>2</sub> emissions (tonnes)	13892
dtype:	int64

```
In [10]: co2['Entity'].unique()
```

array(['Alghanistan', 'Africa', 'Albania', 'Algeria', 'Americas (other)', 'Andorra', 'Angola', 'Antarctic Fisheries', 'Antigua and Barbuda', 'Argentina', 'Armenia', 'Aruba', 'Asia and Pacific (other)', 'Australia', 'Austria', 'Azerbaijan', 'Bahamas', 'Bahrain', 'Bangladesh', 'Barbados', 'Belarus', 'Belgium', 'Belize', 'Benin', 'Bermuda', 'Bhutan', 'Bolivia', 'Bosnia and Herzegovina', 'Botswana', 'Brazil', 'British Virgin Islands', 'Brunei', 'Bulgaria', 'Burkina Faso', 'Burundi', 'Cambodia', 'Cameroon', 'Canada', 'Cape Verde', 'Cayman Islands', 'Central African Republic', 'Chad', 'Chile', 'China', 'Christmas Island', 'Colombia', 'Comoros', 'Republic of the Congo', 'Cook Islands', 'Costa Rica', 'Cote d'Ivoire', 'Croatia', 'Cuba', 'Cyprus', 'Czech Republic', 'Czechoslovakia', 'Democratic Republic of Republic of the Congo', 'Denmark', 'Djibouti', 'Dominica', 'Dominican Republic', 'EU-28', 'Ecuador', 'Egypt', 'El Salvador', 'Equatorial Guinea', 'Eritrea', 'Estonia', 'Ethiopia', 'Europe (other)', 'Faeroe Islands', 'Falkland Islands', 'Fiji', 'Finland', 'France', 'French Guiana', 'French Polynesia', 'Gabon', 'Gambia', 'Georgia', 'Germany', 'Ghana', 'Gibraltar', 'Greece', 'Greenland', 'Grenada', 'Guadeloupe', 'Guatemala', 'Guinea', 'Guinea-Bissau', 'Guyana', 'Haiti', 'Honduras', 'Hong Kong', 'Hungary', 'Iceland', 'India', 'Indonesia', 'International transport', 'Iran', 'Iraq', 'Ireland', 'Israel', 'Italy', 'Jamaica', 'Japan', 'Jordan', 'Kazakhstan', 'Kenya', 'Kiribati', 'Kuwait', 'Kyrgyzstan', 'Kyrgyzstan', 'Laos', 'Latvia', 'Lebanon', 'Lesotho', 'Liberia', 'Lilya', 'Liechtenstein', 'Lithuania', 'Luxembourg', 'Macao', 'Macedonia', 'Madagascar', 'Malawi', 'Malaysia', 'Maldives', 'Mali', 'Malta', 'Marshall Islands', 'Martinique', 'Mauritania', 'Mauritius', 'Mexico', 'Micronesia (country)', 'Middle East', 'Moldova', 'Mongolia', 'Montenegro', 'Montserrat', 'Morocco', 'Mozambique', 'Myanmar', 'Namibia', 'Nauru', 'Nepal', 'Netherlands', 'Niger', 'Nigeria', 'Niue', 'North Korea', 'Norway', 'Oman', 'Pakistan', 'Palau', 'Palestine', 'Panama', 'Papua New Guinea', 'Paraguay', 'Peru', 'Philippines', 'Poland', 'Portugal', 'Qatar', 'Reunion', 'Romania', 'Russia', 'Rwanda', 'Saint Helena', 'Saint Kitts and Nevis', 'Saint Lucia', 'Saint Pierre and Miquelon', 'Saint Vincent and the Grenadines', 'Samoa', 'Sao Tome and Principe', 'Saudi Arabia', 'Senegal', 'Serbia', 'Seychelles', 'Sierra Leone', 'Singapore', 'Sierra Leone', 'South Africa', 'South Korea', 'Saint Maarten (Dutch part)', 'Slovakia', 'Slovenia', 'Solomon Islands', 'Somalia', 'South Africa', 'South Korea', 'South Sudan', 'Spain', 'Sri Lanka', 'Statistical differences', 'Sudan', 'Suriname', 'Swaziland', 'Sweden', 'Switzerland', 'Syria', 'Tajikistan', 'Tanzania', 'Thailand', 'Timor', 'Togo', 'Tonga', 'Trinidad and Tobago', 'Turkmenistan', 'Turks and Caicos Islands', 'Tuvalu', 'Turkey', 'Turkmenistan', 'Turks and Caicos Islands', 'Tuvalu', 'Uganda', 'Ukraine', 'United Arab Emirates', 'United Kingdom', 'United States', 'Uruguay', 'Uzbekistan', 'Vanuatu', 'Venezuela', 'Vietnam', 'Wallis and Futuna Islands', 'World', 'Yemen', 'Zambia', 'Zimbabwe'], dtype=object)

```
In [11]: co2.Year.min()
```

