

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns

1 df = pd.read_csv("/content/netflix-rotten-tomatoes-metacritic-imdb.csv")

1 df.shape
→ (9425, 29)

1 # Split the categorical variables and keep only the first string
2 categorical_columns = ["Genre", "Tags", "Languages", "Director", "Writer", "Actors", "Country Availability"]
3 for col in categorical_columns:
4     df[col] = df[col].str.split(',').str[0]

1 df.info()
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 9425 entries, 0 to 9424
Data columns (total 29 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Title            9425 non-null    object  
 1   Genre            9400 non-null    object  
 2   Tags             9389 non-null    object  
 3   Languages        9266 non-null    object  
 4   Series or Movie  9425 non-null    object  
 5   Hidden Gem Score 9415 non-null    float64 
 6   Country Availability 9414 non-null    object  
 7   Runtime          9424 non-null    object  
 8   Director         7120 non-null    object  
 9   Writer           7615 non-null    object  
 10  Actors           9314 non-null    object  
 11  View Rating     6827 non-null    object  
 12  IMDb Score      9417 non-null    float64 
 13  Rotten Tomatoes Score 5445 non-null    float64 
 14  Metacritic Score 4082 non-null    float64 
 15  Awards Received  5226 non-null    float64 
 16  Awards Nominated For 6376 non-null    float64 
 17  Boxoffice        3754 non-null    float64 
 18  Release Date    9217 non-null    datetime64[ns]
 19  Netflix Release Date 9425 non-null    datetime64[ns]
 20  Production House 4393 non-null    object  
 21  Netflix Link     9425 non-null    object  
 22  IMDb Link        9101 non-null    object  
 23  Summary          9420 non-null    object  
 24  IMDb Votes       9415 non-null    float64 
 25  Image            9425 non-null    object  
 26  Poster           8487 non-null    object  
 27  TMDb Trailer    9425 non-null    object  
 28  Trailer Site     9424 non-null    object  
dtypes: datetime64[ns](2), float64(8), object(19)
memory usage: 2.1+ MB

1 df.head()
```

	Title	Genre	Tags	Languages	Series or Movie	Hidden Gem Score	Country Availability	Runtime	Director	Writer	...	Netflix Release Date	Production House
0	Lets Fight Ghost	Crime	Comedy Programmes	Swedish	Series	4.3	Thailand	< 30 minutes	Tomas Alfredson	John Ajvide Lindqvist	...	2021-03-04	Canal+ Sandre Metronome
1	HOW TO BUILD A GIRL	Comedy	Dramas	English	Movie	7.0	Canada	1-2 hour	Coky Giedroyc	Caitlin Moran	...	2021-03-04	Film + Monument Picture Lionsgate
2	The Con-Heartist	Comedy	Romantic Comedies	Thai	Movie	8.6	Thailand	> 2 hrs	Mez Tharatorn	Pattaranad Bhiboonsawade	...	2021-03-03	Na
3	Gleboka woda	Drama	TV Dramas	Polish	Series	8.7	Poland	< 30 minutes	Nan	Nan	...	2021-03-03	Na
4	Only a Mother	Drama	Social Issue Dramas	Swedish	Movie	8.3	Lithuania	1-2 hour	Alf Sjöberg	Ivar Lo-Johansson	...	2021-03-03	Na

5 rows × 29 columns

1 df.isna().sum()

→ Title	0
Genre	25
Tags	36
Languages	159
Series or Movie	0
Hidden Gem Score	10
Country Availability	11
Runtime	1
Director	2305
Writer	1810
Actors	111
View Rating	2598
IMDb Score	8
Rotten Tomatoes Score	3980
Metacritic Score	5343
Awards Received	4199
Awards Nominated For	3049
Boxoffice	5671
Release Date	208
Netflix Release Date	0
Production House	5032
Netflix Link	0
IMDb Link	324
Summary	5
IMDb Votes	10
Image	0
Poster	938
TMDb Trailer	0
Trailer Site	1
dtype:	int64

## ▼ Missing Values

```

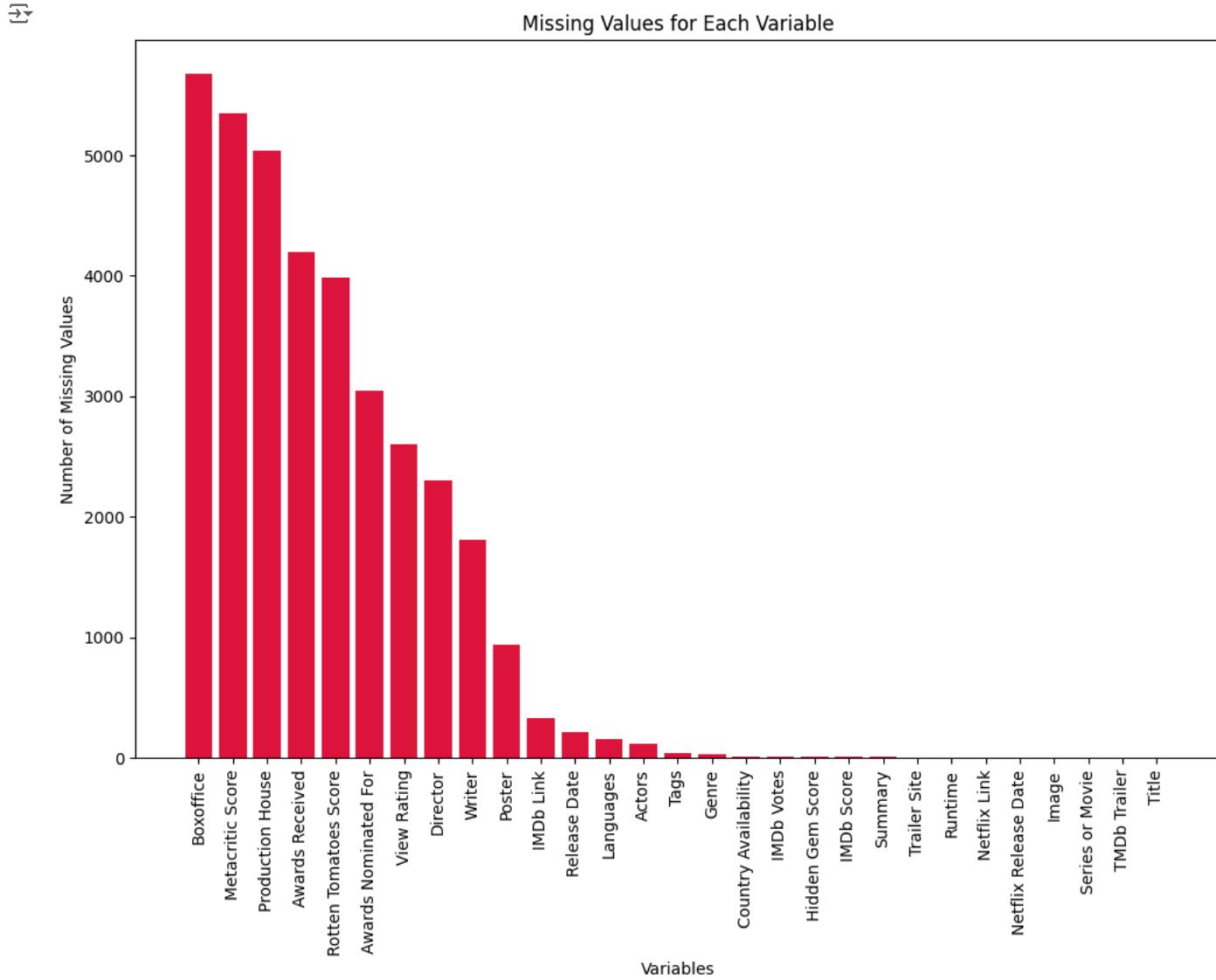
1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 # Define the variables and their corresponding missing value counts
5 variables = ['Title', 'Genre', 'Tags', 'Languages', 'Series or Movie', 'Hidden Gem Score',
6               'Country Availability', 'Runtime', 'Director', 'Writer', 'Actors', 'View Rating',
7               'IMDb Score', 'Rotten Tomatoes Score', 'Metacritic Score', 'Awards Received'],

```

```

8      'Awards Nominated For', 'Boxoffice', 'Release Date', 'Netflix Release Date',
9      'Production House', 'Netflix Link', 'IMDb Link', 'Summary', 'IMDb Votes',
10     'Image', 'Poster', 'TMDb Trailer', 'Trailer Site']
11
12 missing_values = [0, 25, 36, 159, 0, 10, 11, 1, 2305, 1810, 111, 2598, 8, 3980, 5343, 4199, 3049,
13      5671, 208, 0, 5032, 0, 324, 5, 10, 0, 938, 0, 1]
14
15 # Create a DataFrame to store the variables and their missing value counts
16 missing_data_df = pd.DataFrame({'Variable': variables, 'Missing Values': missing_values})
17
18 # Sort the DataFrame by missing value counts in descending order
19 missing_data_df.sort_values(by='Missing Values', ascending=False, inplace=True)
20
21 # Create a bar plot to visualize missing values
22 plt.figure(figsize=(12, 8))
23 plt.bar(missing_data_df['Variable'], missing_data_df['Missing Values'], color='crimson')
24 plt.xticks(rotation=90)
25 plt.xlabel('Variables')
26 plt.ylabel('Number of Missing Values')
27 plt.title('Missing Values for Each Variable')
28 plt.show()
29

```



✓ Columns dropped with more than 50% values missing

```

1 threshold = len(df) * 0.5 # Calculate the threshold for 50% missing values
2 subset_df = df.dropna(thresh=threshold, axis=1)
3
4 # Print the columns dropped

```

```

5 dropped_columns = set(df.columns) - set(subset_df.columns)
6
7 print(dropped_columns)
8
9 df = subset_df

→ {'Boxoffice', 'Metacritic Score', 'Production House'}

```

## ▼ Impute values

```

1 # Calculate the 5% threshold for null values
2 threshold = len(df) * 0.05
3
4 # List of numerical columns (int or float)
5 numerical_columns = ['Hidden Gem Score', 'IMDb Score', 'Rotten Tomatoes Score', 'Awards Received', 'Awards Nominated For',
6
7 # List of categorical columns (object type)
8 categorical_columns = ['Genre', 'Tags', 'Languages', 'Series or Movie', 'Country Availability', 'Runtime', 'Director', 'Writer'
9
10 # Replace null values in numerical columns if the count is less than 5% of the DataFrame
11 for column in numerical_columns:
12     if df[column].isnull().sum() < threshold:
13         df[column].fillna(df[column].mean(), inplace=True)
14
15 # Replace null values in categorical columns if the count is less than 5% of the DataFrame
16 for column in categorical_columns:
17     if df[column].isnull().sum() < threshold:
18         df[column].fillna(df[column].mode().iloc[0], inplace=True)
19
20
21

```

→ <ipython-input-10-78c8223369d0>:13: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df[column].fillna(df[column].mean(), inplace=True)

```

1 df['Director'].fillna('Unknown', inplace=True)
2 df['Writer'].fillna('Unknown', inplace=True)

1 df.drop(columns=['View Rating', 'Poster', 'TMDb Trailer', 'Image', 'Trailer Site', 'Netflix Link', 'IMDb Link', 'Summary'], inplace=True)

1 df['Awards Received'].fillna(0, inplace=True)
2 df['Awards Nominated For'].fillna(0, inplace=True)

