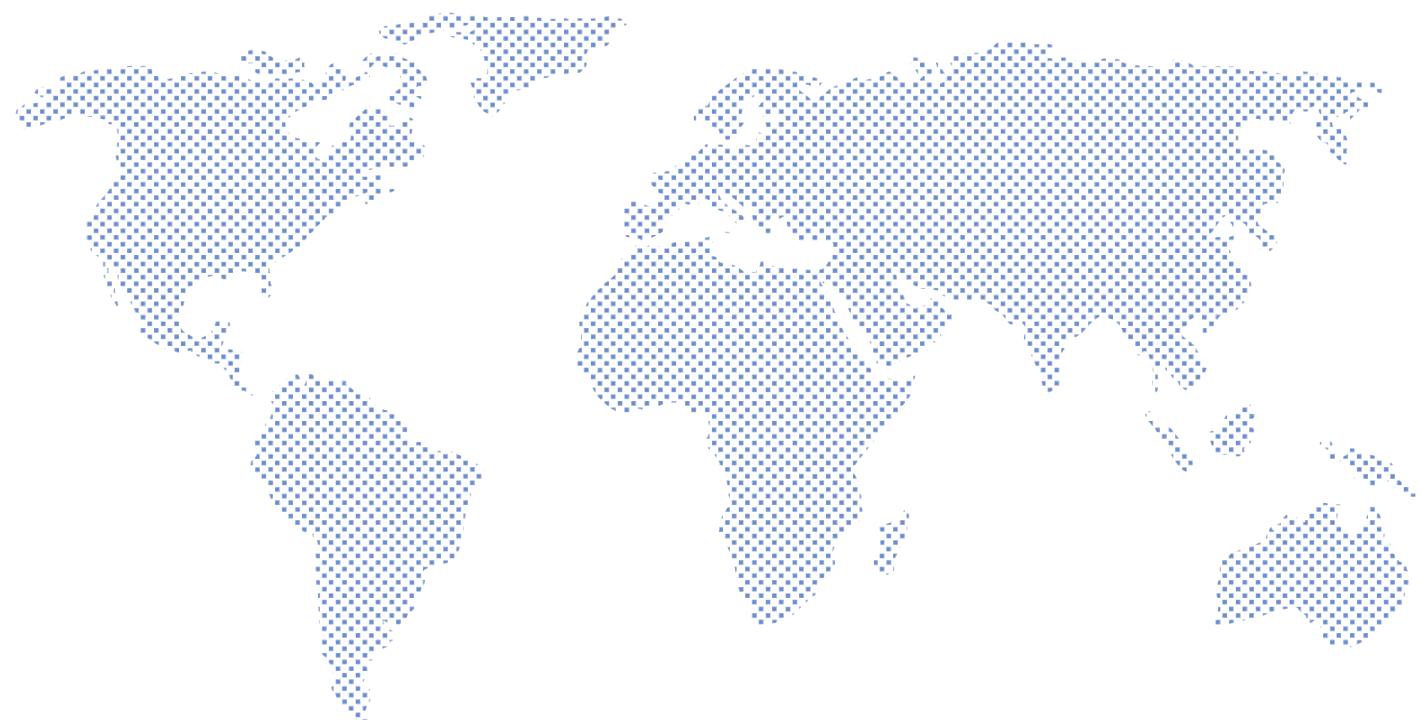


Suicide Rates Overview 1985 to 2016

Compares socio-economic info with suicide rates



Ruozhuo Wang Joyce Jian
Sherry Huang Sherry Huang

CONTENT

Overview	1
Data description	1
Data cleaning	2
Methods used	3
Suicide rates and years	4
Suicide rates and age group	6
Suicide rates and gender	8
Suicide rates and country/location	9
Suicide rates and GDP	11
Reference	14
Python codes	15

OVERVIEW

Our project is to uncover trends and patterns in suicide rates in 1985 to 2016. We will examine relationships between suicide rates and age, gender, climate, GDP; trends in suicide rates over the years; and related questions, as the data admits.

DATA DESCRIPTION

This compiled dataset pulled from four other datasets linked by time and place, and was built to find signals correlated to increased suicide rates among different cohorts globally, across the socio-economic spectrum

<https://www.kaggle.com/russellyates88/suicide-rates-overview-1985-to-2016>.

The datasets includes the following information that we can use to analyze factors associated with suicide rates:

Related factors	Years
	Age Group
	Gender
	Country/Location
	GDP

DATA CLEANING

Country	1985	1986	1987	1988	1989	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016
A	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
B			v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
C	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
D								v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	
E	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
F	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
G								v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	
H								v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	
I	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
J	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
K	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
L	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
M	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v		
N											v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	v	

To obtain reliable datasets for further analysis, we performed data cleaning process by doing the following:

- Dropped unnecessary columns in the DataFrame
- Renamed columns to a more recognizable set of labels
- Removed missing and duplicate data
- Set up several bins (groups) for analysis

Analyzing suicide rates VS Years/Age group/Gender/GDP

Since a lot of data are missing in 1985-1989, and 2015-2016 in the raw dataset, we decided to only use data in 1990-2014 from 38 countries for analysis.

Analyzing suicide rate VS Counties/Location

To have the most countries' continuous data for analysis, we decided to use the data in 2010-2014 from 71 countries for analysis.

METHODS USED

Pandas, Numpy:

Restructured raw data and performed scientific computing for analysis

Matplotlib:

Visualized analysis results and plotted findings using line chart, bar plot, scattered plot and Google Heatmap

API, JSON, Python Requests:

Retrieved coordinates from Google Map using API and imported into the DataFrame

Data Modeling:

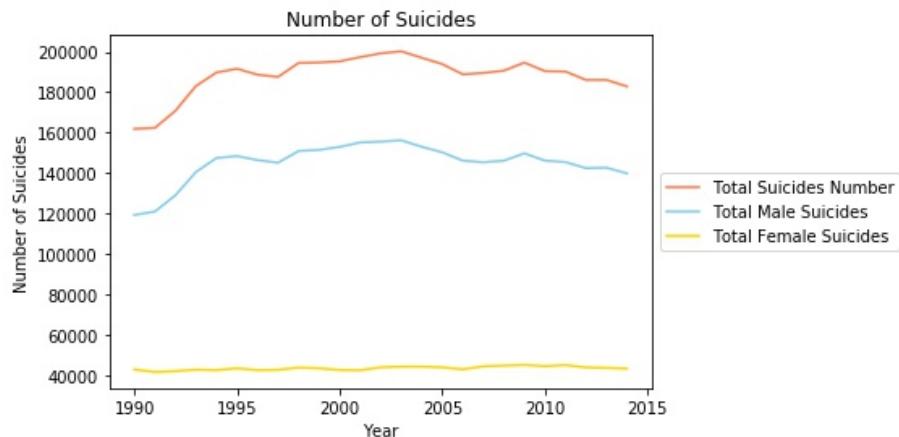
Analyzed output using line of best fit (linear/non-linear regression model) and statistical testing



