

# Global Suicide Statistics

## An Analysis of Critical Variables

Team BSJ - Shannon Houser, Jack McNeilly, Brian Linder

### Introduction

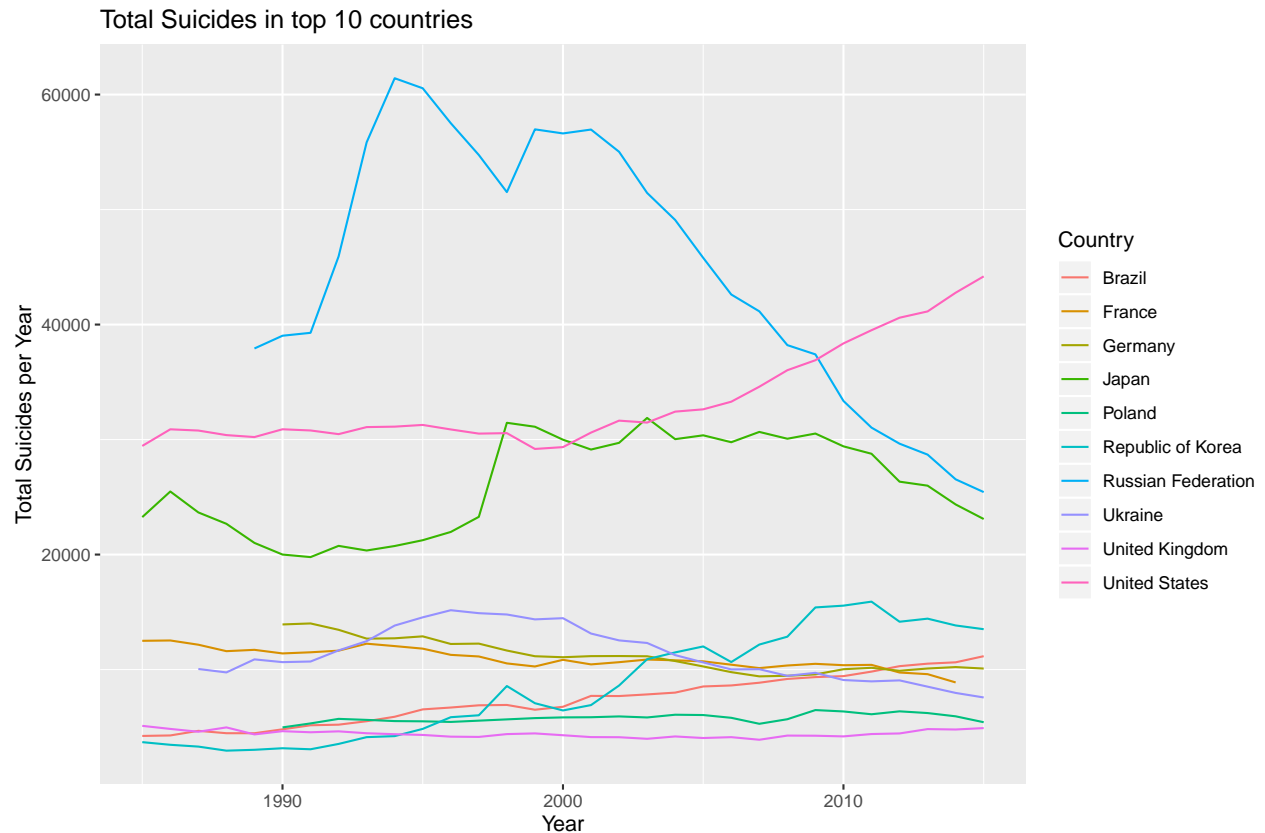
As the Duke community grieves the loss of two of our classmates to suicide in the last 2 weeks, our team plans to analyze general global suicide rates from 1985-2016 in order to see if there are any prevalent factors that might contribute to people taking their own lives. The dataset we have selected compiles data from four distinct datasets that includes information on suicides from over 100 different countries throughout the world. The data compares socio-economic info with suicide rates by country and year. The data is sourced from the World Health Organization, the World Bank and, the United Nations Development Program.

