

Project 2 Analyzing World Happiness Dataset

The "World Happiness Report" consist of six factor in measuring happiness around the world. The first column records the overall rank of the country according to there happiness score in a descending manner. Second column records the region and country. The third column is the happiness score. The remaining column are GDP, Social support, life expectancy, generosity, and corruption.

Context

The data from the world happiness report is collected through the Gallup World Poll. The participants are asked to rate there life on a scale of 1 to 10. The dataset utilize in this instance are the national representative sample from the year 2015 to 2017. Each column represent the extent of its influence toward life evaluation higher than they are in a Dystopia, a theoretical country consist of the world's lowest national averages of each of the six factors.

Aim

The aim for this project would be to answer the following question, which country is the happiest , which factor affects the most on nation happiness, and is happiness in general increasing or decreasing ?

```
In [1]: import pandas as pd
import seaborn as sns
data3 = pd.read_csv('2015.csv')
data2 = pd.read_csv('2016.csv')
data3 = pd.read_csv('2017.csv')
print(data3.isnull().i())
print(data2.isnull().i())
print(data3.isnull().i())

Country Region Happiness Rank Happiness Score Standard Error \
0 False False False False False
1 False False False False False
2 False False False False False
3 False False False False False
4 False False False False False
..
153 False False False False False
154 False False False False False
155 False False False False False
156 False False False False False
157 False False False False False

Economy (GDP per Capita) Family Health (Life Expectancy) Freedom \
0 False False False False False
1 False False False False False
2 False False False False False
3 False False False False False
4 False False False False False
..
153 False False False False False
154 False False False False False
155 False False False False False
156 False False False False False
157 False False False False False

Trust (Government Corruption) Generosity Dystopia Residual
0 False False False False
1 False False False False
2 False False False False
3 False False False False
4 False False False False
..
153 False False False False
154 False False False False
155 False False False False
156 False False False False
157 False False False False

[158 rows x 12 columns]

Country Region Happiness Rank Happiness Score \
0 False False False False
1 False False False False
2 False False False False
3 False False False False
4 False False False False
..
153 False False False False
154 False False False False
155 False False False False
156 False False False False
157 False False False False

Lower Confidence Interval Upper Confidence Interval \
0 False False False False
1 False False False False
2 False False False False
3 False False False False
4 False False False False
..
153 False False False False
154 False False False False
155 False False False False
156 False False False False
157 False False False False

Economy (GDP per Capita) Family Health (Life Expectancy) Freedom \
0 False False False False False
1 False False False False False
2 False False False False False
3 False False False False False
4 False False False False False
..
153 False False False False False
154 False False False False False
155 False False False False False
156 False False False False False
157 False False False False False

Trust (Government Corruption) Generosity Dystopia Residual
0 False False False False
1 False False False False
2 False False False False
3 False False False False
4 False False False False
..
153 False False False False
154 False False False False
155 False False False False
156 False False False False
157 False False False False

[157 rows x 13 columns]

Country Happiness Rank Happiness Score Whisker.high Whisker.low \
0 False False False False False
1 False False False False False
2 False False False False False
3 False False False False False
4 False False False False False
..
150 False False False False False
151 False False False False False
152 False False False False False
153 False False False False False
154 False False False False False

Economy (GDP per Capita) Family Health (Life Expectancy) Freedom \
0 False False False False False
1 False False False False False
2 False False False False False
3 False False False False False
4 False False False False False
..
150 False False False False False
151 False False False False False
152 False False False False False
153 False False False False False
154 False False False False False

Generosity Trust (Government Corruption) Dystopia Residual
0 False False False False
1 False False False False
2 False False False False
3 False False False False
4 False False False False
..
150 False False False False
151 False False False False
152 False False False False
153 False False False False
154 False False False False

[155 rows x 12 columns]
```

Finally we start off by calling for the necessary tools which we'll use later for further analysis. Also I check one more time to make sure there is no null or NA value hidden in the dataset. Which we can see there is no missing value and the data is very complete all around.

```
In [2]: data3.head()

Out[2]: Country Happiness.Rank Happiness.Score Whisker.high Whisker.low Economy.GDP.per.Capita Family Health.Life.Expectancy Freedom Generosity
0 Denmark 2 7.522 7.861728 7.162272 1.403933 1.551122 0.709566 0.650623 0.350230
1 Iceland 3 7.904 7.622000 7.389970 1.400033 1.610974 0.833562 0.627163 0.479540
2 Switzerland 4 7.494 7.561772 7.426227 1.564980 1.510912 0.856131 0.620071 0.290549
3 Finland 5 7.469 7.527542 7.414568 1.445572 1.540247 0.809158 0.617951 0.245483
```

```
In [3]: data3.rename(columns={'Happiness.Rank': 'Happiness.Rank'}, inplace=True)
data3.rename(columns={'Happiness.Score': 'Happiness.Score'}, inplace=True)
data3.rename(columns={'Economy.GDP.per.Capita': 'Economy (GDP per Capita)'}, inplace=True)
data3.rename(columns={'Health.Life.Expectancy': 'Health (Life Expectancy)'}, inplace=True)
data3.rename(columns={'Trust (Government Corruption)': 'Trust (Government Corruption)'}, inplace=True)
data3.rename(columns={'Dystopia.Residual': 'Dystopia.Residual'}, inplace=True)
```

