

Netflix Data Cleaning, Analysis and Visualisation

November 13, 2024

1 Project 2: Netflix Data Cleaning, Anaysis and Visualization

1.0.1 Step 1: Importing necessary Libraries

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud
import warnings
warnings.filterwarnings('ignore', category=FutureWarning)
```

1.0.2 Step 2: Load the dataset

Loading the dataset into the dataframe and inspecting the initial structure of the data

```
[2]: ndata = pd.read_csv('netflix1.csv')
```

```
[3]: ndata.head(5)
```

```
[3]: show_id    type                title    director \
0      s1      Movie      Dick Johnson Is Dead  Kirsten Johnson
1      s3  TV Show      Ganglands      Julien Leclercq
2      s6  TV Show      Midnight Mass      Mike Flanagan
3     s14      Movie  Confessions of an Invisible Girl  Bruno Garotti
4      s8      Movie      Sankofa      Haile Gerima
```

```
country date_added release_year rating duration \
0 United States  9/25/2021      2020  PG-13    90 min
1      France  9/24/2021      2021  TV-MA    1 Season
2 United States  9/24/2021      2021  TV-MA    1 Season
3      Brazil  9/22/2021      2021  TV-PG    91 min
4 United States  9/24/2021      1993  TV-MA   125 min
```

```
listed_in
0 Documentaries
1 Crime TV Shows, International TV Shows, TV Act...
2 TV Dramas, TV Horror, TV Mysteries
3 Children & Family Movies, Comedies
```

4 Dramas, Independent Movies, International Movies

```
[4]: ndata.columns
```

```
[4]: Index(['show_id', 'type', 'title', 'director', 'country', 'date_added',
         'release_year', 'rating', 'duration', 'listed_in'],
         dtype='object')
```

```
[5]: ndata.shape
```

```
[5]: (8790, 10)
```

```
[6]: ndata.describe
```

```
[6]: <bound method NDFrame.describe of
title      director \
0         s1      Movie      Dick Johnson Is Dead  Kirsten Johnson
1         s3      TV Show      Ganglands          Julien Leclercq
2         s6      TV Show      Midnight Mass       Mike Flanagan
3        s14      Movie  Confessions of an Invisible Girl  Bruno Garotti
4         s8      Movie      Sankofa              Haile Gerima
...
8785    s8797      TV Show      Yunus Emre          Not Given
8786    s8798      TV Show      Zak Storm          Not Given
8787    s8801      TV Show      Zindagi Gulzar Hai    Not Given
8788    s8784      TV Show      Yoko              Not Given
8789    s8786      TV Show      YOM              Not Given
```

```
country  date_added  release_year  rating  duration \
0  United States  9/25/2021      2020  PG-13      90 min
1         France  9/24/2021      2021  TV-MA      1 Season
2  United States  9/24/2021      2021  TV-MA      1 Season
3         Brazil  9/22/2021      2021  TV-PG      91 min
4  United States  9/24/2021      1993  TV-MA     125 min
...
8785         Turkey  1/17/2017      2016  TV-PG      2 Seasons
8786  United States  9/13/2018      2016  TV-Y7      3 Seasons
8787         Pakistan  12/15/2016      2012  TV-PG      1 Season
8788         Pakistan  6/23/2018      2016   TV-Y      1 Season
8789         Pakistan  6/7/2018      2016  TV-Y7      1 Season
```

```
listed_in
0      Documentaries
1  Crime TV Shows, International TV Shows, TV Act...
2      TV Dramas, TV Horror, TV Mysteries
3      Children & Family Movies, Comedies
4      Dramas, Independent Movies, International Movies
```

```

...
8785          International TV Shows, TV Dramas
8786          Kids' TV
8787  International TV Shows, Romantic TV Shows, TV ...
8788          Kids' TV
8789          Kids' TV

```

```
[8790 rows x 10 columns]>
```

```
[7]: ndata.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8790 entries, 0 to 8789
Data columns (total 10 columns):
#   Column          Non-Null Count  Dtype
---  -
0   show_id         8790 non-null   object
1   type            8790 non-null   object
2   title           8790 non-null   object
3   director        8790 non-null   object
4   country         8790 non-null   object
5   date_added      8790 non-null   object
6   release_year    8790 non-null   int64
7   rating          8790 non-null   object
8   duration        8790 non-null   object
9   listed_in      8790 non-null   object
dtypes: int64(1), object(9)
memory usage: 686.8+ KB

```

1.0.3 Step 3: Data Cleaning

3.1 Identify columns with null values and treat them accordingly. Dropping rows for the dataset where critical information is missing.

```
[8]: # Check for missing values
ndata.isnull().sum()
```

```

[8]: show_id      0
     type        0
     title       0
     director    0
     country     0
     date_added  0
     release_year 0
     rating      0
     duration    0
     listed_in   0
     dtype: int64

```

No missing values in the data.

```
[9]: # Dropping rows with missing critical information in specific columns
ndata.dropna(subset=['director', 'country'], inplace=True)
```

3.2 Removing and Looking for Duplicates

Looking for duplicates in the column where the value needs to be unique.

