

✓ Introduction

This project investigates the evolving preferences of anime fans within the rapidly expanding global anime market. Once a niche interest, anime has achieved mainstream popularity, driven by its distinctive art styles and compelling narratives. The anime market, valued at USD 31.41 billion in 2023 and projected to reach USD 72.86 billion by 2032 (CAGR 9.8%), necessitates a deeper understanding of consumer trends. This research will analyze data from MyAnimeList, alongside with relevant economic and demographic data from World Bank, to identify key drivers of anime fan preferences. Specifically, this study will explore how factors such as gender, happiness score, country, age, and gdp per capita correlate with user ratings and engagement.

✓ Preperation

Most of the data we use for the machine learning is in google drive. This cell allows us to use the data on Google Drive

```
1 import os
2 from os.path import join
3 from google.colab import drive
4 drive.mount("/content/drive", force_remount=True)
5
6 # to be used a prefix for file I/O
7 path_prefix = "/content/drive/My Drive/Colab/DSA 210/Project/Data"
```

→ Mounted at /content/drive

This cell lists the files in a specified directory and creates two lists based on these filenames.

exceptions_2015 and exceptions_2014 created here will be used in data processing. The year is last year with data on the datalists.

```
1 worldbank_fname_list_csv = os.listdir(path_prefix)
2 worldbank_fname_list = [item[:-4] if item.endswith(".csv") else item for item in os.listdir(path_prefix)]
3
4 exceptions_2015 = ["Renewable electricity output (% of total electricity output).csv",]
5 exceptions_2014 = ["Electric power consumption (kWh per capita).csv"]
6 len(os.listdir(path_prefix)), len(worldbank_fname_list_csv), len(worldbank_fname_list)
```

→ (55, 55, 55)

Check to see if there is any missing data

```
1 n = os.listdir(path_prefix)
2 for item in n:
3     if item not in worldbank_fname_list_csv:
4         print(item)
```

This code block is used to set up authentication for accessing datasets from the Kaggle platform within the Google Colab environment.

```
1 from google.colab import userdata
2 os.environ["KAGGLE_KEY"] = userdata.get('KAGGLE_KEY')
3 os.environ["KAGGLE_USERNAME"] = userdata.get('KAGGLE_USERNAME')
```

This code block imports various Python libraries and modules that are essential for data manipulation, analysis, visualization, and machine learning tasks.

```
1 ### Basic libraries
2 import math
3 import datetime
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 %matplotlib inline
9
10 import pickle
11
12 # Data Preprocessing
13 from scipy.stats import chi2_contingency
14 from scipy import stats
15 from sklearn.preprocessing import LabelEncoder
16
```

```

~
17
18 # Model Training
19 from sklearn.utils import shuffle
20 from sklearn.model_selection import train_test_split, RandomizedSearchCV
21 from sklearn.metrics import mean_absolute_error, root_mean_squared_log_error
22 from sklearn.experimental import enable_halving_search_cv
23 from sklearn.model_selection import HalvingRandomSearchCV
24 from sklearn.model_selection import GridSearchCV
25 from sklearn.ensemble import RandomForestRegressor
26 from sklearn.linear_model import Ridge
27 from sklearn.svm import SVR
28 from sklearn.ensemble import GradientBoostingRegressor
29 from sklearn.linear_model import Lasso
30 from sklearn.linear_model import LinearRegression
31 import xgboost as xgb
32 import lightgbm as lgb
33 import re
34 from sklearn.tree import DecisionTreeRegressor

```

▼ Gathering Data

These code blocks are used to download datasets from Kaggle and then extract its contents.

```

1 !kaggle datasets download -d mathurinache/world-happiness-report
2
3 !unzip "world-happiness-report.zip"

```

→ Dataset URL: <https://www.kaggle.com/datasets/mathurinache/world-happiness-report>
 License(s): CC0-1.0
 Downloading world-happiness-report.zip to /content
 0% 0.00/67.4k [00:00<?, ?B/s]
 100% 67.4k/67.4k [00:00<00:00, 178MB/s]
 Archive: world-happiness-report.zip
 inflating: 2015.csv
 inflating: 2016.csv
 inflating: 2017.csv
 inflating: 2018.csv
 inflating: 2019.csv
 inflating: 2020.csv
 inflating: 2021.csv
 inflating: 2022.csv

```

1 !kaggle datasets download -d azathoth42/myanimelist
2
3 !unzip "myanimelist.zip"

```

→ Dataset URL: <https://www.kaggle.com/datasets/azathoth42/myanimelist>
 License(s): DbCL-1.0
 Downloading myanimelist.zip to /content
 100% 1.91G/1.92G [00:14<00:00, 243MB/s]
 100% 1.92G/1.92G [00:14<00:00, 147MB/s]
 Archive: myanimelist.zip
 inflating: Animelist.csv
 inflating: UserAnimelist.csv
 inflating: UserList.csv
 inflating: anime_cleaned.csv
 inflating: anime_filtered.csv
 inflating: animelists_cleaned.csv
 inflating: animelists_filtered.csv
 inflating: users_cleaned.csv
 inflating: users_filtered.csv

▼ Data Preprocessing

This code block focuses on loading and cleaning dataset of World Happiness Report, into a pandas DataFrame named df.

```

1 df = pd.read_csv('2020.csv')
2 df = df.rename(columns={"Country name": "Country"})
3 df.sort_values(by='Country', ascending=True, inplace=True)
4 df = df[df["Country"] != "Kosovo"].copy()
5 df = df.rename(columns={"Ladder score": "Happiness Score"})
6 df.reset_index(drop=True, inplace=True)
7 df['Country'] = df['Country'].str.lower()
8 df['Country'] = df['Country'].str.strip()
9 df.head()

```

	Country	Regional indicator	Happiness Score	Standard error of ladder score	upperwhisker	lowerwhisker	Logged GDP per capita	Social support	Healthy life expectancy	Freedom to make life choices	Generosity	Pi
0	afghanistan	South Asia	2.5669	0.031311	2.628270	2.505530	7.462861	0.470367	52.590000	0.396573	-0.096429	
1	albania	Central and Eastern Europe	4.8827	0.056116	4.992687	4.772713	9.417931	0.671070	68.708138	0.781994	-0.042309	
2	algeria	Middle East and North Africa	5.0051	0.044236	5.091802	4.918397	9.537965	0.803385	65.905174	0.466611	-0.121105	
3	argentina	Latin America and Caribbean	5.9747	0.053442	6.079446	5.869954	9.810955	0.900568	68.803802	0.831132	-0.194914	
4	armenia	Commonwealth of Independent States	4.6768	0.058595	4.791646	4.561953	9.100476	0.757479	66.750656	0.712018	-0.138780	

```
1 df.info()
```

```
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 152 entries, 0 to 151
Data columns (total 20 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Country         152 non-null    object  
 1   Regional indicator 152 non-null    object  
 2   Happiness Score  152 non-null    float64 
 3   Standard error of ladder score 152 non-null    float64 
 4   upperwhisker    152 non-null    float64 
 5   lowerwhisker    152 non-null    float64 
 6   Logged GDP per capita 152 non-null    float64 
 7   Social support   152 non-null    float64 
 8   Healthy life expectancy 152 non-null    float64 
 9   Freedom to make life choices 152 non-null    float64 
 10  Generosity       152 non-null    float64 
 11  Perceptions of corruption 152 non-null    float64 
 12  Ladder score in Dystopia 152 non-null    float64 
 13  Explained by: Log GDP per capita 152 non-null    float64 
 14  Explained by: Social support   152 non-null    float64 
 15  Explained by: Healthy life expectancy 152 non-null    float64 
 16  Explained by: Freedom to make life choices 152 non-null    float64 
 17  Explained by: Generosity     152 non-null    float64 
 18  Explained by: Perceptions of corruption 152 non-null    float64 
 19  Dystopia + residual    152 non-null    float64 
dtypes: float64(18), object(2)
memory usage: 23.9+ KB
```

This code block iterates through a list of World Bank data filenames and merges the relevant data from each file into the main DataFrame, named df. The goal is to enrich the data with various country-level indicators from the World Bank datasets.

```
1 rows_to_skip = 4
2 for i in range (len(worldbank_fname_list_csv)):
3     fname = worldbank_fname_list_csv[i]
4     df_world_bank = pd.read_csv(join(path_prefix, fname), skiprows=rows_to_skip, sep=',', on_bad_lines='skip')
5
6     if 'Country Name' in df_world_bank.columns:
7         df_world_bank['Country Name'] = df_world_bank['Country Name'].str.lower().str.strip()
8     else:
9         print(f"Warning: 'Country Name' column not found in {fname}")
10        continue
11
12    df_tmp = df[['Country']].copy()
13    if fname in exceptions_2015:
14        df_tmp = df_tmp.merge(df_world_bank[['Country Name', '2015']], left_on='Country', right_on='Country Name', how='left')
15        df[worldbank_fname_list[i]] = df_tmp['2015']
16    elif fname in exceptions_2014:
17        df_tmp = df_tmp.merge(df_world_bank[['Country Name', '2014']], left_on='Country', right_on='Country Name', how='left')
18        df[worldbank_fname_list[i]] = df_tmp['2014']
19    else:
20        df_tmp = df_tmp.merge(df_world_bank[['Country Name', '2020']], left_on='Country', right_on='Country Name', how='left')
21        df[worldbank_fname_list[i]] = df_tmp['2020']
22 df.info()
```

```
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 152 entries, 0 to 151
Data columns (total 75 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Country         152 non-null    object  
 1   Regional indicator 152 non-null    object  
 2   Happiness Score  152 non-null    float64 
 3   Standard error of ladder score 152 non-null    float64 
 4   upperwhisker    152 non-null    float64 
 5   lowerwhisker    152 non-null    float64 
 6   Logged GDP per capita 152 non-null    float64 
 7   Social support   152 non-null    float64 
 8   Healthy life expectancy 152 non-null    float64 
 9   Freedom to make life choices 152 non-null    float64 
 10  Generosity       152 non-null    float64 
 11  Perceptions of corruption 152 non-null    float64 
 12  Ladder score in Dystopia 152 non-null    float64 
 13  Explained by: Log GDP per capita 152 non-null    float64 
 14  Explained by: Social support   152 non-null    float64 
 15  Explained by: Healthy life expectancy 152 non-null    float64 
 16  Explained by: Freedom to make life choices 152 non-null    float64 
 17  Explained by: Generosity     152 non-null    float64 
 18  Explained by: Perceptions of corruption 152 non-null    float64 
 19  Dystopia + residual    152 non-null    float64 
 20  Country Name       152 non-null    object  
 21  2015               152 non-null    float64 
 22  2014               152 non-null    float64 
 23  2020               152 non-null    float64
```

```

-----      -----
0 Country          152 non-null  object
1 Regional indicator 152 non-null  object
2 Happiness Score   152 non-null  float64
3 Standard error of ladder score 152 non-null  float64
4 upperwhisker     152 non-null  float64
5 lowerwhisker     152 non-null  float64
6 Logged GDP per capita 152 non-null  float64
7 Social support    152 non-null  float64
8 Healthy life expectancy 152 non-null  float64
9 Freedom to make life choices 152 non-null  float64
10 Generosity       152 non-null  float64
11 Perceptions of corruption 152 non-null  float64
12 Ladder score in Dystopia 152 non-null  float64
13 Explained by: Log GDP per capita 152 non-null  float64
14 Explained by: Social support 152 non-null  float64
15 Explained by: Healthy life expectancy 152 non-null  float64
16 Explained by: Freedom to make life choices 152 non-null  float64
17 Explained by: Generosity 152 non-null  float64
18 Explained by: Perceptions of corruption 152 non-null  float64
19 Dystopia + residual 152 non-null  float64
20 High-technology exports (current US$) 117 non-null  float64
21 Forest area (% of land area) 130 non-null  float64
22 Government expenditure on education, total (% of GDP) 118 non-null  float64
23 GDP per capita growth (annual %) 129 non-null  float64
24 Urban population (% of total population) 130 non-null  float64
25 Net migration 130 non-null  float64
26 GNI per capita, PPP (current international $) 129 non-null  float64
27 Exports of goods and services (% of GDP) 123 non-null  float64
28 Agriculture, forestry, and fishing, value added (% of GDP) 129 non-null  float64
29 Literacy rate, adult total (% of people ages 15 and above) 28 non-null  float64
30 Access to electricity (% of population) 130 non-null  float64
31 Population growth (annual %) 130 non-null  float64
32 Renewable electricity output (% of total electricity output) 130 non-null  float64
33 Electric power consumption (kWh per capita) 113 non-null  float64
34 Central government debt, total (% of GDP) 51 non-null  float64
35 Renewable internal freshwater resources per capita (cubic meters) 129 non-null  float64
36 Population density (people per sq. km of land area) 130 non-null  float64
37 Time required to get electricity (days) 0 non-null  float64
38 Mobile cellular subscriptions (per 100 people) 130 non-null  float64
39 Fixed broadband subscriptions (per 100 people) 128 non-null  float64
40 Secure Internet servers (per 1 million people) 130 non-null  float64
41 Logistics performance index Overall (1=low to 5=high) 0 non-null  float64
42 Tax revenue (% of GDP) 103 non-null  float64
43 Mortality rate, under-5 (per 1,000 live births) 130 non-null  float64
44 Employment to population ratio, 15+, total (%) (modeled ILO estimate) 130 non-null  float64
45 Death rate, crude (per 1,000 people) 130 non-null  float64
46 Poverty headcount ratio at $2.15 a day (2017 PPP) (% of population) 55 non-null  float64
47 PM2.5 air pollution, mean annual exposure (micrograms per cubic meter) 130 non-null  float64
48 Agricultural land (% of land area) 130 non-null  float64
49 Rural population (% of total population) 130 non-null  float64
50 Income share held by lowest 20% 55 non-null  float64
51 Maternal mortality ratio (modeled estimate, per 100,000 live births) 130 non-null  float64
52 CPIA public sector management and institutions cluster average (1=low to 6=high) 41 non-null  float64

```

```
1 df.head()
```



	Country	Regional indicator	Happiness Score	Standard error of ladder score	upperwhisker	lowerwhisker	Logged GDP per capita	Social support	Healthy life expectancy	Freedom to make life choices	... Contraceptive prevalence any method of married women
0	afghanistan	South Asia	2.5669	0.031311	2.628270	2.505530	7.462861	0.470367	52.590000	0.396573	...
1	albania	Central and Eastern Europe	4.8827	0.056116	4.992687	4.772713	9.417931	0.671070	68.708138	0.781994	...
2	algeria	Middle East and North Africa	5.0051	0.044236	5.091802	4.918397	9.537965	0.803385	65.905174	0.466611	...
3	argentina	Latin America and Caribbean	5.9747	0.053442	6.079446	5.869954	9.810955	0.900568	68.803802	0.831132	...
4	armenia	Commonwealth of Independent States	4.6768	0.058595	4.791646	4.561953	9.100476	0.757479	66.750656	0.712018	...

5 rows × 75 columns

```
1 df_world_bank.head(5)
```

	Country Name	Country Code	Indicator Name	Indicator Code	1960	1961	1962	1963	1964	1965	...	2015	2016
0	aruba	ABW	Fertility rate, total (births per woman)	SP.DYN.TFR.T.IN	4.820000	4.655000	4.471000	4.271000	4.059000	3.842000	...	1.972000	1.953000
1	africa eastern and southern	AFE	Fertility rate, total (births per woman)	SP.DYN.TFR.T.IN	6.723226	6.742124	6.762343	6.778133	6.787821	6.799767	...	4.678619	4.616540
2	afghanistan	AFG	Fertility rate, total (births per woman)	SP.DYN.TFR.T.IN	7.282000	7.284000	7.292000	7.302000	7.304000	7.305000	...	5.405000	5.262000
			Fertility										

This code block iterates through each column in the df DataFrame to identify and handle missing values (represented as NaN). Every column that have more than 50 missing column is dropped.

```

1 for col in df.columns:
2     if df[col].isnull().any():
3         nan_sum = df[col].isnull().sum()
4         print(f"Column '{col}' has {nan_sum} NaN values.")
5         if nan_sum > 50:
6             df.drop(columns=[col], inplace=True)
7             print(f"Column '{col}' has been dropped.")

```

```

→ Column 'High-technology exports (current US$)' has 35 NaN values.
Column 'Forest area (% of land area)' has 22 NaN values.
Column 'Government expenditure on education, total (% of GDP)' has 34 NaN values.
Column 'GDP per capita growth (annual %)' has 23 NaN values.
Column 'Urban population (% of total population)' has 22 NaN values.
Column 'Net migration' has 22 NaN values.
Column 'GNI per capita, PPP (current international $)' has 23 NaN values.
Column 'Exports of goods and services (% of GDP)' has 29 NaN values.
Column 'Agriculture, forestry, and fishing, value added (% of GDP)' has 23 NaN values.
Column 'Literacy rate, adult total (% of people ages 15 and above)' has 124 NaN values.
Column 'Literacy rate, adult total (% of people ages 15 and above)' has been dropped.
Column 'Access to electricity (% of population)' has 22 NaN values.
Column 'Population growth (annual %)' has 22 NaN values.
Column 'Renewable electricity output (% of total electricity output)' has 22 NaN values.
Column 'Electric power consumption (kWh per capita)' has 39 NaN values.
Column 'Central government debt, total (% of GDP)' has 101 NaN values.
Column 'Central government debt, total (% of GDP)' has been dropped.
Column 'Renewable internal freshwater resources per capita (cubic meters)' has 23 NaN values.
Column 'Population density (people per sq. km of land area)' has 22 NaN values.
Column 'Time required to get electricity (days)' has 152 NaN values.
Column 'Time required to get electricity (days)' has been dropped.
Column 'Mobile cellular subscriptions (per 100 people)' has 22 NaN values.
Column 'Fixed broadband subscriptions (per 100 people)' has 24 NaN values.
Column 'Secure Internet servers (per 1 million people)' has 22 NaN values.
Column 'Logistics performance index Overall (1=low to 5=high)' has 152 NaN values.
Column 'Logistics performance index Overall (1=low to 5=high)' has been dropped.
Column 'Tax revenue (% of GDP)' has 49 NaN values.
Column 'Mortality rate, under-5 (per 1,000 live births)' has 22 NaN values.
Column 'Employment to population ratio, 15+, total (%) (modeled ILO estimate)' has 22 NaN values.
Column 'Death rate, crude (per 1,000 people)' has 22 NaN values.
Column 'Poverty headcount ratio at $2.15 a day (2017 PPP) (% of population)' has 97 NaN values.
Column 'Poverty headcount ratio at $2.15 a day (2017 PPP) (% of population)' has been dropped.
Column 'PM2.5 air pollution, mean annual exposure (micrograms per cubic meter)' has 22 NaN values.
Column 'Agricultural land (% of land area)' has 22 NaN values.
Column 'Rural population (% of total population)' has 22 NaN values.
Column 'Income share held by lowest 20%' has 97 NaN values.
Column 'Income share held by lowest 20%' has been dropped.
Column 'Maternal mortality ratio (modeled estimate, per 100,000 live births)' has 22 NaN values.
Column 'CPIA public sector management and institutions cluster average (1=low to 6=high)' has 111 NaN values.
Column 'CPIA public sector management and institutions cluster average (1=low to 6=high)' has been dropped.
Column 'Population in urban agglomerations of more than 1 million (% of total population)' has 52 NaN values.
Column 'Population in urban agglomerations of more than 1 million (% of total population)' has been dropped.
Column 'School enrollment, primary and secondary (gross), gender parity index (GPI)' has 70 NaN values.
Column 'School enrollment, primary and secondary (gross), gender parity index (GPI)' has been dropped.
Column 'Labor force, female (% of total labor force)' has 22 NaN values.
Column 'Literacy rate, youth total (% of people ages 15-24)' has 121 NaN values.
Column 'Literacy rate, youth total (% of people ages 15-24)' has been dropped.
Column 'Population in the largest city (% of urban population)' has 33 NaN values.
Column 'New businesses registered (number)' has 50 NaN values.
Column 'Strength of legal rights index (0=weak to 12=strong)' has 152 NaN values.
Column 'Strength of legal rights index (0=weak to 12=strong)' has been dropped.
Column 'Firms with female participation in ownership (% of firms)' has 142 NaN values.
Column 'Firms with female participation in ownership (% of firms)' has been dropped.
Column 'Cause of death, by injury (% of total)' has 152 NaN values.
Column 'Cause of death, by injury (% of total)' has been dropped.
Column 'Income share held by highest 20%' has 97 NaN values.

```

```
Column 'Income share held by highest 20%' has been dropped.  
Column 'Refugee population by country or territory of asylum' has 29 NaN values.
```

```
1 df.head(5)
```

	Country	Regional indicator	Happiness Score	Standard error of ladder score	upperwhisker	lowerwhisker	Logged GDP per capita	Social support	Healthy life expectancy	Freedom to make life choices	...	Refugee population by country or territory of asylum
0	afghanistan	South Asia	2.5669	0.031311	2.628270	2.505530	7.462861	0.470367	52.590000	0.396573	...	72278
1	albania	Central and Eastern Europe	4.8827	0.056116	4.992687	4.772713	9.417931	0.671070	68.708138	0.781994	...	105
2	algeria	Middle East and North Africa	5.0051	0.044236	5.091802	4.918397	9.537965	0.803385	65.905174	0.466611	...	97651
3	argentina	Latin America and Caribbean	5.9747	0.053442	6.079446	5.869954	9.810955	0.900568	68.803802	0.831132	...	3965
4	armenia	Commonwealth of Independent States	4.6768	0.058595	4.791646	4.561953	9.100476	0.757479	66.750656	0.712018	...	107910

5 rows × 59 columns

```
1 country_list = df['Country'].tolist()
```

This code block focuses on loading user data from MAL user dataset and performing initial processing steps, specifically related to dates and calculating user age.

```
1 users = pd.read_csv('users_cleaned.csv')  
2 users['join_date'] = pd.to_datetime(users['join_date'])  
3 users['last_online'] = pd.to_datetime(users['last_online'])  
4 users['birth_date'] = pd.to_datetime(users['birth_date'])  
5 now = datetime.datetime.now()  
6 users['age'] = (now - users['birth_date']).dt.days / 365.25  
7 users['age'] = users['age'].astype(int)  
8 users.head()
```

	username	user_id	user_watching	user_completed	user_onhold	user_dropped	user_plantowatch	user_days_spent_watching	gen
0	karthiga	2255153	3	49	1	0	0	55.091667	Fem
1	Damonashu	37326	45	195	27	25	59	82.574306	M
2	bskai	228342	25	414	2	5	11	159.483333	M
3	terune_uzumaki	327311	5	5	0	0	0	11.394444	Fem
4	Bas_G	5015094	35	114	6	20	175	30.458333	M

```
1 users.info()
```

```
2 <class 'pandas.core.frame.DataFrame'>  
RangeIndex: 108711 entries, 0 to 108710  
Data columns (total 18 columns):  
 #   Column           Non-Null Count  Dtype     
---  --   
 0   username        108710 non-null   object    
 1   user_id         108711 non-null   int64    
 2   user_watching   108711 non-null   int64    
 3   user_completed  108711 non-null   int64    
 4   user_onhold     108711 non-null   int64    
 5   user_dropped    108711 non-null   int64    
 6   user_plantowatch 108711 non-null   int64    
 7   user_days_spent_watching 108711 non-null   float64  
 8   gender          108711 non-null   object    
 9   location        108706 non-null   object    
 10  birth_date      108711 non-null   datetime64[ns]  
 11  access_rank    0 non-null       float64   
 12  join_date       108711 non-null   datetime64[ns]  
 13  last_online     108711 non-null   datetime64[ns]  
 14  stats_mean_score 108711 non-null   float64
```

```

15 stats_rewatched           108711 non-null  float64
16 stats_episodes            108711 non-null  int64
17 age                       108711 non-null  int64
dtypes: datetime64[ns](3), float64(4), int64(8), object(3)
memory usage: 14.9+ MB

```

This code block checks for and removes two columns, 'access_rank' and 'user_id', from the users DataFrame. These columns are not needed for subsequent analysis or modeling.

```

1 if 'access_rank' in users.columns:
2     users = users.drop(columns=['access_rank'])
3     print("Column 'access_rank' dropped successfully from users.")
4 else:
5     print("Column 'access_rank' not found in users.")
6
7 if 'user_id' in users.columns:
8     users = users.drop(columns=['user_id'])
9     print("Column 'user_id' dropped successfully from users.")
10 else:
11     print("Column 'user_id' not found in users.")

→ Column 'access_rank' dropped successfully from users.
Column 'user_id' dropped successfully from users.

