# Netflix Movies and TV Shows

Haritha Kolli
Instructor: Gahangir Hossain
INFO 5709 DATA VISUALIZATION AND COMMUNICATION
PROJECT PAPER
12/16/2022

## Introduction:

One of the most well-known media and video streaming services is Netflix. They offer more than 8000 movies and TV episodes on their platform, and as of the middle of 2021, they had more than 200 million subscribers worldwide.

This tabular dataset includes listings for all of the Netflix movies and TV episodes, together with information about the actors, directors, ratings, release year, duration, and other factors.

## Dataset:

The Dataset contains 12 features about the TV Shows and Movie content, which is available on the Streaming platform, NETFLIX. This data has been taken from Kaggle.

Dataset Kaggle Link : https://www.kaggle.com/datasets/shivamb/netflix-shows

We have 8807 data entries with the below list of column data:

1. show_id : It's a unique identifier for each movie and TV Show.
2. type : Indicates if it is a TV show or Movie
3. title : Title of the TV Show or Movie
4. director : Director of the TV Show or Movie
5. cast : Actors casted in the TV show or Movie
6. country : Country, which produced the movie or the TV show
7. date_added: Date when the show/ movie been added to the Netflix
8. release_year: When the TV Show or movie got released
9. rating : rating for the TV Show or movie
10. duration : TV show or movie duration in min / number of seasons
11. listed_in : Indicates the genere
12. description : TV Show or Movie summary

```
In [270]: data_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
 3   director      6173 non-null   object
 4   cast          7982 non-null   object
 5   country       7976 non-null   object
 6   date_added    8797 non-null   object
 7   release_year  8807 non-null   int64
 8   rating        8803 non-null   object
 9   duration      8804 non-null   object
 10  listed_in     8807 non-null   object
 11  description   8807 non-null   object
dtypes: int64(1), object(11)
memory usage: 825.8+ KB
```
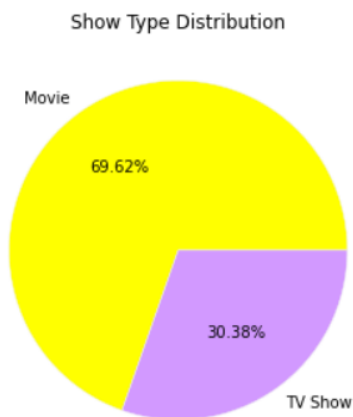
## Tools:

- Python
- Jupyter Notebook
- Libraries Used: pandas, NumPy, datascience, matplotlib, seaborn
- Tableau

## Exploratory Data Analysis:

❖ The below pie chart shows that We have 69.62% of Movie Content and 30.38% of TV Show content available on Netflix based on the data available at hand.

```python
show_type = data_df['type'].value_counts().reset_index()

plt.figure(figsize=(12,5))

plt.pie(show_type['type'], labels=show_type['index'], autopct='%0.2f%%', colors=['yellow','#D299FF'])

plt.title("Show Type Distribution", fontsize=12)

plt.show()
```



Insights: DataSet contains `69.62%` of Movie's data and `30.38%` of TV Show

- Missing Content Check: Created a function 'null_value_check()' to check the null value in the Dataset. We found that 4.08% null content in the data.

```python
def null_value_check(df):

    missing_content = round((df.isnull().sum()/len(df)*100),2)
    print("Missing Content in the dataset:\n",missing_content[missing_content>0])
    print("Total missing %: ",round(missing_content.mean(),2),"%")

null_value_check(data_df)

Missing Content in the dataset:
 director      29.91
cast           9.37
country        9.44
date_added     0.11
rating         0.05
duration       0.03
dtype: float64
Total missing %:   4.08 %
```

- Dropped the null values in the `cast`, `country`, and `date_added` using `dropna()`.
- 'rating' column null values have been replaced with `NA`.
- `duration` null values have been replaced with `0 min`.

- **Feature Generation:**
- Created new features like 'added_year', 'added_month', 'added_weekday' from the existing column `date_added`.
- Marked 'others' for the Country column having more than 1 country in its entry.
- `seasons` column derived from `duration` : `min` values been replaced with 0 and others we have extracted the number from the number of seasons string.
- Created `genre` by considering the first value from the list available in `listed_in` column.

```python
#converting the `date_added` datatype from `object` to `datetime`

data_df['date_added'] = pd.to_datetime(data_df['date_added'])

#New Features from existing feature `date_added`
data_df['added_year'] = data_df['date_added'].dt.year

data_df['added_month'] = data_df['date_added'].dt.month

data_df['added_weekday'] = data_df['date_added'].dt.strftime('%A')

data_df[['date_added','added_weekday','added_month','added_year']].head(2)


'''Extracting Number of seasons from `duration` fetaure'''

data_df['seasons'] = data_df['duration'].apply(lambda x: 0 if 'min' in x else
                                               int(x.split(" ")[0]))


'''Generating new feature `genre` from the `listed_in` feature - considering the first value in the column value'''

data_df['genre'] = data_df['listed_in'].apply(lambda x: x.split(',')[0])

data_df['genre'].value_counts().head()
```