Through further speculation I realize there's a difference in the column name for data after 2016, therefore I rename some of the column name for easier manipulation in the later parts of my code.

```
In [4]: import numpy as np
import matplotlib.pyplot as plt
fig, graph = plt.subplots(nrows=7, ncols=3, figsize=(28, 40))

#Economy vs Score
#2015
graph[0][0].scatter(data3['Economy (GDP per Capita)'], data3['Happiness Score'])
graph[0][0].set_title('2015 GDP vs Score')
graph[0][0].set_xlabel('Economy (GDP per Capita)')
graph[0][0].set_ylabel('Happiness Score')
fit = np.polyfit(data3['Economy (GDP per Capita)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Economy (GDP per Capita)'])
graph[0][0].plot(data3['Economy (GDP per Capita)'], line, color='red')

#2016
graph[0][1].scatter(data3['Economy (GDP per Capita)'], data3['Happiness Score'])
graph[0][1].set_title('2016 GDP vs Score')
graph[0][1].set_xlabel('Economy (GDP per Capita)')
graph[0][1].set_ylabel('Happiness Score')
fit = np.polyfit(data3['Economy (GDP per Capita)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Economy (GDP per Capita)'])
graph[0][1].plot(data3['Economy (GDP per Capita)'], line, color='red')

#2017
graph[0][2].scatter(data3['Economy (GDP per Capita)'], data3['Happiness Score'])
graph[0][2].set_title('2017 GDP vs Score')
graph[0][2].set_xlabel('Economy (GDP per Capita)')
graph[0][2].set_ylabel('Happiness Score')
fit = np.polyfit(data3['Economy (GDP per Capita)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Economy (GDP per Capita)'])
graph[0][2].plot(data3['Economy (GDP per Capita)'], line, color='red')

#Family vs Score
#2015
graph[1][0].scatter(data3['Family'], data3['Happiness Score'])
graph[1][0].set_title('2015 Scatter Plot of Family vs Score')
graph[1][0].set_xlabel('Family')
graph[1][0].set_ylabel('Happiness score')
fit = np.polyfit(data3['Family'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Family'])
graph[1][0].plot(data3['Family'], line, color='red')

#2016
graph[1][1].scatter(data3['Family'], data3['Happiness Score'])
graph[1][1].set_title('2016 Scatter Plot of Family vs Score')
graph[1][1].set_xlabel('Family')
graph[1][1].set_ylabel('Happiness score')
fit = np.polyfit(data3['Family'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Family'])
graph[1][1].plot(data3['Family'], line, color='red')

#2017
graph[1][2].scatter(data3['Family'], data3['Happiness Score'])
graph[1][2].set_title('2017 Scatter Plot of Family vs Score')
graph[1][2].set_xlabel('Family')
graph[1][2].set_ylabel('Happiness score')
fit = np.polyfit(data3['Family'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Family'])
graph[1][2].plot(data3['Family'], line, color='red')

#Health (Life Expectancy) vs Score
#2015
graph[2][0].scatter(data3['Health (Life Expectancy)'], data3['Happiness Score'])
graph[2][0].set_title('2015 of Health vs Score')
graph[2][0].set_xlabel('Health (Life Expectancy)')
graph[2][0].set_ylabel('Happiness score')