```

```

1 us_states = [
2     "Alabama", "Alaska", "Arizona", "Arkansas", "California", "Colorado",
3     "Connecticut", "Delaware", "Florida", "Georgia", "Hawaii", "Idaho",
4     "Illinois", "Indiana", "Iowa", "Kansas", "Kentucky", "Louisiana",
5     "Maine", "Maryland", "Massachusetts", "Michigan", "Minnesota",
6     "Mississippi", "Missouri", "Montana", "Nebraska", "Nevada",
7     "New Hampshire", "New Jersey", "New Mexico", "New York",
8     "North Carolina", "North Dakota", "Ohio", "Oklahoma", "Oregon",
9     "Pennsylvania", "Rhode Island", "South Carolina", "South Dakota",
10    "Tennessee", "Texas", "Utah", "Vermont", "Virginia", "Washington",
11    "West Virginia", "Wisconsin", "Wyoming"
12 ]
13 us_states_lower = [state.lower() for state in us_states]

```

Function to turn location data to country.

```

1 def find_country(input):
2     input = str(input)
3     if input == 'nan':
4         return 'Other'
5     else:
6         for country in country_list:
7             if country in input.lower():
8                 return country
9     for state in us_states_lower:
10        if state in input.lower():
11            return 'United States'
12    return 'Other'

```

This code block applies the above function to determine the country for each user based on their location information.

```

1 users['Country'] = users['location'].apply(find_country)
2 users.info()

→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 108711 entries, 0 to 108710
Data columns (total 17 columns):
 #   Column           Non-Null Count  Dtype  
--- 
0   username          108710 non-null   object 
1   user_watching     108711 non-null   int64  
2   user_completed    108711 non-null   int64  
3   user_onhold       108711 non-null   int64  
4   user_dropped      108711 non-null   int64  
5   user_plantowatch 108711 non-null   int64  
6   user_days_spent_watching 108711 non-null   float64
7   gender            108711 non-null   object 
8   location          108706 non-null   object 
9   birth_date        108711 non-null   datetime64[ns]
10  join_date         108711 non-null   datetime64[ns]
11  last_online       108711 non-null   datetime64[ns]
12  stats_mean_score 108711 non-null   float64
13  stats_rewatched  108711 non-null   float64
14  stats_episodes    108711 non-null   int64  
15  age               108711 non-null   int64  

```

```
16 Country          108711 non-null object
dtypes: datetime64[ns](3), float64(3), int64(7), object(4)
memory usage: 14.1+ MB
```

This code block drops the now unnecessary location column.

```
1 if 'location' in users.columns:
2     users = users.drop(columns=['location'])
3     print("Column 'location' dropped successfully from users.")
4 else:
5     print("Column 'location' not found in users.")

→ Column 'location' dropped successfully from users.
```

Country column is important to merge df to users. The rows with Other value in Country column can't be merged and thus will be dropped.

```
1 other_count = users[users['Country'] == 'Other'].shape[0]
2 print(f"Number of 'Other' values in the 'Country' column: {other_count}")

→ Number of 'Other' values in the 'Country' column: 40904
```

This code block performs several filtering and cleaning steps on the users DataFrame to prepare it for further analysis or merging.

```
1 users = users[users['Country'] != 'Other']
2 users = users[users['user_completed'] != 0]
3 users = users[users['stats_mean_score'] != 0]
4 users = users[users['gender'].isin(['Male', 'Female'])]
5 users.reset_index(drop=True, inplace=True)
6 users.info()

→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 65828 entries, 0 to 65827
Data columns (total 16 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   username        65828 non-null   object  
 1   user_watching   65828 non-null   int64  
 2   user_completed  65828 non-null   int64  
 3   user_onhold     65828 non-null   int64  
 4   user_dropped    65828 non-null   int64  
 5   user_plantowatch 65828 non-null   int64  
 6   user_days_spent_watching 65828 non-null   float64 
 7   gender          65828 non-null   object  
 8   birth_date      65828 non-null   datetime64[ns]
 9   join_date       65828 non-null   datetime64[ns]
 10  last_online     65828 non-null   datetime64[ns]
 11  stats_mean_score 65828 non-null   float64 
 12  stats_rewatched 65828 non-null   float64 
 13  stats_episodes   65828 non-null   int64  
 14  age             65828 non-null   int64  
 15  Country         65828 non-null   object  
dtypes: datetime64[ns](3), float64(3), int64(7), object(3)
memory usage: 8.0+ MB
```

This code block performs a merge operation between two pandas DataFrames: users and df.

```
1 user_country_stats = pd.merge(users, df, on='Country', how='inner')
2 user_country_stats.head()
```



```
username user_watching user_completed user_onhold user_dropped user_plantowatch user_days_spent_watching gender birt
```

0	karthiga	3	49	1	0	0	55.091667	Female	1991
1	Damonashu	45	195	27	25	59	82.574306	Male	1991
2	bskai	25	414	2	5	11	159.483333	Male	1991
3	terune_uzumaki	5	5	0	0	0	11.394444	Female	1991
4	HimeAria	2	87	2	0	28	19.963194	Female	1991

5 rows × 74 columns

This code block is designed to check each column in the user_country_stats DataFrame for missing values (represented as NaN). If a column contains missing values, the code calculates the total number of missing values in that column and prints this information. If the number of missing values is bigger than 5000, the column is dropped.

```
1 for col in user_country_stats.columns:
2     if user_country_stats[col].isnull().any():
3         nan_sum = user_country_stats[col].isnull().sum()
4         print(f"Column '{col}' has {nan_sum} NaN values.")
5         if nan_sum > 5000:
6             user_country_stats.drop(columns=[col], inplace=True)
7             print(f"Column '{col}' has been dropped.")

→ Column 'username' has 1 NaN values.
Column 'High-technology exports (current US$)' has 4748 NaN values.
Column 'Forest area (% of land area)' has 4350 NaN values.
Column 'Government expenditure on education, total (% of GDP)' has 5395 NaN values.
Column 'Government expenditure on education, total (% of GDP)' has been dropped.
Column 'GDP per capita growth (annual %)' has 4350 NaN values.
Column 'Urban population (% of total population)' has 4350 NaN values.
Column 'Net migration' has 4350 NaN values.
Column 'GNI per capita, PPP (current international $)' has 4350 NaN values.
Column 'Exports of goods and services (% of GDP)' has 4419 NaN values.
Column 'Agriculture, forestry, and fishing, value added (% of GDP)' has 4350 NaN values.
Column 'Access to electricity (% of population)' has 4350 NaN values.
Column 'Population growth (annual %)' has 4350 NaN values.
Column 'Renewable electricity output (% of total electricity output)' has 4350 NaN values.
Column 'Electric power consumption (kWh per capita)' has 4381 NaN values.
Column 'Renewable internal freshwater resources per capita (cubic meters)' has 4364 NaN values.
Column 'Population density (people per sq. km of land area)' has 4350 NaN values.
Column 'Mobile cellular subscriptions (per 100 people)' has 4350 NaN values.
Column 'Fixed broadband subscriptions (per 100 people)' has 4350 NaN values.
Column 'Secure Internet servers (per 1 million people)' has 4350 NaN values.
Column 'Tax revenue (% of GDP)' has 6695 NaN values.
Column 'Tax revenue (% of GDP)' has been dropped.
Column 'Mortality rate, under-5 (per 1,000 live births)' has 4350 NaN values.
Column 'Employment to population ratio, 15+, total (%) (modeled ILO estimate)' has 4350 NaN values.
Column 'Death rate, crude (per 1,000 people)' has 4350 NaN values.
Column 'PM2.5 air pollution, mean annual exposure (micrograms per cubic meter)' has 4350 NaN values.
Column 'Agricultural land (% of land area)' has 4350 NaN values.
Column 'Rural population (% of total population)' has 4350 NaN values.
Column 'Maternal mortality ratio (modeled estimate, per 100,000 live births)' has 4350 NaN values.
Column 'Labor force, female (% of total labor force)' has 4350 NaN values.
Column 'Population in the largest city (% of urban population)' has 4662 NaN values.
Column 'New businesses registered (number)' has 23923 NaN values.
Column 'New businesses registered (number)' has been dropped.
Column 'Refugee population by country or territory of asylum' has 5135 NaN values.
Column 'Refugee population by country or territory of asylum' has been dropped.
Column 'Life expectancy at birth, male (years)' has 4350 NaN values.
Column 'Life expectancy at birth, female (years)' has 4350 NaN values.
Column 'Medium and high-tech exports (% manufactured exports)' has 4391 NaN values.
Column 'Military expenditure (% of GDP)' has 4602 NaN values.
Column 'Inflation, consumer prices (annual %)' has 4903 NaN values.
Column 'Population, total' has 4350 NaN values.
Column 'GDP per capita (current US$)' has 4350 NaN values.
Column 'Unemployment, total (% of total labor force) (modeled ILO estimate)' has 4350 NaN values.
Column 'Fertility rate, total (births per woman)' has 4350 NaN values.
```

This code block drops the rows with NaN value in most of the columns.

```

1 if 'Inflation, consumer prices (annual %)' in user_country_stats.columns:
2     user_country_stats.dropna(subset=['Inflation, consumer prices (annual %)'], inplace=True)
3     print("Rows with missing values in 'Inflation, consumer prices (annual %)' dropped.")
4 else:
5     print("'Inflation, consumer prices (annual %)' column not found in the DataFrame.")
6
7 user_country_stats.info()

```

→ Rows with missing values in 'Inflation, consumer prices (annual %)' dropped.
<class 'pandas.core.frame.DataFrame'>
Index: 60925 entries, 0 to 65827
Data columns (total 70 columns):

#	Column	Non-Null Count	Dtype
0	username	60924	non-null object
1	user_watching	60925	non-null int64
2	user_completed	60925	non-null int64
3	user_onhold	60925	non-null int64
4	user_dropped	60925	non-null int64
5	user_plantowatch	60925	non-null int64
6	user_days_spent_watching	60925	non-null float64
7	gender	60925	non-null object
8	birth_date	60925	non-null datetime64[ns]
9	join_date	60925	non-null datetime64[ns]
10	last_online	60925	non-null datetime64[ns]
11	stats_mean_score	60925	non-null float64
12	stats_rewatched	60925	non-null float64
13	stats_episodes	60925	non-null int64
14	age	60925	non-null int64
15	Country	60925	non-null object
16	Regional indicator	60925	non-null object
17	Happiness Score	60925	non-null float64
18	Standard error of ladder score	60925	non-null float64
19	upperwhisker	60925	non-null float64
20	lowerwhisker	60925	non-null float64
21	Logged GDP per capita	60925	non-null float64
22	Social support	60925	non-null float64
23	Healthy life expectancy	60925	non-null float64
24	Freedom to make life choices	60925	non-null float64
25	Generosity	60925	non-null float64
26	Perceptions of corruption	60925	non-null float64
27	Ladder score in Dystopia	60925	non-null float64
28	Explained by: Log GDP per capita	60925	non-null float64
29	Explained by: Social support	60925	non-null float64
30	Explained by: Healthy life expectancy	60925	non-null float64
31	Explained by: Freedom to make life choices	60925	non-null float64
32	Explained by: Generosity	60925	non-null float64
33	Explained by: Perceptions of corruption	60925	non-null float64
34	Dystopia + residual	60925	non-null float64
35	High-technology exports (current US\$)	60527	non-null float64
36	Forest area (% of land area)	60925	non-null float64
37	GDP per capita growth (annual %)	60925	non-null float64
38	Urban population (% of total population)	60925	non-null float64
39	Net migration	60925	non-null float64
40	GNI per capita, PPP (current international \$)	60925	non-null float64
41	Exports of goods and services (% of GDP)	60872	non-null float64
42	Agriculture, forestry, and fishing, value added (% of GDP)	60925	non-null float64
43	Access to electricity (% of population)	60925	non-null float64
44	Population growth (annual %)	60925	non-null float64
45	Renewable electricity output (% of total electricity output)	60925	non-null float64
46	Electric power consumption (kWh per capita)	60894	non-null float64
47	Renewable internal freshwater resources per capita (cubic meters)	60911	non-null float64
48	Population density (people per sq. km of land area)	60925	non-null float64
49	Mobile cellular subscriptions (per 100 people)	60925	non-null float64
50	Fixed broadband subscriptions (per 100 people)	60925	non-null float64
51	Secure Internet servers (per 1 million people)	60925	non-null float64

This code block is used to identify and count missing values within each column of the user_country_stats DataFrame.

```

1 for col in user_country_stats.columns:
2     if user_country_stats[col].isnull().any():
3         nan_sum = user_country_stats[col].isnull().sum()
4         print(f"Column '{col}' has {nan_sum} NaN values.")

→ Column 'username' has 1 NaN values.
Column 'High-technology exports (current US$)' has 398 NaN values.
Column 'Exports of goods and services (% of GDP)' has 53 NaN values.
Column 'Electric power consumption (kWh per capita)' has 31 NaN values.
Column 'Renewable internal freshwater resources per capita (cubic meters)' has 14 NaN values.
Column 'Population in the largest city (% of urban population)' has 312 NaN values.
Column 'Medium and high-tech exports (% manufactured exports)' has 41 NaN values.
Column 'Military expenditure (% of GDP)' has 252 NaN values.

```

This code block drops the column named username in the user_country_stats DataFrame.

```

1 if 'username' in user_country_stats.columns:
2     user_country_stats.dropna(subset=['username'], inplace=True)
3     print("Rows with missing values in 'username' dropped.")
4 else:
5     print("'username' column not found in the DataFrame.")

```

→ Rows with missing values in 'username' dropped.

This code block focuses on addressing any remaining missing values within the numerical columns of the user_country_stats DataFrame.

```

1 numeric_columns = user_country_stats.select_dtypes(include=np.number).columns
2 user_country_stats[numeric_columns] = user_country_stats[numeric_columns].fillna(user_country_stats[numeric_columns].mean())

```

This code block is used to systematically check each column in the user_country_stats DataFrame for missing values.

```

1 for col in user_country_stats.columns:
2     if user_country_stats[col].isnull().any():
3         nan_sum = user_country_stats[col].isnull().sum()
4         print(f"Column '{col}' has {nan_sum} NaN values.")

```

This code block is used to remove three columns from the user_country_stats DataFrame: 'upperwhisker', 'lowerwhisker', and 'Standard error of ladder score'. These columns are considered unnecessary for the subsequent analysis or machine learning tasks.

```

1 if 'upperwhisker' in user_country_stats.columns:
2     user_country_stats = user_country_stats.drop(columns=['upperwhisker'])
3     print("Column 'upperwhisker' dropped successfully from user_country_stats.")
4 else:
5     print("Column 'upperwhisker' not found in user_country_stats.")
6
7 if 'lowerwhisker' in user_country_stats.columns:
8     user_country_stats = user_country_stats.drop(columns=['lowerwhisker'])
9     print("Column 'lowerwhisker' dropped successfully from user_country_stats.")
10 else:
11     print("Column 'lowerwhisker' not found in user_country_stats.")
12
13 if 'Standard error of ladder score' in user_country_stats.columns:
14     user_country_stats = user_country_stats.drop(columns=['Standard error of ladder score'])
15     print("Column 'Standard error of ladder score' dropped successfully from user_country_stats.")
16 else:
17     print("Column 'Standard error of ladder score' not found in user_country_stats.")

```

→ Column 'upperwhisker' dropped successfully from user_country_stats.
 Column 'lowerwhisker' dropped successfully from user_country_stats.
 Column 'Standard error of ladder score' dropped successfully from user_country_stats.

This code block focuses on converting categorical data in two columns of user_country_stats into a numerical format.

```

1 label_encoder = LabelEncoder()
2 user_country_stats['Regional indicator'] = label_encoder.fit_transform(user_country_stats['Regional indicator'])
3 user_country_stats['Country'] = label_encoder.fit_transform(user_country_stats['Country'])

```

This line of code converts the categorical 'gender' column into a numerical representation, creating a new column called gender_numeric.

```
1 user_country_stats['gender_numeric'] = user_country_stats['gender'].map({'Male': 0, 'Female': 1})
```

```
1 user_country_stats.info()
```

→ <class 'pandas.core.frame.DataFrame'>
 Index: 60924 entries, 0 to 65827
 Data columns (total 68 columns):

#	Column	Non-Null Count	Dtype
0	username	60924	object
1	user_watching	60924	non-null int64
2	user_completed	60924	non-null int64
3	user_onhold	60924	non-null int64
4	user_dropped	60924	non-null int64
5	user_plantowatch	60924	non-null int64
6	user_days_spent_watching	60924	non-null float64
7	gender	60924	non-null object
8	birth_date	60924	non-null datetime64[ns]
9	join_date	60924	non-null datetime64[ns]
10	last_online	60924	non-null datetime64[ns]
11	stats_mean_score	60924	non-null float64
12	stats_rewatched	60924	non-null float64
13	stats_episodes	60924	non-null int64

14	age	60924	non-null	int64
15	Country	60924	non-null	int64
16	Regional indicator	60924	non-null	int64
17	Happiness Score	60924	non-null	float64
18	Logged GDP per capita	60924	non-null	float64
19	Social support	60924	non-null	float64
20	Healthy life expectancy	60924	non-null	float64
21	Freedom to make life choices	60924	non-null	float64
22	Generosity	60924	non-null	float64
23	Perceptions of corruption	60924	non-null	float64
24	Ladder score in Dystopia	60924	non-null	float64
25	Explained by: Log GDP per capita	60924	non-null	float64
26	Explained by: Social support	60924	non-null	float64
27	Explained by: Healthy life expectancy	60924	non-null	float64
28	Explained by: Freedom to make life choices	60924	non-null	float64
29	Explained by: Generosity	60924	non-null	float64
30	Explained by: Perceptions of corruption	60924	non-null	float64
31	Dystopia + residual	60924	non-null	float64
32	High-technology exports (current US\$)	60924	non-null	float64
33	Forest area (% of land area)	60924	non-null	float64
34	GDP per capita growth (annual %)	60924	non-null	float64
35	Urban population (% of total population)	60924	non-null	float64
36	Net migration	60924	non-null	float64
37	GNI per capita, PPP (current international \$)	60924	non-null	float64
38	Exports of goods and services (% of GDP)	60924	non-null	float64
39	Agriculture, forestry, and fishing, value added (% of GDP)	60924	non-null	float64
40	Access to electricity (% of population)	60924	non-null	float64
41	Population growth (annual %)	60924	non-null	float64
42	Renewable electricity output (% of total electricity output)	60924	non-null	float64
43	Electric power consumption (kWh per capita)	60924	non-null	float64
44	Renewable internal freshwater resources per capita (cubic meters)	60924	non-null	float64
45	Population density (people per sq. km of land area)	60924	non-null	float64
46	Mobile cellular subscriptions (per 100 people)	60924	non-null	float64
47	Fixed broadband subscriptions (per 100 people)	60924	non-null	float64
48	Secure Internet servers (per 1 million people)	60924	non-null	float64
49	Mortality rate, under-5 (per 1,000 live births)	60924	non-null	float64
50	Employment to population ratio, 15+, total (%) (modeled ILO estimate)	60924	non-null	float64
51	Death rate, crude (per 1,000 people)	60924	non-null	float64
52	PM2.5 air pollution, mean annual exposure (micrograms per cubic meter)	60924	non-null	float64

Conduct Exploratory Data Analysis

```

1 numerical_features = ['age',
2                         'stats_mean_score',
3                         'Happiness Score',
4                         'user_days_spent_watching',
5                         'user_completed',
6                         'GDP per capita (current US$)',
7                         'Logged GDP per capita',
8                         'Social support',
9                         'Healthy life expectancy',
10                        'Dystopia + residual',
11                        'Freedom to make life choices',
12                        'Generosity',
13                        'Perceptions of corruption',
14                        'user_watching',
15                        'user_plantowatch']
16
17 interesting_pairs = [
18     ('Happiness Score', 'Logged GDP per capita'),
19     ('Happiness Score', 'Healthy life expectancy'),
20     ('stats_mean_score', 'user_completed'),
21     ('user_days_spent_watching', 'user_completed'),
22 ]

```

This code block focuses on visualizing the distribution of the gender column in the user_country_stats DataFrame using a pie chart.

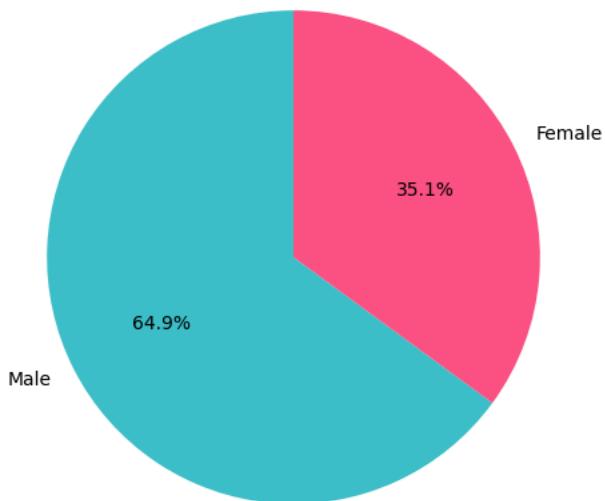
```

1 categorical_feature = 'gender'
2 my_colors = ['#3FC1C9', '#FC5185']
3
4 gender_counts = user_country_stats[categorical_feature].value_counts()
5 plt.figure(figsize=(8, 6))
6 plt.pie(gender_counts, labels=gender_counts.index, autopct='%.1f%%', startangle=90, colors=my_colors)
7 plt.title(f'Distribution of {categorical_feature}')
8 plt.show()

```

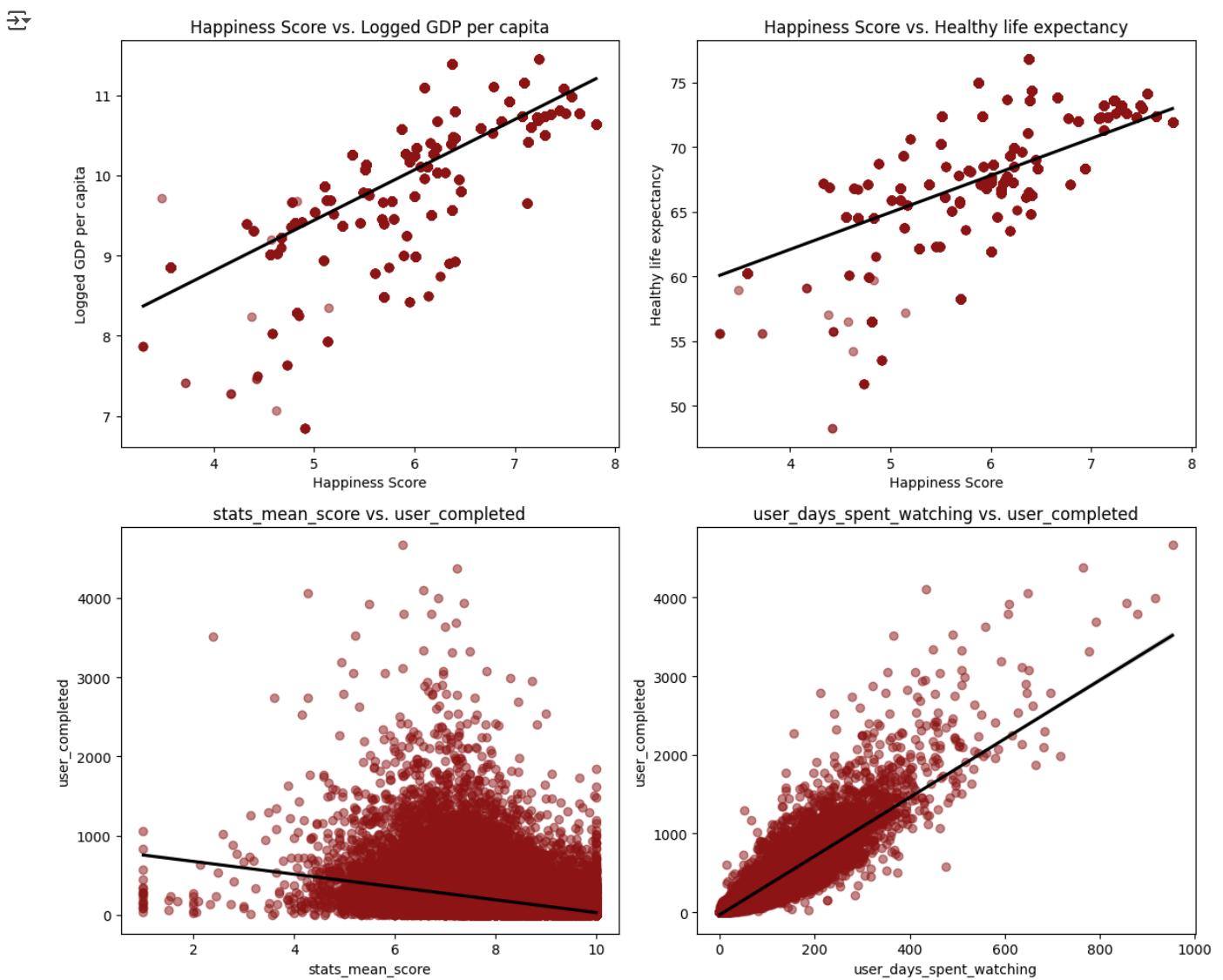


Distribution of gender



This code block focuses on visualizing the relationships between specific pairs of numerical features identified as "interesting" earlier in the notebook. It uses scatter plots with regression lines to show the general trend between the features.

```
1 num_pairs = len(interesting_pairs)
2 num_cols = 2
3 num_rows = math.ceil(num_pairs / num_cols)
4
5 fig, axes = plt.subplots(num_rows, num_cols, figsize=(12, 10))
6 axes = axes.flatten()
7
8 for i, pair in enumerate(interesting_pairs):
9     feature1, feature2 = pair
10    sns.regplot(x=feature1, y=feature2, data=user_country_stats, scatter_kws={'alpha': 0.5, 'color': '#8E1616'}, line_kws={'color':
11        axes[i].set_title(f'{feature1} vs. {feature2}')
12
13 for i in range(num_pairs, num_rows * num_cols):
14     fig.delaxes(axes[i])
15
16 plt.tight_layout()
17 plt.show()
```