```
[10]: # Remove duplicates based on 'title', 'director', and 'show_id'
ndata.drop_duplicates(subset=['title', 'director', 'show_id'], inplace=True)
```

```
[11]: # Verify that duplicates have been removed
print("Number of duplicates after removal:", ndata.duplicated(subset=['title', 'director', 'show_id']).sum())
```

Number of duplicates after removal: 0

3.3 Dropping Columns that are not required for analysis.

```
[12]: # Dropping 'show_id' and 'description' columns
ndata.drop(columns=['show_id', 'description'], inplace=True, errors='ignore')
```

```
[13]: print("Remaining columns:\n", ndata.columns)
```

Remaining columns:

```
Index(['type', 'title', 'director', 'country', 'date_added', 'release_year',
       'rating', 'duration', 'listed_in'],
      dtype='object')
```

```
[14]: ndata.head(2)
```

```
[14]:
```

	type	title	director	country	date_added	\
0	Movie	Dick Johnson Is Dead	Kirsten Johnson	United States	9/25/2021	
1	TV Show	Ganglands	Julien Leclercq	France	9/24/2021	

	release_year	rating	duration	\
0	2020	PG-13	90 min	
1	2021	TV-MA	1 Season	

	listed_in
0	Documentaries
1	Crime TV Shows, International TV Shows, TV Act...

3.4 Data Type Conversion

```
[15]: # Converting 'date_added' to datetime datatype.
ndata['date_added'] = pd.to_datetime(ndata['date_added'], errors='coerce')
```

```
[16]: ndata['date_added'].dtype
```

```
[16]: dtype('<M8[ns]')
```

3.5 Splitting column

```
[17]: # Splitting 'listed_in' column into individual genres
      ndata['genres'] = ndata['listed_in'].apply(lambda x: x.split(', '))
```

```
[18]: print(ndata[['listed_in', 'genres']].head())
```

```

               listed_in \
0               Documentaries
1  Crime TV Shows, International TV Shows, TV Act...
2               TV Dramas, TV Horror, TV Mysteries
3               Children & Family Movies, Comedies
4  Dramas, Independent Movies, International Movies

               genres
0               [Documentaries]
1  [Crime TV Shows, International TV Shows, TV Ac...
2               [TV Dramas, TV Horror, TV Mysteries]
3               [Children & Family Movies, Comedies]
4  [Dramas, Independent Movies, International Mov...
```

Exploding the data can affect certain analyses, such as: Counting the total number of unique titles (Movies and TV Shows), Analyzing the distribution of content types (Movies vs. TV Shows), Any analysis that requires the unique count of titles rather than genre-based analysis.

So to avoid this let's make a copy of the dataset ('ndata') and name it as ('data_exploded'). We will Keep the original dataset ('ndata') for queries that require unique titles and the exploded DataFrame ('data_exploded') only for genre-specific analysis.

```
[19]: # Creating a separate exploded DataFrame for genre analysis
      data_exploded = ndata.copy()
```

```
[20]: data_exploded = data_exploded.explode('genres')
```

```
[21]: data_exploded.head()
```

```
[21]:
```

	type	title	director	country	date_added	\
0	Movie	Dick Johnson Is Dead	Kirsten Johnson	United States	2021-09-25	
1	TV Show	Ganglands	Julien Leclercq	France	2021-09-24	
1	TV Show	Ganglands	Julien Leclercq	France	2021-09-24	
1	TV Show	Ganglands	Julien Leclercq	France	2021-09-24	
2	TV Show	Midnight Mass	Mike Flanagan	United States	2021-09-24	

	release_year	rating	duration	\
0	2020	PG-13	90 min	
1	2021	TV-MA	1 Season	
1	2021	TV-MA	1 Season	

```

1          2021  TV-MA  1 Season
2          2021  TV-MA  1 Season

```

```

                                listed_in          genres
0                                Documentaries    Documentaries
1  Crime TV Shows, International TV Shows, TV Act...    Crime TV Shows
1  Crime TV Shows, International TV Shows, TV Act...  International TV Shows
1  Crime TV Shows, International TV Shows, TV Act...    TV Action & Adventure
2                                TV Dramas, TV Horror, TV Mysteries    TV Dramas

```

```
[19]: ndata.head()
```

```

[19]:      type      title      director      country \
0  Movie      Dick Johnson Is Dead  Kirsten Johnson  United States
1  TV Show      Ganglands  Julien Leclercq      France
2  TV Show      Midnight Mass  Mike Flanagan  United States
3  Movie  Confessions of an Invisible Girl  Bruno Garotti      Brazil
4  Movie      Sankofa  Haile Gerima  United States

```

```

      date_added  release_year  rating  duration \
0  2021-09-25      2020  PG-13    90 min
1  2021-09-24      2021  TV-MA    1 Season
2  2021-09-24      2021  TV-MA    1 Season
3  2021-09-22      2021  TV-PG    91 min
4  2021-09-24      1993  TV-MA    125 min

```