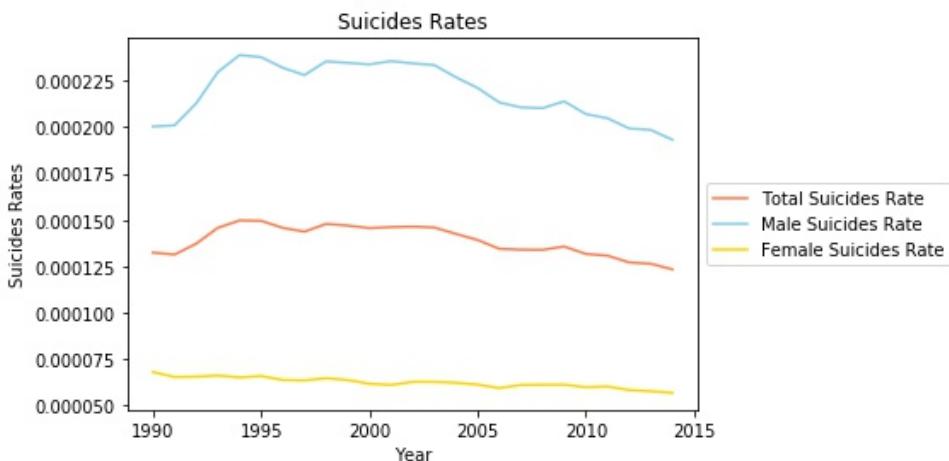
YEARS

Suicide rates and years

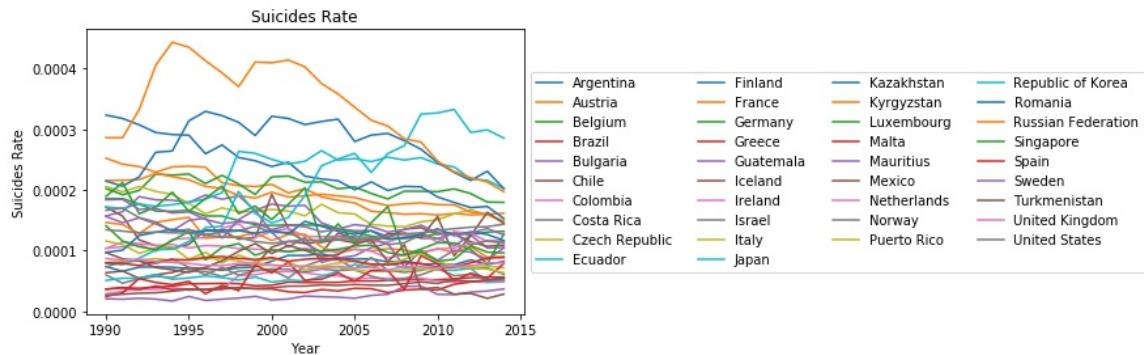
The total number of suicides and number of male suicides fluctuates over the years, but they were all increasing in the long term. However, number of female suicides were almost the same every year.



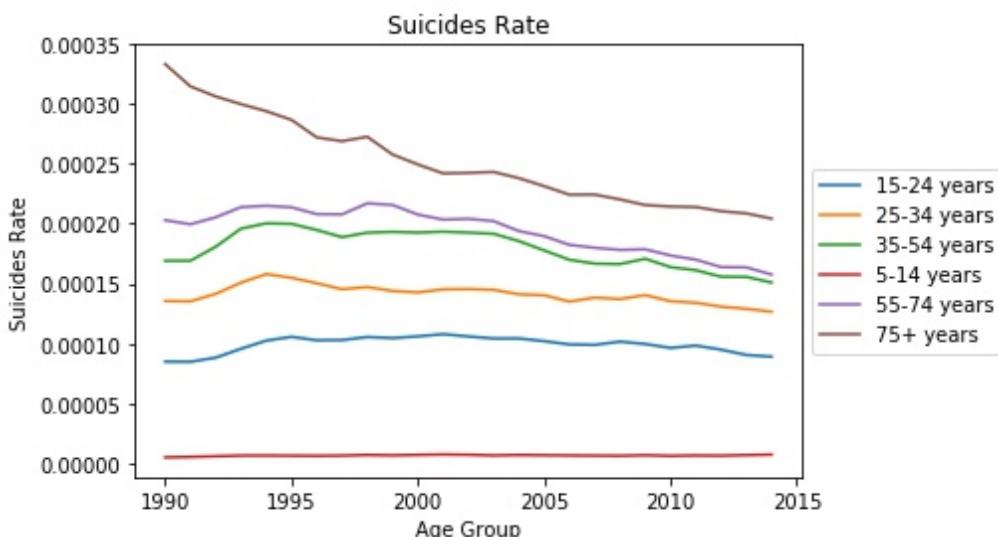
Although the total number of suicides were increasing, the total suicide rates, as well as both gender suicide rates were decreasing from 1990 to 2014. However, the suicide rates of female was relatively consistent over the years. The fluctuation of total suicide rates was mainly caused by suicide rates of male.



YEARS



Among 38 countries in the dataset, there was no clear trend of suicide rates across countries from 1990 to 2014.

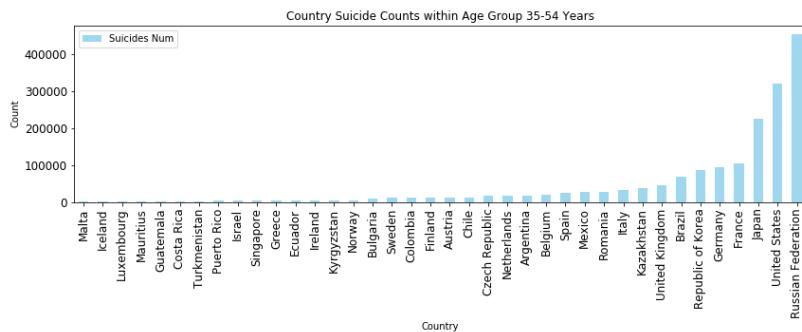
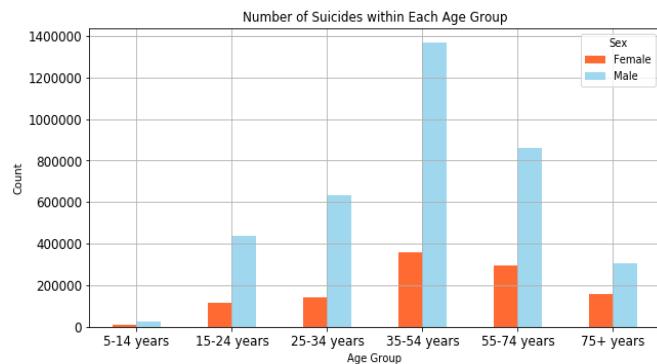


The suicides rate of all age groups were decreasing over 1990-2014. Age group of 75+ years old dropped the most, which might be resulted from improved healthcare and pension system.

AGE GROUP

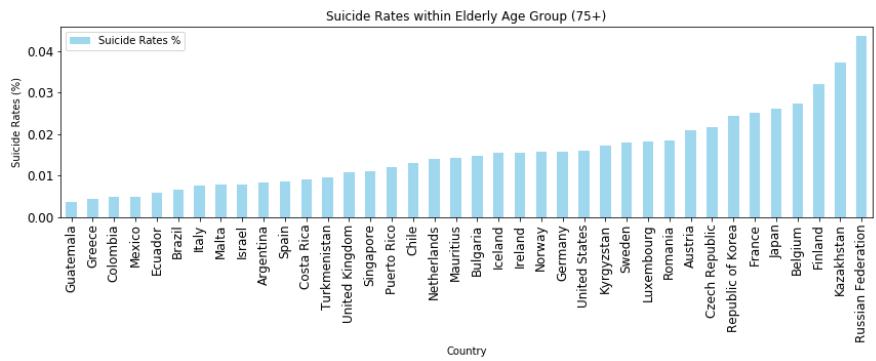
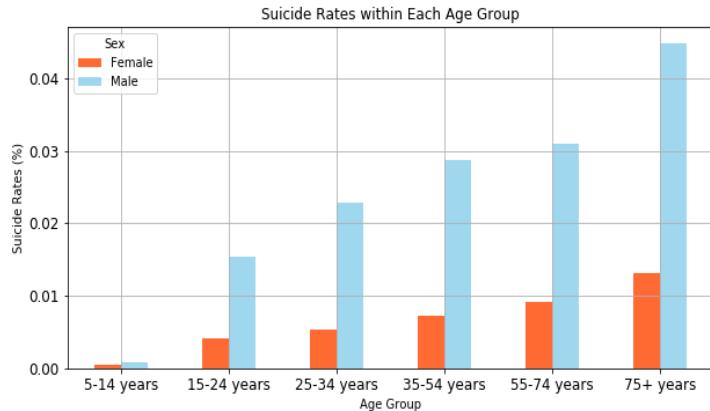
Suicide rates and age group

Among all the age groups, female have much less total number of suicides comparing to that of male, and age group of 35-54 years old has the highest number of suicides in both genders. After that, as age increases, the number of suicides decreases. This large number of suicides within middle age group is mainly contributed by the following countries: Japan, United States, and the Russian Federation.



Suicide rates reveal a different trend: the rates are the highest among elderly group (75+ years) for both genders. And the country with the highest elderly suicide rate is still Russia Federation.

