

DSC 365 Final Project

George Tzimas Maxwell Ruther Gulbanu Madiyarova
Nitheesh Samiappan Rohith Reddy Patlolla Anwesh Ramesh

March 2024

1 Introduction

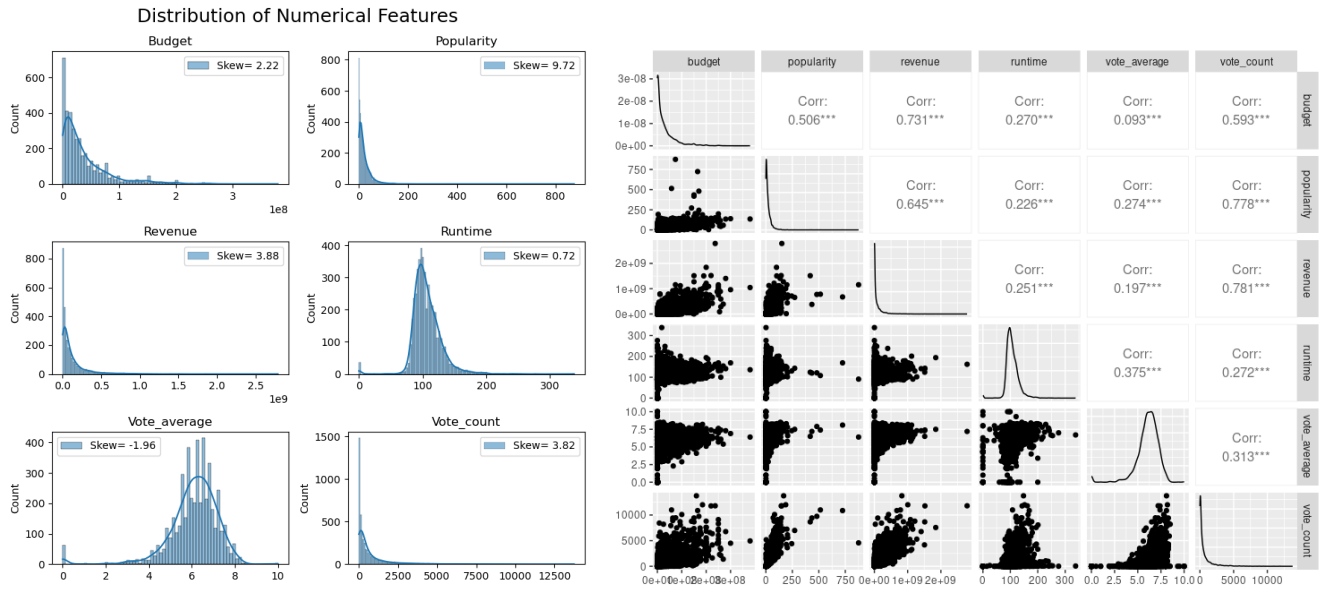
Our dataset includes variables such as Title, Runtime, Genre, Revenue, Budget, Profit, Director, Actor, Production Company, Awards, Oscar Awards, Season, Year, Release Date and etc.... Some of the Key variables for drawing conclusions are revenue, profit, genre, budget, director, actor, and awards. These variables provide insights into financial performance, popular genres, critical acclaim, and potential correlations between creative talents and movie success.

Movies have been a cornerstone of entertainment ever since their introduction in the beginning of the 20th century. From the early black-and-white films to the modern CGI blockbusters of today, they have changed and evolved through technological advancements and cultural changes. One thing that they have maintained throughout these years is the type of message and theme they are trying to disseminate to the audience. These themes and messages can be broadly categorized into **genres**, with each genre trying to impart various sentiments, perspectives, and emotional experiences to its audience.

In this project, we will be exploring movies through the lens of genre in order to unveil what genres have had the biggest impact in terms of **revenue** and **popularity**. We will also take a look into temporal changes within the 21st century, focusing on the movies released between 2001 and 2015. The primary dataset is provided by Kaggle [[The Movie Database \(TMDB\), 2017](#)], although additional features were added either through the use of TMDB's API or through feature engineering.

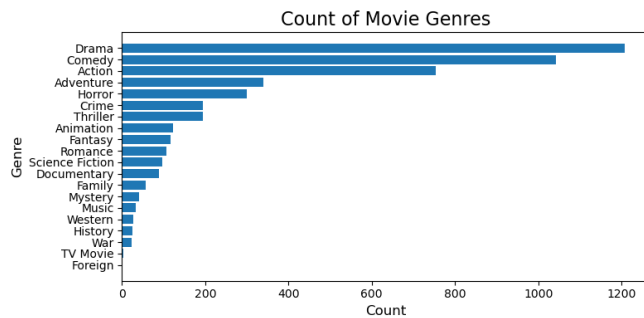
2 Exploratory Data Analysis

The first stage of exploration involved looking over all of the main numeric features in the dataset and their distributions. One of our initial hypotheses was that some of these features may be correlated with each other, which could provide some valuable insight as to where to proceed from there.

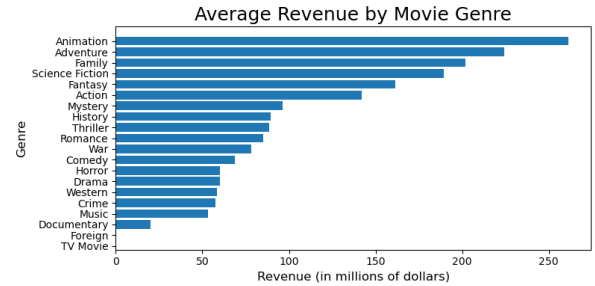


((a)) Distribution of numeric features.

((b)) Correlation plot of numeric features.



((c)) Movie count by genre.



((d)) Average revenue (non-inflation-adjusted) by genre.

Figure 1: Early stages of exploratory data analysis.

The relatively large positive correlation between popularity and revenue (0.65) as well as the non-uniformly distributed nature of average revenue by genre led us down the path of exploring all these different metrics in terms of genre, with the hope of unveiling some interesting insights.

3 Visualizations

3.1 Revenue by Genre, Broadly

Our first visualization is a horizontal bar chart that shows the average revenue for each genre. When viewing its [interactive version](#), one can hover the mouse over a bar to display the poster of the top-earning movie of that genre.



