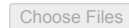


Import the CSV file

```

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

```

 all-weeks-countries.csv

- **all-weeks-countries.csv**(text/csv) - 6847558 bytes, last modified: 8/28/2022 - 100% done

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	

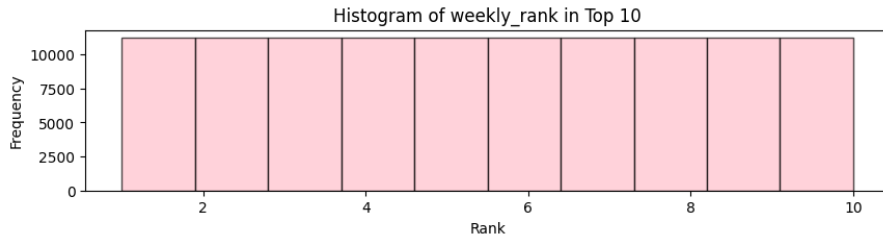
Preliminary analysis

```

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='pink')
5 plt.title('Histogram of weekly_rank in Top 10')
6 plt.xlabel('Rank')
7 plt.ylabel('Frequency')
8

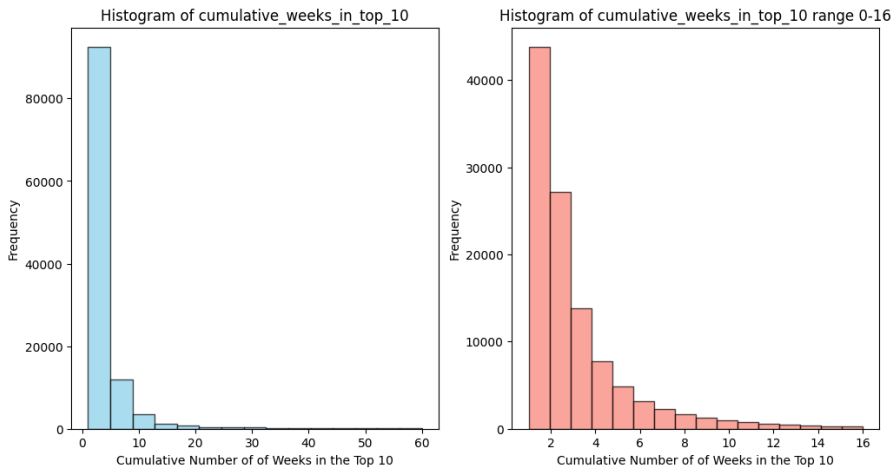
```

Text(0, 0.5, 'Frequency')