1 # Calculate the mean of 'Rotten Tomatoes Score'
2 rotten_tomatoes_mean = df['Rotten Tomatoes Score'].mean()
3
4 # Replace null values in 'Rotten Tomatoes Score' column with its mean
5 df['Rotten Tomatoes Score'].fillna(rotten_tomatoes_mean, inplace=True)

1 df = df.dropna(subset=['Release Date'])

1 df.info()

```

→ <class 'pandas.core.frame.DataFrame'>  
Int64Index: 9217 entries, 0 to 9423  
Data columns (total 18 columns):  

#	Column	Non-Null Count	Dtype
0	Title	9217 non-null	object
1	Genre	9217 non-null	object
2	Tags	9217 non-null	object
3	Languages	9217 non-null	object
4	Series or Movie	9217 non-null	object
5	Hidden Gem Score	9217 non-null	float64
6	Country Availability	9217 non-null	object
7	Runtime	9217 non-null	object
8	Director	9217 non-null	object
9	Writer	9217 non-null	object
10	Actors	9217 non-null	object

```

11 IMDb Score      9217 non-null   float64
12 Rotten Tomatoes Score 9217 non-null   float64
13 Awards Received    9217 non-null   float64
14 Awards Nominated For 9217 non-null   float64
15 Release Date       9217 non-null   datetime64[ns]
16 Netflix Release Date 9217 non-null   datetime64[ns]
17 IMDb Votes         9217 non-null   float64
dtypes: datetime64[ns](2), float64(6), object(10)
memory usage: 1.3+ MB

```

```

1 df['Release Date'] = pd.to_datetime(df['Release Date'])
2
3 # Find the last date in the 'Release Date' column
4 last_date = df['Release Date'].max()
5
6 # Calculate the number of days from 'Release Date' till the last date
7 df['Days_Since_Movie_Released'] = (last_date - df['Release Date']).dt.days

```

→ <ipython-input-17-8fb3f2ad1857>:1: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df['Release Date'] = pd.to\_datetime(df['Release Date'])  
<ipython-input-17-8fb3f2ad1857>:7: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df['Days\_Since\_Movie\_Released'] = (last\_date - df['Release Date']).dt.days

## Extracting year and month for further analysis of release date

```

1 # Extracting year and month from 'Release Date'
2 df['Release Year'] = df['Release Date'].dt.year
3 df['Release Month'] = df['Release Date'].dt.month

```

→ <ipython-input-18-36b4f31718c4>:2: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df['Release Year'] = df['Release Date'].dt.year  
<ipython-input-18-36b4f31718c4>:3: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df['Release Month'] = df['Release Date'].dt.month

```

1 df['Netflix Release Date'] = pd.to_datetime(df['Netflix Release Date'])
2
3 # Find the last date in the 'Release Date' column
4 last_date = df['Netflix Release Date'].max()
5
6 # Calculate the number of days from 'Release Date' till the last date
7 df['Days_Since_Movie_Released_on_Netflix'] = (last_date - df['Netflix Release Date']).dt.days

```

→ <ipython-input-19-caa7da3ce69c>:1: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df['Netflix Release Date'] = pd.to\_datetime(df['Netflix Release Date'])  
<ipython-input-19-caa7da3ce69c>:7: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df['Days\_Since\_Movie\_Released\_on\_Netflix'] = (last\_date - df['Netflix Release Date']).dt.days

```

1 df.info()
2
3 → <class 'pandas.core.frame.DataFrame'>
Int64Index: 9217 entries, 0 to 9423
Data columns (total 22 columns):
 #   Column          Non-Null Count  Dtype  

```

```
---- ----
 0 Title          9217 non-null  object
 1 Genre          9217 non-null  object
 2 Tags           9217 non-null  object
 3 Languages      9217 non-null  object
 4 Series or Movie 9217 non-null  object
 5 Hidden Gem Score 9217 non-null  float64
 6 Country Availability 9217 non-null  object
 7 Runtime         9217 non-null  object
 8 Director        9217 non-null  object
 9 Writer          9217 non-null  object
 10 Actors         9217 non-null  object
 11 IMDb Score    9217 non-null  float64
 12 Rotten Tomatoes Score 9217 non-null  float64
 13 Awards Received 9217 non-null  float64
 14 Awards Nominated For 9217 non-null  float64
 15 Release Date   9217 non-null  datetime64[ns]
 16 Netflix Release Date 9217 non-null  datetime64[ns]
 17 IMDb Votes     9217 non-null  float64
 18 Days_Since_Movie_Released 9217 non-null  int64
 19 Release Year   9217 non-null  int64
 20 Release Month  9217 non-null  int64
 21 Days_Since_Movie_Released_on_Netflix 9217 non-null  int64
dtypes: datetime64[ns](2), float64(6), int64(4), object(10)
memory usage: 1.6+ MB
```

```
1 df.drop(['Release Date','Netflix Release Date'],axis=1,inplace=True)
```

→ <ipython-input-21-1f093ac1e93f>:1: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df.drop(['Release Date','Netflix Release Date'],axis=1,inplace=True)

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

→ <class 'pandas.core.frame.DataFrame'>  
Int64Index: 9217 entries, 0 to 9423  
Data columns (total 20 columns):  
 # Column Non-Null Count Dtype  
---  
 0 Title 9217 non-null object  
 1 Genre 9217 non-null object  
 2 Tags 9217 non-null object  
 3 Languages 9217 non-null object  
 4 Series or Movie 9217 non-null object  
 5 Hidden Gem Score 9217 non-null float64  
 6 Country Availability 9217 non-null object  
 7 Runtime 9217 non-null object  
 8 Director 9217 non-null object  
 9 Writer 9217 non-null object  
 10 Actors 9217 non-null object  
 11 IMDb Score 9217 non-null float64  
 12 Rotten Tomatoes Score 9217 non-null float64  
 13 Awards Received 9217 non-null float64  
 14 Awards Nominated For 9217 non-null float64  
 15 IMDb Votes 9217 non-null float64  
 16 Days\_Since\_Movie\_Released 9217 non-null int64  
 17 Release Year 9217 non-null int64  
 18 Release Month 9217 non-null int64  
 19 Days\_Since\_Movie\_Released\_on\_Netflix 9217 non-null int64  
dtypes: float64(6), int64(4), object(10)
memory usage: 1.5+ MB

## Checking Outliers

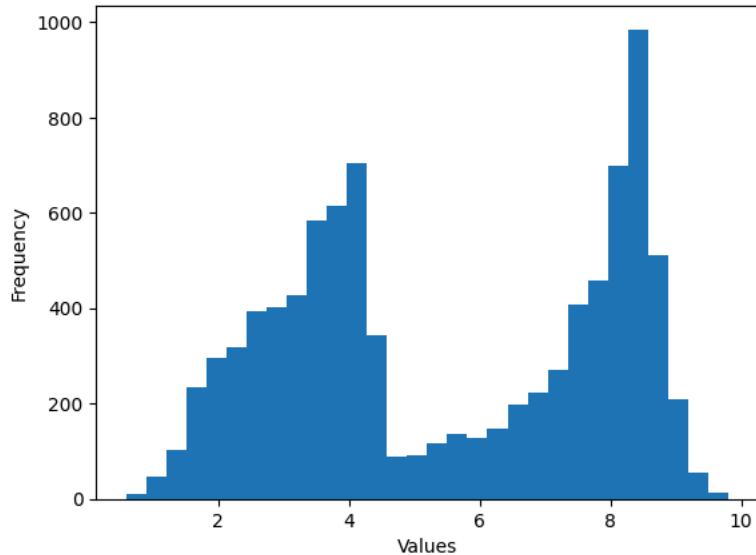
```
1 numerical_columns = ["Hidden Gem Score", "IMDb Score", "Rotten Tomatoes Score",
2                      "Awards Received", "Awards Nominated For", "IMDb Votes","Days_Since_Movie_Released","Days_Since_Movie_Rele
3
4