AGE GROUP



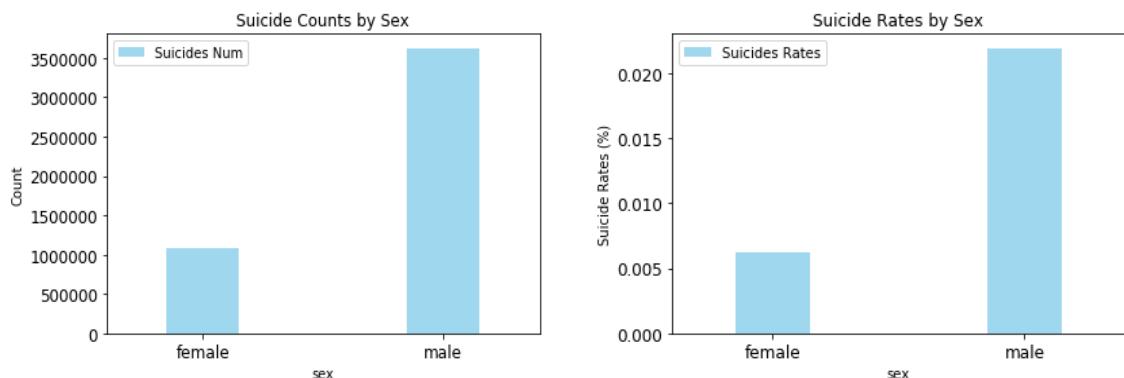
Conclusion

Either from suicide rates or suicide counts perspective, suicide is correlated with age.

GENDER

Suicide rates and gender

Over 1990-2014, the number of suicides for female in all countries is significantly less than that of male (1081921 vs 3625386, t -test statistics ≈ -18.77 , p -value $\approx 1.458e76 \approx 0$). This difference is also in suicide rates (0.006% vs 0.022% , t -test statistics ≈ -51.72 , p -value = 0).



The male to female suicides ratio is highest among 25-34 years old: for each female who committed suicide, there would be 4 male counterparts. Some possible risk factors for high suicides counts within middle age group might be: disrupted marital status, depression, addictive disorder, etc.

sex	female	male	male_to_female_ratio
age			
5-14 years	11519	24501	2
15-24 years	113466	435073	3
25-34 years	143549	631340	4
35-54 years	355754	1368261	3
55-74 years	297580	859581	2
75+ years	160053	306630	1

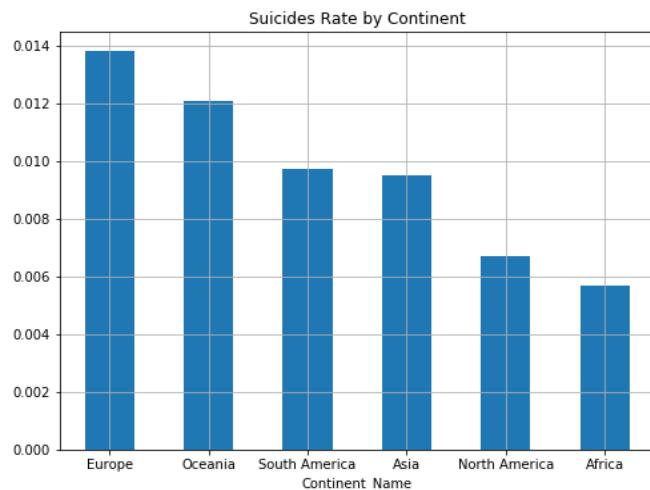
Conclusion

Either from suicide rates or suicide counts perspective, suicide is correlated with gender, with female having significantly less suicides counts comparing to that of male.

COUNTRY

Suicide rates and country

Continent



Europe suffers the highest rate of suicides (0.0138%) among the six continents, which is more than twice of that of Africa.



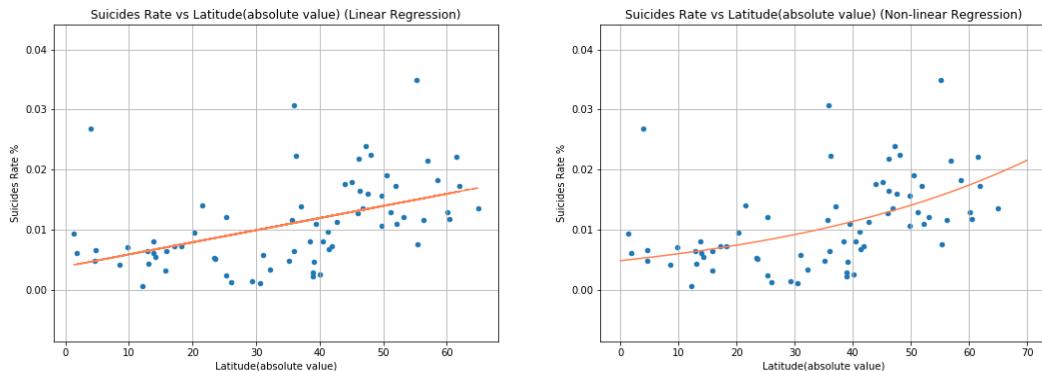
Google Heatmap by Suicide Rate (Europe)

COUNTRY

Latitude (absolute value):

The correlation coefficient between the absolute value of latitude and suicide rate, also known as R-value, is 0.4727, which means they have a moderate positive relationship. In other words, with countries located further away from equator suicide rate sees a remarkable increasing trend.

Comparison using linear and non-linear model



Both linear and non-linear regression analysis have been conducted to describe this trend. Based on the R-value, the exponential model better explains the data. Thus, it can be concluded that the increasing rate also grows with latitude increasing.

- **Linear model:** $\text{Suicide Rate} = 0.0002017 * \text{Latitude}(abs) + 0.003861$

$R\text{-square} = 0.2235$

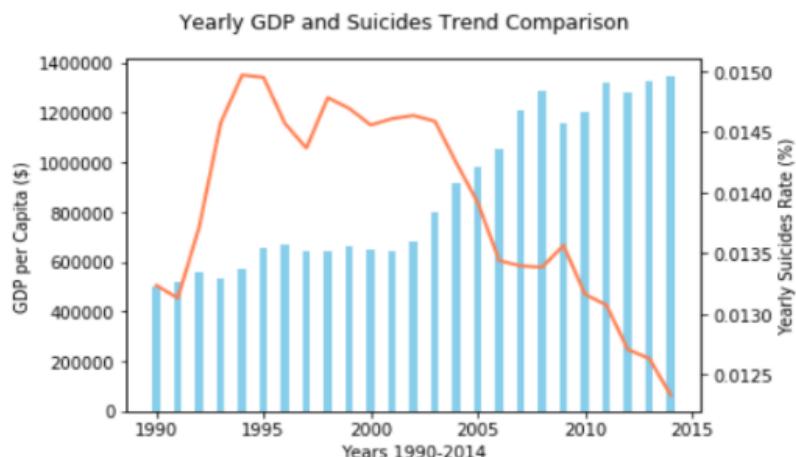
- **Non-linear model:** $\text{Suicide Rate} = 0.004839 * \exp^{(0.02134 * \text{Latitude}(abs))}$

$R\text{-square} = 0.2492$

GDP

Suicide and economic factor (GDP)

*Overall trending comparison of suicide and GDP
trending of 1990-2014*



In general the two series' value trend in an opposite direction over time

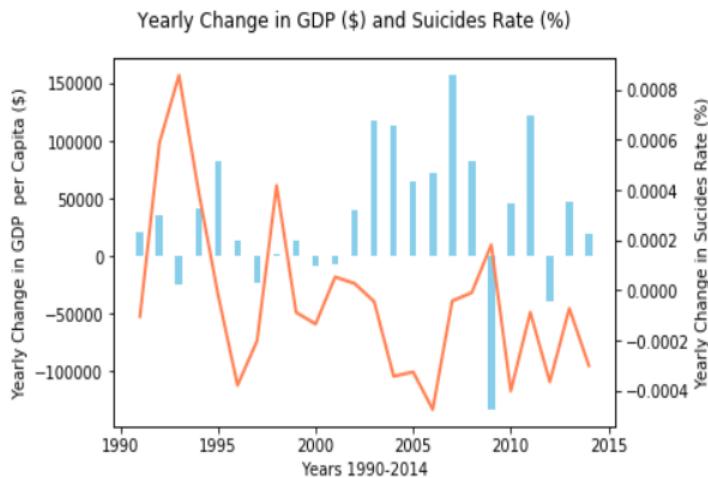
- GDP always increases
- Suicide rate augments first and then shows a continuous downward slope

However yearly data is not consistently matching each other's growth/decrease for the two variables.

Time series difference comparison

To have a better understanding of the evolution of suicide rate and GDP, instead of taking in consideration of each year's value of the two variables directly, we calculate the 1-period-shift difference ($D_1 = V_{\text{year2}} - V_{\text{year1}}$) of the two series and then compare the difference trending of the study year range.