```

                                listed_in \
0                                Documentaries
1  Crime TV Shows, International TV Shows, TV Act...
2                                TV Dramas, TV Horror, TV Mysteries
3                                Children & Family Movies, Comedies
4  Dramas, Independent Movies, International Movies

```

```

                                genres
0                                [Documentaries]
1  [Crime TV Shows, International TV Shows, TV Ac...
2                                [TV Dramas, TV Horror, TV Mysteries]
3                                [Children & Family Movies, Comedies]
4  [Dramas, Independent Movies, International Mov...

```

1.0.4 Step 4: Exploratory Data Analysis (EDA)

4.1 Content Type Distribution (Movies vs. TV Shows)

```
[20]: type_counts = ndata['type'].value_counts()
```

```
[106]: fig, ax = plt.subplots(1, 2, figsize=(16, 6))
```

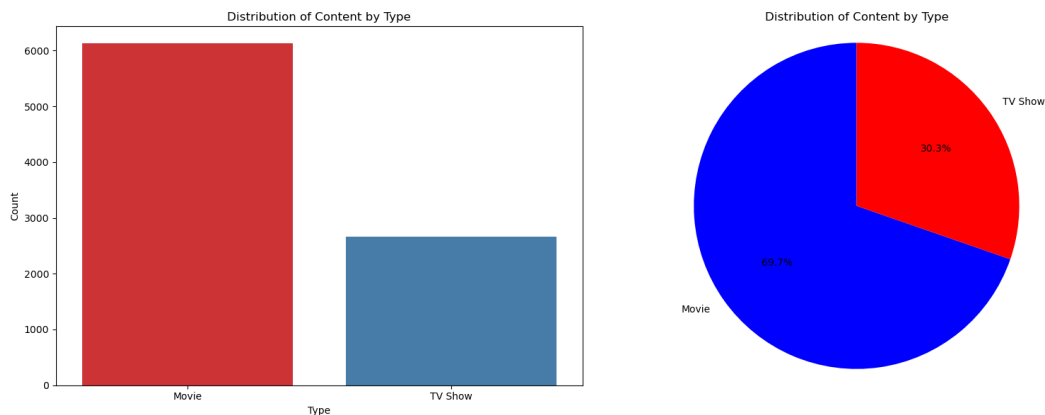
```

# First subplot: Bar plot
sns.barplot(x=type_counts.index, y=type_counts.values, palette='Set1', ax=ax[0])
ax[0].set_title('Distribution of Content by Type')
ax[0].set_xlabel('Type')
ax[0].set_ylabel('Count')

# Second subplot: Pie chart
ax[1].pie(type_counts, labels=type_counts.index, autopct='%0.1f%%',
          ↪startangle=90, colors=['blue', 'red'])
ax[1].set_title('Distribution of Content by Type')
ax[1].axis('equal')

plt.tight_layout()
plt.show()

```



4.2 Most Common Genres (we'll use "data_exploded")

```

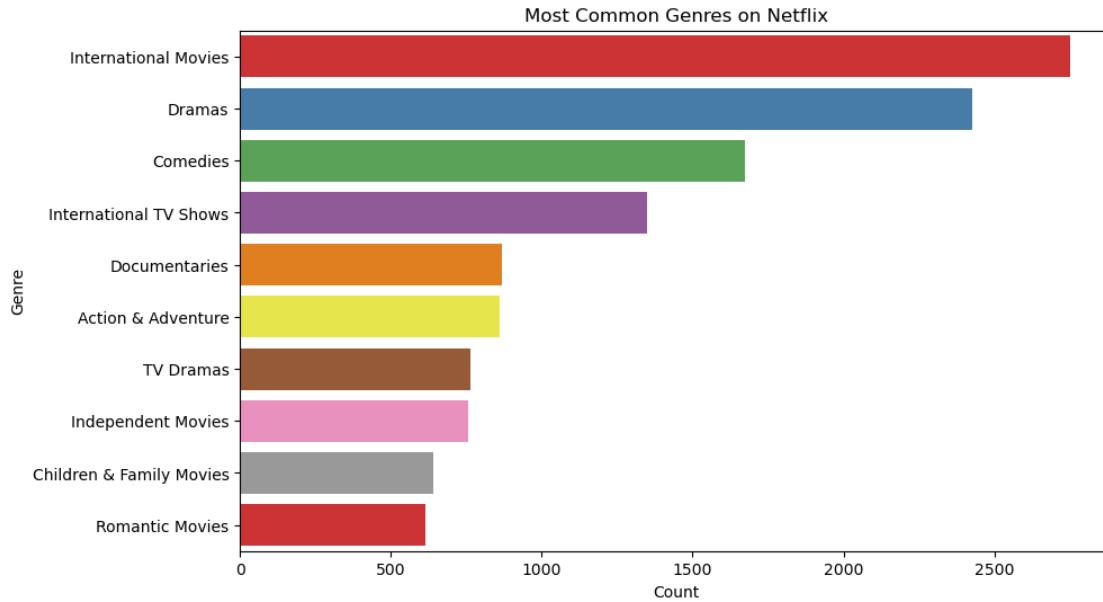
[46]: all_genres = sum(ndata['genres'], [])
      genre_counts = pd.Series(all_genres).value_counts().head(10)

```