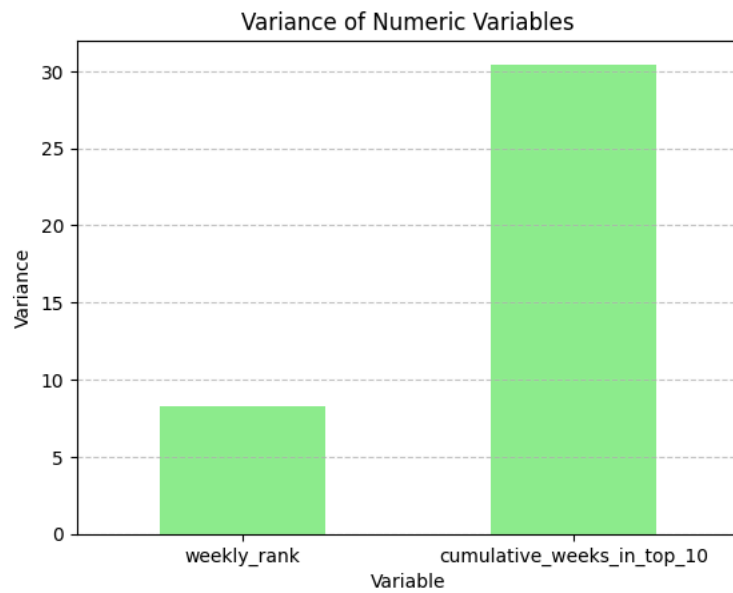
```
1 fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(12, 6))
2
3 axes[0].hist(top10['cumulative_weeks_in_top_10'], bins=15, edgecolor='black', alpha=0.7, color='skyblue')
4 axes[0].set_title('Histogram of cumulative_weeks_in_top_10')
5 axes[0].set_xlabel('Cumulative Number of of Weeks in the Top 10')
6 axes[0].set_ylabel('Frequency')
7
8 axes[1].hist(top10['cumulative_weeks_in_top_10'], bins=16, range=(1, 16), edgecolor='black', alpha=0.7, color='salmon')
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')

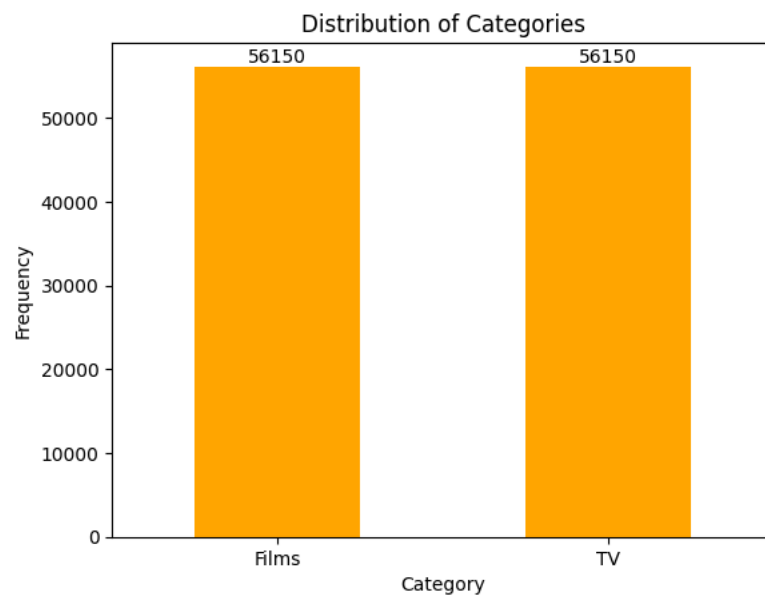


```
1 variances = top10.var()
2
3 variances.plot(kind='bar', color='lightgreen')
4 plt.title('Variance of Numeric Variables')
5 plt.xlabel('Variable')
6 plt.ylabel('Variance')
7 plt.xticks(rotation=0)
8 plt.grid(axis='y', linestyle='--', alpha=0.7)
9 plt.show()
```

```
<ipython-input-5-4c8b6dacff89>:1: FutureWarning: The default value of numeric_on  
variances = top10.var()
```



```
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()
```

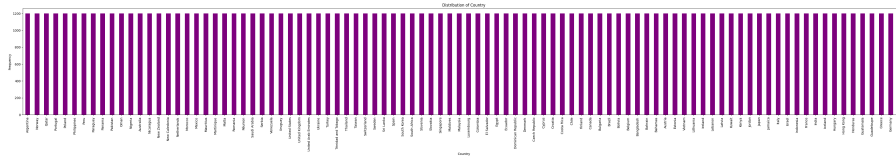


```

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(method='spearman')
2 print(spearman_corr)

```

