



Objective:

- Import and clean the Netflix data set using pandas in python.
- Analyze data trends and distributions using summary statistics and visualizations in python.
- Create charts like bar charts pie charts using Matplotlib and Seaborn.
- Extract key trends and patterns, such as top and most common genre.



Importing Libraries:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import plotly.express as px
import seaborn as sns
from wordcloud import WordCloud
```



Data Analysis:

```
data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8790 entries, 0 to 8789
Data columns (total 10 columns):
 # Column
                  Non-Null Count Dtype
    show id
                  8790 non-null
                  8790 non-null
                                 object
    type
    title
                  8790 non-null
                                  object
     director
                  8790 non-null
                                 object
     country
                  8790 non-null
                                 object
     date added
                  8790 non-null
                                 object
    release year 8790 non-null
                                  int64
     rating
                  8790 non-null
                                 object
     duration
                  8790 non-null
                                  object
                  8790 non-null object
    listed in
dtypes: int64(1), object(9)
memory usage: 686.8+ KB
```



#Prints first 5 rows of the dataset
data.head()

26]

listed_in	duration	rating	release_year	date_added	country	director	title	type	show_id	
Documentaries	90 min	PG-13	2020	9/25/2021	United States	Kirsten Johnson	Dick Johnson Is Dead	Movie	s1	0
Crime TV Shows, International TV Shows, TV Act	1 Season	TV-MA	2021	9/24/2021	France	Julien Leclercq	Ganglands	TV Show	s 3	
TV Dramas, TV Horror, TV Mysteries	1 Season	TV-MA	2021	9/24/2021	United States	Mike Flanagan	Midnight Mass	TV Show	s6	2
Children & Family Movies, Comedies	91 min	TV-PG	2021	9/22/2021	Brazil	Bruno Garotti	Confessions of an Invisible Girl	Movie	s14	3
Dramas, Independent Movies, International Movies	125 min	TV-MA	1993	9/24/2021	United States	Haile Gerima	Sankofa	Movie	s8	4

Data Analysis:

```
print("Top Ratings Counts:")
   data.rating.value_counts()
Top Ratings Counts:
rating
TV-MA
           3205
TV-14
           2157
TV-PG
            861
            799
PG-13
            490
TV-Y7
TV-Y
             306
            287
TV-G
             220
             41
TV-Y7-FV
NC-17
Name: count, dtype: int64
```

Data Cleaning:Adjusting Date Format

```
# Convert 'date_added' to datetime
    data['date_added'] = pd.to_datetime(data['date_added'], errors='coerce')

# Print the first few rows of the DataFrame to check the conversion
    print(data[['date_added']].head())

[28]

...    date_added
    0 2021-09-25
    1 2021-09-24
    2 2021-09-24
    3 2021-09-22
    4 2021-09-24
```

Data Cleaning:Checking for Null Values







Data Cleaning:

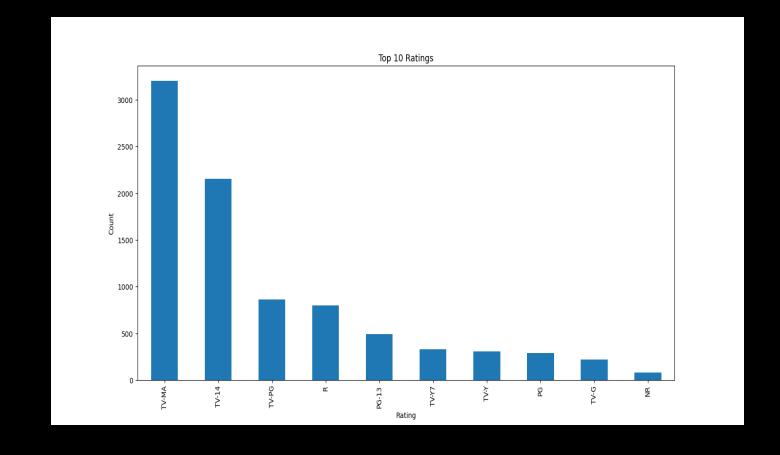
```
# Data Cleaning
   print("Missing Values Before Cleaning:")
   data.isnull().sum()
Missing Values Before Cleaning:
show id
type
title
director
country
date added
                4386
release year
rating
duration
listed in
dtype: int64
```