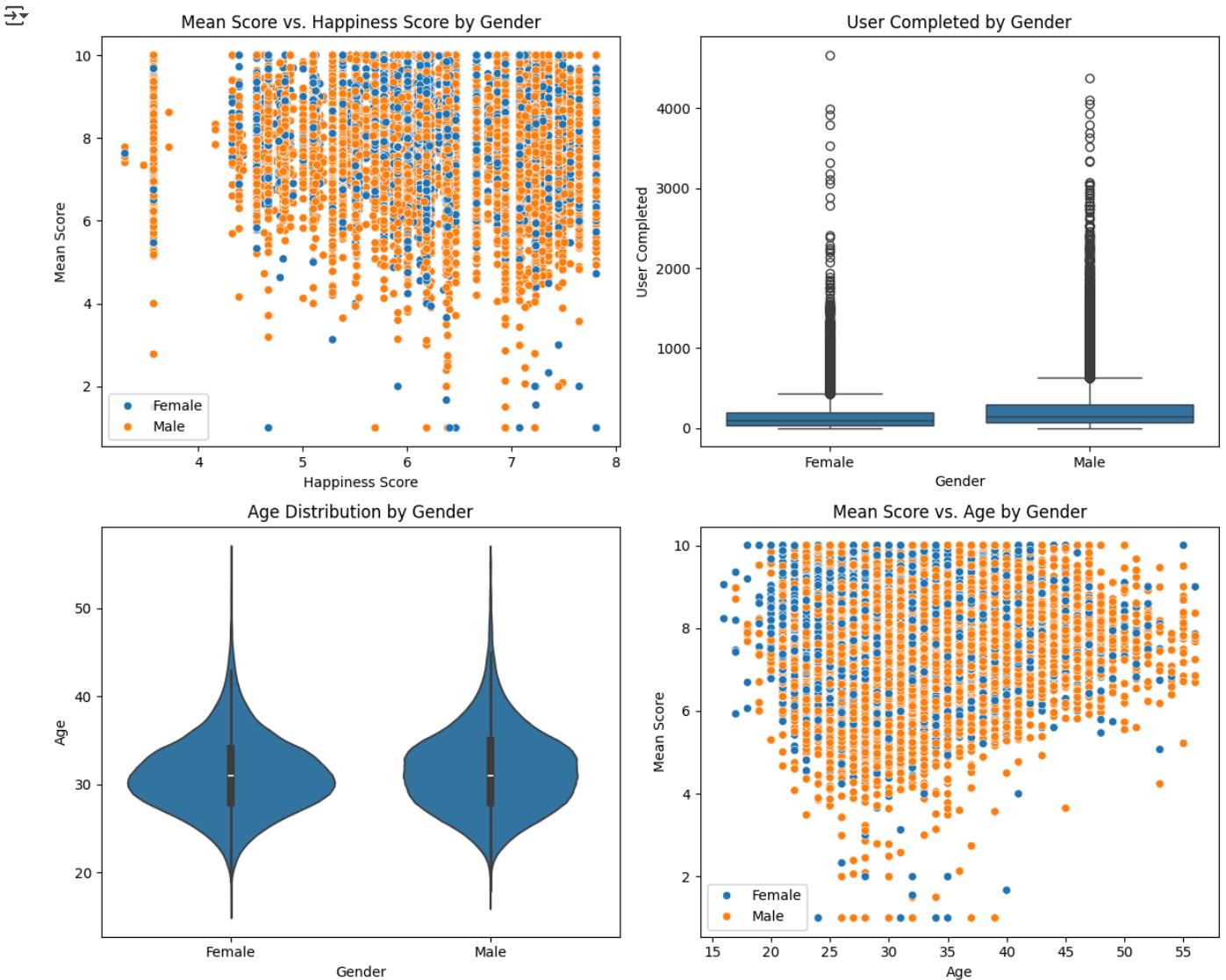
This code block uses the Seaborn and Matplotlib libraries to create a figure with multiple plots, specifically a 2x2 grid of subplots. Each subplot visualizes the relationship between different features from the user_country_stats DataFrame, differentiated by gender.

```

1 fig, axes = plt.subplots(2, 2, figsize=(12, 10))
2 sns.scatterplot(x='Happiness Score', y='stats_mean_score', hue='gender', data=user_country_stats, ax=axes[0, 0])
3 axes[0, 0].set_title('Mean Score vs. Happiness Score by Gender')
4 axes[0, 0].set_xlabel('Happiness Score')
5 axes[0, 0].set_ylabel('Mean Score')
6 axes[0, 0].legend()
7
8 sns.boxplot(x='gender', y='user_completed', data=user_country_stats, ax=axes[0, 1])
9 axes[0, 1].set_title('User Completed by Gender')
10 axes[0, 1].set_xlabel('Gender')
11 axes[0, 1].set_ylabel('User Completed')
12
13 sns.violinplot(x='gender', y='age', data=user_country_stats, ax=axes[1, 0])
14 axes[1, 0].set_title('Age Distribution by Gender')
15 axes[1, 0].set_xlabel('Gender')
16 axes[1, 0].set_ylabel('Age')
17
18 sns.scatterplot(x='age', y='stats_mean_score', hue='gender', data=user_country_stats, ax=axes[1, 1])
19 axes[1, 1].set_title('Mean Score vs. Age by Gender')
20 axes[1, 1].set_xlabel('Age')
21 axes[1, 1].set_ylabel('Mean Score')
22 axes[1, 1].legend()
23
24

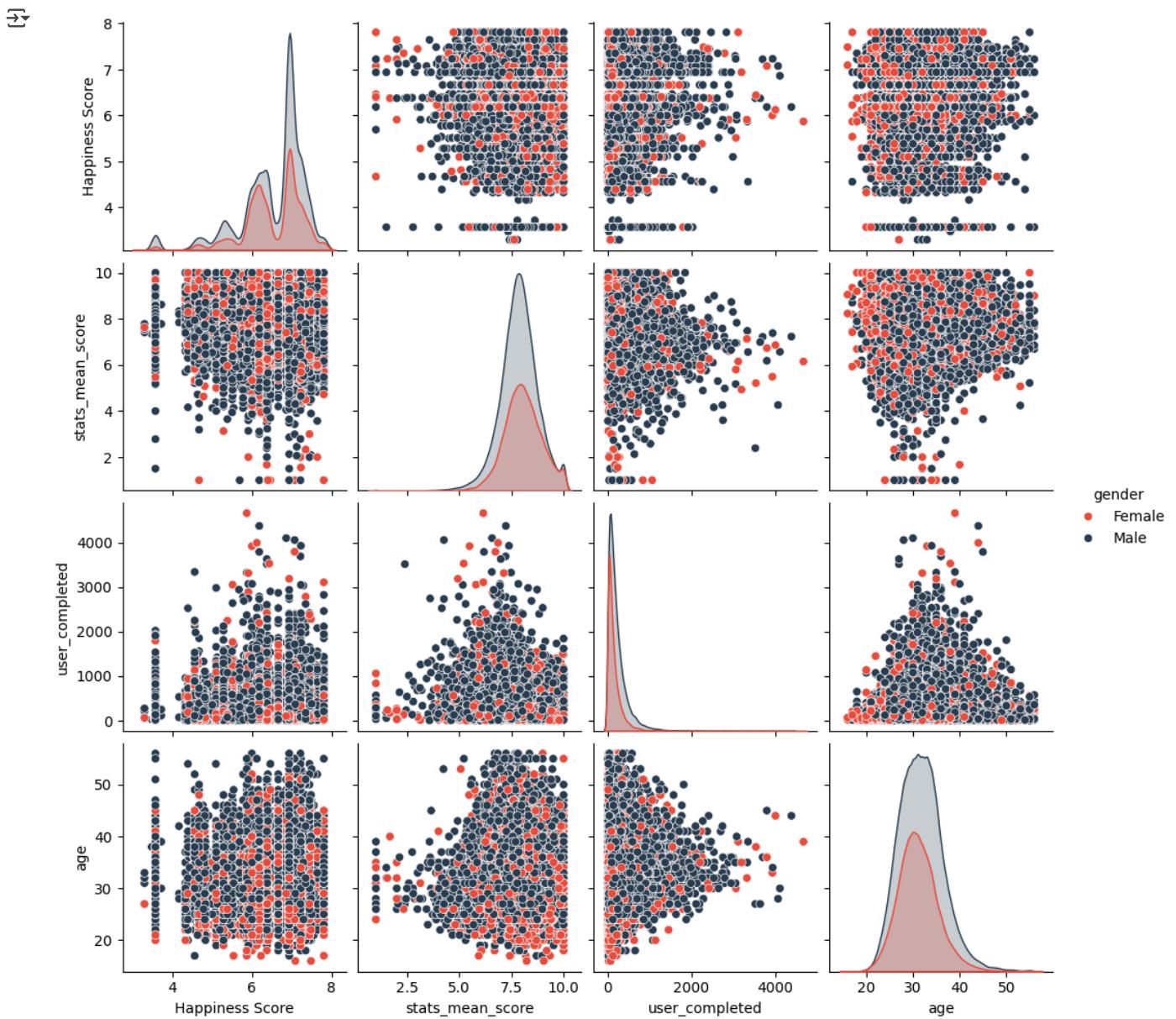
```

```
25 plt.tight_layout()
26 plt.show();
```



This code block generates a pair plot using the seaborn library. The resulting pair plot will show scatter plots for every combination of the specified numerical variables ('Happiness Score', 'stats_mean_score', 'user_completed', and 'age').

```
1 sns.pairplot(user_country_stats,
2             vars=['Happiness Score', 'stats_mean_score', 'user_completed', 'age'],
3             hue='gender',
4             palette={'Male': '#2C3E50', 'Female': '#E74C3C'})
5 plt.show()
```



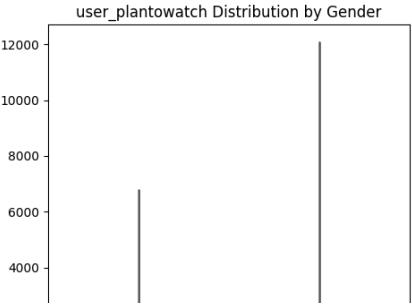
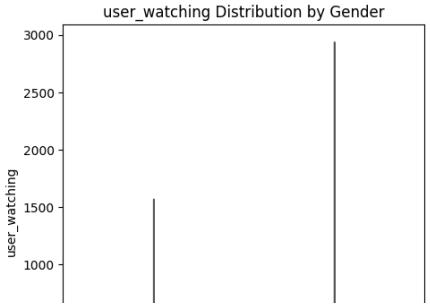
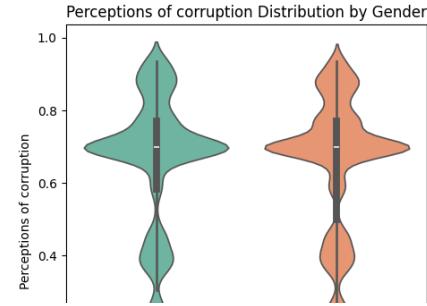
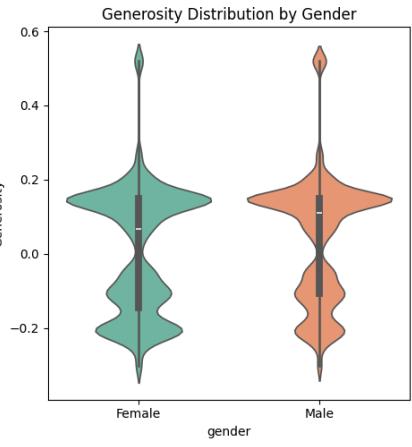
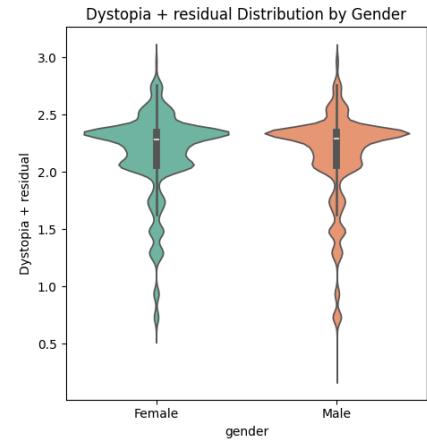
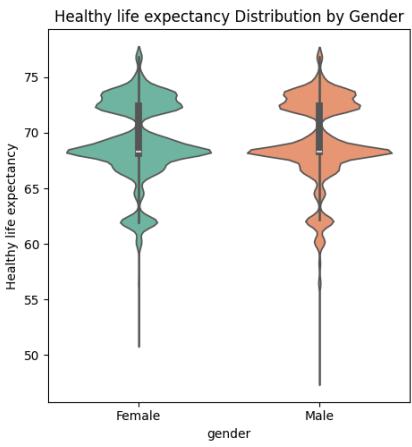
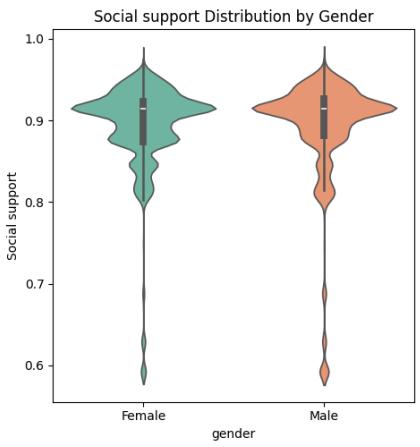
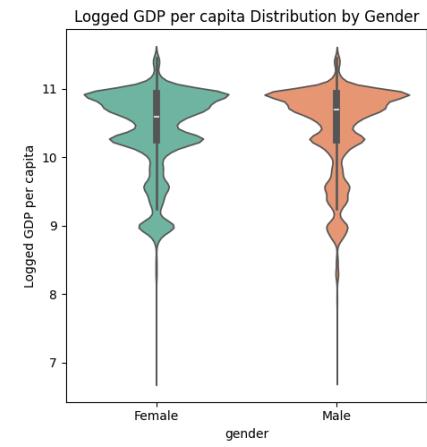
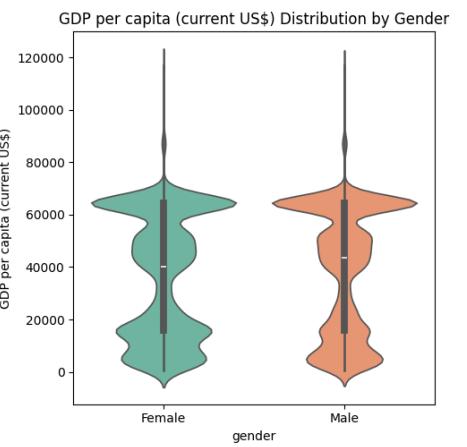
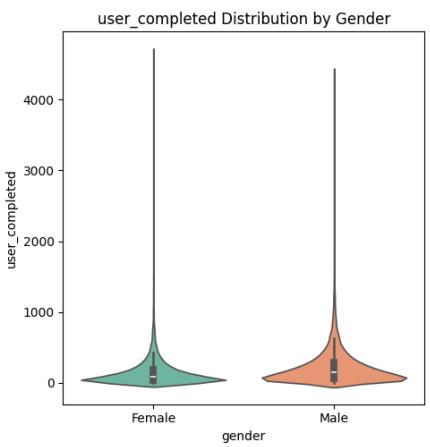
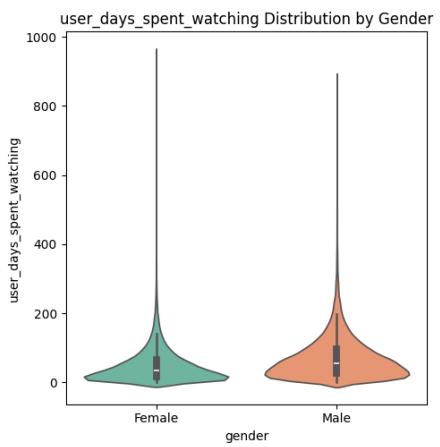
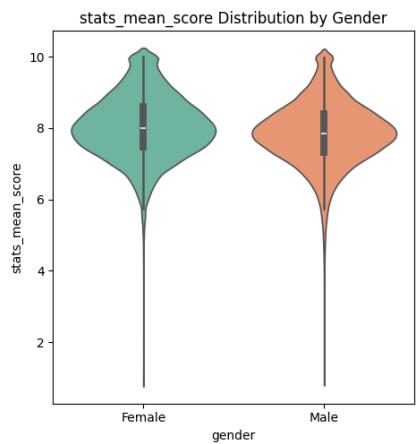
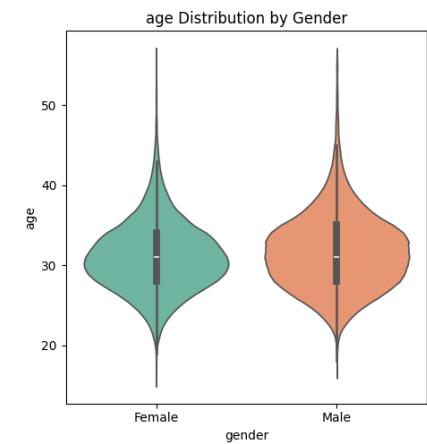
This code block generates a series of violin plots, one for each numerical feature in the numerical_features list. Each plot shows the distribution of that feature for both male and female users, allowing for a visual comparison of how these distributions differ by gender.

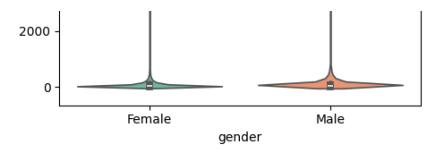
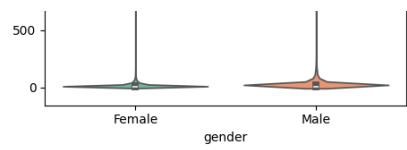
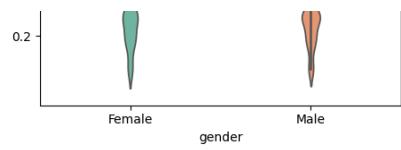
```

1 num_features = len(numerical_features)
2 num_cols = 3
3 num_rows = math.ceil(num_features / num_cols)
4
5 fig, axes = plt.subplots(num_rows, num_cols, figsize=(15, 5 * num_rows))
6
7 axes = axes.flatten()
8
9 for i, feature in enumerate(numerical_features):
10     sns.violinplot(x='gender', y=feature, data=user_country_stats, ax=axes[i], hue='gender', palette="Set2", legend=False)
11     axes[i].set_title(f'{feature} Distribution by Gender')
12
13 for i in range(num_features, num_rows * num_cols):
14     fig.delaxes(axes[i])
15
16 plt.tight_layout()
17 plt.show()

```

[]

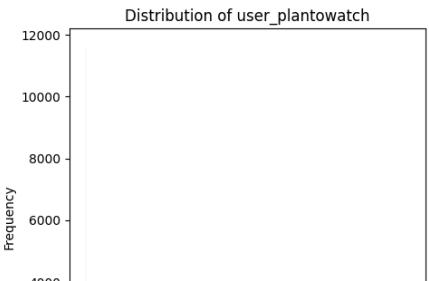
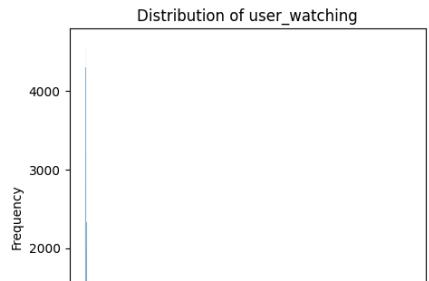
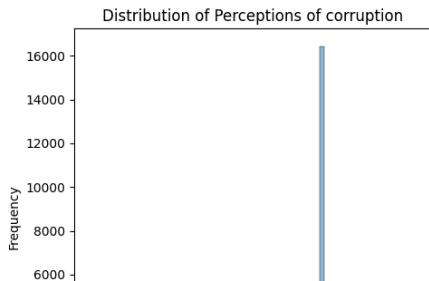
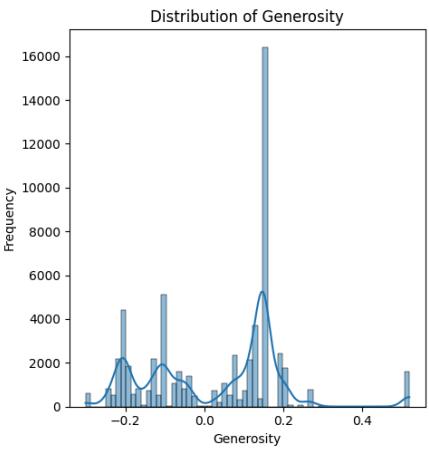
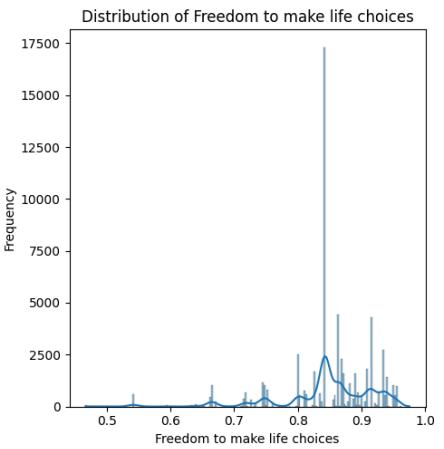
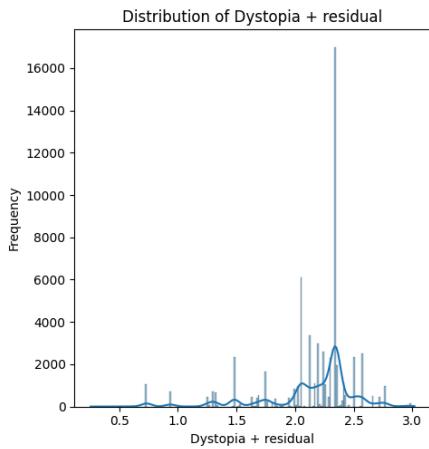
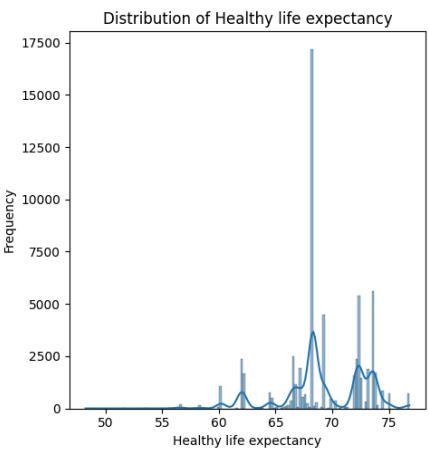
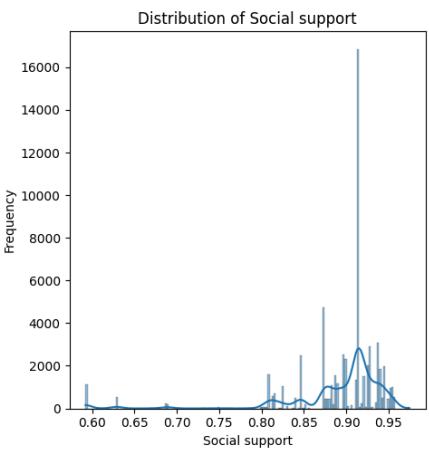
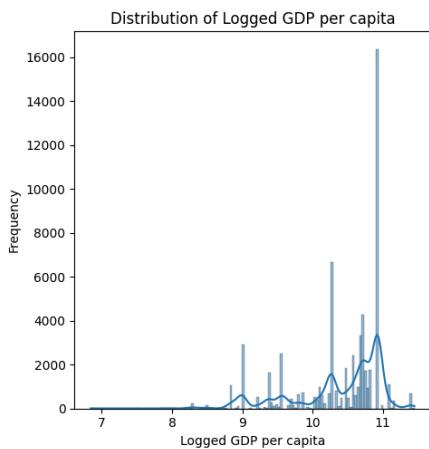
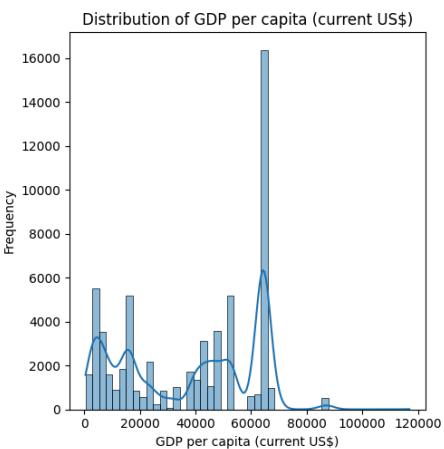
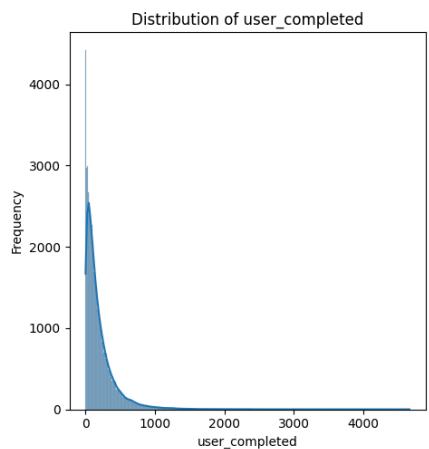
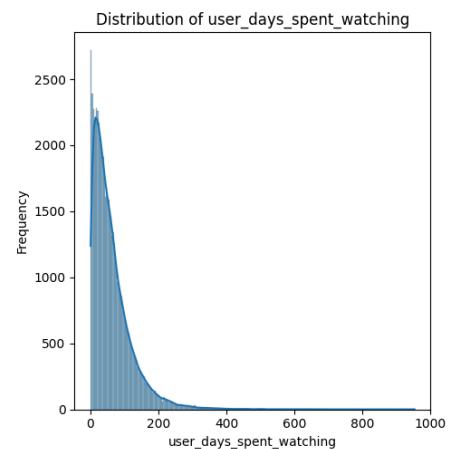
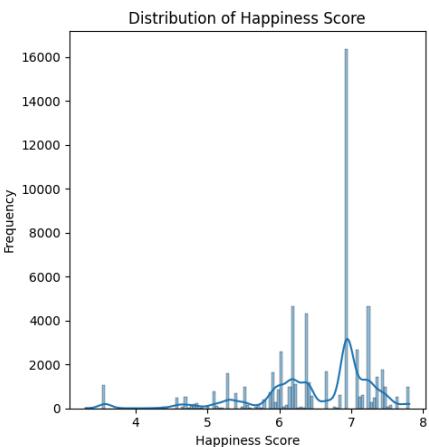
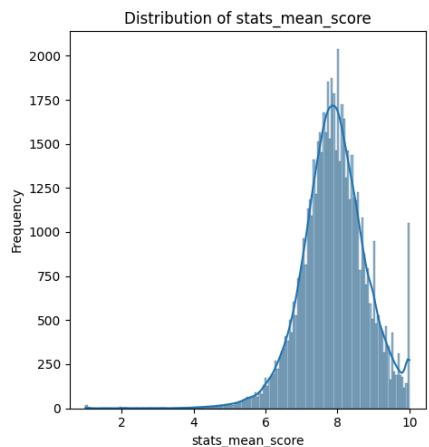
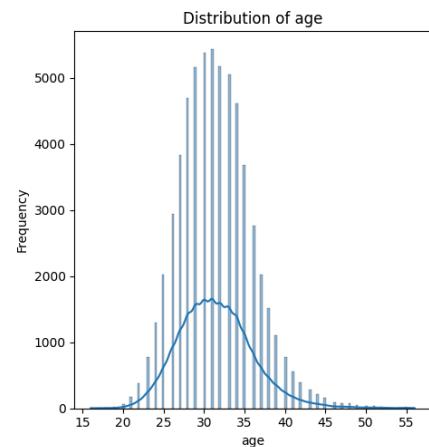