Out[11]: 1751

```
In [12]: co2.Year.max()
```

Out[12]: 2017

```
In [13]: co2['Annual CO2 emissions (tonnes)'].max()
```

Out[13]: 36153261645.0

```
In [14]: co2['Annual CO2 emissions (tonnes)'].min()
```

Out[14]: -625522256.7

```
In [15]: condition=co2['Annual CO2 emissions (tonnes)']>=0
```

```
df=co2[condition]
```

```
In [17]: df['Annual CO2 emissions (tonnes)'].min()
```

Out[17]: 0.0

```
In [18]: df['Annual CO2 emissions (tonnes)'].max()
```

Out[18]: 36153261645.0

```
In [19]: df.info()
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 20813 entries, 0 to 20812  
Data columns (total 4 columns):  
Entity 20813 non-null object  
Code 18646 non-null object  
Year 20813 non-null int64  
Annual CO<sub>2</sub> emissions (tonnes) 20813 non-null float64  
dtypes: float64(1), int64(1), object(2)  
memory usage: 813.0 KB

## 2. Data Cleaning

Data cleaning is the time consuming process. This is the most important steps of data science life cycle. Here we have to perform function on missing value.

```
In [20]: co2.isnull().sum()
```

Out[20]:

Entity	0
Code	2207
Year	0
Annual CO <sub>2</sub> emissions (tonnes)	0
dtype:	int64

```
In [21]: drop_col = co2.drop(['Code'], axis=1)
```

```
In [22]: drop_col.head()
```

Out[22]:

	Entity	Year	Annual CO <sub>2</sub> emissions (tonnes)
0	Alghanistan	1949	14656.00
1	Alghanistan	1950	84272.00
2	Alghanistan	1951	91600.00
3	Alghanistan	1952	91600.00
4	Alghanistan	1953	106256.00

```
In [23]: drop_col.isnull().sum()
```

Out[23]:

Entity	0
Year	0
Annual CO <sub>2</sub> emissions (tonnes)	0
dtype:	int64

## 3. Visualisation

Data visualisation is the process through which we can give data to meaning full insight.

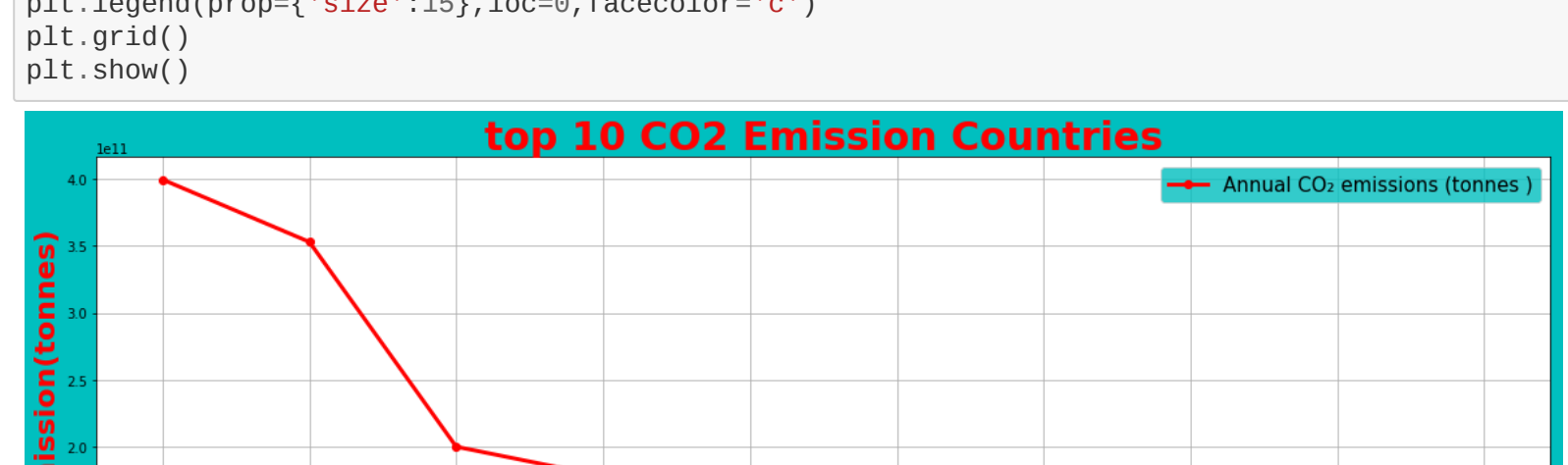
```
In [24]: df.groupby('Entity')[['Annual CO2 emissions (tonnes)']].sum().sort_values(ascending=True)[1:11]
```

Out[24]:

Entity	153888.00
Antarctic Fisheries	153888.00
Taxila	246022.58
Niue	271578.58
Saint Helena	340388.79
Wallis and Futuna Islands	429573.15
Christmas Island	1330822.00
Nontserrat	1539754.12
Kiribati	1913894.59
Liechtenstein	2080665.44
Cook Islands	2944872.31
Name: Annual CO <sub>2</sub> emissions (tonnes), dtype: float64	

