

Enter your Name in the Next Cell

Derek Cheung (RUID 204005689)

Grading Rubric

Score: / Max(0, 20 - Total Deductions)

Content Area	Deduction	Times Deducted	Check	Comments
Abstract				
Missing	5		[ ]	
Insufficient/Wrong Focus	1		[ ]	
Data Dictionary (Metadata)				
Missing	5		[ ]	
Insufficient/Wrong Form or Wording	1		[ ]	
Graphs				
Missing	5		[ ]	
Missing Title	1		[ ]	
Missing/Wrong Labels	1		[ ]	
Pre-Lab				
Missing	5		[ ]	
Insufficient/Wrong Answer	2		[ ]	
No/Incorrect/Insufficient Model Specification	2		[ ]	
No/Incorrect Statistical Hypothesis Statement	2		[ ]	
Post-Lab				
Missing	5		[ ]	
Insufficient/Wrong Answer	2		[ ]	
Correlations				
Missing	5		[ ]	
Insufficient/Wrong Analysis	2		[ ]	
Missing Graph	2		[ ]	
Estimations				
Missing	5		[ ]	
No or incorrect discussion/interpretation of...				
Hypothesis tests and p-values	2		[ ]	
R <sup>2</sup>	2		[ ]	
F-Statistic	2		[ ]	
Multicollinearity/VIF	2		[ ]	
Heteroskedasticity/Test	2		[ ]	
Autocorrelation/Test	2		[ ]	
No/insufficient model selection	2		[ ]	
Elasticities				
Missing	5		[ ]	
Incorrect Interpretation	2		[ ]	
Missing Summary Table	2		[ ]	
Model Portfolio				
Missing	5		[ ]	
General Comments:				

Contents

- 1. [Collaboration Policy](#)
- 2. [Introduction](#)
  - A. [Purpos](#)
  - B. [Assignment](#)
  - C. [Problem Solution](#)
- 3. [Documentation](#)
  - A. [Abstract](#)
  - B. [Data Dictionary](#)
  - 4. [Pre-lab](#)
  - 5. [Tasks](#)
  - 6. [Post-lab Analysis and Conclusions](#)

Collaboration Policy

[Back to Contents](#)

The submitted assignment must be your work.

There is to be no collaboration for this assignment.

Introduction

[Back to Contents](#)

Purpose

[Back to Contents](#)

The purpose of this lab is to allow you to:

- identify and define your own economic problem;
- collect your own data;
- analyze your data using the tools you learned;
- estimate a multiple linear regression model in Pandas;
- interpret key statistics;
- identify shortcomings in the proposed linear model;
- summarize the regression output;
- estimate elasticities and judge their reasonableness;
- build a model portfolio;
- interpret the model results.

Assignment

[Back to Contents](#)

You have to define your own problem for this lab. There are two rules:

- There must be a minimum of five (5) legitimate independent variables. Dummies for a concept variable are not legitimate but the concept variable is legitimate.
- You cannot repeat a problem discussed in class, used in a lesson, a tutorial, or was in a previous lab.

Problem Statement

[Back to Contents](#)

State your problem, why it's a problem, and what you expect to show.

The Chinese Economy Party (CCP) has announced it's economic stimulus plan in response to the Covid-19 pandemic. The CCP has put in place plans to increase government spending in order to build new and improve existing infrastructure. This plan will stimulate the short term economy by creating jobs while being able to benefit in the long term, from a more efficient infrastructure after the pandemic is over. These projects have caused the demand for raw materials (example being construction metals) to increase drastically and has led the prices for scrap metal to record levels. Global shipping companies, being some of the most devastated companies from the Covid-19 pandemic, are looking to cash in on these scrap metal prices by selling some of their older and less fuel efficient industrial cargo ships (ships that would have otherwise continued to sail for much longer periods of time) for a short term injection of cash to hopefully be able to ride out the remainder of the economic slowdown we have experienced since March of 2020. This may cause a problem for the future global economy, as we have become so reliant on the merchant marine fleet, "the unsung heroes of globalization". This decrease in shipping capacity may indicate a slow down in globalization as countries look to become more self sufficient, moving forward after the pandemic. This would then have huge impacts on the global economy as the productivity and efficiency offered by specialization will definitely be lessened.

In this lab we look to explore this relationship between global gross domestic product and the total deadweight tons of the global merchant fleet (total shipping capacity) and try to capture the affect total shipping capacity will have on the global economy. To do this we will estimate a complex regression model that will also take into account world population, global unemployment, global energy consumption, and world literacy rate.

The idea for this lab is from Economics Explained's most recent Youtube video uploaded on 12/03/2020. Here is a link to the video <https://www.youtube.com/watch?v=spg-2uQp-37w>.

