# Analyzing Happiness on a Global Scale
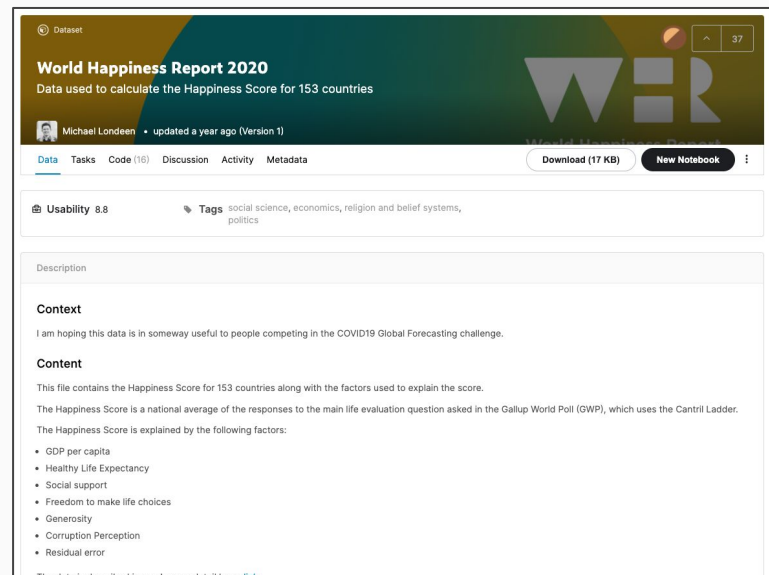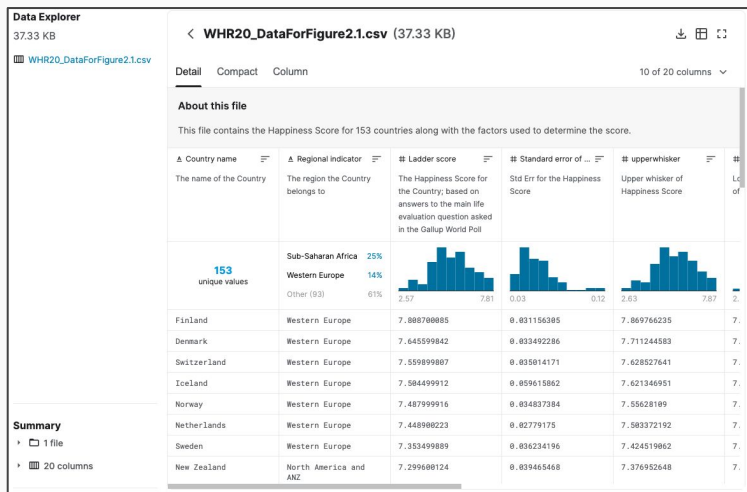
By Valerie Andrade

# Motivation & Summary

- Personally speaking, I'm very curious about anything social science-related and therefore was searching for datasets in this area
  - I'm a big fan of the concept of trying to measure different aspects of human development and seeing data compared against other countries, regions, indicators, etc.

- My goal was to come as close as possible to finding and analyzing datasets that measured some kind of human development data against other specific criteria such as country statistics/figures

# Dataset #1 - World Happiness Report

- Gallup is a global analytics and research firm that since 2005 has produced a "World Happiness Report" using data from a survey called the Gallup World Poll

- This dataset contains a "Happiness Score" for 153 countries and uses data collected from 2017-2019. It was released Feb. 2020

- The Happiness Score is a national average of the responses to the main life evaluation question asked in the Gallup World Poll (GWP), which uses the Cantril Ladder Scale.

- Source: https://www.kaggle.com/londeen/world-happiness-report-2020?select=WHR20_DataForFigure2.1.csv
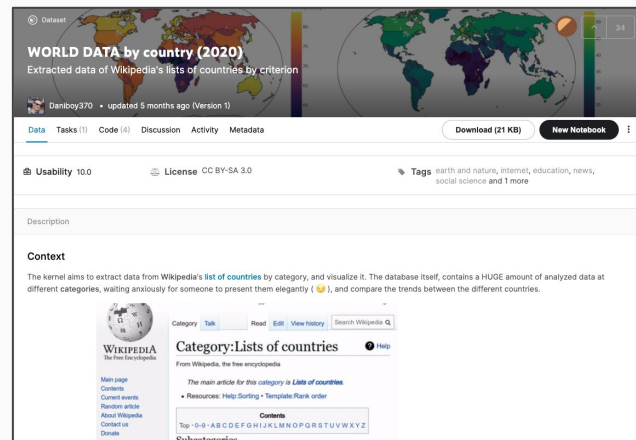
# Dataset #1 - World Happiness Report