```

[47]: plt.figure(figsize=(10, 6))
      sns.barplot(x=genre_counts.index, y=genre_counts.values, palette='Set1')
      plt.title('Most Common Genres on Netflix')
      plt.xlabel('Genre')
      plt.ylabel('Count')
      plt.show()

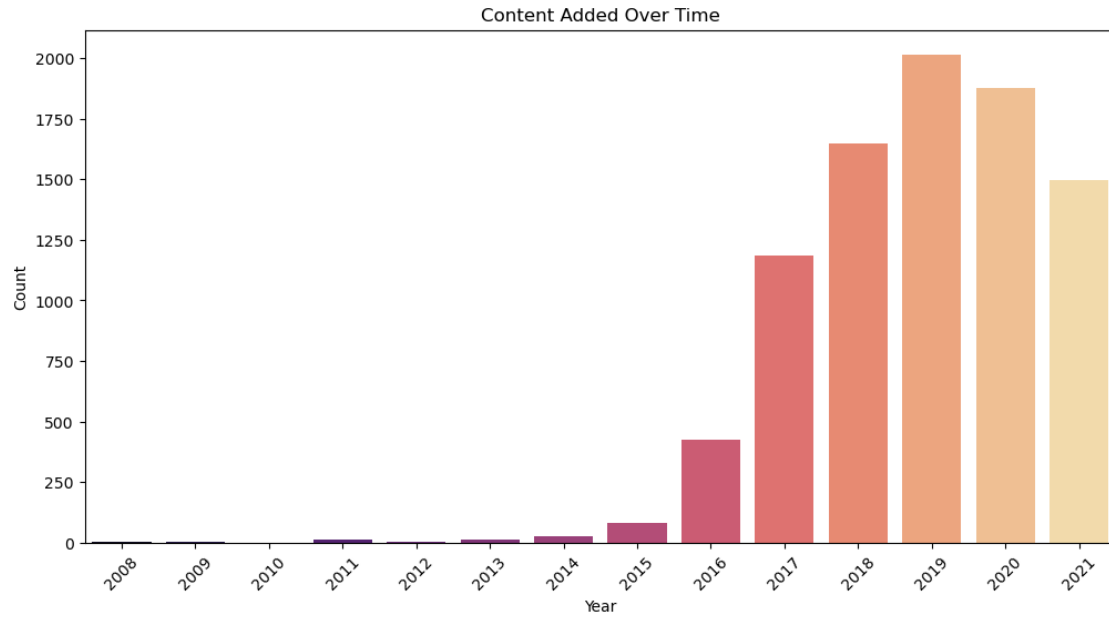
```



4.3 Content Added Over Time

```
[48]: ndata['year_added'] = ndata['date_added'].dt.year
```

```
[49]: plt.figure(figsize=(12, 6))
sns.countplot(x='year_added', data=ndata, palette='magma')
plt.title('Content Added Over Time')
plt.xlabel('Year')
plt.ylabel('Count')
plt.xticks(rotation=45)
plt.show()
```

4.4 Top 10 Directors with most Titles

```
[36]: top_directors = ndata['director'].value_counts().head(10)
```

```
[37]: plt.figure(figsize=(10, 6))
sns.barplot(x=top_directors.values, y=top_directors.index, palette='magma')
plt.title('Top 10 Directors with the Most Titles')
plt.xlabel('Number of Titles')
plt.ylabel('Director')
plt.show()
```


4.6 Ratings

```
[55]: ndata['rating'].value_counts()
```

```
[55]: rating
      TV-MA      3205
      TV-14      2157
      TV-PG       861
      R          799
      PG-13       490
      TV-Y7       333
      TV-Y        306
      PG          287
      TV-G        220
      NR          79
      G           41
      TV-Y7-FV     6
      NC-17        3
      UR           3
      Name: count, dtype: int64
```

```
[67]: import matplotlib.cm as cm
```

```
[65]: ratings = ndata['rating'].value_counts().reset_index()
      ratings.columns = ['rating', 'count']
```