fit = np.polyfit(data3['Health (Life Expectancy)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Health (Life Expectancy)'])
graph[2][0].plot(data3['Health (Life Expectancy)'], line, color='red')

#2016
graph[2][1].scatter(data3['Health (Life Expectancy)'], data3['Happiness Score'])
graph[2][1].set_title('2016 of Health vs Score')
graph[2][1].set_xlabel('Health (Life Expectancy)')
graph[2][1].set_ylabel('Happiness score')
fit = np.polyfit(data3['Health (Life Expectancy)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Health (Life Expectancy)'])
graph[2][1].plot(data3['Health (Life Expectancy)'], line, color='red')

#2017
graph[2][2].scatter(data3['Health (Life Expectancy)'], data3['Happiness Score'])
graph[2][2].set_title('2017 of Health vs Score')
graph[2][2].set_xlabel('Health (Life Expectancy)')
graph[2][2].set_ylabel('Happiness score')
fit = np.polyfit(data3['Health (Life Expectancy)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Health (Life Expectancy)'])
graph[2][2].plot(data3['Health (Life Expectancy)'], line, color='red')

#Freedom vs Score
graph[3][0].scatter(data3['Freedom'], data3['Happiness Score'])
graph[3][0].set_title('2015 of Freedom vs Score')
graph[3][0].set_xlabel('Freedom (Life Expectancy)')
graph[3][0].set_ylabel('Happiness score')
fit = np.polyfit(data3['Freedom'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Freedom'])
graph[3][0].plot(data3['Freedom'], line, color='red')

graph[3][1].scatter(data3['Freedom'], data3['Happiness Score'])
graph[3][1].set_title('2016 of Freedom vs Score')
graph[3][1].set_xlabel('Freedom (Life Expectancy)')
graph[3][1].set_ylabel('Happiness score')
fit = np.polyfit(data3['Freedom'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Freedom'])
graph[3][1].plot(data3['Freedom'], line, color='red')

graph[3][2].scatter(data3['Freedom'], data3['Happiness Score'])
graph[3][2].set_title('2017 of Freedom vs Score')
graph[3][2].set_xlabel('Freedom (Life Expectancy)')
graph[3][2].set_ylabel('Happiness score')
fit = np.polyfit(data3['Freedom'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Freedom'])
graph[3][2].plot(data3['Freedom'], line, color='red')

#Trust (Government Corruption) vs Score
graph[4][0].scatter(data3['Trust (Government Corruption)'], data3['Happiness Score'])
graph[4][0].set_title('2015 of Trust (Government Corruption) vs Score')
graph[4][0].set_xlabel('Trust (Government Corruption)')
graph[4][0].set_ylabel('Happiness score')
fit = np.polyfit(data3['Trust (Government Corruption)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Trust (Government Corruption)'])
graph[4][0].plot(data3['Trust (Government Corruption)'], line, color='red')

graph[4][1].scatter(data3['Trust (Government Corruption)'], data3['Happiness Score'])
graph[4][1].set_title('2016 of Trust (Government Corruption) vs Score')
graph[4][1].set_xlabel('Trust (Government Corruption)')
graph[4][1].set_ylabel('Happiness score')
fit = np.polyfit(data3['Trust (Government Corruption)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Trust (Government Corruption)'])