Documentation

[Back to Contents](#)

Appropriate documentation.

Abstract

[Back to Contents](#)

In this lab we explore the relationship world population, total deadweight tons of the global merchant fleet (total shipping capacity), total global energy consumption, total global unemployment rate, and world literacy rates has on global gross domestic product. Specifically, we try to create an OLS model that tries to predict global gross domestic product based on world population, total deadweight tons of the global merchant fleet (total shipping capacity), total global energy consumption, total global unemployment rate, and world literacy rates. We ultimately conclude that the fit is statistically significant with total global energy consumption having the highest affect on global gross domestic product, and this is further reinforced by the consistent statistical significance of total global energy consumption across all models developed. This leads us to conclude that total global energy consumption is definitely related to global gross domestic product in some degree.

Data Dictionary

[Back to Contents](#)

A scatter plot with 'Total Deadweight Tons of Global Merchant Fleet (in hundred of millions tons)' on the x-axis and 'Global Energy Consumption' on the y-axis. The x-axis ranges from 6 to 20, and the y-axis ranges from 30 to 40. There are approximately 15 data points showing a positive correlation.

Total Deadweight Tons (hundred of millions tons)	Global Energy Consumption
6.5	28.5
6.8	29.5
7.2	30.5
7.5	31.5
7.8	32.5
8.2	33.5
8.5	34.5
8.8	35.5
9.2	36.5
9.5	37.5
9.8	38.5
10.2	39.5
10.5	40.5
10.8	41.5
11.2	42.5
11.5	43.5

#correlation 3: Energy vs GDP

```
ax = sns.relplot(x = "Energy", y = "GDP", data = df)
ax.set(title = "Global Energy Consumption vs Global Gross Domestic Product",
        xlabel = "Global Energy Consumption",
        ylabel = "Global Gross Domestic Product");
```

A scatter plot titled 'Global Energy Consumption vs Global Gross Domestic Product'. The x-axis is 'Global Energy Consumption' (ranging from 350 to 650) and the y-axis is 'Global Gross Domestic Product' (ranging from 30 to 90). The plot shows a strong positive correlation between the two variables.

Global Energy Consumption	Global Gross Domestic Product
350	25
360	28
370	30
380	32
390	34
400	36
410	38
420	40
430	42
440	44
450	46
460	48
470	50
480	52
490	54
500	56
510	58
520	60
530	62
540	64
550	66
560	68
570	70
580	72
590	74
600	76
610	78
620	80
630	82
640	84
650	86

Note: The GGDP and Energy values for 2019 are projected values.

Pre-lab

[Back to Contents](#)

Data description, testable hypotheses, statistical hypotheses.

Type of Data

This data is secondary, quantitative, macro time series data collected from a variety of sources. It is quantitative and macro because it contains only numerical data, aggregated over the entire world. Additionally, the data is collected starting from 1991 till 2019.

Testable Hypothesis

A potential testable hypothesis is:

We expect the global gross domestic product to increase as total shipping capacity, world population, energy consumption, and world literacy rises, and when global unemployment rates decline. This is because, intuitively, high shipping capacity, energy consumption, and world literacy rates combined with a low unemployment rate should indicate a higher production output. Additionally as the global population rises, we would require more output in order to sustain growth and we have more labor resources to apply towards productive output.

Tentative Specific Model

A potential specific model takes the form of:

$$GGDP = \beta_0 + (\beta_1 \times PopulationScaled) + (\beta_2 \times DWTONsScaled) + (\beta_3 \times Energy) + (\beta_4 \times Unemployment) + (\beta_5 \times Literacy) + noise$$

where the parameters  $\beta_0$  represents the intercept,  $\beta_1$  represents the affect global population has on global gross domestic product,  $\beta_2$  represents the affect total shipping capacity has on global gross domestic product,  $\beta_3$  represents the affect total energy consumption has on global gross domestic product,  $\beta_4$  represents the affect global unemployment has on global gross domestic product,  $\beta_5$  represents the affect world literacy rates has on global gross domestic product, and noise represents the random variations or error associated with statistical models.

Statistical Hypothesis

The statistical hypothesis corresponding to the testable hypothesis in regards to the parameters of the model takes the form of:

- $H_0: \beta_1 = 0$  Ha:  $\beta_1 > 0$
- $H_0: \beta_2 = 0$  Ha:  $\beta_2 > 0$
- $H_0: \beta_3 = 0$  Ha:  $\beta_3 > 0$
- $H_0: \beta_4 = 0$  Ha:  $\beta_4 < 0$
- $H_0: \beta_5 = 0$  Ha:  $\beta_5 > 0$

where the alternative of  $\beta_1$ ,  $\beta_2$ ,  $\beta_3$ ,  $\beta_4$ , and  $\beta_5$  states that global population, total shipping capacity, total energy consumption, and world literacy rate has a positive relationship with global gross domestic product, and the alternative of  $\beta_4$  states that global unemployment has a negative relationship with global gross domestic product. This follows from our testable hypothesis where we concluded, from intuition, that as the world's population increases we must produce more to sustain the population which leads to increases in productivity and more energy consumption and more shipping capacity. Therefore any increases in unemployment should signal productive slowdowns which would then lead to lower global gross domestic product.