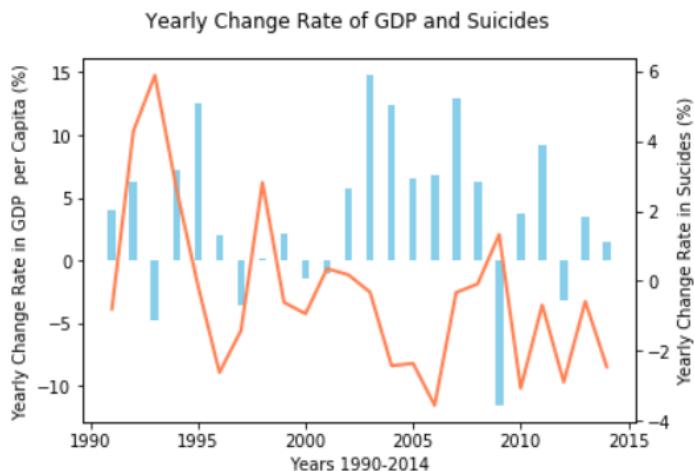
GDP



From the graph representation, the two variables scale very differently and evolution has no matching pattern during the studied period.

Time series difference in % change comparison

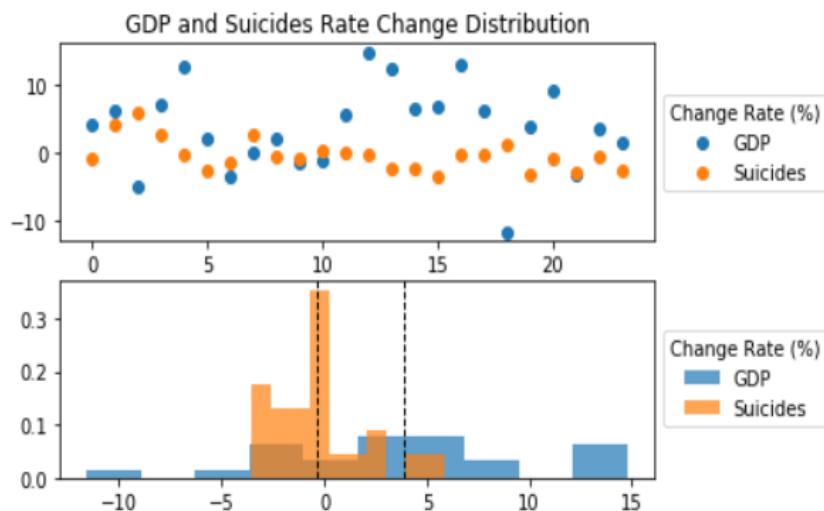
The % change also shows a similar graphic which indicates no significant relationship between the two variables being studied.



GDP

Evolution rate (%) Distribution

We plot the two series of % change rate year over year respectively on a scatterplot and also on a histogram plot.



Both distribution plots tell no similarity, however in order to make sure of the non-existent correlation of suicide rate and GDP figure, we pull out a T-test on the following hypotheses.

Null Hypothesis (H_0): GDP is affecting suicide rates

Alternative Hypothesis (H_1): GDP has no influence on suicide rates

T-test result gives a P value which equals 0.07841685 (>0.05)

Thus statistically saying, the suicides growth is not significantly related to GDP growth to the extent of a confidence level of 95%.

REFERENCE

United Nations Development Program. (2018). Human development index (HDI).
<http://hdr.undp.org/en/indicators/137506>

World Bank. (2018). World development indicators: GDP (current US\$) by country:1985 to 2016.

<http://databank.worldbank.org/data/source/world-development-indicators#>

[Szamil]. (2017). *Suicide in the Twenty-First Century* [dataset].
<https://www.kaggle.com/szamil/suicide-in-the-twenty-first-century/notebook>

World Health Organization. (2018). Suicide prevention.
http://www.who.int/mental_health/suicide-prevention/en/

PYTHON CODES

Data cleaning 1990 - 2014

```
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
```

```
import os
import csv
import pandas as pd
from config import gkey
import gmmaps
gmmaps.configure(api_key=gkey)
import json
import requests
import numpy as np

file = "Data/master.csv"
suicide_data = pd.read_csv(file, encoding = "ISO-8859-1")
suicide_data.head()

sub_year = suicide_data.loc[(suicide_data["year"]>=1990)&(suicide_data["year"]<=2014),:]
sub_year.head()

country_list = pd.DataFrame(sub_year["country"].value_counts())
country_list = country_list.loc[country_list["country"]== 300,:]
country_list.head()

names = country_list.index

country_chosen = ['Sweden', 'Finland', 'Greece', 'Luxembourg', 'Guatemala', 'Netherlands',
'Turkmenistan', 'Spain', 'Bulgaria', 'Ecuador', 'Brazil', 'Puerto Rico',
'Argentina', 'Israel', 'France', 'Germany', 'Russian Federation',
'Singapore', 'Republic of Korea', 'Iceland', 'Mexico', 'Czech Republic',
'Romania', 'Chile', 'Mauritius', 'Malta', 'United States', 'Italy',
'Kazakhstan', 'Belgium', 'Colombia', 'Costa Rica', 'Austria',
'United Kingdom', 'Ireland', 'Kyrgyzstan', 'Japan', 'Norway']

sub_year = sub_year.set_index("country")
sub_dataset = sub_year.loc[country_chosen,:]
sub_dataset.head()

sub_dataset.to_csv("Data/sub_dataset.csv", encoding = "utf-8", index = True, header = True)
```

PYTHON CODES

Data cleaning 2010 - 2014

```
1 # coding: utf-8
2 import os
3 import csv
4 import pandas as pd
5
6 #from config import gkey
7 file = "data/master2.csv"
8 suicide_data = pd.read_csv(file, encoding = "ISO-8859-1")
9 suicide_data.head()
10
11 group_by_year = suicide_data.groupby(["year"])
12 country_count = group_by_year["country"].nunique()
13 country_count
14
15 sub_year = suicide_data.loc[(suicide_data["year"]>=2010)&(suicide_data["year"]<=2014),:]
16 sub_year.head()
17
18 country_list = pd.DataFrame(sub_year["country"].value_counts())
19 country_list = country_list.loc[country_list["country"]== 60,:]
20 country_list.head()
21
22 names = country_list.index
23
24 country_chosen = ['Mexico', 'Costa Rica', 'Norway', 'Italy', 'France', 'Argentina',
25     'Turkmenistan', 'Cuba', 'Hungary', 'Sweden', 'Australia', 'Croatia',
26     'Mauritius', 'Chile', 'Armenia', 'Ireland', 'Grenada', 'Thailand',
27     'Ecuador', 'Saint Vincent and Grenadines', 'Cyprus', 'Luxembourg',
28     'Seychelles', 'Puerto Rico', 'Belize', 'Bulgaria', 'Brazil',
29     'Uzbekistan', 'Serbia', 'Bahrain', 'El Salvador', 'Latvia', 'Belgium',
30     'Czech Republic', 'Israel', 'Switzerland', 'Slovenia', 'Austria',
31     'Kyrgyzstan', 'United Kingdom', 'Kazakhstan', 'Turkey', 'Qatar',
32     'Colombia', 'South Africa', 'Panama', 'Portugal', 'Japan', 'Germany',
33     'Republic of Korea', 'Kuwait', 'Lithuania', 'Estonia', 'Iceland',
34     'United States', 'Spain', 'Finland', 'Denmark', 'Suriname', 'Poland',
35     'Guatemala', 'Georgia', 'Russian Federation', 'Netherlands', 'Romania',
36     'Greece', 'Malta', 'Paraguay', 'Nicaragua', 'Saint Lucia', 'Singapore']
37
38 sub_year = sub_year.set_index("country")
39 sub_dataset = sub_year.loc[country_chosen,:]
40 sub_dataset.head()
41
42 sub_dataset.to_csv("Data/sub_dataset_2010-2014.csv", encoding = "utf-8", index = True, header = True)
43
44 len(country_chosen)
45
```

