# Import the CSV file

```
1 from google.colab import files
2 import pandas as pd
3
4 uploaded = files.upload()
```

Choose Files  no files selected    Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving all-weeks-countries.csv to all-weeks-countries.csv

```
1 file_name = list(uploaded.keys())[0]
2 top10 = pd.read_csv(file_name)
3
4 top10 = pd.DataFrame(top10)
5
6 top10.info()
7 top10.describe()
```
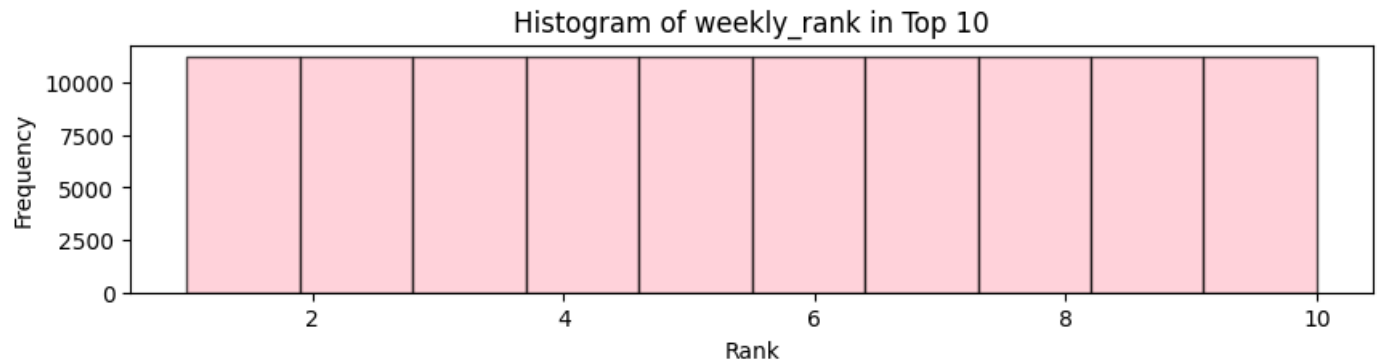
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 112300 entries, 0 to 112299
Data columns (total 8 columns):
 #   Column                      Non-Null Count   Dtype
---  ------                      --------------   -----
 0   country_name                112300 non-null  object
 1   country_iso2                112300 non-null  object
 2   week                        112300 non-null  object
 3   category                    112300 non-null  object
 4   weekly_rank                 112300 non-null  int64
 5   show_title                  112300 non-null  object
 6   season_title                54668 non-null   object
 7   cumulative_weeks_in_top_10  112300 non-null  int64
dtypes: int64(2), object(6)
memory usage: 6.9+ MB
```

|       | weekly_rank   | cumulative_weeks_in_top_10 |
|-------|---------------|----------------------------|
| count | 112300.000000 | 112300.000000              |
| mean  | 5.500000      | 3.468281                   |
| std   | 2.872294      | 5.518189                   |
| min   | 1.000000      | 1.000000                   |
| 25%   | 3.000000      | 1.000000                   |
| 50%   | 5.500000      | 2.000000                   |
| 75%   | 8.000000      | 3.000000                   |
| max   | 10.000000     | 60.000000                  |

## Preliminary analysis

```
1 import matplotlib.pyplot as plt
2
3 plt.figure(figsize=(10, 2))
4 plt.hist(top10['weekly_rank'], bins=10, edgecolor='black', alpha=0.7, color='p
5 plt.title('Histogram of weekly_rank in Top 10')
6 plt.xlabel('Rank')
7 plt.ylabel('Frequency')
8
```
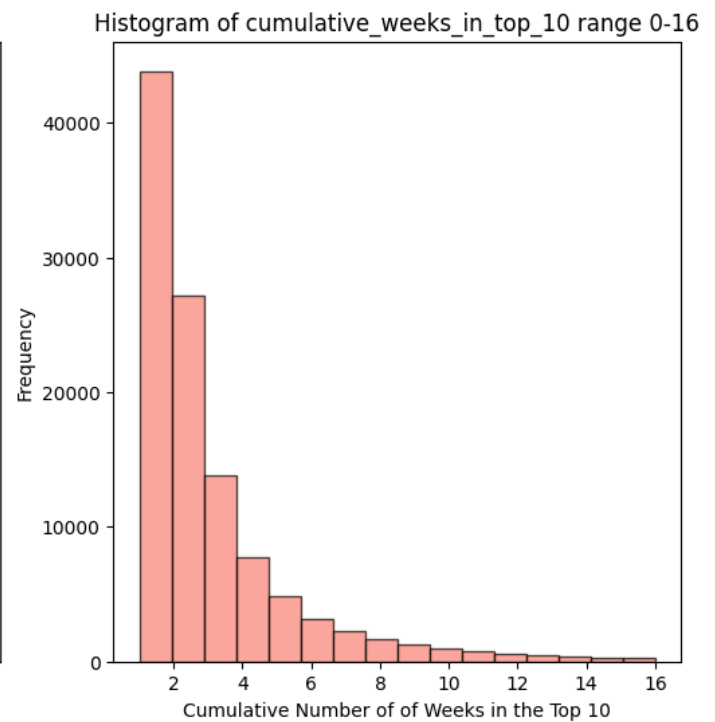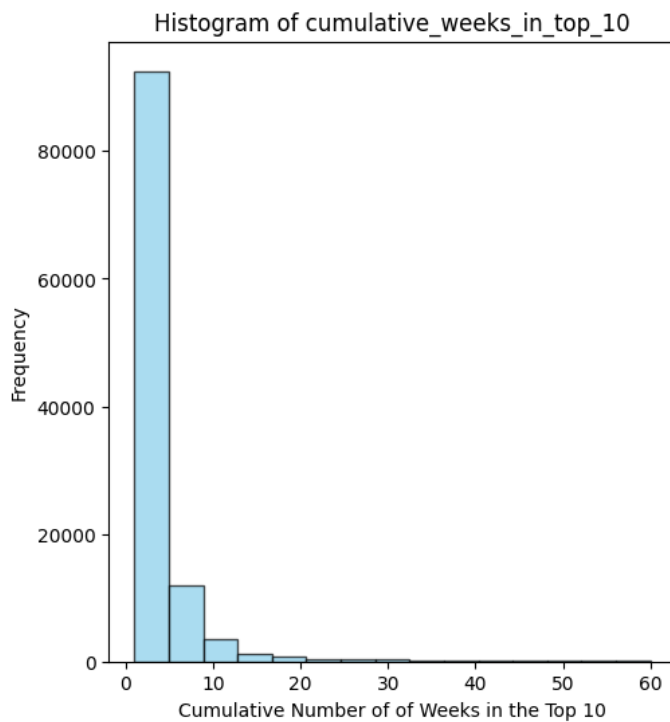
Text(0, 0.5, 'Frequency')

```
1 fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(12, 6))
2
3 axes[0].hist(top10['cumulative_weeks_in_top_10'], bins=15, edgecolor='black',
4 axes[0].set_title('Histogram of cumulative_weeks_in_top_10')
5 axes[0].set_xlabel('Cumulative Number of of Weeks in the Top 10')
6 axes[0].set_ylabel('Frequency')
7
8 axes[1].hist(top10['cumulative_weeks_in_top_10'], bins=16, range=(1, 16), edge
9 axes[1].set_title('Histogram of cumulative_weeks_in_top_10 range 0-16')
10 axes[1].set_xlabel('Cumulative Number of of Weeks in the Top 10')
11 axes[1].set_ylabel('Frequency')
12
```
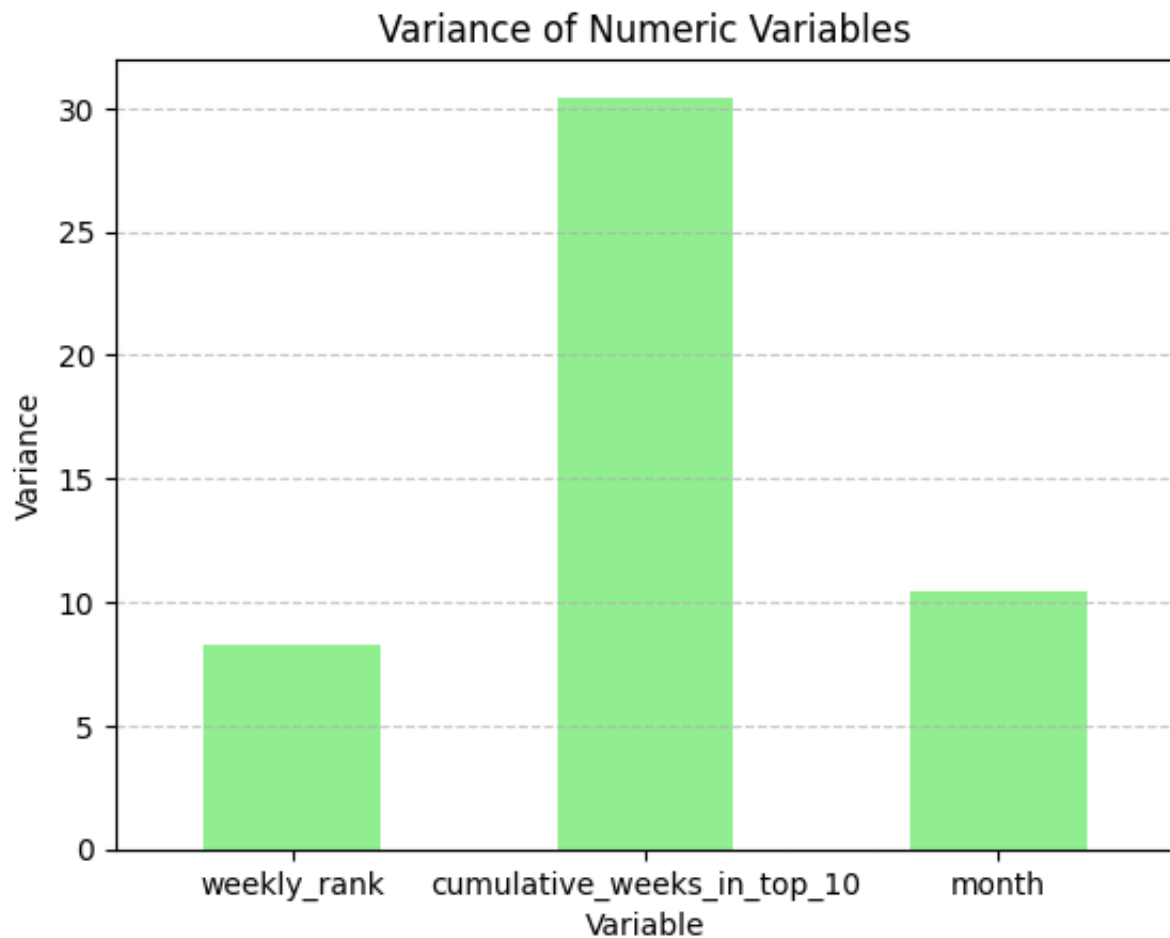
Text(0, 0.5, 'Frequency')

