

EDA ON NETFLIX MOVIES AND TV SHOWS DATASET

Presented By : Harsh Mishra

Guide By : Naina Devi

College : Techno India University



AGENDA:

- To understand the composition and trends of Netflix's content library.
- To analyze movies and TV shows based on country, genre, release year, and ratings.
- To visualize data insights through graphs and charts.
- To identify patterns in Netflix's content strategy



INTRODUCTION :

This project performs **Exploratory Data Analysis (EDA)** on Netflix's Movies and TV Shows dataset. The dataset includes information like title, director, cast, country, release year, rating, duration, and type of content. The goal is to gain meaningful insights into Netflix's content distribution and trends. Through data cleaning, visualization, and analysis, the project helps understand what kind of content Netflix produces most, how it has evolved over time, and which countries contribute the most to its library.

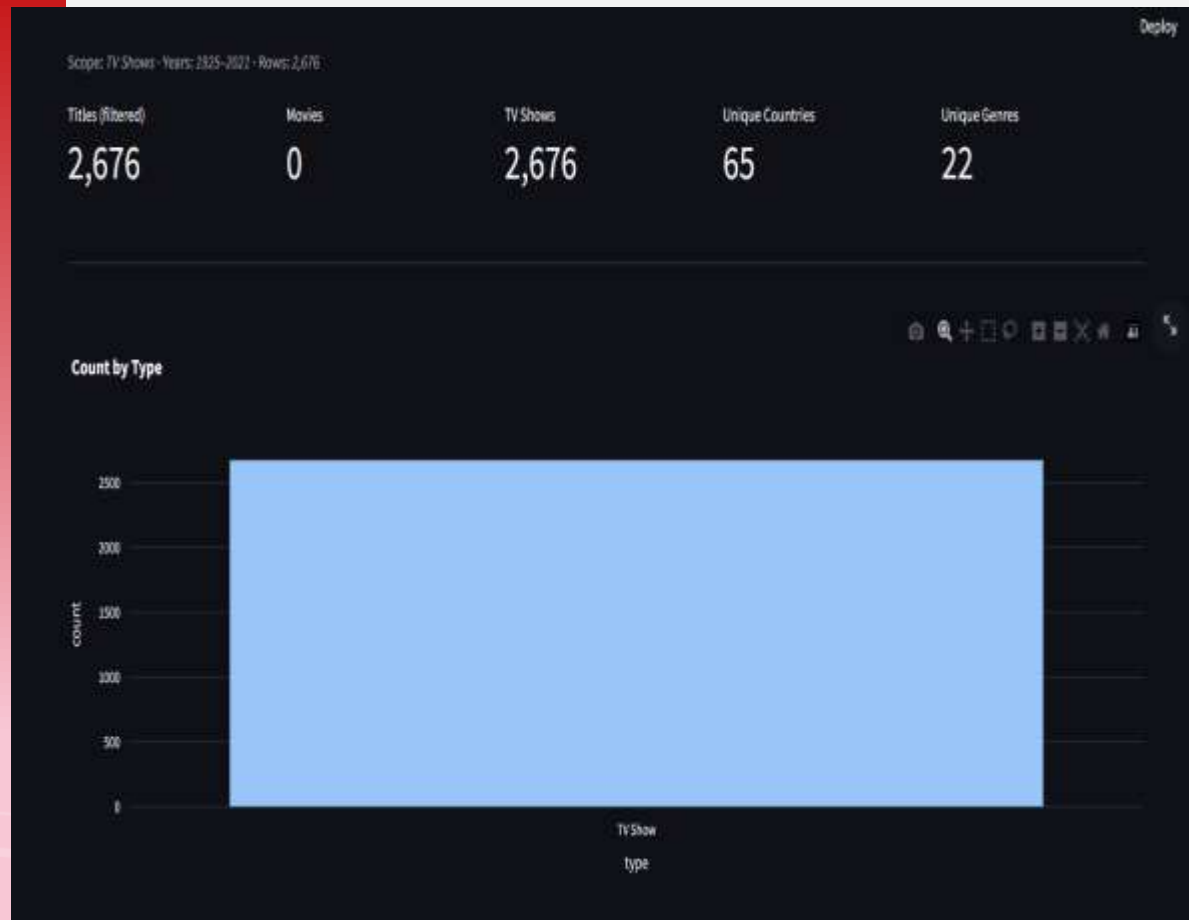


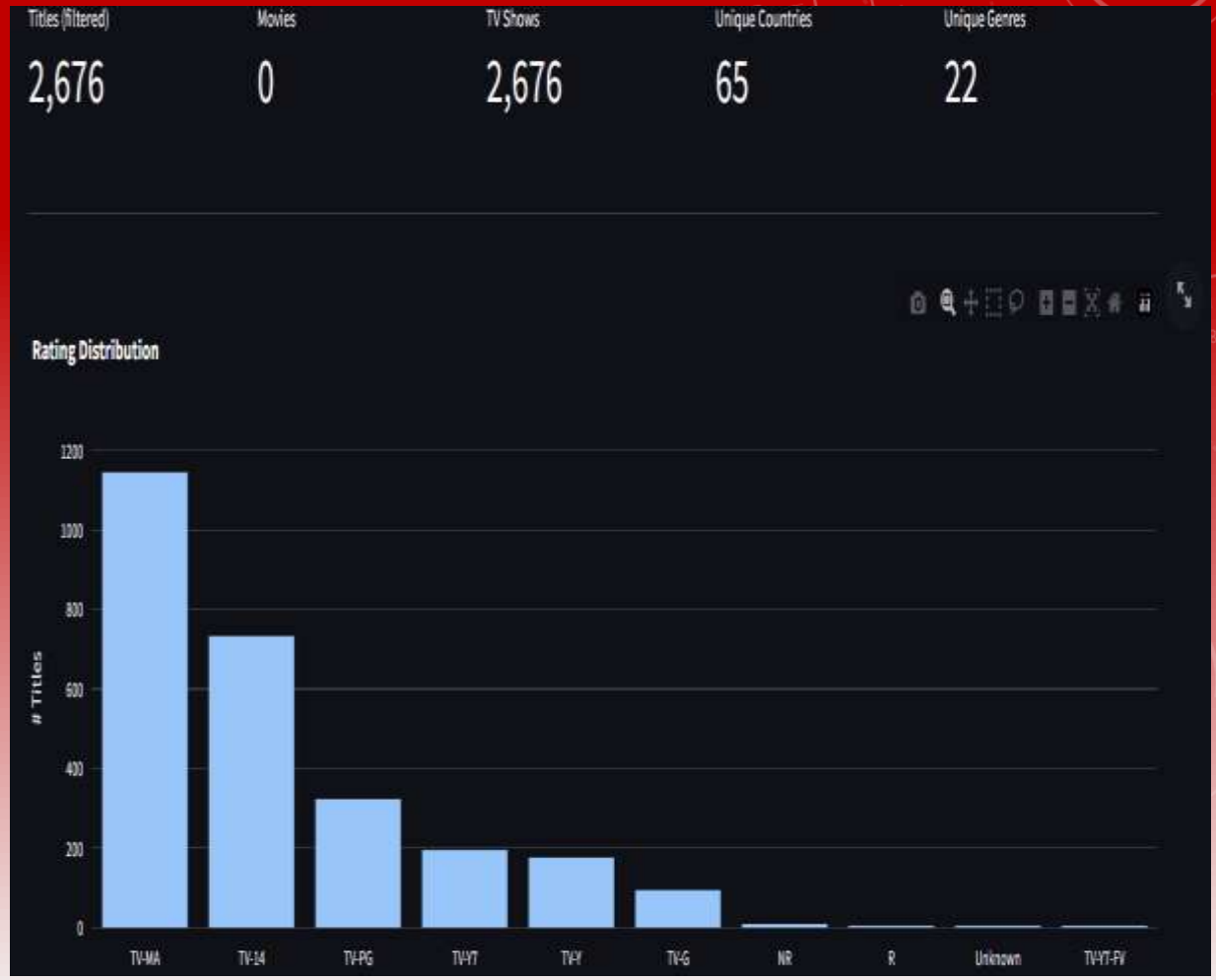
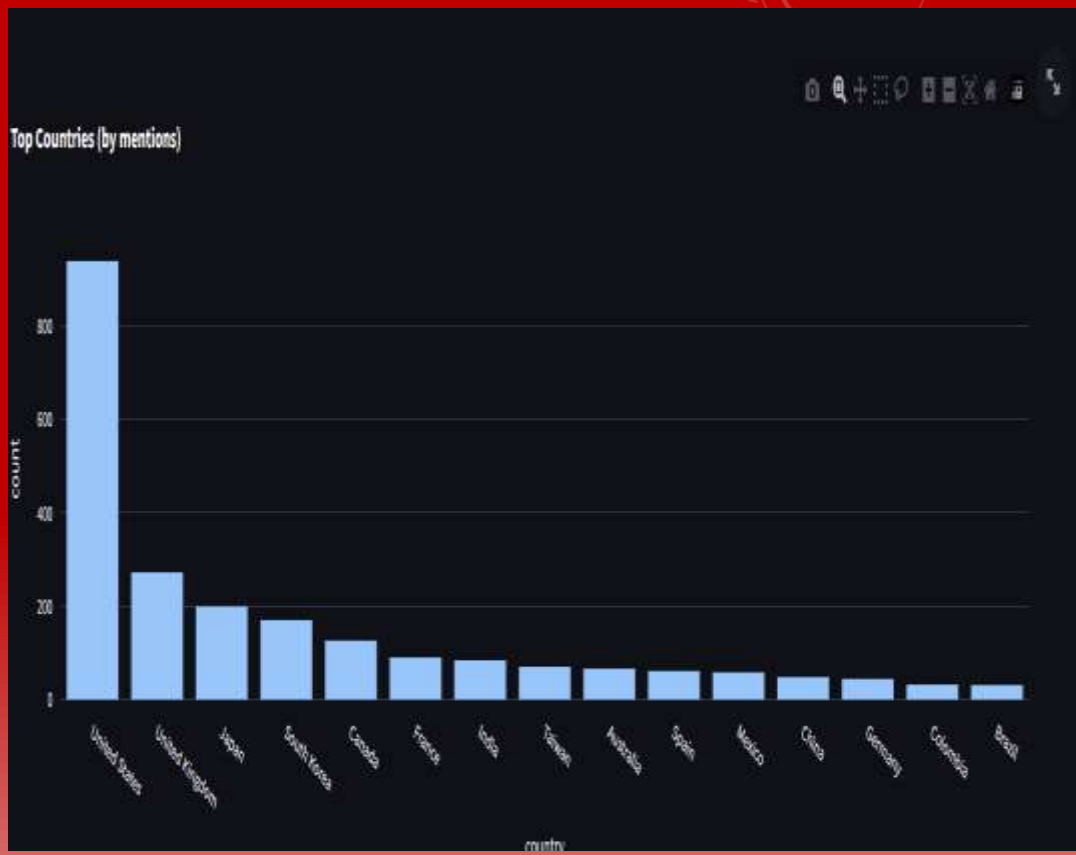
LIBRARIES AND TECHNOLOGIES USED