The statistical hypothesis corresponding to the testable hypothesis in regards to the comparison of the naive and sophisticated model takes the form of:

- $H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0$
- Ha: at least one parameter ( $\beta_1$ ,  $\beta_2$ ,  $\beta_3$ ,  $\beta_4$ ,  $\beta_5$ )  $\neq 0$

where the null hypothesis says its naive (or restricted) model is better, and the alternative hypothesis says the sophisticated (or unrestricted) model is better.

Tasks

[Back to Contents](#)

Data import, data examination, model estimation, elasticity calculations, portfolio construction. Be sure to include all you learned this semester.

Data Import

```
In [1]: #import libs
import numpy as np
import pandas as pd
import statsmodels.api as sm
import statsmodels.formula.api as smf
from statsmodels.iolib.summary2 import summary_col
from statsmodels.stats.api import anova_lm
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: #import data and set row index to year
df = pd.read_excel('lab7Data.xlsx', index_col = 'Year')
df.head()
```

Out [2]:

	GGDP	Population	DWTONs	Energy	Unemployment	Literacy	PopulationScaled	DWTONsScaled
Year								
1991	87.698	771348100	1.988305e+06	620.000	5.395	86.30	7.713468	19.883051
1992	86.357	763109100	1.937777e+06	598.006	5.392	86.30	7.631091	19.377770
1993	81.229	754789525	1.868174e+06	583.929	5.570	86.13	7.547859	18.681738
1994	76.336	746402249	1.811526e+06	575.594	5.670	85.89	7.464022	18.115264
1995	75.199	737979719	1.753092e+06	568.778	5.638	85.60	7.379797	17.530919

```
In [3]: #drop redundant columns
df = df[['GGDP', 'PopulationScaled', 'DWTONsScaled', 'Energy', 'Unemployment', 'Literacy']]
df.head()
```

Out [3]:

	GGDP	PopulationScaled	DWTONsScaled	Energy	Unemployment	Literacy
Year						
1991	87.698	7.713468	19.883051	620.000	5.395	86.30
1992	86.357	7.631091	19.377770	598.006	5.392	86.30
1993	81.229	7.547859	18.681738	583.929	5.570	86.13
1994	76.336	7.464022	18.115264	575.594	5.670	85.89
1995	75.199	7.379797	17.530919	568.778	5.638	85.60

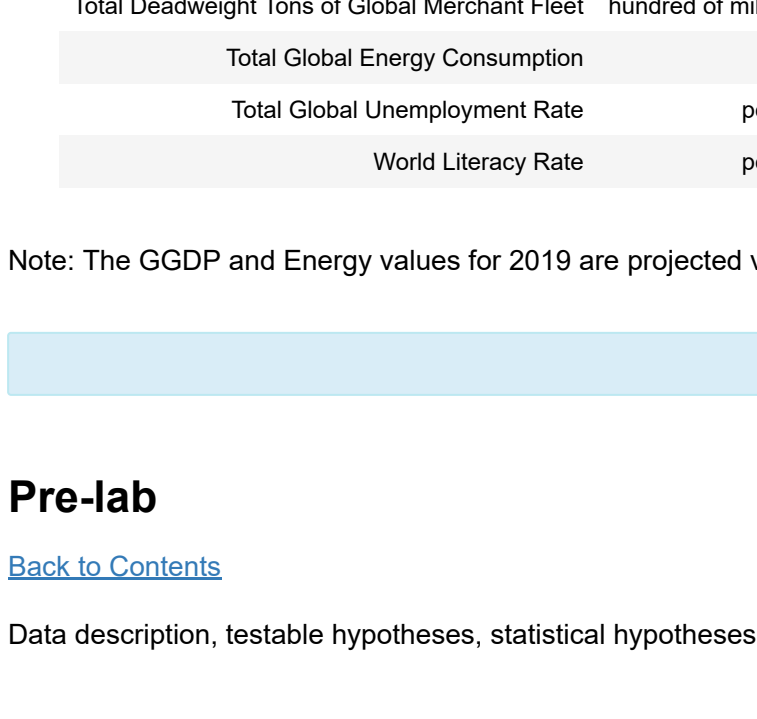
```
In [4]: #describe stats
df.describe().T
```