```
1 variances = top10.var()
2
3 variances.plot(kind='bar', color='lightgreen')
4 plt.title('Variance of Numeric Variables')
5 plt.xlabel('Variable')
6 plt.ylabel('Variance')
7 plt.xticks(rotation=0)
8 plt.grid(axis='y', linestyle='--', alpha=0.7)
9 plt.show()
```
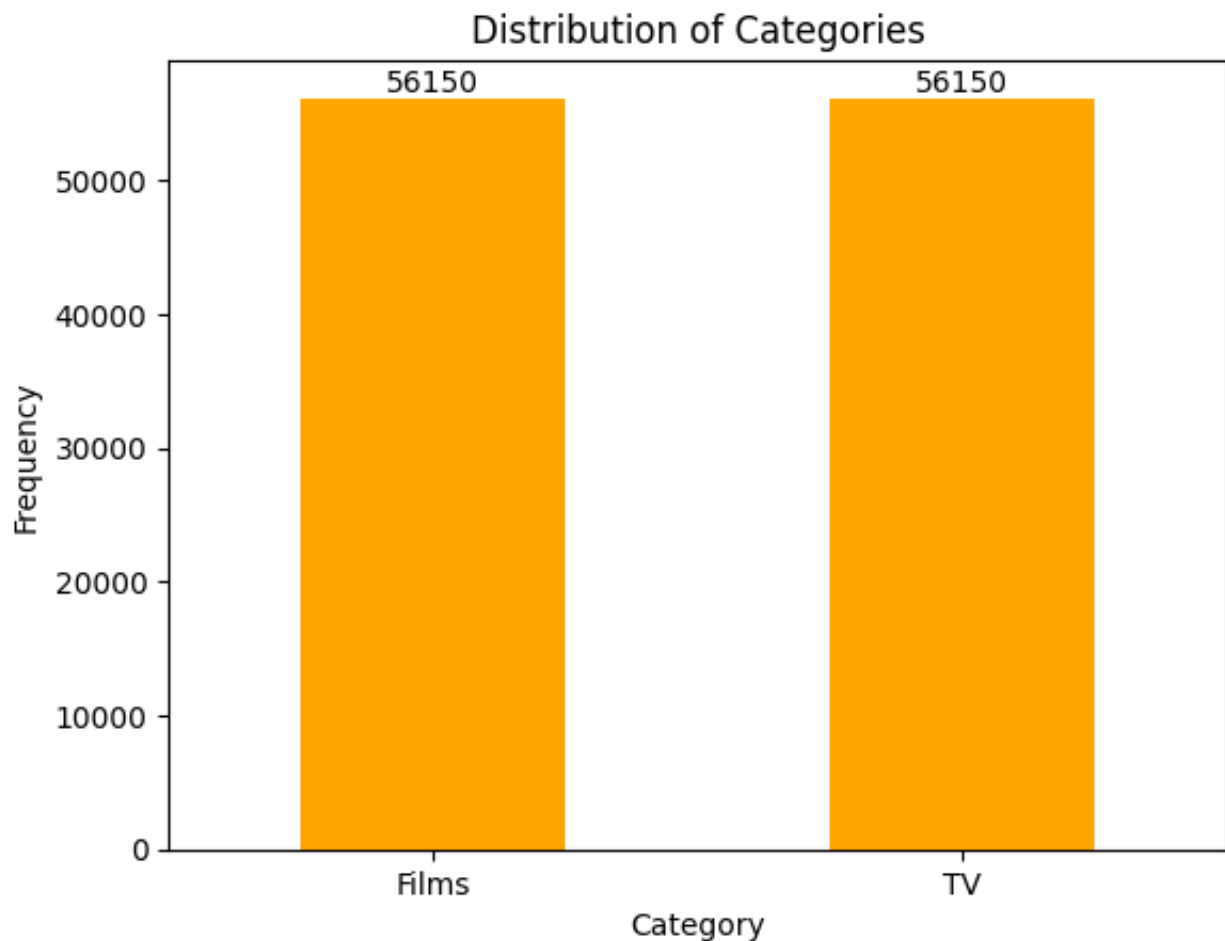
```
<ipython-input-6-4c8b6dacff89>:1: FutureWarning: The default value of numeric_
    variances = top10.var()
```



Variance of Numeric Variables

```
1 category_counts = top10['category'].value_counts()
2 category_counts.plot(kind='bar', color='orange')
3
4 for i, count in enumerate(category_counts):
5     plt.text(i, count, str(count), ha='center', va='bottom')
6
7 plt.title('Distribution of Categories')
8 plt.xlabel('Category')
9 plt.ylabel('Frequency')
10 plt.xticks(rotation=0)
11 plt.show()
```
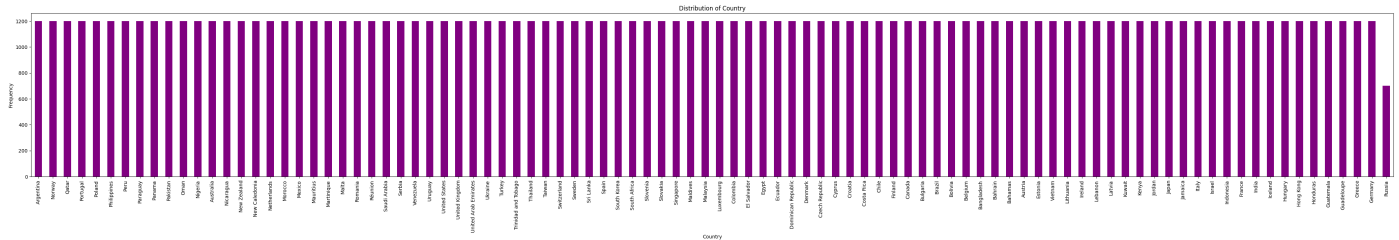
```
1 plt.figure(figsize=(50, 6))
2 top10['country_name'].value_counts().plot(kind='bar', color='purple')
3 plt.title('Distribution of Countries')
4 plt.xlabel('Country')
5 plt.ylabel('Frequency')
6
7 print(top10['country_name'].nunique())
8
9 plt.title('Distribution of Country')
10 plt.xticks(rotation=90)
11 plt.show()
```

94



```
1 spearman_corr = top10[['weekly_rank', 'cumulative_weeks_in_top_10']].corr(meth
2 print(spearman_corr)
```

```
                            weekly_rank   cumulative_weeks_in_top_10
weekly_rank                    1.000000                     0.028064
cumulative_weeks_in_top_10     0.028064                     1.000000
```

```
1 plt.figure(figsize=(10, 10))
2 plt.scatter(top10['weekly_rank'], top10['cumulative_weeks_in_top_10'], color='
3 plt.title('Weekly Rank vs. Cumulative Weeks in Top 10')
4 plt.xlabel('Weekly Rank')
5 plt.ylabel('Cumulative Weeks in Top 10')
6 plt.grid(True)
7 plt.show()
```



Weekly Rank vs. Cumulative Weeks in Top 10

```
1 plt.figure(figsize=(10, 3))
2 plt.scatter(top10['cumulative_weeks_in_top_10'], top10['category'], color='gol
3 plt.xlabel('Cumulative Weeks in Top 10')
4 plt.ylabel('Category')
5 plt.title('Category vs. Cumulative Weeks in Top 10')
6 plt.xticks(rotation=90)
7 plt.grid(axis='x')
8 plt.show()
```

```
1 plt.figure(figsize=(25, 6))
2 plt.scatter(top10['country_name'], top10['cumulative_weeks_in_top_10'], color=
3 plt.xlabel('Country')
4 plt.ylabel('Cumulative Weeks in Top 10')
5 plt.title('Scatter Plot of Country vs Cumulative Weeks in Top 10')
6 plt.xticks(rotation=90)
7 plt.grid(True)
8 plt.show()
```



## Initial Results and Code

Does date correlate with shows that have larger cumulative weeks in Top 10?