Our goal is to examine these different socio-economic, location, and gender factors to gain insight regarding how the variables of the dataset impact increased suicide rates. Each observation corresponds to the number of suicides that occurred in a certain country and within a certain age and gender group. The variables include country, year, sex, age group, count of suicides, population, suicide rate, country-year composite key, HDI for year, gdp for year, gdp per capita, and generation.

### Data Analysis Plan

#### Summary Statistics and Visualizations

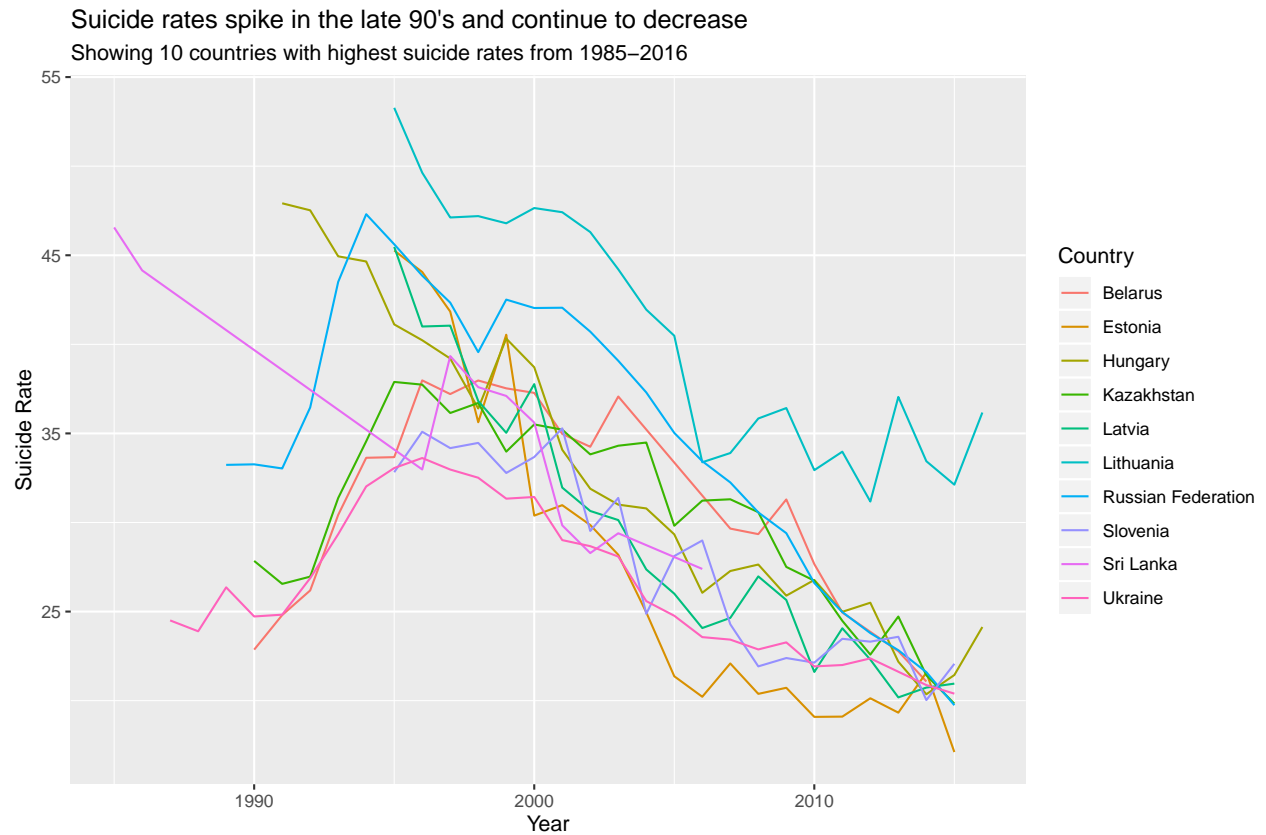
```
# A tibble: 10 x 2
  country      total_suicides
  <chr>         <dbl>
1 Russian Federation 1209742
2 United States    1034013
3 Japan            806902
4 France           329127
5 Ukraine          319950
6 Germany          291262
7 Republic of Korea 261730
8 Brazil           226613
9 Poland           139098
10 United Kingdom   136805
```