❖ Top 10 countries content-wise available on the Netflix

```python
plt.figure(figsize=(15,7))

fig = sns.barplot(data = top10_Countries, x='country', y='title', hue='type', ci=None, palette=['lightblue','yellow'])

for patch in fig.patches:

    x= patch.get_x() + 0.05

    height = patch.get_height()
    if height>0:
        val = ("{0:.2f}".format(patch.get_height()/len(data_df)*100)+"%")
        fig.text(x, height + 0.5 , val)

plt.legend(loc="upper center")

plt.xlabel("country")

plt.ylabel("Count of Movies & TV Shows")

plt.title("Top 10 Country content-wise Number of TV Shows & Movies on Netflix", fontsize=14)

plt.show()
```
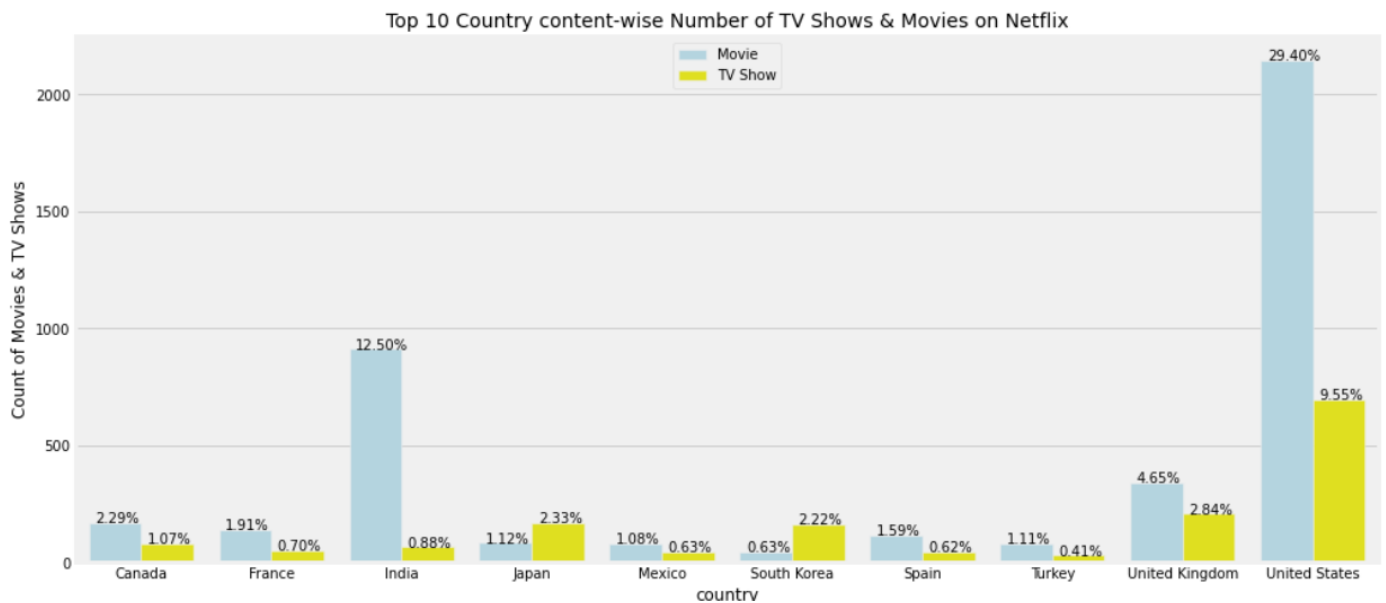


Top 10 Country content-wise Number of TV Shows & Movies on Netflix

**Insights:**

- "United States" is the top country, whose content is most available on Netflix.
- "India" stands 2nd place in the content on Netflix.
- There are very few Indian TV Shows available on Netflix as per the Data, we have at hand.
- "UK" stands 3rd place content-wise on Netflix.
- And also "South Korea" have a very low movie content of `0.63%` on Netflix.

❖ Top 10 genres of content available on Netflix? And What are the different types of
  genres available in the American Content on Netflix? What are the top 3 genres for US
  content?

```python
plt.figure(figsize=(20,7))

draw_circle = plt.Circle((0,0),0.8, color='white')

plt.subplot(1,2,1)

fig1 = plt.pie(top10_genre['title'], labels=top10_genre['genre'], autopct="%0.2f%%", colors = my_palette)

plt.title("Content Available on Netflix")

add_circle_patch = plt.gcf()

add_circle_patch.gca().add_artist(draw_circle)


#Barplot to show the American genres - i.e `United States` content on Netflix

plt.subplot(1,2,2)

UnitedStates_Content = data_df[data_df['country'] == 'United States'].groupby(by='genre').agg({'title':'count'}).sort_values('ti

fig2 = sns.barplot(data = UnitedStates_Content, x = 'genre', y='title', palette = my_palette, ci=None)

fig2.bar_label(fig2.containers[0])

fig2.set_title("American Netflix Content Genres.")

fig2.set_xticklabels(UnitedStates_Content['genre'],rotation = 90)

plt.show()
```



Content Available on Netflix          American Netflix Content Genres.