Movie Genres by Average Film Revenue (2001-2015)

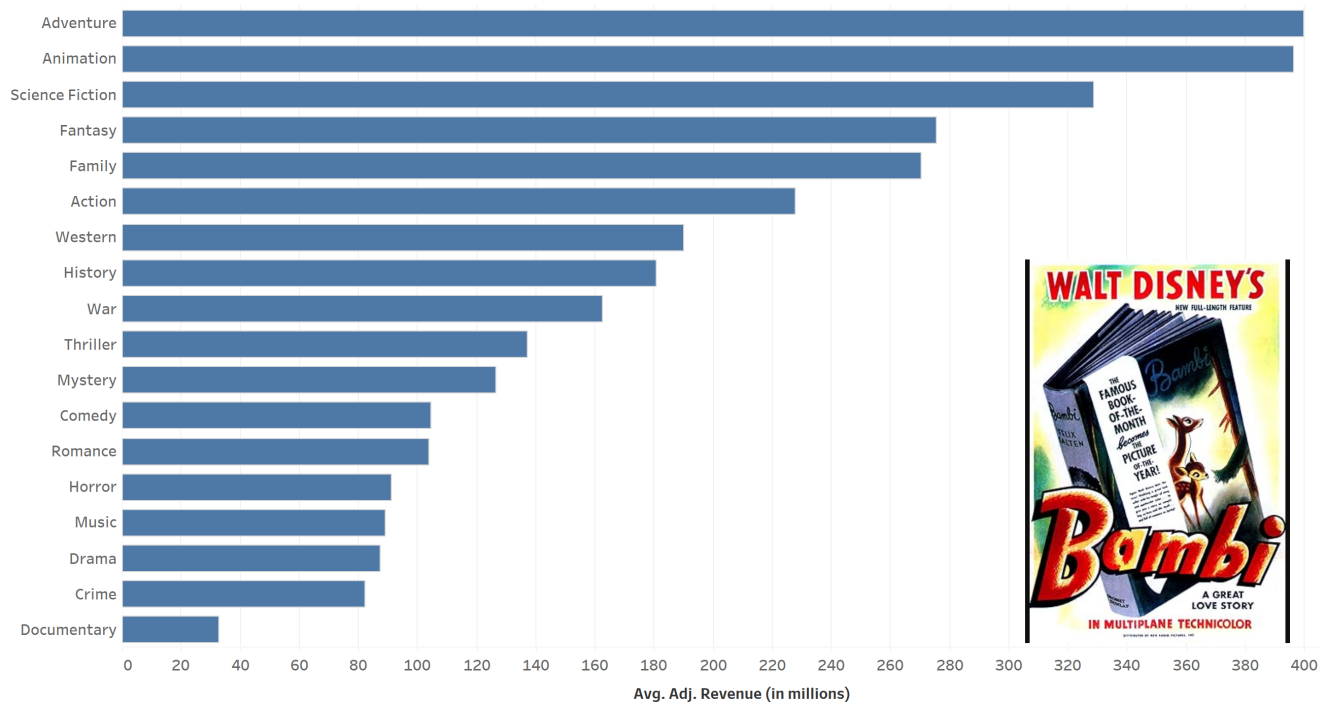


Figure 2: Average revenue by movie genre from 2001 to 2015.

The bar chart was instrumental in giving us an initial sense of what movie genres might be the most successful (as well as some examples of specific films that might define that success.) We here favored averages instead of sums for the visualized measure of revenue, because the averages afford a better view of smaller genres' successes.

Prominently featured are several of the highest-average earning genres: Adventure, Animation, Science Fiction, and Family. These become a special focus of our analysis later on, in a section on genres that thrive in the warmer months.

To further explore how revenue distributes by genre, we next analyzed it year-to-year. In the below [stacked area chart](#), yearly revenue sums are visualized in a a time series analysis. Five genres were the top-grossing of the observed period by quite some margin: Action, Adventure, Drama, Comedy, and Animation. We binned all other genres into "Other". This stacked area chart has the double utility of illustrating both how total revenue has changed annually, as well as part of how that revenue breaks down across genres.

Action and Adventure are the super stars of the bunch. These are followed by Comedy and Drama, which are neck-and-neck. Animation earns the least of these, but does well to make it into this top 5 because it consist of far fewer films, as noted in our next graphic, the mosaic plot.

The below mosaic plot is another time series piece, covering the same top genres as in area chart. Encoded in the width of these mosaic tiles is the number of films for that year.

Its most striking takeaway might be that the distribution of genres stays pretty consistent throughout the period (save for 2007, when the share of "Other" genres briefly spikes.) Given our interest in how genre-specific revenue totals might change over time, this relative consistency of