```

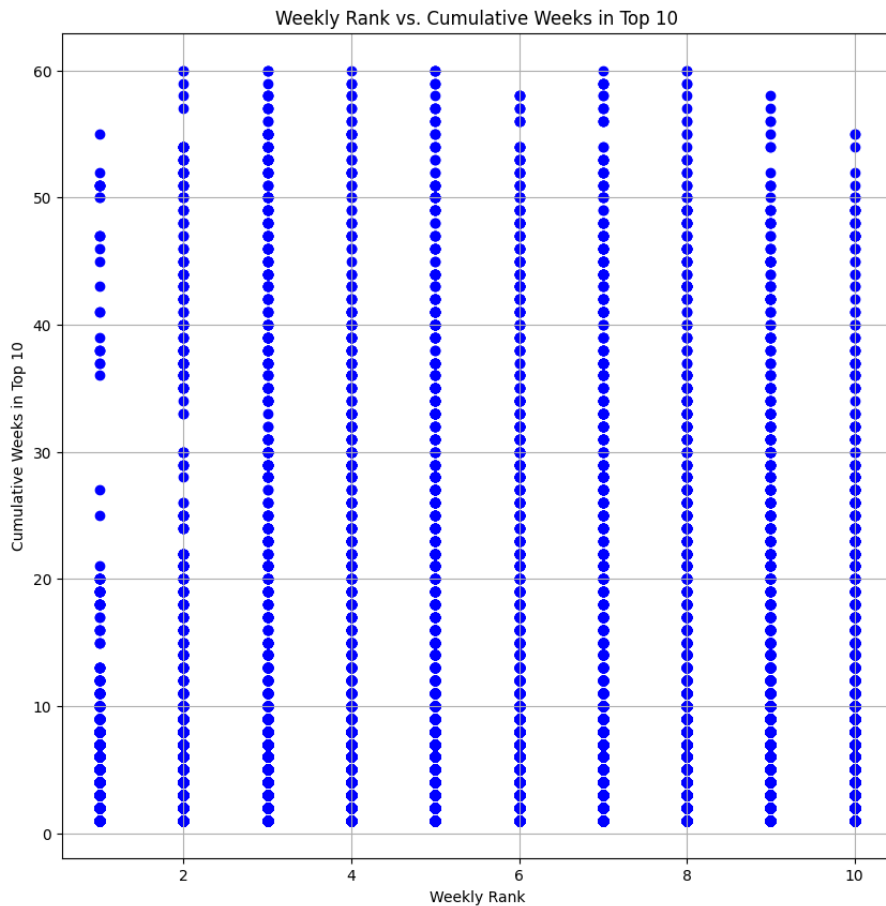
              weekly_rank  cumulative_weeks_in_top_10
weekly_rank              1.000000                0.028064
cumulative_weeks_in_top_10  0.028064                1.000000

```

```

1 plt.figure(figsize=(10, 10))
2 plt.scatter(top10['weekly_rank'], top10['cumulative_weeks_in_top_10'], color='blue')
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()

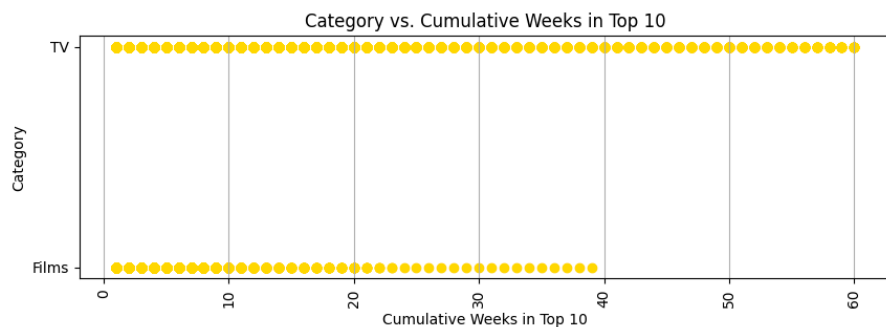
```



```

1 plt.figure(figsize=(10, 3))
2 plt.scatter(top10['cumulative_weeks_in_top_10'], top10['category'], color='gold')
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()

```