Out [4]:

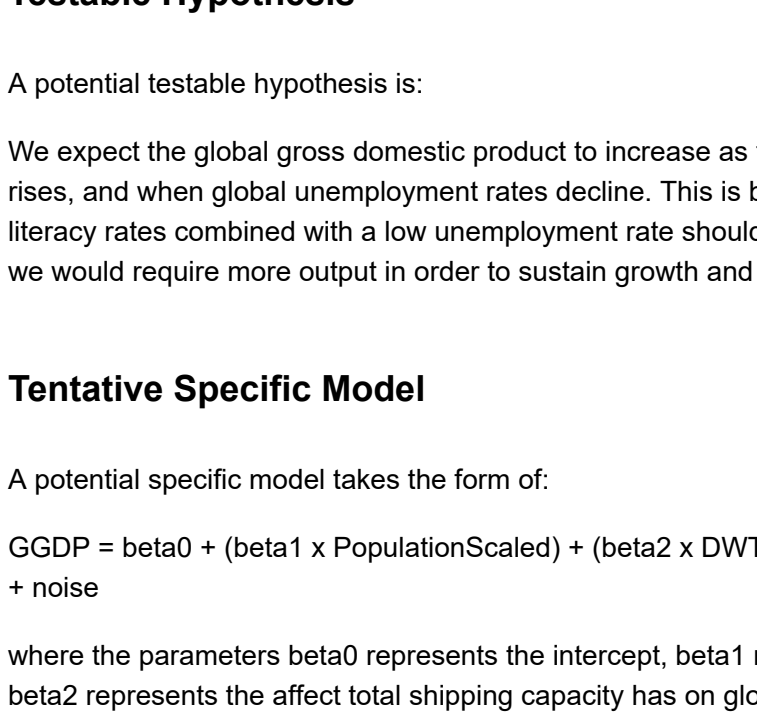
	count	mean	std	min	25%	50%	75%	max
GGDP	29.0	51.86228	21.842192	23.967000	31.573000	47.517000	75.146000	87.698000
PopulationScaled	29.0	6.564847	0.696986	5.414289	5.984794	6.541907	7.125828	7.713468
DWTONsScaled	29.0	11.290307	4.520854	6.512820	7.728012	9.074743	15.374839	19.883051
Energy	29.0	466.414103	87.538966	351.696000	382.070000	465.007000	549.766000	620.000000
Unemployment	29.0	5.645241	0.330158	4.758000	5.430000	5.670000	5.854000	6.189000
Literacy	29.0	81.777366	3.609607	74.910000	80.210000	82.410000	84.730000	86.300000

Data Exploration

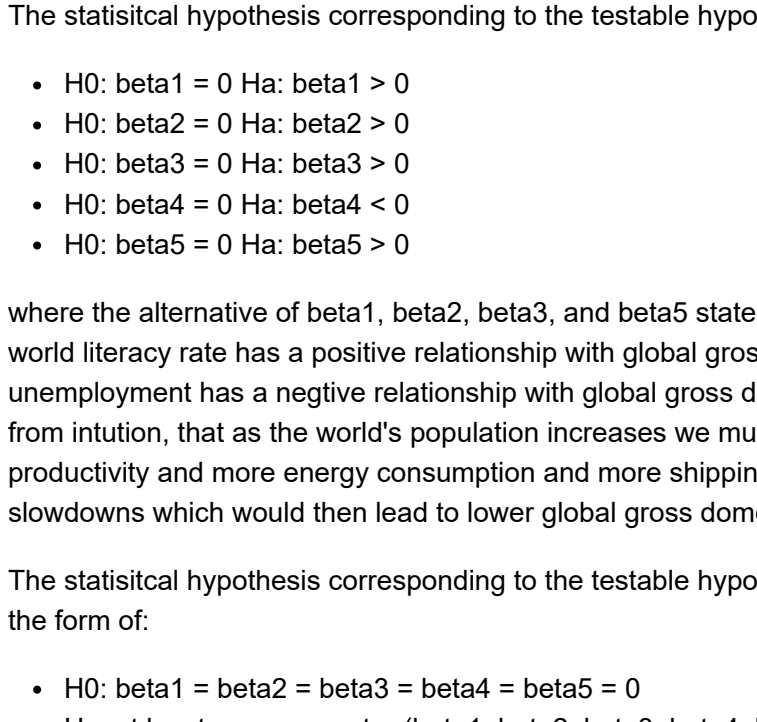
```
In [5]: #graph 1: GGDP over time
ax = sns.lineplot(y = 'GGDP', x = df.index, data = df);
ax.set(title = 'Global Gross Domestic Product over Time (1991-2019)',
        xlabel = 'Year',
        ylabel = 'Global Gross Domestic Product');
```