1 import matplotlib.pyplot as plt
2
3 # Assuming your data is stored in a DataFrame called 'df' and 'column' represents the column of interest
4 dict_skew = {}
5 for i in numerical_columns:
6     plt.hist(df[i], bins=30)
```

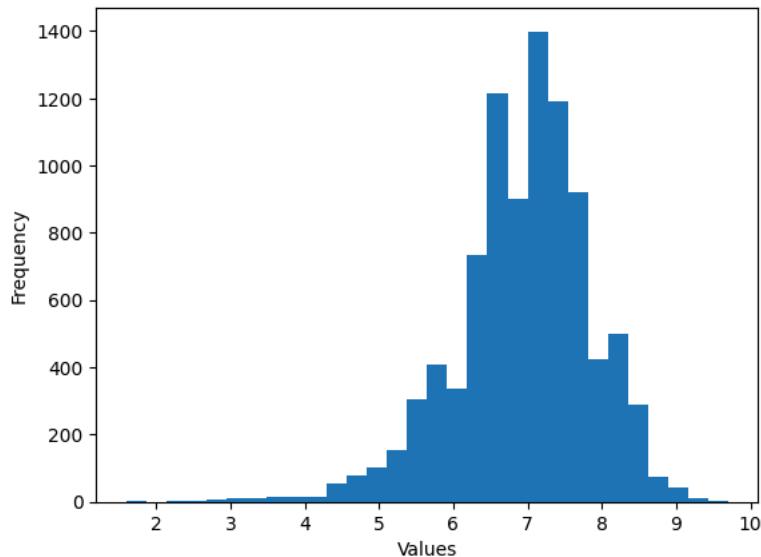
```
7 plt.xlabel('Values')
8 plt.ylabel('Frequency')
9 plt.title(i)
10 plt.show()
11 data = df[i]
12 skewness = np.mean((data - np.mean(data))**3) / np.power(np.var(data), 1.5)
13 dict_skew[i] = skewness
14
15 for i,k in dict_skew.items():
16     print(f"Pearson's First Coefficient of Skewness for {i}:", k)
17
18
```



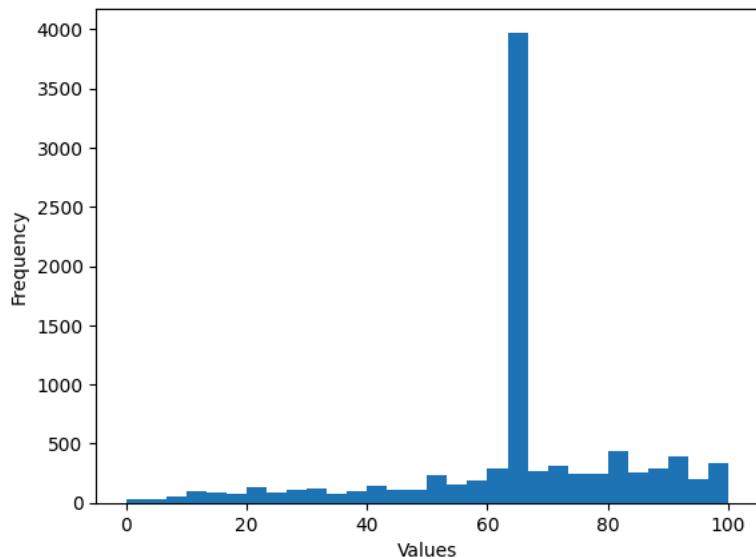
Hidden Gem Score



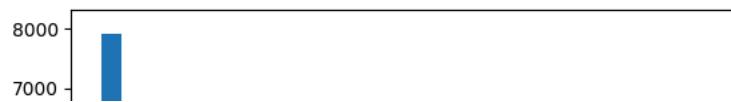
IMDb Score

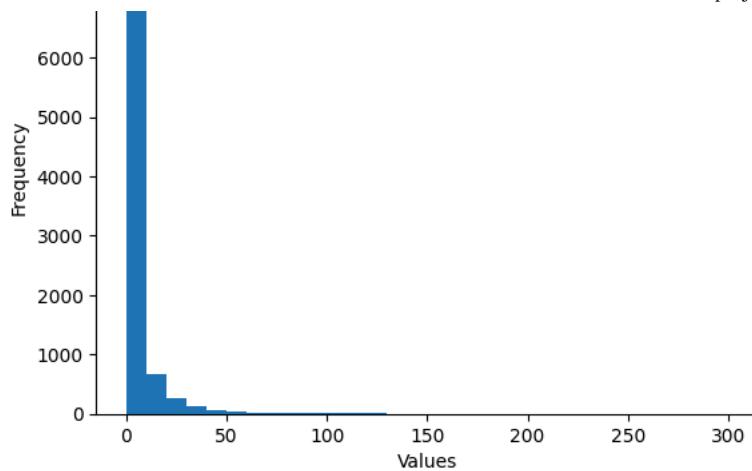


Rotten Tomatoes Score

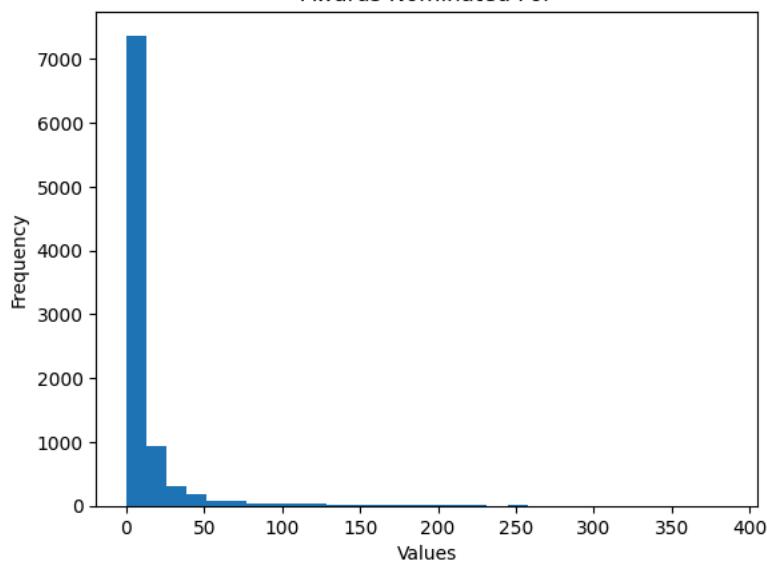


Awards Received

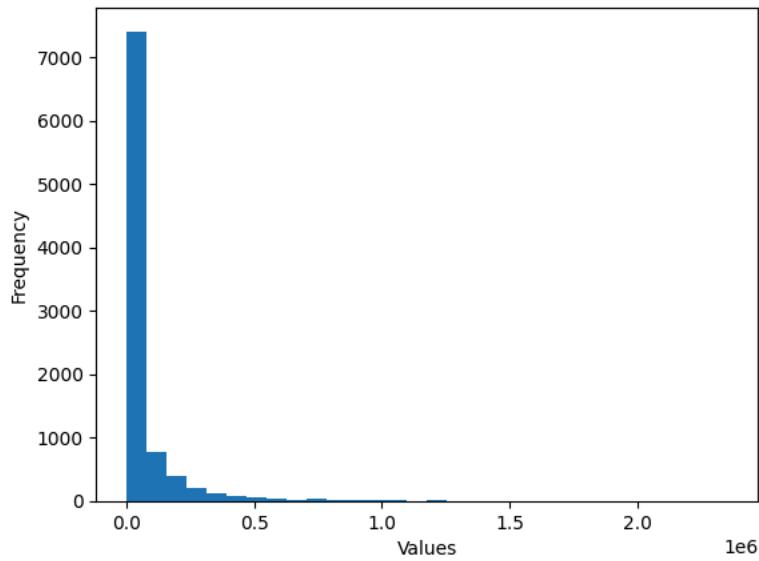




Awards Nominated For

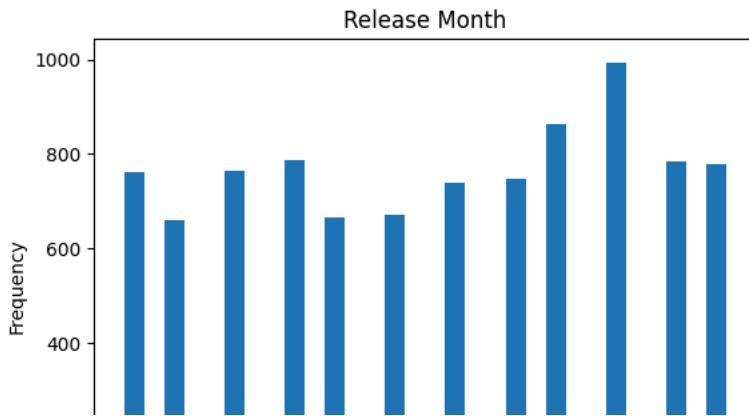
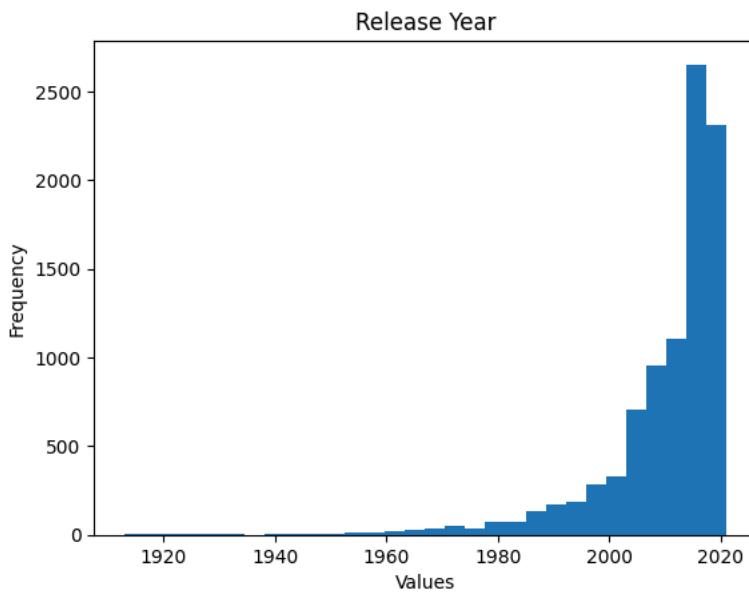
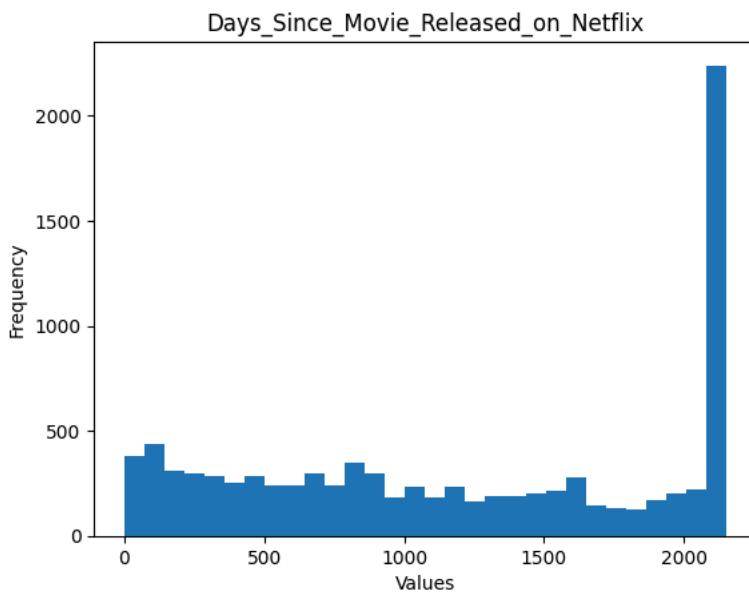
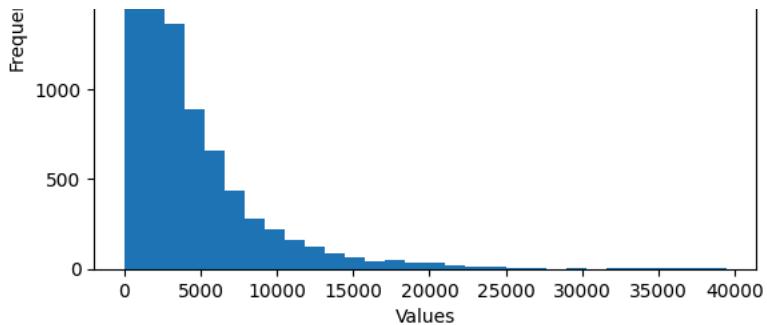


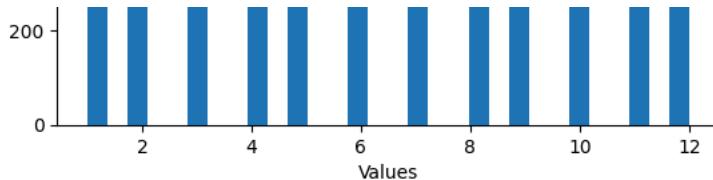
IMDb Votes



Days\_Since\_Movie\_Released







Pearson's First Coefficient of Skewness for Hidden Gem Score: -0.006706672020348805  
Pearson's First Coefficient of Skewness for IMDb Score: -0.6414454765727965  
Pearson's First Coefficient of Skewness for Rotten Tomatoes Score: -0.8508410923195496  
Pearson's First Coefficient of Skewness for Awards Received: 7.473282616099965  
Pearson's First Coefficient of Skewness for Awards Nominated For: 5.810237735721319  
Pearson's First Coefficient of Skewness for IMDb Votes: 5.526193276102253  
Pearson's First Coefficient of Skewness for Days\_Since\_Movie\_Released: 2.7643856653018335  
Pearson's First Coefficient of Skewness for Days\_Since\_Movie\_Released\_on\_Netflix: -0.04782662642652016  
Pearson's First Coefficient of Skewness for Release Year: -2.7610321481500946  
Pearson's First Coefficient of Skewness for Release Month: -0.10273635829957133

- ✓ Data is not normally distributed, so log transformation is applied stabilize the variance and reduce the impact of outliers.