```

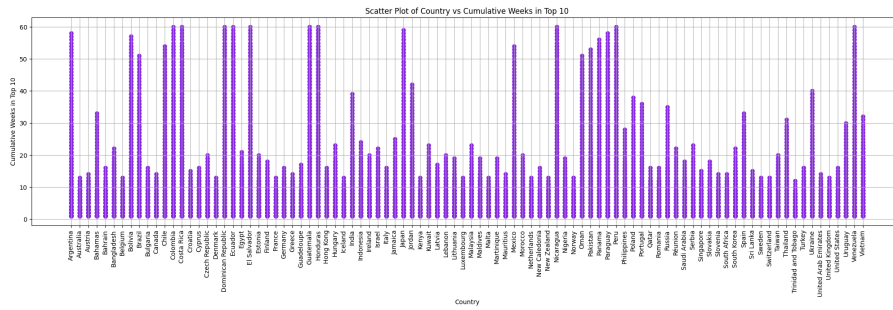
1 plt.figure(figsize=(25, 6))
2 plt.scatter(top10['country_name'], top10['cumulative_weeks_in_top_10'], color='blueviolet')
3 plt.xlabel('Country')

```

```

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_10'].mean()
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.values, color = "limegreen")
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 statistics to
16 #determine the significance of the relationship between month and cumulative weeks in top 10
17 from scipy.stats import kruskal
18
19 data_by_month = [group.values for name, group in top10.groupby('month')['cumulative_weeks_in_top_10']]
20 print(data_by_month)
21
22 kruskal_result = kruskal(*data_by_month)
23 print(kruskal_result)

```

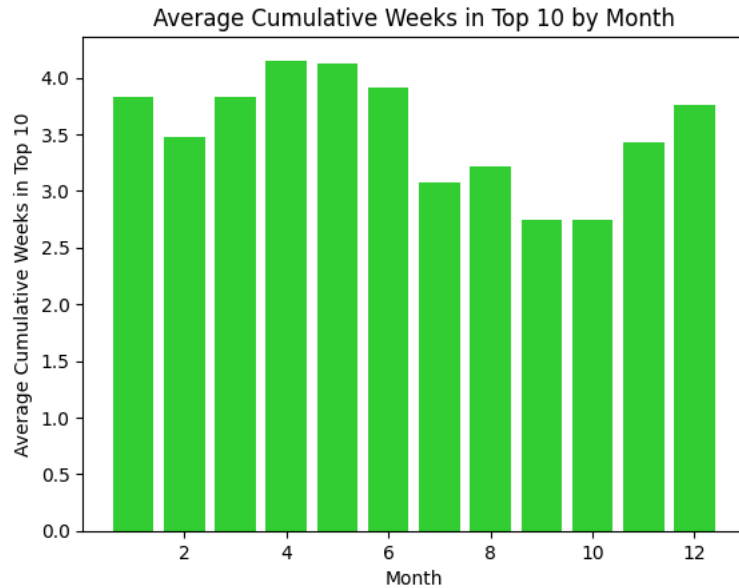
month

```

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



```

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

```

```

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 relationships and if signifcant,
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-statistic: 10.100335003904233
p-value: 3.560190557251464e-24

t-test between month 1 and month 12:
t-statistic: 0.9766721429985747
p-value: 0.328745455099556

t-test between month 1 and month 11:
t-statistic: 5.6946452915759895
p-value: 1.2568120167180884e-08

t-test between month 1 and month 10:
t-statistic: 17.878055935482745
p-value: 6.756436741720882e-71

t-test between month 1 and month 9:
t-statistic: 16.708487267559864
p-value: 3.5808423797454773e-62

t-test between month 12 and month 11:
t-statistic: 4.703586237878578
p-value: 2.5790393231842418e-06

t-test between month 12 and month 10:
t-statistic: 17.0504933145285
p-value: 1.2016725795091587e-64

t-test between month 12 and month 9:
t-statistic: 16.187596910486995
p-value: 1.9232420848985566e-58

t-test between month 11 and month 10:
t-statistic: 12.650184174461105
p-value: 1.6333181222166085e-36

t-test between month 11 and month 9:
t-statistic: 12.162578595087783
p-value: 7.093872323563874e-34

t-test between month 10 and month 9:
t-statistic: -0.11913295030528707
p-value: 0.9051714405437928