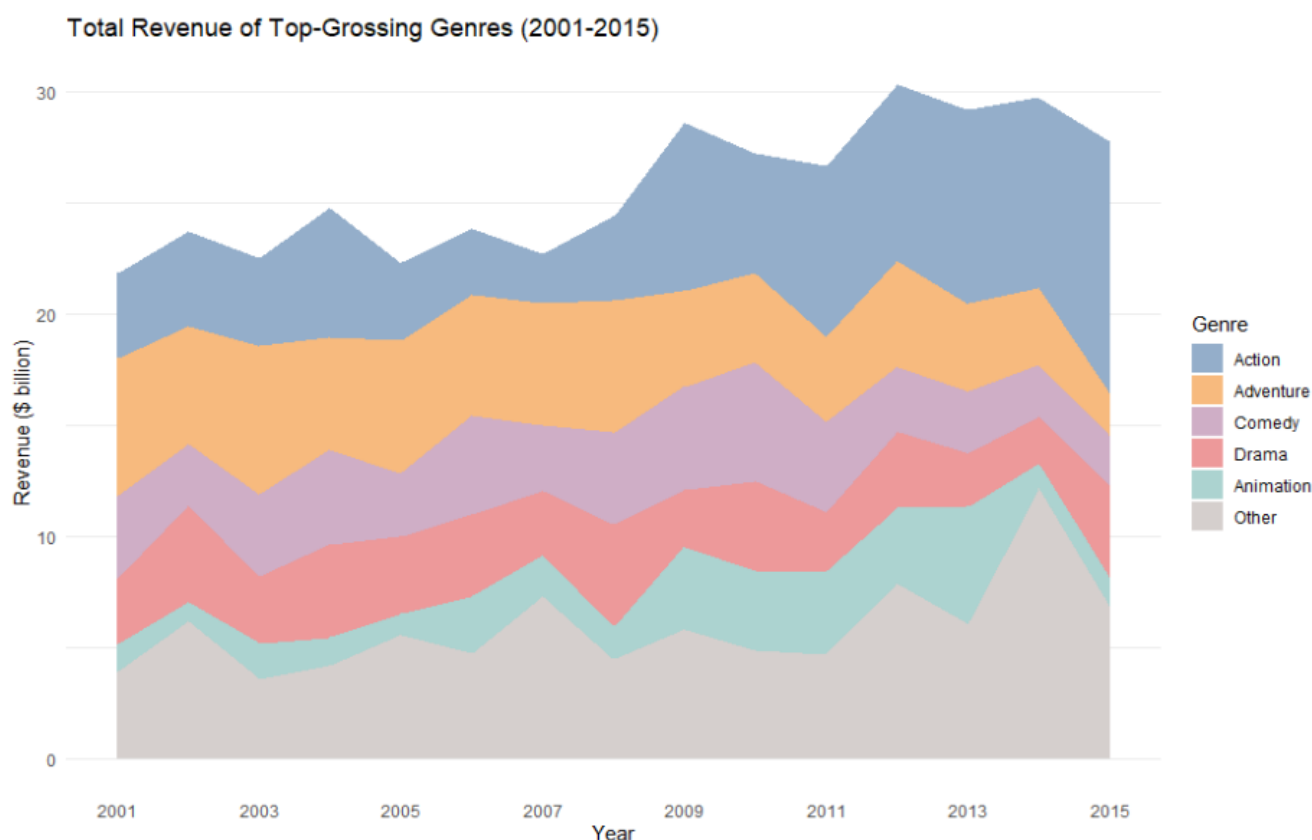


Figure 3: Stacked area chart featuring revenue by year, emphasizing top-selling genres. The above categories are sorted in descending order of revenue total.

sample count was useful context in our analysis. The insight suggests that any jump in a genre-specific portion of a revenue total is unlikely to be a result of a jump in sample count of that genre's films. Instead, such a boost in genre's revenue might likelier owe to above-average success of its films.

Other useful takeaways from this mosaic plot might be how Comedy and Drama form the biggest shares of the sample, though they are third and fourth in terms of their all-time shares of the revenue (behind Action and Adventure.) Also, first visualized here is the small sample size of Animation that was previously mentioned.

3.2 Identifying and examining seasonal boosts in some genres' revenues

We again take a look at revenue by genre over time, but now we instead focus on the month of the year in which films were released, leaving out the year. By scanning the below heatmap to visualize the revenue averages of many genres for each month all-time, one can easily identify patterns in this statistic. Most striking to us was how some genres appear to enjoy hot streaks in the warmer months, in Spring and Summer. We thus noted the best "warm season" genres as Adventure, Animation, Science Fiction, and Family.

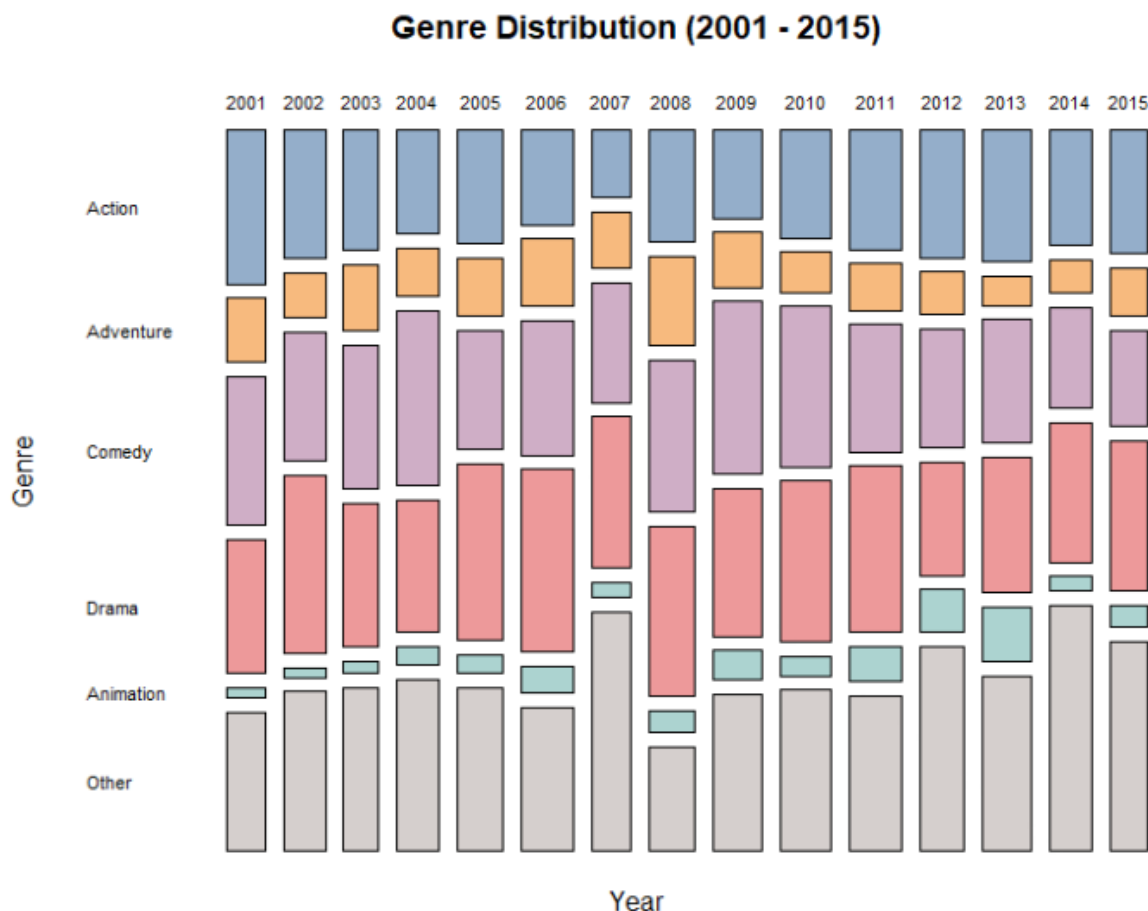


Figure 4: Mosaic plot showing the count of movies by genre from 2001 to 2015.

