

World Happiness Report - Analysis

Introduction

For my data analysis, I decided to examine data from the World Happiness Report from 2015 to 2023. It is based on a global survey in which the population rates their own lives on a scale from 0 to 10. The reason I chose this data is that there are several key factors that influence this value. In my analysis, these values are as follows:

- 1. GDP per capita
- 2. Social support
- 3. Healthy life expectancy
- 4. Freedom to make life choices
- 5. Generosity
- 6. Perceptions of corruption

The data comes from various sources/years and is divided into different CSV files. These are getting merged and analyzed.

Research Questions

- 1. Which variables have the strongest influence on the happiness score?
- 2. Has the coronavirus pandemic had an impact on the happiness score?
- 3. Is it possible to predict a country's happiness score by using a Linear Regression model?
- 4. Are there regional differences in terms of the happiness score?

Key limitations

There are limitations in my data research that affect the quality of the data.

On the one hand, the perception of one's own well-being is very subjective. This means that a certain number of points can be interpreted differently (either as good or bad). If one person rates their quality of life as a "7," the same quality of life may be a "5" for another person. In addition, this data consists of six factors. The problem behind this is that important data is missing, such as environmental quality or inequality. I will address other limitations in the course of my research.

Data Preparation and Merging

To meet the requirements of my data analysis, I will first merge all tables. I do this by adding a new column in which I insert the year of the data record. The result is a CSV file with 1367 rows.

```
import pandas as pd
import glob
import os
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
file_paths = sorted(glob.glob("*.csv"))
all_dfs = []
for file in file_paths:
    temp_df = pd.read_csv(file)
    year_val = int(os.path.basename(file).split('.')[1].split('.')[0])
    temp_df['Year'] = year_val
    all_dfs.append(temp_df)
df = pd.concat(all_dfs, ignore_index=True)
df = df.fillna(df.mean(numeric_only=True))
display(df)
```

country	region	happiness_score	gdp_per_capita	social_support	healthy_life_expectancy	freedom_to_make_life_choices	generosity	perceptions_of_corruption	Year	
0	Switzerland	Western Europe	7.587	1.39651	1.34951	0.94143	0.66557	0.29678	2015	
1	Iceland	Western Europe	7.661	1.30232	1.40223	0.94784	0.62977	0.43630	0.14546	2015
2	Denmark	Western Europe	7.577	1.32548	1.36058	0.87464	0.64986	0.34139	0.48307	2015
3	Norway	Western Europe	7.522	1.49900	1.33095	0.88521	0.68973	0.34699	0.36003	2015
4	Canada	North America and ANZ	7.427	1.32629	1.32261	0.80563	0.62939	0.45811	0.30979	2015
...
1362	Congo (Kinshasa)	Sub-Saharan Africa	3.207	0.53100	0.79400	0.15050	0.37500	0.18300	0.06800	2023
1363	Zimbabwe	Sub-Saharan Africa	3.204	0.57800	0.88100	0.06800	0.36300	0.11200	0.11700	2023
1364	Sierra Leone	Sub-Saharan Africa	3.158	0.67000	0.54000	0.09200	0.37100	0.19300	0.05150	2023
1365	Lebanon	Middle East and North Africa	2.392	1.41700	0.47600	0.39800	0.32300	0.06150	0.07200	2023
1366	Afghanistan	South Asia	1.859	0.64500	0.00000	0.08700	0.00000	0.09300	0.05900	2023

1367 rows x 10 columns

```
In [1]: df.info()
```

To ensure memory usage, I have a file with 100K rows, but with large amounts of data, hardware performance can suffer. To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df = df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

General Statistic (Distribution and Average Happiness Score over the Years) + Key questions

The main topic of my data analysis is the "Happiness Score". Accordingly, I decided to take a closer look at this value first. The lower plots show the frequencies of all values of the happiness_score. From this, I can draw the following conclusions:

Mean (AVG): The global average across all years is 5.44. Range (Min/Max): The global distribution shows enormous differences. While most countries lie between approximately 4.3 and 6.5, there are also those that fall significantly outside the range. The gap between happy and unhappy ranges from Afghanistan in 2023 with a value of 1.86 to 7.84, which was recorded by Finland in 2021.

```
In [2]: df.describe()
```

```
happiness_score    5.4400000000000005
gdp_per_capita      10614.000000000001
social_support      0.6910000000000001
healthy_life_expectancy 78.74000000000001
freedom_to_make_life_choices 0.4610000000000001
generosity          0.12800000000000002
perceptions_of_corruption 0.12800000000000002
Year                2015.00000000000001
```

memory usage: 106.1 MB + 100K = 106.1 MB

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1 MB.

To ensure that there are no errors have occurred when merging all the data. Since I now have 10 columns (0-9), this means that the new column for specifying the individual date record (Year) is included. This is also indicated in the output under "9". The so-called "non-null count" is also very relevant. This ensures that my merged file does not contain any missing data. To guarantee this, I first used "df.replace(np.nan,np.nan,only=True)" to fill in all missing values with the calculated average value. This allowed me to significantly improve my data quality, as otherwise my regression models would not work (NaN throws an error code). In addition, I can identify the data types of the individual columns (under "Dtypes"). "float64" stands for decimal numbers (happiness_score, gdp_per_capita...); "object" usually stands for text (country, region); and "int64" stands for integers (Year). Last but not least, you can see how much memory is used by the data in my example. 936.1 MB + 100K = 106.1