```
1 #Calculate the average cumulative_weeks_in_top_10 for each month
2 top10['week'] = pd.to_datetime(top10['week'])
3 top10['month'] = top10['week'].dt.month
4
5 average_weeks_top_10_monthly = top10.groupby('month')['cumulative_weeks_in_top
6 print(average_weeks_top_10_monthly)
7
```

```
 8 #view distribution of the average cumulative_weeks_in_top_10 for each month
 9 plt.bar(average_weeks_top_10_monthly.index, average_weeks_top_10_monthly.value
10 plt.xlabel('Month')
11 plt.ylabel('Average Cumulative Weeks in Top 10')
12 plt.title('Average Cumulative Weeks in Top 10 by Month')
13 plt.show()
14
15 #The data is non-normally distributed thus we must use non-parametric statisti
16 #determine the significance of the relationship between month and cumulative we
17 from scipy.stats import kruskal
18
19 data_by_month = [group.values for name, group in top10.groupby('month')['cumul
20 print(data_by_month)
21
22 kruskal_result = kruskal(*data_by_month)
23 print(kruskal_result)
```

```
month
1     3.827766
2     3.475399
3     3.834274
4     4.151882
5     4.122366
6     3.916801
7     3.080856
8     3.216422
9     2.750665
10    2.745319
11    3.424202
12    3.754388
Name: cumulative_weeks_in_top_10, dtype: float64
```



Average Cumulative Weeks in Top 10 by Month

```
[array([ 2,  3,  1, ..., 21,  1,  3]), array([1, 2, 1, ..., 3, 1, 3]), array([
KruskalResult(statistic=798.6599409998698, pvalue=3.674882741765935e-164)
```

```
1 #view distribution of select months (samples) to see if samples are normally d
2 january_data = top10[top10['month'] == 1]['cumulative_weeks_in_top_10']
3
4 # Plot histogram for the distribution of cumulative weeks in the top 10 for Ja
5 plt.hist(january_data, bins=10, color='red')
6 plt.xlabel('Cumulative Weeks in Top 10')
7 plt.ylabel('Frequency')
8 plt.title('Distribution of Cumulative Weeks in Top 10 for January')
9 plt.show()
```



```
1 from itertools import combinations
2 from scipy.stats import ttest_ind
3
4 months = top10['month'].unique()
5
6 #Assuming the dataset is sufficiently large, we can perform t-tests on the
7 #average cumualtive weeks in top 10 between each month
8
9 ''' from these results we can infer which months have significant relationship
```

```
10 whether there is a TV consumption trend'''
11
12 for month1, month2 in combinations(months, 2):
13     data_month1 = top10[top10['month'] == month1]['cumulative_weeks_in_top_10'
14     data_month2 = top10[top10['month'] == month2]['cumulative_weeks_in_top_10'
15
16     t_statistic, p_value = ttest_ind(data_month1, data_month2)
17
18     print(f"t-test between month {month1} and month {month2}:")
19     print("t-statistic:", t_statistic)
20     print("p-value:", p_value)
21     print()
```

```
t-test between month 8 and month 7:
t-statistic: 2.2826475001277426
p-value: 0.022457711439931027

t-test between month 8 and month 6:
t-statistic: -8.345107815884798
p-value: 7.52103500068677e-17

t-test between month 8 and month 5:
t-statistic: -11.433355372990585
p-value: 3.4038920819053535e-30

t-test between month 8 and month 4:
t-statistic: -11.723393867829902
p-value: 1.1965784751814039e-31

t-test between month 8 and month 3:
t-statistic: -8.15789419545575
p-value: 3.5866989719926787e-16

t-test between month 8 and month 2:
t-statistic: -3.6105253323563606
p-value: 0.0003062414208590005

t-test between month 8 and month 1:
t-statistic: -9.479776916863766
p-value: 2.7728148746402483e-21

t-test between month 8 and month 12:
t-statistic: -8.000897214186281
p-value: 1.294415927300952e-15

t-test between month 8 and month 11:
t-statistic: -3.2330023507875816
p-value: 0.0012267304830323638
```

```
t-test between month 8 and month 10:
t-statistic: 8.483800036374486
p-value: 2.3016125552583263e-17

t-test between month 8 and month 9:
t-statistic: 7.776981868394483
p-value: 7.745411536924857e-15

t-test between month 7 and month 6:
t-statistic: -9.456778357490943
p-value: 3.454243130848794e-21

t-test between month 7 and month 5:
t-statistic: -12.607706363093271
p-value: 2.4437547961980922e-36

t-test between month 7 and month 4:
t-statistic: -12.612229835380301
p-value: 2.351690185059873e-36

t-test between month 7 and month 3:
t-statistic: -9.239488276367874
p-value: 2.6741742418225577e-20
```

Are shows that make the top 10 in the United States more likely to make the top 10 in other countries? Western countries? English speaking countries?

```
1  # A list of every country for which top 10 data is included in the dataset
2  countries = top10["country_name"].unique()
3  print(countries)
4
5  print("\n")
6
7  num_countries = len(countries)
8  print(num_countries)
```

```
['Argentina' 'Australia' 'Austria' 'Bahamas' 'Bahrain' 'Bangladesh'
 'Belgium' 'Bolivia' 'Brazil' 'Bulgaria' 'Canada' 'Chile' 'Colombia'
 'Costa Rica' 'Croatia' 'Cyprus' 'Czech Republic' 'Denmark'
 'Dominican Republic' 'Ecuador' 'Egypt' 'El Salvador' 'Estonia' 'Finland'
 'France' 'Germany' 'Greece' 'Guadeloupe' 'Guatemala' 'Honduras'
 'Hong Kong' 'Hungary' 'Iceland' 'India' 'Indonesia' 'Ireland' 'Israel'
 'Italy' 'Jamaica' 'Japan' 'Jordan' 'Kenya' 'Kuwait' 'Latvia' 'Lebanon'
 'Lithuania' 'Luxembourg' 'Malaysia' 'Maldives' 'Malta' 'Martinique'
 'Mauritius' 'Mexico' 'Morocco' 'Netherlands' 'New Caledonia'
 'New Zealand' 'Nicaragua' 'Nigeria' 'Norway' 'Oman' 'Pakistan' 'Panama'
 'Paraguay' 'Peru' 'Philippines' 'Poland' 'Portugal' 'Qatar' 'Romania'
 'Russia' 'Réunion' 'Saudi Arabia' 'Serbia' 'Singapore' 'Slovakia'
 'Slovenia' 'South Africa' 'South Korea' 'Spain' 'Sri Lanka' 'Sweden'
 'Switzerland' 'Taiwan' 'Thailand' 'Trinidad and Tobago' 'Turkey'
 'Ukraine' 'United Arab Emirates' 'United Kingdom' 'United States'
 'Uruguay' 'Venezuela' 'Vietnam']
```

```
94
```

```python
1  #The shows/movies that make the Top 10 in the United States
2  us_shows = top10[top10['country_name'] == 'United States']
3
4  #count the frequency of how many times each show_title appears in the Top 10 fo
5  us_show_frequency = us_shows.groupby('show_title').size()
6  us_show_frequency_sorted = us_show_frequency.sort_values(ascending=False)
7
8  #the top 10 most frequently occuring shows/movies that make the top 10 in the l
9  us_top10 = us_show_frequency_sorted.head(10)
10 print(us_top10)
11
```

```
show_title
CoComelon         52
Stranger Things   43
Ozark             23
Manifest          20
All American      18
Virgin River      15
Bridgerton        12
You               11
The Witcher       11
Squid Game        11
dtype: int64
```

```
1  #The shows/movies that make the Top 10 in Canada
2  can_shows = top10[top10['country_name'] == 'Canada']
3
4  #count the frequency of how many times each show_title appears in the Top 10 f
5  can_show_frequency = can_shows.groupby('show_title').size()
6  can_show_frequency_sorted = can_show_frequency.sort_values(ascending=False)
7
8  #the top 10 most frequently occuring shows/movies that make the top 10 in Cana
9  can_top10 = can_show_frequency_sorted.head(10)
10 print(can_top10)
```

```
   show_title
   Stranger Things    37
   Ozark              21
   Blindspot          21
   Manifest           14
   Young Sheldon      13
   Maid               12
   Bridgerton         12
   The Witcher        11
   Love Is Blind      11
   You                11
   dtype: int64
```

```
1  '''Count the number of times Top 10 US shows/movies make the Top 10 in other c
2
3  top10_frequency_by_country = {}
4
5  #Compute frequency of each show in the top 10 for each country that is not the
6  for country in top10['country_name'].unique():
7      if country != 'United States':
8          country_shows = top10[top10['country_name'] == country]
9          country_show_frequency = country_shows['show_title'].value_counts().he
10         top10_frequency_by_country[country] = country_show_frequency
11
12 #DataFrame to store contingency table
13 contingency_table = pd.DataFrame(index=us_top10.index, columns=top10_frequency_
14
15 #Fill the contingency table with frequency counts
16 for country, country_frequency in top10_frequency_by_country.items():
17     for show in us_top10.index:
18         # Fill in the frequency count for each show in the top 10 for the curr
19         contingency_table.loc[show, country] = country_frequency.get(show, 0)
20
21 #Fill the contingency table for the US
```