## Insights:

- 'Dramas' genre takes the 1st place on Netflix TV Shows & Movies.
- `Comedies` is 2nd most genre available on Netflix.
- `Action & Adventure` and `International TV Shows` takes 3rd and 4th place on Netflix.
- Top 3 genres for United States Content are - `Dramas, Comedies, Action & Adventure.

❖ On which weekday does Netflix releases new content? How many releases does it do on that weekday?

```
weekday_releases = data_df.groupby(by=['added_weekday','type']).agg({'title':'count'}).reset_index()

weekday_releases
plt.figure(figsize=(12,7))

fig = sns.barplot(data = weekday_releases , x = 'title', y='added_weekday', palette = my_palette, ci=None)

fig.bar_label(fig.containers[0])

fig.set_title("Netflix Content Realses on Weekdays")

fig.set_xlabel("Number of Releases")

fig.set_ylabel("Weekday")

plt.show()
```



Netflix Content Realses on Weekdays

Insights:

- Netflix releases most of its new content on Friday.
- Based on the data available, it did 1005 releases on Friday.

This dataset with the above features like country, date_added, and type, title, seasons, duration, release year, and genre, rating gives us more information to answer the below hypothesis and I find this is best for explaining them with evidence.

# Hypothesis:

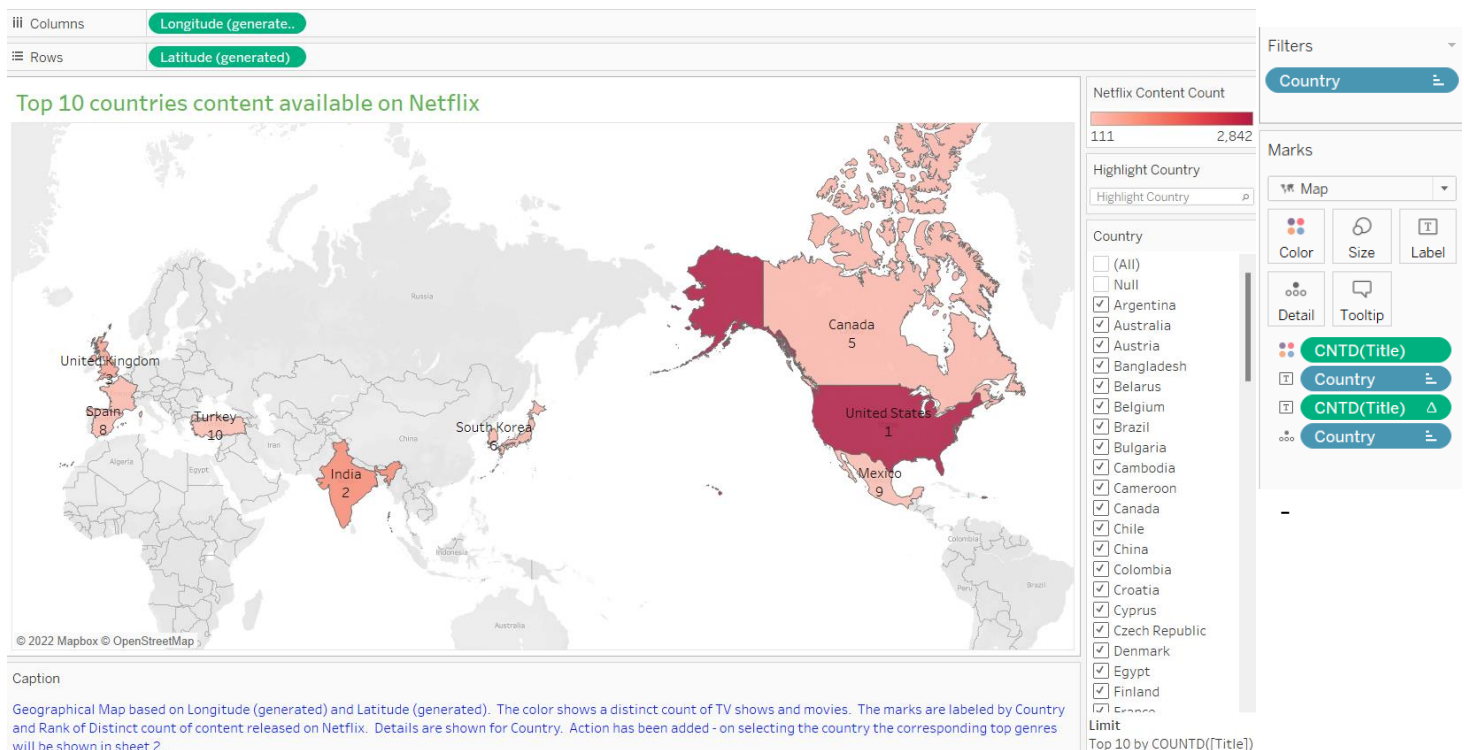1. What are the top 10 countries' content available on Netflix? And What are the top 10 genres produced by United Stated on Netflix? What are the Top 10 genres of content available on Netflix?
2. What month is most suitable for the content to be released on Netflix? As we know, a smaller number of releases means more chances that people may view your content. see which weekday most of the Netflix releases happen in that month and What is the topmost-rated content released in that month?
3. Which rating for the TV Show content is more available on Netflix? And What is the top TV Show with more seasons? And analyze the trend line for both TV Shows and Movies along the date added on Netflix.