```

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

```

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 for the Unites States
5 us_show_freq = 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 occurring shows/movies that make the top 10 in the United States
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 for Canada
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 occurring shows/movies that make the top 10 in Canada
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 countries'''
2
3 top10_frequency_by_country = {}
4
5 #Compute frequency of each show in the top 10 for each country that is not the US
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().head(10)
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_by_country.keys())
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 current country
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, 0)
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]

```

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.sum()
14    expected_frequencies = expected_frequencies[:, contingency_table.columns.get_loc(country)]
15
16    chi2_stat, p_value, _, _ = chi2_contingency([observed_frequencies, expected_frequencies])
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 that makes
28 the Top 10 in the US is likely to make the Top 10 in the listed country. Possibly, the US market may
29 influence the markets in these other countries.'''
30
31 significant_countries = [country for country, p_value in p_values.items() if p_value < 0.01]
32 print("List of Significant Countries", significant_countries)

```

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
 Nigeria: 0.0019520571297548723
 Norway: 3.034455861014874e-05
 Oman: 0.005268137544786315
 Pakistan: 0.0003658332893763554
 Panama: 0.01326424107608904
 Paraguay: 0.015316351845605206
 Peru: 0.0206159800194768
 Philippines: 0.015008299565895349
 Poland: 0.03306371588803032
 Portugal: 0.0017040677828814136
 Qatar: 0.0005453804935978162
 Romania: 0.0037403045996279014
 Russia: 6.145670330520167e-16
 Réunion: 0.033139093941267496
 Saudi Arabia: 1.9679226758956863e-05
 Serbia: 3.818540781175382e-07
 Singapore: 0.029370495562954648
 Slovakia: 1.0532719195131897e-05
 Slovenia: 0.008692652331488748
 South Africa: 0.04623109664188473
 South Korea: 0.004584812871061249
 Spain: 0.005239982027083498
 Sri Lanka: 0.018355977793873757
 Sweden: 0.004804139909884401
 Switzerland: 0.05878084007392329
 Taiwan: nan
 Thailand: 0.039870176610579594
 Trinidad and Tobago: 0.0007869445378556762
 Turkey: 1.0088783514926642e-05
 Ukraine: 8.185745238206755e-16
 United Arab Emirates: 0.0008636492796893472
 United Kingdom: 4.8371130152199034e-06
 Uruguay: 0.004888806769799585
 Venezuela: 0.026938079134912037
 Vietnam: 0.012766005229262795
 United States: 1.919391971704617e-17

List of Significant Countries ['Bahamas', 'Bahrain', 'Bangladesh', 'Bolivia', 'Brazil', 'Canada', 'Chile', 'Colombia', 'Czech', 'Denmark', 'Ecuador', 'Egypt', 'France', 'Germany', 'Greece', 'Hong Kong', 'Hungary', 'India', 'Indonesia', 'Israel', 'Italy', 'Japan', 'Korea', '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', 'Uruguay', 'Venezuela', 'Vietnam', 'United States']
 /usr/local/lib/python3.10/dist-packages/scipy/stats/contingency.py:134: RuntimeWarning: invalid value encountered in divide
 expected = reduce(np.multiply, margsums) / observed.sum() ** (d - 1)