```
In [25]: plt.figure(figsize=(20,7), facecolor='c')
plt.plot(df.groupby('Entity')['Annual CO2 emissions (tonnes)'].sum().sort_values(ascending=True)[1:11].color='r', linewidth=3, marker='o')
```

plt.title('Least 10 CO<sub>2</sub> Emission Countries', fontsize=30, color='r', fontweight='bold')  
plt.xlabel('Countries', fontsize=22, color='r', fontweight='bold')  
plt.ylabel('CO<sub>2</sub> Emission (tonnes)', fontsize=22, color='r', fontweight='bold')  
plt.legend(prop={'size':15}, loc=6, facecolor='c')  
plt.grid()
plt.show()



```
In [26]: df.groupby('Entity')[['Annual CO2 emissions (tonnes)']].sum().sort_values(ascending=False)[1:11]
```

Out[26]:

Entity	3.993783e+11
United States	3.993783e+11
EU-28	3.285535e+11
China	2.981256e+11
Asia and Pacific (other)	1.789573e+11
Europe (other)	1.574562e+11
Russia	1.065931e+11
Americas (other)	9.86416e+10
Germany	9.85656e+10
United Kingdom	7.79716e+10
Japan	6.230461e+10
Name: Annual CO <sub>2</sub> emissions (tonnes), dtype: float64	

```
In [27]: plt.figure(figsize=(20,7), facecolor='c')
plt.plot(df.groupby('Entity')['Annual CO2 emissions (tonnes)'].sum().sort_values(ascending=False)[1:11].color='r', linewidth=3, marker='o')
```

plt.title('top 10 CO<sub>2</sub> Emission Countries', fontsize=30, color='r', fontweight='bold')  
plt.xlabel('Countries', fontsize=22, color='r', fontweight='bold')  
plt.ylabel('CO<sub>2</sub> Emission (tonnes)', fontsize=22, color='r', fontweight='bold')  
plt.legend(prop={'size':15}, loc=6, facecolor='c')  
plt.grid()
plt.show()



## Group the year in ratio 50

```
In [28]: x1 = df['Year']<=1800
p1 = df[x1]
```

```
In [29]: x2 = df['Year']>=1850
y2 = df['Year']>=1800
p2 = df[x2&y2]
```

```
In [30]: x3 = df['Year']<=1900
y3 = df['Year']>=1850&1900
p3 = df[x3&y3]
```

```
In [31]: x4 = df['Year']<=1950
y4 = df['Year']>=1900
p4 = df[x4&y4]
```

```
In [32]: x5 = df['Year']<=2000
y5 = df['Year']>=1950
p5 = df[x5&y5]
```

```
In [33]: x6 = df['Year']>=2000
p6 = df[x6]
```

## 1. From year 1750 to 1800

```
In [34]: p1.groupby('Entity')[['Annual CO2 emissions (tonnes)']].sum().sort_values(ascending=False)[1:11]
```

Out[34]:

Entity	771850912.0
Poland	766523456.0
Germany	4920752.0
United States	252216.0
Canada	58624.0
Americas (other)	58624.0
Statistical differences	0.0
International transport	0.0
Name: Annual CO <sub>2</sub> emissions (tonnes), dtype: float64	