**Hypothesis 1:** What are the top 10 countries' content available on Netflix? And What are the top 10 genres produced by United Stated on Netflix? What are the Top 10 genres of content available on Netflix?

- **Geographical Map:**

- The below graph is achieved by using the latitude and longitude values of the attribute `Country` and the contrast of the color indicates the number of TV Shows and Movies content on Netflix.
- The `Rank` has been used from the `Quick Table Calculations` to display the rank on the map for the countries.



Geographical Map based on Longitude (generated) and Latitude (generated). The color shows a distinct count of TV shows and movies. The marks are labeled by Country and Rank of Distinct count of content released on Netflix. Details are shown for Country. Action has been added - on selecting the country the corresponding top genres will be shown in sheet 2.

**Filter** : Top 10 by `Title` - count(Distinct)
**Color** : CNTD(Title) – Number of TV Shows and Movie Content
**Label** : Country , CNTD(Title) – Rank from Quick Table Calculations

- **Action** has been added to the sheet – on click of any country the corresponding top genres in the next sheet.

<span style="color:green">Insights:</span>

- Top 10 Countries are the United States, India, United Kingdom, Japan, Canada, South Korea, France, Spain, Mexico, Turkey
- `United States` ranks **1** in the content on Netflix.
- `India` takes the **2nd** position in terms of content available on Netflix.

- **Bubble chart:**
- This chart is created using the attribute/dimension – Genre and measure – CNTD(title) i.e Number of TV Shows and Movie content available for each genre of the selected country in the previous sheet.
- This is an interactive chart as we have added an action – to filter out the results for the genre based on the country chosen.
- It shows the top genres of each country's content available on Netflix.
  **Filter** : Genre – top 10 by field 'title` - count(distinct)
  **Label** : Genre , CNTD(Title) – 'percentile of total' from Quick Table Calculations.
  **Color** : Genre
  **Size** : CNTD(Title)



Genre and % of Total Distinct count of Title. The color shows details about the Genre. Size shows distinct Netflix content including TV Shows and Movies. The marks are labeled by Genre and % of the Total Netflix content. The data is filtered on Action (Country), which keeps 1 member. The view is filtered on Genre, which keeps the top 10 of 36 members.

- For **United States** Netflix content – the top 10 genres are – Dramas, Comedies, Action & Adventure, Children & Family Movies, Documentaries, Stand-up Comedy, Kid's TV, Horror Movies, Crime TV Shows, and International TV Shows.
- The **topmost available genre** is **Dramas** which is **18.85%** for US content**.**
- **Comedies** with **16.48%** take the **2nd** place and `**Action and Adventure**` takes the **3rd** place for the US Netflix content.

- **Bar Chart:**
- The below bar chart is created using the attribute called `genre` and `country` attributes as filters to show the top 10 genres and their corresponding top 10 country content and the `CNTD(Country)` content for each genre to display as a label.
- The distinct colors indicate the top 10 countries.
- Highlights have been added for the attributes genre and country to highlight the specific country and genre at a time on select.
- To display the percentage of the country content – 'percentile of total' has been selected from the quick table calculations.



Top 10 Genres with top 10 Country TV Shows and Movie Content

Caption

Distinct Netflix Content ( TV Shows & Movies ) for each Genre. The color shows details about the Country. The marks are labeled by % of the Total Count of Country Netflix Content in each of the top 10 genres available on Netflix. The view is filtered by Genre and Country. The Genre filter keeps the top 10 of 36 members. The Country filter keeps the top 10 of 84 members.

Insights:

- Top 10 genres from the above bar chart are – Dramas, Comedies, Action & Adventure, Children & Family Movies, Documentaries, International TV Shows, Stand-up comedy, Kid's TV, Horror Movies, and Crime TV Shows.

- '**South Korea**` stands top **with 62.56%** in producing `**International TV Shows**` content on Netflix.
- `India` stands 2[nd] next to the United States in producing top 5 genre content on Netflix as shown in the above chart.