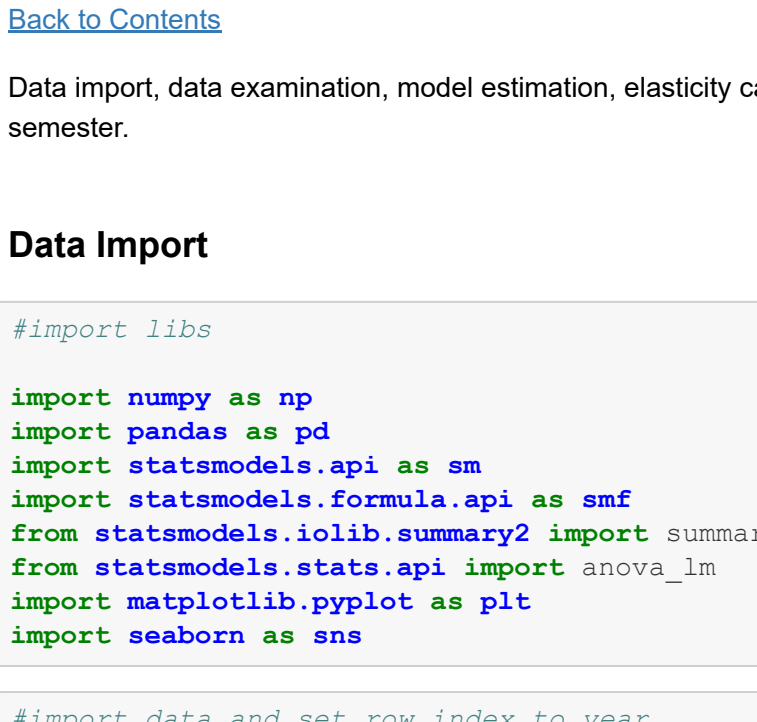
```
In [6]: #graph 2: Population over time
ax = sns.lineplot(y = 'PopulationScaled', x = df.index, data = df);
ax.set(title = 'Global Population over Time (1991-2019)',
        xlabel = 'Year',
        ylabel = 'Global Population (in billions)');
```



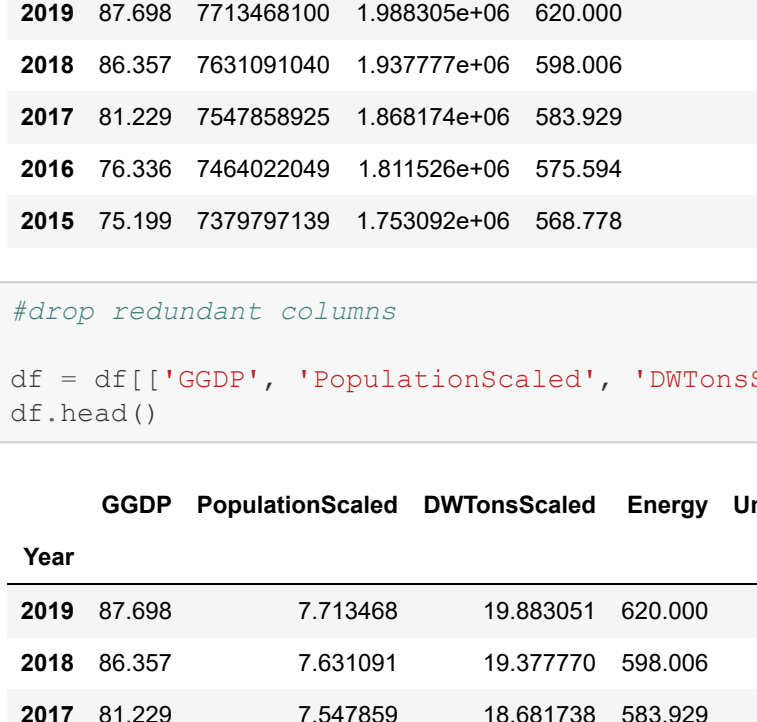
```
In [7]: #graph 3: DWTONsScaled over time
ax = sns.lineplot(y = 'DWTONsScaled', x = df.index, data = df);
ax.set(title = 'Total Deadweight Tons of Global Merchant Fleet (total shipping capctiy)\n over Time (1991-2019)',
        xlabel = 'Year',
        ylabel = 'Total Deadweight Tons of Global Merchant Fleet\n (in hundred of millions tons)');
```