```
22 for show in us_top10.index:
23     contingency_table.loc[show, 'United States'] = us_show_frequency.get(show,
24
25 print(contingency_table)
```

|                | Argentina | Australia | Austria | Bahamas | Bahrain | Bangladesh | \ |
|----------------|-----------|-----------|---------|---------|---------|------------|---|
| show_title     |           |           |         |         |         |            |   |
| CoComelon      | 0         | 0         | 0       | 17      | 0       | 0          |   |
| Stranger Things| 24        | 41        | 40      | 30      | 36      | 44         |   |
| Ozark          | 0         | 14        | 0       | 13      | 0       | 0          |   |
| Manifest       | 15        | 26        | 20      | 16      | 0       | 21         |   |
| All American   | 0         | 0         | 0       | 0       | 0       | 0          |   |
| Virgin River   | 0         | 13        | 0       | 12      | 0       | 0          |   |
| Bridgerton     | 0         | 14        | 13      | 11      | 13      | 0          |   |
| You            | 0         | 12        | 15      | 0       | 11      | 0          |   |
| The Witcher    | 0         | 0         | 12      | 0       | 12      | 0          |   |
| Squid Game     | 11        | 11        | 12      | 0       | 16      | 22         |   |

|                | Belgium | Bolivia | Brazil | Bulgaria | ... | Thailand | \ |
|----------------|---------|---------|--------|----------|-----|----------|---|
| show_title     |         |         |        |          | ... |          |   |
| CoComelon      | 0       | 0       | 0      | 0        | ... | 0        |   |
| Stranger Things| 41      | 27      | 31     | 46       | ... | 21       |   |
| Ozark          | 0       | 0       | 0      | 0        | ... | 0        |   |
| Manifest       | 15      | 0       | 0      | 20       | ... | 0        |   |
| All American   | 0       | 0       | 0      | 0        | ... | 0        |   |
| Virgin River   | 12      | 0       | 0      | 0        | ... | 0        |   |
| Bridgerton     | 14      | 0       | 0      | 15       | ... | 0        |   |
| You            | 12      | 0       | 0      | 15       | ... | 0        |   |
| The Witcher    | 11      | 0       | 0      | 15       | ... | 0        |   |
| Squid Game     | 11      | 0       | 10     | 16       | ... | 0        |   |

|                | Trinidad and Tobago | Turkey | Ukraine | United Arab Emirates | \ |
|----------------|---------------------|--------|---------|----------------------|---|
| show_title     |                     |        |         |                      |   |
| CoComelon      | 0                   | 0      | 0       | 0                    |   |
| Stranger Things| 34                  | 34     | 46      | 34                   |   |
| Ozark          | 12                  | 0      | 28      | 10                   |   |
| Manifest       | 16                  | 0      | 0       | 0                    |   |
| All American   | 0                   | 0      | 0       | 0                    |   |
| Virgin River   | 12                  | 0      | 0       | 0                    |   |
| Bridgerton     | 13                  | 0      | 35      | 13                   |   |
| You            | 13                  | 19     | 21      | 10                   |   |
| The Witcher    | 11                  | 0      | 27      | 11                   |   |
| Squid Game     | 0                   | 16     | 0       | 14                   |   |

|                | United Kingdom | Uruguay | Venezuela | Vietnam | United States |
|----------------|----------------|---------|-----------|---------|---------------|
| show_title     |                |         |           |         |               |
| CoComelon      | 0              | 0       | 0         | 0       | 52.0          |
| Stranger Things| 46             | 21      | 19        | 25      | 43.0          |
| Ozark          | 17             | 0       | 0         | 0       | 23.0          |

```
Manifest                              0      16      15       0      20.0
All American                          0       0       0       0      18.0
Virgin River                         12       0       0       0      15.0
Bridgerton                           13       0       0       0      12.0
You                                  13      12       0       0      11.0
The Witcher                           0       0       0       0      11.0
Squid Game                           10       0       0       0      11.0

[10 rows x 94 columns]
```

```python
1  '''Calculate whether the results in the contingency table above are significant
2
3  import numpy as np
4  from scipy.stats import chi2_contingency
5
6  p_values = {}
7
8  for country in contingency_table.columns:
9      # Extract observed frequencies for the current country
10     observed_frequencies = contingency_table[country].values.astype(float)
11     row_totals = contingency_table.sum(axis=1)
12     column_totals = contingency_table.sum(axis=0)
13     expected_frequencies = np.outer(row_totals, column_totals) / row_totals.su
14     expected_frequencies = expected_frequencies[:, contingency_table.columns.g
15
16     chi2_stat, p_value, _, _ = chi2_contingency([observed_frequencies, expecte
17
18     p_values[country] = p_value
19
20
21 for country, p_value in p_values.items():
22     print(f"{country}: {p_value}")
23
24
25 print("\n")
26
27 '''Here, significant countries suggests that in these countries, a show/movie
28 the Top 10 in the US is likely to make the Top 10 in the listed country. Possi
29 influence the markets in these other countries.'''
30
31 significant_countries = [country for country, p_value in p_values.items() if p
32 print("List of Significant Countries", significant_countries)
```

```
Argentina: 0.028787031633686188
Australia: 0.014589738469597569
```

```
Austria: 0.09908073729779852
Bahamas: 1.6428018500900734e-07
Bahrain: 0.00040544203937159415
Bangladesh: 0.00016475360623926946
Belgium: 0.075575420324457381
Bolivia: 0.00702701050954094235
Brazil: 0.002305190808831503
Bulgaria: 0.05274111160975641
Canada: 3.7241957280345492e-06
Chile: 0.0027174447198394522
Colombia: 0.004851091032492417
Costa Rica: 0.022772288597190448
Croatia: 0.020101224321794453
Cyprus: 0.036197197212430726
Czech Republic: 1.9325766371444175e-05
Denmark: 0.007930065789479997
Dominican Republic: 0.0206159800194768
Ecuador: 0.03946046413725036
Egypt: 2.577880163511151e-05
El Salvador: 0.06337385755013583
Estonia: 0.00018414750490193212
Finland: 0.03380490569942349
France: 0.023836926402612604
Germany: 0.0017473424514078122
Greece: 0.021830932614843312
Guadeloupe: 0.02470583816573554
Guatemala: 0.04796088964421329
Honduras: 0.07860712593283108
Hong Kong: 0.02204444827856596
Hungary: 1.5212543936266608e-06
Iceland: 0.0004255930132056846
India: 6.65118228025686e-05
Indonesia: 0.012185316340483103
Ireland: 4.3212020022413026e-12
Israel: 0.002209538283476711
Italy: 0.0011210599199247232
Jamaica: 0.05537084550764695
Japan: 0.0019156756229678359
Jordan: 2.9884026444445108e-05
Kenya: 0.0505726389320982
Kuwait: 0.00012518788664030034
Latvia: 0.0036623239495632374
Lebanon: 0.00014233708594856714
Lithuania: 0.01013552331735906
Luxembourg: 0.03257002315611079
Malaysia: 0.014188477714176969
Maldives: 0.004153869248331972
Malta: 0.0008714551433422026
Martinique: 0.034957038135512285
Mauritius: 0.015399724013805837
```

```
Mexico: 0.0027843972455754208
Morocco: 0.000611412614606106
Netherlands: 5.587868575947894e-05
New Caledonia: 0.021102821293318457
New Zealand: 1.9388276955265418e-05
Nicaragua: 0.04409396244769613
```

```
 1  print("Number of significant countries: ", len(significant_countries))
 2
 3  #English language countries as per The University of Tennessee Knoxville
 4  #Note this list does not include the UK, Australia but I have added them in
 5
 6  #strong limitation is the definition of english-language country
 7  #not every country is listed and not every country was checked against the 94
 8  '''https://gradschool.utk.edu/future-students/office-of-graduate-admissions/ap
 9  admission-requirements/testing-requirements/countries-with-english-as-official-
10
11  english_speaking_countries = [
12      'Anguilla', 'Antigua and Barbuda', 'Bahamas', 'Barbados', 'Belize', 'Belgi
13      'British Virgin Islands', 'Burundi', 'Cameroon', 'Canada', 'Cayman Islands
14      'Dominica', 'Fiji', 'Gambia', 'Ghana', 'Grenada', 'Guyana', 'Hong Kong', '
15      'Liberia', 'Malawi', 'Malta', 'Marshall Islands', 'Micronesia', 'Namibia',
16      'Nigeria', 'Niue', 'Norfolk Island', 'Northern Mariana Islands', 'Pakistan
17      'Philippines', 'Pitcairn Islands', 'Rwanda', 'Saint Kitts and Nevis', 'Sai
18      'Sierra Leone', 'Singapore', 'Sint Maarten', 'Solomon Islands', 'Somalia',
19      'Swaziland', 'Tanzania', 'Tonga', 'Trinidad and Tobago', 'Turks and Caicos
20      'Zimbabwe', 'United Kingdom', 'Australia']
21
22  # Of the significant countries, which are english-speaking?
23  num_significant_english_speaking_countries = 0
24
25  for country in significant_countries:
26      if country in english_speaking_countries:
27          num_significant_english_speaking_countries = num_significant_english_sp
28
29  print("Number of significant english speaking countries: ", num_significant_en
```

```
Number of significant countries:  52
Number of significant english speaking countries:  10
```