```
data.drop_duplicates(inplace=True)
   data['date added'] = data['date added'].fillna(pd.to datetime('today'))
   print("Missing Values After Cleaning:")
   data.isnull().sum()
Missing Values After Cleaning:
show id
               0
type
title
director
country
date added
release year
rating
duration
listed in
dtype: int64
```



Top 10 ratings

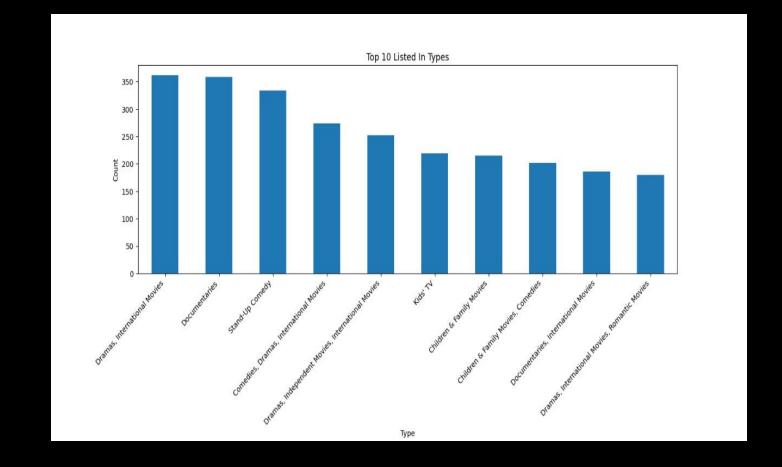
```
# Keep 'rating' as a categorical variable
data['rating'] = data['rating'].astype(str)
data['Years']=data['release_year']
year = data['rating'].value_counts()
year[:10].plot(kind='bar')
plt.title('Top 10 Ratings')
plt.xlabel('Rating')
plt.ylabel('Count')
plt.show()
```





Top 10 Listed in Genre Types

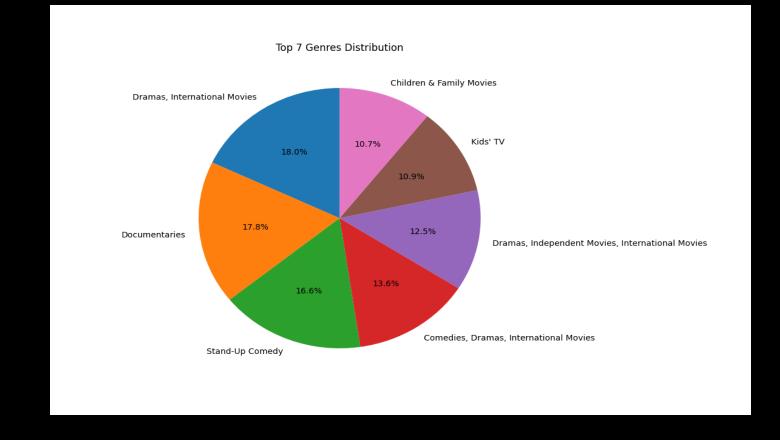
```
#Top 10 Genre Lists
top_10_types = data['listed_in'].value_counts()
top_10_types[:10].plot(kind='bar')
plt.title('Top 10 Listed In Genre Types')
plt.xlabel('Type of Genres')
plt.xticks(rotation=45, ha='right')
plt.ylabel('Count')
plt.show()
```





Top 7 Genres Distribution

```
#Top 7 Genre Distribution
genre_list = data['listed_in'].value_counts()
genre_list[:7].plot(kind='pie', autopct='%1.1f%%', startangle=90)
plt.title('Top 7 Genres Distribution')
plt.ylabel('')
plt.show()
```

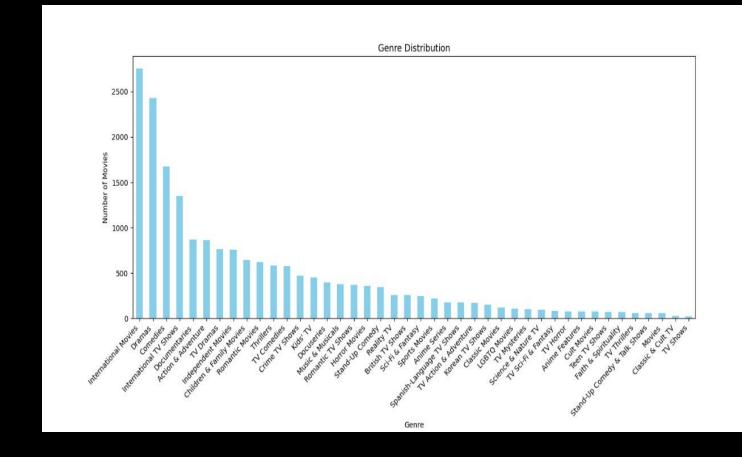




Genre Distribution