**Hypothesis 2:** What month is most suitable for the content to be released on Netflix? As we know, a smaller number of releases means more chances that people may view your content. see which weekday most of the Netflix releases happen in that month and
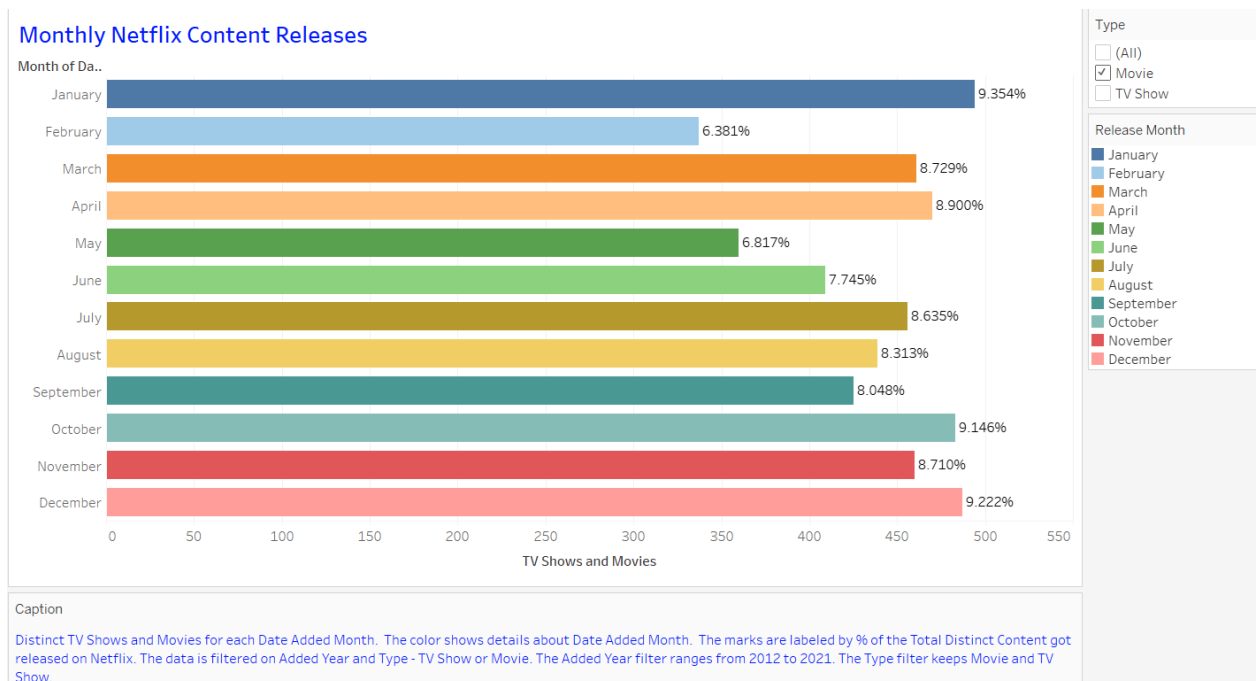What is the topmost-rated content released in that month?

- **Bar Chart:** The horizontal bar chart below shows the percentile of the number of movies got released on Netflix in each month during the year 2012 – 2021.
- This chart has been created using the dimension – `Date_Added` on expanding it to month from the `rows` dropdown. And count the distinct title as columns.
- And a **filter** has been created to filter only movie content from the available data.
- The distinct **color** has shown each month.
- for the **label** – we have selected CNTD(Title) and `percentile of total` from `Quick Table Calculations`.



**Insights:**

- **`February` and `May` are the best months to release any new movie content** as less content has been released in those months from the above bar chart. More chance of high views from the audience.

- High amount of movie content has been released in the months of January, December, October, April, March, and November.

- **Funnel Chart:**

The funnel chart has been considered to show the linear flow of the movie content releases in a specific month on weekdays in decreasing order.

- **Filter** was added to select only the movie content from the attribute `type`.
- **Action** :'Add Filter' action has been created in the previous bar chart, on click of any month – the corresponding weekday's Netflix content release percentage has been shown.
- **Color** : the distinct color has been assigned to each weekday from the `Added_Weekday` attribute.
- **Size** : To have variance in the sizes – we have considered the measure – NetflixContent.csv(Count) of each weekday
- **Label** : `Added_weekday` and the count of netflixcontent.csv measure was considered.
- **Sort** : sorting the results in the descending order of the Netflix Movie content number.



**Weekday Movie Releases on Netflix on Specific Month**

Caption

Count of Movie Releases on Netflix in a Specific month chosen in the previous chart. The color shows details about Weekday, on which the Movie got released. Size shows the count of Netflix movie content. The marks are labeled by the count of NetflixContent and Added Weekday. The data is filtered on Action (MONTH(Date Added)), which keeps 1 member.

Insights:

- When the `February` month was selected - `Monday` has least number of movie releases on Netflix.

- When `May` month was selected - `Sunday` and `Wednesday` have the less number of movie releases on Netflix. And it's the best time to release the movie content in the May.

- **Donut Chart:** It is similar to the pie chart – this type of chart is selected to represent the top 10 movie ratings and their corresponding proportions of the movie releases on Netflix.
- `**Rating**` dimension has been used to show the top 10 ratings and a filter has been added to pick the top 10 ratings using the Top 10 – by 'Added Month` - distinct count.
- **Action :**  is added to pick the specific month's content while displaying the ratings for the movie content.
- **Color :** distinct color represents the different ratings available in the data.
- **Size :** count of the measure ( CNT(NetflixContent.csv)) been used to differentiate the sizes.
- **Label :**  Rating and the 'percentile of  total' from the `quick table calculations` is used for the CNT(NetflixContent.csv) measure  to display the labels.