```
 1  '''Can we use classification to predict whether Top 10 US shows make the Top 1
 2
 3  from sklearn.model_selection import train_test_split
 4  from sklearn.tree import DecisionTreeClassifier, plot_tree
```

```python
 5 from sklearn.metrics import confusion_matrix, classification_report
 6 import matplotlib.pyplot as plt
 7 import pandas as pd
 8 from sklearn.preprocessing import OneHotEncoder
 9
10 us_shows = top10[top10['country_name'] == 'United States']
11
12 # Determine target variable indicating whether the show made it to the top 10
13 target_variable = (us_shows['show_title'].isin(top10[top10['country_name'] ==
14
15 # Combine features and target variable into a dataset
16 prepared_dataset = pd.concat([us_shows[['show_title', 'category', 'weekly_rank
17 prepared_dataset.columns = ['show_title', 'category', 'weekly_rank', 'country_
18
19 print(prepared_dataset)
20
21 prepared_dataset['category'] = prepared_dataset['category'].map({'Films': 0, '
22 X = prepared_dataset[['weekly_rank', 'category']]
23 y = prepared_dataset['is_top_10_canada']
24
25 # Split the data into training and testing sets
26 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, rando
27
28 # Instantiate the decision tree classifier
29 classifier = DecisionTreeClassifier(max_depth=4, random_state=42)
30
31 # Train the decision tree classifier on the training data
32 classifier.fit(X_train, y_train)
33
34 # Make predictions on the testing data
35 y_pred = classifier.predict(X_test)
36
37 # Print confusion matrix and classification report
38 conf_matrix = confusion_matrix(y_test, y_pred)
39 class_report = classification_report(y_test, y_pred)
40
41 print("Confusion Matrix:")
42 print(conf_matrix)
43 print("\nClassification Report:")
44 print(class_report)
45
46 # Visualize the decision tree
47 plt.figure(figsize=(15, 10))
48 plot_tree(classifier, feature_names=X.columns, class_names=["Not in Canada Top
49 plt.show()
```

50

```
                                       show_title category  weekly_rank  \
107500                                   Day Shift    Films            1
107501                              Look Both Ways    Films            2
107502  Untold: The Girlfriend Who Didn't Exist    Films            3
107503                                   Uncharted    Films            4
107504                                      Sing 2    Films            5
...                                            ...      ...          ...
108695                                   CoComelon       TV            6
108696                         Newly Rich, Newly Poor       TV            7
108697                                 Sweet Tooth       TV            8
108698                                   CoComelon       TV            9
108699                                   CoComelon       TV           10

           country_name  is_top_10_canada
107500  United States                 1
107501  United States                 1
107502  United States                 1
107503  United States                 0
107504  United States                 0
...              ...               ...
108695  United States                 0
108696  United States                 0
108697  United States                 0
108698  United States                 0
108699  United States                 0

[1200 rows x 5 columns]
Confusion Matrix:
[[ 35  48]
 [ 26 131]]

Classification Report:
              precision    recall  f1-score   support

           0       0.57      0.42      0.49        83
           1       0.73      0.83      0.78       157

    accuracy                           0.69       240
   macro avg       0.65      0.63      0.63       240
weighted avg       0.68      0.69      0.68       240
```
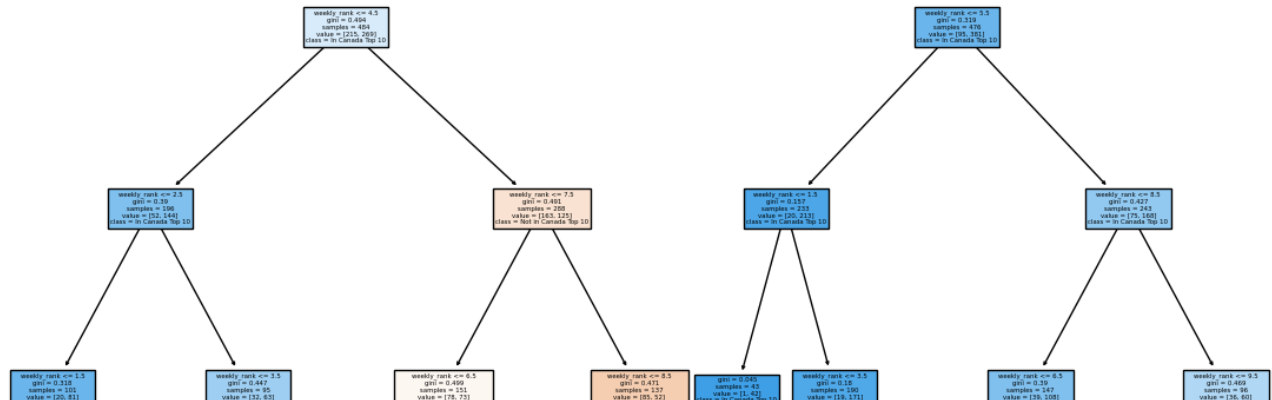
```
1  '''Can we use classification to predict whether Top 10 US shows make the Top 1
2  Using cross validation and different random state'''
3  from sklearn.model_selection import train_test_split, cross_val_score
4  from sklearn.tree import DecisionTreeClassifier, plot_tree
5  from sklearn.metrics import confusion_matrix, classification_report
6  import matplotlib.pyplot as plt
7  import pandas as pd
8  from sklearn.preprocessing import OneHotEncoder
9
10 # Assuming top10 and us_shows are defined earlier
11
12 # Determine target variable indicating whether the show made it to the top 10
13 target_variable = (us_shows['show_title'].isin(top10[top10['country_name'] ==
14
15 # Combine features and target variable into a dataset
16 prepared_dataset = pd.concat([us_shows[['show_title', 'category', 'weekly_rank
17 prepared_dataset.columns = ['show_title', 'category', 'weekly_rank', 'country_
18
19 prepared_dataset['category'] = prepared_dataset['category'].map({'Films': 0, '
20 X = prepared_dataset[['weekly_rank', 'category']]
21 y = prepared_dataset['is_top_10_canada']
22
23 # Split the data into training and testing sets
24 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, rando
25
26 # Instantiate the decision tree classifier
27 classifier = DecisionTreeClassifier(max_depth=4, random_state=0)
28
29 # Perform cross-validation
30 cv_scores = cross_val_score(classifier, X_train, y_train, cv=5)
31
32 print("Cross-validation Scores:", cv_scores)
33 print("Mean CV Score:", cv_scores.mean())
34
```

```
35 # Train the decision tree classifier on the training data
36 classifier.fit(X_train, y_train)
37
38 # Make predictions on the testing data
39 y_pred = classifier.predict(X_test)
40
41 # Print confusion matrix and classification report
42 conf_matrix = confusion_matrix(y_test, y_pred)
43 class_report = classification_report(y_test, y_pred)
44
45 print("Confusion Matrix:")
46 print(conf_matrix)
47 print("\nClassification Report:")
48 print(class_report)
49
50 # Visualize the decision tree
51 plt.figure(figsize=(15, 10))
52 plot_tree(classifier, feature_names=X.columns, class_names=["Not in Canada Top
53 plt.show()
```

```
Cross-validation Scores: [0.72916667 0.70833333 0.68229167 0.734375   0.671875
Mean CV Score: 0.7052083333333333
Confusion Matrix:
[[ 35  46]
 [ 33 126]]

Classification Report:
              precision    recall  f1-score   support

           0       0.51      0.43      0.47        81
           1       0.73      0.79      0.76       159

    accuracy                           0.67       240
   macro avg       0.62      0.61      0.62       240
weighted avg       0.66      0.67      0.66       240
```
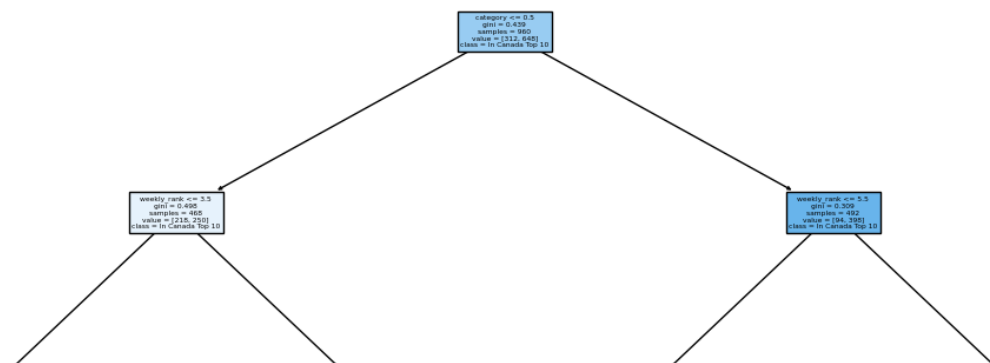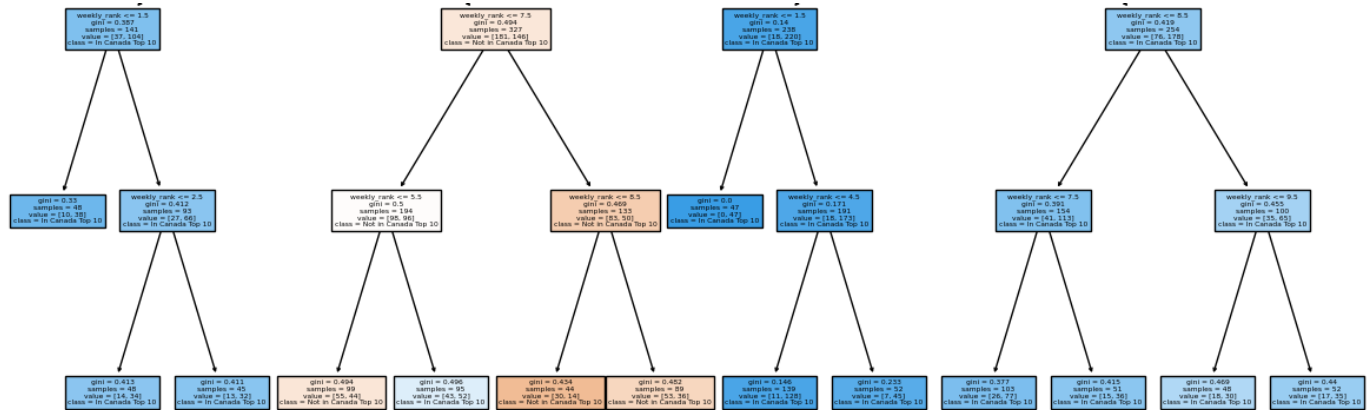
```
1 '''Can we use classification to predict whether Top 10 US shows make the Top 1(
2
3 from sklearn.model_selection import train_test_split
4 from sklearn.tree import DecisionTreeClassifier, plot_tree
5 from sklearn.metrics import confusion_matrix, classification_report
6 import matplotlib.pyplot as plt
7 import pandas as pd
8 from sklearn.preprocessing import OneHotEncoder
9
10 us_shows = top10[top10['country_name'] == 'United States']
11
12 # Determine target variable indicating whether the show made it to the top 10
13 target_variable = (us_shows['show_title'].isin(top10[top10['country_name'] ==
14
15 # Combine features and target variable into a dataset
16 prepared_dataset = pd.concat([us_shows[['show_title', 'category', 'weekly_rank
```