```
# Count the occurrences of each genre
genre_counts = genres_series.value_counts()

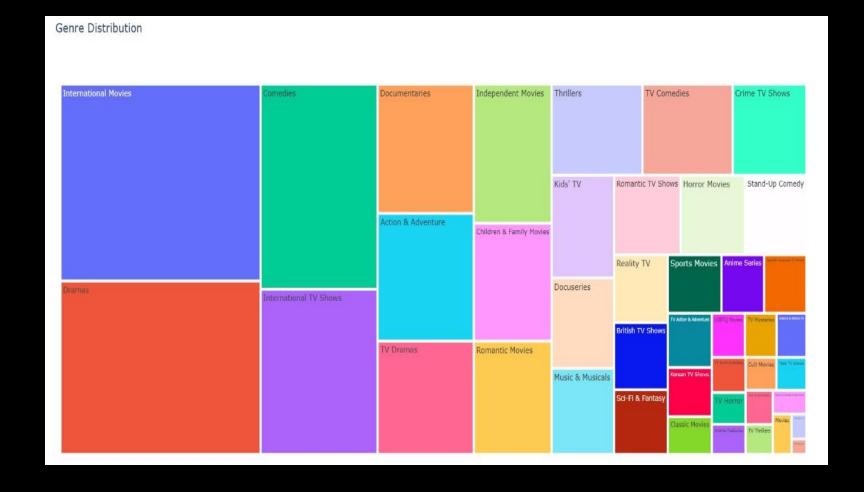
# Plot a bar chart for the genre distribution
genre_counts.plot(kind='bar', figsize=(10, 6), color='skyblue')
plt.title('Genre Distribution')
plt.xlabel('Genre')
plt.ylabel('Number of Movies')
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for
plt.show()
```





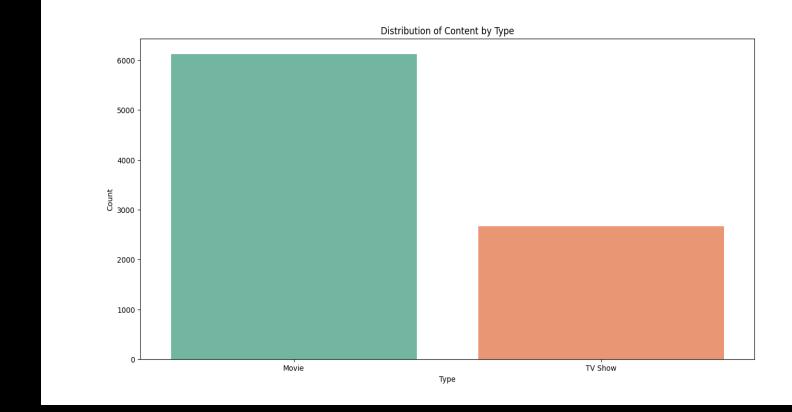
Tree Map of Genre Distribution