```
[77]: cmap = plt.colormaps.get_cmap('magma')
      colors = [cmap(i) for i in np.linspace(0, 1, len(ratings))]
```

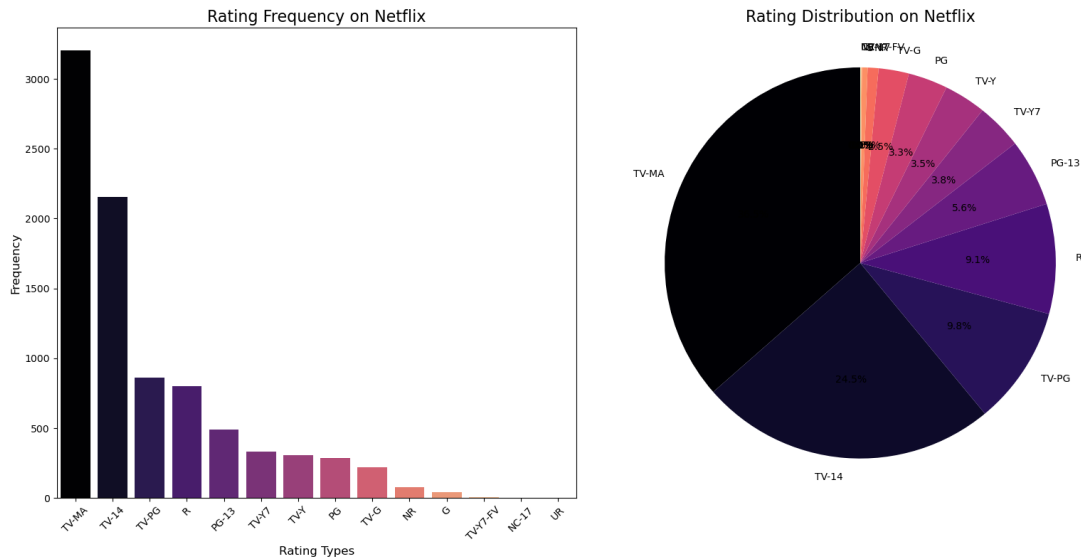
```
[78]: fig, axes = plt.subplots(1, 2, figsize=(16, 8))

      # Bar Chart (left side)
      sns.barplot(ax=axes[0], x=ratings['rating'], y=ratings['count'], palette=colors)
      axes[0].set_title('Rating Frequency on Netflix', fontsize=16)
      axes[0].set_xlabel('Rating Types', fontsize=12)
      axes[0].set_ylabel('Frequency', fontsize=12)
      axes[0].tick_params(axis='x', rotation=45)

      # Pie Chart (right side)
      axes[1].pie(ratings['count'], labels=ratings['rating'], autopct='%1f%%',
                  ↪startangle=90, colors=colors)
      axes[1].set_title('Rating Distribution on Netflix', fontsize=16)
      axes[1].axis('equal') # Ensures the pie chart is a circle

      plt.tight_layout()
```

```
plt.show()
```



4.7 Top 10 Movies and TV shows producing countries

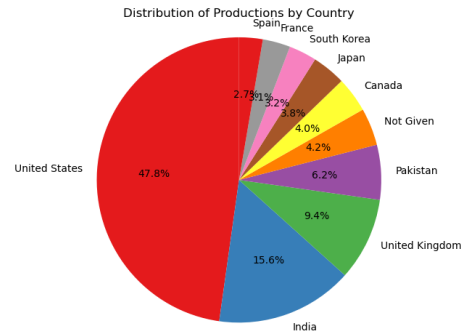
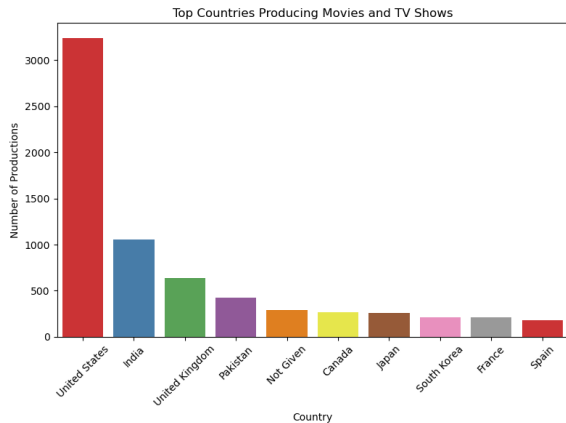
```
[108]: country_counts = ndata['country'].value_counts().head(10)
```

```
[109]: # Create a figure with two subplots (1 row, 2 columns)
fig, ax = plt.subplots(1, 2, figsize=(16, 6))

# First subplot: Bar plot for top countries
sns.barplot(x=country_counts.index, y=country_counts.values, palette='Set1',
            ax=ax[0])
ax[0].set_title('Top Countries Producing Movies and TV Shows')
ax[0].set_xlabel('Country')
ax[0].set_ylabel('Number of Productions')
ax[0].tick_params(axis='x', rotation=45)

# Second subplot: Pie chart for top countries
ax[1].pie(country_counts, labels=country_counts.index, autopct='%0.1f%%',
           startangle=90, colors=sns.color_palette('Set1',
           n_colors=len(country_counts)))
ax[1].set_title('Distribution of Productions by Country')
ax[1].axis('equal')

# Display the plots
plt.tight_layout()
plt.show()
```



4.8 Montly Release

```
[110]: # Extracting month from the date_added
ndata['month'] = ndata['date_added'].dt.month_name()
```