```

1 # import numpy as np
2
3 # # Assuming your data is stored in a DataFrame called 'df'
4 # # Numeric columns
5 # numerical_columns = ["Hidden Gem Score", "IMDb Score", "Rotten Tomatoes Score",
6 #                      "Awards Received", "Awards Nominated For", "IMDb Votes",
7 #                      "Days_Since_Movie_Released", "Release Year"]
8
9 # # Apply log transformation to columns with positive skewness
10 # log_transform_columns = ["Awards Received", "IMDb Votes",
11 #                           "Days_Since_Movie_Released", "Release Year"]
12
13 # for column in log_transform_columns:
14 #     df[column] = np.log(df[column])
15
16 # # Apply square root transformation to column with negative skewness
17 # square_root_transform_column = "IMDb Score"
18 # df[square_root_transform_column] = np.sqrt(df[square_root_transform_column])
19

1 # # Calculate skewness before transformation
2 skewness_before = df[numerical_columns].skew()
3
4 # # Apply log transformation to columns with positive skewness
5 log_transform_columns = ["Days_Since_Movie_Released", "Release Year"]
6
7 # for column in log_transform_columns:
8 #     # Add a small positive constant (1e-6) to the data to avoid zero or negative values
9 #     df[column] = np.log(df[column] + 1e-6)
10 #
11 # # Apply square root transformation to column with negative skewness
12 # square_root_transform_column = "IMDb Score"
13 # # Adding a small constant to handle any negative values (e.g., 1e-6)
14 # df[square_root_transform_column] = np.sqrt(df[square_root_transform_column] + 1e-6)
15
16 # # Calculate skewness after transformation
17 skewness_after = df[numerical_columns].skew()
18
19 # # Display skewness before and after transformation
20 # print("Skewness Before Transformation:")
21 # print(skewness_before)
22
23 # print("\nSkewness After Transformation:")
24 # print(skewness_after)
25
26 # # Plot histograms before and after transformation
27 # plt.figure(figsize=(12, 6))
28 # for i, column in enumerate(numerical_columns, 1):
29 #     plt.subplot(2, 4, i)
30 #     plt.hist(df[column], bins=30)
31 #     plt.xlabel('Values')
32 #     plt.ylabel('Frequency')
33 #     plt.title(column)
34
35 # plt.tight_layout()
36 # plt.show()

```

## Descriptive analysis questions

- ✓ What are the most common genres and languages among movies and web series available on Netflix?

<sup>1</sup>  
2 # Get the most common genres and languages

```
3 most_common_genres = df['Genre'].value_counts().head(10)
4 most_common_languages = df['Languages'].value_counts().head(10)
5
6 print("Most Common Genres:")
7 print(most_common_genres)
8
9 print("\nMost Common Languages:")
10 print(most_common_languages)
11
12 # Create bar chart for most common genres
13 plt.figure(figsize=(10, 6))
14 most_common_genres.plot(kind='bar', color='crimson')
15 plt.title('Top 10 Most Common Genres on Netflix')
16 plt.xlabel('Genre')
17 plt.ylabel('Frequency')
18 plt.xticks(rotation=45)
19 plt.tight_layout()
20 plt.show()
21
22
23 # Create bar chart for most common languages
24 plt.figure(figsize=(10, 6))
25 most_common_languages.plot(kind='bar', color='black')
26 plt.title('Top 10 Most Common Languages on Netflix')
27 plt.xlabel('Language')
28 plt.ylabel('Frequency')
29 plt.xticks(rotation=45)
30 plt.tight_layout()
31 plt.show()
32
33
34
35 plt.show()
```

## ↳ Most Common Genres:

Comedy	2071
Drama	1801
Action	1651
Animation	1119
Documentary	767
Crime	556
Biography	386
Adventure	294
Horror	203
Reality-TV	57

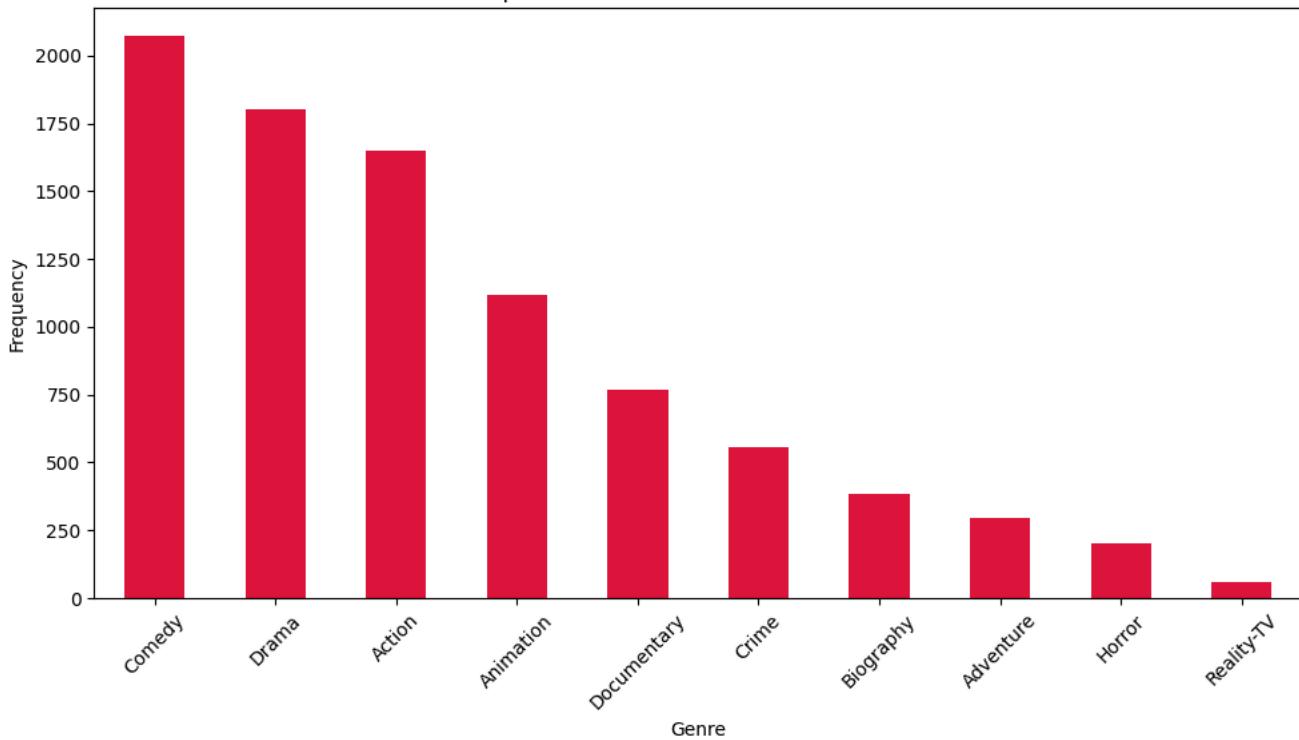
Name: Genre, dtype: int64

## Most Common Languages:

English	5498
Japanese	918
Korean	488
Hindi	278
Spanish	260
French	236
Mandarin	157
German	142
Cantonese	120
Portuguese	77

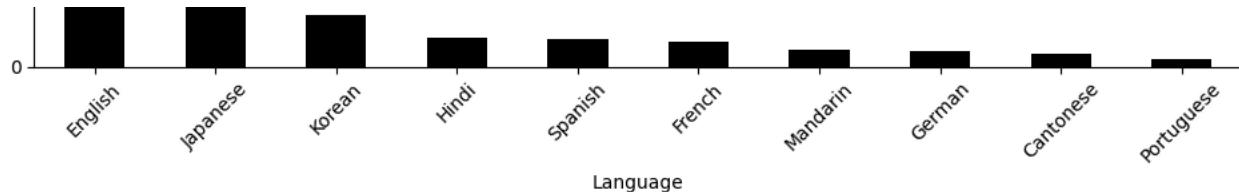
Name: Languages, dtype: int64

Top 10 Most Common Genres on Netflix



Top 10 Most Common Languages on Netflix

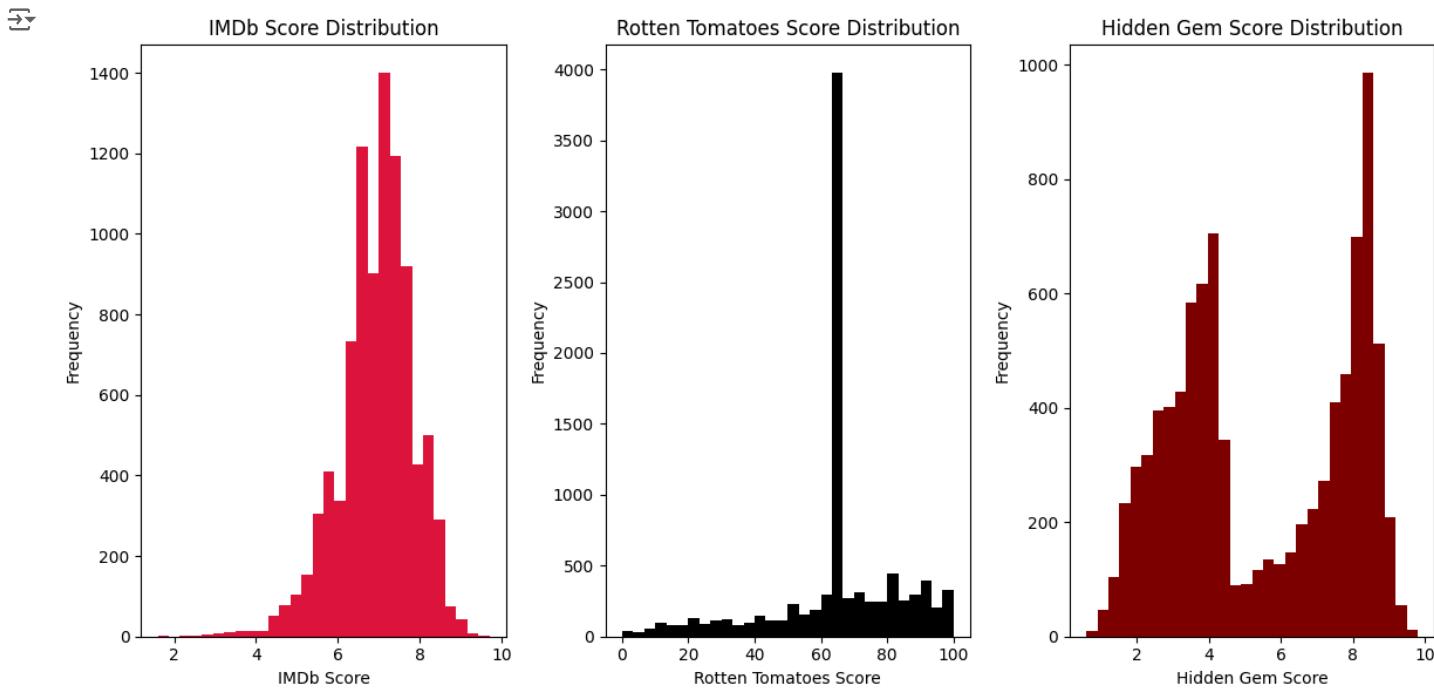




1 Start coding or [generate](#) with AI.

## ▼ New Section

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 plt.figure(figsize=(12, 6))
5
6 # IMDb Score distribution
7 plt.subplot(1, 3, 1)
8 plt.hist(df['IMDb Score'], bins=30, color='crimson')
9 plt.xlabel('IMDb Score')
10 plt.ylabel('Frequency')
11 plt.title('IMDb Score Distribution')
12
13 # Rotten Tomatoes Score distribution
14 plt.subplot(1, 3, 2)
15 plt.hist(df['Rotten Tomatoes Score'], bins=30, color='black')
16 plt.xlabel('Rotten Tomatoes Score')
17 plt.ylabel('Frequency')
18 plt.title('Rotten Tomatoes Score Distribution')
19
20 # Hidden Gem Score distribution
21 plt.subplot(1, 3, 3)
22 plt.hist(df['Hidden Gem Score'], bins=30, color='maroon')
23 plt.xlabel('Hidden Gem Score')
24 plt.ylabel('Frequency')
25 plt.title('Hidden Gem Score Distribution')
26
27 plt.tight_layout()
28 plt.show()
29
```



```

1 import scipy.stats as stats
2
3 # Assuming your data is stored in a DataFrame called 'df'
4 # Assuming 'IMDb Score', 'Rotten Tomatoes Score', and 'Hidden Gem Score' are the columns with scores
5
6 # Calculate the skewness of IMDb Score, Rotten Tomatoes Score, and Hidden Gem Score
7 imdb_skewness = stats.skew(df['IMDb Score'])
8 rt_skewness = stats.skew(df['Rotten Tomatoes Score'])
9 hg_skewness = stats.skew(df['Hidden Gem Score'])
10
11 print("Skewness of IMDb Score:", imdb_skewness)
12 print("Skewness of Rotten Tomatoes Score:", rt_skewness)
13 print("Skewness of Hidden Gem Score:", hg_skewness)

→ Skewness of IMDb Score: -0.6414454765727965
Skewness of Rotten Tomatoes Score: -0.8508410923195496
Skewness of Hidden Gem Score: -0.006706672020348805