```

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 countries in the dataset
8 '''https://gradschool.utk.edu/future-students/office-of-graduate-admissions/applying-to-graduate-school/
9 admission-requirements/testing-requirements/countries-with-english-as-official-language/'''
10
11 english_speaking_countries = [
12     'Anguilla', 'Antigua and Barbuda', 'Bahamas', 'Barbados', 'Belize', 'Belgium', 'Bermuda', 'Botswana',
13     'British Virgin Islands', 'Burundi', 'Cameroon', 'Canada', 'Cayman Islands', 'Christmas Island', 'Cook Islands',
14     'Dominica', 'Fiji', 'Gambia', 'Ghana', 'Grenada', 'Guyana', 'Hong Kong', 'India', 'Ireland', 'Jersey', 'Kenya',
15     'Liberia', 'Malawi', 'Malta', 'Marshall Islands', 'Micronesia', 'Namibia', 'New Zealand',
16     'Nigeria', 'Niue', 'Norfolk Island', 'Northern Mariana Islands', 'Pakistan', 'Palau', 'Papua New Guinea',
17     'Philippines', 'Pitcairn Islands', 'Rwanda', 'Saint Kitts and Nevis', 'Saint Lucia', 'Samoa', 'Seychelles',
18     'Sierra Leone', 'Singapore', 'Sint Maarten', 'Solomon Islands', 'Somalia', 'South Africa', 'South Sudan', 'Sudan',
19     'Swaziland', 'Tanzania', 'Tonga', 'Trinidad and Tobago', 'Turks and Caicos Islands', 'Tuvalu', 'Uganda', 'Zambia',
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_speaking_countries + 1
28
29 print("Number of significant english speaking countries: ", num_significant_english_speaking_countries)

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 10 in Canada?'''
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 in Canada
13 target_variable = (us_shows['show_title'].isin(top10[top10['country_name'] == 'Canada']['show_title'])).astype(int)
14
15 # Combine features and target variable into a dataset
16 prepared_dataset = pd.concat([us_shows[['show_title', 'category', 'weekly_rank', 'country_name']], target_variable], axis=1)
17 prepared_dataset.columns = ['show_title', 'category', 'weekly_rank', 'country_name', 'is_top_10_canada']
18
19 print(prepared_dataset)
20
21 prepared_dataset['category'] = prepared_dataset['category'].map({'Films': 0, 'TV': 1})
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, random_state=42)
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 10", "In Canada Top 10"], filled=True)
49 plt.show()
50
```

```

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:

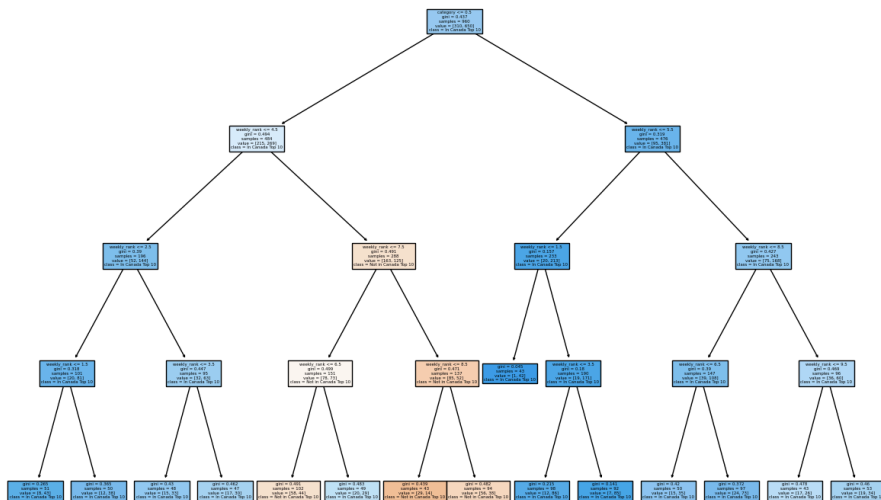
```

[[ 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 10 in Canada?
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 in Canada
13 target_variable = (us_shows['show_title'].isin(top10[top10['country_name'] == 'Canada']['show_title'])).astype(int)
14
15 # Combine features and target variable into a dataset
16 prepared_dataset = pd.concat([us_shows[['show_title', 'category', 'weekly_rank', 'country_name']], target_variable], axis=1)
17 prepared_dataset.columns = ['show_title', 'category', 'weekly_rank', 'country_name', 'is_top_10_canada']
18
19 prepared_dataset['category'] = prepared_dataset['category'].map({'Films': 0, 'TV': 1})
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, random_state=0)
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 10", "In Canada Top 10"], filled=True)
53 plt.show()