- The Happiness Score is explained by the following factors:
  - GDP per capita
  - Healthy Life Expectancy
  - Social support
  - Freedom to make life choices
  - Generosity
  - Corruption Perception
  - Residual error

- Their survey included questions for life evaluations including the Cantril Ladder. Example:
  - **"Are you satisfied or dissatisfied with your freedom to choose what you do with your life?** Please imagine a ladder, with steps numbered from 0 at the bottom to 10 at the top. The top of the ladder represents the best possible life for you and the bottom of the ladder represents the worst possible life for you. On which step of the ladder would you say you personally feel you stand at this time?"

# Dataset #2 - World Data by Country

- This dataset includes extracted global data from Wikipedia by 9 different categories. For this project I used the following:
  - **Fertility Rate in 2018 for 202 countries (via World Bank)**
    - "The fertility rate is the expected number of children born per woman in her child-bearing years."
  - **Median Age in 2018 for 224 countries (via CIA World Factbook)**
    - "Median age is the age that divides a population into two numerically equally sized groups - that is, half the people are younger than this age and half are older."
  - **Urbanization population rate in 2019 for 213 countries (via World Bank)**
    - "Urban population describes the percentage of the total population living in urban areas, as defined by the country."
  - **Population Growth Rate from 2015 - 2020 for 208 countries (via UN)**

# Questions

- Question 1:
  - Are any of the following 4 variables (Fertility Rate, Median Age, Urbanization Rate and Population Growth Rate) correlated to a country's Happiness Score (from per Gallup World Poll)?

- Question 1a:
  - For each variable, how strong, moderate or weak is the correlation?

# Data Cleanup & Exploration

- Merged 5 datasets into 1 main dataframe
- Removed unrelated columns from Happiness Score dataset
- Resolved duplicate country issue that created multiple rows
- Had to go through the CSVs to resolve country names that were inconsistent with another and reimported
- Renamed columns
- Confirmed no missing values
- Grouped countries by region indicators (groupby)
- Noticed about 79 countries from World Data dataset were dropped as Happiness Score dataset only measured 153 countries

# Data Cleanup & Exploration - Examples

## Merged 5 datasets into 1 main dataframe



## Made new DF containing relevant columns



## Renamed column names and reordered them

# Data Cleanup & Exploration - Examples

## Made new DF to groupby Region & Country

```
In [27]: region_groupby = main_df.groupby('Region', as_index = True)['Country']
         region_counts = pd.DataFrame(region_groupby.value_counts())
         region_counts
```

Out[27]:

| | | Country |
|---|---|---|
| Region | Country | |
| Central and Eastern Europe | Albania | 1 |
| | Bosnia and Herzegovina | 1 |
| | Bulgaria | 1 |
| | Croatia | 1 |
| | Czech Republic | 1 |
| | Estonia | 1 |
| | Hungary | 1 |
| | Latvia | 1 |
| | Lithuania | 1 |
| | Montenegro | 1 |

## Made another new DF to groupby Region & get the .mean() for each column

```
In [28]: region_group = main_df.groupby("Region")

         region_df = pd.DataFrame(region_group.mean())
         region_df2 = region_df.sort_values(by='Happiness Score (0 - 10)', ascending = False)
         region_df2
```

Out[28]:

| | Happiness Score (0 - 10) | Fertility Rate in 2018 (births/woman) | Median Age in 2018 | Urbanization Rate in 2019 (%) | Population Growth: 2015-2020 (%) |
|---|---|---|---|---|---|
| Region | | | | | |
| North America and ANZ | 7.173525 | 1.650000 | 39.225000 | 84.170000 | 1.027500 |
| Western Europe | 6.967405 | 1.535000 | 41.530000 | 80.710050 | 0.580500 |
| Latin America and Caribbean | 5.981786 | 2.219048 | 28.571429 | 72.448143 | 1.249524 |
| Central and Eastern Europe | 5.875664 | 1.592857 | 41.764286 | 63.249286 | -0.315714 |
| East Asia | 5.566740 | 1.740000 | 39.840000 | 80.395800 | 0.686000 |
| Southeast Asia | 5.517788 | 2.087500 | 29.162500 | 53.313875 | 1.512500 |
| Commonwealth of Independent States | 5.358342 | 2.193333 | 33.375000 | 55.670000 | 0.835000 |
| Middle East and North | | | | | |