```

## Correlation between scores

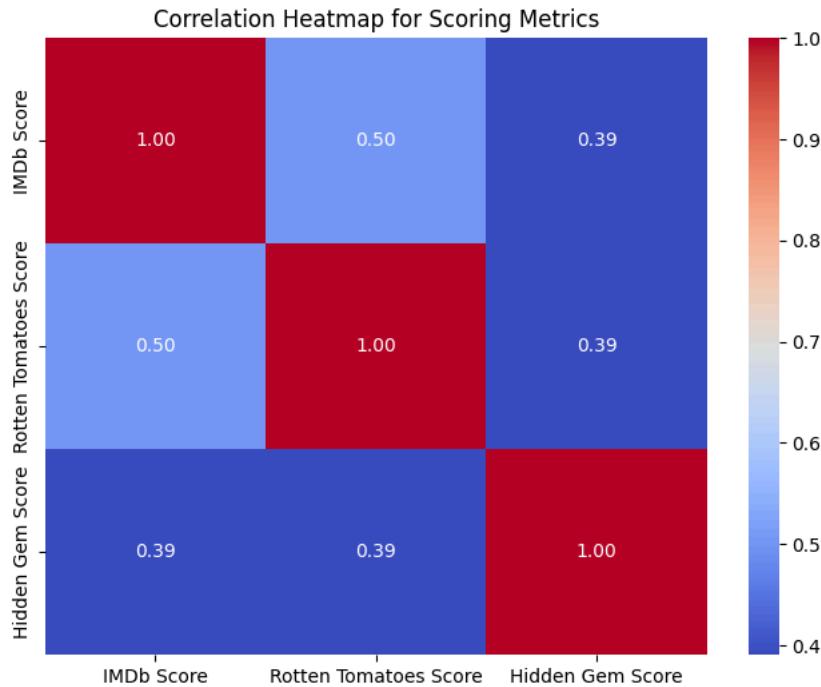
```

1 # Select the relevant columns
2 scores_df = df[['IMDb Score', 'Rotten Tomatoes Score', 'Hidden Gem Score']]
3
4 # Calculate the correlation matrix
5 correlation_matrix = scores_df.corr()
6
7 # Display the correlation matrix
8 print("Correlation Matrix:")
9 print(correlation_matrix)
10
11 # Visualize the correlation matrix as a heatmap
12 import seaborn as sns
13 import matplotlib.pyplot as plt
14
15 plt.figure(figsize=(8, 6))
16 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
17 plt.title("Correlation Heatmap for Scoring Metrics")
18 plt.show()

```

### Correlation Matrix:

	IMDb Score	Rotten Tomatoes Score	Hidden Gem Score
IMDb Score	1.000000	0.502980	0.390769
Rotten Tomatoes Score	0.502980	1.000000	0.390403
Hidden Gem Score	0.390769	0.390403	1.000000



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

→ <class 'pandas.core.frame.DataFrame'>  
Int64Index: 9217 entries, 0 to 9423  
Data columns (total 20 columns):  
 # Column Non-Null Count Dtype   
---  
 0 Title 9217 non-null object   
 1 Genre 9217 non-null object   
 2 Tags 9217 non-null object   
 3 Languages 9217 non-null object   
 4 Series or Movie 9217 non-null object   
 5 Hidden Gem Score 9217 non-null float64  
 6 Country Availability 9217 non-null object   
 7 Runtime 9217 non-null object   
 8 Director 9217 non-null object   
 9 Writer 9217 non-null object   
 10 Actors 9217 non-null object   
 11 IMDb Score 9217 non-null float64  
 12 Rotten Tomatoes Score 9217 non-null float64  
 13 Awards Received 9217 non-null float64  
 14 Awards Nominated For 9217 non-null float64  
 15 IMDb Votes 9217 non-null float64  
 16 Days\_Since\_Movie\_Released 9217 non-null int64  
 17 Release Year 9217 non-null int64  
 18 Release Month 9217 non-null int64  
 19 Days\_Since\_Movie\_Released\_on\_Netflix 9217 non-null int64  
dtypes: float64(6), int64(4), object(10)  
memory usage: 1.5+ MB

```
1 df.to_csv('/content/Final_Netflix_Data_For_Dashboard.csv')
```

## Exploratory data analysis

```
1 import pandas as pd  

2  

3 # Assuming your data is stored in a DataFrame called 'df'  

4 # Assuming 'IMDb Votes' is the relevant column  

5  

6 # Step 1: Calculate the average IMDb votes  

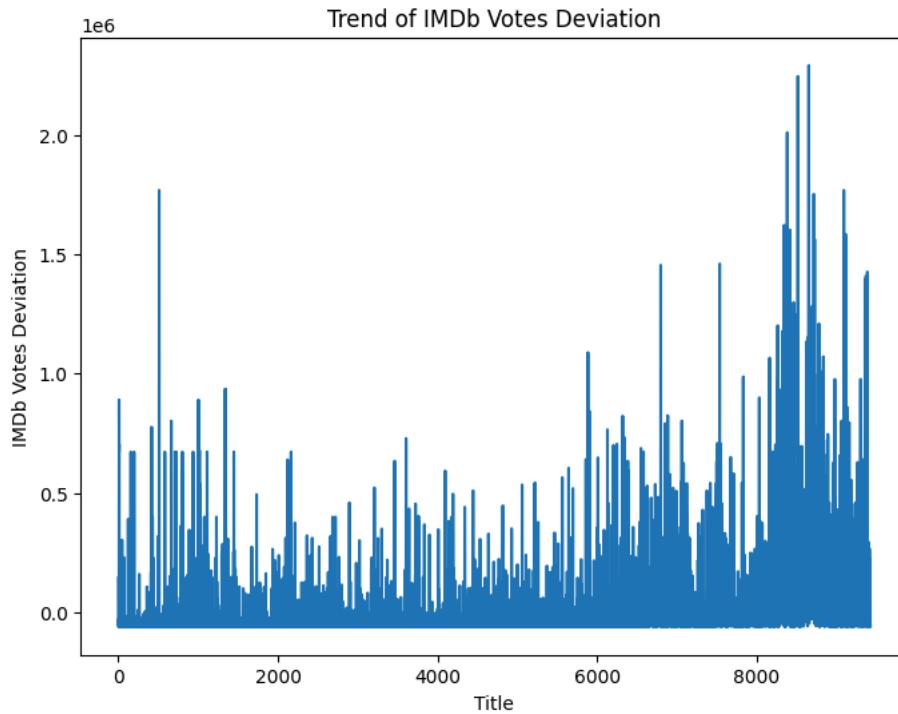
7 avg_imdb_votes = df['IMDb Votes'].mean()  

8
```

```
9 # Step 2: Create the new column by subtracting each IMDb vote record from the average IMDb votes
10 df['IMDb Votes Deviation'] = df['IMDb Votes'] - avg_imdb_votes
11
12 # Step 3: Examine the trend in the new column
13 import matplotlib.pyplot as plt
14
15 plt.figure(figsize=(8, 6))
16 plt.plot(df['IMDb Votes Deviation'])
17 plt.xlabel('Title')
18 plt.ylabel('IMDb Votes Deviation')
19 plt.title('Trend of IMDb Votes Deviation')
20 plt.show()
21
```

→ <ipython-input-33-cc3ac611d136>:10: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)



```
1 df.head(2).transpose()
```



0

1

<b>Title</b>	Lets Fight Ghost	HOW TO BUILD A GIRL
<b>Genre</b>	Crime	Comedy
<b>Tags</b>	Comedy Programmes	Dramas
<b>Languages</b>	Swedish	English
<b>Series or Movie</b>	Series	Movie
<b>Hidden Gem Score</b>	4.3	7.0
<b>Country Availability</b>	Thailand	Canada
<b>Runtime</b>	< 30 minutes	1-2 hour
<b>Director</b>	Tomas Alfredson	Coky Giedroyc
<b>Writer</b>	John Ajvide Lindqvist	Caitlin Moran
<b>Actors</b>	Lina Leandersson	Cleo
<b>IMDb Score</b>	7.9	5.8
<b>Rotten Tomatoes Score</b>	98.0	79.0
<b>Awards Received</b>	74.0	1.0
<b>Awards Nominated For</b>	57.0	0.0
<b>IMDb Votes</b>	205926.0	2838.0
<b>Days_Since_Movie_Released</b>	4718	553
<b>Release Year</b>	2008	2020
<b>Release Month</b>	12	5
<b>Days_Since_Movie_Released_on_Netflix</b>	0	0
<b>IMDb Votes Deviation</b>	144509.472088	-58578.527912

```

1 import pandas as pd
2
3
4 # Assuming your data is stored in a DataFrame called 'df'
5 # Assuming 'IMDb Score' and 'IMDb Votes' are the relevant columns
6
7 # Step 1: Calculate the average IMDb votes
8 avg_imdb_votes = df['IMDb Votes'].mean()
9
10 # Step 2: Create the IMDb Votes Deviation column
11 df['IMDb Votes Deviation'] = (df['IMDb Votes'] - avg_imdb_votes)/len(df['IMDb Votes'])
12
13 # Step 3: Define a function to calculate the weighted IMDb Score
14 def weighted_imdb_score(row):
15     weight_imdb_votes = row['IMDb Votes Deviation']
16     return (row['IMDb Score'] * (1 + weight_imdb_votes))
17
18 # Step 4: Apply the weighted_imdb_score function to create the weighted IMDb Score column
19 df['Weighted IMDb Score'] = df.apply(weighted_imdb_score, axis=1)
20
21 # Display the DataFrame with the new weighted IMDb Score column
22 print(df[['IMDb Score', 'IMDb Votes', 'Weighted IMDb Score']].head())
23

```



	IMDb Score	IMDb Votes	Weighted IMDb Score
0	7.9	205926.0	131.760782
1	5.8	2838.0	-31.061827
2	7.4	131.0	-41.803961
3	7.5	47.0	-42.437231
4	6.7	88.0	-37.880790