## Top 10 Rating Movies got released in a specific month.

**Filters**
- Rating
- Type: Movie
- Action (MONTH(D.. ⊘

**Marks**

All

⊘ AGG(AVG(0))

⊘ Pie ▾

| Color | Size | Label |
| Detail | Tooltip | Angle |

- ∷ Rating
- ▷ CNT(NetflixCo..
- ⊘ CNT(NetflixCo..
- T CNT(Netflix.. △
- T Rating

○ AGG(avg(0))

**Action (MONTH(Date A...**
- ☐ (All)
- ☐ January
- ☑ February
- ☐ March
- ☐ April
- ☐ May
- ☐ June
- ☐ July
- ☐ August
- ☐ September
- ☐ October
- ☐ November
- ☐ December

**Type**
- ☐ (All)
- ☑ Movie
- ☐ TV Show

**Rating**
- ■ NR
- ■ PG
- ■ PG-13
- ■ R
- ■ TV-14
- ■ TV-G
- ■ TV-MA
- ■ TV-PG
- ■ TV-Y
- ■ TV-Y7

CNT(NetflixContent.csv)
333

*Chart labels:* 2.10% TV-Y7 · 4.50% PG · 4.80% PG-13 · 10.21% TV-PG · 13.51% R · 39.34% TV-MA · 21.92% TV-14 · 1.50% TV-G

**Caption**

Created a Donut chart using the dual axis - AVG(0) and avg(0). For pane AVG(0): Color shows details about the Rating. Size shows a count of NetflixContent.csv. The marks are labeled by % of Total Count of Netflix movie Content and Rating. The data is filtered on Type and Action (MONTH(Date Added)). The Type filter keeps Movie. The Action (MONTH(Date Added)) filter keeps 1 member. The view is filtered on Rating, which keeps 10 of 18 members.

## Top 10 Rating Movies got released in a specific month.

**Action (MONTH(Date A...**
- ☐ (All)
- ☐ January
- ☐ February
- ☐ March
- ☐ April
- ☑ May
- ☐ June
- ☐ July
- ☐ August
- ☐ September
- ☐ October
- ☐ November
- ☐ December

**Type**
- ☐ (All)
- ☑ Movie
- ☐ TV Show

**Rating**
- ■ NR
- ■ PG
- ■ PG-13
- ■ R
- ■ TV-14
- ■ TV-G
- ■ TV-MA
- ■ TV-PG
- ■ TV-Y
- ■ TV-Y7

CNT(NetflixContent.csv)
357

*Chart labels:* 7.28% TV-PG · 1.12% TV-Y7 · 0.56% NR · 5.60% PG · 5.32% PG-13 · 12.89% R · 24.37% TV-14 · 1.68% TV-G · 40.34% TV-MA

**Caption**

Created a Donut chart using the dual axis - AVG(0) and avg(0). For pane AVG(0): Color shows details about the Rating. Size shows a count of NetflixContent.csv. The marks are labeled by % of Total Count of Netflix movie Content and Rating. The data is filtered on Type and Action (MONTH(Date Added)). The Type filter keeps Movie. The Action (MONTH(Date Added)) filter keeps 1 member. The view is filtered on Rating, which keeps 10 of 18 members.

**Insights:**

- In `February` month, the most released movie content was `TV-MA` rated – which means the content is intended for adults and maybe not for children under 17.

- `TV-14` is the 2nd most-rated movie content released in the month of February. This type of content is not for children under 14 years old.
- From the below donut chart, In `May` as well – `TV-MA` and `TV-14` are the high rated movie content available on Netflix.

