

Data Analysis: Netflix

1. Aims, objectives and background

1.1 Introduction

Data analysis has been used by many companies for a long time. With the main objective of looking for trends and answering questions using data visualization techniques.

As one of the top streaming service in 2021, Netflix has been a staple to many households. From TV shows to movies, there are over 17,000 titles worldwide. As an avid fan of Netflix, I have spent countless hours binging on the latest hits. Using this opportunity, I want to explore Netflix and its content.

Being a huge company, analysing its data will allow for informed decision-making based on facts. Reduce costs by looking for improvement opportunities, trends, and patterns in their data and plan their strategies accordingly, provide a better user recommendations by analysing their interest and behaviours.

1.2 Aims and objectives

Through this project, I would like to find out:

- How much content there actually is
- When did Netflix peaked
- Look for any trends and relationships

The aims for this project proposal are to:

- Find a dataset that allows for data cleaning and analysis
- Clean the data to make it usable
- Visualizing the data using graph plotting techniques
- Analyzing the data for trends and answering questions

1.3 Data

1.3.1 Data acquisition and reliability

After doing some research, I have found out the official Netflix API has been shutdown in 2014. This lead me to searching for datasets that are available on the internet. Which brought me to Kaggle, an online community platform for data scientists, that allows users to collaborate with other users, find and publish datasets.

Based on the posting, the dataset was regularly updated up till September 2021 and based on the usability rating it scored a 10/10, showing this can be a trusted and reliable dataset.

1.3.2 Data type and suitability

The dataset is in the CSV format allowing it to be easily used with Pandas. It consists of the main information such as title, country, data added etc, this allows us to answer the questions that I have posed.

1.3.3 Time period

The years provided in the dataset span from 1925 – 2021. However, having only a small quantity of titles (<10) in that particular year, it would not make a significant impact on the changes of the data. Hence, we will be focusing mainly on 2000 – 2021.

1.4 Ethical consideration

1.4.1 copyrights

Based on Kaggle's terms, I am allowed to use it as long as it is for personal and non-commercial. For specific copyrights, they will be displayed on the dataset page.

After looking through, I found the license on the dataset's page, CC0: Public Domain, which states "The person who associated a work with this deed has dedicated the work to the public domain by waiving all of his or her rights to the work worldwide under copyright law, including all related and neighboring rights, to the extent allowed by law. You can copy, modify, distribute and perform the work, even for commercial purposes, all without asking permission." Thus, the data is safe to use.

1.4.2 Reusage of data

The legality of the dataset has been confirmed as of writing. Anyone planning on using the dataset provided may refer to the link in the references for the terms and conditions to prevent any infringement of the law.

1.4.2 Potential impacts

As this project is used for personal and educational use, any conclusions from this project will be based on assumptions after looking at the graphs and should not be used as a professional analytical judgement.

2. Importing Data and libraries

importing the required libraries for data visualization

```
In [1]: import os
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
import plotly.graph_objects as pgo
```

Applying error handling methods, in case a file is not available.

```
In [2]: try:
        netflix = pd.read_csv('netflix_titles.csv')
    except FileNotFoundError:
        print("File does not exist")
```

Original dataset with missing values

```
In [3]: netflix
```

```
Out [3]:
```

	show_id	type	click	director	cast	country	date_added	release_year	rating	duration	listed_in	descripti
0	s1	Movie	Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG-13	90 min	Documentaries	As I father the end li film
1	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabulane, Thabane...	South Africa	September 24, 2021	2021	TV-MA	2 Seasons	International TV Shows, TV Dramas, TV Mysteries	AF cross paths a party Cape T
2	s3	TV Show	Ganglands	Julien Leclercq	Samii Bouajila, Tracy Goudas, Sami Joui, Nab...	NaN	September 24, 2021	2021	TV-MA	1 Season	Crime TV Shows, Korean TV Shows, TV Act...	To prot his fan from power drug
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA	1 Season	Docueries, Reality TV	Fe flirtatious and tal down ank
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA	2 Seasons	International TV Shows, TV Shows, TV ...	In a cty cent cent known
...
8802	s8803	Movie	Zodiac	David Fincher	Mark Ruffalo, Jake Gyllenhaal, Robert Downey J...	United States	November 20, 2019	2007	R	158 min	Cult Movies, Dramas, Thrillers	A politi cartooni a crim repor and i
8803	s8804	TV Show	Zombie Dumb	NaN	NaN	NaN	July 1, 2019	2018	TV-Y7	2 Seasons	Kids' TV Shows, TV Comedies	While livi alone i life form town sur
8804	s8805	Movie	Zombieland	Ruben Fleischer	Eisenberg, Jesse Harnesse, Woody Ha...	United States	November 1, 2019	2009	R	88 min	Comedies, Horror Movies	Looking survi world tak over by z
8805	s8806	Movie	Zoom	Peter Hewitt	Tim Allen, Courteney Cox, Chevy Chase, Kate Ma...	United States	January 11, 2020	2006	PG	88 min	Children & Family Movies, Comedies	Dragg from civil life form superher
8806	s8807	Movie	Zubaan	Mozez Singh	Vicky Kaushal, Sarah-Jane Dias, Raaghav Chan...	India	March 2, 2019	2015	TV-14	111 min	Dramas, Movies, Music & Musicals	A scrap boy work way ir a t

8807 rows x 12 columns

3. Data cleaning and processing

Before plotting any graphs, I need to ensure that the dataset does not have any missing data values by:

- scanning through the data set and calculating the percentage of missing values
- Filling up missing data with as much relevant details
- Removing missing data will be dropped

3.1 Calculation of missing data in percentage

```
In [4]: for i in netflix.columns:
        nullPercentage = netflix[i].isna().sum() / len(netflix)*100
        if nullPercentage > 0 :
            print(f"{i} NaN rate is {round(nullPercentage,2)}%")
```

director NaN rate is 29.91%

country NaN rate is 9.44%

date_added NaN rate is 0.11%

rating NaN rate is 0.05%

duration NaN rate is 0.03%

There are 6 columns that have missing values.

- Director and cast's information is missing a big portion, although it is not important for the purpose of data analysis, it is vital so that it won't not be dropped later on.
- Country's information is vital for our data analysis, hence, it is vital to fill it up with the most repeated country. This will not impact the data analysis results later on.
- date_added, rating and duration will be dropped as their amount will have little impact on the dataset.

```
In [5]: netflix["cast"].fillna("No Data", inplace=True)
netflix["director"].fillna("No Data", inplace=True)
netflix["country"].fillna(netflix["country"].mode()[0])
netflix.dropna(inplace=True)
netflix.drop_duplicates(inplace=True) #drop duplicate rows
```

Double checking all missing values has been filled or removed.

```
In [6]: for i in netflix.columns:
        nullPercentage = netflix[i].isna().sum() / len(netflix)*100
        if nullPercentage > 0 :
            print(f"{i} NaN rate is {round(nullPercentage,2)}%")
        else:
            print("no missing values")
```

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

no missing values