graph[4][1].plot(data3['Trust (Government Corruption)'], line, color='red')

graph[4][2].scatter(data3['Trust (Government Corruption)'], data3['Happiness Score'])
graph[4][2].set_title('2017 of Trust (Government Corruption) vs Score')
graph[4][2].set_xlabel('Trust (Government Corruption)')
graph[4][2].set_ylabel('Happiness score')
fit = np.polyfit(data3['Trust (Government Corruption)'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Trust (Government Corruption)'])
graph[4][2].plot(data3['Trust (Government Corruption)'], line, color='red')

#Generosity vs Score
graph[5][0].scatter(data3['Generosity'], data3['Happiness Score'])
graph[5][0].set_title('2015 of Generosity vs Score')
graph[5][0].set_xlabel('Generosity')
graph[5][0].set_ylabel('Happiness score')
fit = np.polyfit(data3['Generosity'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Generosity'])
graph[5][0].plot(data3['Generosity'], line, color='red')

graph[5][1].scatter(data3['Generosity'], data3['Happiness Score'])
graph[5][1].set_title('2016 of Generosity vs Score')
graph[5][1].set_xlabel('Generosity')
graph[5][1].set_ylabel('Happiness score')
fit = np.polyfit(data3['Generosity'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Generosity'])
graph[5][1].plot(data3['Generosity'], line, color='red')

graph[5][2].scatter(data3['Generosity'], data3['Happiness Score'])
graph[5][2].set_title('2017 of Generosity vs Score')
graph[5][2].set_xlabel('Generosity')
graph[5][2].set_ylabel('Happiness score')
fit = np.polyfit(data3['Generosity'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Generosity'])
graph[5][2].plot(data3['Generosity'], line, color='red')

#Dystopia Residual vs Score
graph[6][0].scatter(data3['Dystopia Residual'], data3['Happiness Score'])
graph[6][0].set_title('2015 of Dystopia Residual vs Score')
graph[6][0].set_xlabel('Dystopia Residual')
graph[6][0].set_ylabel('Happiness score')
fit = np.polyfit(data3['Dystopia Residual'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Dystopia Residual'])
graph[6][0].plot(data3['Dystopia Residual'], line, color='red')

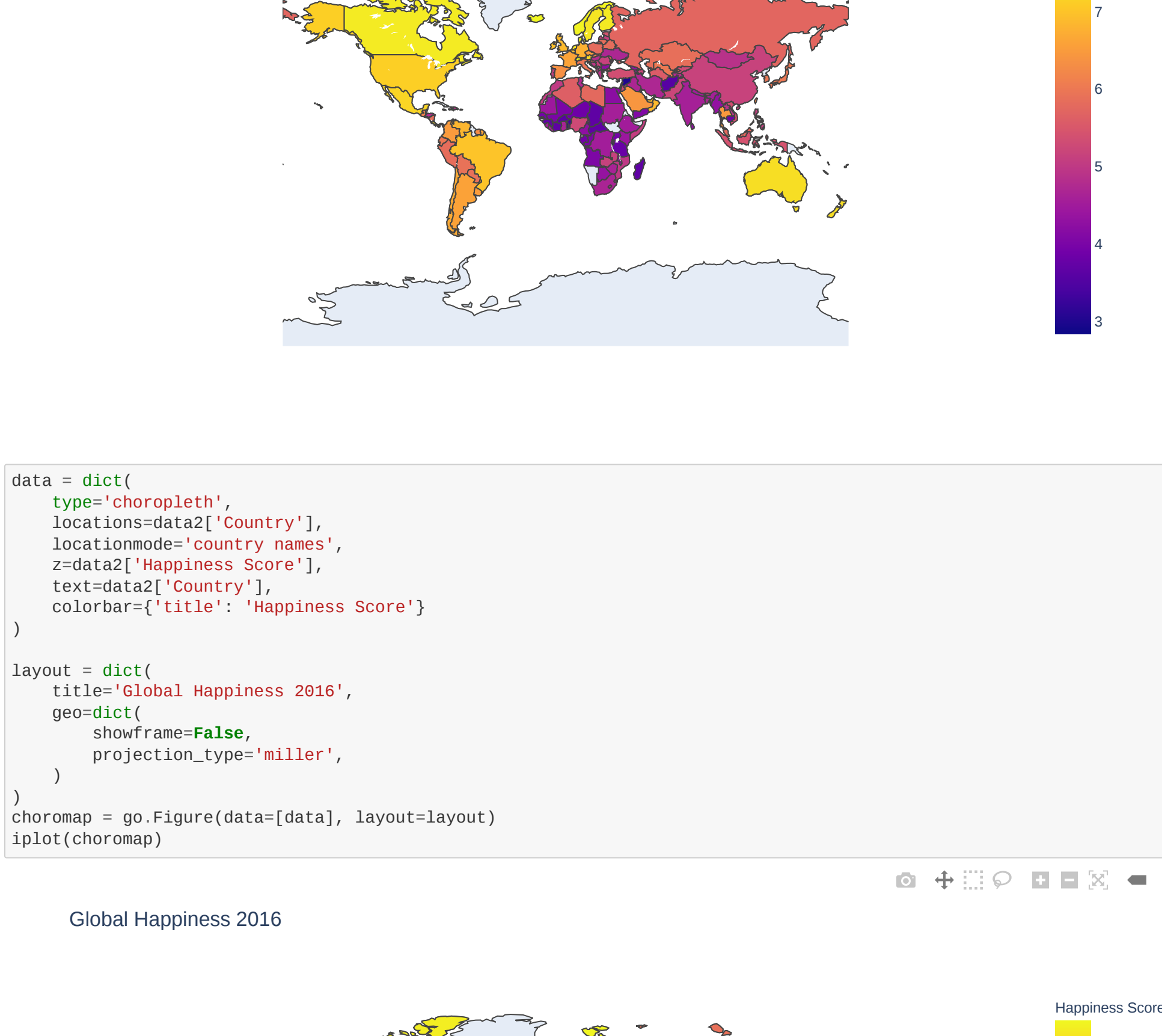
graph[6][1].scatter(data3['Dystopia Residual'], data3['Happiness Score'])
graph[6][1].set_title('2016 of Dystopia Residual vs Score')
graph[6][1].set_xlabel('Dystopia Residual')
graph[6][1].set_ylabel('Happiness score')
fit = np.polyfit(data3['Dystopia Residual'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Dystopia Residual'])
graph[6][1].plot(data3['Dystopia Residual'], line, color='red')

graph[6][2].scatter(data3['Dystopia Residual'], data3['Happiness Score'])
graph[6][2].set_title('2017 of Dystopia Residual vs Score')
graph[6][2].set_xlabel('Dystopia Residual')
graph[6][2].set_ylabel('Happiness score')
fit = np.polyfit(data3['Dystopia Residual'], data3['Happiness Score'], 1)
line = np.polyval(fit, data3['Dystopia Residual'])
graph[6][2].plot(data3['Dystopia Residual'], line, color='red')

plt.tight_layout()
plt.show()
```