```python
17 prepared_dataset.columns = ['show_title', 'category', 'weekly_rank', 'country_
18
19 print(prepared_dataset)
20
21 prepared_dataset['category'] = prepared_dataset['category'].map({'Films': 0, '
22 X = prepared_dataset[['weekly_rank', 'category']]
23 y = prepared_dataset['is_top_10_Argentina']
24
25 # Split the data into training and testing sets
26 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, randor
27
28 # Instantiate the decision tree classifier
29 classifier = DecisionTreeClassifier(max_depth=4, random_state=42)
30
31 # Train the decision tree classifier on the training data
32 classifier.fit(X_train, y_train)
33
34 # Make predictions on the testing data
35 y_pred = classifier.predict(X_test)
36
37 # Print confusion matrix and classification report
38 conf_matrix = confusion_matrix(y_test, y_pred)
39 class_report = classification_report(y_test, y_pred)
40
41 print("Confusion Matrix:")
42 print(conf_matrix)
43 print("\nClassification Report:")
44 print(class_report)
45
46 # Visualize the decision tree
47 plt.figure(figsize=(15, 10))
48 plot_tree(classifier, feature_names=X.columns, class_names=["Not in Argentina
49 plt.show()
50
```

```
                                    show_title category  weekly_rank  \
    107500                          Day Shift    Films            1
    107501                     Look Both Ways    Films            2
    107502  Untold: The Girlfriend Who Didn't Exist    Films            3
    107503                          Uncharted    Films            4
    107504                             Sing 2    Films            5
    ...                                    ...      ...          ...
    108695                          CoComelon       TV            6
    108696             Newly Rich, Newly Poor       TV            7
    108697                        Sweet Tooth       TV            8
    108698                          CoComelon       TV            9
```

```
108699                                        CoComelon        TV              10

           country_name  is_top_10_Argentina
107500   United States                      1
107501   United States                      1
107502   United States                      0
107503   United States                      0
107504   United States                      0
...                ...                    ...
108695   United States                      0
108696   United States                      1
108697   United States                      0
108698   United States                      0
108699   United States                      0

[1200 rows x 5 columns]
Confusion Matrix:
[[80 40]
 [49 71]]

Classification Report:
              precision    recall  f1-score   support

           0       0.62      0.67      0.64       120
           1       0.64      0.59      0.61       120

    accuracy                           0.63       240
   macro avg       0.63      0.63      0.63       240
weighted avg       0.63      0.63      0.63       240
```
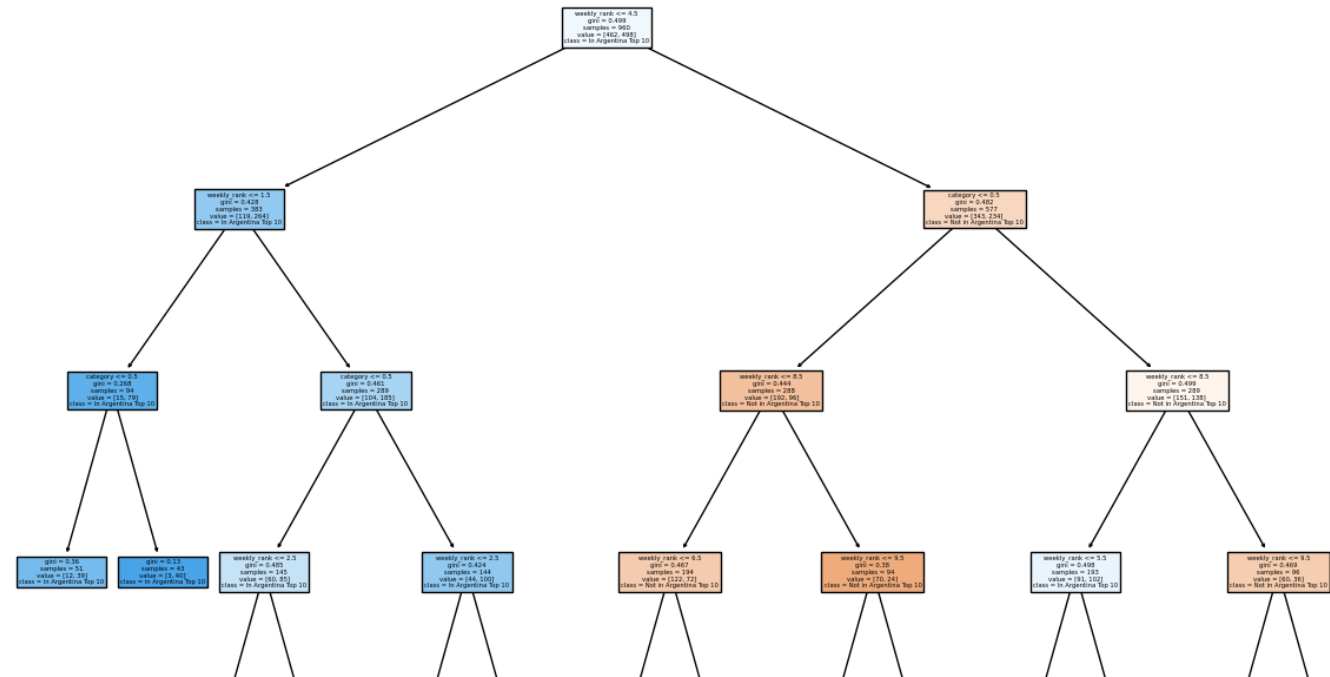
Do TV shows with more seasons make the top 10 list more often? Have larger number of cumulative weeks in Top 10?

```
1 #subset data to include rows where season_title is included
2 subset = top10[top10['season_title'].notnull()]
3 print(subset.info())
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 54668 entries, 10 to 112299
Data columns (total 9 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   country_name              54668 non-null  object
 1   country_iso2              54668 non-null  object
 2   week                      54668 non-null  datetime64[ns]
 3   category                  54668 non-null  object
 4   weekly_rank               54668 non-null  int64
 5   show_title                54668 non-null  object
 6   season_title              54668 non-null  object
 7   cumulative_weeks_in_top_10  54668 non-null  int64
 8   month                     54668 non-null  int64
dtypes: datetime64[ns](1), int64(3), object(5)
memory usage: 4.2+ MB
None
```

```
1 #view format of season_title entries
2 print(subset['season_title'].head(60))
3
4 #add column season_number of the numeric part of season_title
5 subset['season_number'] = subset['season_title'].str.extract(r'(\d+)')
6
7 '''some season_title entries do not include a number'''
8 print(subset['season_number'].head(60))
```

```
10          Pasión de Gavilanes: Season 2
11                  Another Self: Season 1
12          Pasión de Gavilanes: Season 1
13                      Manifest: Season 1
14                  The Sandman: Season 1
15    Extraordinary Attorney Woo: Season 1
16                     High Heat: Season 1
17                      Manifest: Season 2
18                      Manifest: Season 3
19            Never Have I Ever: Season 3
30          Pasión de Gavilanes: Season 2
```