<ipython-input-35-c2d5213326ba>:11: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
df['IMDb Votes Deviation'] = (df['IMDb Votes'] - avg\_imdb\_votes)/len(df['IMDb Votes'])  
<ipython-input-35-c2d5213326ba>:19: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)  
`df['Weighted IMDb Score'] = df.apply(weighted_imdb_score, axis=1)`

```
1 df['Weighted IMDb Score'].value_counts()
```

```
→ 480.350121    12
-39.666893      9
-13.883281      8
126.057154      6
323.660372      6
...
-40.385250      1
-40.431063      1
-36.042922      1
-43.386771      1
-44.027245      1
Name: Weighted IMDb Score, Length: 8701, dtype: int64
```

```
1 # Step 5: Calculate the correlation between 'Awards Nominated For' and 'Weighted IMDb Score' for 'Hit' titles
2 correlation = df['Awards Nominated For'].corr(df['Weighted IMDb Score'])
3
4 print("Correlation between 'Awards Nominated For' and 'Weighted IMDb Score' for 'Hit' titles:", correlation)
```

```
→ Correlation between 'Awards Nominated For' and 'Weighted IMDb Score' for 'Hit' titles: 0.5090434323505878
```

```
1 # # Step 1: Calculate the 99th percentile of 'Awards Nominated For'
2 percentile_95 = np.percentile(df['Awards Nominated For'], 95)
3
4 # # Step 2: Remove values greater than 100
5 df = df[df['Awards Nominated For'] <= percentile_95]
6
7
```

```
1 df['Awards Nominated For'].value_counts()
```

```
→ 0.0      2895
1.0      1006
2.0      694
3.0      534
4.0      421
...
248.0    1
165.0    1
112.0    1
263.0    1
138.0    1
Name: Awards Nominated For, Length: 201, dtype: int64
```

```
1 df['Weighted IMDb Score'].value_counts()
```

```
→ 480.350121    12
-39.666893      9
-13.883281      8
126.057154      6
323.660372      6
...
-40.385250      1
-40.431063      1
-36.042922      1
-43.386771      1
-44.027245      1
Name: Weighted IMDb Score, Length: 8701, dtype: int64
```

```
1 positive_count = df[df['Weighted IMDb Score'] > threshold].shape[0]
2 negative_count = df[df['Weighted IMDb Score'] <= threshold].shape[0]
3
4 print("Number of Positive Values:", positive_count)
5 print("Number of Negative Values:", negative_count)
```

```
→ Number of Positive Values: 134
Number of Negative Values: 9083
```

Double-click (or enter) to edit

```

1 import pandas as pd
2 import numpy as np
3
4 # Assuming your data is stored in a DataFrame called 'df'
5 # Assuming 'Awards Nominated For', 'Weighted IMDb Score', and 'Target' are relevant columns
6
7 # Step 1: Calculate the 80th percentile of 'Awards Nominated For'
8 percentile_80 = np.percentile(df['Awards Nominated For'], 80)
9 print('vscdfgh',percentile_80)
10
11 percentile_70 = np.percentile(df['Awards Nominated For'], 60)
12 print('70%',percentile_70)
13
14 # Step 2: Define the success_label function
15 def success_label(nominated_for, weighted_imdb_score, threshold,percentile_70):
16     if nominated_for > threshold and weighted_imdb_score > 0:
17         return 'Super Hit'
18     elif nominated_for > threshold or weighted_imdb_score > 0:
19         return 'Moderate Hit'
20     else:
21         return 'Not Hit'
22
23 # Step 3: Apply the success_label function to create the 'Target' variable
24 df['Target'] = df.apply(lambda row: success_label(row['Awards Nominated For'],
25                                                 row['Weighted IMDb Score'],
26                                                 percentile_80,percentile_70), axis=1)
27
28 # Check the value counts of the 'Target' variable to ensure it's balanced
29 print(df['Target'].value_counts())
30

```

```

→ vscdfgh 13.0
70% 4.0
Not Hit      6241
Moderate Hit 1909
Super Hit    1067
Name: Target, dtype: int64
<ipython-input-42-40b46964ec3c>:24: 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/stable/user\\_guide/indexing.html#returning-a-view](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view)

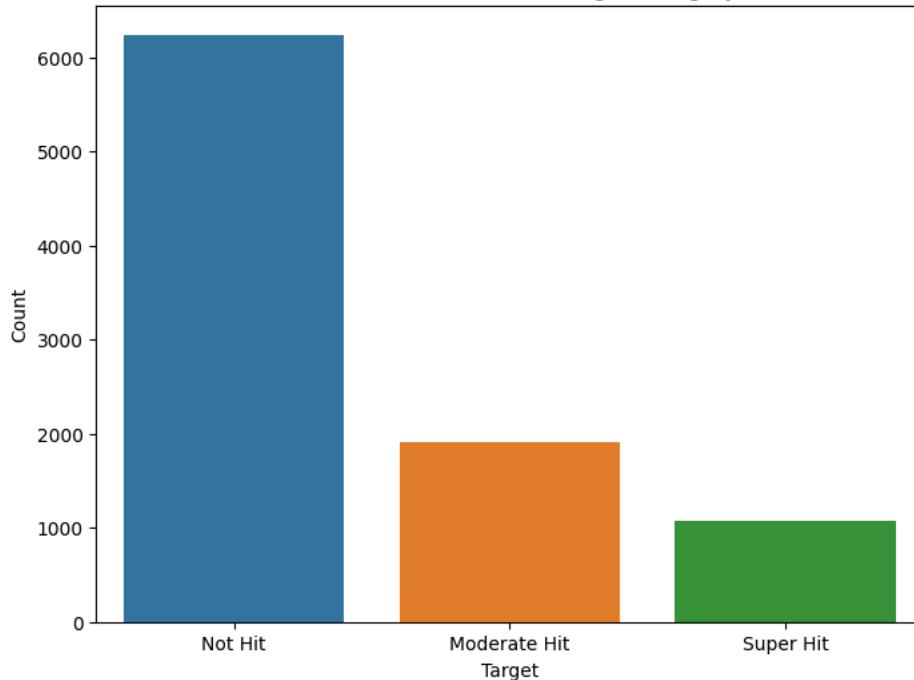
```

1 # Calculate value counts for the 'Target' variable
2 target_value_counts = df['Target'].value_counts()
3
4 # Create a bar chart for the 'Target' variable
5 plt.figure(figsize=(8, 6))
6 sns.barplot(x=target_value_counts.index, y=target_value_counts.values)
7 plt.xlabel('Target')
8 plt.ylabel('Count')
9 plt.title('Number of Titles in Each Target Category')
10 plt.show()

```



Number of Titles in Each Target Category

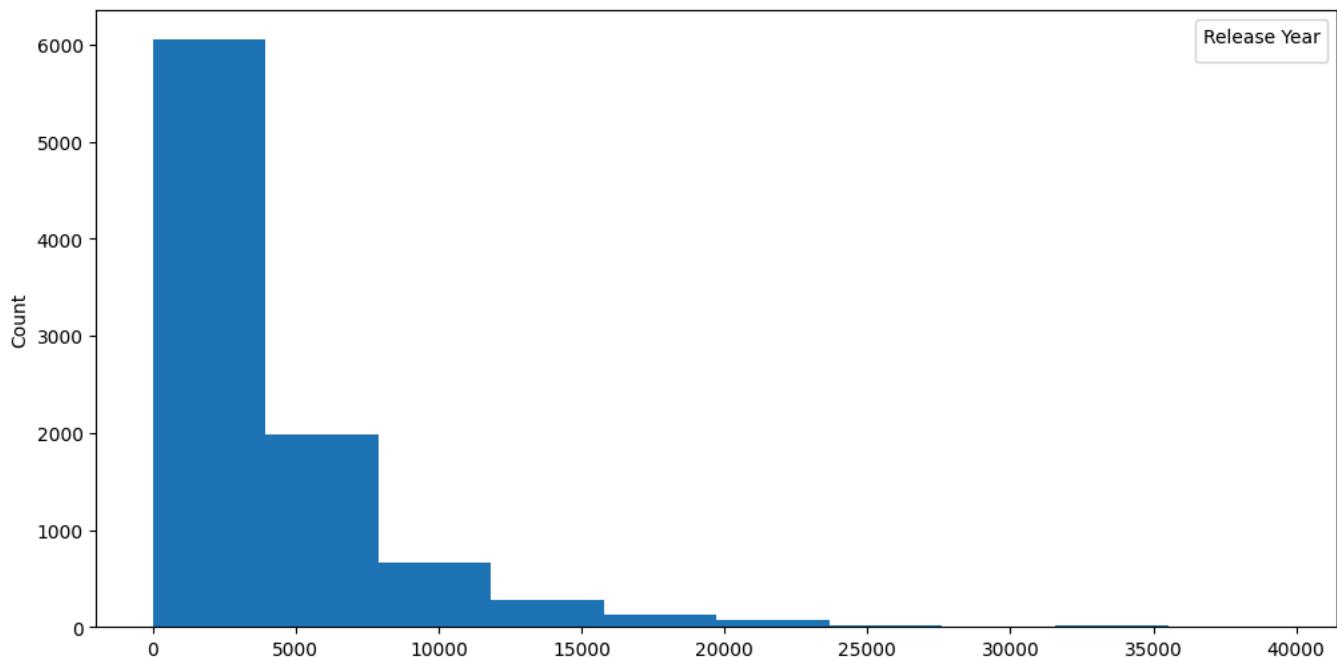


```

1 plt.figure(figsize=(12, 6))
2 plt.hist(x='Days_Since_Movie_Released', data=df, bins=10)
3 # plt.xlabel('Release Year')
4 plt.ylabel('Count')
5 # plt.title('Count Plot of Release Year with Target as the Legend')
6 plt.legend(title='Release Year', loc='upper right') # Customize the legend title and location
7 plt.show()

```

WARNING:matplotlib.legend:No artists with labels found to put in legend. Note that artists whose label start with an underscore



```

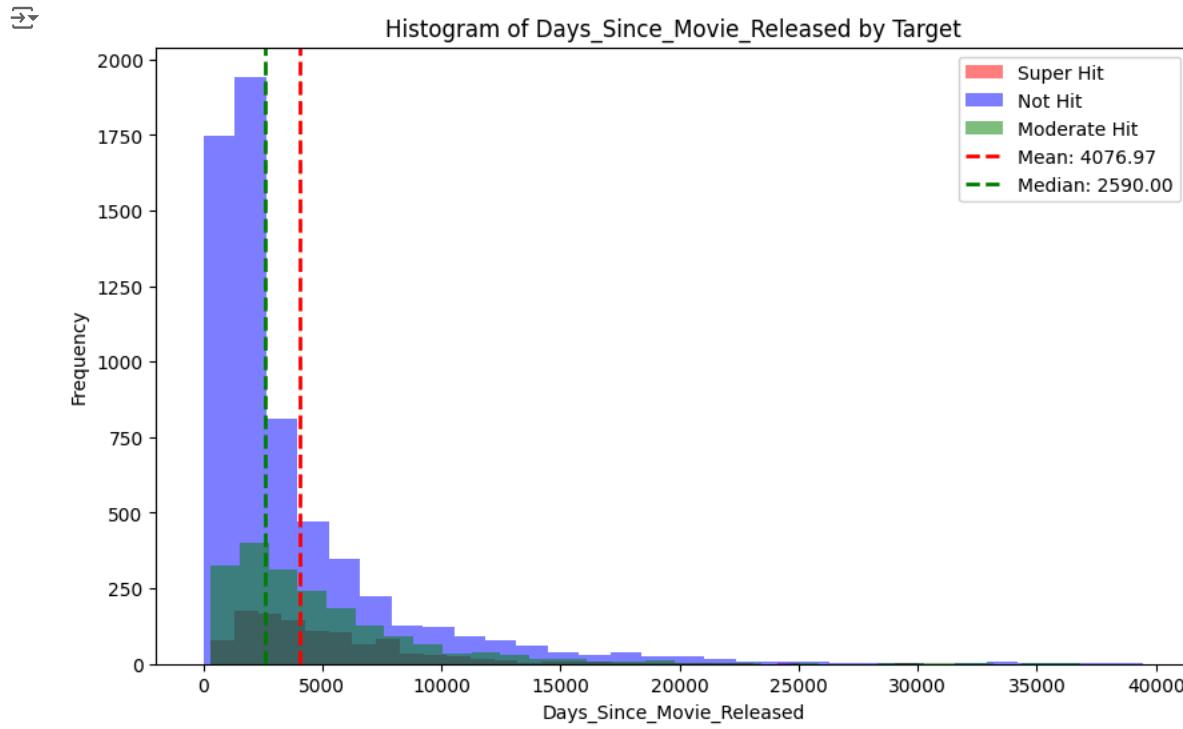
1 # Assuming you have a DataFrame named 'df' and want to plot a histogram of a column 'data_column'
2 data_column = 'Days_Since_Movie_Released'
3
4 categorical_column = 'Target' # Replace with the name of your categorical column
5
6 # Get unique categories in the categorical column and assign colors to each category
7 unique_categories = df[categorical_column].unique()
8 custom_colors = ['red', 'blue', 'green'] # Add more colors as needed

```