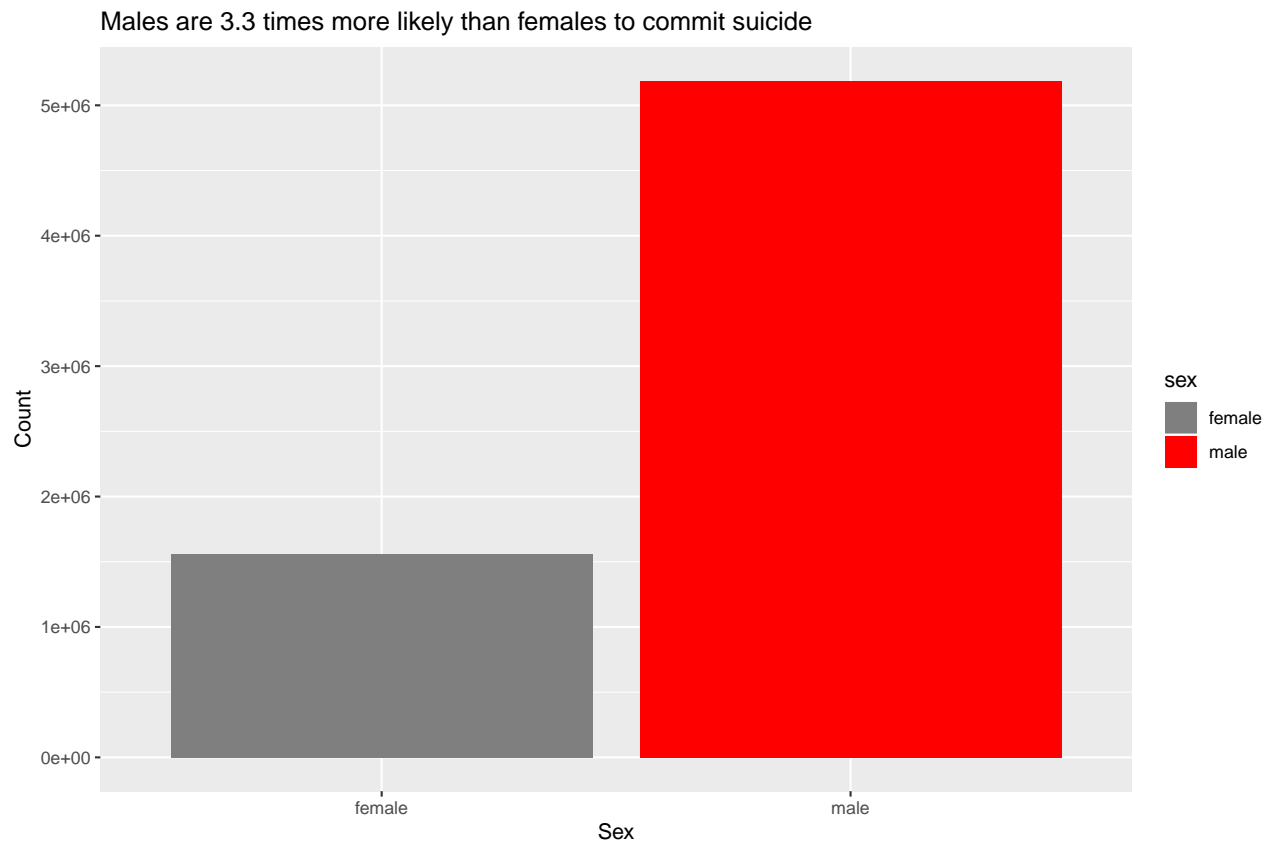
The top ten countries in terms of total suicides seem to be those that are most populated. We then decided to graph how the total number of suicides in these countries has changed over the years. We decided that this is not really helpful information and went on to explore further.

```
# A tibble: 10 x 2
  country      rate_suicide
  <chr>         <dbl>
1 Lithuania    40.4
2 Sri Lanka    35.3
3 Russian Federation 34.9
4 Hungary      32.8
5 Belarus      31.1
6 Kazakhstan   30.5
7 Latvia       29.3
8 Slovenia     27.8
9 Estonia      27.3
10 Ukraine     26.6
```