```
[112]: # Count the number of releases per month
monthly_counts = ndata['month'].value_counts().reindex(
    ['January', 'February', 'March', 'April', 'May', 'June', 'July', 'August',
    ↪ 'September', 'October', 'November', 'December']
).fillna(0) # Ensures months with no data are shown as 0
```

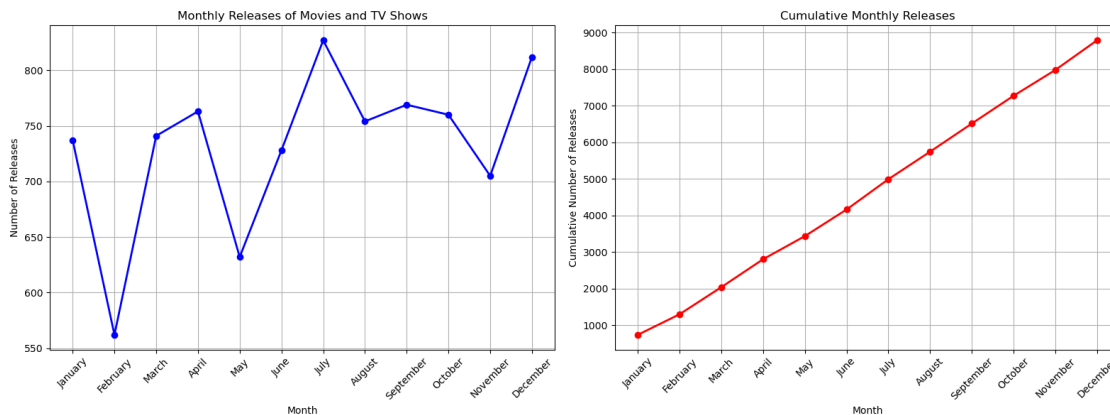
```
[113]: fig, ax = plt.subplots(1, 2, figsize=(16, 6))

# First subplot: Line plot for monthly releases
ax[0].plot(monthly_counts.index, monthly_counts.values, marker='o', color='b',
    ↪ linestyle='-', linewidth=2, markersize=6)
ax[0].set_title('Monthly Releases of Movies and TV Shows')
ax[0].set_xlabel('Month')
ax[0].set_ylabel('Number of Releases')
ax[0].tick_params(axis='x', rotation=45)
ax[0].grid(True)

# Second subplot: Line plot for the cumulative distribution of releases
ax[1].plot(monthly_counts.index, monthly_counts.cumsum(), marker='o',
    ↪ color='r', linestyle='-', linewidth=2, markersize=6)
ax[1].set_title('Cumulative Monthly Releases')
ax[1].set_xlabel('Month')
ax[1].set_ylabel('Cumulative Number of Releases')
ax[1].tick_params(axis='x', rotation=45)
ax[1].grid(True)

plt.tight_layout()
```

```
plt.show()
```



1.0.5 Step 5: Feature Engineering - Creating new features for better analysis.

5.1 Extracting the Duration in Minutes

The 'duration' column contains either the duration of a movie in minutes (e.g., "90 min") or the number of seasons for a TV show (e.g., "1 Season"). We will extract the duration in minutes for movies and convert it to an integer.

```
[79]: import re

# Creating a new column 'duration_minutes' with extracted duration for movies
ndata['duration_minutes'] = ndata['duration'].apply(lambda x: int(re.
    ↪findall(r'\d+', x)[0]) if 'min' in x else None)

# Display the new column
print(ndata[['title', 'type', 'duration', 'duration_minutes']].head())
```

	title	type	duration	duration_minutes
0	Dick Johnson Is Dead	Movie	90 min	90.0
1	Ganglands	TV Show	1 Season	NaN
2	Midnight Mass	TV Show	1 Season	NaN
3	Confessions of an Invisible Girl	Movie	91 min	91.0
4	Sankofa	Movie	125 min	125.0

5.2 Extracting Year, Month, and Day

```
[82]: ndata['year_added'] = ndata['date_added'].dt.year
ndata['month_added'] = ndata['date_added'].dt.month
ndata['day_added'] = ndata['date_added'].dt.day

print(ndata[['date_added', 'year_added', 'month_added', 'day_added']].head())
```

	date_added	year_added	month_added	day_added
0	2021-09-25	2021	9	25
1	2021-09-24	2021	9	24
2	2021-09-24	2021	9	24
3	2021-09-22	2021	9	22
4	2021-09-24	2021	9	24

5.3 Creating Binary Columns for Content Type

```
[83]: ndata['is_movie'] = ndata['type'].apply(lambda x: 1 if x == 'Movie' else 0)
      ndata['is_tv_show'] = ndata['type'].apply(lambda x: 1 if x == 'TV Show' else 0)
      print(ndata[['type', 'is_movie', 'is_tv_show']].head())
```

	type	is_movie	is_tv_show
0	Movie	1	0
1	TV Show	0	1
2	TV Show	0	1
3	Movie	1	0
4	Movie	1	0

```
[84]: ndata.shape
```

```
[84]: (8790, 16)
```

5.4 Extracting the Number of Seasons for TV Shows

```
[86]: ndata['num_seasons'] = ndata.apply(lambda x: int(re.findall(r'\d+', x[
      ↪x['duration'])[0]) if 'Season' in x['duration'] else 0, axis=1)
      print(ndata[['title', 'type', 'duration', 'num_seasons']].head())
```

	title	type	duration	num_seasons
0	Dick Johnson Is Dead	Movie	90 min	0
1	Ganglands	TV Show	1 Season	1
2	Midnight Mass	TV Show	1 Season	1
3	Confessions of an Invisible Girl	Movie	91 min	0
4	Sankofa	Movie	125 min	0

```
[87]: ndata.shape
```

```
[87]: (8790, 17)
```

5.5 Creating a Popularity Metric Using Ratings

Creating a popularity score based on the rating category (e.g., TV-MA, PG-13). Higher maturity ratings might correlate with more popular or trending content.