To start to dig into this phenomenon, we took the totals of sample count, revenue, and popularity for all summer releases all-time, then calculated the shares of those that our "warm season" genres enjoy. The shares are illustrated in the below bar chart, which puts these genres' shares of the various totals side-by-side for easy comparison.

What we found aligned with the high-average streaks of the heatmap: these "warm season" genres enjoy outsize proportions of the Spring and Summertime revenue and popularity. For example, the share of popularity that Adventure enjoys is double its share of the sample count.

To underscore the phenomenon with some contrast, we also featured the "Thriller" genre, which in the heatmap showed consistent revenue averages throughout all seasons. Its indifference to warmer months might well be shown here, as its revenue share is proportionally even with its count.

Soon we will dive into a more numerically gritty analysis of genre-level revenue, but first we take one last look at these Spring and Summer genres. In the above bar chart, we determined that these genres enjoy outsize shares of summertime revenue and popularity. But do these genres

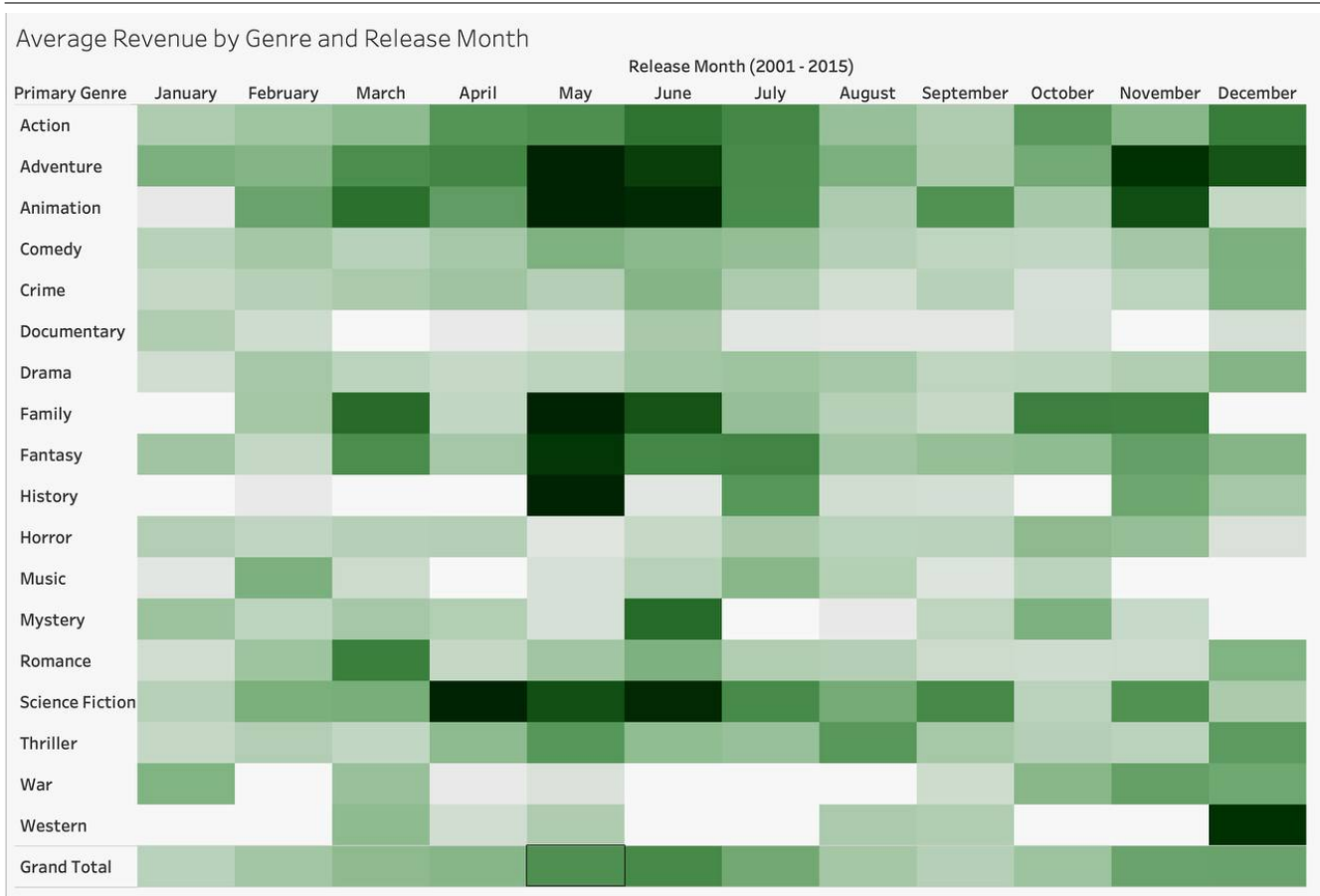


Figure 5: Heatmap of average revenues by genre.

also enjoy higher rates of releasing a summer blockbuster?

First, we define a blockbuster as a movie that places in the top quartile of revenue. Second, we divide films into "Cooler season" and "Warmer Season" bins specific to their genre. With those attributes created, we were able to make the below Sankey. Making this season attribute a subcategory of the genre (one might notice that the genre is named parenthetically in the season value) was an important step, as Sankey's might thrive when most the attributes are hierarchical. This hierarchical structure ensures that the rightmost attribute can be interpreted as an outcome, a conventional feature of Sankeys.

Whether a genre enjoys a warm-season boost in blockbusting is illustrated by the how the blockbusting portions differ between its warmer and cooler season sections. For example, Adventure is shown to provide the biggest share of blockbusters of these genres. However, the proportion of blockbusters to "not"-busters doesn't substantially differ between the warmer and cooler seasons. So Adventure's blockbusting does not get a warm-season boost like its revenue and popularity shares did.