- **Programming Language:** Python
- **Libraries Used:**
 - **Numpy** for data manipulation
 - **Matplotlib**
 - **Sckit learn** for predictive modeling
 - **Joblib** for model serialization
- **Deployment and Interface:**
 - **Streamlit** for rapid development
 - **Render** for cloud deployment
- **Dataset Source:**
 - Kaggle
 - Mymoviedb.csv provides rich features for Netflix shows EDA
 - Development in **Jupyter Notebook** and version control with **GitHub**

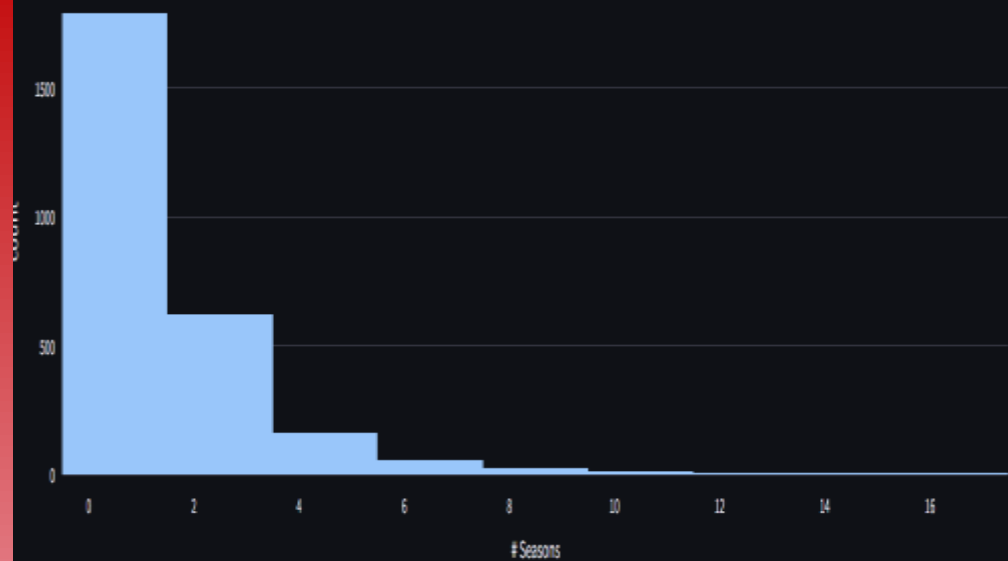


PROJECT WORKING:

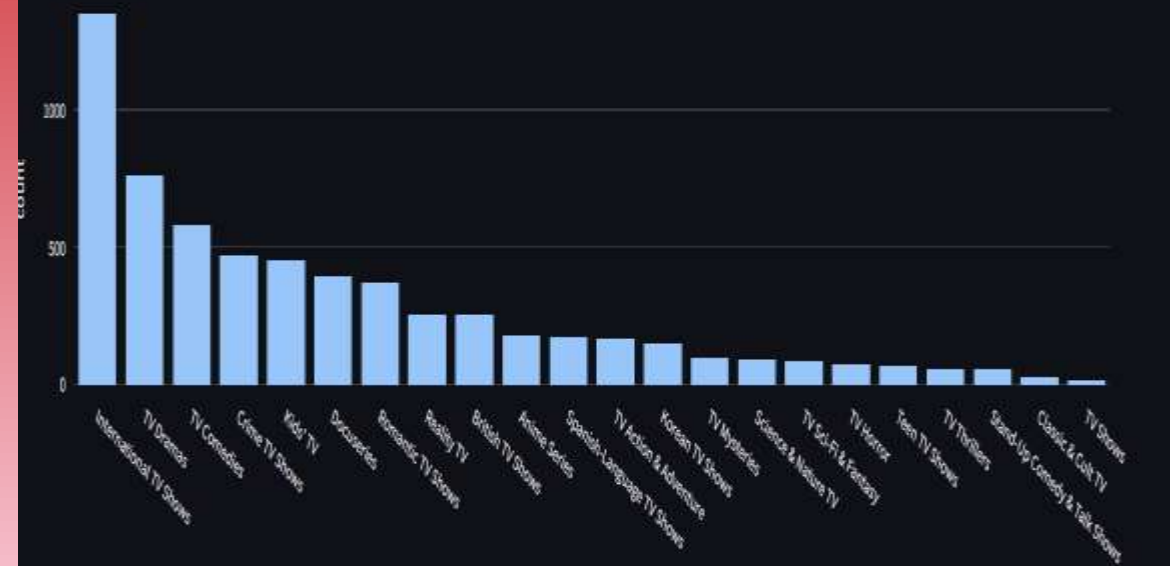




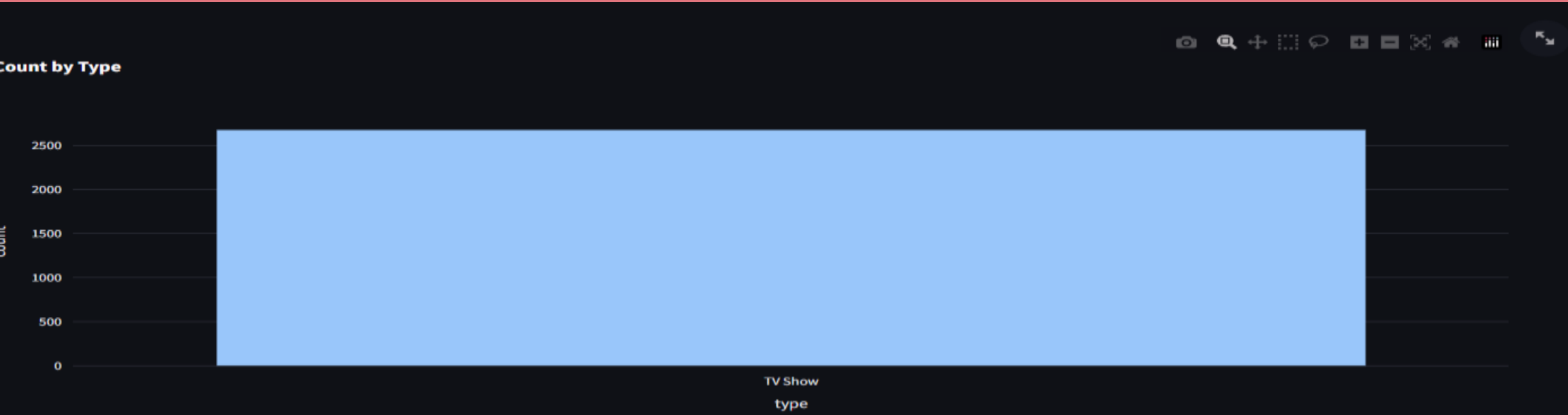
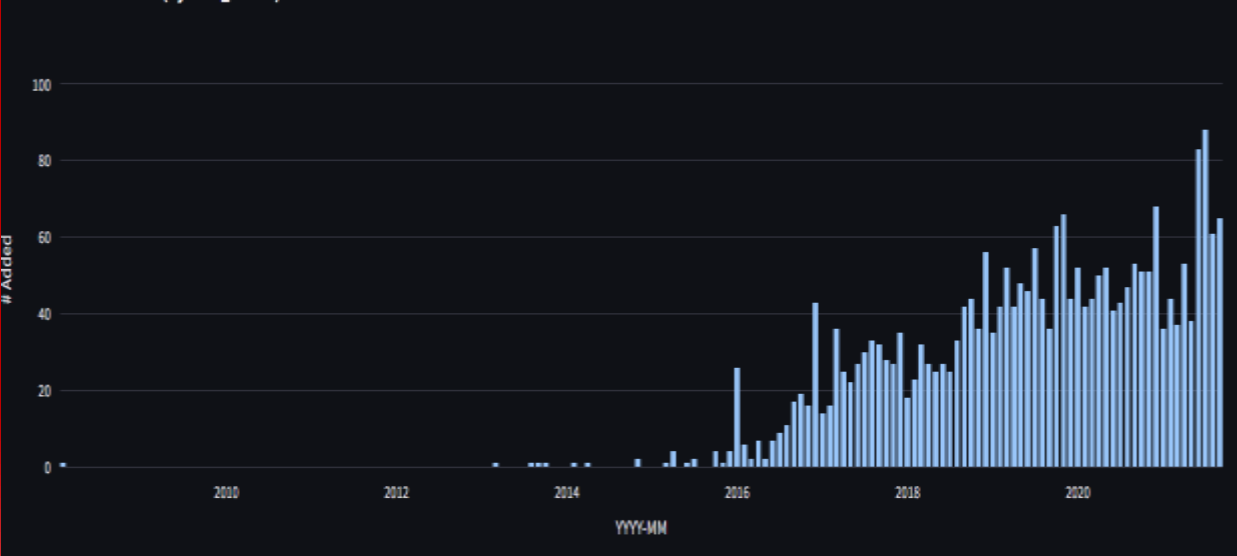
TV Show Seasons Histogram



Top Genres



Month-wise Additions (by date_added)



APPLICATIONS:

- Helps Netflix and other streaming services understand global content trends.
- Useful for data analysts and entertainment researchers.
- Supports decision-making in content acquisition and production.
- Educational use for learning data visualization and analytics.



PROBLEM FACED AND ITS SOLUTIONS:

CHALLENGES

Missing Values:

Many entries in *director*, *cast*, and *country* columns were missing.

Duplicate Records:

Some titles appeared more than once in the dataset.

Unformatted Dates and Durations:

The *date_added* and *duration* columns had inconsistent formats.