Here I utilize a scatterplot to see the relationship between happiness score and the different categories of the dataset. A interesting trend is the correlation of GDP and Family with the Happiness score. The higher the happiness score the higher the GDP and Family score.

```
In [5]: sns.heatmap(data3.corr(), vmax=.8, contour=True)
plt.show()
```

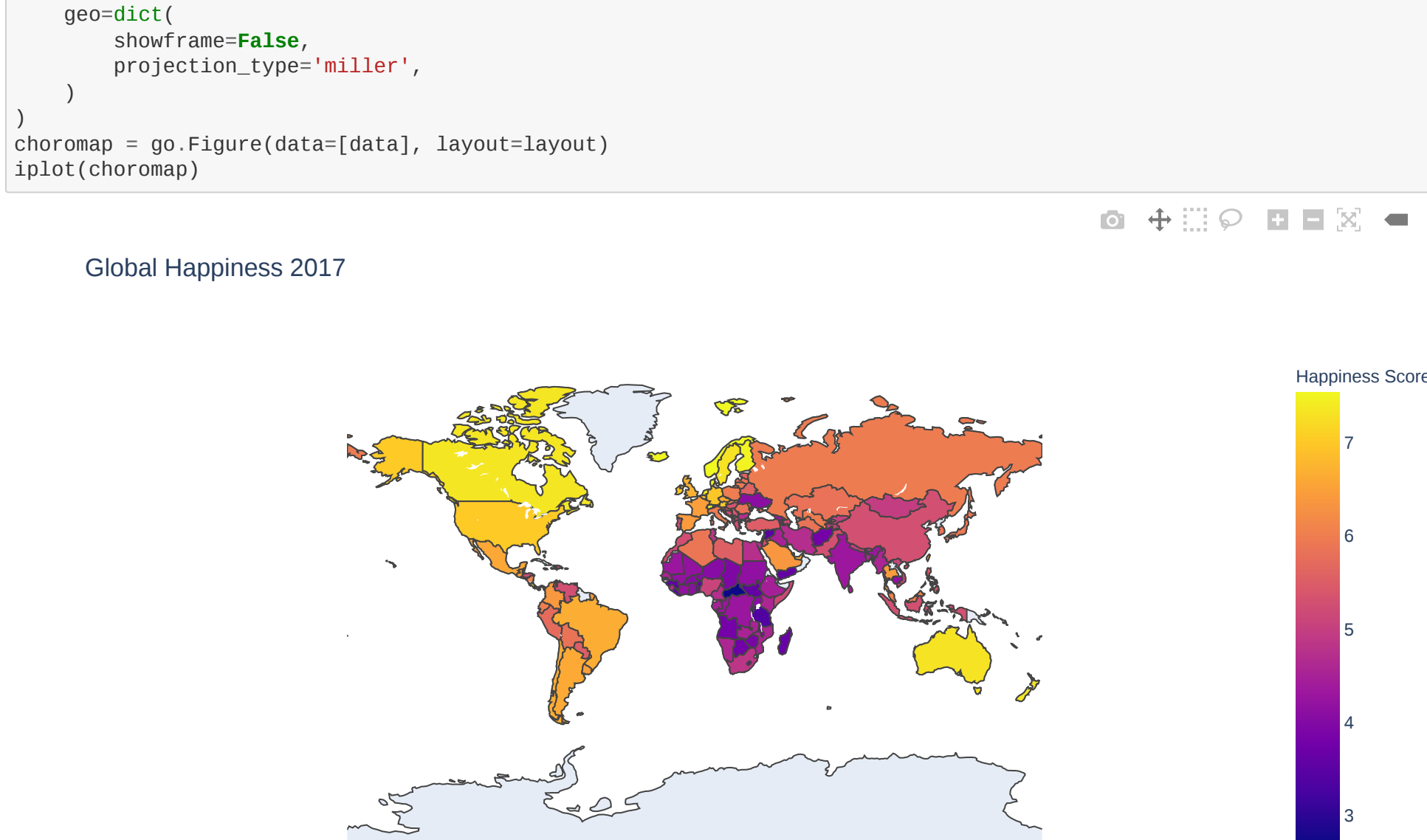


Here utilizing the `corr()` function I've deeper into the correlation each category have with the other. Which we can see the a near perfect score between GDP, Family, Life Expectancy and Happiness score.