In this vein, Science Fiction and Animation *do* enjoy a blockbuster boost, while the Family genre might actually be negatively affected.



Some Genres Thrive in Warmer Weather

Shares by genre of various totals from Spring and Summer releases, years 2001 to 2015.

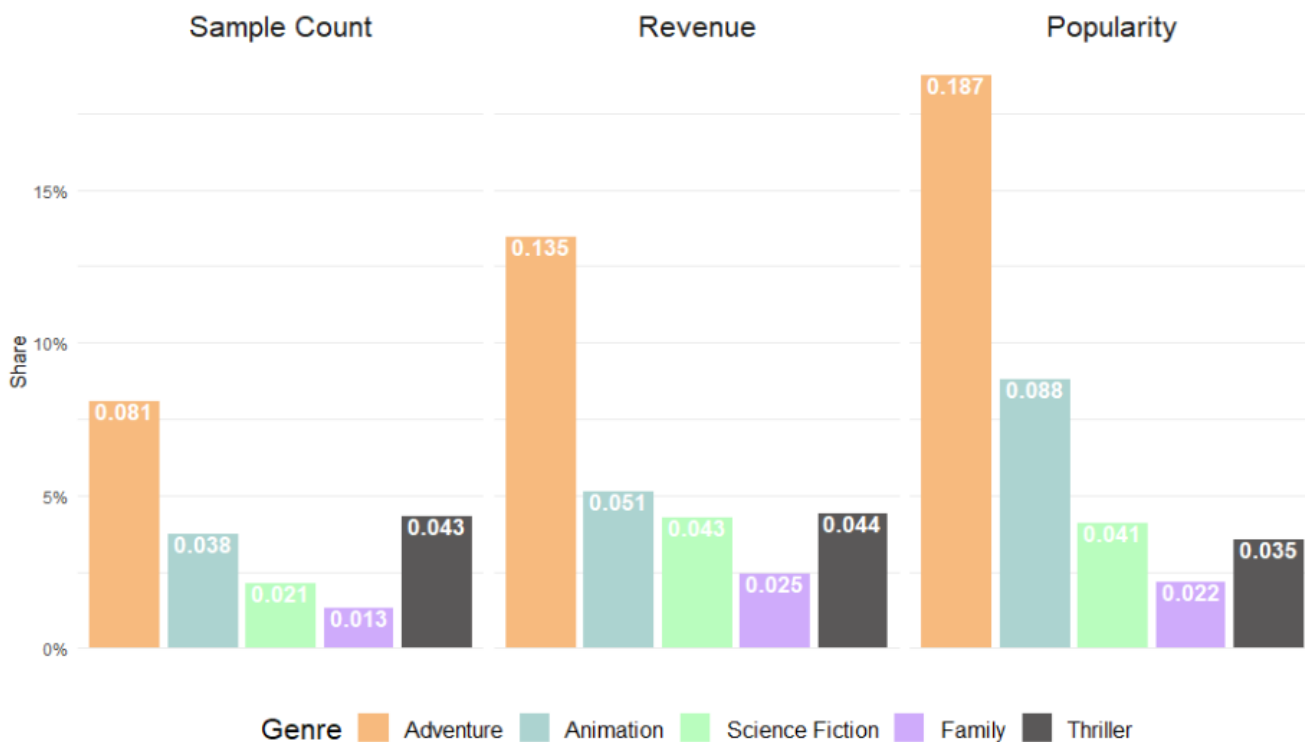
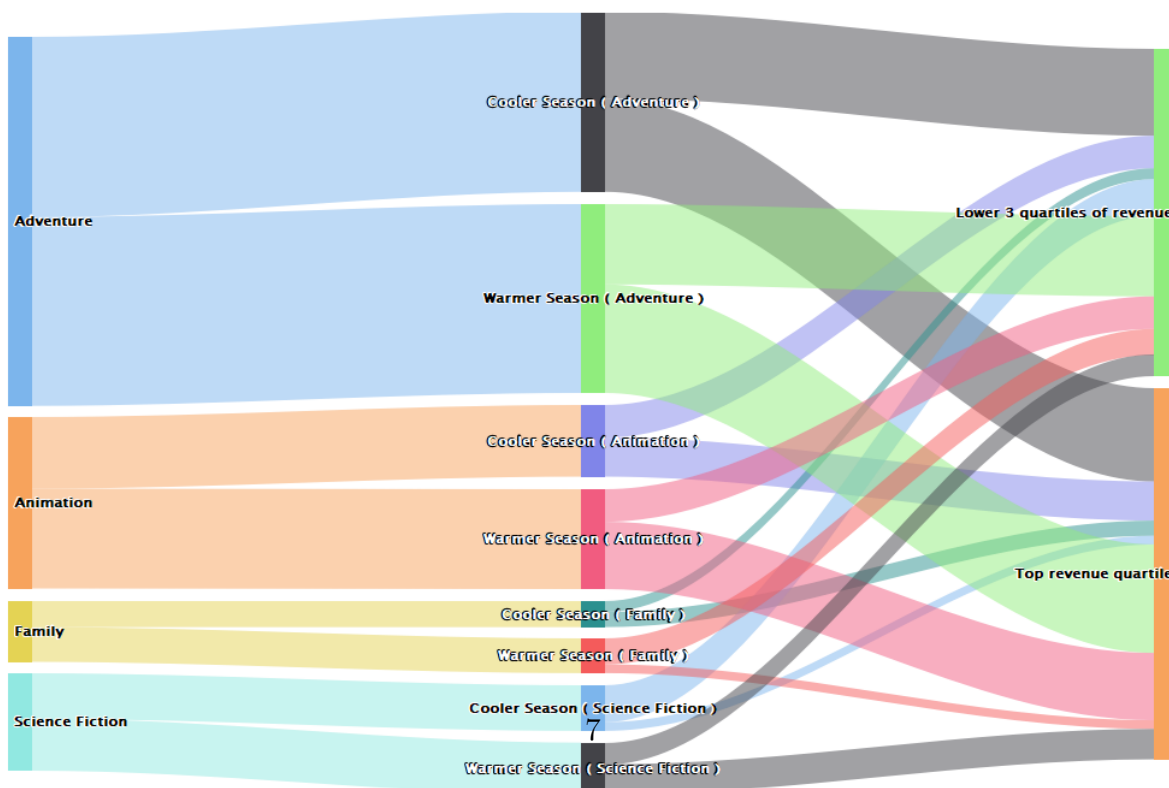


Figure 6: Genres observed in the preceding heatmap as having high average revenues during Spring and Summer are here showcased as having outsized shares of Revenue and Popularity in that warmer period. Thriller is not one of those genres, but is featured here beside the summer heavy-hitters contrast.