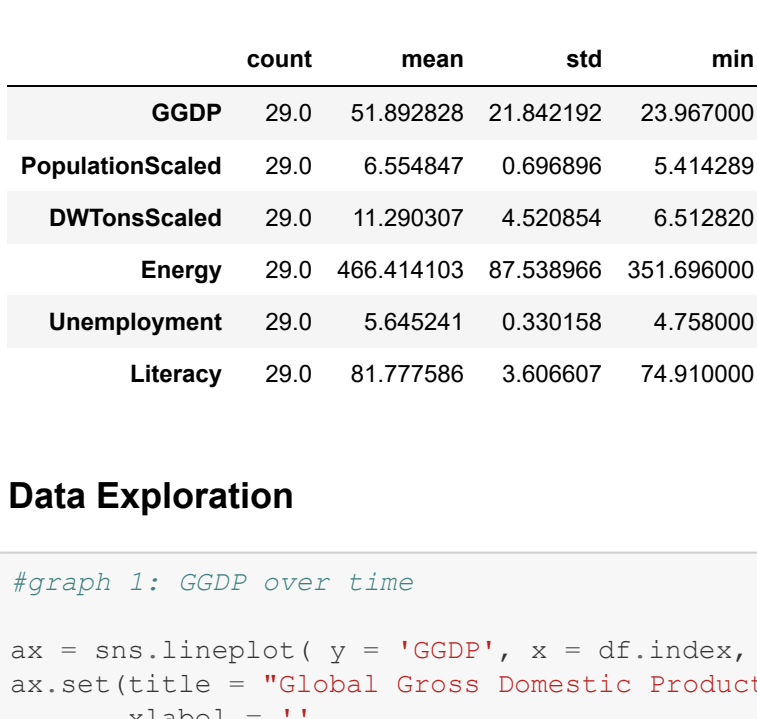
```
In [8]: #graph 4: Energy over time
ax = sns.lineplot(y = 'Energy', x = df.index, data = df);
ax.set(title = 'Global Energy Consumption over Time (1991-2019)',
        xlabel = 'Year',
        ylabel = 'Global Energy Consumption');
```



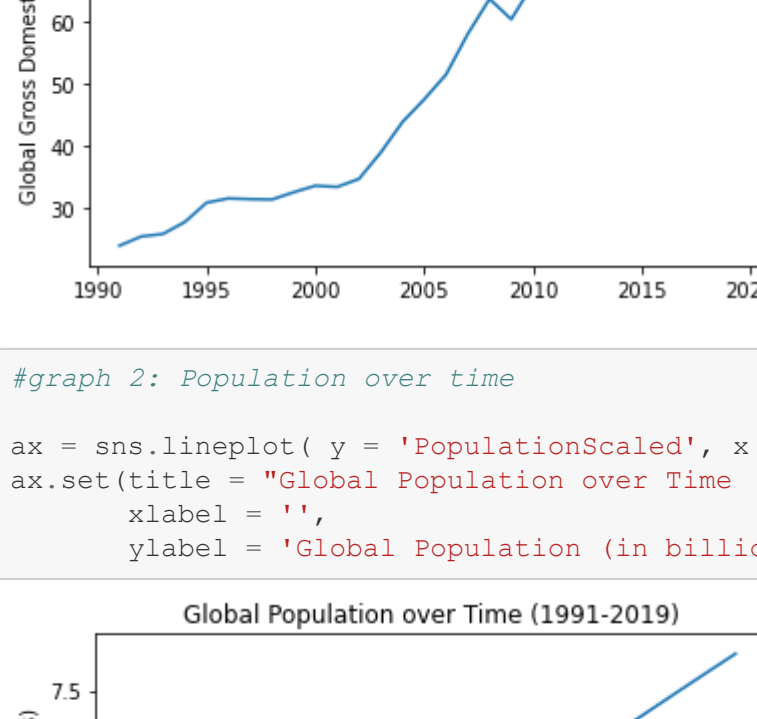
```
In [9]: #graph 5: Unemployment over time
ax = sns.lineplot(y = 'Unemployment', x = df.index, data = df);
ax.set(title = 'Global Unemployment over Time (1991-2019)',
        xlabel = 'Year',
        ylabel = 'Global Unemployment Rate');
```