PYTHON CODES

Data analysis

```
1 ## Import Dependencies
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import numpy as np
5 import scipy.stats as stats
6
7 from config import gkey
8 import gmmaps
9 gmmaps.configure(api_key=gkey)
10 import json
11 import requests
12
13 from scipy.stats import linregress
14 from scipy.optimize import curve_fit
15
16 # ## Import Cleaned Dataset
17 data = pd.read_csv("data/sub_dataset.csv")
18 data.head()
19
20 # ## Analyse Year Factor
21 # Group by Year
22 groupby_yr=data.groupby("year")
23 year_sum = groupby_yr["suicides_no"].sum()
24 year_sum.head()
25
26 # Total suicides number and suicides number by Gender Over 1990-2014
27
28 male = data.loc[(data["sex"] == "male")]
29 female = data.loc[(data["sex"] == "female")]
30
31 male_sum_yr = male.groupby(['year'])['suicides_no'].sum()
32 female_sum_yr = female.groupby(['year'])['suicides_no'].sum()
33
34 total_population=data.groupby(['year'])['population'].sum()
35 male_population=male.groupby(['year'])['population'].sum()
36 female_population=female.groupby(['year'])['population'].sum()
37
38 total_rate=year_sum/total_population
39 male_rate=male_sum_yr/male_population
40 female_rate=female_sum_yr/female_population
41
42 # Plot Number of Suicides by Gender
43
44 total_suicides_no = year_sum.plot(kind='line', color="coral", label="Total Suicides Number")
45 male_sum_yr_plot = male_sum_yr.plot(kind='line', color="skyblue", label="Total Male Suicides")
```

PYTHON CODES

```
46 female_sum_yr_plot = female_sum_yr.plot(kind='line', color="gold", label='Total Female Suicides')
47
48
49 total_suicides_no.set_xlabel("Year")
50 total_suicides_no.set_ylabel("Number of Suicides")
51 plt.title("Number of Suicides")
52
53 plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
54 plt.savefig("Images/num_suicides.jpg",bbox_inches = "tight")
55 plt.show()
56
57 # Plot Suicide Rates by Gender
58
59 total_rate_plot = total_rate.plot(kind='line', color="coral",label="Total Suicides Rate")
60 male_rate_plot = male_rate.plot(kind='line', color="skyblue", label="Male Suicides Rate")
61 female_rate_plot = female_rate.plot(kind='line', color="gold", label='Female Suicides Rate')
62
63 total_rate_plot.set_xlabel("Year")
64 total_rate_plot.set_ylabel("Suicides Rates")
65 plt.title("Suicides Rates")
66
67 plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
68 plt.savefig("Images/rate_suicides.jpg",bbox_inches = "tight")
69 plt.show()
70
71 # Suicides Rates VS Country Over 1990–2014
72
73 country_year=data.groupby(['year','country'])['suicides_no'].sum()
74 country_population=data.groupby(['year','country'])['population'].sum()
75 country_rate=country_year/country_population
76
77 # Plot Suicide Rates By Country
78
79 country_rate_plot= country_rate.unstack().plot(kind='line', figsize = (10,6))
80
81 country_rate_plot.set_xlabel("Year")
82 country_rate_plot.set_ylabel("Suicides Rate")
83 plt.title("Suicides Rate")
84
85 plt.legend(loc='center left', bbox_to_anchor=(1, 0.5),
86 |           ncol=4)
87 plt.savefig("Images/country_year.jpg",bbox_inches = "tight")
88 plt.show()
89
90 # Specifically look at Suicide Rates in Mexico
```

PYTHON CODES

```
91
92     mexico = data.loc[(data["country"] == "Mexico")]
93     mex_suicides = mexico.groupby(['year'])['suicides_no'].sum()
94     mex_population=mexico.groupby(['year'])['population'].sum()
95
96     mex_rate=mex_suicides/mex_population
97     mex_plot = mex_rate.plot(kind='line', color="gold", label='Mexico Suicides Rate')
98
99     total_suicides_no.set_xlabel("Year")
100    total_suicides_no.set_ylabel("Suicides Rate")
101    plt.title("Mexico Suicides Rate")
102
103    plt.legend(loc='best', bbox_to_anchor=(1, 0.5))
104    plt.savefig("Images/mex_suicides.jpg",bbox_inches='tight')
105    plt.show()
106
107    # Suicide Rates by Age Over 1990–2014
108
109    age_year=data.groupby(['year','age'])['suicides_no'].sum()
110    age_population=data.groupby(['year','age'])['population'].sum()
111    age_rate=age_year/age_population
112
113    # Plot Suicide Rates by Age
114
115    age_rate_plot= age_rate.unstack().plot(kind='line')
116
117    age_rate_plot.set_xlabel("Age Group")
118    age_rate_plot.set_ylabel("Suicides Rate")
119    plt.title("Suicides Rate")
120
121    plt.legend(loc='center left', bbox_to_anchor=(1, 0.5),
122               | ncol=1)
123    plt.savefig("Images/age_year.jpg",bbox_inches = "tight")
124    plt.show()
125
126
127    # ## Analyse Age/Gender Factor
128
129    # Group by Age,Gender
130
131    sub_data = data.loc[:,["country","year","sex","age","suicides_no","population","suicides/100k pop"]]
132    groupby_age_sex = sub_data.groupby(["age","sex"])
133    groupby_age_sex.count()
```

PYTHON CODES

```
135 # Create New DataFrame
136 suicide_sum = pd.DataFrame({"suicides_sum": groupby_age_sex["suicides_no"].sum()})
137 # Unstack the DataFrame to create bar chart
138 suicide_sum = suicide_sum.unstack()
139 suicide_sum
140
141 # Create an index column to sort age group
142
143 ind = [2, 3, 4, 1, 5, 6]
144 suicide_sum["index_col"] = ind
145 suicide_sum = suicide_sum.sort_values("index_col")['suicides_sum']
146 suicide_sum
147
148 # Plot Number of Suicides by Age Group
149
150 suicide_sum.plot(kind="bar", rot = 0, figsize = (10,5),
151 |           | color = ["orangered","skyblue"], fontsize = 12, alpha = 0.8)
152
153 plt.legend(["Female","Male"],title="Sex")
154 plt.title("Number of Suicides within Each Age Group")
155 plt.xlabel("Age Group")
156 plt.ylabel("Count")
157 plt.grid()
158 plt.savefig("Images/age_group_counts.png", bbox_inches = "tight")
159 plt.show()
160
161 # Suicide rate for Each Age Group
162
163 suicide_rate = pd.DataFrame({"suicide_rates_%": groupby_age_sex["suicides_no"].sum()/groupby_age_sex["population"].sum()*100})
164 suicide_rate = suicide_rate.unstack()
165
166 # Sort by Age Group
167
168 ind = [2, 3, 4, 1, 5, 6]
169 suicide_rate["index_col"] = ind
170 suicide_rate = suicide_rate.sort_values("index_col")["suicide_rates_%"]
171 suicide_rate
172
173 # Plot Suicides Rates by Age Group
174
175 suicide_rate.plot(kind="bar", rot = 0, figsize = (10,5),
176 |           | color = ["orangered","skyblue"], fontsize = 12, alpha = 0.8)
177
178 plt.legend(["Female","Male"],title="Sex")
179 plt.title("Suicide Rates within Each Age Group")
```

PYTHON CODES

```
180 plt.xlabel("Age Group")
181 plt.ylabel("Suicide Rates (%)")
182 plt.grid()
183 plt.savefig("Images/age_group_rates.png", bbox_inches = "tight")
184 plt.show()
185
186 # Number of Suicides by Gender
187
188 sex_group = sub_data.groupby(["sex"])
189 sum_by_sex = pd.DataFrame(sex_group["suicides_no"].sum())
190 sum_by_sex
191
192 # Plot Number of Suicides by Gender
193
194 plt.figure()
195
196 sum_by_sex.plot(kind = "bar", color = "skyblue", width = 0.3, alpha = 0.8, rot = 0, fontsize = 12)
197 plt.title("Suicide Counts by Sex")
198 plt.ylabel("Count")
199 plt.legend(["Suicides Num"], loc="best")
200 plt.savefig("Images/sex_group_counts.png", bbox_inches = "tight")
201 plt.show()
202
203 # Suicide Rates by Gender
204
205 rate_by_sex = pd.DataFrame({"suicide_rates_%": sex_group["suicides_no"].sum()/sex_group["population"].sum()*100})
206 rate_by_sex
207
208 # Plot Suicide Rates by Gender
209
210 plt.figure()
211
212 rate_by_sex.plot(kind = "bar", color = "skyblue", width = 0.3, alpha = 0.8, rot = 0, fontsize = 12)
213 plt.title("Suicide Rates by Sex")
214 plt.ylabel("Suicide Rates (%)")
215 plt.legend(["Suicides Rates"], loc="best")
216 plt.savefig("Images/sex_group_rates.png", bbox_inches = "tight")
217 plt.show()
218
219 # Perform two sample T-Test on female and male data
220
221 df_female = sub_data.loc[sub_data["sex"] == "female", :]
222 df_male = sub_data.loc[sub_data["sex"] == "male", :]
223
```