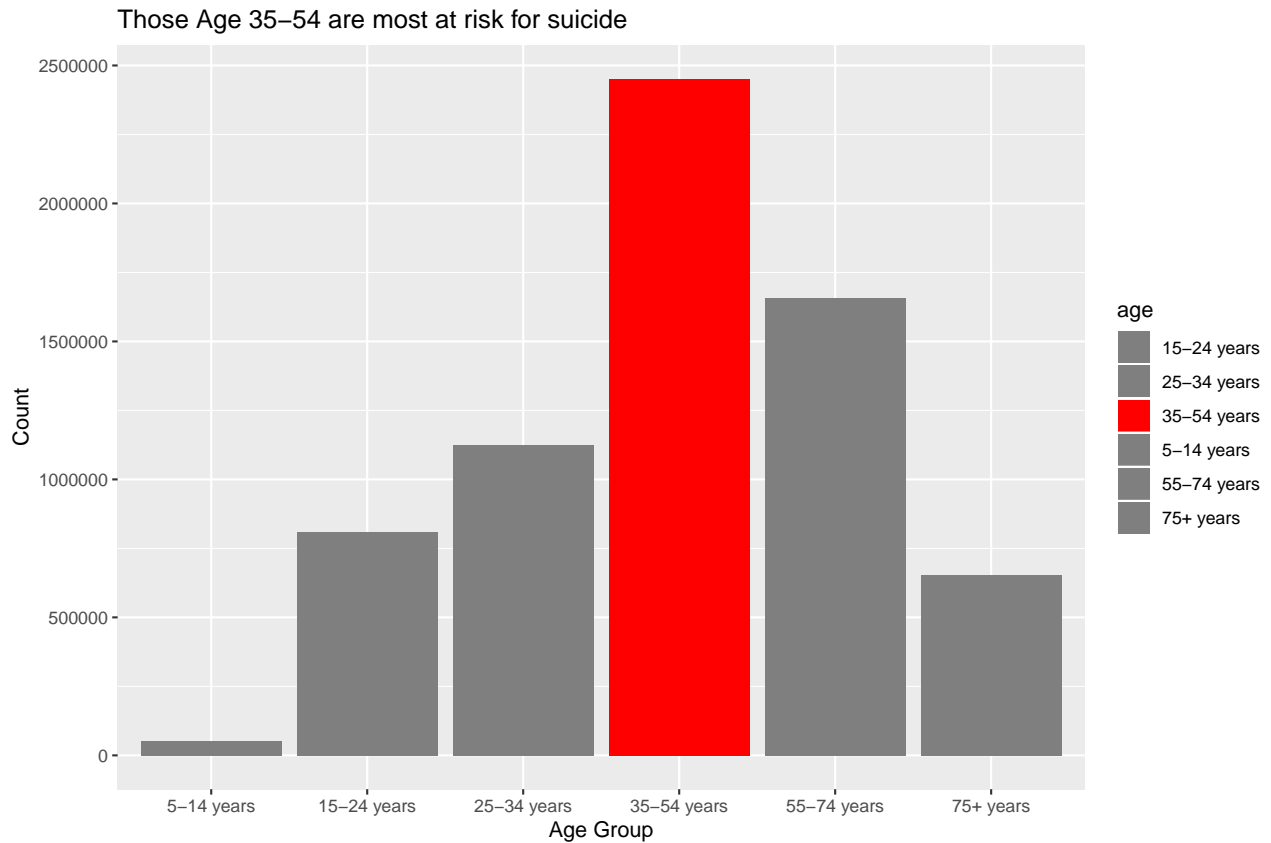
When we examined the suicide rate (per 100k people), we found a strong correlation between geopolitical circumstances as 9/10 of the top 10 countries for suicide rate were part of the ex- Soviet Union. They are all Eastern European countries that may share history, religions, wars, etc. that we are unable to currently predict.