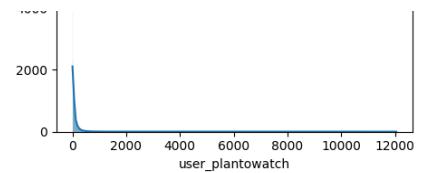
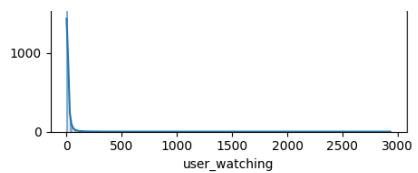
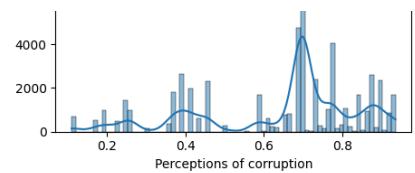


This code block focuses on visualizing the distribution of several numerical features from the user_country_stats DataFrame using histograms.

```
1 num_rows = 5
2 num_cols = 3
3
4 fig, axes = plt.subplots(num_rows, num_cols, figsize=(15, 25))
5
6 axes = axes.flatten()
7
8 for i, feature in enumerate(numerical_features):
9     sns.histplot(user_country_stats[feature], kde=True, ax=axes[i])
10    axes[i].set_title(f'Distribution of {feature}')
11    axes[i].set_xlabel(feature)
12    axes[i].set_ylabel('Frequency')
13
14 plt.tight_layout()
15
16 plt.show()
```

[]



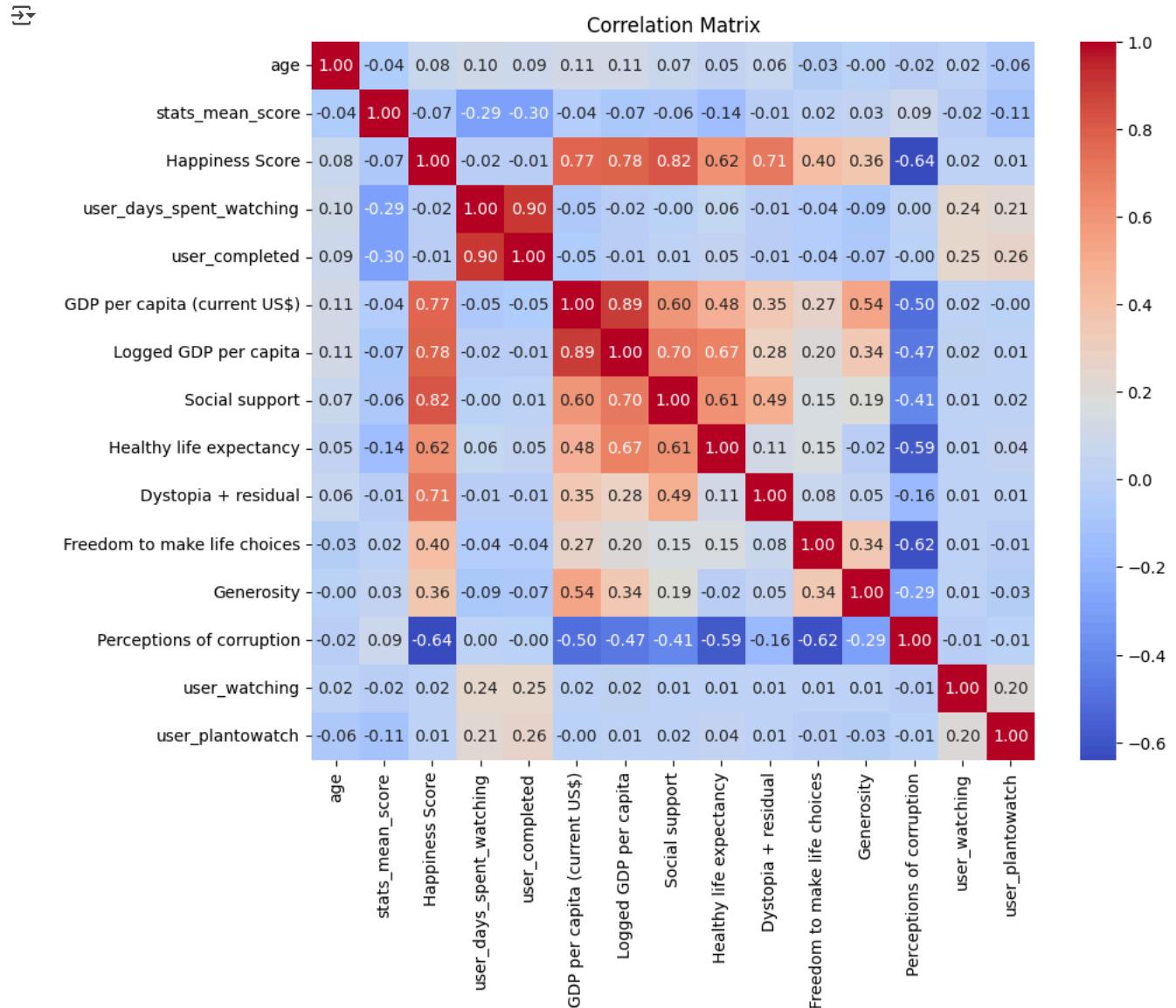


This code block calculates and visualizes the correlation matrix for the numerical features in the user_country_stats DataFrame.

```

1 correlation_matrix = user_country_stats[numerical_features].corr()
2
3 plt.figure(figsize=(10, 8))
4 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
5 plt.title('Correlation Matrix')
6 plt.show()

```



▼ Hypothesis Tests

The threshold for the hypotheses is $\alpha = 0.05$. If the p-value is less than or equal to α , we reject the null hypothesis and conclude that the variables are likely dependent. If the p-value is greater than α , we fail to reject the null hypothesis and conclude that the variables are likely independent.

H1: Users with a higher 'Mean Score' tend to have higher engagement metrics, such as a larger number of 'User Completed' anime and more 'User Days Spent Watching', suggesting satisfaction leads to deeper engagement

```

1 correlation_with_completed = user_country_stats['stats_mean_score'].corr(user_country_stats['user_completed'])
2 correlation_with_days_spent = user_country_stats['stats_mean_score'].corr(user_country_stats['user_days_spent_watching'])
3
4 print(f"Correlation between Mean Score and User Completed: {correlation_with_completed}")
5 print(f"Correlation between Mean Score and User Days Spent Watching: {correlation_with_days_spent}")

```

→ Correlation between Mean Score and User Completed: -0.2983295936486229
Correlation between Mean Score and User Days Spent Watching: -0.28560726935590597

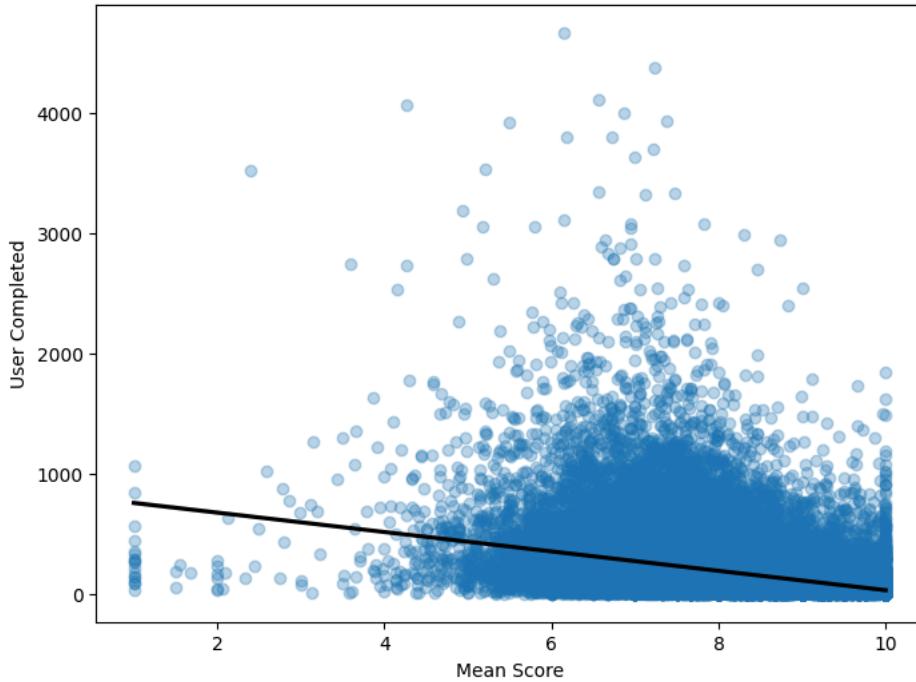
```

1 plt.figure(figsize=(8, 6))
2 sns.regplot(x='stats_mean_score', y='user_completed', data=user_country_stats, scatter_kws={'alpha':0.3}, line_kws={'color': 'black'}
3 plt.title('Mean Score vs. User Completed')
4 plt.xlabel('Mean Score')
5 plt.ylabel('User Completed')
6 plt.show()
7
8 plt.figure(figsize=(8, 6))
9 sns.regplot(x='stats_mean_score', y='user_days_spent_watching', data=user_country_stats, scatter_kws={'alpha':0.3}, line_kws={'color': 'black'}
10 plt.title('Mean Score vs. User Days Spent Watching')
11 plt.xlabel('Mean Score')
12 plt.ylabel('User Days Spent Watching')
13 plt.show()

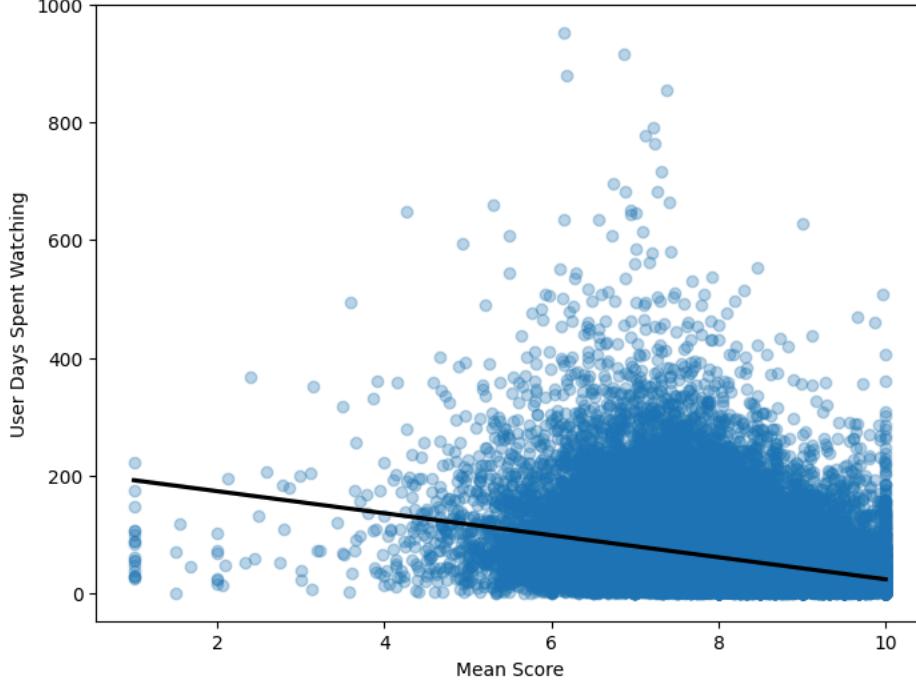
```

⤵

Mean Score vs. User Completed



Mean Score vs. User Days Spent Watching



```

1 tmp = user_country_stats.sample(n=500, random_state=42).copy()
2
3 tmp['mean_score_category'] = pd.qcut(tmp['stats_mean_score'], 3, labels=['low', 'medium', 'high'])
4 tmp['completed_category'] = pd.qcut(tmp['user_completed'], 3, labels=['low', 'medium', 'high'])
5 tmp['days_spent_category'] = pd.qcut(tmp['user_days_spent_watching'], 3, labels=['low', 'medium', 'high'])
6
7 contingency_table_completed = pd.crosstab(tmp['mean_score_category'], tmp['completed_category'])
8
9 contingency_table_days_spent = pd.crosstab(tmp['mean_score_category'], tmp['days_spent_category'])

```

```

10
11 chi2_stat_completed, p_val_completed, dof_completed, expected_completed = chi2_contingency(contingency_table_completed)
12
13 chi2_stat_days_spent, p_val_days_spent, dof_days_spent, expected_days_spent = chi2_contingency(contingency_table_days_spent)
14
15 print(f"Chi-square test for Mean Score and User Completed:")
16 print(f"Chi-square statistic: {chi2_stat_completed}")
17 print(f"P-value: {p_val_completed}")
18
19 print(f"\nChi-square test for Mean Score and User Days Spent Watching:")
20 print(f"Chi-square statistic: {chi2_stat_days_spent}")
21 print(f"P-value: {p_val_days_spent}")

→ Chi-square test for Mean Score and User Completed:
Chi-square statistic: 53.40157794100588
P-value: 7.022420420867169e-11

Chi-square test for Mean Score and User Days Spent Watching:
Chi-square statistic: 39.017057813558736
P-value: 6.910159230681829e-08

```

Conclusion: We reject the null hypothesis and conclude that there is a significant association

H2: Users from countries with higher 'GDP per capita (current US\$)' tend to watch more anime ('User Completed', 'User Watched Episodes') and potentially have higher 'Mean Scores', possibly reflecting better access or more leisure time/income.

```

1 correlation_with_completed = user_country_stats['GDP per capita (current US$)'].corr(user_country_stats['user_completed'])
2 correlation_with_days_spent = user_country_stats['GDP per capita (current US$)'].corr(user_country_stats['user_days_spent_watching'])
3 correlation_with_mean_score = user_country_stats['GDP per capita (current US$)'].corr(user_country_stats['stats_mean_score'])
4
5 print(f"Correlation between GDP per capita and User Completed: {correlation_with_completed}")
6 print(f"Correlation between GDP per capita and User Days Spent Watching: {correlation_with_days_spent}")
7 print(f"Correlation between GDP per capita and Mean Score: {correlation_with_mean_score}")

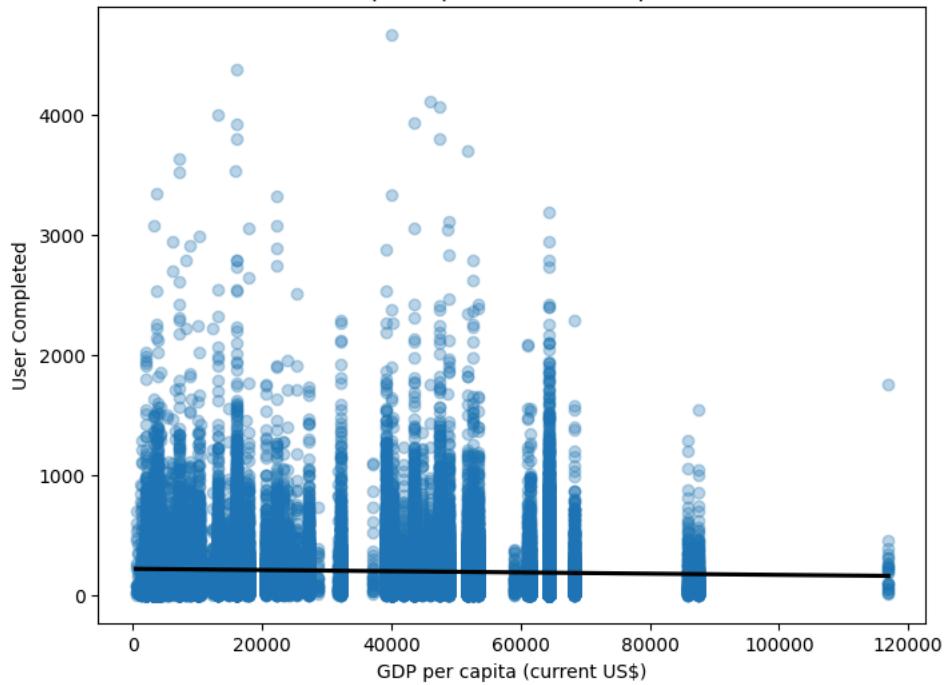
→ Correlation between GDP per capita and User Completed: -0.049227594442320624
Correlation between GDP per capita and User Days Spent Watching: -0.05220467832304506
Correlation between GDP per capita and Mean Score: -0.035823753105141014

1 plt.figure(figsize=(8, 6))
2 sns.regplot(x='GDP per capita (current US$)', y='user_completed', data=user_country_stats, scatter_kws={'alpha':0.3}, line_kws={'co:
3 plt.title('GDP per capita vs. User Completed')
4 plt.xlabel('GDP per capita (current US$)')
5 plt.ylabel('User Completed')
6 plt.show()
7
8 plt.figure(figsize=(8, 6))
9 sns.regplot(x='GDP per capita (current US$)', y='user_days_spent_watching', data=user_country_stats, scatter_kws={'alpha':0.3}, line_
10 plt.title('GDP per capita vs. User Days Spent Watching')
11 plt.xlabel('GDP per capita (current US$)')
12 plt.ylabel('User Days Spent Watching')
13 plt.show()
14
15 plt.figure(figsize=(8, 6))
16 sns.regplot(x='GDP per capita (current US$)', y='stats_mean_score', data=user_country_stats, scatter_kws={'alpha':0.3}, line_kws={'c
17 plt.title('GDP per capita vs. Mean Score')
18 plt.xlabel('GDP per capita (current US$)')
19 plt.ylabel('Mean Score')
20 plt.show()

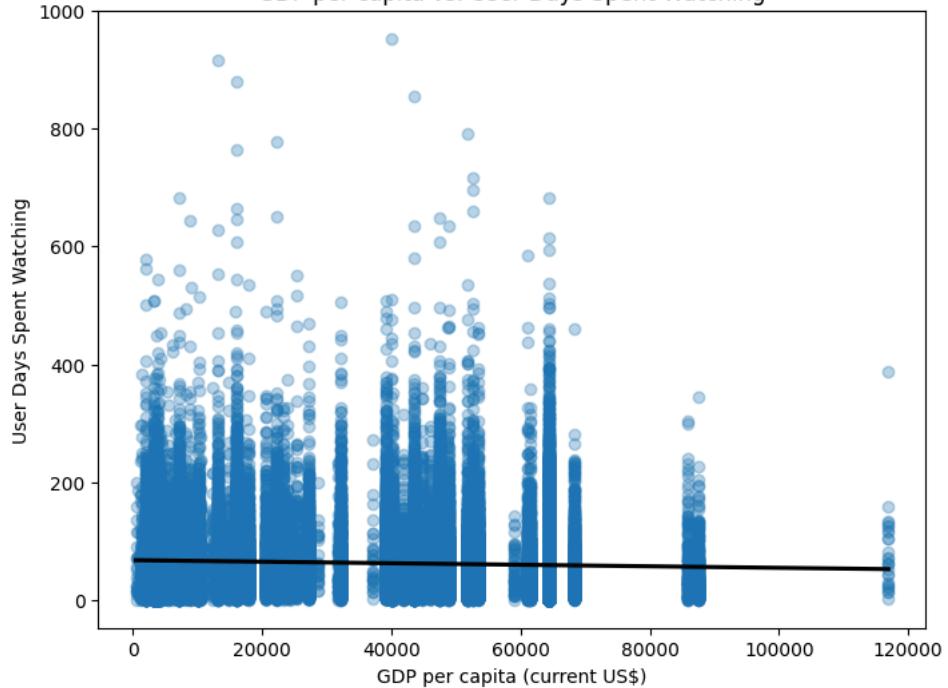
```

[]