```

9 category_colors = custom_colors[:len(unique_categories)]
10
11 plt.figure(figsize=(10, 6))
12 for i, category in enumerate(unique_categories):
13     category_data = df[df[categorical_column] == category][data_column]
14     plt.hist(category_data, bins=30, alpha=0.5, color=category_colors[i], label=category)
15
16 # Calculate mean and median for the overall data
17 mean_val = df[data_column].mean()
18 median_val = df[data_column].median()
19
20 # Add mean and median lines to the histogram
21 plt.axvline(mean_val, color='red', linestyle='dashed', linewidth=2, label=f'Mean: {mean_val:.2f}')
22 plt.axvline(median_val, color='green', linestyle='dashed', linewidth=2, label=f'Median: {median_val:.2f}')
23
24 plt.xlabel(data_column)
25 plt.ylabel('Frequency')
26 plt.title(f'Histogram of {data_column} by {categorical_column}')
27 plt.legend()
28 plt.show()

```



## Cluster analysis

```

1 df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 9217 entries, 0 to 9423
Data columns (total 23 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Title            9217 non-null    object  
 1   Genre            9217 non-null    object  
 2   Tags             9217 non-null    object  
 3   Languages        9217 non-null    object  
 4   Series or Movie  9217 non-null    object  
 5   Hidden Gem Score 9217 non-null    float64 
 6   Country Availability 9217 non-null    object  
 7   Runtime          9217 non-null    object  
 8   Director         9217 non-null    object  
 9   Writer           9217 non-null    object  
 10  Actors           9217 non-null    object  
 11  IMDb Score       9217 non-null    float64 
 12  Rotten Tomatoes Score 9217 non-null    float64 
 13  Awards Received  9217 non-null    float64 
 14  Awards Nominated For 9217 non-null    float64 
 15  IMDb Votes        9217 non-null    float64 
 16  Days_Since_Movie_Released 9217 non-null    int64

```

```

17 Release Year          9217 non-null  int64
18 Release Month         9217 non-null  int64
19 Days_Since_Movie_Released_on_Netflix 9217 non-null  int64
20 IMDb Votes Deviation 9217 non-null  float64
21 Weighted IMDb Score   9217 non-null  float64
22 Target                 9217 non-null  object
dtypes: float64(8), int64(4), object(11)
memory usage: 1.7+ MB

```

```
1 df_copy = df.copy()
```

```
1 df = df_copy.copy()
```

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

```

→ <class 'pandas.core.frame.DataFrame'>
Int64Index: 9217 entries, 0 to 9423
Data columns (total 23 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Title            9217 non-null    object  
 1   Genre             9217 non-null    object  
 2   Tags              9217 non-null    object  
 3   Languages         9217 non-null    object  
 4   Series or Movie   9217 non-null    object  
 5   Hidden Gem Score  9217 non-null    float64 
 6   Country Availability 9217 non-null  object  
 7   Runtime            9217 non-null    object  
 8   Director           9217 non-null    object  
 9   Writer             9217 non-null    object  
 10  Actors             9217 non-null    object  
 11  IMDb Score        9217 non-null    float64 
 12  Rotten Tomatoes Score 9217 non-null  float64 
 13  Awards Received   9217 non-null    float64 
 14  Awards Nominated For 9217 non-null  float64 
 15  IMDb Votes         9217 non-null    float64 
 16  Days_Since_Movie_Released 9217 non-null  int64  
 17  Release Year       9217 non-null    int64  
 18  Release Month      9217 non-null    int64  
 19  Days_Since_Movie_Released_on_Netflix 9217 non-null  int64  
 20  IMDb Votes Deviation 9217 non-null  float64 
 21  Weighted IMDb Score   9217 non-null  float64 
 22  Target              9217 non-null  object  
dtypes: float64(8), int64(4), object(11)
memory usage: 1.7+ MB

```

```
1 df_cluster = df.drop(['Genre', 'Languages', 'Series or Movie', 'Country Availability', 'Runtime',
2                         'Director', 'Writer', 'Actors', 'Tags', 'Title', 'Target'], axis=1)
```

```
1 df_cluster.info()
```

```

→ <class 'pandas.core.frame.DataFrame'>
Int64Index: 9217 entries, 0 to 9423
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Hidden Gem Score 9217 non-null    float64 
 1   IMDb Score        9217 non-null    float64 
 2   Rotten Tomatoes Score 9217 non-null  float64 
 3   Awards Received   9217 non-null    float64 
 4   Awards Nominated For 9217 non-null  float64 
 5   IMDb Votes         9217 non-null    float64 
 6   Days_Since_Movie_Released 9217 non-null  int64  
 7   Release Year       9217 non-null    int64  
 8   Release Month      9217 non-null    int64  
 9   Days_Since_Movie_Released_on_Netflix 9217 non-null  int64  
 10  IMDb Votes Deviation 9217 non-null  float64 
 11  Weighted IMDb Score   9217 non-null  float64 
dtypes: float64(8), int64(4)
memory usage: 936.1 KB

```

```

1 from sklearn.preprocessing import StandardScaler
2 from sklearn.cluster import KMeans
3
4 # Drop the columns that are not suitable for clustering (e.g., non-numeric columns)
5
6
7 # Standardize the numerical features

```

```

8 scaler = StandardScaler()
9 df_scaled = scaler.fit_transform(df_cluster)
10
11 # Apply K-Means clustering with k=3
12 kmeans = KMeans(n_clusters=3, random_state=42)
13 cluster_labels = kmeans.fit_predict(df_scaled)
14
15 # Add cluster labels back to the DataFrame
16 df['Cluster'] = cluster_labels
17
18 # Get the cluster centroids
19 cluster_centroids = pd.DataFrame(kmeans.cluster_centers_, columns=df_cluster.columns)
20
21 # Analyze the cluster centroids to identify important features
22 print("Cluster Centroids:")
23 # print(cluster_centroids)
24
25
26


```

1 cluster\_centroids

	Hidden Gem Score	IMDb Score	Rotten Tomatoes Score	Awards Received	Awards Nominated For	IMDb Votes	Days_Since_Movie_Released	Release Year	Release Month	Days_Since_Movie_R
0	-0.912693	-0.556360	-0.633800	-0.100028	-0.080789	0.171243		0.676031	-0.675439	-0.009020
1	0.532708	0.225876	0.269915	-0.150893	-0.182679	-0.338700		-0.383218	0.383205	-0.008167
2	-0.615271	0.964709	0.928227	2.768927	3.060504	3.351275		0.301794	-0.305903	0.175956

1 Start coding or generate with AI.

```

1 # Look at the distribution of features within each cluster
2 for i in range(3):
3     print(f"Cluster {i}:")
4     print(df[df['Cluster'] == i].describe(include='all'))


```



	IMDb	Votes	Deviation	Weighted	IMDb	Score	Target	Cluster
mean	2006.602837	7.304965			1678.959811			
std	10.276887	3.593589			713.728277			
min	1960.000000	1.000000			0.000000			
25%	2002.000000	5.000000			1374.000000			
50%	2009.000000	7.000000			2151.000000			
75%	2014.000000	11.000000			2151.000000			
max	2020.000000	12.000000			2151.000000			

	IMDb	Votes	Deviation	Weighted	IMDb	Score	Target	Cluster
count		423.000000		423.000000		423	423.0	
unique		Nan		Nan		2	Nan	
top		Nan		Nan	Super	Hit	Nan	
freq		Nan		Nan		399	Nan	
mean	53.696806		436.581293			Nan	2.0	
std	37.770202		332.925967			Nan	0.0	
min	-5.351473		-35.246932			Nan	2.0	
25%	29.746010		235.826651			Nan	2.0	
50%	46.471571		359.182037			Nan	2.0	
75%	68.679665		511.312742			Nan	2.0	
max	248.755612		2322.727188			Nan	2.0	

[11 rows x 24 columns]

```
1 df.drop(['Title'],axis=1,inplace=True)
```

```
1 # Analyze each cluster
2 for i in range(df['Cluster'].nunique()):
3     print(f"Cluster {i}:")

