

NETFLIX MOVIE RECOMMENDATION SYSTEM



Table of content

- Abstract
- Description of Dataset
- Data Exploration
- Data Cleaning
- Data visualization
- Recommendation System

Abstract

My project will be the Netflix movies and TV Show dataset, Netflix is an application that keeps growing exponentially whole around the world and it is the most famous streaming platform. It given a large number of movies and series available on the platform, it is a perfect opportunity to flex our data manipulation skills and dive into the entertainment industry.

The goal of this project was to Recommends content based on movie description. Here I would recommend Movie based on Movie titles only. Similar movies would have similar names thus having a high cosine similarity. Then the movies that are most likely to be similar are recommended.

I generates TF-IDF matrix and finds cosine similarity of each movie with other movies and displays top 10 similar movies, and to analysis and visualization of content is available in different countries, analysis of Actors / Directors and find interesting insights.

Dataset

This dataset contains data collected from Netflix of different TV shows and movies from the year 2008 to 2021.

- type: Gives information about 2 different unique values one is TV Show and another is Movie
- title: Gives information about the title of Movie or TV Show
- director: Gives information about the director who directed the Movie or TV Show
- cast: Gives information about the cast who plays role in Movie or TV Show
- release_year: Gives information about the year when Movie or TV Show was released
- rating: Gives information about the Movie or TV Show are in which category (eg like the movies are only for students, or adults, etc)
- duration: Gives information about the duration of Movie or TV Show
- listed_in: Gives information about the genre of Movie or TV Show
- description: Gives information about the description of Movie or TV Show

EDA

Importing library

```
import pandas as pd
import numpy as np
import plotly.graph_objects as go
import matplotlib.pyplot as plt
import plotly.graph_objects as go
import pandas as pd
import seaborn as sns
```

See the Netflix data:

```
In [4]: # imports a CSV file (netflix_titles) to DataFrame format
df = pd.read_csv("netflix_titles.csv")
df.head(10)
```

Out[4]:

	show_id	type	title	director	cast	country	date_added	release_year	rating	duration	listed_in	description
0	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG- 13	90 min	Documentaries	As her father nears the end of

Data Exploration

Exploring the data

```
Data Exploration

In [5]: #print a summary of a DataFrame
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8807 entries, 0 to 8806
Data columns (total 12 columns):
#   Column             Non-Null Count  Dtype
---  --
0   show_id             8807 non-null   object
1   type                8807 non-null   object
2   title               8807 non-null   object
```

Finding if the dataset contains null values

```
In [6]: #find the nulls value
df.isnull().sum()

Out[6]: show_id      0
        type         0
        title        0
        director    2634
        cast        825
        country     831
        date_added   10
```

Finding how many unique values are there in the dataset

```
In [7]: #find unique value
df.nunique()

Out[7]: show_id      8807
        type         2
        title      8807
        director    4528
        cast       7692
        country     748
        date_added  1767
        release_year 74
        rating      17
```

Data Cleaning

Replacing null values with 'NULLs'

```
Data Cleaning

In [8]: #Replace null values with Null word
df['country'].fillna('Null',inplace=True)
df['rating'].fillna('Null',inplace=True)
df.isnull().sum().sum()

Out[8]: 3472
```

Remove missing value

```
df['rating'].fillna('Null',inplace=True)
df.isnull().sum().sum()

Out[8]: 3472

In [56]: df = df.dropna( how='any',subset=['cast', 'director'])

In [57]: df.isnull().sum()
```

Converting into a proper date-time format

```
dtype: int64

In [10]: #Converting into date-time format and adding two more features year and month.

df['date_added'] = pd.to_datetime(df['date_added'])
df['year_added'] = df['date_added'].dt.year
df['month_added'] = df['date_added'].dt.month

In [65]: #Correlation between the features
```

Data Visualization

Chart 1

Viewing the correlation between the features.

```
In [65]: #Correlation between the features
month_year_df = df.groupby('year_added')['month_added'].value_counts().unstack().fillna(0).T

plt.figure(figsize=(11,8))
sns.heatmap(month_year_df, linewidths=0.025, cmap="Greens")
plt.title("Content Heatmap")
plt.ylabel("Month")
plt.xlabel("Year")
plt.show()
```