Seasonal blockbusting rate is distinct from seasonal average

Though these genres might all register higher average revenues during the warmer months, their proportions of blockbusters in that time don't necessarily improve.





3.3 Assessing correlation between revenue and popularity

For movies, earning high revenue means selling lots of tickets. In this section, we explore whether selling lots of tickets coincides with receiving popular acclaim. We do this by determining the linear correlation between revenue and popularity.

An issue we encountered in this dataset was that the values of revenue were stubbornly concentrated in a certain area. We thought to remove outliers to get a better look at this area, but what we identified to be outliers felt too meaningful and too numerous to remove. To address this overplotting, we created two visualizations which each tackle the issue differently. Featured below is the first approach where we break the problem down by genre by creating faceted scatterplots. Each plot reflects one of the 5 top-grossing genres or "Other", as were visualized previously, and features a correlation value specific to that genre.

The correlations vary widely here. For Animation, the correlation between revenue and popularity is very strong at 0.81. This result aligns with what we've already determined about the Animation genre's outsize success. In contrast, Comedy shows the lowest correlation of these, 0.55. This might align with how senses of humor might vary widely, person-to-person.

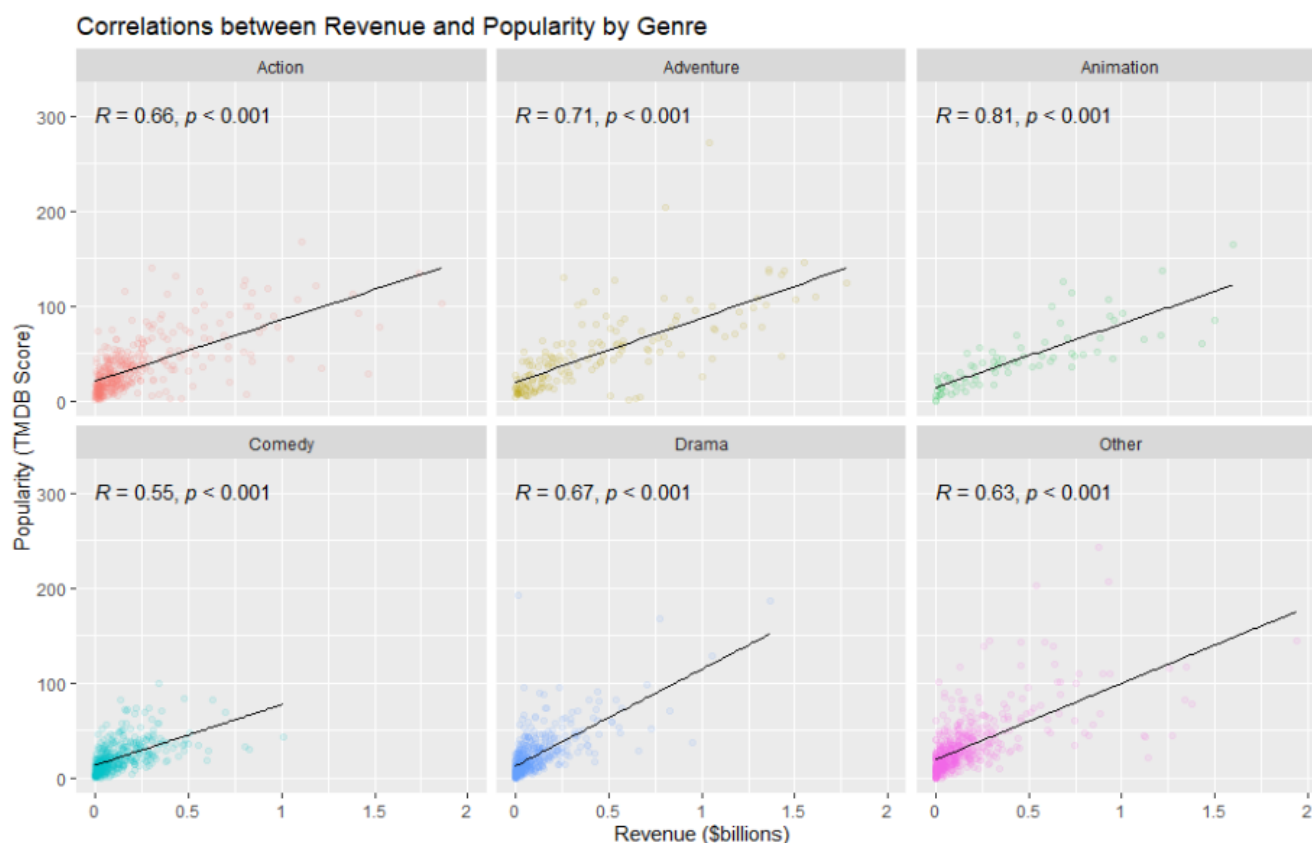


Figure 8: Do higher revenues imply higher popularity scores? Maybe partly, but less than one might expect. And it varies a lot by genre.

Next is our second analysis of the linear correlation between revenue and popularity, looks at the data overall. No longer separating for genre, this linear correlation was determined to be 0.58

, a moderately positive correlation.

In addition to assessing correlation, we here also examine the density of points that so vexed us in previous analyses. Leveraging the alpha value was not as helpful in the second scatterplot as it had been in the first, so in the bottom half of the image, we have "zoomed in" on the area that contains the gross majority of points. This zooming in is done on the area indicated by the cyan rectangle, in the bottom-left of the upper plot.

It is in this area that the 2D-density plot for this same set of points is located. (This is why the second plot is also surrounded by a cyan rectangle; we seek to indicate that these rectangles portray the very same range of revenue and popularity.) Thus, this density plot informs us that the gross majority of films pull in less than \$200 million in revenue and show less than 40 for their TMDB popularity scores.