```
In [35]: plt.style.use('ggplot')
plt.figure(figsize=(15,5))
p1.groupby('Entity')[['Annual CO2 emissions (tonnes)']].sum().sort_values(ascending=False)[1:11]
```

plt.plot(kind='line', linewidth=3, marker='o')  
plt.title('CO<sub>2</sub> Emission from 1750 to 1800', color='r', fontweight='bold', fontsize=20)  
plt.xlabel('Countries', fontsize=15, fontweight='bold', color='k')  
plt.ylabel('CO<sub>2</sub> Emission (tonnes)', fontsize=15, fontweight='bold', color='k')  
plt.legend(prop={'size':20}, loc=1)  
plt.show()



## 2. From year 1800 to 1850

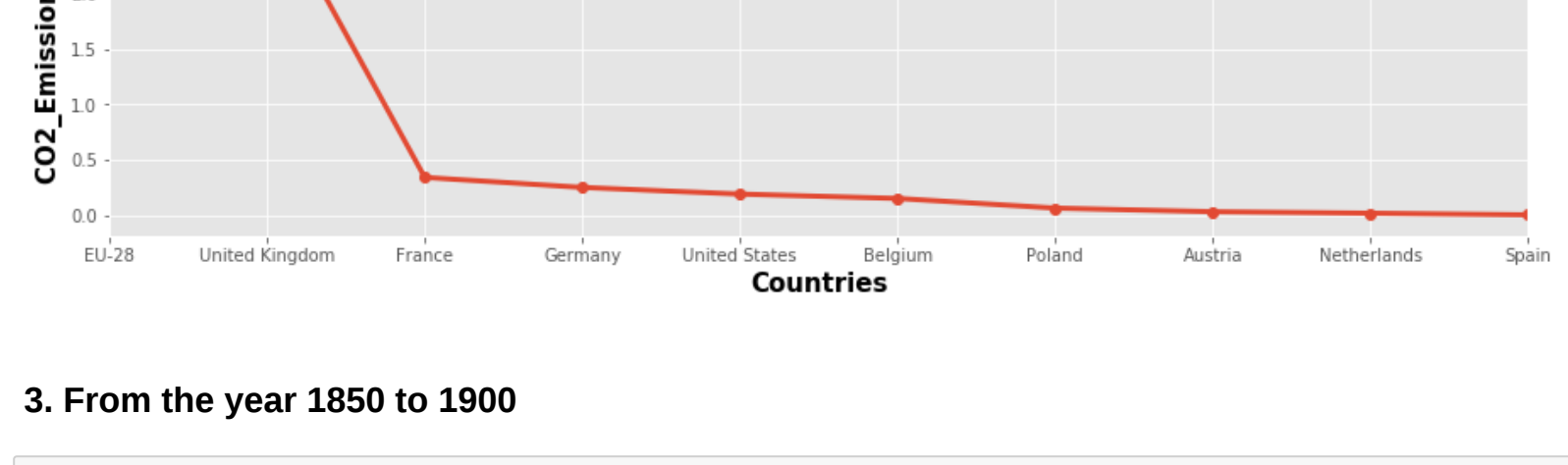
```
In [36]: p2.groupby('Entity')[['Annual CO2 emissions (tonnes)']].sum().sort_values(ascending=False)[1:11]
```

Out[36]:

Entity	3.784469e+09
EU-28	2.927456e+09
United Kingdom	3.404222e+08
France	2.497398e+08
Germany	1.897734e+08
United States	1.504951e+08
Belgium	6.295436e+07
Poland	3.064936e+07
Austria	1.655448e+07
Netherlands	2.748072e+06
Name: Annual CO <sub>2</sub> emissions (tonnes), dtype: float64	

```
In [37]: plt.style.use('ggplot')
plt.figure(figsize=(15,5))
p2.groupby('Entity')[['Annual CO2 emissions (tonnes)']].sum().sort_values(ascending=False)[1:11]
```

plt.plot(kind='line', marker='o', linewidth=3)  
plt.title('CO<sub>2</sub> Emission from 1800 to 1850', color='k', fontweight='bold', fontsize=20)  
plt.xlabel('Countries', fontsize=15, fontweight='bold', color='k')  
plt.ylabel('CO<sub>2</sub> Emission (tonnes)', fontsize=15, fontweight='bold', color='k')  
plt.legend(prop={'size':20}, loc=1)  
plt.show()



## 3. From the year 1850 to 1900

```
In [38]: p3.groupby('Entity')[['Annual CO2 emissions (tonnes)']].sum().sort_values(ascending=False)[1:11]
```

Out[38]:

Entity	2.701552e+10
EU-28	2.394048e+10
United States	1.047446e+10
Germany	6.136036e+09
France	3.370841e+09
Europe (other)	1.776498e+09
Belgium	1.387546e+09
Poland	8.433759e+08
Czechoslovakia	5.288383e+08
Austria	5.288383e+08
Name: Annual CO <sub>2</sub> emissions (tonnes), dtype: float64	