# Data Cleanup & Exploration - **Issues**

## 8 duplications of Guinea found



## Resolved by making new DF & dropping Guinea

# Data Cleanup & Exploration - **Issues**

**Found 2 rows of Russia with different values for Fertility Rates**



| | Country | Regional indicator | Ladder score | Standard error of ladder score | upperwhisker | lowerwhisker | Logged GDP per capita | Social support | Healthy life expectancy | Freedom to make life choices | ... | Explained by: Healthy life expectancy | Explained by: Freedom to make life choices | Explained b Generosi |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 116 | Russia | Commonwealth of Independent States | 5.546 | 0.03961 | 5.623635 | 5.468365 | 10.128872 | 0.903151 | 64.100456 | 0.729893 | ... | 0.680446 | 0.3995 | 0.09904 |
| 117 | Russia | Commonwealth of Independent States | 5.546 | 0.03961 | 5.623635 | 5.468365 | 10.128872 | 0.903151 | 64.100456 | 0.729893 | ... | 0.680446 | 0.3995 | 0.09904 |

2 rows × 25 columns

**Resolved by dropping the lower value. Used the code below:**

- duplicate_row = df1.loc[(df1['Country'] == "Russia") & (df1['Fertility'] == 1.60)].index
- df2 = df1.drop(duplicate_row)
- df2

# Data Analysis - Fertility Rate vs Happiness Score
*Scatter Plot & Linear Regression*



Fertility Rate (2018) vs Happiness Score (2017-2019) by Country

y = -0.59x + 7.04

- The correlation coefficient for these two variables is R=-0.6618.

- The scatter plot indicates a moderate negative linear association between a country's Happiness Score and its average Fertility Rates.

- There appears to be a moderate relationship between the two variables.

- It's interesting to see the data show us that for many countries the Happiness Score decreases as the Fertility Rate increases.

# Data Analysis - Average Fertility Rate by Region
## *Bar Chart*



Fertility Rate by Region (2018)

| Region | Fertility Rate (birth/woman) |
|---|---|
| North America and ANZ | 1.65 |
| Western Europe | 1.54 |
| Latin America and Caribbean | 2.22 |
| Central and Eastern Europe | 1.59 |
| East Asia | 1.74 |
| Southeast Asia | 2.09 |
| Commonwealth of Independent States | 2.19 |
| Middle East and North Africa | 2.59 |
| South Asia | 2.6 |
| Sub-Saharan Africa | 4.38 |

# Data Analysis - Median Age vs Happiness Score
## *Scatter Plot & Linear Regression*



Median Age (2018) vs Happiness Score (2017-2019) by Country

$y = 0.08x + 2.95$

- The correlation coefficient for these two variables is R=0.6753.

- The scatter plot indicates a moderate positive linear association between a country's Happiness Score and its Median Age.

- There appears to be a moderate relationship between the two variables.

- It's interesting to see the data show us that for many countries, the Happiness Score increases as the Median Age increases.

# Data Analysis - Average Median Age by Region
## *Bar Chart*



Average Median Age by Region (2018)

# Data Analysis - Urbanization vs Happiness Score
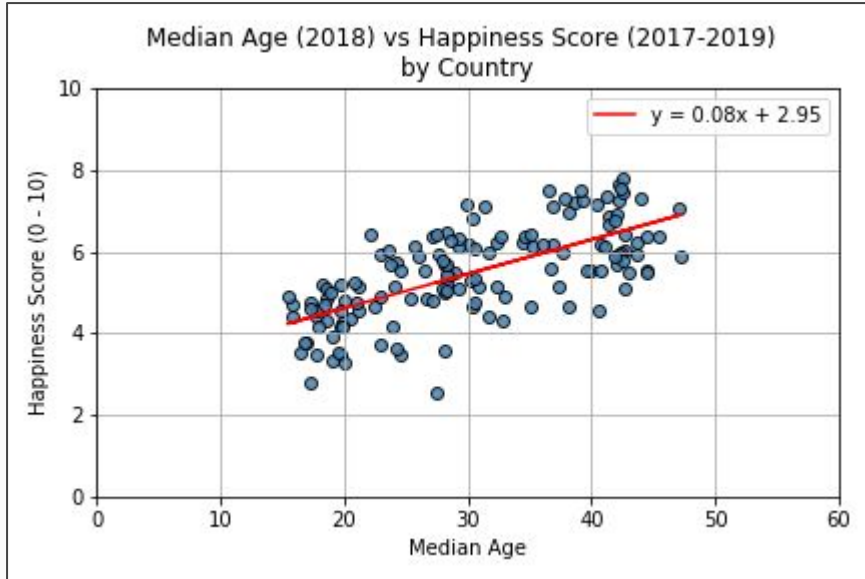*Scatter Plot & Linear Regression*