PYTHON CODES

```
224 stats.ttest_ind(df_female["suicides_no"], df_male["suicides_no"], equal_var=False)
225
226 stats.ttest_ind(df_female["suicides/100k pop"], df_male["suicides/100k pop"], equal_var=False)
227
228 # Male to Female Suicides Ratio
229
230 suicide_sum["male_to_female_ratio"] = suicide_sum["male"]//suicide_sum["female"]
231 suicide_sum
232
233 # # Location
234
235 # get data from 2010 to 2014
236 file_2010_2014 = "data/sub_dataset_2010-2014.csv"
237 data_2010_2014 = pd.read_csv(file_2010_2014, encoding = "ISO-8859-1")
238 data_2010_2014.head()
239
240 # calculate total number of suicides and population of each country during this 5 years
241 rate_2010_2014 = data_2010_2014.groupby(["country"])
242 rate_2010_2014 = rate_2010_2014.sum()
243 rate_2010_2014.head()
244
245 # add a column and calculate total suicides rate of each country during this 5 years
246 rate_2010_2014["suicides_rate_%"] = rate_2010_2014["suicides_no"]/rate_2010_2014["population"]*100
247 rate_2010_2014 = rate_2010_2014[["suicides_no","population","suicides_rate_%"]]
248
249 # add 2 columns for lat and lng
250 rate_2010_2014["lat"]= ""
251 rate_2010_2014["lng"]= ""
252 rate_2010_2014.head()
253
254 # Loop through the countries and run a lat/long search for each city
255 # save lat and lng data to rate_2010_2014
256 for index, row in rate_2010_2014.iterrows():
257
258     country = index
259     target_url = ('https://maps.googleapis.com/maps/api/geocode/json?
260                 | address={0}&key={1}').format(country, gkey)
261
262     try:
263         response = requests.get(target_url).json()
264         lat = response['results'][0]["geometry"]["location"]["lat"]
265         lng = response['results'][0]["geometry"]["location"]["lng"]
266
267         rate_2010_2014.loc[index, "lat"] = lat
268         rate_2010_2014.loc[index, "lng"] = lng
```

PYTHON CODES

```
269     except :
270         print(country)
271
272 rate_2010_2014.head()
273
274 # reset the index
275 rate_general_2010_2014 = rate_2010_2014.reset_index()
276 rate_general_2010_2014.head()
277
278 # save the df to csv file
279 rate_general_2010_2014.to_csv("data/dataset_attd_2010_2014.csv", encoding = "utf-8", index = True, header = True)
280
281 # get data from dataset_attd_2010_2014.csv
282 file_att = "data/dataset_attd_2010_2014.csv"
283 rate_general_2010_2014 = pd.read_csv(file_att, encoding = "ISO-8859-1")
284 rate_general_2010_2014.head()
285
286 # sort data by suicides rate
287 rate_general_2010_2014_rank = rate_general_2010_2014.sort_values("suicides_rate_%", ascending=False)
288 rate_general_2010_2014_rank.head()
289
290 # save the top 3 to "top3"
291 rate_general_2010_2014_rank = rate_general_2010_2014_rank.reset_index(drop = True)
292 rate_general_2010_2014_rank
293 top3 = rate_general_2010_2014_rank.iloc[0:3,:]
294 top3
295
296 # save the bottom 3 to "bottom3"
297 bottom3 = rate_general_2010_2014_rank.iloc[68:71,:]
298 bottom3
299
300 # combine the top3 and bottom 3 as top_bottom
301 top_bottom = pd.concat([top3,bottom3])
302 top_bottom
303
304 # plot the top 3 and bottom 3
305 plt.figure(figsize=(10,6))
306 top_bottom_plot = plt.bar(top_bottom['country'], top_bottom["suicides_rate_%"], color='r', alpha=0.5, align="center")
307
308 plt.xlabel("Countries")
309 plt.ylabel("Suicides Rate %")
310 plt.title("Top 3 and Bottom 3")
311
312 top_bottom_plot[3].set_color('b')
313 top_bottom_plot[4].set_color('b')
```

PYTHON CODES

```
314 top_bottom_plot[5].set_color('b')
315
316 plt.savefig("Images/Location_top_bottom.png")
317
318 # Plot Heatmap
319 locations_71 = rate_general_2010_2014[["lat", "lng"]]
320 suicides_rate_71 = rate_general_2010_2014["suicides_rate_%"].astype(float)
321
322 fig = gmaps.figure()
323
324 # Create heat layer
325 heat_layer = gmaps.heatmap_layer(locations_71, weights=suicides_rate_71 * 1000000,dissipating=False, max_intensity=10000,point_radius=1)
326
327 # Add layer
328 fig.add_layer(heat_layer)
329
330 # read data from "Countries of the world"
331 file_region_71 = "Data/Countries of the world.xls"
332 region_71 = pd.read_excel(file_region_71, encoding = "ISO-8859-1")
333 region_71.head()
334
335 # remove space in column "country" of "region_71" and "Country" of "rate_general_2010_2014"
336 region_71['Country'] = region_71["Country"].str.strip()
337 rate_general_2010_2014['country'] = rate_general_2010_2014["country"].str.strip()
338
339 # merge "region_71" and "rate_general_2010_2014"
340 data_merged_71 = pd.merge(rate_general_2010_2014,region_71, left_on = "country",right_on = "Country", how = "inner")
341 data_merged_71.head()
342
343 # check which country is left after merge
344 list1 = data_merged_71["country"].tolist()
345 list2 = rate_general_2010_2014["country"].tolist()
346
347 for x in list2:
348     if x not in list1:
349         print(x)
350
351 # read data from "country-and-continent.csv"
352 file_continent = "Data/country-and-continent.csv"
353 continent = pd.read_csv(file_continent, encoding = "ISO-8859-1")
354 continent.head()
355
356 # split the "Country_Name" of "continent" and save to new
357 new= continent["Country_Name"].str.split(", ", n = 1, expand = True)
```

PYTHON CODES

```
359 # add the first column of "new" to "CountryName" in "continent"
360 continent["CountryName"] = new[0]
361 continent.head()
362
363 # Merge "continent" with "data_merged_71"
364 data_complete = pd.merge(data_merged_71, continent, left_on = "country", right_on = "CountryName", how = "inner")
365
366 # list the columns of "data_complete"
367 data_complete.columns
368
369 # remove unuseful columns
370 data_complete = data_complete[['country', 'CountryName', 'suicides_no', 'population', 'suicides_rate_%', 'lat', 'lng',
371     'Country', 'Region', 'Population', 'Area sq. mi.',
372     'Pop. Density per sq. mi.', 'Coastline coast/area ratio',
373     'Climate', 'Agriculture', 'Industry',
374     'Service', 'Continent_Name','Country_Name']]
375
376 data_complete.head(3)
377
378 # Delete the country that show twice
379 data_complete = data_complete.drop(data_complete[(data_complete['country'] == "Armenia") & (data_complete['Continent_Name'] == "Asia")].index)
380 data_complete = data_complete.drop(data_complete[(data_complete['country'] == "Cyprus") & (data_complete['Continent_Name'] == "Asia")].index)
381 data_complete = data_complete.drop(data_complete[(data_complete['country'] == "Turkey") & (data_complete['Continent_Name'] == "Europe")].index)
382 data_complete = data_complete.drop(data_complete[(data_complete['country'] == "Kazakhstan") & (data_complete['Continent_Name'] == "Europe")].index)
383 data_complete = data_complete.drop(data_complete[(data_complete['country'] == "Russian Federation") & (data_complete['Continent_Name'] == "Asia")].index)
384 data_complete = data_complete.drop(data_complete[(data_complete['country'] == "Georgia") & (data_complete['Continent_Name'] == "Asia")].index)
385
386 # Check if each country shows only once
387 data_complete["country"].value_counts()
388
389 # Check which country is left after merge
390 list1 = data_complete["country"].tolist()
391 list2 = data_merged_71["country"].tolist()
392
393 for x in list2:
394     if x not in list1:
395         print(x)
396
397 # save the df to csv file
398 data_complete.to_csv("Data/data_complete.csv", encoding = "utf-8", index = True, header = True)
399
400 # classify countries by continent
401 continent_count = pd.DataFrame(data_complete["Continent_Name"].value_counts())
402 continent_count = continent_count.rename(columns={"Continent_Name": "Continent"})
```