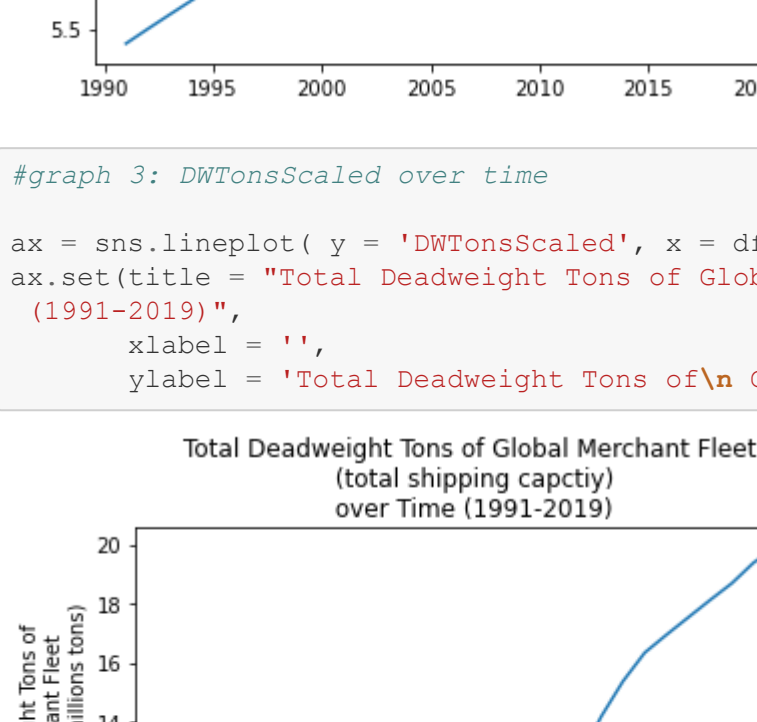
```
In [10]: #graph 6: Literacy over time
ax = sns.lineplot(y = 'Literacy', x = df.index, data = df);
ax.set(title = 'Global Literacy Rate over Time (1991-2019)',
        xlabel = 'Year',
        ylabel = 'Global Literacy Rate');
```



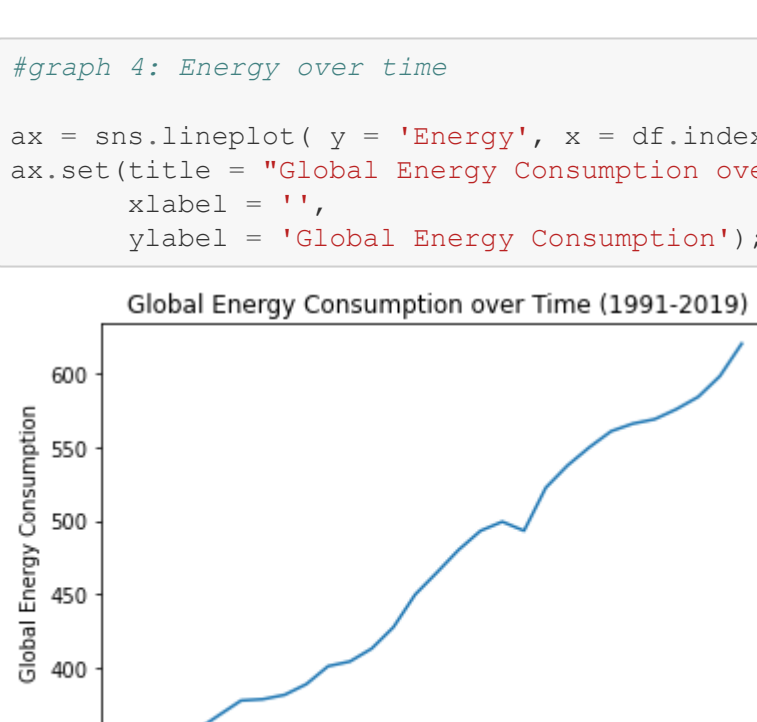
```
In [11]: #correlation 1: Population vs GGDP
ax = sns.relplot(x = 'PopulationScaled', y = 'GGDP', data = df)
ax.set(title = 'Global Population vs Global Gross Domestic Product',
        xlabel = 'Global Population (in billions)',
        ylabel = 'Global Gross Domestic Product');
```



```
In [12]: #correlation 2: DWTONsScaled vs GGDP
ax = sns.relplot(x = 'DWTONsScaled', y = 'GGDP', data = df)
ax.set(title = 'Total Deadweight Tons of Global Merchant Fleet vs Global Gross Domestic Product',
        xlabel = 'Total Deadweight Tons of Global Merchant Fleet\n (in hundred of millions tons)',
        ylabel = 'Global Gross Domestic Product');
```