Chart 2

Movie duration over years

```
In [12]: # Create the years and durations Lists
Year = [2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020]
Duration = [103, 101, 99, 100, 100, 95, 95, 96, 93, 90]

# Create a dictionary with the two Lists
movie_dict = {'years': Year, 'durations': Duration}

line = plt.plot(Year, Duration , linewidth = 1 , color='g')

plt.title('Netflix Movie Durations 2011-2020',fontsize = 20,style='italic',weight='bold',rotation=0,color='purple');
plt.xlabel('YEARS',fontsize = 8, weight = 'bold',color='green');
plt.ylabel('MOVIE DURATIONS',fontsize = 8, weight = 'bold',color='green');

plt.show()
```

Chart3

Pie chart to show the percentage of content type

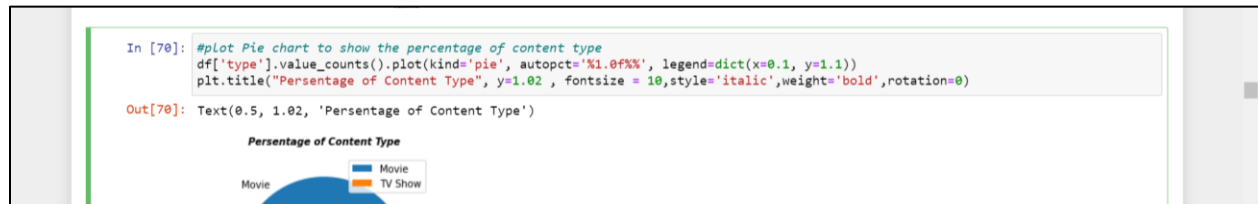


Chart 4

Bar charts of the number of Movies per Country.(descending)



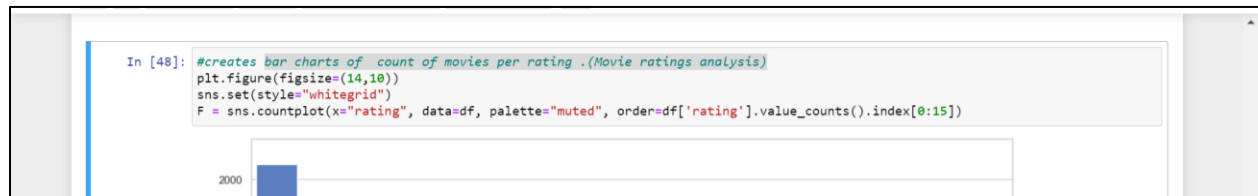
Chart 5

The content added for Movies compared to TV Shows over years



Chart 6

Bar charts of count of movies per rating .(Movie ratings analysis)



Recommendation System (TF-IDF)

CONTENT-BASED FILTERING



```
Recommendation System

In [59]: from sklearn.feature_extraction.text import TfidfVectorizer

#removing stopwords
tf = TfidfVectorizer(stop_words='english')

#Construct the required TF-IDF matrix by fitting and transforming
# the data
tf_matrix = tf.fit_transform(df['description'])

#output the shape of tfidf_matrix
tf_matrix.shape

Out[59]: (5522, 14500)

In [61]: #Linear Kernel
from sklearn.metrics.pairwise import linear_kernel

#Cosine similarity
cosine_sim = linear_kernel(tf_matrix, tf_matrix)
```

There are about 14500 words described for the 5522 movies in this dataset.

```
In [61]: #Linear Kernel
from sklearn.metrics.pairwise import linear_kernel

#Cosine similarity
cosine_sim = linear_kernel(tf_matrix, tf_matrix)

In [60]: indices = pd.Series(df.index, index = df['title']).drop_duplicates()

In [62]: indices
Out[62]: title
My Little Pony: A New Generation    1
Sankofa                             2
The Starling                       3
Je Suis Karl                       4
Confessions of an Invisible Girl    5
...
```

- Initialize Tfidf vectorizer & fit into title column
- Calculate cosine similarity
- Convert all titles into a Series associated with movie index numbers
- Function that gets movie recommendations based on the cosine similarity score of movie titles

Tools

- Numpy and Pandas for data manipulation
- Matplotlib and Seaborn for plotting