```
In [6]: import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd

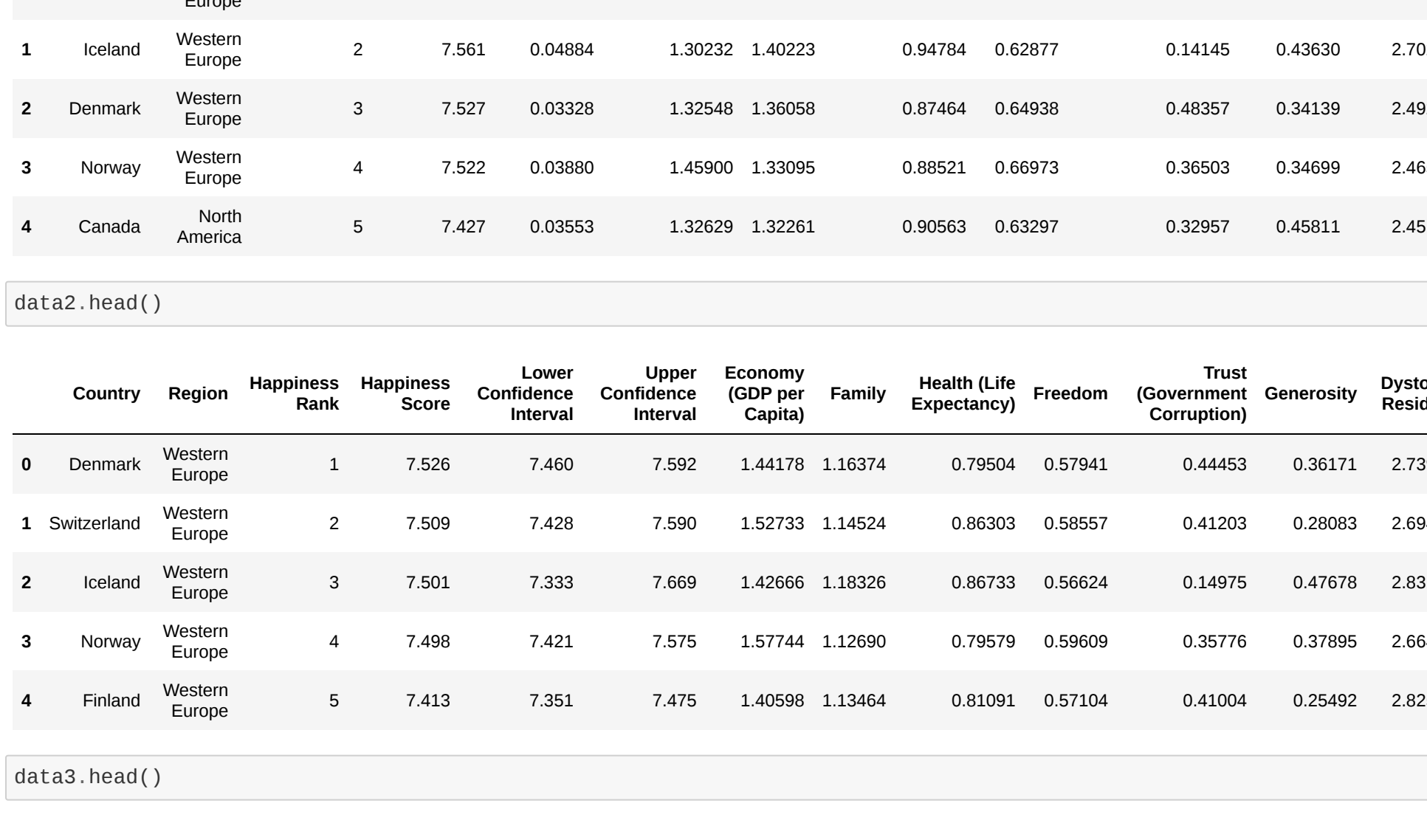
data = dict(
    type='choropleth',
    locations=data3['Country'],
    locationsmode='country names',
    z=data3['Happiness Score'],
    text=data3['Country'],
    colorbar={'title': 'Happiness Score'})

layout = dict(
    title='Global Happiness 2015',
    geo=dict(
        showframe=False,
        projection_type='miller',
    ),
)
choromap = go.Figure(data=[data], layout=layout)
iplot(choromap)
```