```

Cross-validation Scores: [0.72916667 0.70833333 0.68229167 0.734375 0.671875

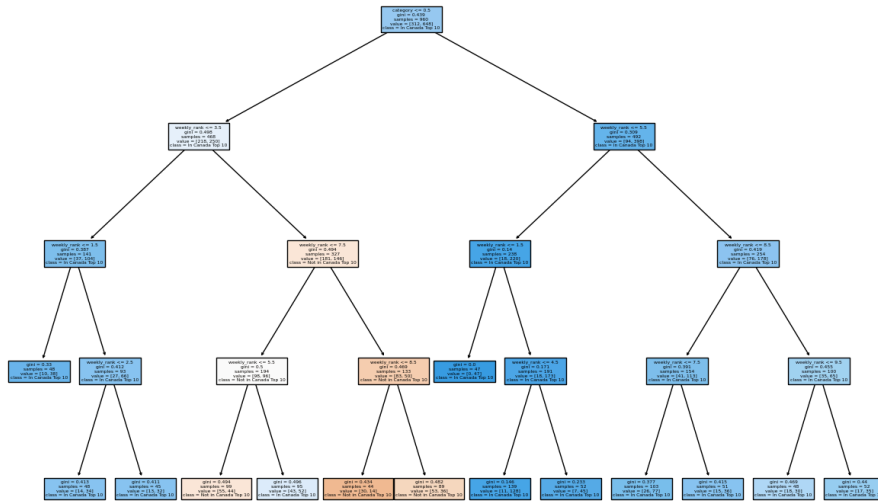
Mean CV Score: 0.7052083333333333

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 10 in Argentina?'''
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 in Argentina
13 target_variable = (us_shows['show_title'].isin(top10[top10['country_name'] == 'Argentina']['show_title'])).astype(int)
14
15 # Combine features and target variable into a dataset
16 prepared_dataset = pd.concat([us_shows[['show_title', 'category', 'weekly_rank', 'country_name']], target_variable], axis=1)
17 prepared_dataset.columns = ['show_title', 'category', 'weekly_rank', 'country_name', 'is_top_10_Argentina']
18
19 print(prepared_dataset)
20
21 prepared_dataset[['category']] = prepared_dataset[['category']].map({'TV-14': 0, 'TV-13': 1})

```



```
21 prepared_dataset['category'] = prepared_dataset['category'].map(lambda x: 0 if x == 'Not in Argentina Top 10' else 1)
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, random_state=42)
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 Top 10", "In Argentina Top 10"], filled=True)
49 plt.show()
50
```

```

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:

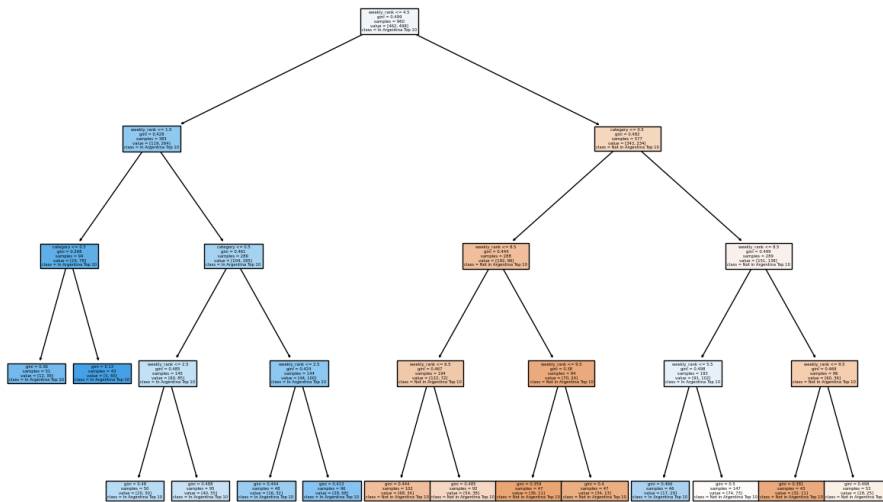
```

[[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))

```