```
31                      Another Self: Season 1
32                        The Sandman: Season 1
33              Pasión de Gavilanes: Season 1
34                            Manifest: Season 1
35       Extraordinary Attorney Woo: Season 1
36                                Alba: Season 1
37                            Manifest: Season 2
38                            Manifest: Season 3
39                 Never Have I Ever: Season 3
50              Pasión de Gavilanes: Season 2
51                            Manifest: Season 1
52                      Another Self: Season 1
53              Pasión de Gavilanes: Season 1
54                                Alba: Season 1
55           Keep Breathing: Limited Series
56                        The Sandman: Season 1
57                            Manifest: Season 2
58                    Virgin River: Season 4
59                            Manifest: Season 3
70              Pasión de Gavilanes: Season 2
71                                Alba: Season 1
72                            Manifest: Season 1
73                    Virgin River: Season 4
74              Pasión de Gavilanes: Season 1
75                            Manifest: Season 2
76                          Stranger Things 4
77           Keep Breathing: Limited Series
78                            Manifest: Season 3
79                      Another Self: Season 1
90                                Alba: Season 1
91              Pasión de Gavilanes: Season 2
92                            Manifest: Season 1
93                    Virgin River: Season 4
94                          Stranger Things 4
95                   Resident Evil: Season 1
96              Pasión de Gavilanes: Season 1
97                            Manifest: Season 2
98        Café con aroma de mujer: Season 1
99                    Rebelde Way: Temporada 1
110                         Stranger Things 4
111                               Alba: Season 1
112                  Resident Evil: Season 1
113              The Longest Night: Season 1
114                           Manifest: Season 1
115             Pasión de Gavilanes: Season 1
116                          Capitani: Season 2
117       Café con aroma de mujer: Season 1
118                         Stranger Things 2
```

```
1 #subset data to include rows where season_number is included
2 subset2 = subset[subset['season_number'].notnull()]
3 print(subset2.info())
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 50640 entries, 10 to 112299
Data columns (total 10 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   country_name              50640 non-null  object
 1   country_iso2              50640 non-null  object
 2   week                      50640 non-null  datetime64[ns]
 3   category                  50640 non-null  object
 4   weekly_rank               50640 non-null  int64
 5   show_title                50640 non-null  object
 6   season_title              50640 non-null  object
 7   cumulative_weeks_in_top_10  50640 non-null  int64
 8   month                     50640 non-null  int64
 9   season_number             50640 non-null  object
dtypes: datetime64[ns](1), int64(3), object(6)
memory usage: 4.2+ MB
None
```

```
1  '''The method of isolating season_number includes many errors.
2  For our sake we will look at only entries where season number is 10 or smaller
3
4  #many errors exist when isolating season number from season_title
5  unique_season_numbers = subset2['season_number'].unique()
6  print(unique_season_numbers)
7
8  #convert season_number to int type
9  subset2.loc[:, 'season_number'] = subset2['season_number'].fillna(-5).astype(i
10
11 subset3 = subset2.loc[(subset2['season_number'] > 0) & (subset2['season_number
12
13 unique_season_numbers_again = subset3['season_number'].unique()
14 print(unique_season_numbers_again)
```

```
   ['2' '1' '3' '4' '6' '5' '42' '81' '17' '100' '8' '99' '7' '15' '11' '9'
    '14' '245' '101' '1988' '18' '24' '10' '2011' '13' 'ر' '12' '56' '892'
    '2045' '2020' '20' '97' '2022' '2021' '800' '1867' '60']
   [ 2  1  3  4  6  5  8  7  9 10]
   <ipython-input-29-4c5af1bdabe4>:9: SettingWithCopyWarning:
   A value is trying to be set on a copy of a slice from a DataFrame.
   Try using .loc[row_indexer,col_indexer] = value instead

   See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/st
     subset2.loc[:, 'season_number'] = subset2['season_number'].fillna(-5).astype
   <ipython-input-29-4c5af1bdabe4>:9: DeprecationWarning: In a future version, `d
     subset2.loc[:, 'season_number'] = subset2['season_number'].fillna(-5).astype
```
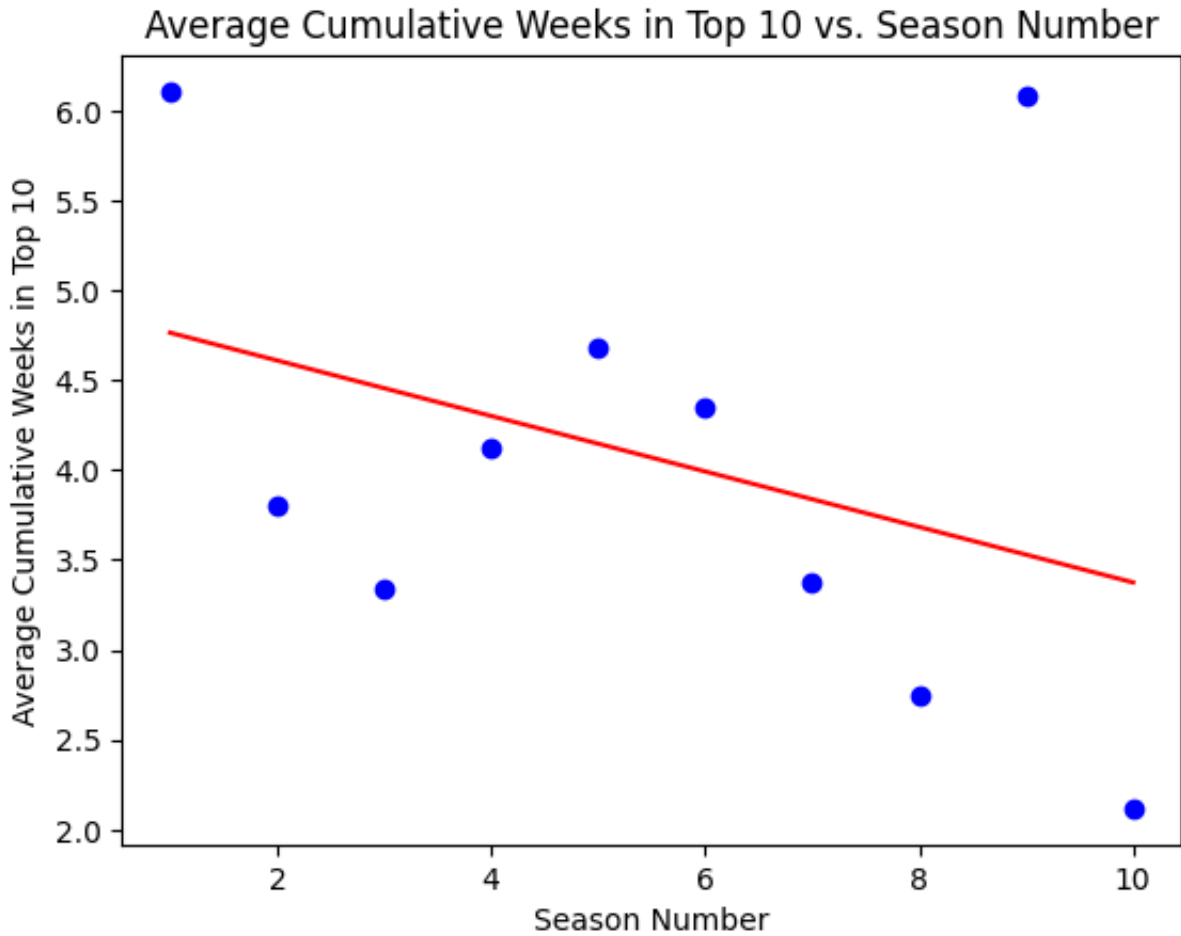
```
1  '''subset3 is our dataset now'''
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from scipy.stats import pearsonr
5
6  #average cumulative weeks in top 10 for each season number 1 through 10
7  average_cumulative_weeks = subset3.groupby('season_number')['cumulative_weeks_
8  average_cumulative_weeks.columns = ['season_number', 'average_cumulative_weeks
9  print(average_cumulative_weeks)
10
11 #initiliaze variables for statistical analysis
12 season_number = average_cumulative_weeks['season_number']
13 average_cumulative_weeks_in_top_10 = average_cumulative_weeks['average_cumulat
14
15 '''In this context season number can be numeric instead of rank because
16 we want to see if higher season number shows higher average cumulative weeks i
17
```
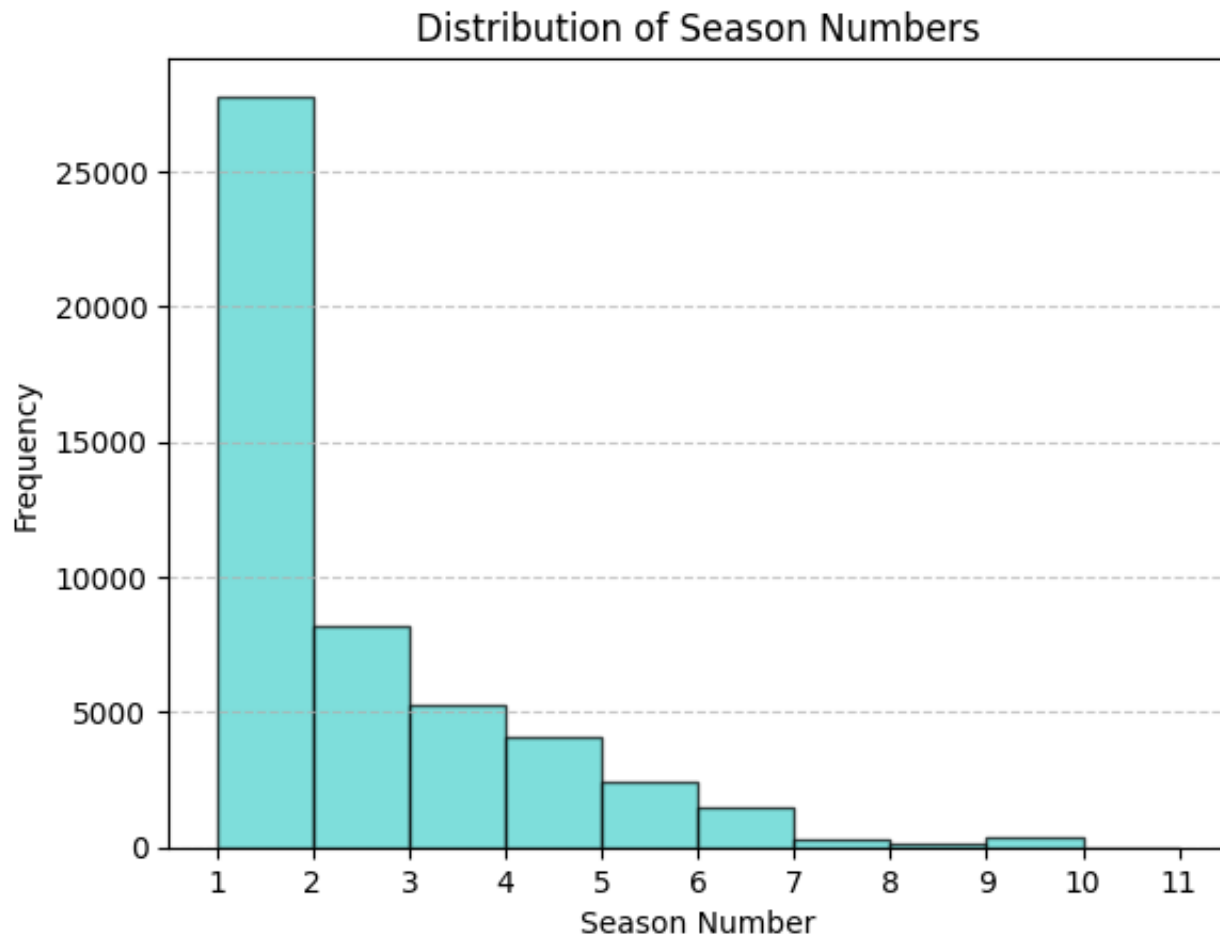
```
18 correlation_coefficient, p_value = pearsonr(season_number, average_cumulative_
19 print("Pearson correlation coefficient:", correlation_coefficient)
20 print("P-value:", p_value)
21
22 plt.scatter(season_number, average_cumulative_weeks_in_top_10, color='blue', l
23 plt.plot(season_number, np.poly1d(np.polyfit(season_number, average_cumulative
24 plt.xlabel('Season Number')
25 plt.ylabel('Average Cumulative Weeks in Top 10')
26 plt.title('Average Cumulative Weeks in Top 10 vs. Season Number')
27 plt.show()
```