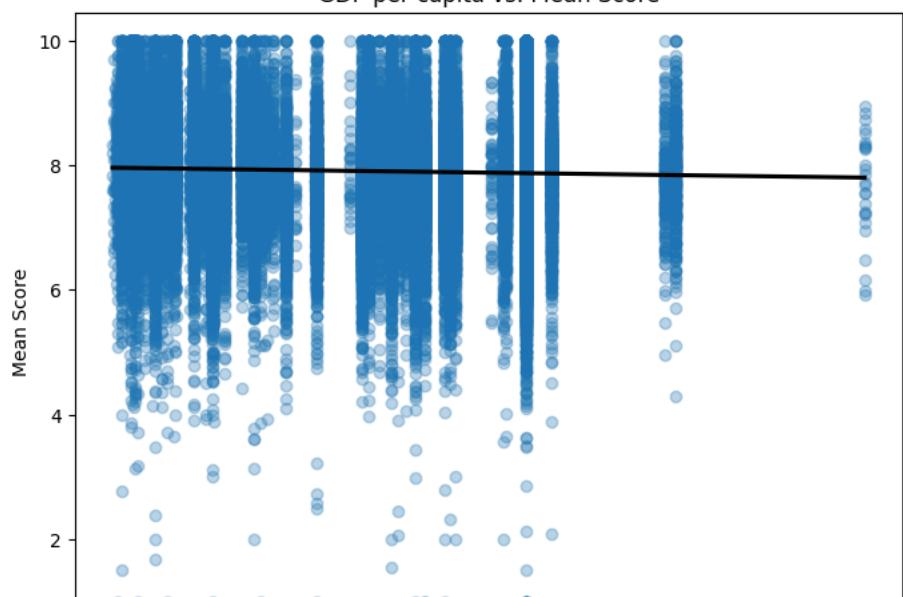
GDP per capita vs. User Completed

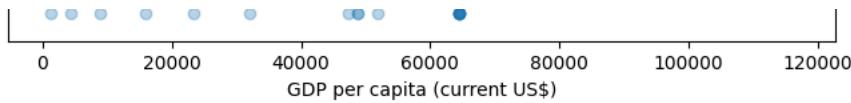


GDP per capita vs. User Days Spent Watching



GDP per capita vs. Mean Score





```

1 tmp = user_country_stats.sample(n=500, random_state=42).copy()
2
3 tmp['gdp_category'] = pd.qcut(tmp['GDP per capita (current US$)'], 3, labels=['low', 'medium', 'high'])
4
5 tmp['completed_category'] = pd.qcut(tmp['user_completed'], 3, labels=['low', 'medium', 'high'])
6
7 tmp['days_spent_category'] = pd.qcut(tmp['user_days_spent_watching'], 3, labels=['low', 'medium', 'high'])
8
9 tmp['mean_score_category'] = pd.qcut(tmp['stats_mean_score'], 3, labels=['low', 'medium', 'high'])
10
11 contingency_table_completed = pd.crosstab(tmp['gdp_category'], tmp['completed_category'])
12
13 contingency_table_days_spent = pd.crosstab(tmp['gdp_category'], tmp['days_spent_category'])
14
15 contingency_table_mean_score = pd.crosstab(tmp['gdp_category'], tmp['mean_score_category'])
16
17 chi2_stat_completed, p_val_completed, dof_completed, expected_completed = chi2_contingency(contingency_table_completed)
18
19 chi2_stat_days_spent, p_val_days_spent, dof_days_spent, expected_days_spent = chi2_contingency(contingency_table_days_spent)
20
21 chi2_stat_mean_score, p_val_mean_score, dof_mean_score, expected_mean_score = chi2_contingency(contingency_table_mean_score)
22
23 print(f"Chi-square test for GDP and User Completed:")
24 print(f"Chi-square statistic: {chi2_stat_completed}")
25 print(f"P-value: {p_val_completed}")
26
27 print(f"\nChi-square test for GDP and User Days Spent Watching:")
28 print(f"Chi-square statistic: {chi2_stat_days_spent}")
29 print(f"P-value: {p_val_days_spent}")
30
31 print(f"\nChi-square test for GDP and Mean Score:")
32 print(f"Chi-square statistic: {chi2_stat_mean_score}")
33 print(f"P-value: {p_val_mean_score}")

➡ Chi-square test for GDP and User Completed:
Chi-square statistic: 5.383837618156164
P-value: 0.25013050340042614

Chi-square test for GDP and User Days Spent Watching:
Chi-square statistic: 3.1269498625861756
P-value: 0.5368099282784534

Chi-square test for GDP and Mean Score:
Chi-square statistic: 8.533135573192155
P-value: 0.07388917746679664

```

Conclusion: We fail to reject the null hypothesis and conclude that there is no significant association

H3: Higher 'Fixed broadband subscriptions', 'Mobile cellular subscriptions', or 'Secure Internet servers (per 1 million people)' correlate with higher 'User Completed' counts for users in that country.

```

1 correlation_with_days_broadband = user_country_stats['Fixed broadband subscriptions (per 100 people)'].corr(user_country_stats['stats_mean_score'])
2 correlation_with_mobile = user_country_stats['Mobile cellular subscriptions (per 100 people)'].corr(user_country_stats['stats_mean_score'])
3 correlation_with_server = user_country_stats['Secure Internet servers (per 1 million people)'].corr(user_country_stats['stats_mean_score'])
4
5 print(f"Correlation between Fixed broadband subscriptions and Mean Score: {correlation_with_days_broadband}")
6 print(f"Correlation between Mobile cellular subscriptions and Mean Score: {correlation_with_mobile}")
7 print(f"Correlation between Secure Internet servers and Mean Score: {correlation_with_server}")

➡ Correlation between Fixed broadband subscriptions and Mean Score: -0.10001294168824629
Correlation between Mobile cellular subscriptions and Mean Score: 0.009402702266438736
Correlation between Secure Internet servers and Mean Score: 0.005194191434505355

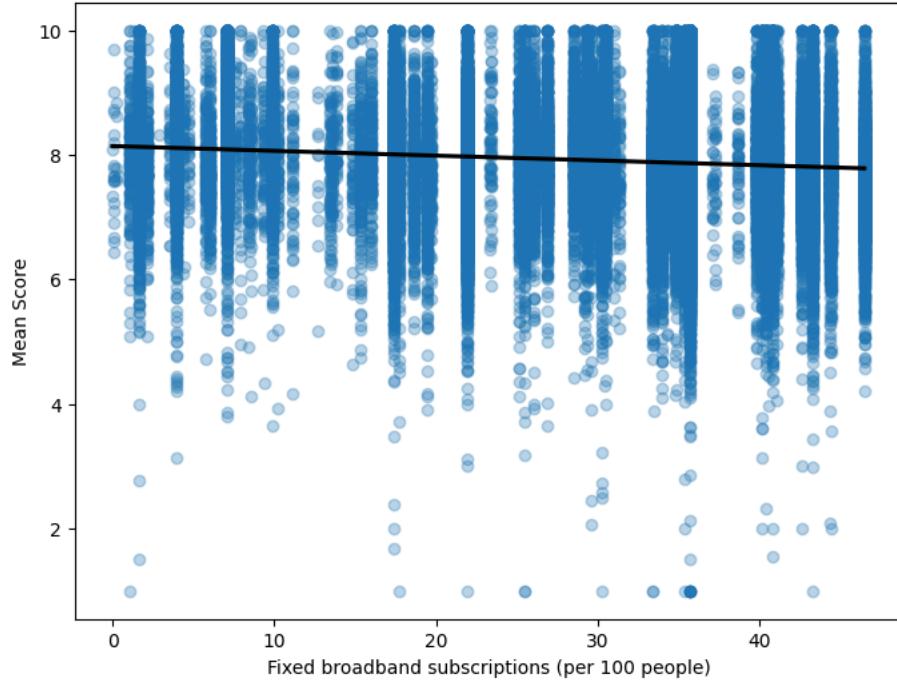
1 plt.figure(figsize=(8, 6))
2 sns.regplot(x='Fixed broadband subscriptions (per 100 people)', y='stats_mean_score', data=user_country_stats, scatter_kws={'alpha': 0.5})
3 plt.title('Fixed broadband subscriptions (per 100 people) vs. Mean Score')
4 plt.xlabel('Fixed broadband subscriptions (per 100 people)')
5 plt.ylabel('Mean Score')
6 plt.show()
7
8 plt.figure(figsize=(8, 6))
9 sns.regplot(x='Mobile cellular subscriptions (per 100 people)', y='stats_mean_score', data=user_country_stats, scatter_kws={'alpha': 0.5})
10 plt.title('Mobile cellular subscriptions (per 100 people) vs. Mean Score')
11 plt.xlabel('Mobile cellular subscriptions (per 100 people)')

```

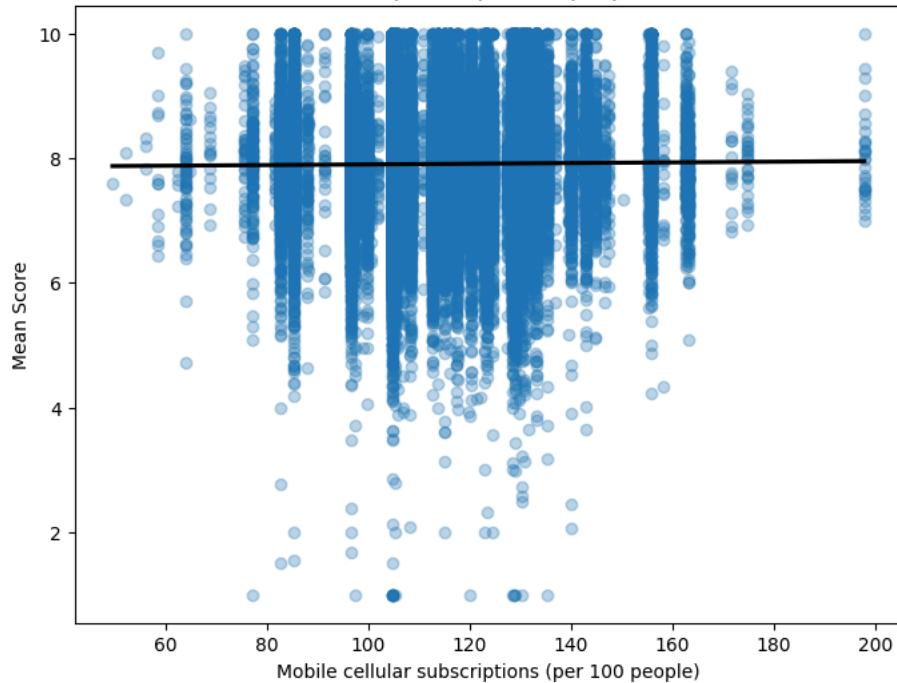
```
12 plt.ylabel('Mean Score')
13 plt.show()
14
15 plt.figure(figsize=(8, 6))
16 sns.regplot(x='Secure Internet servers (per 1 million people)', y='stats_mean_score', data=user_country_stats, scatter_kws={'alpha':
17 plt.title('Secure Internet servers (per 1 million people) vs. Mean Score')
18 plt.xlabel('Secure Internet servers (per 1 million people)')
19 plt.ylabel('Mean Score')
20 plt.show()
```

[X]

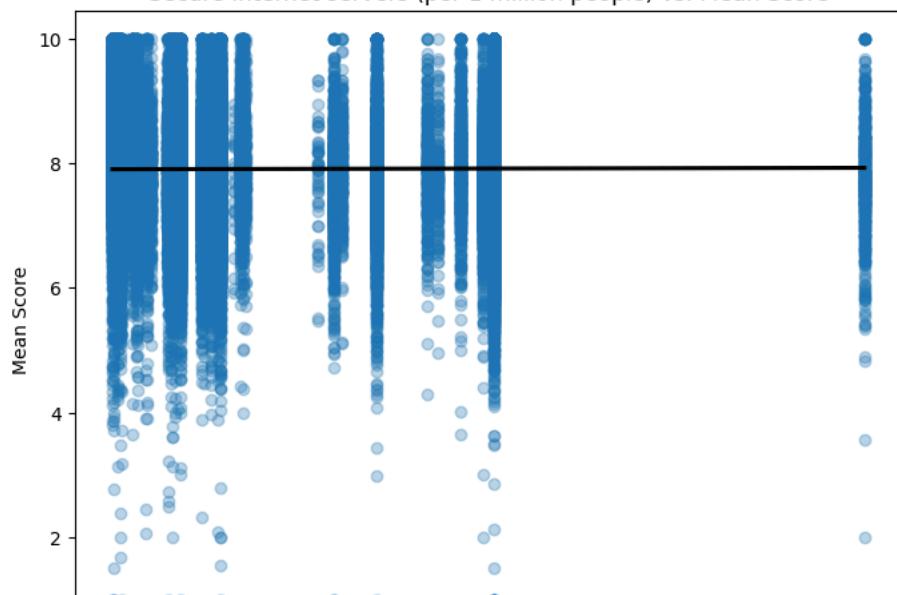
Fixed broadband subscriptions (per 100 people) vs. Mean Score

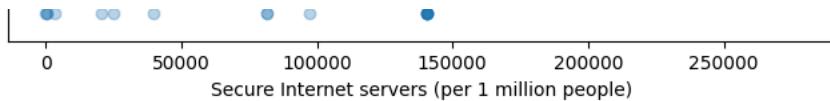


Mobile cellular subscriptions (per 100 people) vs. Mean Score



Secure Internet servers (per 1 million people) vs. Mean Score





```

1 tmp = user_country_stats.sample(n=500, random_state=42).copy()
2
3 tmp['broadband_category'] = pd.qcut(tmp['Fixed broadband subscriptions (per 100 people)'], 3, labels=['low', 'medium', 'high'])
4 tmp['mobile_category'] = pd.qcut(tmp['Mobile cellular subscriptions (per 100 people)'], 3, labels=['low', 'medium', 'high'])
5 tmp['server_category'] = pd.qcut(tmp['Secure Internet servers (per 1 million people)'], 3, labels=['low', 'medium', 'high'])
6 tmp['mean_score_category'] = pd.qcut(tmp['stats_mean_score'], 3, labels=['low', 'medium', 'high'])
7
8 contingency_table_broadband = pd.crosstab(tmp['broadband_category'], tmp['mean_score_category'])
9
10 contingency_table_mobile = pd.crosstab(tmp['mobile_category'], tmp['mean_score_category'])
11
12 contingency_table_server = pd.crosstab(tmp['server_category'], tmp['mean_score_category'])
13
14 chi2_stat_broadband, p_val_broadband, dof_broadband, expected_broadband = chi2_contingency(contingency_table_broadband)
15
16 chi2_stat_mobile, p_val_mobile, dof_mobile, expected_mobile = chi2_contingency(contingency_table_mobile)
17
18 chi2_stat_server, p_val_server, dof_server, expected_server = chi2_contingency(contingency_table_server)
19
20 print(f"Chi-square test for Fixed broadband subscriptions (per 100 people) and Mean Score:")
21 print(f"Chi-square statistic: {chi2_stat_broadband}")
22 print(f"P-value: {p_val_broadband}")
23
24 print(f"\nChi-square test for Mobile cellular subscriptions (per 100 people) and Mean Score:")
25 print(f"Chi-square statistic: {chi2_stat_mobile}")
26 print(f"P-value: {p_val_mobile}")
27
28 print(f"\nChi-square test for Secure Internet servers (per 1 million people) and Mean Score:")
29 print(f"Chi-square statistic: {chi2_stat_server}")
30 print(f"P-value: {p_val_server}")

→ Chi-square test for Fixed broadband subscriptions (per 100 people) and Mean Score:
Chi-square statistic: 9.162235472910053
P-value: 0.05716984705773421

Chi-square test for Mobile cellular subscriptions (per 100 people) and Mean Score:
Chi-square statistic: 3.653328246357328
P-value: 0.45495040686979304

Chi-square test for Secure Internet servers (per 1 million people) and Mean Score:
Chi-square statistic: 9.243525206249089
P-value: 0.055292535695738776

```

Conclusion: We fail to reject the null hypothesis and conclude that there is no significant association

H4: Higher national happiness correlates with higher average 'User Mean Score' from users in that country.

```

1 correlation = user_country_stats['Happiness Score'].corr(user_country_stats['stats_mean_score'])
2
3 print(f"Correlation between 'Happiness Score and Mean Score: {correlation}")

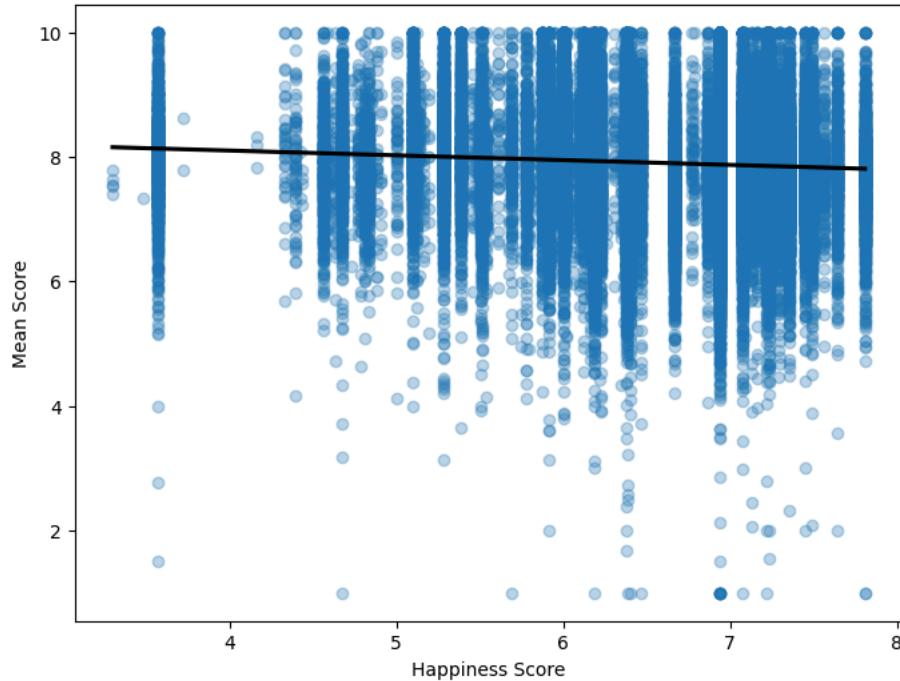
→ Correlation between 'Happiness Score and Mean Score: -0.06687589793370768

1 plt.figure(figsize=(8, 6))
2 sns.regplot(x='Happiness Score', y='stats_mean_score', data=user_country_stats, scatter_kws={'alpha':0.3}, line_kws={'color': 'black'})
3 plt.title('Happiness Score vs. Mean Score')
4 plt.xlabel('Happiness Score')
5 plt.ylabel('Mean Score')
6 plt.show()

```



Happiness Score vs. Mean Score



```
1 tmp = user_country_stats.sample(n=500, random_state=42).copy()
2
3 tmp['happiness_category'] = pd.qcut(tmp['Happiness Score'], 3, labels=['low', 'medium', 'high'])
4 tmp['mean_score_category'] = pd.qcut(tmp['stats_mean_score'], 3, labels=['low', 'medium', 'high'])
5
6 contingency_table = pd.crosstab(tmp['happiness_category'], tmp['mean_score_category'])
7
8 chi2_stat, p_val, dof, expected = chi2_contingency(contingency_table)
9
10 print(f"Chi-square test for Happiness Score and Mean Score:")
11 print(f"Chi-square statistic: {chi2_stat}")
12 print(f"P-value: {p_val}")

→ Chi-square test for Happiness Score and Mean Score:
Chi-square statistic: 5.702573456003083
P-value: 0.22248859870279736
```

Conclusion: We fail to reject the null hypothesis and conclude that there is no significant association

H5: Male and female users show differences in their anime watching habits, such as the total time spent watching anime or the average scores they give to anime.

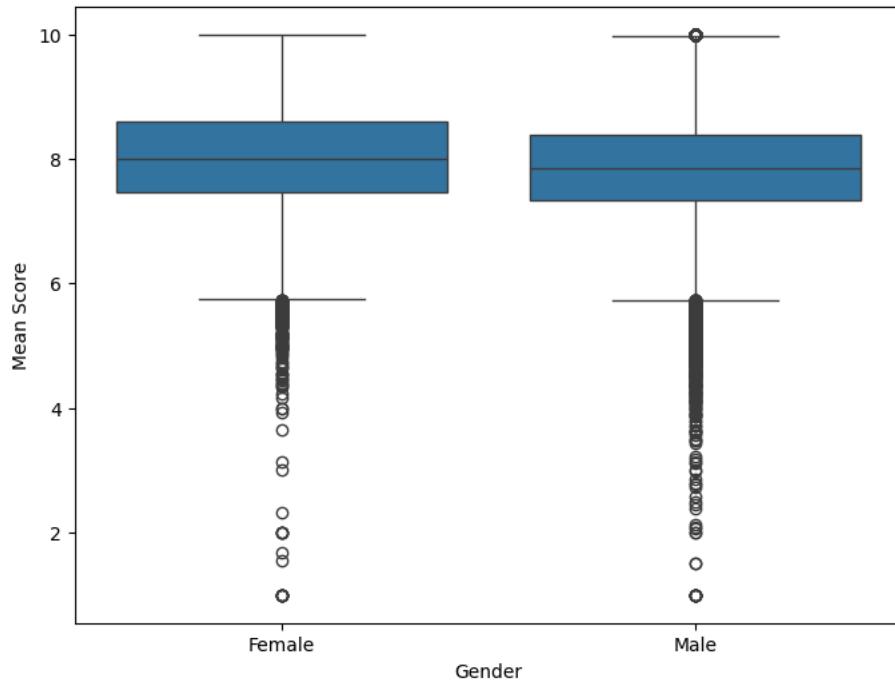
```
1 correlation_with_days_spent_watching = user_country_stats['user_days_spent_watching'].corr(user_country_stats['gender_numeric'])
2 correlation_with_mean_score = user_country_stats['stats_mean_score'].corr(user_country_stats['gender_numeric'])
3
4 print(f"Correlation between Gender and User Days Spent Watching: {correlation_with_days_spent_watching}")
5 print(f"Correlation between Gender and Mean Score: {correlation_with_mean_score}")

→ Correlation between Gender and User Days Spent Watching: -0.17631070389647777
Correlation between Gender and Mean Score: 0.09952580830135989
```

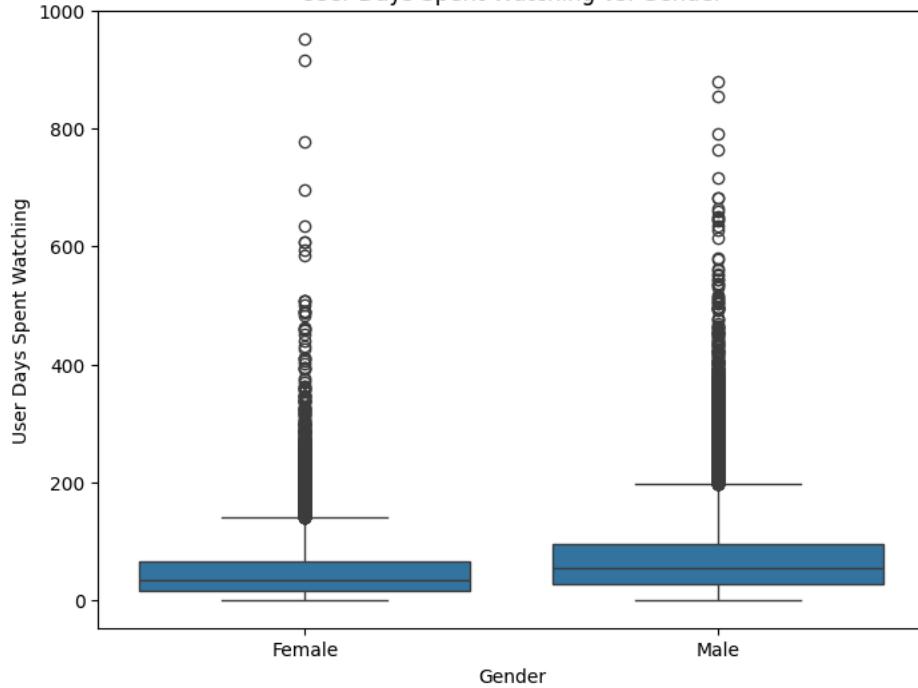
```
1 plt.figure(figsize=(8, 6))
2 sns.boxplot(x='gender', y='stats_mean_score', data=user_country_stats)
3 plt.title('Mean Score vs. Gender')
4 plt.xlabel('Gender')
5 plt.ylabel('Mean Score')
6 plt.show()
7
8 plt.figure(figsize=(8, 6))
9 sns.boxplot(x='gender', y='user_days_spent_watching', data=user_country_stats)
10 plt.title('User Days Spent Watching vs. Gender')
11 plt.xlabel('Gender')
12 plt.ylabel('User Days Spent Watching')
13 plt.show()
```

[x]

Mean Score vs. Gender



User Days Spent Watching vs. Gender



```

1 tmp = user_country_stats.sample(n=500, random_state=42).copy()
2
3 tmp['mean_score_category'] = pd.qcut(tmp['stats_mean_score'], 3, labels=['low', 'medium', 'high'])
4 tmp['days_spent_watching_category'] = pd.qcut(tmp['user_days_spent_watching'], 3, labels=['low', 'medium', 'high'])
5
6 contingency_table_mean = pd.crosstab(tmp['gender'], tmp['mean_score_category'])
7 contingency_table_days_watching = pd.crosstab(tmp['gender'], tmp['days_spent_watching_category'])
8
9 chi2_stat_mean, p_val_mean, dof_mean, expected_mean = chi2_contingency(contingency_table_mean)
10
11 chi2_stat_days_watching, p_val_days_watching, dof_days_watching, expected_days_watching = chi2_contingency(contingency_table_days_watching)
12
13 print(f"\nChi-square test for Mean Score and Gender:")
14 print(f"Chi-square statistic: {chi2_stat_mean}")
15 print(f"P-value: {p_val_mean}")
16
17 print(f"\nChi-square test for User Days Spent Watching and Gender:")
18 print(f"Chi-square statistic: {chi2_stat_days_watching}")
19 print(f"P-value: {p_val_days_watching}")

```

[x] Chi-square test for Mean Score and Gender:
 Chi-square statistic: 0.8739665445468503
 P-value: 0.6459822371344817

```
Chi-square test for User Days Spent Watching and Gender:  
Chi-square statistic: 12.10587008445579  
P-value: 0.002350951727874106
```

Conclusion: We fail to reject the null hypothesis for average score and conclude that there is no significant association between gender and the average score users give to anime. However, we reject the null hypothesis for user days spent watching and conclude that there is a significant association between gender and the total time spent watching anime.