```
[88]: # Defining a dictionary for rating scores
      rating_scores = {
          'TV-MA': 5, 'R': 5, 'NC-17': 4,
          'TV-14': 4, 'PG-13': 3, 'PG': 2,
          'TV-PG': 2, 'TV-Y7': 1, 'TV-Y': 1,
```

```
'G': 1, 'NR': 3, 'UR': 3, 'TV-G': 2
}
```

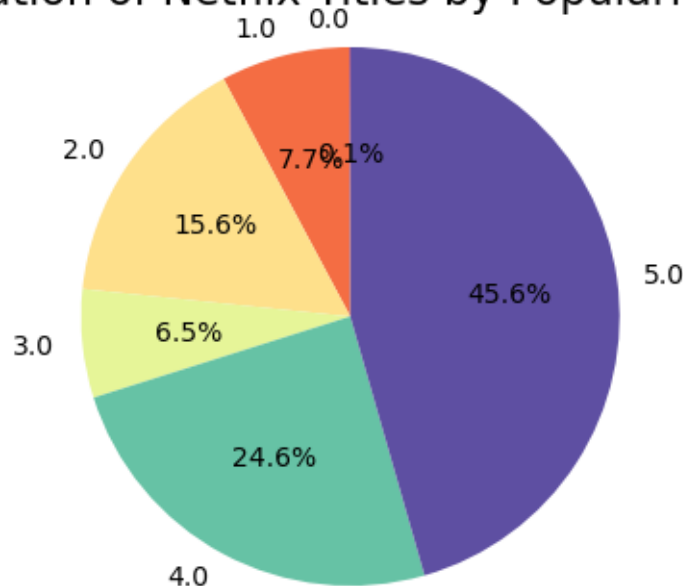
```
[89]: ndata['popularity_score'] = ndata['rating'].map(rating_scores).fillna(0)
print(ndata[['title', 'rating', 'popularity_score']].head())
```

	title	rating	popularity_score
0	Dick Johnson Is Dead	PG-13	3.0
1	Ganglands	TV-MA	5.0
2	Midnight Mass	TV-MA	5.0
3	Confessions of an Invisible Girl	TV-PG	2.0
4	Sankofa	TV-MA	5.0

```
[92]: popularity_distribution = ndata['popularity_score'].value_counts().sort_index()
cmap = plt.colormaps.get_cmap('Spectral')
colors = [cmap(i) for i in np.linspace(0, 1, len(popularity_distribution))]

# Create the pie chart
plt.figure(figsize=(6, 4))
plt.pie(popularity_distribution, labels=popularity_distribution.index,
        autopct='%.1f%%', startangle=90, colors=colors)
plt.title('Distribution of Netflix Titles by Popularity Score', fontsize=16)
plt.axis('equal') # Ensures the pie chart is a circle
plt.show()
```

Distribution of Netflix Titles by Popularity Score



The pie chart will display the percentage distribution of Netflix titles based on their popularity_score

5.6 One-Hot Encode the rating Column

We can one-hot encode the rating column to prepare it for machine learning part.

```
[93]: ndata = pd.get_dummies(ndata, columns=['rating'], prefix='rating')
```

```
[94]: ndata.columns
```

```
[94]: Index(['type', 'title', 'director', 'country', 'date_added', 'release_year',  
        'duration', 'listed_in', 'genres', 'year_added', 'duration_minutes',  
        'month_added', 'day_added', 'is_movie', 'is_tv_show', 'num_seasons',  
        'popularity_score', 'rating_G', 'rating_NC-17', 'rating_NR',  
        'rating_PG', 'rating_PG-13', 'rating_R', 'rating_TV-14', 'rating_TV-G',  
        'rating_TV-MA', 'rating_TV-PG', 'rating_TV-Y', 'rating_TV-Y7',  
        'rating_TV-Y7-FV', 'rating_UR'],  
        dtype='object')
```

```
[95]: ndata.shape
```

```
[95]: (8790, 31)
```

1.0.6 Step 6: Machine Learning - Building a Simple Recommendation Model

For this project, I'm creating a basic content-based recommendation model using the genres information.

6.1 Vectorize the Genres and Combine with Additional Features

Let's use TF-IDF Vectorization for genres and concatenate it with numerical features.