Urbanization Rate (2019) vs Happiness Score (2017-2019) by Country
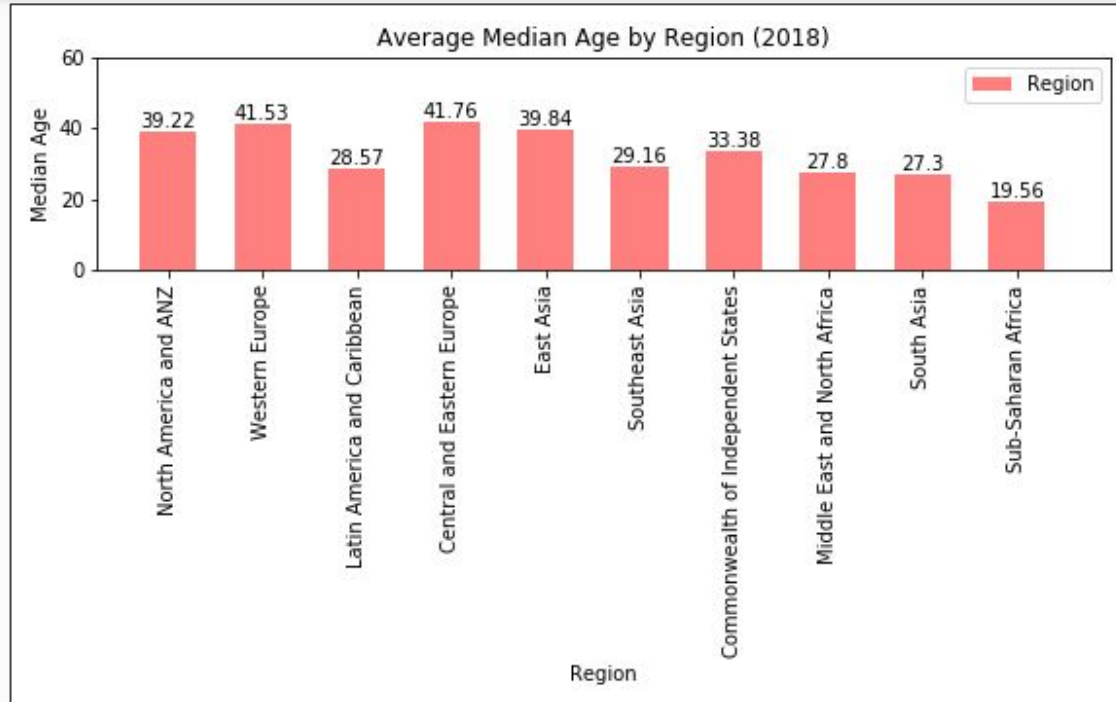
- The correlation coefficient for these two variables is R=0.6676.

- The scatter plot indicates a moderate positive linear association between a country's Happiness Score and Urbanization Population in 2019.

- There appears to be a moderate relationship between the two variables.

- It's interesting to see the data show us that for many countries the Happiness Score is higher when more of the population is urbanized.
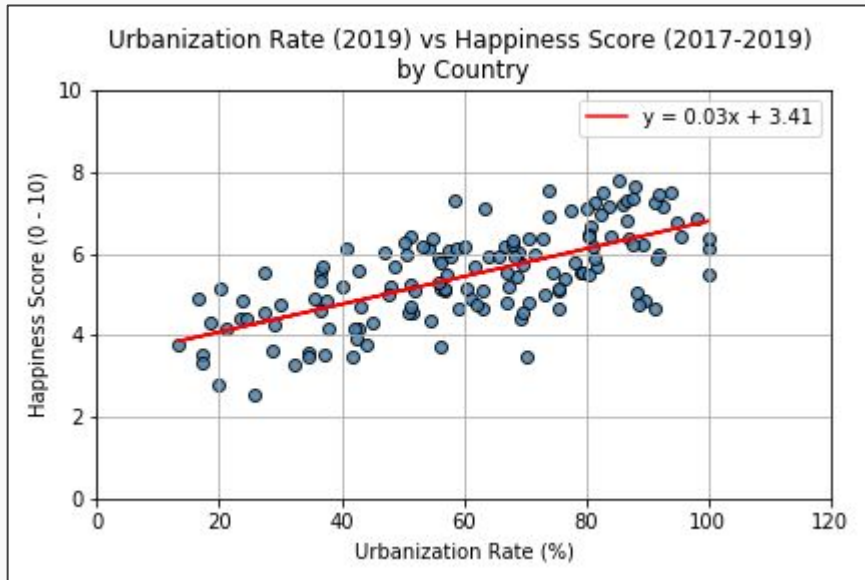
# Data Analysis - Avg. Urbanization Rate by Region
*Bar Chart*