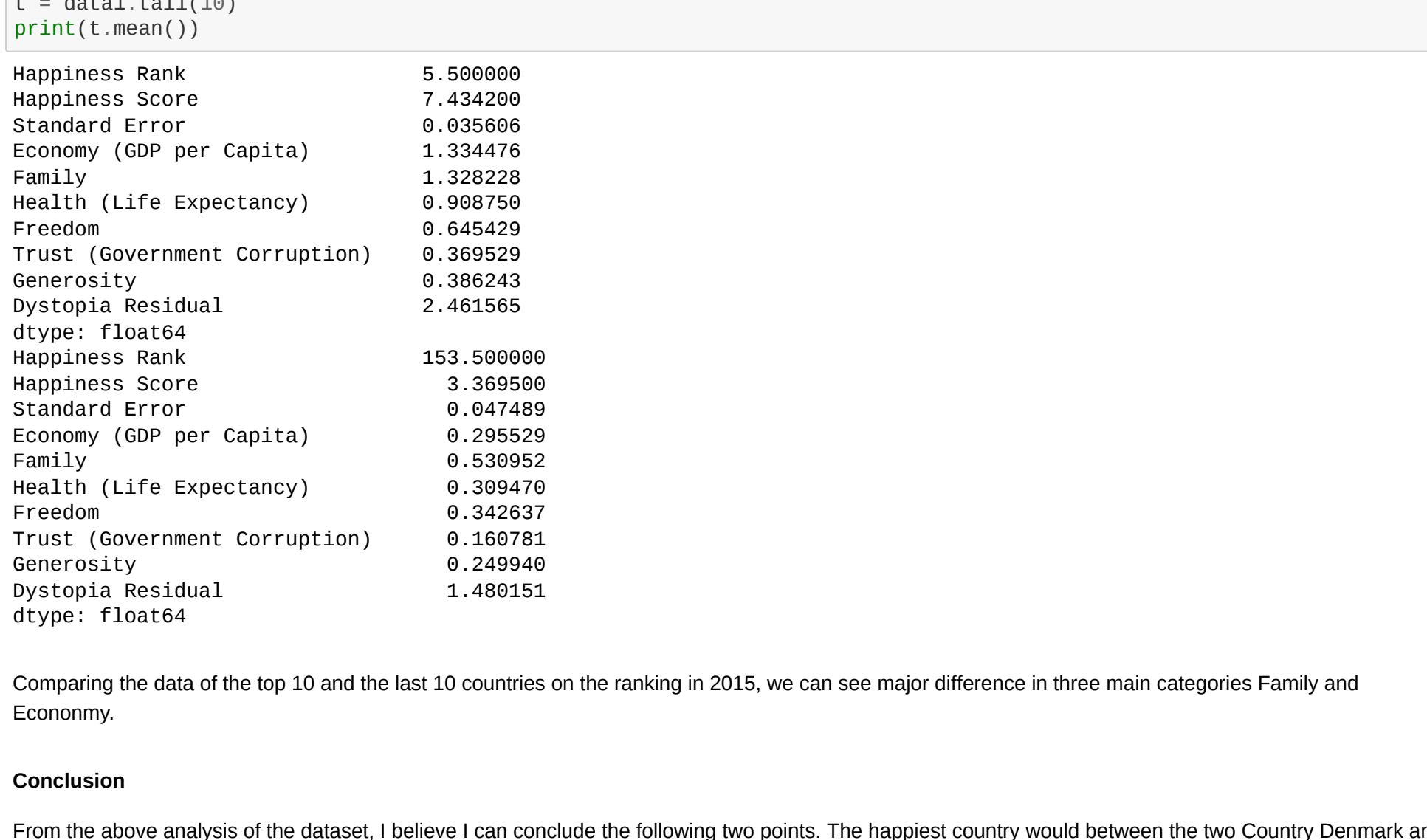
```
In [7]: data = dict(
    type='choropleth',
    locations=data3['Country'],
    locationsmode='country names',
    z=data3['Happiness Score'],
    text=data3['Country'],
    colorbar={'title': 'Happiness Score'})

layout = dict(
    title='Global Happiness 2016',
    geo=dict(
        showframe=False,
        projection_type='miller',
    ),
)
choromap = go.Figure(data=[data], layout=layout)
iplot(choromap)
```



```
In [8]: data = dict(
    type='choropleth',
    locations=data3['Country'],
    locationsmode='country names',
    z=data3['Happiness Score'],
    text=data3['Country'],
    colorbar={'title': 'Happiness Score'})

layout = dict(
    title='Global Happiness 2017',
    geo=dict(
        showframe=False,
        projection_type='miller',
    ),
)
choromap = go.Figure(data=[data], layout=layout)
iplot(choromap)
```