```
# Count the occurrences of each genre
genre_counts = genres_series.value_counts().reset_index()
genre_counts.columns = ['Genre', 'Count']
fig = px.treemap(genre_counts, path=['Genre'], values='Count', title='Genre Distribution')
fig.show()
```





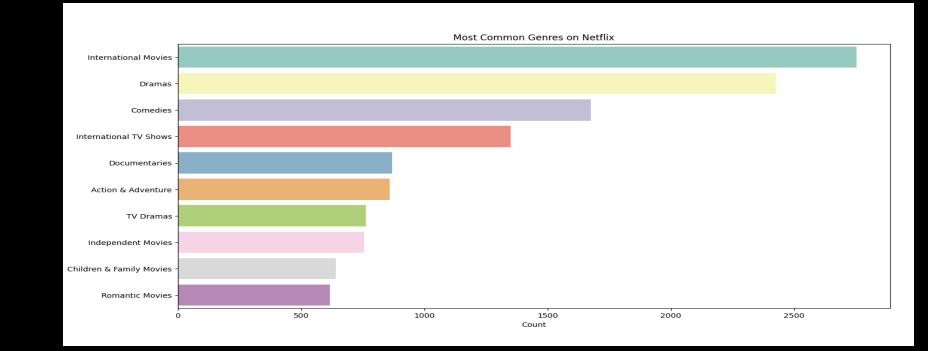
Distribution of Content by Type





Most common Genres on Netflix

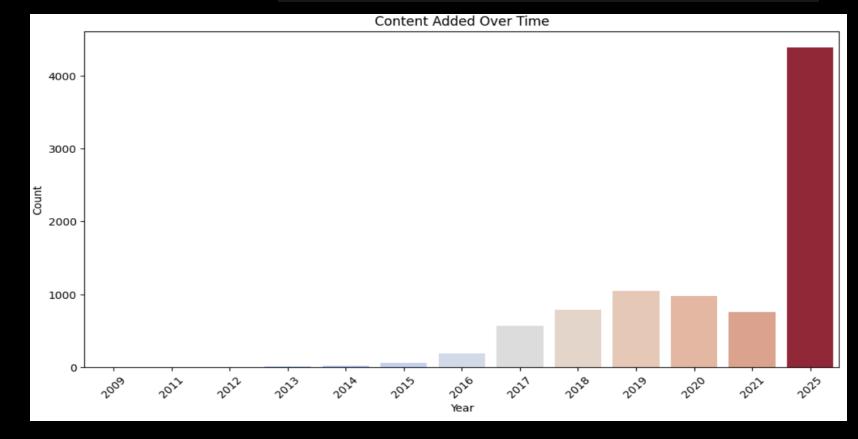
Most Common Genres data['genres'] = data['listed_in'].apply(lambda x: x.split(', ')) all genres = sum(data['genres'], []) genre counts = pd.Series(all genres).value counts().head(10) plt.figure(figsize=(10, 6)) # Explicitly assigning 'y' to 'hue' to avoid the deprecation warning sns.barplot(x=genre counts.values, y=genre_counts.index, hue=genre counts.values, # Assign 'y' variable to 'hue' palette='Set3', dodge=False) plt.title('Most Common Genres on Netflix') plt.xlabel('Count') plt.ylabel('Genre') # Remove the legend since it's unnecessary in this case plt.legend([], [], frameon=False) plt.show()





Content Added Over Time

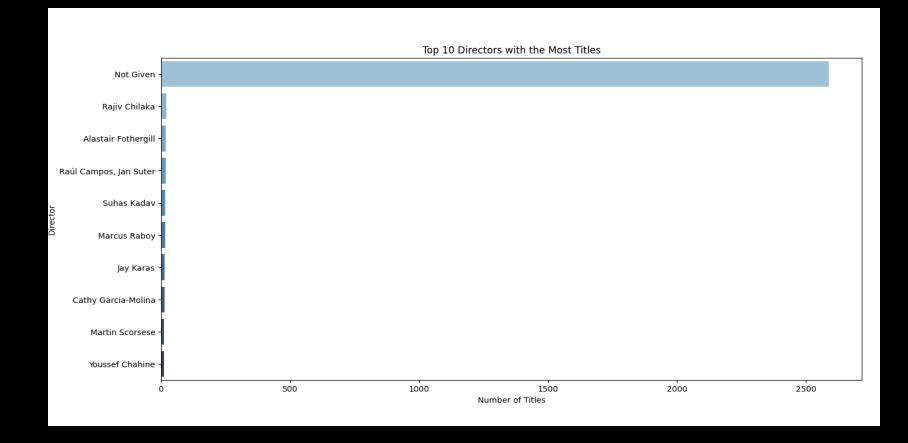
```
# Content Added Over Time
plt.figure(figsize=(12, 6))
# Explicitly assigning 'x' to 'hue' to avoid the deprecation warning
sns.countplot(x='year_added', data=data, palette='coolwarm', hue='year_added')
# Add title and labels
plt.title('Content Added Over Time')
plt.xlabel('Year')
plt.ylabel('Year')
plt.ylabel('Count')
plt.xticks(rotation=45)
# Remove the legend since it's unnecessary in this case
plt.legend([], [], frameon=False)
plt.show()
```





Top 10 Directors with Most Titles

```
# Top 10 Directors with the Most Titles
plt.figure(figsize=(10, 6))
# Plotting the bar chart without hue
sns.barplot(x=top directors.values, y=top directors.index)
# Add title and labels
plt.title('Top 10 Directors with the Most Titles')
plt.xlabel('Number of Titles')
plt.ylabel('Director')
plt.show()
```





Word Cloud of Movie Titles

```
Baby Game
Rise
                                                Legend
                               Woman
Present
   Fire
                                 _aBad
```



Key Insights:

- Most Netflix content falls under a few dominant genres such as Drama, Comedy, Documentaries.
- Netflix adds content at varying rates over the years, highlighting trends in expansion.
- Movies dominate over other TV shows in Netflix's content library.
- The top directors, Rajiv Chilaka, Alastair Fothergill, Raul Campos, Jan Sutas, and others, with the most Netflix content were identified.