Cluttered Visualizations:

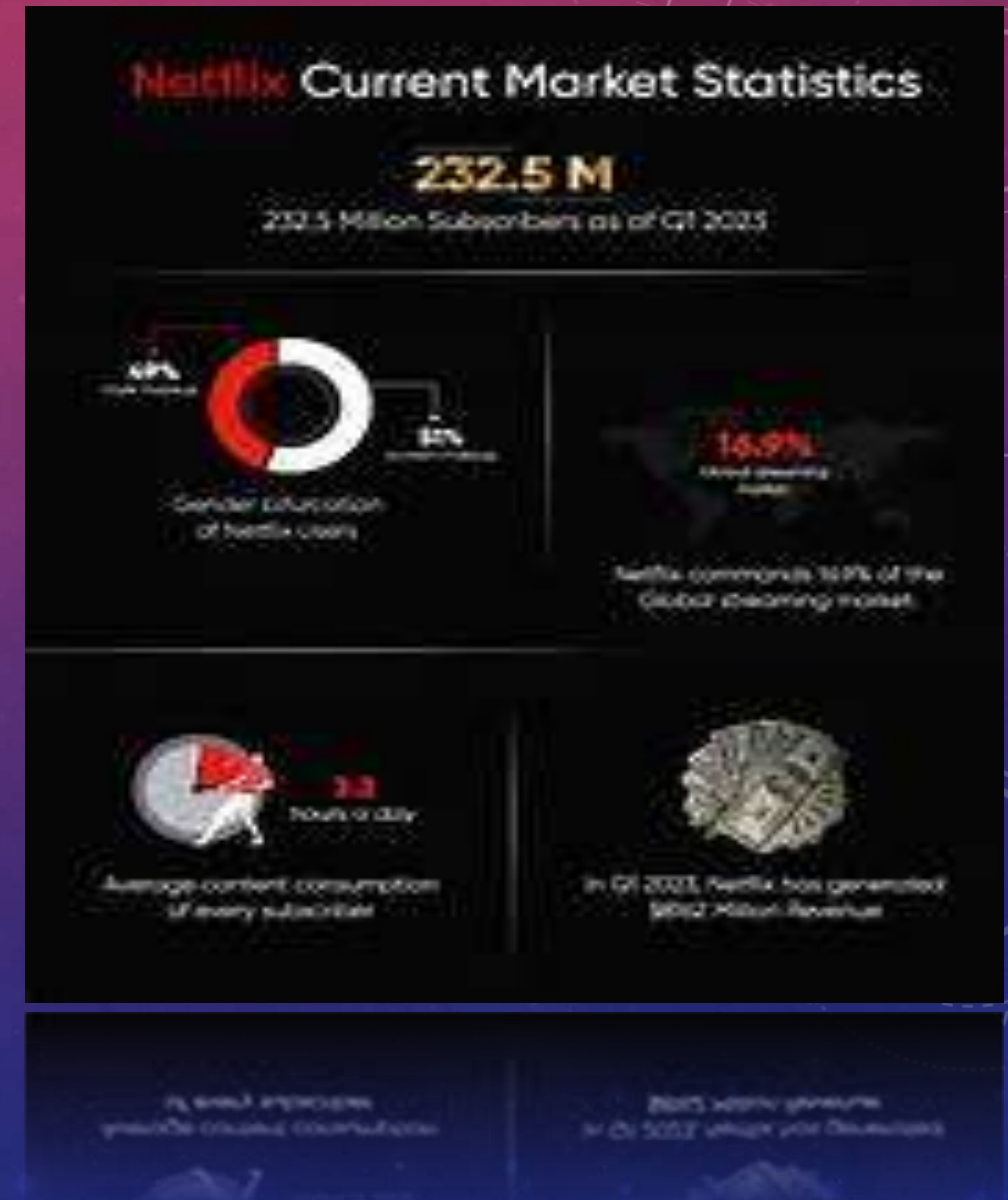
Graphs were overcrowded with too many categories.

SOLUTIONS

- Filled missing values with “Unknown” or removed incomplete rows using Pandas.
- Used `drop_duplicates()` to clean the data and ensure accuracy.
- Converted them to proper formats using `pd.to_datetime()` and string operations.
- Displayed only top 10 values for better clarity and presentation

FUTURE SCOPE:

- Extend analysis using machine learning to predict content popularity.
- Include IMDb ratings and viewer engagement data for deeper insights.
- Create an interactive dashboard using Power BI or Tableau.
- Automate updates with Netflix's latest data through APIs.



PROJECT LINK :

- ❑ <https://github.com/MishraHarsh25/EDA-on-Netflix>
- ❑ <https://eda-on-netflix-sgse.onrender.com/>

THANK YOU