4
5     # Identify Patterns and Similarities
6     print("Common properties:")
7     print(df[df['Cluster'] == i].mode().iloc[0]) # prints the mode of each feature within the cluster
8
9     # Feature Importance
10    # In KMeans, all features contribute equally to the distance calculation. However, you can look at the mode or mean value
11
12    # Understanding Audience Segments
13    print("Audience segment properties:")
14    df[df['Cluster'] == i].describe() # prints summary statistics of each feature within the cluster
```

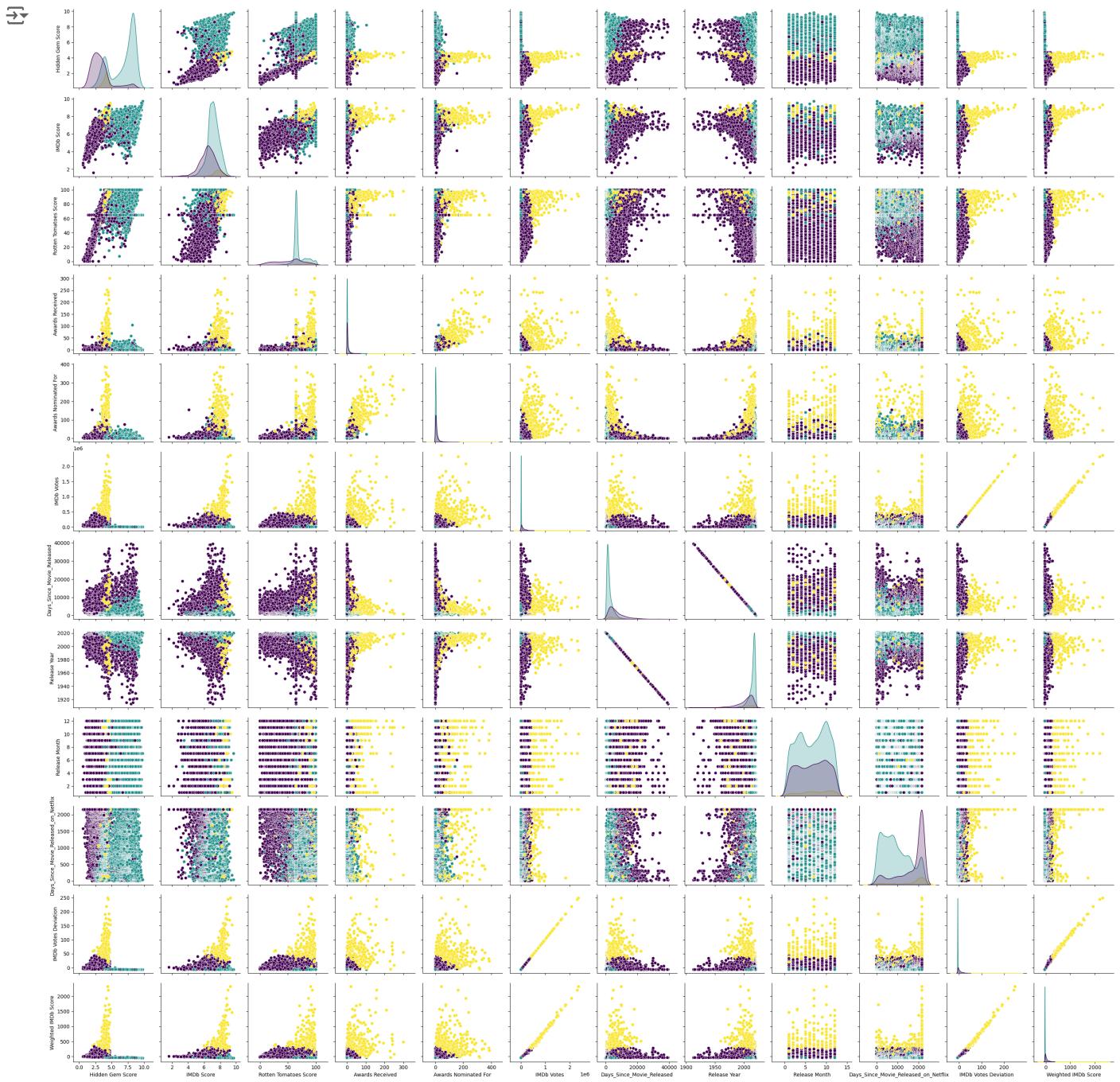
	Not Hit
Target	0
Cluster	
Name: 0, dtype: object	
Audience segment properties:	
Cluster 1:	
Common properties:	
Genre	Comedy
Tags	Dramas
Languages	English
Series or Movie	Movie
Hidden Gem Score	8.3
Country Availability	Japan
Runtime	1-2 hour
Director	Unknown
Writer	Unknown
Actors	Jackie Chan
IMDb Score	6.8
Rotten Tomatoes Score	64.691276

1/9/25, 12:02 AM

Netflix\_project.ipynb - Colab

```
writer          George Lucas
Actors          Ewan McGregor
IMDb Score      7.6
Rotten Tomatoes Score 64.691276
Awards Received   26.0
Awards Nominated For 69.0
IMDb Votes        733336.0
Days_Since_Movie_Released 8213
Release Year      2013
Release Month     12
Days_Since_Movie_Released_on_Netflix 2151
IMDb Votes Deviation 72.900019
Weighted IMDb Score 480.350121
Target           Super Hit
Cluster          2
Name: 0, dtype: object
Audience segment properties:
```

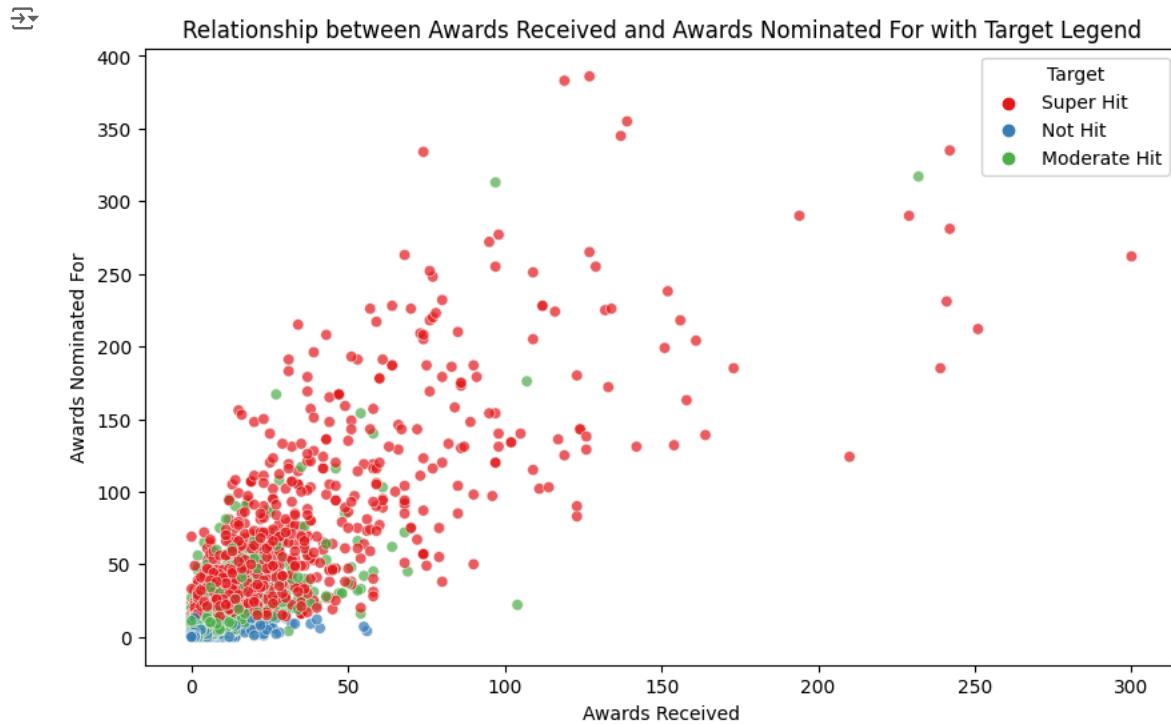
```
1 # Visualize the clusters
2 sns.pairplot(df, hue='Cluster', diag_kind='kde', palette='viridis')
3 plt.show()
```



```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Scatter plot with legend based on 'Target' variable
5 plt.figure(figsize=(10, 6))
6 sns.scatterplot(x='Awards Received', y='Awards Nominated For', hue='Target', data=df, palette='Set1', alpha=0.7)
7 plt.xlabel('Awards Received')
8 plt.ylabel('Awards Nominated For')
9 plt.title('Relationship between Awards Received and Awards Nominated For with Target Legend')
10 plt.legend(title='Target', loc='upper right')
11 plt.show()
12

```



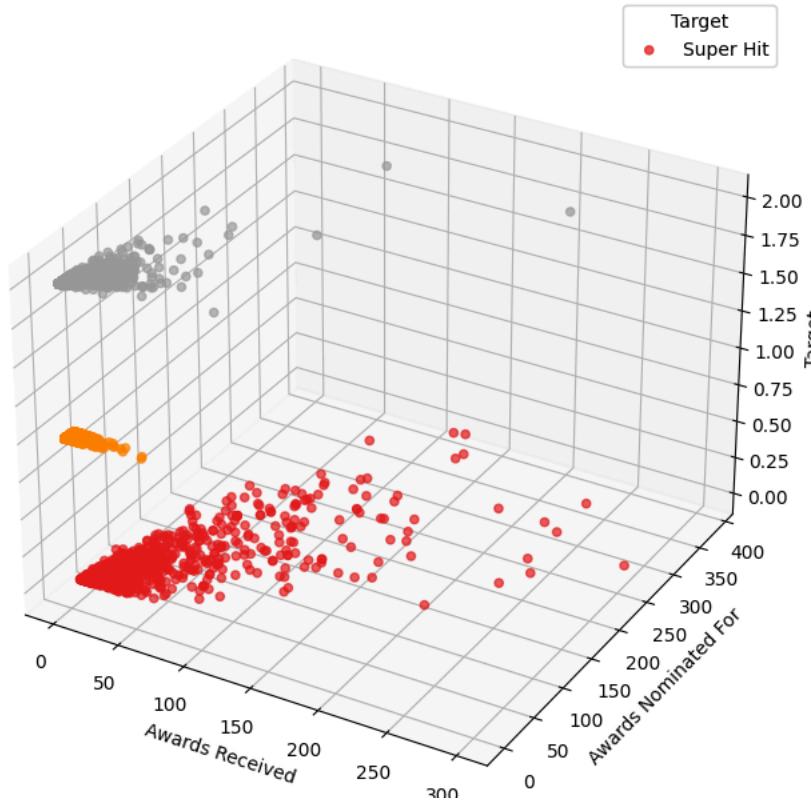
```

1 import matplotlib.pyplot as plt
2 from mpl_toolkits.mplot3d import Axes3D
3 import pandas as pd
4
5 # Convert 'Target' column to numeric category codes
6 df['Target_code'] = pd.factorize(df['Target'])[0]
7
8 # Create a mapping of numeric category codes to 'Target' labels
9 target_mapping = {code: label for code, label in zip(df['Target_code'], df['Target'])}
10
11 # Create a 3D scatter plot
12 fig = plt.figure(figsize=(10, 8))
13 ax = fig.add_subplot(111, projection='3d')
14
15 # Scatter plot with color based on 'Target' variable
16 scatter = ax.scatter(df['Awards Received'], df['Awards Nominated For'], df['Target_code'], c=df['Target_code'], cmap='Set1',
17
18 # Set labels for x, y, and z axes
19 ax.set_xlabel('Awards Received')
20 ax.set_ylabel('Awards Nominated For')
21 ax.set_zlabel('Target')
22
23 # Add legend with 'Target' labels
24 target_labels = [target_mapping[code] for code in sorted(df['Target_code'].unique())]
25 legend = ax.legend(*scatter.legend_elements(), labels=target_labels, title='Target', loc='upper right')
26 ax.add_artist(legend)
27
28 # Show the 3D scatter plot
29 plt.title('3D Scatter Plot: Awards Received vs. Awards Nominated For with Target')
30 plt.show()
31

```

```
→ <ipython-input-59-22f0547d105e>:25: UserWarning: You have mixed positional and keyword arguments, some input may be discarded
  legend = ax.legend(*scatter.legend_elements(), labels=target_labels, title='Target', loc='upper right')
```

### 3D Scatter Plot: Awards Received vs. Awards Nominated For with Target



## Predictive analysis

### ✓ Logistic regression with all columns

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

```
→ <class 'pandas.core.frame.DataFrame'>
Int64Index: 9217 entries, 0 to 9423
Data columns (total 24 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Genre            9217 non-null    object  
 1   Tags             9217 non-null    object  
 2   Languages        9217 non-null    object  
 3   Series or Movie  9217 non-null    object  
 4   Hidden Gem Score 9217 non-null    float64 
 5   Country Availability 9217 non-null    object  
 6   Runtime          9217 non-null    object  
 7   Director         9217 non-null    object  
 8   Writer           9217 non-null    object  
 9   Actors           9217 non-null    object  
 10  IMDb Score       9217 non-null    float64 
 11  Rotten Tomatoes Score 9217 non-null    float64 
 12  A               9217 non-null    object  
 13  B               9217 non-null    object  
 14  C               9217 non-null    object  
 15  D               9217 non-null    object  
 16  E               9217 non-null    object  
 17  F               9217 non-null    object  
 18  G               9217 non-null    object  
 19  H               9217 non-null    object  
 20  I               9217 non-null    object  
 21  J               9217 non-null    object  
 22  K               9217 non-null    object  
 23  L               9217 non-null    object 
```