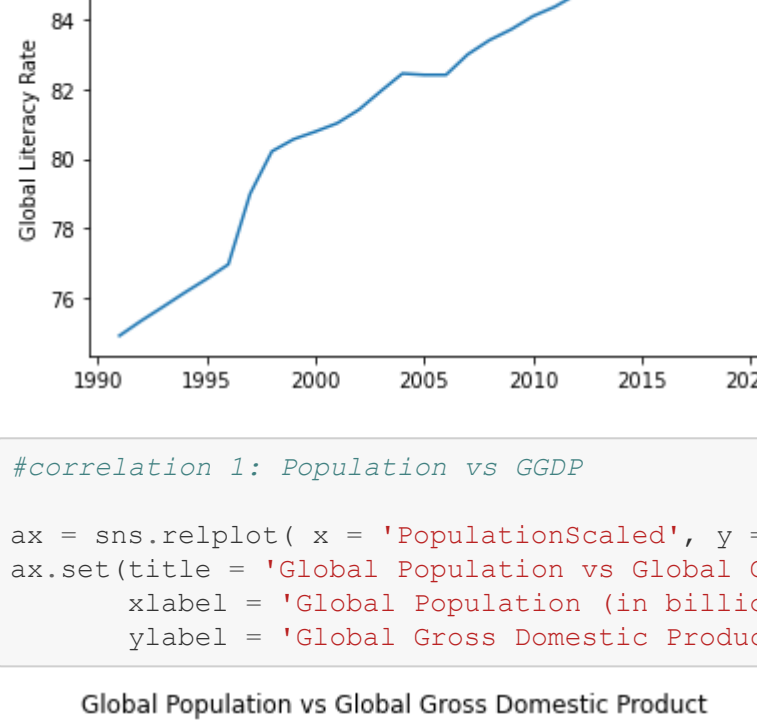
```
In [13]: #correlation 3: Energy vs GGDP
ax = sns.relplot(x = 'Energy', y = 'GGDP', data = df)
ax.set(title = 'Global Energy Consumption vs Global Gross Domestic Product',
        xlabel = 'Global Energy Consumption',
        ylabel = 'Global Gross Domestic Product');
```



```
In [14]: #correlation 4: Unemployment vs GGDP
ax = sns.relplot(x = 'Unemployment', y = 'GGDP', data = df)
ax.set(title = 'Global Unemployment Rate vs Global Gross Domestic Product',
        xlabel = 'Global Unemployment Rate',
        ylabel = 'Global Gross Domestic Product');
```



```
In [15]: #correlation 5: Literacy vs GGDP
ax = sns.relplot(x = 'Literacy', y = 'GGDP', data = df)
ax.set(title = 'Global Literacy Rate vs Global Gross Domestic Product',
        xlabel = 'Global Literacy Rate',
        ylabel = 'Global Gross Domestic Product');
```



Model(s) Estimation

In [18]: #regression 1; all variables used

```
mod1 = smf.ols(formula = "GGDP ~ PopulationScaled + DWTONsScaled + Energy + Unemployment + Literacy", d
reg1 = mod1.fit()
print(reg1.summary())
```

```
#correlation 3: Energy vs GDP
ax = sns.relplot(x = 'Energy', y = 'GDP', data = df)
ax.set(title = 'Global Energy Consumption vs Global Gross Domestic Product',
        xlabel = 'Global Energy Consumption',
        ylabel = 'Global Gross Domestic Product');
```

```
#correlation 4: Unemployment vs GDP
ax = sns.relplot(x = 'Unemployment', y = 'GDP', data = df)
ax.set(title = 'Global Unemployment Rate vs Global Gross Domestic Product',
        xlabel = 'Global Unemployment Rate',
        ylabel = 'Global Gross Domestic Product');
```

Notes:  
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.  
[2] The condition number is large, 3.08e+04, this might indicate that there are strong multicollinearity or other numerical problems.

In [19]: #regression 2: DWTONsScaled only

```
mod2 = smf.ols(formula = "GGDP ~ DWTONsScaled", data = df)
reg2 = mod2.fit()
print(reg2.summary())
```

```

#correlation matrix
df.corr()

```

	GGDP	PopulationScaled	DWTONsScaled	Energy	Unemployment	Literacy
GGDP	1.000000	0.976163	0.967675	0.992881	0.089585	0.913640
PopulationScaled	0.976163	1.000000	0.951171	0.990096	0.236587	0.970719
DWTONsScaled	0.967675	0.951171	1.000000	0.963606	0.027839	0.863517
Energy	0.992881	0.990096	0.963606	1.000000	0.135322	0.940330
Unemployment	0.089585	0.236587	0.027839	0.135322	1.000000	0.401828
Literacy	0.913640	0.970719	0.863517	0.940330	0.401828	1.000000

```

# graph of correlation matrix
ax = sns.heatmap(df.corr(), annot = True, cmap = "coolwarm").set_title('Heatmap of the Correlation Matrix')

```

	GGDP	PopulationScaled	DWTONsScaled	Energy	Unemployment	Literacy
GGDP	1	0.98	0.97	0.99	0.09	0.91
PopulationScaled	0.98	1	0.95	0.99	0.24	0.97
DWTONsScaled	0.97	0.95	1	0.96	0.02	0.86
Energy	0.99	0.99	0.96	1	0.14	0.94
Unemployment	0.09	0.24	0.02	0.14	1	0.4
Literacy	0.91	0.97	0.86	0.94	0.4	1