Through the three choropleth I was able to find which states have the highest score. Which most of them are Canada, Australia, Finland, Norway, Iceland. All of these countries have a high GDP which may mean there are more factors which affecting the happiness scores.

```
In [10]: data3.head()

Out[10]: Country Region Happiness Rank Happiness Score Standard Error Economy (GDP per Capita) Family Health (Life Expectancy) Freedom Trust (Government Corruption) Generosity Dystopia Residual
0 Switzerland Western Europe 1 7.587 0.03411 1.38951 1.34951 0.94143 0.60557 0.41978 0.29878 2.51738
1 Denmark Western Europe 2 7.561 0.04084 1.30232 1.40223 0.94784 0.62877 0.41415 0.43630 2.70201
2 Iceland Western Europe 3 7.951 0.03028 1.32548 1.36058 0.87464 0.64028 0.44357 0.34139 2.49204
3 Norway Western Europe 4 7.522 0.03880 1.45800 1.31055 0.85521 0.66973 0.35053 0.34659 2.45531
4 Canada North America 5 7.427 0.03953 1.32629 1.32661 0.95653 0.63297 0.32957 0.45811 2.45176
```

```
In [11]: data3.head()

Out[11]: Country Region Happiness Rank Happiness Score Lower Confidence Interval Upper Confidence Interval Economy (GDP per Capita) Family Health (Life Expectancy) Freedom Trust (Government Corruption) Generosity Dystopia Residual
0 Denmark Western Europe 1 7.526 7.465 7.582 1.44118 1.53734 0.795264 0.57944 0.44453 0.36211 2.73839
1 Switzerland Western Europe 2 7.509 7.428 7.590 1.52733 1.44524 0.86363 0.58557 0.41203 0.28083 2.69463
2 Iceland Western Europe 3 7.951 7.733 7.669 1.42666 1.38280 0.86733 0.56624 0.44769 0.28137
3 Finland Western Europe 4 7.458 7.421 7.475 1.57744 1.12926 0.79579 0.59609 0.38776 0.37895 2.66455
4 Norway Western Europe 5 7.413 7.351 7.475 1.40599 1.13454 0.81591 0.67104 0.41004 0.25482 2.82590
```

```
In [12]: data3.head()

Out[12]: Country Happiness Rank Happiness Score Whisker.high Whisker.low Economy (GDP per Capita) Family Health (Life Expectancy) Freedom Generosity Trust (Government Corruption) Dystopia Residual
0 Norway 1 7.537 7.594445 7.479556 1.615483 1.533524 0.796667 0.65423 0.350212 0.319564 2.277027
1 Denmark 2 7.522 7.581728 7.462272 1.402383 1.551122 0.709566 0.620007 0.350230 0.400770 2.313707
2 Iceland 3 7.951 7.622000 7.389970 1.400033 1.610974 0.833562 0.627163 0.479540 0.153827 2.227215
3 Switzerland 4 7.494 7.561772 7.426227 1.564980 1.510912 0.856131 0.620071 0.290549 0.367007 2.272716
4 Finland 5 7.469 7.527542 7.414568 1.445572 1.540247 0.809158 0.617951 0.245483 0.382602 2.430182
```

Through observing the data of the top five countries from the year 2015 to 2017, we can see Western Europe are always in the top five. In the top five, Denmark and Iceland being the long-standing place of number two and number three.

```
In [16]: h = data3.head(10)
print(h.mean())
print(h.std())

Happiness Rank 5.590000
Happiness Score 6.434209
Standard Error 0.035696
Economy (GDP per Capita) 1.334478
Family 1.228428
Health (Life Expectancy) 0.908759
Freedom 0.645429
Trust (Government Corruption) 0.369429
Generosity 0.380243
Dystopia Residual 2.499480
dtype: float64
```

Comparing the data of the top 10 and the last 10 countries on the ranking in 2015, we can see major difference in three main categories Family and Economy.

Conclusion

From the above analysis of the dataset, I believe I can conclude the following two points. The happiest country would be between the two Country Denmark and Iceland. Which have the highest happiness score above 7.5. Also Western Europe is definitely the happiest place to live in, if we see its high rank in the period from 2015-2017.

In regards to the factor which affects happiness the most, I believe it's completely certain there are only certain factors affecting the overall score. Yet from my analysis there are certain factors which seem noticeable when looking at countries with high happiness score over 7, which are GDP, Family, and Life expectancy. These three factors are always accompanied when a high ranking country on the happiness scale is mentioned.

Reference: <https://www.kaggle.com/datasets/steffenworld/happiness> <https://www.oecd.org/dataoecd/30/46/39206949.pdf> <https://data360.org/2020/09/08/2020-world-happiness-report/> <https://www.oecd.org/dataoecd/30/46/39206949.pdf> <https://data360.org/2020/09/08/2020-world-happiness-report/>