Graphing the average suicide rates over time in the ten countries with the highest suicide rates, we can see that the suicide rates peaked in the late 1990's and have continued to decrease since.



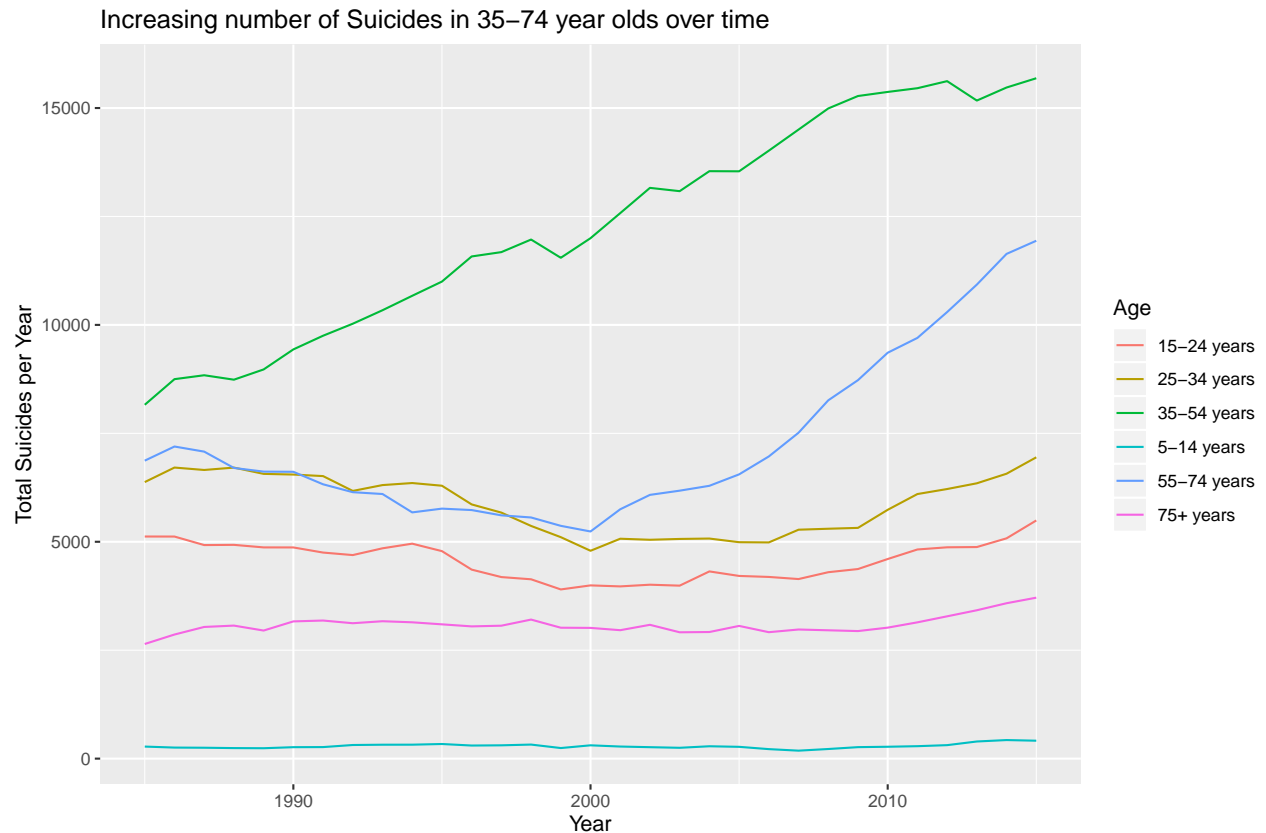
From this data visualization, it is obvious that sex is probably a very important variable when predicting suicide. Men are more than 3.3 times more likely than women to commit suicide.