```
[96]: from sklearn.feature_extraction.text import TfidfVectorizer  
      from sklearn.preprocessing import MinMaxScaler  
      from sklearn.metrics.pairwise import cosine_similarity
```

```
[97]: # Initialize the TF-IDF Vectorizer for genres  
      tfidf = TfidfVectorizer(stop_words='english')  
      tfidf_matrix = tfidf.fit_transform(ndata['listed_in'])
```

```
[100]: # Normalize numerical features  
       scaler = MinMaxScaler()  
       numerical_features = ndata[['popularity_score', 'is_movie', 'is_tv_show']]  
       normalized_features = scaler.fit_transform(numerical_features)
```

```
[101]: # Concatenate TF-IDF matrix with numerical features  
       import scipy.sparse as sp  
       feature_matrix = sp.hstack([tfidf_matrix, normalized_features])
```

```
[102]: print("Feature Matrix Shape:", feature_matrix.shape)
```

Feature Matrix Shape: (8790, 47)

6.2 Computing Cosine Similarity and Build the Recommendation Function

Let's use cosine similarity to find similar titles based on the combined feature set.

```
[ ]: # Compute cosine similarity matrix
cosine_sim = cosine_similarity(feature_matrix, feature_matrix)

# Function to get recommendations
def get_enhanced_recommendations(title, data=ndata, cosine_sim=cosine_sim):
    # Get the index of the title
    idx = data[data['title'] == title].index[0]

    # Get similarity scores for this title
    sim_scores = list(enumerate(cosine_sim[idx]))

    # Sort titles based on similarity scores
    sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)

    # Get the indices of the top 5 most similar titles
    sim_indices = [i[0] for i in sim_scores[1:6]]

    # Return the top 5 similar titles
    return data['title'].iloc[sim_indices]

# Example usage
print("Enhanced Recommendations for 'Stranger Things':\n",
      get_enhanced_recommendations('Stranger Things'))
```

```
[103]: # Compute cosine similarity matrix
cosine_sim = cosine_similarity(feature_matrix, feature_matrix)
```

```
[104]: # Function to get recommendations
def get_enhanced_recommendations(title, data=ndata, cosine_sim=cosine_sim):
    # Geting the index of the title
    idx = data[data['title'] == title].index[0]

    # Geting similarity scores for this title
    sim_scores = list(enumerate(cosine_sim[idx]))

    # Sorting titles based on similarity scores
    sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)

    # Geting the indices of the top 5 most similar titles
    sim_indices = [i[0] for i in sim_scores[1:6]]
```

```
# Returning the top 5 similar titles
return data['title'].iloc[sim_indices]
```

```
[105]: print("Enhanced Recommendations for 'Stranger Things':\n",  
        ↪get_enhanced_recommendations('Stranger Things'))
```

```
Enhanced Recommendations for 'Stranger Things':  
7633      Stranger Things  
7446      Nightflyers  
8485      Helix  
6613      Manifest  
8122      The Vampire Diaries  
Name: title, dtype: object
```

```
[115]: ndata.to_csv('E:/New folder/enhanced_netflix_data.csv', index=False)
```

1.0.7 With the steps outlined above, the project is now fully completed according to the requirements mentioned by Unified Mentor. Let's go through a quick summary to confirm that every part of the project has been addressed

2 Project Summary

2.0.1 Data Import and Initial Inspection:

1. Imported the dataset into Jupyter Notebook.
2. Inspected the initial structure and verified the column types.

2.0.2 Data Cleaning:

1. Checked for missing values (none found in this case).
2. Removed duplicates based on title, director, and show_id.
3. Dropped unnecessary columns (show_id and description).
4. Converted date_added to datetime format.
5. Split the listed_in column into individual genres and created an exploded DataFrame (data_exploded) for genre-specific analysis.

2.0.3 Exploratory Data Analysis (EDA):

1. Performed content type distribution (Movies vs. TV Shows).
2. Analyzed the most common genres using the exploded DataFrame.
3. Explored trends of content added over time.
4. Identified the top 10 directors with the most titles.
5. Created a word cloud of movie titles.
6. Visualized rating distribution using bar charts and pie charts.
7. Analyzed Top countries producing movies and TV shows
8. Visualized Monthly release of TV shows and movies

2.0.4 Feature Engineering:

1. Created new features, such as genres (count of genres per title) and duration_minutes (duration in minutes for movies).
2. Extracted year, month, and day from date_added.
3. Created binary columns for content type (is_movie, is_tv_show).
4. Extracted the number of seasons for TV shows.
5. Created a popularity metric based on ratings.
6. One-hot encoded the rating column for machine learning.

2.0.5 Machine Learning:

1. Built a basic content-based recommendation model using TF-IDF vectorization for genres.
2. Enhanced the recommendation model by including additional features (popularity score, binary columns for content type, and number of genres).
3. Provided a function to generate recommendations based on the enhanced model.

2.1 Conclusion:

This project provided valuable insights into the dataset through data cleaning, exploratory analysis, and feature engineering. The content-based recommendation model was successfully implemented, offering personalized suggestions based on genres and additional features. Moving forward, the model can be further optimized by incorporating additional data, testing with different algorithms, and evaluating performance through advanced metrics. This work serves as a foundation for building a more robust recommendation system that could be applied to larger datasets or deployed in real-time applications.