

## PROJECT: INVESTIGATING NETFLIX MOVIES



1 hidden cell

**Netflix!** What started in 1997 as a DVD rental service has since exploded into one of the largest entertainment and media companies.

Given the large number of movies and series available on the platform, it is a perfect opportunity to flex your exploratory data analysis skills and dive into the entertainment industry. Our friend has also been brushing up on their Python skills and has taken a first crack at a CSV file containing Netflix data. They believe that the average duration of movies has been declining. Using your friends initial research, you'll delve into the Netflix data to see if you can determine whether movie lengths are actually getting shorter and explain some of the contributing factors, if any.

You have been supplied with the dataset `netflix_data.csv`, along with the following table detailing the column names and descriptions. This data does contain null values and some outliers, but handling these is out of scope for the project. Feel free to experiment after submitting!

## The data

### netflix\_data.csv

| Column                    | Description                     |
|---------------------------|---------------------------------|
| <code>show_id</code>      | The ID of the show              |
| <code>type</code>         | Type of show                    |
| <code>title</code>        | Title of the show               |
| <code>director</code>     | Director of the show            |
| <code>cast</code>         | Cast of the show                |
| <code>country</code>      | Country of origin               |
| <code>date_added</code>   | Date added to Netflix           |
| <code>release_year</code> | Year of Netflix release         |
| <code>duration</code>     | Duration of the show in minutes |
| <code>description</code>  | Description of the show         |
| <code>genre</code>        | Show genre                      |

```
# Importing pandas and matplotlib
import pandas as pd
import matplotlib.pyplot as plt

# Read in the Netflix CSV as a DataFrame
netflix_df = pd.read_csv("netflix_data.csv")

# Subset the DataFrame for type "Movie"
netflix_subset = netflix_df[netflix_df["type"] == "Movie"]

# Select only the columns of interest
netflix_movies = netflix_subset[["title", "country", "genre", "release_year",
                                "duration"]]

# Filter for durations shorter than 60 minutes
short_movies = netflix_movies[netflix_movies.duration < 60]

# Define an empty list
colors = []

# Iterate over rows of netflix_movies
for label, row in netflix_movies.iterrows():
    if row["genre"] == "Children":
        colors.append("red")
    elif row["genre"] == "Documentaries":
        colors.append("blue")
    elif row["genre"] == "Stand-Up":
        colors.append("green")
    else:
        colors.append("black")

# Inspect the first 10 values in your list
colors[:10]

# Set the figure style and initialize a new figure
fig = plt.figure(figsize=(12,8))

# Create a scatter plot of duration versus release_year
plt.scatter(netflix_movies.release_year, netflix_movies.duration, c=colors)

# Create a title and axis labels
plt.title("Movie Duration by Year of Release")
plt.xlabel("Release year")
plt.ylabel("Duration (min)")

# Show the plot
plt.show()
```

```
# Are we certain that movies are getting shorter?
```

```
answer = "no"
```