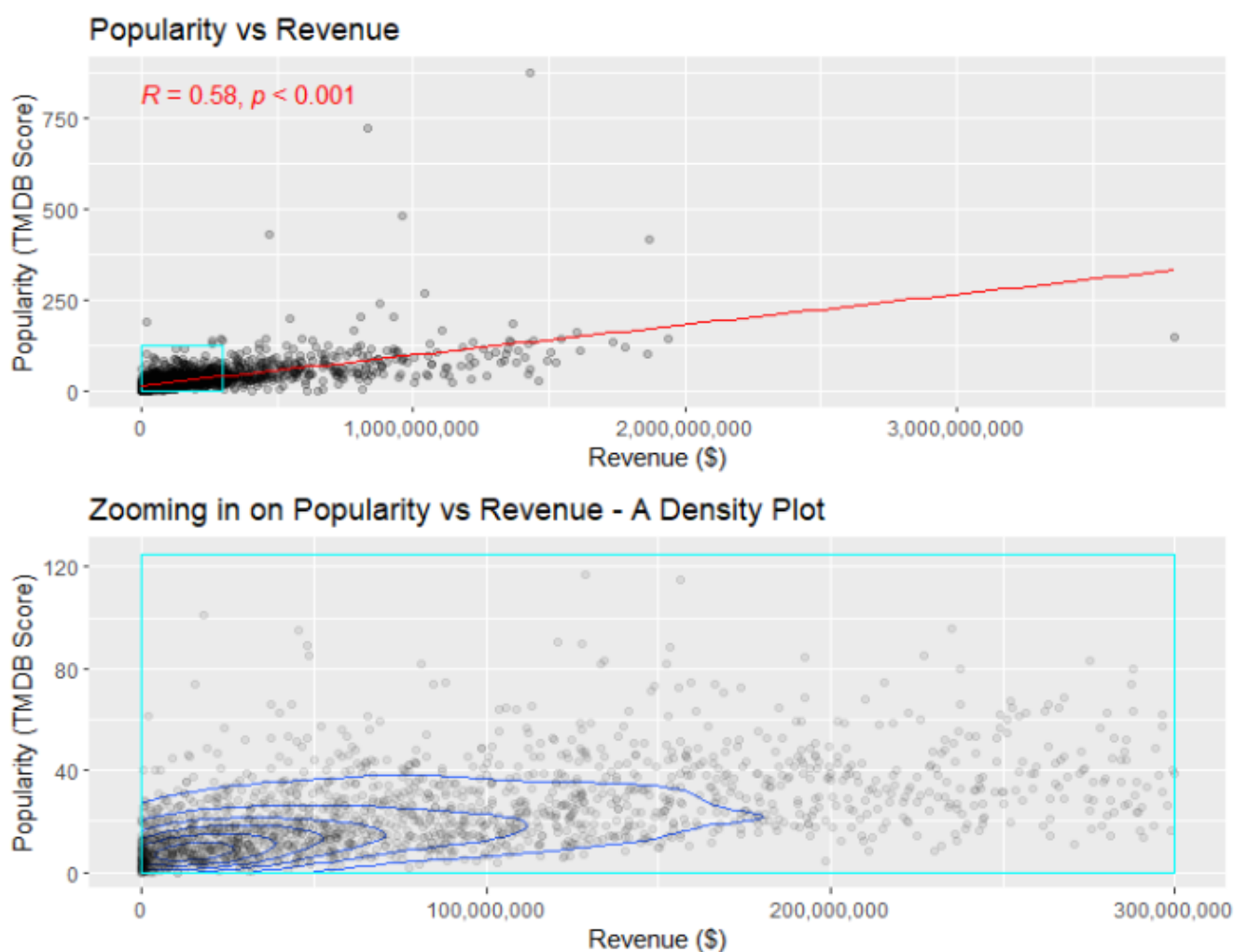


Figure 9: Putting aside genre, what does the relationship between revenue and popularity look like? In our density plot, notice how incredible concentrated the points are in that area, which is a small portion of the upper plot.



4 Analysis and Discussion

The main results of our analysis are the following:

1. Despite being comprised of relatively few films, the Animation genre is in the top 5 for earning revenue. And its spring- and summertime share of the popularity scores is more than double that of its sample count.
2. The average revenue across all genres is high in late Spring and early Summer, but this is especially so for the genres of Adventure, Animation, Family, and Science Fiction.
 - (a) Though these genres might enjoy higher average revenues in that time, they do not all also enjoy higher rates of releasing hit movies.
3. The Animation genre shows the strongest link between revenue and popularity, a strongly positive correlation of 0.81 . This is far stronger than the correlation of 0.58 shown in the greater dataset, which is only moderately positive.

References

[The Movie Database (TMDB), 2017] The Movie Database (TMDB) (2017). TMDB 5000 movie dataset. <https://www.kaggle.com/datasets/tmdb/tmdb-movie-metadata>. Accessed: 2024-03-18.

A Appendix A: Individual Reports

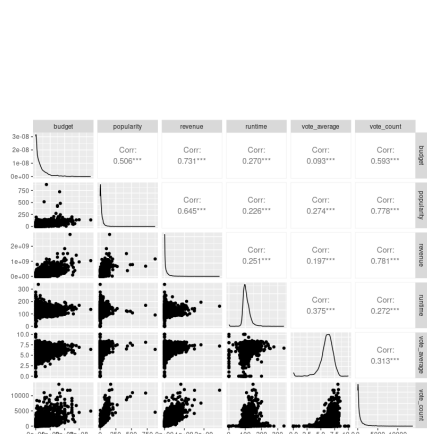
A.1 George Tzimas

- Performed the initial data cleaning process, transforming the dataset from JSON into a relational table.
- Performed feature engineering by extracting relevant information from both csv files provided by Kaggle, and also appending new features through the use of the TMDB and OMDB APIs.
- Created some preliminary visualizations of different aspects of the dataset to give us some ideas as to what direction we want to take this project towards.
- As the group liaison, I was in charge of scheduling weekly meetings through Zoom to go over next steps for project milestones and/or directions for our initial drafts.
- For the in-class presentation, I created interactive visualizations that could be directly embedded into PowerPoint.
- Created the outlines and drafts for the project milestones and the final paper submission through overleaf