H6: Users from different regional indicators (e.g., Western Europe, Sub-Saharan Africa) exhibit distinct anime watching behaviors, such as mean score, watching time, or completed anime.

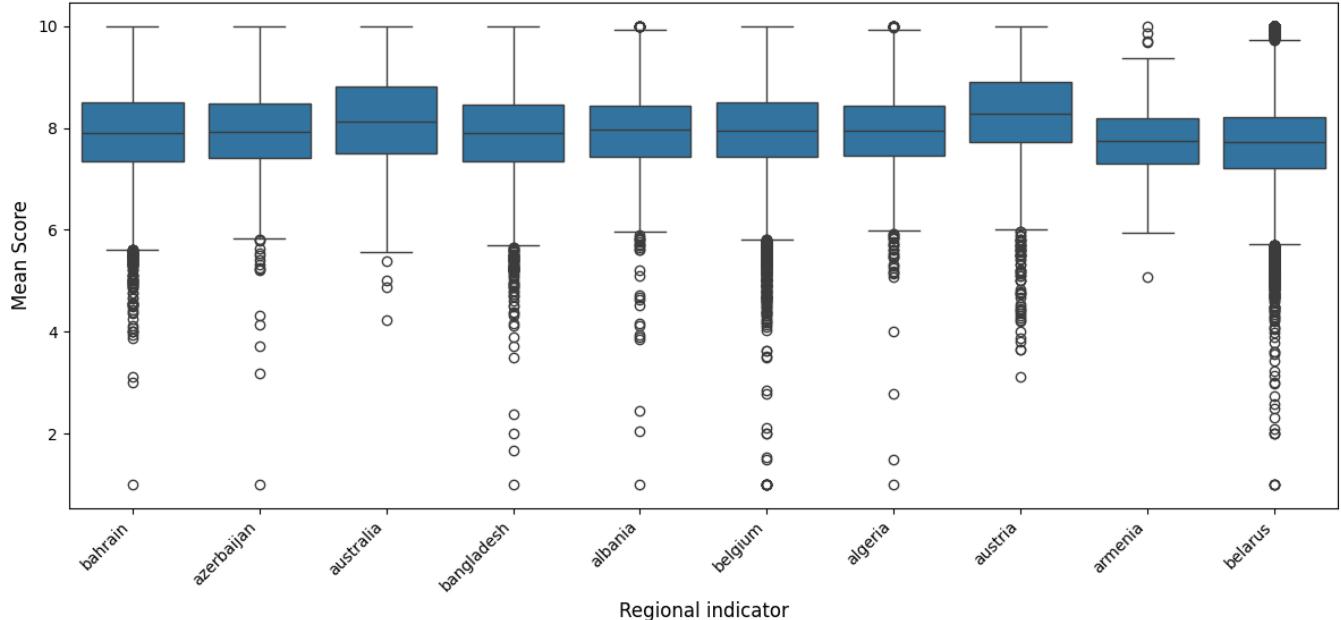
```
1 correlation_with_days_spent_watching = user_country_stats['user_days_spent_watching'].corr(user_country_stats['Regional indicator'])  
2 correlation_with_mean_score = user_country_stats['stats_mean_score'].corr(user_country_stats['Regional indicator'])  
3  
4 print(f"Correlation between Regional indicator and User Days Spent Watching: {correlation_with_days_spent_watching}")  
5 print(f"Correlation between Regional indicator and Mean Score: {correlation_with_mean_score}")
```

→ Correlation between Regional indicator and User Days Spent Watching: 0.011590359946842354
Correlation between Regional indicator and Mean Score: -0.0521957966286115

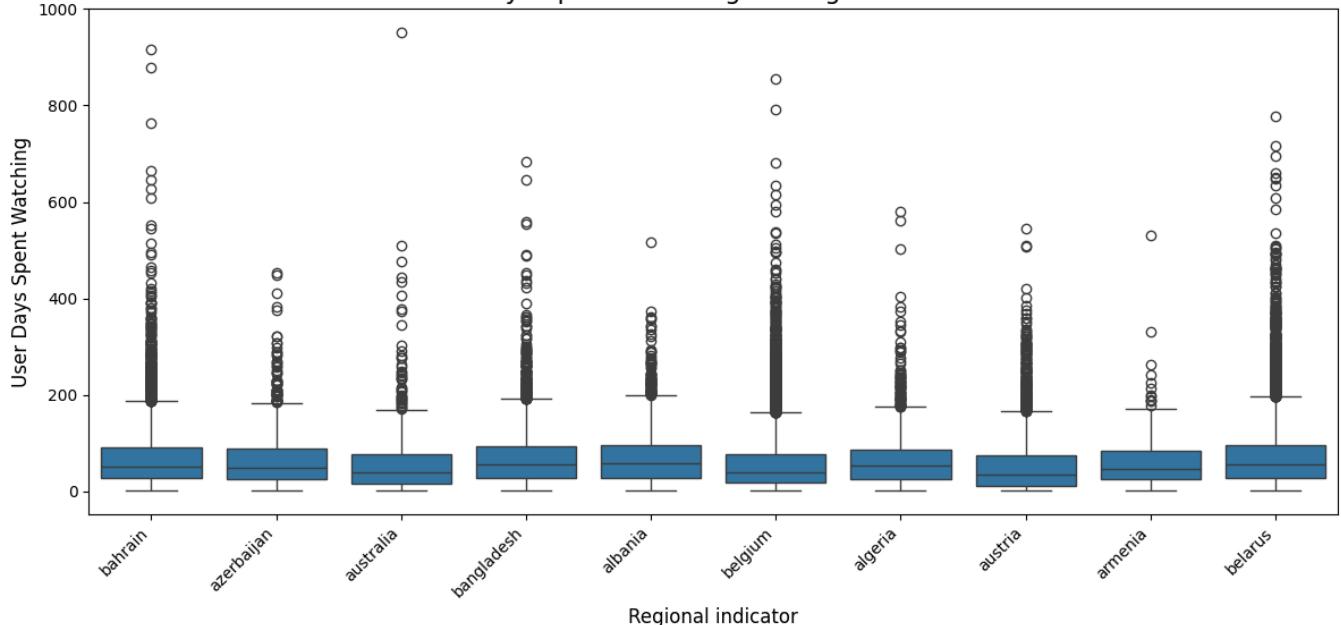
```
1 original_labels = label_encoder.inverse_transform(user_country_stats['Regional indicator'].unique())  
2  
3 plt.figure(figsize=(12, 6))  
4 sns.boxplot(x='Regional indicator', y='stats_mean_score', data=user_country_stats)  
5 plt.title('Mean Score vs. Regional indicator', fontsize=16)  
6 plt.xlabel('Regional indicator', fontsize=12)  
7 plt.ylabel('Mean Score', fontsize=12)  
8 plt.xticks(ticks=np.arange(len(original_labels)), labels=original_labels, rotation=45, ha='right')  
9 plt.tight_layout()  
10 plt.show()  
11  
12 plt.figure(figsize=(12, 6))  
13 sns.boxplot(x='Regional indicator', y='user_days_spent_watching', data=user_country_stats)  
14 plt.title('User Days Spent Watching vs. Regional indicator', fontsize=16)  
15 plt.xlabel('Regional indicator', fontsize=12)  
16 plt.ylabel('User Days Spent Watching', fontsize=12)  
17 plt.xticks(ticks=np.arange(len(original_labels)), labels=original_labels, rotation=45, ha='right')  
18 plt.tight_layout()  
19 plt.show()
```

[

Mean Score vs. Regional indicator



User Days Spent Watching vs. Regional indicator



```

1 tmp = user_country_stats.sample(n=500, random_state=42).copy()
2
3 tmp['mean_score_category'] = pd.qcut(tmp['stats_mean_score'], 3, labels=['low', 'medium', 'high'])
4 tmp['days_spent_watching_category'] = pd.qcut(tmp['user_days_spent_watching'], 3, labels=['low', 'medium', 'high'])
5
6 contingency_table_mean = pd.crosstab(tmp['Regional indicator'], tmp['mean_score_category'])
7 contingency_table_days_watching = pd.crosstab(tmp['Regional indicator'], tmp['days_spent_watching_category'])
8
9 chi2_stat_mean, p_val_mean, dof_mean, expected_mean = chi2_contingency(contingency_table_mean)
10
11 chi2_stat_days_watching, p_val_days_watching, dof_days_watching, expected_days_watching = chi2_contingency(contingency_table_days_watching)
12
13 print(f"Chi-square test for Mean Score and Regional indicator:")
14 print(f"Chi-square statistic: {chi2_stat_mean}")
15 print(f"P-value: {p_val_mean}")
16
17 print(f"\nChi-square test for User Days Spent Watching and Regional indicator:")

```

```
18 print(f"Chi-square statistic: {chi2_stat_days_watching}")
19 print(f"P-value: {p_val_days_watching}")
```

```
→ Chi-square test for Mean Score and Regional indicator:
Chi-square statistic: 30.314558157556036
P-value: 0.03449826829611721
```

```
Chi-square test for User Days Spent Watching and Regional indicator:
Chi-square statistic: 39.6734092927294
P-value: 0.0023117573694563078
```

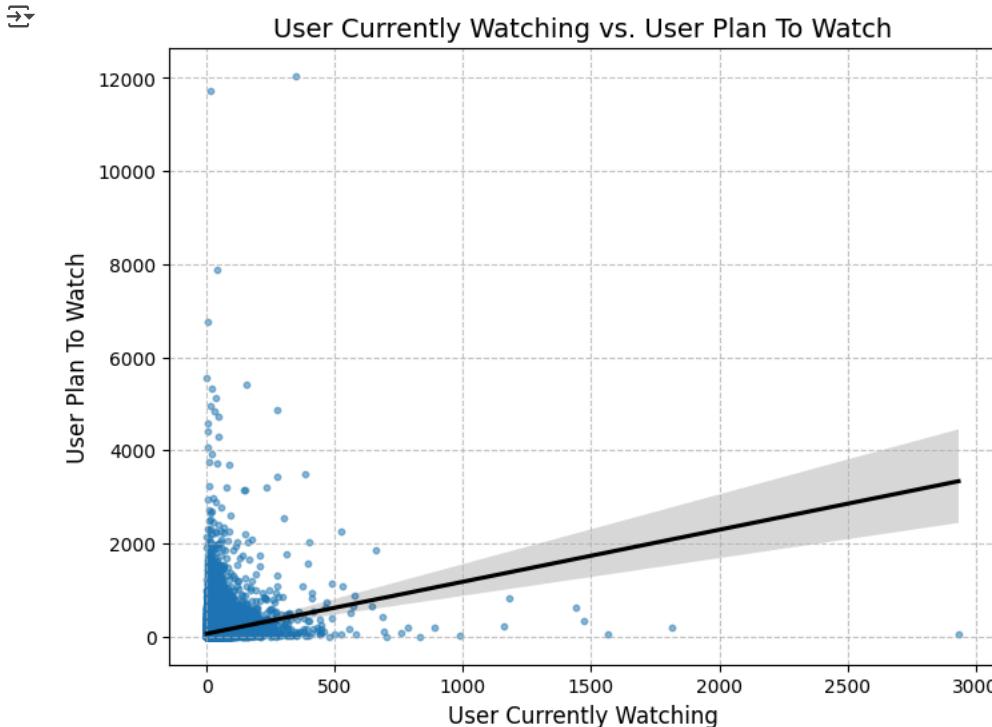
Conclusion: We reject the null hypothesis and conclude that there is a significant association

H7: There's a positive correlation between the number of anime a user is currently watching and the number they plan to watch, suggesting active viewers tend to plan future viewing.

```
1 correlation = user_country_stats['user_watching'].corr(user_country_stats['user_plantowatch'])
2
3 print(f"Correlation between User Watching and User Plan to Watch: {correlation}")
```

```
→ Correlation between User Watching and User Plan to Watch: 0.20253436731703003
```

```
1 plt.figure(figsize=(8, 6))
2 sns.regplot(x='user_watching', y='user_plantowatch', data=user_country_stats, scatter_kws={'alpha':0.5, 's': 10}, line_kws={'color':
3 plt.title('User Currently Watching vs. User Plan To Watch', fontsize=14)
4 plt.xlabel('User Currently Watching', fontsize=12)
5 plt.ylabel('User Plan To Watch', fontsize=12)
6 plt.grid(True, linestyle='--', alpha=0.7)
7 plt.show()
```



```
1 tmp = user_country_stats.sample(n=500, random_state=42).copy()
2
3 tmp['user_watching_category'] = pd.qcut(tmp['user_watching'], 3, labels=['low', 'medium', 'high'])
4 tmp['user_plantowatch_category'] = pd.qcut(tmp['user_plantowatch'], 3, labels=['low', 'medium', 'high'])
5
6 contingency_table = pd.crosstab(tmp['user_watching_category'], tmp['user_plantowatch_category'])
7
8 chi2_stat, p_val, dof, expected = chi2_contingency(contingency_table)
9
10 print(f"Chi-square test for User Currently Watching and User Plan To Watch:")
11 print(f"Chi-square statistic: {chi2_stat}")
12 print(f"P-value: {p_val}")
```

```
→ Chi-square test for User Currently Watching and User Plan To Watch:
Chi-square statistic: 76.58364062372213
P-value: 9.212465740454514e-16
```

Conclusion: We reject the null hypothesis and conclude that there is a significant association

✓ Preparation For Machine Learning

```
1 user_country_stats.info()

→ <class 'pandas.core.frame.DataFrame'>
Index: 60924 entries, 0 to 65827
Data columns (total 68 columns):
 #   Column           Non-Null Count Dtype
 ---  -- 
 0   username        60924 non-null  object
 1   user_watching   60924 non-null  int64
 2   user_completed  60924 non-null  int64
 3   user_onhold    60924 non-null  int64
 4   user_dropped   60924 non-null  int64
 5   user_plantowatch 60924 non-null  int64
 6   user_days_spent_watching 60924 non-null  float64
 7   gender          60924 non-null  object
 8   birth_date      60924 non-null  datetime64[ns]
 9   join_date       60924 non-null  datetime64[ns]
10  last_online     60924 non-null  datetime64[ns]
11  stats_mean_score 60924 non-null  float64
12  stats_rewatched 60924 non-null  float64
13  stats_episodes   60924 non-null  int64
14  age              60924 non-null  int64
15  Country          60924 non-null  int64
16  Regional indicator 60924 non-null  int64
17  Happiness Score  60924 non-null  float64
18  Logged GDP per capita 60924 non-null  float64
19  Social support   60924 non-null  float64
20  Healthy life expectancy 60924 non-null  float64
21  Freedom to make life choices 60924 non-null  float64
22  Generosity       60924 non-null  float64
23  Perceptions of corruption 60924 non-null  float64
24  Ladder score in Dystopia 60924 non-null  float64
25  Explained by: Log GDP per capita 60924 non-null  float64
26  Explained by: Social support 60924 non-null  float64
27  Explained by: Healthy life expectancy 60924 non-null  float64
28  Explained by: Freedom to make life choices 60924 non-null  float64
29  Explained by: Generosity 60924 non-null  float64
30  Explained by: Perceptions of corruption 60924 non-null  float64
31  Dystopia + residual 60924 non-null  float64
32  High-technology exports (current US$) 60924 non-null  float64
33  Forest area (% of land area) 60924 non-null  float64
34  GDP per capita growth (annual %) 60924 non-null  float64
35  Urban population (% of total population) 60924 non-null  float64
36  Net migration 60924 non-null  float64
37  GNI per capita, PPP (current international $) 60924 non-null  float64
38  Exports of goods and services (% of GDP) 60924 non-null  float64
39  Agriculture, forestry, and fishing, value added (% of GDP) 60924 non-null  float64
40  Access to electricity (% of population) 60924 non-null  float64
41  Population growth (annual %) 60924 non-null  float64
42  Renewable electricity output (% of total electricity output) 60924 non-null  float64
43  Electric power consumption (kWh per capita) 60924 non-null  float64
44  Renewable internal freshwater resources per capita (cubic meters) 60924 non-null  float64
45  Population density (people per sq. km of land area) 60924 non-null  float64
46  Mobile cellular subscriptions (per 100 people) 60924 non-null  float64
47  Fixed broadband subscriptions (per 100 people) 60924 non-null  float64
48  Secure Internet servers (per 1 million people) 60924 non-null  float64
49  Mortality rate, under-5 (per 1,000 live births) 60924 non-null  float64
50  Employment to population ratio, 15+, total (%) (modeled ILO estimate) 60924 non-null  float64
51  Death rate, crude (per 1,000 people) 60924 non-null  float64
52  PM2.5 air pollution, mean annual exposure (micrograms per cubic meter) 60924 non-null  float64
```

This code block removes two columns, `username` and `gender`, from the `user_country_stats` DataFrame.

```
1 user_country_stats = user_country_stats.drop(labels="username", axis=1)
2 user_country_stats = user_country_stats.drop(labels="gender", axis=1)
```

This code block transforms the date columns in the `user_country_stats` DataFrame into a numerical representation.

```
1 user_country_stats['join_date'] = (user_country_stats['join_date'] - pd.Timestamp('1970-01-01')).dt.days
2 user_country_stats['last_online'] = (user_country_stats['last_online'] - pd.Timestamp('1970-01-01')).dt.days
3 user_country_stats['birth_date'] = (user_country_stats['birth_date'] - pd.Timestamp('1970-01-01')).dt.days
```

The purpose of this function is to evaluate the performance of a machine learning model on both training and testing datasets. It calculates several common evaluation metrics and returns them in a dictionary.

```
1 def show_scores(model,
2                  train_features,
3                  train_labels,
```

```

4         test_features,
5         test_labels):
6
7     np.random.seed(42)
8     train_preds = model.predict(X=train_features)
9     test_preds = model.predict(X=test_features)
10
11    train_preds_clipped = np.maximum(0, train_preds)
12    test_preds_clipped = np.maximum(0, test_preds)
13
14    scores = {"Training MAE": mean_absolute_error(y_true=train_labels,
15                                              y_pred=train_preds),
16              "Test MAE": mean_absolute_error(y_true=test_labels,
17                                              y_pred=test_preds),
18              "Training RMSLE": root_mean_squared_log_error(y_true=train_labels,
19                                              y_pred=train_preds_clipped),
20              "Test RMSLE": root_mean_squared_log_error(y_true=test_labels,
21                                              y_pred=test_preds_clipped),
22              "Training R^2": model.score(X=train_features,
23                                           y=train_labels),
24              "Test R^2": model.score(X=test_features,
25                                           y=test_labels)}
26
27    return scores

```

▼ Machine Learning Model (Mean Score)

This code block prepares the sample data for training a machine learning model to predict the 'stats_mean_score'.

```

1 sample = user_country_stats.sample(n=1000, random_state=42)
2
3 x_sample = sample.drop(labels="stats_mean_score", axis=1)
4 y_sample = sample["stats_mean_score"]

```

This code block splits the data into training and testing sets.

```
1 x_sample_train, x_sample_test, y_sample_train, y_sample_test = train_test_split(x_sample, y_sample, test_size=0.2, random_state=42)
```

This code block splits the main dataset into features and the target variable, and then further divides these into training and testing sets.

```

1 X = user_country_stats.drop('stats_mean_score', axis=1)
2 y = user_country_stats['stats_mean_score']
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

This code block focuses on cleaning the column names of training and testing feature sets, X_train and X_test.

```

1 X_train.columns = [re.sub('[^A-Za-z0-9_]+', '_', col) for col in X_train.columns]
2 X_test.columns = [re.sub('[^A-Za-z0-9_]+', '_', col) for col in X_test.columns]

```

This code block trains a baseline machine learning model, specifically a Random Forest Regressor, on a sample of the sample data and then evaluates its performance.

```

1 np.random.seed(42)
2
3 base_model_sample = RandomForestRegressor(n_jobs=-1)
4
5 base_model_sample.fit(X=x_sample_train, y=y_sample_train)
6
7 base_model_sample_scores = show_scores(model=base_model_sample,
8                                         train_features=x_sample_train,
9                                         train_labels=y_sample_train,
10                                        test_features=x_sample_test,
11                                        test_labels=y_sample_test)
12
13 print(f"Model score on {len(sample)} samples:")
14 base_model_sample_scores

```

→ Model score on 1000 samples:

```

{'Training MAE': 0.24717012500000024,
 'Test MAE': 0.6344839999999998,
 'Training RMSLE': 0.038674498687953,
 'Test RMSLE': 0.1017272274203456,

```

```
'Training R^2': 0.8779475785877052,
'Test R^2': -0.0063945933095461704}
```

This code block tunes the hyperparameters of a Random Forest Regressor model using a technique called Randomized Search Cross-Validation.