PYTHON CODES

```
405 # Plot the target dataset in terms of which countries are included
406 continent_plot = continent_count.plot(kind="bar", figsize=(9,6), alpha=0.7)
407 plt.xticks(rotation = 0)
408 plt.title("Number of Countries by Continents")
409 plt
410
411 # Save the figure
412 plt.savefig("Images/Location_countries_analyzed.png")
413
414 # Plot bar chart to show mean suicides rate of each continent
415 group_by_continent = data_complete.groupby(data_complete["Continent_Name"])
416 mean_continent = group_by_continent["suicides_rate_%"].mean()
417 mean_continent.sort_values(inplace=True, ascending=False)
418 conti_plot = mean_continent.plot(kind="bar", grid=True, figsize=(8,6),
419 | | | title="Suicides Rate by Continent")
420 plt.xticks(rotation = 0)
421
422 # Save the figure
423 plt.savefig("Images/Location_continent.png")
424
425 # plot scatter chart show relation between suicides rate and population density
426 area_plot = data_complete.plot(kind="scatter", x="Pop. Density per sq. mi.", y="suicides_rate_%", grid=True, figsize=(8,6),
427 | | | title="Suicides Rate vs Population Density")
428 plt.ylabel("Suicides Rate %")
429 plt.xlim(-0.75, 610)
430
431 # Calculate correlation coefficient
432 cc_popden = data_complete["suicides_rate_%"].corr(data_complete["Pop. Density per sq. mi."])
433 print(f" The R-value is {cc_popden}")
434
435 # Save the figure
436 plt.savefig("Images/Location_PopulationDensity.png")
437
438 # plot scatter chart show relation between suicides rate and coast/area ratio
439 coast_plot = data_complete.plot(kind="scatter", x="Coastline coast/area ratio", y="suicides_rate_%", grid=True, figsize=(8,6),
440 | | | title="Suicides Rate vs Coast/Area Ratio")
441 plt.xlim(-0.75, 12)
442 plt.ylabel("Suicides Rate %")
443 plt.xlabel("Coast/Area Ratio")
444
445 # Calculate correlation coefficient
446 cc_coast_area = data_complete["suicides_rate_%"].corr(data_complete["Coastline coast/area ratio"])
447 print(f" The R-value is {cc_coast_area}")
```

PYTHON CODES

```
449 # Save the figure
450 plt.savefig("Images/Location_CoastAreaRatio.png")
451
452 # plot scatter chart show relation between suicides rate and Industry Percentage
453 agri_plot = data_complete.plot(kind="scatter", x="Industry", y="suicides_rate_%", grid=True, figsize=(8,6),
454 |           title="Suicides Rate vs Industry Percentage")
455 plt.ylabel("Suicides Rate %")
456 plt.xlabel("Industry Percentage")
457
458 # Calculate correlation coefficient
459 cc_IndPercentage = data_complete["suicides_rate_%"].corr(data_complete["Industry"])
460 print(f" The R-value is {cc_IndPercentage}")
461
462 # Save the figure
463 plt.savefig("Images/Location_Industry.png")
464
465 # Calculate the average suicides rate for each climate type
466 group_by_climate = data_complete.groupby(data_complete["Climate"])
467 mean_climate = group_by_climate["suicides_rate_%"].mean()
468
469 # plot scatter chart show relation between suicides rate and climate
470 ax = data_complete.plot(kind="scatter", x="Climate", y="suicides_rate_%", grid=True, figsize=(8,6),
471 |           title="Suicides Rate vs Climate")
472
473 # Add a line show mean value
474 mean_climate.plot(kind = "line",ax = ax,color='red',linestyle = "-.")
475 ax.set_xlim(bottom = 0.0,ymax = 0.033)
476 ax.set_xlim(xmin=0.8, xmax=4.2)
477 plt.ylabel("Suicides Rate %")
478
479 # Calculate correlation coefficient
480 cc_climate = data_complete["suicides_rate_%"].corr(data_complete["Climate"])
481 print(f" The R-value is {cc_climate}")
482
483 # Save the figure
484 plt.savefig("Images/Location_Climate.png")
485
486 # convert data of lat to float
487 data_complete["lat"] = data_complete["lat"].astype(float)
488
489 # plot scatter chart show relation between suicides rate and latitude
490 ax = data_complete.plot(kind="scatter", x="lat", y="suicides_rate_%", grid=True, figsize=(8,6),
491 |           title="Suicides Rate vs Latitude (Linear Regression)")
492
```

PYTHON CODES

```
493 # Linear regression
494 (slope, intercept, r_value, p_value, std_err) = linregress(data_complete['lat'], data_complete['suicides_rate_%'])
495 fit = slope * data_complete['lat'] + intercept
496 ax.plot(data_complete['lat'],fit,"b--",color = "coral")
497 plt.ylabel("Suicides Rate %")
498 plt.xlabel("Latitude")
499
500 # get R-value
501 print(f" The R-value is {r_value}")
502
503 # Save the figure
504 plt.savefig("Images/Location_lat_linear.png")
505
506 # Non-Linear regression
507 ax = data_complete.plot(kind="scatter", x="lat", y="suicides_rate_%", grid=True, figsize=(8,6),
508 |           title="Suicides Rate vs Latitude (Non-linear Regression)")
509
510 def model(z, a, b):
511     return a * np.exp(-b * z)
512
513 x_axis = data_complete["lat"]
514 y_axis = data_complete["suicides_rate_%"]
515
516 popt, pcov = curve_fit(model, x_axis, y_axis, p0=(5, 0.1))
517
518 # prepare some data for a plot
519 xx = np.linspace(-40, 70)
520 yy = model(xx, *popt)
521
522 ax.plot(xx,yy, '-',color = "coral")
523
524 r_value = model(x_axis, *popt) - y_axis
525 model_y = model(x_axis, *popt)
526
527 plt.ylabel("Suicides Rate %")
528 plt.xlabel("Latitude")
529
530 # Save the figure
531 plt.savefig("Images/Location_lat_expo.png")
532
533 print(popt)
534 r_value.head()
```

PYTHON CODES

```
536 # exponential model analysis
537 # calculate the mean value for "lat" and "suicides_rate_%"
538 data_analysis = pd.DataFrame(x_axis)
539 data_analysis["suicides_rate_%"] = data_complete["suicides_rate_%"]
540 x_mean = data_analysis["lat"].mean()
541 y_mean = data_analysis["suicides_rate_%"].mean()
542
543 # exponential model analysis
544 # calculate r-square
545 data_analysis["model_y"] = model_y
546 data_analysis["r_value"] = r_value
547
548 data_analysis["(y-y_mean)2"] = (data_analysis["suicides_rate_%"] - y_mean)**2
549 ss_tot = data_analysis["(y-y_mean)2"].sum()
550
551 data_analysis["(y_model-y_mean)2"] = (data_analysis["model_y"] - y_mean)**2
552 ss_reg = data_analysis["(y_model-y_mean)2"].sum()
553
554 data_analysis["(r_value)2"] = (data_analysis["r_value"])**2
555 ss_res = data_analysis["(r_value)2"].sum()
556
557 r_square = 1 - (ss_res / ss_tot)
558 print(f" The R-square for this model is {r_square}")
559
560 data_analysis.head()
561
562 # linear model analysis
563 data_analysis_1 = pd.DataFrame(x_axis)
564 data_analysis_1["suicides_rate_%"] = data_complete["suicides_rate_%"]
565
566 data_analysis_1["model_y"] = fit
567 data_analysis_1["r_value"] = data_analysis_1["model_y"] - data_analysis_1["suicides_rate_%"]
568
569 data_analysis_1["(y-y_mean)2"] = (data_analysis_1["suicides_rate_%"] - y_mean)**2
570 ss_tot = data_analysis_1["(y-y_mean)2"].sum()
571
572 data_analysis_1["(y_model-y_mean)2"] = (data_analysis_1["model_y"] - y_mean)**2
573 ss_reg = data_analysis_1["(y_model-y_mean)2"].sum()
574
575 data_analysis_1["(r_value)2"] = (data_analysis_1["r_value"])**2
576 ss_res = data_analysis_1["(r_value)2"].sum()
577
578 r_square = 1 - (ss_res / ss_tot)
579 print(f" The R-square for this model is {r_square}")
```