```
     season_number   average_cumulative_weeks_in_top_10
0             1                              6.106364
1             2                              3.799366
2             3                              3.333964
3             4                              4.119763
4             5                              4.675585
5             6                              4.346075
6             7                              3.378277
7             8                              2.740964
8             9                              6.077586
9            10                              2.117647
Pearson correlation coefficient: -0.35851553225279903
P-value: 0.3090139342720427
```



Average Cumulative Weeks in Top 10 vs. Season Number

```
1 '''We can view a histogram of season_number to determine which number of seaso
2
3 plt.hist(subset3['season_number'], bins=range(1, 12), edgecolor='black', color
4 plt.xlabel('Season Number')
5 plt.ylabel('Frequency')
6 plt.title('Distribution of Season Numbers')
7 plt.xticks(range(1, 12))
8 plt.grid(axis='y', linestyle='--', alpha=0.7)
9 plt.show()
```
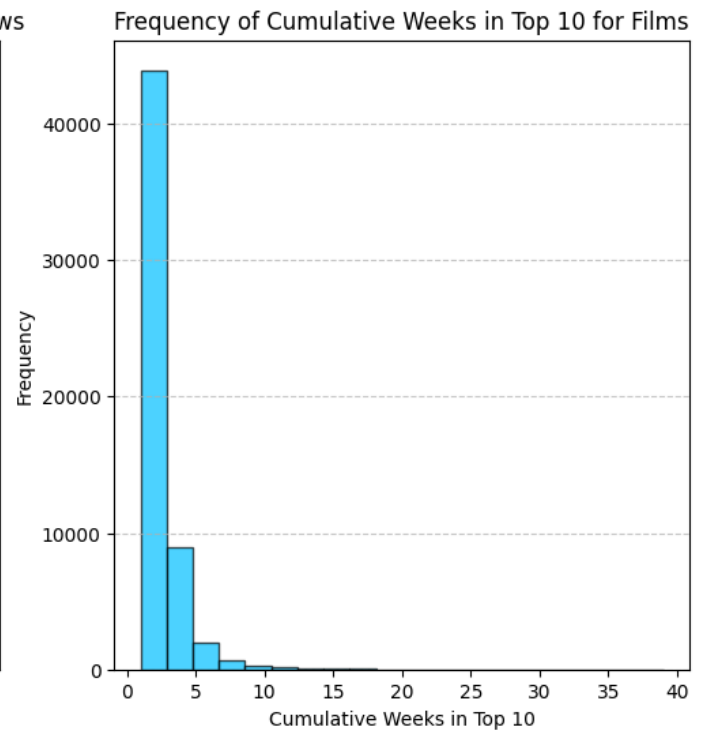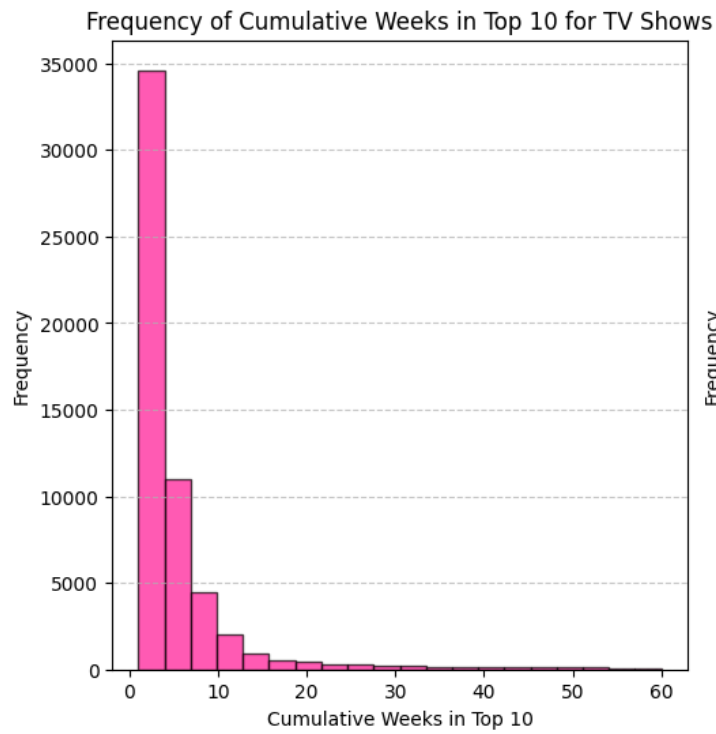


Are TV shows or movies more likely to spend a longer amount of cumulative time on the Top 10 list?

RECALL: Some TV shows made the Top 10 list for up to 60 consecutive weeks, whereas, movies only made the list for up to 40 weeks (yellow scatterplot above)

```
1 '''lets look at the cumulative weeks on top 10 frequency distribution for TV sl
```

```
 2
 3 tv_shows_subset = top10[top10['category'] == 'TV']
 4
 5 movies_shows_subset = top10[top10['category'] == 'Films']
 6
 7 fig, axes = plt.subplots(1, 2, figsize=(12, 6))
 8
 9 axes[0].hist(tv_shows_subset['cumulative_weeks_in_top_10'], bins=20, edgecolor:
10 axes[0].set_xlabel('Cumulative Weeks in Top 10')
11 axes[0].set_ylabel('Frequency')
12 axes[0].set_title('Frequency of Cumulative Weeks in Top 10 for TV Shows')
13 axes[0].grid(axis='y', linestyle='--', alpha=0.7)
14
15 axes[1].hist(movies_shows_subset['cumulative_weeks_in_top_10'], bins=20, edgec
16 axes[1].set_xlabel('Cumulative Weeks in Top 10')
17 axes[1].set_ylabel('Frequency')
18 axes[1].set_title('Frequency of Cumulative Weeks in Top 10 for Films')
19 axes[1].grid(axis='y', linestyle='--', alpha=0.7)
```

Frequency of Cumulative Weeks in Top 10 for TV Shows          Frequency of Cumulative Weeks in Top 10 for Films

```python
1  #Calculate average cumulative weeks in top 10 for films and TV categories
2  avg_cumulative_weeks_tv = tv_shows_subset['cumulative_weeks_in_top_10'].mean()
3  print(avg_cumulative_weeks_tv)
4
5  avg_cumulative_weeks_films = movies_shows_subset['cumulative_weeks_in_top_10']
6  print(avg_cumulative_weeks_films)
7
8  #samples are unpaired and non-normally distributed so we will perform wilcoxon
9  from scipy.stats import mannwhitneyu
10
11 statistic, p_value = mannwhitneyu(tv_shows_subset['cumulative_weeks_in_top_10'
12
13 print("Wilcoxon Rank Sum test statistic:", statistic)
14 print("P-value:", p_value)
```

```
4.936420302760463
2.000142475512021
Wilcoxon Rank Sum test statistic: 2152311801.0
P-value: 0.0
```