Average Urbanization Rate by Region (2019)

# Data Analysis - Pop. Growth vs Happiness Score
## *Scatter Plot & Linear Regression*



Population Growth Rate (2015-2020) vs Happiness Score (2017-2019) by Country

y = -0.47x + 6.14
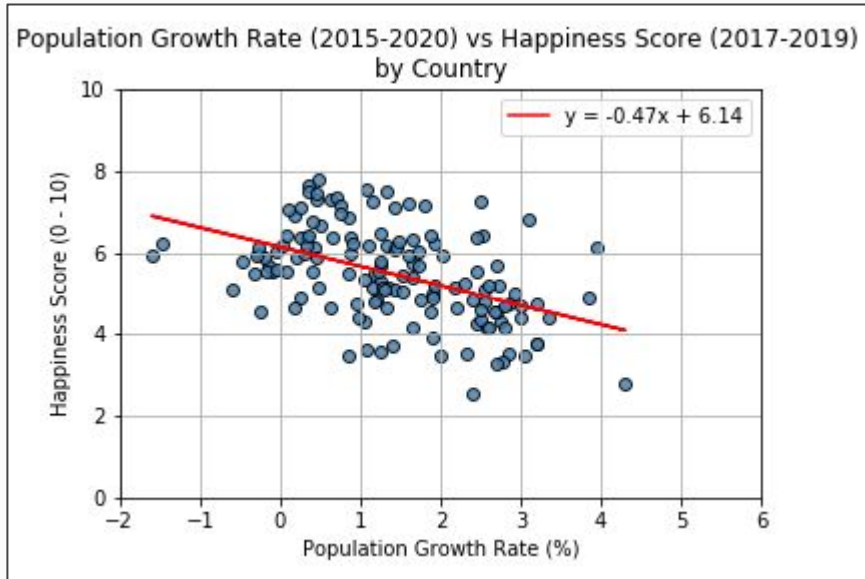
- The correlation coefficient for these two variables is R=-0.4712.

- The scatter plot indicates a semi-moderate negative linear association between a country's Happiness Score and Population Growth Rate (estimated 2015-2020).

- There appears to be a semi-moderate relationship between the two variables.

- It's interesting to see the data show us that for some countries the Happiness Score decreases as the Population Growth Rate increases.

# Data Analysis - Avg. Pop. Growth Rate by Region
*Bar Chart*



Average Population Growth Rate by Region (2015-2020)
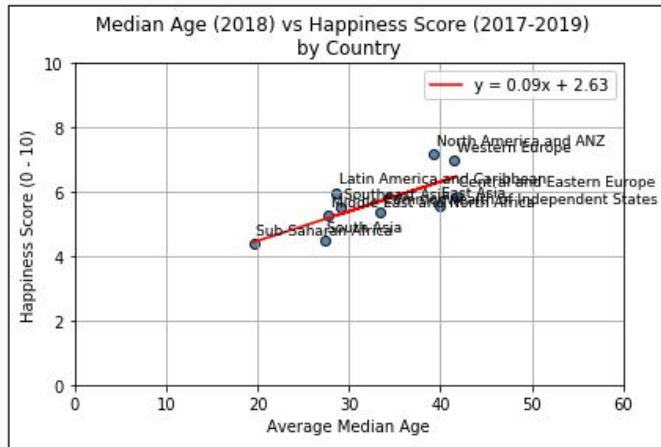
# Findings & Conclusion

- I expected to find stronger correlations than what I did.

- I believe there is a moderate correlation between the World Happiness Score and 3 out of the 4 variables, which were: Fertility Rate, Median Age, and Urbanization Population Rate.

- The correlation between World Happiness Score and Population Growth Rate was the weakest out of the 4 variables, though still semi-moderate.

# Final Thoughts

- Difficulties that arose mainly came during the plotting of the bar charts (as the region names were very long and overlapped each other) and displaying the data in a way that makes sense to others. Another issue came during comparison of variables with different years (though they overlapped). The data analysis might not be as accurate as I would've hoped for because of this.

- I also tried very hard to make scatter plots by region with the names of the region next to each plot, but the labels overlapped too much and wasn't readable (as shown to the right).

- Overall, this was an interesting topic to analyze. It made me curious about the endless possibility of variables I could compare against.



Example of code I would try harder to crack next time

THANK YOU