This visualization shows us that age may also have a strong influence on the likeliness of someone to commit suicide. From the last two visualizations, we may suggest the possibility that middle-aged men are most at risk for suicide. We may also want to examine this relationship more and see brainstorm what life factors make middle-aged men more likely than any other group to commit suicide.

```
# A tibble: 6 x 2
  generation tot_gen
  <chr>      <dbl>
1 Boomers    2284498
2 Silent     1781744
3 Generation X 1532804
4 Millennials  623459
5 G.I. Generation 510009
6 Generation Z  15906
```

This shows the total number of suicides per generation. This has a lot to do with age group and thus is redundant; however, it may help us to better understand what kinds of life circumstances outside of the data these people may have faced to lead them to commit suicide.



From this visualization we can see that over time, the number of global suicides for those between the ages of 35–74 have increased the most drastically. The other age groups seem to be roughly stable; however, it does appear that all other age groups are increasing at the very end of the graph. It would be interesting to see if this sad trend continued past the last year of this study's data collection, 2016.

# A tibble: 101 x 2

	country	avg_gdp
	<chr>	<dbl>
1	United States	1.05e13
2	Japan	4.34e12
3	Germany	2.74e12
4	United Kingdom	1.82e12
5	France	1.78e12
6	Italy	1.48e12
7	Brazil	1.02e12
8	Canada	9.13e11
9	Russian Federation	8.84e11
10	Spain	8.57e11