PYTHON CODES

```
584 # # Economy Analysis (GDP)
585
586 # Read source data again for GDP
587 file = "data/sub_dataset.csv"
588 datagdp = pd.read_csv(file)
589 datagdp.head()
590
591 # sorting by country, year and sex
592 datagdp.sort_values(["country","year","sex"],inplace = True)
593
594 # # dropping ALL duplicate rows
595 datagdp.drop_duplicates(subset =['country','gdp_for_year ($ ','gdp_per_capita ($ ')],keep = 'first', inplace = True)
596
597 # Sum up suicide_no per year
598 grouped_per_year = datagdp.groupby('year')
599
600 # Sum up GDP per year
601 gdp = grouped_per_year['gdp_per_capita ($').sum()
602
603 # x axis is year
604 yeardata = datagdp['year'].unique()
605
606 # Suicides Group by Year (from earlier analysis)
607 groupby_yr=data.groupby("year")
608 year_sum = groupby_yr["suicides_no"].sum()
609 year_sum.head()
610
611 total_rate=year_sum/total_population
612
613 # Suicides rate (%)
614 suicides = total_rate*100
615
616 # Plot GDP and Suicide Rate yearly
617
618 width = .35 # width of a bar
619
620 m1_t = pd.DataFrame({
621   'gdp' : gdp,
622   'suicides' : suicides})
623
624 m1_t = m1_t.reset_index()
625
626
627 fig, ax = plt.subplots()
```

PYTHON CODES

```
628 fig.suptitle('Yearly GDP and Suicides Trend Comparison')
629 ax1 = plt.gca()
630 ax2 = ax1.twinx() # instantiate a second axes that shares the same x-axis
631
632 ax1.bar(x =m1_t['year'], height = m1_t['gdp'], width = width, color ='skyblue')
633
634 ax2.plot(m1_t['year'], m1_t['suicides'] ,color ='Coral', lw=2, ls='--')
635
636 ax.set(xlabel='Years 1990-2014', ylabel='GDP per Capita ($)')
637 ax2.set(ylabel='Yearly Suicides Rate (%)')
638
639 plt.savefig('Images/Yearly GDP and Suicides Trend Comparison.png',bbox_inches='tight')
640
641 # Yearly Change analysis of GDP and Suicides Rate
642
643 #Find First discrete difference of series Suicides and GDP
644 # Shift 1 period to calculate difference ---> yearly difference (GDP growth/decrease)
645 s_diff = suicides.diff()
646 g_diff = gdp.diff()
647
648
649 m2_t = pd.DataFrame({
650     'gdp_diff' : g_diff,
651     'suicides_diff' : s_diff})
652
653 m2_t = m2_t.reset_index()
654 m2_t.head()
655
656 fig, ax1 = plt.subplots()
657 fig.suptitle('Yearly Change in GDP ($) and Suicides Rate (%) ')
658
659 ax = plt.gca()
660 ax2 = ax1.twinx() # instantiate a second axes that shares the same x-axis
661
662 ax1.bar(x =m2_t['year'], height = m2_t['gdp_diff'], width = width, color ='skyblue')
663 ax2.plot(m2_t['year'], m2_t['suicides_diff'], color ='Coral', lw=2, ls='--')
664
665 ax.set(xlabel='Years 1990-2014', ylabel='Yearly Change in GDP per Capita ($)')
666 ax2.set(ylabel='Yearly Change in Suicides Rate (%)')
667
668 plt.savefig('Images/Yearly Change in GDP ($) and Suicides Rate (%).png',bbox_inches='tight')
669 plt.show()
670
```

PYTHON CODES

```
671 # Yearly Change Rate analysis of GDP and Suicides Rate
672
673 # Ratio calculation
674 ggrowth = g_diff*100/gdp
675 sgrowth = s_diff*100/suicides
676
677
678 m3_t = pd.DataFrame({
679     'gdp_growth' : g_diff*100/gdp,
680     'suicides_growth' : s_diff*100/suicides})
681
682 m3_t = m3_t.reset_index()
683 m3_t.head()
684
685 fig, ax1 = plt.subplots()
686 fig.suptitle('Yearly % Change of GDP and Suicides ')
687
688 ax = plt.gca()
689 ax2 = ax1.twinx() # instantiate a second axes that shares the same x-axis
690
691 ax1.bar(x =m3_t['year'], height = m3_t['gdp_growth'], width = width, color ='skyblue')
692 ax2.plot(m3_t['year'], m3_t['suicides_growth'], color ='Coral', lw=2, ls='--')
693
694 ax.set(xlabel='Years 1990–2014', ylabel='Yearly % Change in GDP per Capita (%)')
695 ax2.set(ylabel='Yearly % Change in Suicides (%)')
696
697 plt.savefig('Images/Yearly Change Rate of GDP and Suicides.png',bbox_inches='tight')
698
699 plt.show()
700
701 # Check m3_t table
702 m3_t.head()
703
704 # Finding: first value is invalid since no prior year for first year 1990
705
706 # drop first value which is invalid (missing value)
707 sgrowth= sgrowth.drop(index=[1990])
708 ggrowth = ggrowth.drop(index=[1990])
709
710 # Find statistics for suicides
711 v1 = np.var(sgrowth)
712
713 std1 = np.std(sgrowth)
714
```

PYTHON CODES

```
715 mean1 = np.mean(sgrowth)
716 print('Statistics for Suicides change rate data:')
717 print('-----')
718 print(f'Variance: {v1}  Mean: {mean1} and ')
719 print(f'Standard Deviation: {std1}')
720
721 # Find statistics for GDP
722 v2 = np.var(growth)
723
724 std2 = np.std(growth)
725
726 mean2 = np.mean(growth)
727 print('Statistics for GDP change rate data:')
728 print('-----')
729 print(f'Variance: {v2}  Mean: {mean2} and ')
730 print(f'Standard Deviation: {std2}')
731
732 # Coefficient of Variation = Standard deviation / mean
733 coef1 = std1/ mean1
734
735 coef2 = std2 /mean2
736 print(f'Coefficient of Suicides change rate: {coef1}  and Coefficient of GDP change rate: {coef2} ')
737
738
739 # ### Coefficient difference is evident and important
740 # * Opposite sign (resulting from negative mean of Suicides Rate change rate)
741 # * Different absolute value : -7.131614 and 1.60759356
742 #
743 # ```However when coefficient is a negative value , its representation and comparison power with
744 # |another positive mean population is limited, thus it is eliminated from final presentation```
745
746 # # Null hypothesis:
747 # Suicides evolution is not significantly related to GDP evolution.
748 #
749 # # Hypothesis:
750 # Suicides evolution is significantly influenced by GDP evolution.
751
752 fig, ax = plt.subplots()
753
754 # Scatter Plot of Data
755 plt.subplot(2,1,1)
756 plt.scatter(range(len(growth)), growth, label="GDP")
757 plt.scatter(range(len(sgrowth)), sgrowth, label="Suicides")
```

PYTHON CODES

```
758 # Create a legend
759 plt.legend(title="Change Rate (%)", fancybox=True, loc='center left', bbox_to_anchor=(1, 0.5))
760 plt.title('GDP and Suicides Rate Change Distribution')
761
762 # Histogram Plot of Data
763 plt.subplot(2,1,2)
764 plt.hist(growth, 10, density=True, alpha=0.7, label="GDP")
765 plt.hist(sgrowth, 10, density=True, alpha=0.7, label="Suicides")
766 plt.axvline(growth.mean(), color='k', linestyle='dashed', linewidth=1)
767 plt.axvline(sgrowth.mean(), color='k', linestyle='dashed', linewidth=1)
768
769 # Create a legend
770 plt.legend(title="Change Rate (%)", fancybox=True, loc='center left', bbox_to_anchor=(1, 0.5))
771 #----> define legend box location figure , same for subplot1
772
773 plt.savefig('Images/GDP and Suicides Rate Change Distribution.png',bbox_inches='tight')
```