**Hypothesis 3:** What are the top 3 rated contents for the TV Show that is more available on Netflix? And What is the top TV Show with more seasons? And analyze the trend line for both TV Shows and Movies along the date added on Netflix.
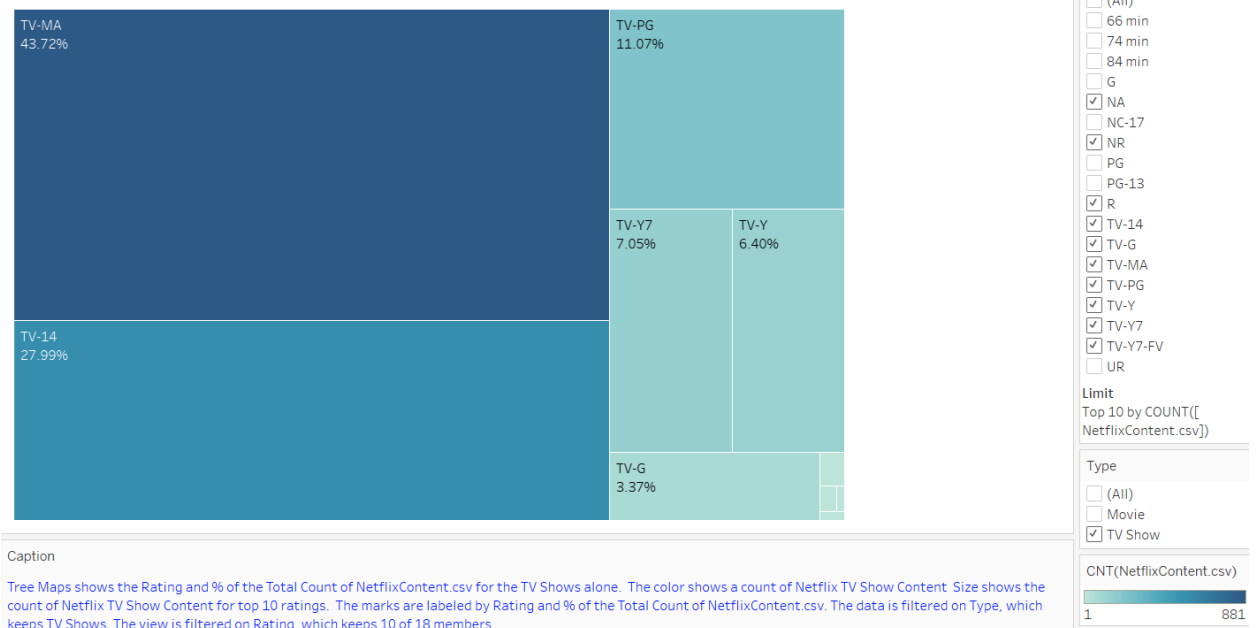
- **Tree Map:**
  This graphical representation is chosen to display the rated content specifically for the TV shows in a hierarchy.
- this chart is created using the `rating` dimension with a filter on content `type` to select only the TV show content.
- Label : 'Rating' and the count of the measure `NetflixContent.csv` is used .
- Size  : the size of the treemap is based on the count of `NetflixContent.csv` measure.
- Color : the darker the color represents the high TV Show content for that rating on Netflix. And the lighter the color – less content for that rating is on Netflix.
- To filter out the top 10 ratings – we selected the top 10 by the measure '`NetflixContent.csv` count
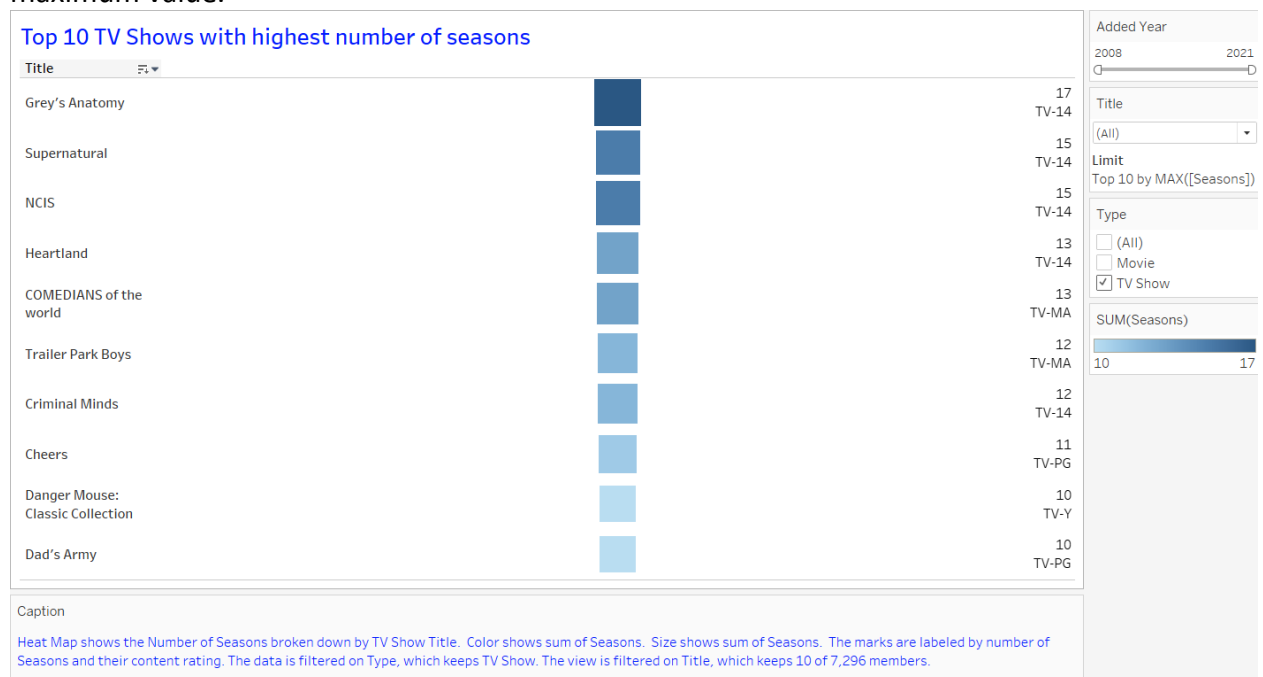


**Top 10 Rated content Analysis for the TV Shows on Netflix**



Caption

Tree Maps shows the Rating and % of the Total Count of NetflixContent.csv for the TV Shows alone.  The color shows a count of Netflix TV Show Content  Size shows the count of Netflix TV Show Content for top 10 ratings.  The marks are labeled by Rating and % of the Total Count of NetflixContent.csv. The data is filtered on Type, which keeps TV Shows. The view is filtered on Rating, which keeps 10 of 18 members.