# ... with 91 more rows

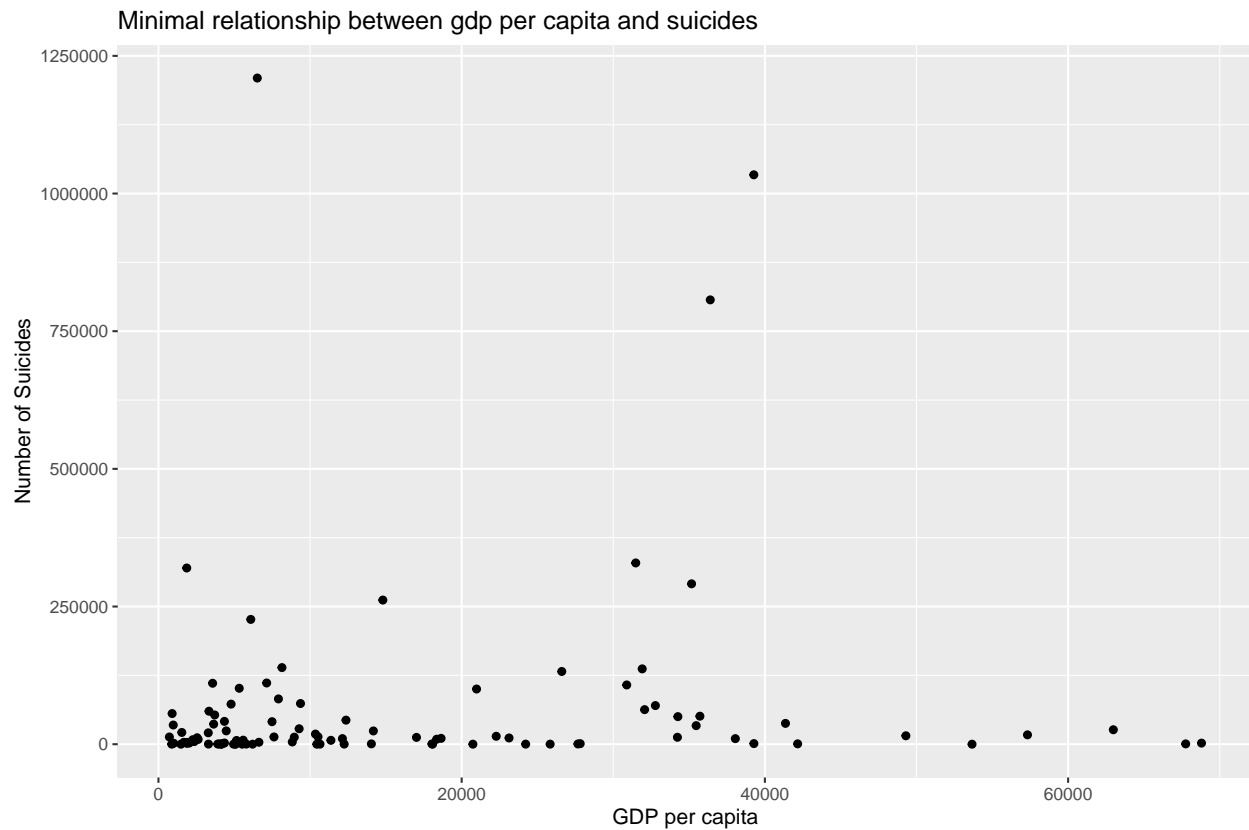
# A tibble: 101 x 2

	country	avg_gdp_capita
	<chr>	<dbl>
1	Luxembourg	68798.
2	Qatar	67756.
3	Switzerland	62982.
4	Norway	57320.
5	San Marino	53664.

```

6 Denmark 49300.
7 United Arab Emirates 42162
8 Sweden 41358.
9 Iceland 39275.
10 United States 39270.
# ... with 91 more rows

```



## Planning

The response variables we will test are the numbers of suicides and the number of suicides per 100k people. The explanatory variables we will examine are different age groups, sex, years, countries, and the socioeconomic status of each country, including their Human Development Index (HDI), growth domestic product (GDP), and GDP per capita.

In addition to observing each individual explanatory variable's impact on the response variables, we will examine how the following combination of explanatory variables and the corresponding result on the response variables:

- Age groups faceted by sex
- HDI with age
- HDI with sex
- HDI with sex and age
- GDP with age
- GDP with sex
- GDP with sex and age
- the above combinations faceted by time period (years)
- the above combinations for Each individual country
- the above combinations for World regions including continents and sub-regions of each continent

In our analysis of the dataset we plan to use statistical methods and tools in R including, linear modeling, regression modeling, a combination of visualization techniques, and null hypothesis testing.

In our analysis of the dataset we plan to use statistical methods and tools in R including, linear modeling, regression modeling, a combination of visualization techniques, and null hypothesis testing.

From our preliminary analysis, we believe that trying to find predictors of suicide using modeling techniques would be a good place to start. We believe that such variables as sex and age may have large impacts on the response variables of total suicide numbers and mean suicide rates. We also believe that gdp and gdp per capita may not play as large of a role as people may think. Instead, perhaps geopolitical factors that are outside of our datasets scope play a large part in suicide determinants.

Furthermore, we plan to explore how these factors have changed over time and if the changes are statistically significant. For example, we will explore whether the total number of suicides in the US has significantly changed between 1985 and 2016, and compare these changes with comparable nations within the data.

## Glimpse of Data

Observations: 27,820

Variables: 12

```
$ country      <chr> "Albania", "Albania", "Albania", "Albania", "A...
$ year         <dbl> 1987, 1987, 1987, 1987, 1987, 1987, 1987, 1987...
$ sex          <chr> "male", "male", "female", "male", "male", "fem...
$ age          <chr> "15-24 years", "35-54 years", "15-24 years", "...
$ suicides_no  <dbl> 21, 16, 14, 1, 9, 1, 6, 4, 1, 0, 0, 0, 2, 17, ...
$ population   <dbl> 312900, 308000, 289700, 21800, 274300, 35600, ...
$ `suicides/100k pop` <dbl> 6.71, 5.19, 4.83, 4.59, 3.28, 2.81, 2.15, 1.56...
$ `country-year` <chr> "Albania1987", "Albania1987", "Albania1987", "...
$ `HDI for year` <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA...
$ `gdp_for_year ($)` <dbl> 2156624900, 2156624900, 2156624900, 2156624900...
$ `gdp_per_capita ($)` <dbl> 796, 796, 796, 796, 796, 796, 796, 796, 796, 7...
$ generation   <chr> "Generation X", "Silent", "Generation X", "G.I..."
```