```
1 np.random.seed(42)
2
3 rf_grid = {"n_estimators": np.arange(10, 200, 10),
4             "max_depth": [None, 10, 20],
5             "min_samples_split": np.arange(2, 10, 1),
6             "min_samples_leaf": np.arange(1, 10, 1),
7             "max_features": [0.5, 1.0, "sqrt"], }
8
9 sample_rs_model = RandomizedSearchCV(estimator=RandomForestRegressor(),
10                                         param_distributions=rf_grid,
11                                         n_iter=20,
12                                         cv=3,
13                                         verbose=3)
14
15 sample_rs_model.fit(X=x_sample_train,
16                       y=y_sample_train)
```

```

Fitting 3 folds for each of 20 candidates, totalling 60 fits
[CV 1/3] END max_depth=10, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=130; score=0.145 total time=
[CV 2/3] END max_depth=10, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=130; score=0.110 total time=
[CV 3/3] END max_depth=10, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=130; score=0.158 total time=
[CV 1/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=7, n_estimators=60; score=0.124 total time=
[CV 2/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=7, n_estimators=60; score=0.130 total time=
[CV 3/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=7, n_estimators=60; score=0.156 total time=
[CV 1/3] END max_depth=10, max_features=0.5, min_samples_leaf=9, min_samples_split=5, n_estimators=140; score=0.128 total time=
[CV 2/3] END max_depth=10, max_features=0.5, min_samples_leaf=9, min_samples_split=5, n_estimators=140; score=0.134 total time=
[CV 3/3] END max_depth=10, max_features=0.5, min_samples_leaf=9, min_samples_split=5, n_estimators=140; score=0.172 total time=
[CV 1/3] END max_depth=10, max_features=0.5, min_samples_leaf=8, min_samples_split=3, n_estimators=50; score=0.113 total time=
[CV 2/3] END max_depth=10, max_features=0.5, min_samples_leaf=8, min_samples_split=3, n_estimators=50; score=0.127 total time=
[CV 3/3] END max_depth=10, max_features=0.5, min_samples_leaf=8, min_samples_split=3, n_estimators=50; score=0.173 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=7, min_samples_split=7, n_estimators=140; score=0.144 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=7, min_samples_split=7, n_estimators=140; score=0.128 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=7, min_samples_split=7, n_estimators=140; score=0.134 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=3, n_estimators=180; score=0.144 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=3, n_estimators=180; score=0.132 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=3, n_estimators=180; score=0.155 total time=
[CV 1/3] END max_depth=10, max_features=1.0, min_samples_leaf=2, min_samples_split=7, n_estimators=160; score=0.139 total time=
[CV 2/3] END max_depth=10, max_features=1.0, min_samples_leaf=2, min_samples_split=7, n_estimators=160; score=0.110 total time=
[CV 3/3] END max_depth=10, max_features=1.0, min_samples_leaf=2, min_samples_split=7, n_estimators=160; score=0.178 total time=
[CV 1/3] END max_depth=10, max_features=1.0, min_samples_leaf=6, min_samples_split=3, n_estimators=150; score=0.127 total time=
[CV 2/3] END max_depth=10, max_features=1.0, min_samples_leaf=6, min_samples_split=3, n_estimators=150; score=0.131 total time=
[CV 3/3] END max_depth=10, max_features=1.0, min_samples_leaf=6, min_samples_split=3, n_estimators=150; score=0.165 total time=
[CV 1/3] END max_depth=None, max_features=0.5, min_samples_leaf=4, min_samples_split=2, n_estimators=110; score=0.127 total time=
[CV 2/3] END max_depth=None, max_features=0.5, min_samples_leaf=4, min_samples_split=2, n_estimators=110; score=0.122 total time=
[CV 3/3] END max_depth=None, max_features=0.5, min_samples_leaf=4, min_samples_split=2, n_estimators=110; score=0.160 total time=
[CV 1/3] END max_depth=10, max_features=0.5, min_samples_leaf=3, min_samples_split=2, n_estimators=190; score=0.149 total time=
[CV 2/3] END max_depth=10, max_features=0.5, min_samples_leaf=3, min_samples_split=2, n_estimators=190; score=0.125 total time=
[CV 3/3] END max_depth=10, max_features=0.5, min_samples_leaf=3, min_samples_split=2, n_estimators=190; score=0.172 total time=
[CV 1/3] END max_depth=10, max_features=1.0, min_samples_leaf=1, min_samples_split=7, n_estimators=120; score=0.152 total time=
[CV 2/3] END max_depth=10, max_features=1.0, min_samples_leaf=1, min_samples_split=7, n_estimators=120; score=0.115 total time=
[CV 3/3] END max_depth=10, max_features=1.0, min_samples_leaf=1, min_samples_split=7, n_estimators=120; score=0.184 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=5, min_samples_split=6, n_estimators=90; score=0.161 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=5, min_samples_split=6, n_estimators=90; score=0.122 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=5, min_samples_split=6, n_estimators=90; score=0.152 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=20; score=0.118 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=20; score=0.120 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=20; score=0.135 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=160; score=0.135 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=160; score=0.108 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=160; score=0.154 total time=
[CV 1/3] END max_depth=20, max_features=0.5, min_samples_leaf=1, min_samples_split=8, n_estimators=10; score=0.105 total time=
[CV 2/3] END max_depth=20, max_features=0.5, min_samples_leaf=1, min_samples_split=8, n_estimators=10; score=0.111 total time=
[CV 3/3] END max_depth=20, max_features=0.5, min_samples_leaf=1, min_samples_split=8, n_estimators=10; score=0.150 total time=
[CV 1/3] END max_depth=None, max_features=1.0, min_samples_leaf=3, min_samples_split=2, n_estimators=140; score=0.145 total time=
[CV 2/3] END max_depth=None, max_features=1.0, min_samples_leaf=3, min_samples_split=2, n_estimators=140; score=0.128 total time=
[CV 3/3] END max_depth=None, max_features=1.0, min_samples_leaf=3, min_samples_split=2, n_estimators=140; score=0.157 total time=
[CV 1/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=2, n_estimators=100; score=0.142 total time=
[CV 2/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=2, n_estimators=100; score=0.133 total time=
[CV 3/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=2, n_estimators=100; score=0.164 total time=
[CV 1/3] END max_depth=20, max_features=1.0, min_samples_leaf=7, min_samples_split=7, n_estimators=10; score=0.100 total time=
[CV 2/3] END max_depth=20, max_features=1.0, min_samples_leaf=7, min_samples_split=7, n_estimators=10; score=0.067 total time=
[CV 3/3] END max_depth=20, max_features=1.0, min_samples_leaf=7, min_samples_split=7, n_estimators=10; score=0.164 total time=
[CV 1/3] END max_depth=10, max_features=sqrt, min_samples_leaf=1, min_samples_split=7, n_estimators=150; score=0.163 total time=
[CV 2/3] END max_depth=10, max_features=sqrt, min_samples_leaf=1, min_samples_split=7, n_estimators=150; score=0.147 total time=
[CV 3/3] END max_depth=10, max_features=sqrt, min_samples_leaf=1, min_samples_split=7, n_estimators=150; score=0.158 total time=
[CV 1/3] END max_depth=None, max_features=1.0, min_samples_leaf=8, min_samples_split=2, n_estimators=20; score=0.100 total time=
[CV 2/3] END max_depth=None, max_features=1.0, min_samples_leaf=8, min_samples_split=2, n_estimators=20; score=0.116 total time=
[CV 3/3] END max_depth=None, max_features=1.0, min_samples_leaf=8, min_samples_split=2, n_estimators=20; score=0.180 total time=

```

RandomizedSearchCV



`best_estimator_:`
RandomForestRegressor

RandomForestRegressor

```
1 sample_rs_model.best_params_
```

```

2 { 'n_estimators': np.int64(150),
  'min_samples_split': np.int64(7),
  'min_samples_leaf': np.int64(1),
  'max_features': 'sqrt',
  'max_depth': 10}

```

```

1 sample_rs_model_scores = show_scores(model=sample_rs_model,
2                                     train_features=x_sample_train,
3                                     train_labels=y_sample_train,
4                                     test_features=x_sample_test,
5                                     test_labels=y_sample_test)
6 sample_rs_model_scores

```

```

7 {'Training MAE': 0.4282388818111595,
8  'Test MAE': 0.5979108009171518,
9 }

```

```
'Training RMSLE': 0.06718667252319614,  
'Test RMSLE': 0.08965091362916926,  
'Training R^2': 0.626504320648437,  
'Test R^2': 0.175379145651261}
```

This code block performs a Grid Search to find the best hyperparameters for a RandomForestRegressor model using the sample training data.

```
1 np.random.seed(42)
2
3 rf_grid = {"n_estimators": np.arange(50, 150, 50),
4             "max_depth": [None, 10, 20],
5             "min_samples_split": np.arange(2, 8, 2),
6             "min_samples_leaf": np.arange(2, 6, 2),
7             "max_features": [0.5, "sqrt"],
8             "oob_score": [True, False]
9 }
10
11 sample_gs_model = GridSearchCV(estimator=RandomForestRegressor(n_jobs=-1),
12                                 param_grid=rf_grid,
13                                 cv=3,
14                                 verbose=3,
15                                 n_jobs=-1)
16
17 sample_gs_model.fit(X=x_sample_train,
18                      y=y_sample_train)
```

```
→ Fitting 3 folds for each of 144 candidates, totalling 432 fits
  ▶ GridSearchCV
    ⓘ ⓘ
  ▶ best_estimator_:
    RandomForestRegressor
      ▶ RandomForestRegressor ⓘ
```

```
1 sample_gs_model.best_params_
```

```
→ {'max_depth': 10,  
 'max_features': 0.5,  
 'min_samples_leaf': np.int64(4),  
 'min_samples_split': np.int64(2),  
 'n_estimators': np.int64(100),  
 'oob_score': False}
```

```
1 sample_gs_model_scores = show_scores(model=sample_rs_model,
2                                         train_features=x_sample_train,
3                                         train_labels=y_sample_train,
4                                         test_features=x_sample_test,
5                                         test_labels=y_sample_test)
6 sample_gs_model_scores
```

```
→ {'Training MAE': 0.4282388818111595,  
   'Test MAE': 0.5979108009171518,  
   'Training RMSLE': 0.06718667252319614  
   'Test RMSLE': 0.08965091362916926,  
   'Training R^2': 0.626504320648437,  
   'Test R^2': 0.175379145651261}
```

This code block sets up and runs a `HalvingRandomSearchCV` for a `RandomForestRegressor` model. The goal is to find the best combination of hyperparameters for the model efficiently.

```
16 sample_hrs_model.fit(X=x_sample_train,  
17                      y=y_sample_train)
```


HalvingRandomSearchCV



▶ RandomForestRegressor [?](#)



```

1 sample_hrs_model.best_params_
2
3 ➔ {'oob_score': False,
4     'n_estimators': np.int64(90),
5     'min_samples_split': np.int64(9),
6     'min_samples_leaf': np.int64(7),
7     'max_features': 'sqrt',
8     'max_depth': 10}

1 sample_hrs_model_scores = show_scores(model=sample_hrs_model,
2                                         train_features=x_sample_train,
3                                         train_labels=y_sample_train,
4                                         test_features=x_sample_test,
5                                         test_labels=y_sample_test)
6 sample_hrs_model_scores

7 ➔ {'Training MAE': 0.5158075035008145,
8     'Test MAE': 0.5807193498483294,
9     'Training RMSLE': 0.08305704803787987,
10    'Test RMSLE': 0.08779701839845318,
11    'Training R^2': 0.431613061749733,
12    'Test R^2': 0.21029860671427314}

```

This code block trains a baseline machine learning model, specifically a Random Forest Regressor, on the full training dataset and then evaluates its performance.

```

1 base_model = RandomForestRegressor(n_jobs=-1)
2
3 base_model.fit(X=X_train, y=y_train)
4
5 base_model_scores = show_scores(model=base_model,
6                                 train_features=X_train,
7                                 train_labels=y_train,
8                                 test_features=X_test,
9                                 test_labels=y_test)
10 base_model_scores

7 ➔ {'Training MAE': 0.2205965304991896,
8     'Test MAE': 0.5887088469429627,
9     'Training RMSLE': 0.038302907056537425,
10    'Test RMSLE': 0.09367731487555116,
11    'Training R^2': 0.8946714395636954,
12    'Test R^2': 0.2453869189700515}

```

This code block focuses on tuning the hyperparameters of a RandomForestRegressor model using Randomized Search Cross-Validation.

```

1 np.random.seed(42)
2
3 rf_grid = {"n_estimators": np.arange(10, 200, 10),
4            "max_depth": [None, 10, 20],
5            "min_samples_split": np.arange(2, 10, 1),
6            "min_samples_leaf": np.arange(1, 10, 1),
7            "max_features": [0.5, 1.0, "sqrt"], }
8
9 rs_model = RandomizedSearchCV(estimator=RandomForestRegressor(),
10                                param_distributions=rf_grid,
11                                n_iter=20,
12                                cv=3,
13                                verbose=3)
14
15 rs_model.fit(X_train, y_train)

```

```

Fitting 3 folds for each of 20 candidates, totalling 60 fits
[CV 1/3] END max_depth=10, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=130; score=0.246 total time=
[CV 2/3] END max_depth=10, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=130; score=0.247 total time=
[CV 3/3] END max_depth=10, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=130; score=0.235 total time=
[CV 1/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=7, n_estimators=60; score=0.257 total time=
[CV 2/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=7, n_estimators=60; score=0.261 total time=
[CV 3/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=7, n_estimators=60; score=0.251 total time=
[CV 1/3] END max_depth=10, max_features=0.5, min_samples_leaf=9, min_samples_split=5, n_estimators=140; score=0.262 total time=
[CV 2/3] END max_depth=10, max_features=0.5, min_samples_leaf=9, min_samples_split=5, n_estimators=140; score=0.266 total time=
[CV 3/3] END max_depth=10, max_features=0.5, min_samples_leaf=9, min_samples_split=5, n_estimators=140; score=0.257 total time=
[CV 1/3] END max_depth=10, max_features=0.5, min_samples_leaf=8, min_samples_split=3, n_estimators=50; score=0.259 total time=
[CV 2/3] END max_depth=10, max_features=0.5, min_samples_leaf=8, min_samples_split=3, n_estimators=50; score=0.258 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=7, min_samples_split=7, n_estimators=140; score=0.256 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=7, min_samples_split=7, n_estimators=140; score=0.257 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=7, min_samples_split=7, n_estimators=140; score=0.249 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=3, n_estimators=180; score=0.256 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=3, n_estimators=180; score=0.258 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=3, n_estimators=180; score=0.248 total time=
[CV 1/3] END max_depth=10, max_features=1.0, min_samples_leaf=2, min_samples_split=7, n_estimators=160; score=0.259 total time=
[CV 2/3] END max_depth=10, max_features=1.0, min_samples_leaf=2, min_samples_split=7, n_estimators=160; score=0.261 total time=
[CV 3/3] END max_depth=10, max_features=1.0, min_samples_leaf=2, min_samples_split=7, n_estimators=160; score=0.253 total time=
[CV 1/3] END max_depth=10, max_features=1.0, min_samples_leaf=6, min_samples_split=3, n_estimators=150; score=0.259 total time=
[CV 2/3] END max_depth=10, max_features=1.0, min_samples_leaf=6, min_samples_split=3, n_estimators=150; score=0.262 total time=
[CV 3/3] END max_depth=10, max_features=1.0, min_samples_leaf=6, min_samples_split=3, n_estimators=150; score=0.256 total time=
[CV 1/3] END max_depth=None, max_features=0.5, min_samples_leaf=4, min_samples_split=2, n_estimators=110; score=0.257 total time=
[CV 2/3] END max_depth=None, max_features=0.5, min_samples_leaf=4, min_samples_split=2, n_estimators=110; score=0.262 total time=
[CV 3/3] END max_depth=None, max_features=0.5, min_samples_leaf=4, min_samples_split=2, n_estimators=110; score=0.255 total time=
[CV 1/3] END max_depth=10, max_features=0.5, min_samples_leaf=3, min_samples_split=2, n_estimators=190; score=0.261 total time=
[CV 2/3] END max_depth=10, max_features=0.5, min_samples_leaf=3, min_samples_split=2, n_estimators=190; score=0.264 total time=
[CV 3/3] END max_depth=10, max_features=0.5, min_samples_leaf=3, min_samples_split=2, n_estimators=190; score=0.256 total time=
[CV 1/3] END max_depth=10, max_features=1.0, min_samples_leaf=1, min_samples_split=7, n_estimators=120; score=0.257 total time=
[CV 2/3] END max_depth=10, max_features=1.0, min_samples_leaf=1, min_samples_split=7, n_estimators=120; score=0.259 total time=
[CV 3/3] END max_depth=10, max_features=1.0, min_samples_leaf=1, min_samples_split=7, n_estimators=120; score=0.253 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=5, min_samples_split=6, n_estimators=90; score=0.254 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=5, min_samples_split=6, n_estimators=90; score=0.258 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=5, min_samples_split=6, n_estimators=90; score=0.247 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=20; score=0.240 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=20; score=0.236 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=3, min_samples_split=8, n_estimators=20; score=0.230 total time=
[CV 1/3] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=160; score=0.253 total time=
[CV 2/3] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=160; score=0.256 total time=
[CV 3/3] END max_depth=20, max_features=sqrt, min_samples_leaf=2, min_samples_split=2, n_estimators=160; score=0.246 total time=
[CV 1/3] END max_depth=20, max_features=0.5, min_samples_leaf=1, min_samples_split=8, n_estimators=10; score=0.204 total time=
[CV 2/3] END max_depth=20, max_features=0.5, min_samples_leaf=1, min_samples_split=8, n_estimators=10; score=0.206 total time=
[CV 3/3] END max_depth=20, max_features=0.5, min_samples_leaf=1, min_samples_split=8, n_estimators=10; score=0.206 total time=
[CV 1/3] END max_depth=None, max_features=1.0, min_samples_leaf=3, min_samples_split=2, n_estimators=140; score=0.251 total time=
[CV 2/3] END max_depth=None, max_features=1.0, min_samples_leaf=3, min_samples_split=2, n_estimators=140; score=0.253 total time=
[CV 3/3] END max_depth=None, max_features=1.0, min_samples_leaf=3, min_samples_split=2, n_estimators=140; score=0.248 total time=
[CV 1/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=2, n_estimators=100; score=0.259 total time=
[CV 2/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=2, n_estimators=100; score=0.264 total time=
[CV 3/3] END max_depth=None, max_features=0.5, min_samples_leaf=6, min_samples_split=2, n_estimators=100; score=0.255 total time=
[CV 1/3] END max_depth=20, max_features=1.0, min_samples_leaf=7, min_samples_split=7, n_estimators=10; score=0.223 total time=
[CV 2/3] END max_depth=20, max_features=1.0, min_samples_leaf=7, min_samples_split=7, n_estimators=10; score=0.223 total time=
[CV 3/3] END max_depth=20, max_features=1.0, min_samples_leaf=7, min_samples_split=7, n_estimators=10; score=0.225 total time=
[CV 1/3] END max_depth=10, max_features=sqrt, min_samples_leaf=1, min_samples_split=7, n_estimators=150; score=0.245 total time=
[CV 2/3] END max_depth=10, max_features=sqrt, min_samples_leaf=1, min_samples_split=7, n_estimators=150; score=0.248 total time=
[CV 3/3] END max_depth=10, max_features=sqrt, min_samples_leaf=1, min_samples_split=7, n_estimators=150; score=0.235 total time=
[CV 1/3] END max_depth=None, max_features=1.0, min_samples_leaf=8, min_samples_split=2, n_estimators=20; score=0.242 total time=
[CV 2/3] END max_depth=None, max_features=1.0, min_samples_leaf=8, min_samples_split=2, n_estimators=20; score=0.246 total time=
[CV 3/3] END max_depth=None, max_features=1.0, min_samples_leaf=8, min_samples_split=2, n_estimators=20; score=0.241 total time=

```

RandomizedSearchCV



best_estimator_:
RandomForestRegressor

RandomForestRegressor

```
1 rs_model.best_params_
```

```

2 { 'n_estimators': np.int64(140),
  'min_samples_split': np.int64(5),
  'min_samples_leaf': np.int64(9),
  'max_features': 0.5,
  'max_depth': 10}

```

```

1 rs_model_scores = show_scores(model=rs_model,
2                                train_features=X_train,
3                                train_labels=y_train,
4                                test_features=X_test,
5                                test_labels=y_test)
6 rs_model_scores

```

```

7 {'Training MAE': 0.5473572853514472,
8  'Test MAE': 0.5809276716238158,
9 }

```

```
'Training RMSLE': 0.08870563578612232,
'Test RMSLE': 0.09300200957031936,
'Training R^2': 0.35659503446317753,
'Test R^2': 0.25901433056358125}
```

This code block takes the best hyperparameters found during a previous randomized search on a sample of the data (`sample_rs_model.best_params_`) and uses them to train a new RandomForestRegressor model on the full training dataset (`X_train, y_train`).

```
1 np.random.seed(42)
2
3 rf_best_model_rs = RandomForestRegressor(**sample_rs_model.best_params_, n_jobs=-1)
4
5 rf_best_model_rs.fit(X=X_train,
6                      y=y_train)
7
8 rf_best_model_rs_scores = show_scores(model=rf_best_model_rs,
9                                       train_features=X_train,
10                                      train_labels=y_train,
11                                      test_features=X_test,
12                                      test_labels=y_test)
13
14 rf_best_model_rs_scores
```

→ {'Training MAE': 0.5624191487918332,
 'Test MAE': 0.5876372043698553,
 'Training RMSLE': 0.09061482159477974,
 'Test RMSLE': 0.09415714889156918,
 'Training R^2': 0.32594390438005383,
 'Test R^2': 0.24094986223080406}

This code block trains a machine learning model, specifically a Random Forest Regressor, using the optimal hyperparameters found during a previous Grid Search performed on a sample of the data. It then evaluates the performance of this trained model.

```
1 np.random.seed(42)
2
3 rf_best_model_gs = RandomForestRegressor(**sample_gs_model.best_params_, n_jobs=-1)
4
5 rf_best_model_gs.fit(X=X_train,
6                      y=y_train)
7
8 rf_best_model_gs_scores = show_scores(model=rf_best_model_gs,
9                                       train_features=X_train,
10                                      train_labels=y_train,
11                                      test_features=X_test,
12                                      test_labels=y_test)
13
14 rf_best_model_gs_scores
```

→ {'Training MAE': 0.5435080531798863,
 'Test MAE': 0.5814392662242237,
 'Training RMSLE': 0.08762451259715064,
 'Test RMSLE': 0.09304934505106814,
 'Training R^2': 0.3703261958364007,
 'Test R^2': 0.2581559797832569}

This code block trains a machine learning model, specifically a Random Forest Regressor, using the optimal hyperparameters found during a previous Halving Randomized Search performed on a sample of the data. It then evaluates the performance of this trained model.

```
1 np.random.seed(42)
2
3 rf_best_model_hrs = RandomForestRegressor(**sample_hrs_model.best_params_, n_jobs=-1)
4
5 rf_best_model_hrs.fit(X=X_train,
6                      y=y_train)
7
8 rf_best_model_hrs_scores = show_scores(model=rf_best_model_hrs,
9                                       train_features=X_train,
10                                      train_labels=y_train,
11                                      test_features=X_test,
12                                      test_labels=y_test)
13
14 rf_best_model_hrs_scores
```

→ {'Training MAE': 0.565585096031433,
 'Test MAE': 0.5875540100160795,
 'Training RMSLE': 0.09170307199110485,
 'Test RMSLE': 0.09412788456318996,
 'Training R^2': 0.3115352189683698,
 'Test R^2': 0.24132504820335732}

This code block is setting up and running a HalvingRandomSearchCV for a RandomForestRegressor model. The goal is to efficiently find the best combination of hyperparameters for the model by exploring a range of possible values.

```
1 np.random.seed(42)
2
3 rf_grid = {"n_estimators": np.arange(50, 300, 25),
4             "max_depth": [10, 20, 30, None],
5             "min_samples_split": np.arange(5, 20, 2),
6             "min_samples_leaf": np.arange(3, 15, 2),
7             "max_features": [0.5, 0.7, 1.0, "sqrt"],
8             "bootstrap": [True, False],
9             "oob_score": [True, False]}
10
11 halving_rs_model = HalvingRandomSearchCV(estimator=RandomForestRegressor(n_jobs=-1),
12                                             param_distributions=rf_grid,
13                                             cv=5,
14                                             verbose=3,
15                                             n_jobs=-1)
16
17 halving_rs_model.fit(X_train, y_train)
```

```
    n_iterations: 8
    n_required_iterations: 8
    n_possible_iterations: 8
    min_resources_: 10
    max_resources_: 48739
    aggressive_elimination: False
    factor: 3
    -----
    iter: 0
    n_candidates: 4873
    n_resources: 10
    Fitting 5 folds for each of 4873 candidates, totalling 24365 fits
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_validation.py:528: FitFailedWarning:
    6055 fits failed out of a total of 24365.
    The score on these train-test partitions for these parameters will be set to nan.
    If these failures are not expected, you can try to debug them by setting error_score='raise'.

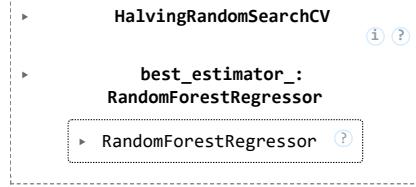
    Below are more details about the failures:
    -----
    6055 fits failed with the following error:
    Traceback (most recent call last):
      File "/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_validation.py", line 866, in _fit_and_score
        estimator.fit(X_train, y_train, **fit_params)
      File "/usr/local/lib/python3.11/dist-packages/sklearn/base.py", line 1389, in wrapper
        return fit_method(estimator, *args, **kwargs)
        ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
      File "/usr/local/lib/python3.11/dist-packages/sklearn/ensemble/_forest.py", line 448, in fit
        raise ValueError("Out of bag estimation only available if bootstrap=True")
    ValueError: Out of bag estimation only available if bootstrap=True

        warnings.warn(some_fits_failed_message, FitFailedWarning)
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
        warnings.warn(
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are
    -4.27193078e-04 -4.44089210e-17]
        warnings.warn(
    -----
    iter: 1
    n_candidates: 1625
    n_resources: 30
    Fitting 5 folds for each of 1625 candidates, totalling 8125 fits
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_validation.py:528: FitFailedWarning:
    2005 fits failed out of a total of 8125.
    The score on these train-test partitions for these parameters will be set to nan.
    If these failures are not expected, you can try to debug them by setting error_score='raise'.