- When analyzing the TV Shows ratings – from the above tree map, 43.72% of `TV-MA`, rated content is more available on Netflix.
- `TV-14`, `TV-PG` rated content stands 2nd and 3rd in terms of TV Shows on Netflix.

- **Heat Map:** Heatmap is chosen to display the top TV Shows with the highest number of seasons.
- This chart has been created using the `title` and the `season` attributes along with a filter `type` to select only the TV Show content.
- Size : the number of seasons indicate the size of the rectangles
- Label : `rating` and the number of seasons are displayed in labels.
- Color : the intensity of the color depends on the number of seasons.
- And the top 10 tv shows are selected using the filter on the top 10 of `seasons` with maximum value.
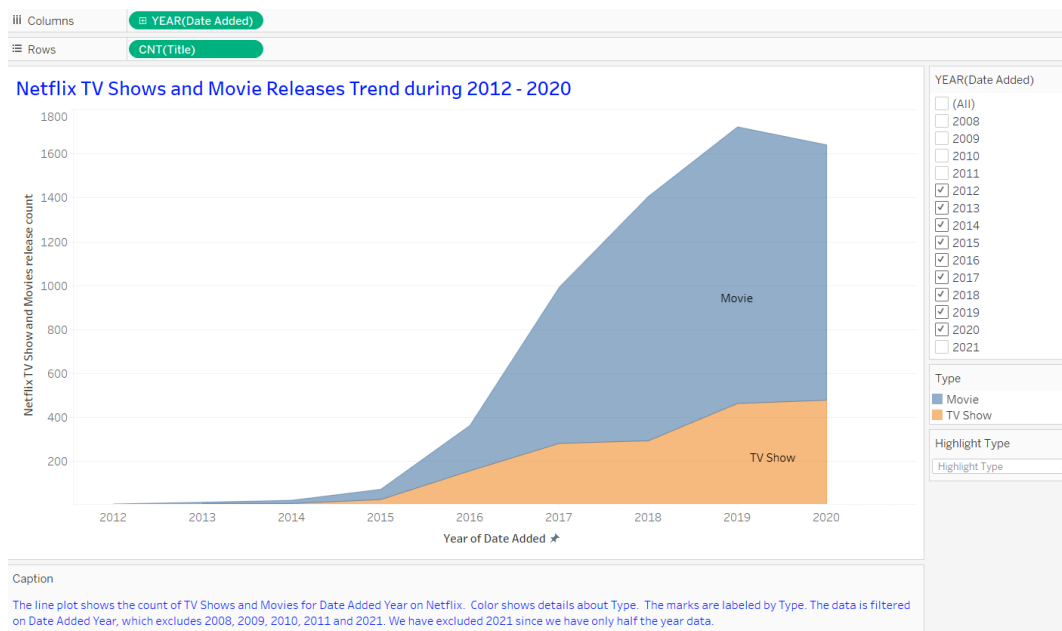


Caption

Heat Map shows the Number of Seasons broken down by TV Show Title. Color shows sum of Seasons. Size shows sum of Seasons. The marks are labeled by number of Seasons and their content rating. The data is filtered on Type, which keeps TV Show. The view is filtered on Title, which keeps 10 of 7,296 members.

Insights :

- From the below heatmap, `**Grey's Anatomy**` is the big TV Show with the highest number of **seasons** with a value of **17** and its **'TV-14'** rated content – which means not appropriate for children under the age of 14.
- `**Supernatural**` and `**NICS**` stands next with 15 seasons and both are TV-14 rated content.

- **Lineplot:**
- The below line plot is chosen to display the trend of the number of TV Shows and Movies released on the Netflix OTT platform during the period of 2012 – 2020. We have excluded the 2021 year as we have only half year data, which we cant use for any conclusions.
- The below plot is created using the dimensions – year of `Date_Added` and the `title`.
- Label : Netflix content `type` is shown in labels
- Color : Color differentiates the type of content between TV shows and Movie.
- Highlighter : can highlight the specific Netflix content type trend line.



Insights:

- From the above line plot, the number of content releases has been increasing from year to year with a negligible decline in between.

## Conclusion:

In this project paper, I have analyzed the various features showing what type of content is available on Netflix. To what extent it is appropriate for the children? And Which countries stand on top in producing content on video streaming platforms like Netflix? What are the major genres available to watch for the users?

Always new content is going to release on Netflix- It's one of the most trending video streaming platforms. Most of the content on Netflix is produced by the United States, India, the UK, and other countries. And the top genres available are Dramas, Comedies, Action & Adventure. And ~40% of the content is most suitable for the adults and Less content is available for the Kid's comparatively to the other genres.

## References:

- *Netflix Movie and TV Shows content from the Kaggle*
  *https://www.kaggle.com/datasets/shivamb/netflix-shows*
- *For creating a Donut chart - https://kb.tableau.com/articles/issue/creating-donut-charts*