    Below are more details about the failures:
    -----
    2005 fits failed with the following error:
    Traceback (most recent call last):
      File "/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_validation.py", line 866, in _fit_and_score
        estimator.fit(X_train, y_train, **fit_params)
      File "/usr/local/lib/python3.11/dist-packages/sklearn/base.py", line 1389, in wrapper
        return fit_method(estimator, *args, **kwargs)
        ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
      File "/usr/local/lib/python3.11/dist-packages/sklearn/ensemble/_forest.py", line 448, in fit
        raise ValueError("Out of bag estimation only available if bootstrap=True")
    ValueError: Out of bag estimation only available if bootstrap=True

        warnings.warn(some_fits_failed_message, FitFailedWarning)
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
    -0.01605811]
        warnings.warn(
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are
    3.41004882e-01 2.90062201e-01]
        warnings.warn(
    -----
    iter: 2
    n_candidates: 542
    n_resources: 90
    Fitting 5 folds for each of 542 candidates, totalling 2710 fits
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
    -0.00157941]
        warnings.warn(
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are
    7.06175891e-01 7.62732700e-01]
        warnings.warn(
    -----
    iter: 3
    n_candidates: 181
    n_resources: 270
    Fitting 5 folds for each of 181 candidates, totalling 905 fits
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
        warnings.warn(
    /usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are
    7.25171879e-01 7.30219253e-01]
        warnings.warn(
    -----
    iter: 4
    n_candidates: 61
```

```
n_resources: 810
Fitting 5 folds for each of 61 candidates, totalling 305 fits
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are r
  5.94897933e-01  5.76805513e-01]
  warnings.warn(
-----
iter: 5
n_candidates: 21
n_resources: 2430
Fitting 5 folds for each of 21 candidates, totalling 105 fits
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are r
  5.65216244e-01  5.52050012e-01]
  warnings.warn(
-----
iter: 6
n_candidates: 7
n_resources: 7290
Fitting 5 folds for each of 7 candidates, totalling 35 fits
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are r
  4.93965097e-01  6.15766271e-01]
  warnings.warn(
-----
iter: 7
n_candidates: 3
n_resources: 21870
Fitting 5 folds for each of 3 candidates, totalling 15 fits
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the test scores are r
  warnings.warn(
/usr/local/lib/python3.11/dist-packages/sklearn/model_selection/_search.py:1108: UserWarning: One or more of the train scores are r
  5.76771625e-01  4.24070309e-01]
  warnings.warn(
```



```

1 halving_rs_model.best_params_
→ {'oob_score': False,
  'n_estimators': np.int64(275),
  'min_samples_split': np.int64(17),
  'min_samples_leaf': np.int64(7),
  'max_features': 0.5,
  'max_depth': 30,
  'bootstrap': True}

1 halving_rs_model_scores = show_scores(model=halving_rs_model,
2                                         train_features=X_train,
3                                         train_labels=y_train,
4                                         test_features=X_test,
5                                         test_labels=y_test)
6 halving_rs_model_scores

→ {'Training MAE': 0.43239679105730994,
  'Test MAE': 0.5809039460064926,
  'Training RMSLE': 0.07387080889863588,
  'Test RMSLE': 0.09281565264812022,
  'Training R^2': 0.5763891058035406,
  'Test R^2': 0.26164274423707845}

```

This code block trains a Ridge Regression model. It then evaluates the trained model's performance using show_scores function and prints the resulting scores.

```

1 ridge_model = Ridge()
2
3 ridge_model.fit(X_train, y_train)
4
5 ridge_scores = show_scores(model=ridge_model,
6                             train_features=X_train,
7                             train_labels=y_train,
8                             test_features=X_test,
9                             test_labels=y_test)
10
11 print("Ridge Regression Scores:")
12 ridge_scores

→ Ridge Regression Scores:
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_ridge.py:215: LinAlgWarning: Ill-conditioned matrix (rcond=2.93359e-27)
    return linalg.solve(A, Xy, assume_a="pos", overwrite_a=True).T
{'Training MAE': 0.6231777443912474,
 'Test MAE': 0.6186031282587662,
 'Training RMSLE': 0.10036177894638101,
 'Test RMSLE': 0.09854144256145399,
 'Training R^2': 0.17322049052005628,
 'Test R^2': 0.16280710973535883}

```

This code block trains and evaluates a Support Vector Regression (SVR) model with a Radial Basis Function (RBF) kernel.

```

1 svr_model = SVR(kernel="rbf")
2
3 svr_model.fit(X_train, y_train)
4
5 svr_scores = show_scores(model=svr_model,
6                          train_features=X_train,
7                          train_labels=y_train,
8                          test_features=X_test,
9                          test_labels=y_test)
10
11 print("SVR (RBF Kernel) Scores:")
12 svr_scores

→ SVR (RBF Kernel) Scores:
{'Training MAE': 0.6881027560787162,
 'Test MAE': 0.6842749615898589,
 'Training RMSLE': 0.10811746913004266,
 'Test RMSLE': 0.1062625410687166,
 'Training R^2': 0.008715864813863994,
 'Test R^2': 0.007762053219004805}

```

This code block trains and evaluates a Gradient Boosting Regressor.

```

1 gbr_model = GradientBoostingRegressor(random_state=42)
2

```

```

3 gbr_model.fit(X_train, y_train)
4
5 gbr_scores = show_scores(model=gbr_model,
6                           train_features=X_train,
7                           train_labels=y_train,
8                           test_features=X_test,
9                           test_labels=y_test)
10
11 print("Gradient Boosting Regressor Scores:")
12 gbr_scores
→ Gradient Boosting Regressor Scores:
{'Training MAE': 0.5751491075187253,
 'Test MAE': 0.5798634603002869,
 'Training RMSLE': 0.09284064173529244,
 'Test RMSLE': 0.09280764947156771,
 'Training R^2': 0.2898418118980568,
 'Test R^2': 0.262580555628648}

```

This code block focuses on training and evaluating a Lasso Regression model.

```

1 lasso_model = Lasso(random_state=42)
2
3 lasso_model.fit(X_train, y_train)
4
5 lasso_scores = show_scores(model=lasso_model,
6                            train_features=X_train,
7                            train_labels=y_train,
8                            test_features=X_test,
9                            test_labels=y_test)
10
11 print("Lasso Regression Scores:")
12 lasso_scores
→ Lasso Regression Scores:
{'Training MAE': 0.632726509690792,
 'Test MAE': 0.6292975515719712,
 'Training RMSLE': 0.10120086801639935,
 'Test RMSLE': 0.09936992092910946,
 'Training R^2': 0.15270338476070877,
 'Test R^2': 0.14448708096918006}

```

This code block focuses on training and evaluating a Linear Regression model.

```

1 linear_model = LinearRegression()
2
3 linear_model.fit(X_train, y_train)
4
5 linear_scores = show_scores(model=linear_model,
6                            train_features=X_train,
7                            train_labels=y_train,
8                            test_features=X_test,
9                            test_labels=y_test)
10
11 print("Linear Regression Scores:")
12 linear_scores
→ Linear Regression Scores:
{'Training MAE': 0.6245620575467123,
 'Test MAE': 0.6199371872971146,
 'Training RMSLE': 0.10063056015042172,
 'Test RMSLE': 0.09867805024108317,
 'Training R^2': 0.16845493968968595,
 'Test R^2': 0.16032442404375724}

```

This code block trains and evaluates a XGBoost (Extreme Gradient Boosting) Regressor model.

```

1 xgb_model = xgb.XGBRegressor(objective='reg:squarederror', random_state=42)
2
3 xgb_model.fit(X_train, y_train)
4
5 xgb_scores = show_scores(model=xgb_model,
6                           train_features=X_train,
7                           train_labels=y_train,
8                           test_features=X_test,
9                           test_labels=y_test)
10
11 print("XGBoost Regressor Scores:")
12 xgb_scores

```

```

☒ XGBoost Regressor Scores:
{'Training MAE': 0.4955529241429421,
 'Test MAE': 0.5951668821028818,
 'Training RMSLE': 0.0790799696419051,
 'Test RMSLE': 0.09462947333056976,
 'Training R^2': 0.47928815972210703,
 'Test R^2': 0.22986016432730916}

```

This code block trains and evaluates a LightGBM Regressor model.

```

1 base_lgb_model = lgb.LGBMRegressor(objective='regression', random_state=42)
2
3 base_lgb_model .fit(X_train, y_train)
4
5 base_lgb_scores = show_scores(model=base_lgb_model ,
6                               train_features=X_train,
7                               train_labels=y_train,
8                               test_features=X_test,
9                               test_labels=y_test)
10
11 print("LightGBM Regressor Scores:")
12 base_lgb_scores

☒ [LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.019830 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true` .
[LightGBM] [Info] Total Bins 7622
[LightGBM] [Info] Number of data points in the train set: 48739, number of used features: 64
[LightGBM] [Info] Start training from score 7.906696
LightGBM Regressor Scores:
{'Training MAE': 0.5525830835730449,
 'Test MAE': 0.5788805861720081,
 'Training RMSLE': 0.08900097029380527,
 'Test RMSLE': 0.09255008213444889,
 'Training R^2': 0.34836009160766324,
 'Test R^2': 0.2661600517020939}

```

This code block trains and evaluates a Decision Tree Regressor model.

```

1 dt_model = DecisionTreeRegressor(random_state=42)
2
3 dt_model.fit(X_train, y_train)
4
5 dt_scores = show_scores(model=dt_model,
6                         train_features=X_train,
7                         train_labels=y_train,
8                         test_features=X_test,
9                         test_labels=y_test)
10
11 print("Decision Tree Regression Scores:")
12 dt_scores

☒ Decision Tree Regression Scores:
{'Training MAE': 2.186778767804326e-19,
 'Test MAE': 0.8620016413623306,
 'Training RMSLE': 0.0,
 'Test RMSLE': 0.1367951210369189,
 'Training R^2': 1.0,
 'Test R^2': -0.5571233874357828}

```

Because LightGBM Regressor produced the best score, we will try to tune it to get an even better score.

This code block focuses on tuning the hyperparameters of a LightGBM Regressor model using Randomized Search Cross-Validation.

```

1 np.random.seed(42)
2
3 lgb_grid = {
4     'n_estimators': [100, 200, 300],
5     'learning_rate': [0.01, 0.1, 0.2],
6     'num_leaves': [31, 50, 100],
7     'max_depth': [-1, 10, 20],
8     'min_child_samples': [20, 50, 100],
9     'subsample': [0.8, 1.0],
10    'colsample_bytree': [0.8, 1.0],
11 }
12
13 tuned_lgb_model = RandomizedSearchCV(estimator=lgb.LGBMRegressor(random_state=42),
14                                         param_distributions=lgb_grid,
15                                         n_iter=20,

```

```

16                                         cv=3,
17                                         verbose=3,
18                                         n_jobs=-1)
19
20 tuned_lgb_model.fit(X=X_train, y=y_train)

→ Fitting 3 folds for each of 20 candidates, totalling 60 fits
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.056658 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 7622
[LightGBM] [Info] Number of data points in the train set: 48739, number of used features: 64
[LightGBM] [Info] Start training from score 7.906696
    ▶ RandomizedSearchCV
        ⓘ ⓘ
    ▶ best_estimator_:
        LGBMRegressor
            ▶ LGBMRegressor

```

1 tuned_lgb_model.best_params_

```

→ {'subsample': 1.0,
    'num_leaves': 31,
    'n_estimators': 100,
    'min_child_samples': 100,
    'max_depth': 20,
    'learning_rate': 0.1,
    'colsample_bytree': 0.8}

1 tuned_lgb_scores = show_scores(model=tuned_lgb_model,
2                                     train_features=X_train,
3                                     train_labels=y_train,
4                                     test_features=X_test,
5                                     test_labels=y_test)
6 tuned_lgb_scores

→ {'Training MAE': 0.555199268874275,
    'Test MAE': 0.5779629516602062,
    'Training RMSLE': 0.08976384401293586,
    'Test RMSLE': 0.09243194321941373,
    'Training R^2': 0.33928446621819053,
    'Test R^2': 0.26818849514884713}

```

This code block is used to organize and compare the performance scores of various machine learning models that were trained and evaluated in the preceding cells.

```

1 base_model_scores["model_name"] = "RandomForestRegressor (Base)"
2 rs_model_scores["model_name"] = "RandomForestRegressor (Tuned)"
3 ridge_scores["model_name"] = "Ridge"
4 svr_scores["model_name"] = "SVR"
5 gbr_scores["model_name"] = "GradientBoostingRegressor"
6 lasso_scores["model_name"] = "Lasso"
7 linear_scores["model_name"] = "LinearRegression"
8 xgb_scores["model_name"] = "XGBoost"
9 base_lgb_scores["model_name"] = "LightGBM (Base)"
10 tuned_lgb_scores["model_name"] = "LightGBM (Tuned)"
11 dt_scores["model_name"] = "DecisionTreeRegressor"
12 rf_best_model_rs_scores["model_name"] = "RandomForestRegressor (Random Search Tuned)"
13 rf_best_model_gs_scores["model_name"] = "RandomForestRegressor (Grid Search Tuned)"
14 halving_rs_model_scores["model_name"] = "RandomForestRegressor (Halving Random Search Tuned)"
15
16
17 all_model_scores = [base_model_scores,
18                     rs_model_scores,
19                     rf_best_model_rs_scores,
20                     rf_best_model_gs_scores,
21                     halving_rs_model_scores,
22                     ridge_scores,
23                     svr_scores,
24                     gbr_scores,
25                     lasso_scores,
26                     linear_scores,
27                     xgb_scores,
28                     base_lgb_scores,
29                     tuned_lgb_scores,
30                     dt_scores]
31
32 model_comparison_df = pd.DataFrame(all_model_scores).sort_values(by="Test R^2", ascending=True)
33 model_comparison_df

```

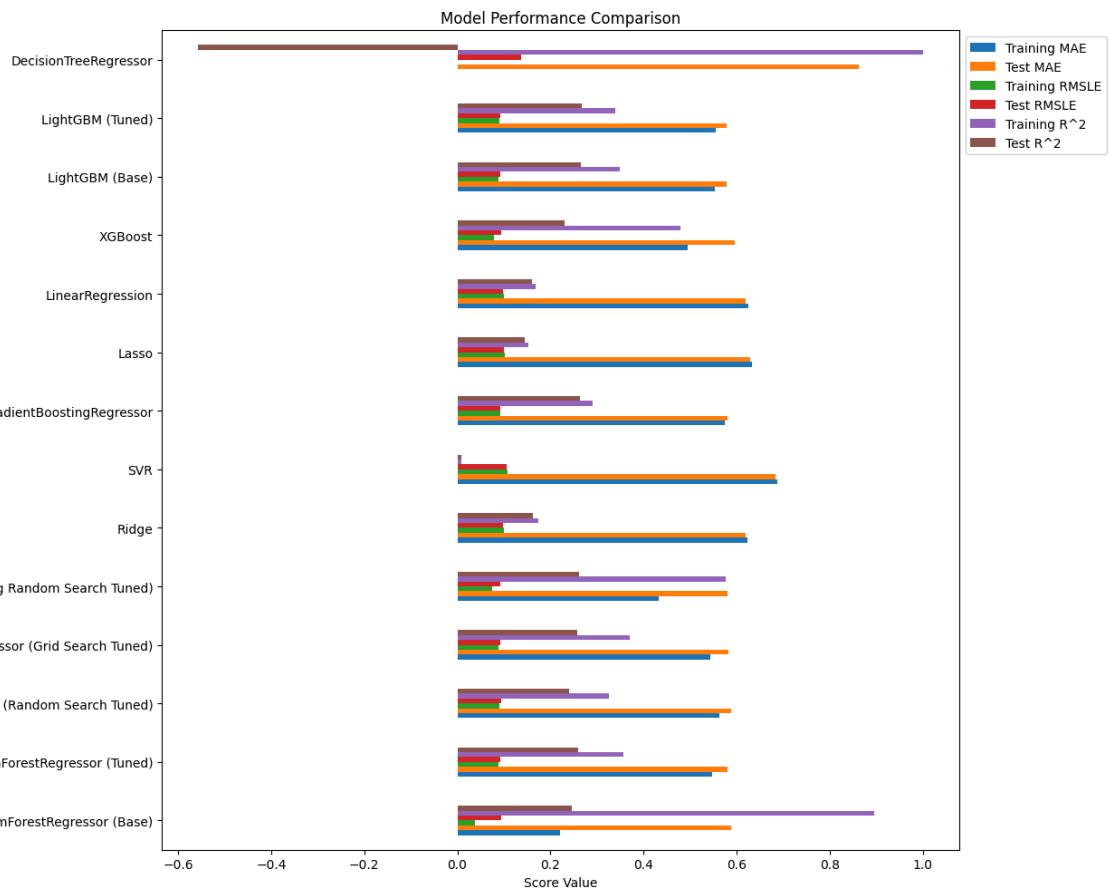
	Training MAE	Test MAE	Training RMSLE	Test RMSLE	Training R^2	Test R^2	model_name
13	2.186779e-19	0.862002	0.000000	0.136795	1.000000	-0.557123	DecisionTreeRegressor
6	6.881028e-01	0.684275	0.108117	0.106263	0.008716	0.007762	SVR
8	6.327265e-01	0.629298	0.101201	0.099370	0.152703	0.144487	Lasso
9	6.245621e-01	0.619937	0.100631	0.098678	0.168455	0.160324	LinearRegression
5	6.231777e-01	0.618603	0.100362	0.098541	0.173220	0.162807	Ridge
10	4.955529e-01	0.595167	0.079080	0.094629	0.479288	0.229860	XGBoost
2	5.624191e-01	0.587637	0.090615	0.094157	0.325944	0.240950	RandomForestRegressor (Random Search Tuned)
0	2.205965e-01	0.588709	0.038303	0.093677	0.894671	0.245387	RandomForestRegressor (Base)
3	5.435081e-01	0.581439	0.087625	0.093049	0.370326	0.258156	RandomForestRegressor (Grid Search Tuned)
1	5.473573e-01	0.580928	0.088706	0.093002	0.356595	0.259014	RandomForestRegressor (Tuned)
4	4.323968e-01	0.580904	0.073871	0.092816	0.576389	0.261643	RandomForestRegressor (Halving Random Search T...)
7	5.751491e-01	0.579863	0.092841	0.092808	0.289842	0.262581	GradientBoostingRegressor
11	5.525831e-01	0.578881	0.089001	0.092550	0.348360	0.266160	LightGBM (Base)
12	5.551993e-01	0.577963	0.089764	0.092432	0.339284	0.268188	LightGBM (Tuned)

This code block focuses on visualizing the performance of the various machine learning models trained in the previous steps. It creates a bar chart to easily compare the scores of each model.

```

1 model_scores = pd.DataFrame({
2     "RandomForestRegressor (Base)": base_model_scores,
3     "RandomForestRegressor (Tuned)": rs_model_scores,
4     "RandomForestRegressor (Random Search Tuned)": rf_best_model_rs_scores,
5     "RandomForestRegressor (Grid Search Tuned)": rf_best_model_gs_scores,
6     "RandomForestRegressor (Halving Random Search Tuned)": halving_rs_model_scores,
7     "Ridge": ridge_scores,
8     "SVR": svr_scores,
9     "GradientBoostingRegressor": gbr_scores,
10    "Lasso": lasso_scores,
11    "LinearRegression": linear_scores,
12    "XGBoost": xgb_scores,
13    "LightGBM (Base)": base_lgb_scores,
14    "LightGBM (Tuned)": tuned_lgb_scores,
15    "DecisionTreeRegressor": dt_scores
16 })
17
18 model_scores = model_scores.transpose()
19
20 fig, ax = plt.subplots(figsize=(15, 10))
21 model_scores.plot(kind="barh", ax=ax)
22 ax.set_title("Model Performance Comparison")
23 ax.set_xlabel("Score Value")
24 ax.set_ylabel("Model")
25 plt.legend(loc="upper left", bbox_to_anchor=(1, 1))
26 plt.tight_layout()
27 plt.show()

```



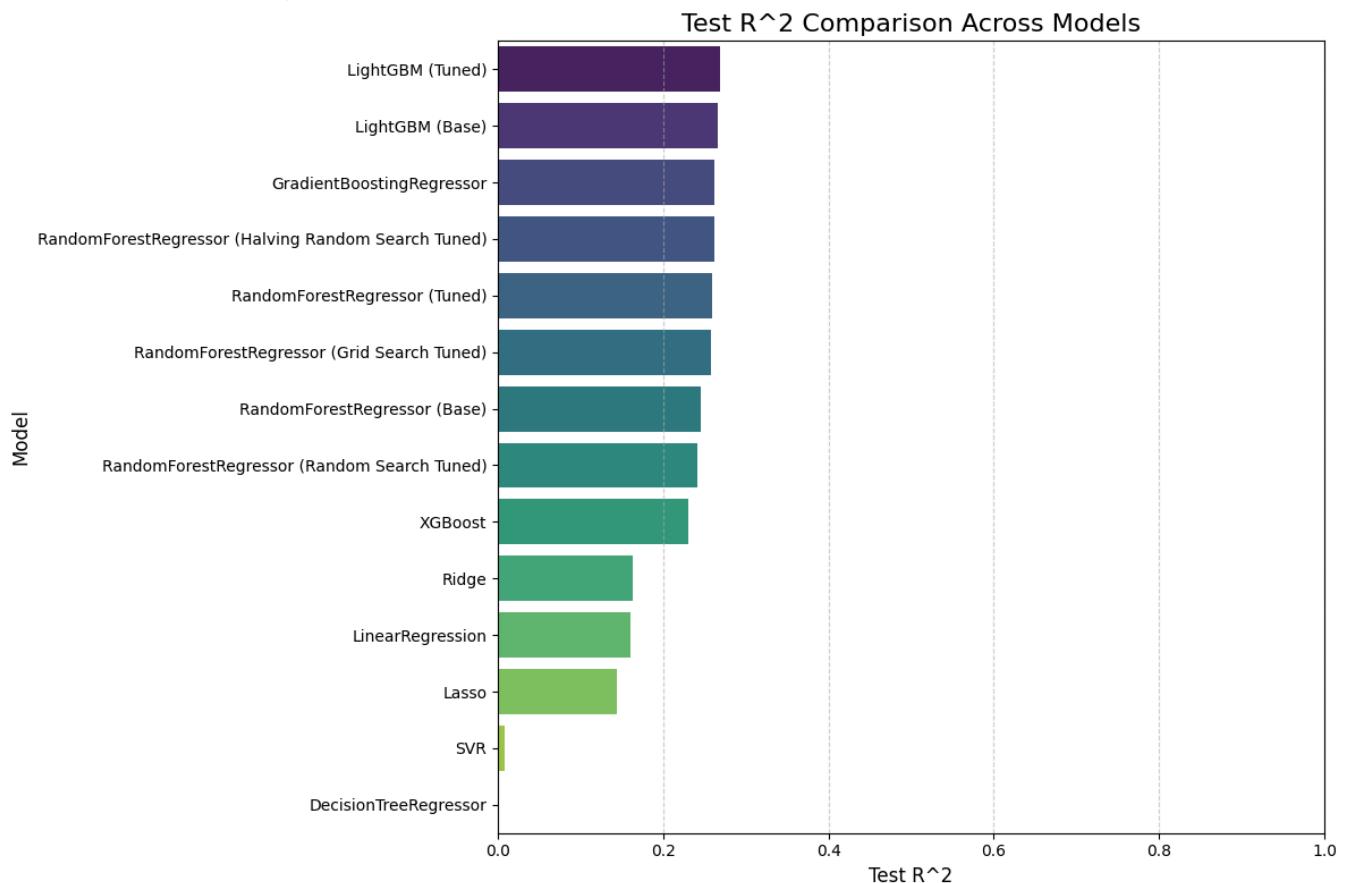
This code block focuses on visualizing the performance of the machine learning models trained and evaluated previously. It specifically creates a horizontal bar chart to compare the Test R² score for each model.

```
1 model_comparison_df_sorted = model_comparison_df.sort_values(by="Test R^2", ascending=False)
2
3 plt.figure(figsize=(12, 8))
4 sns.barplot(x="Test R^2", y="model_name", data=model_comparison_df_sorted, palette="viridis")
5 plt.title("Test R^2 Comparison Across Models", fontsize=16)
6 plt.xlabel("Test R^2", fontsize=12)
7 plt.ylabel("Model", fontsize=12)
8 plt.xlim(0, 1)
9 plt.grid(axis='x', linestyle='--', alpha=0.6)
10 plt.tight_layout()
11 plt.show()
```

```

[4]: <ipython-input-107-43a2ba204d6a>:4: FutureWarning:
      Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set `le
      sns.barplot(x="Test R^2", y="model_name", data=model_comparison_df_sorted, palette="viridis")

```



▼ Machine Learning Model (Total Completed)

This code block prepares the sample data for training a machine learning model to predict the 'user_completed'.

```

1 sample = user_country_stats.sample(n=1000, random_state=42)
2
3 x_sample = sample.drop(labels="user_completed", axis=1)
4 y_sample = np.log1p(sample["user_completed"])

```

This code block splits the data into training and testing sets.

```
1 x_sample_train, x_sample_test, y_sample_train, y_sample_test = train_test_split(x_sample, y_sample, test_size=0.2, random_state=42)
```

This code block splits the main dataset into features and the target variable, and then further divides these into training and testing sets.

```

1 X = user_country_stats.drop('user_completed', axis=1)
2 y = np.log1p(user_country_stats['user_completed'])
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

This code block focuses on cleaning the column names of training and testing feature sets, X_train and X_test.

```

1 X_train.columns = [re.sub('[^A-Za-z0-9_]+', '_', col) for col in X_train.columns]
2 X_test.columns = [re.sub('[^A-Za-z0-9_]+', '_', col) for col in X_test.columns]

```