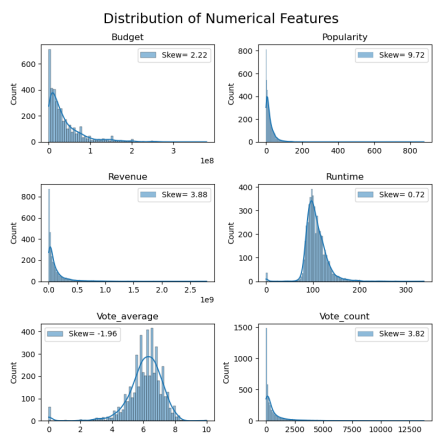


This project challenged me to incorporate many of the techniques learned in the class towards my visualizations. I found myself putting a lot more emphasis on the encoding methods used for each visualization, the visual appeal of the plots in terms of axes tick marks, axes labels, and legends where appropriate. Another valuable thing is the ability to criticize my own work as well as the work of my collaborators more effectively. In addition, coordinating group goals and delegating responsibilities has taught me a lot about what it takes to effectively work within the confines of a group and not on an individual basis. Finally, the challenge of creating a data-driven story from beginning to end has been great learning experience.

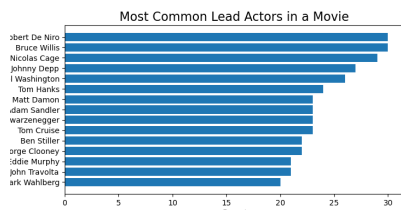
One of the most challenging and humbling aspects that I've learned through this project is how much material ends up being discarded as we go through all the stages from the initial drafting process to the final outputs. The amount of revision and adjustment needed to make everything cohesive and flow together has been eye-opening, and I think it is important to learn about it early on during the academic years so that you are prepared when pursuing a full-time career in an environment that requires team collaboration and meeting deadlines.



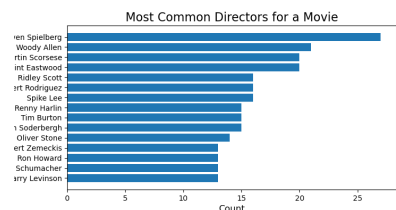
((a)) Correlation plot.



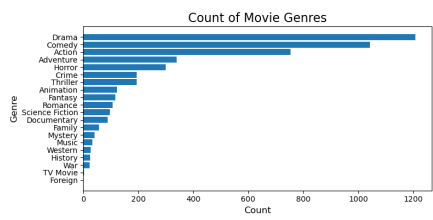
((b)) Distribution plot.



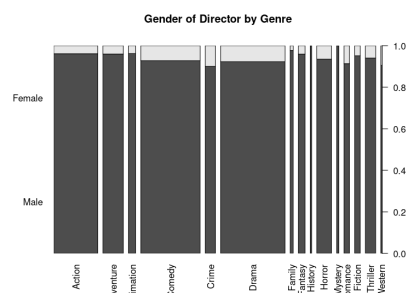
((c)) Top 20 actors.



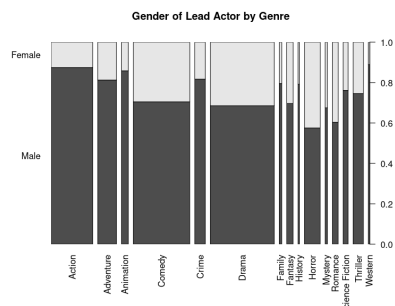
((d)) Top 20 directors.



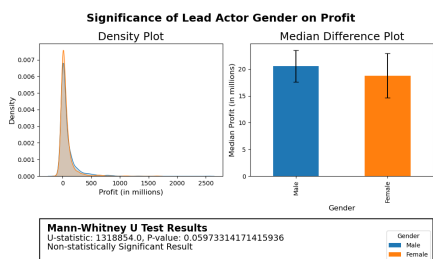
((e)) Top 20 movie genres.



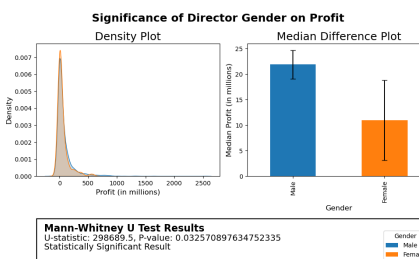
((f)) Gender proportions in actors.



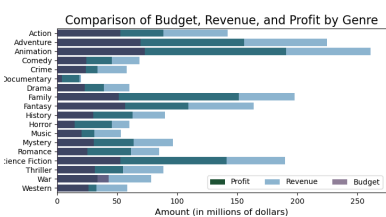
((g)) Gender proportions in directors.



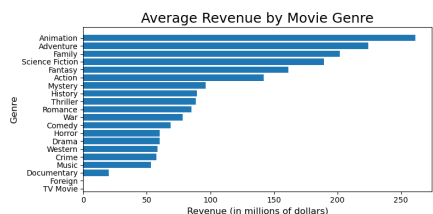
((h)) Statistical test on profit by actor gender.



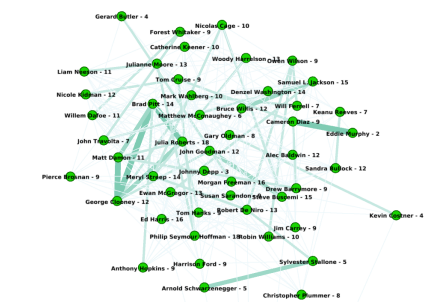
((i)) Statistical test on profit by director gender.



((j)) Stacked bar chart of average budget, revenue, profit by genre.



((k)) Average revenue by genre.



((l)) Network graph of actors and counts of collaborations.



A.2 Maxwell Ruther

- Wrote the "Visualization" and "Analysis & Discussion" sections of this report.
- Fully reworked central visualizations to address a core lack of coherent data message.
- Actively collaborated with group liaison to meet each project milestone promptly and thoughtfully.
- Persisted in engaging the group to determine and implement consistent steps in preprocessing, which included matters of time range, applying inflation adjustment, feature engineering, and sharing the resulting dataset.

I contributed the Summer Genre barchart, the Sankey, and both scatterplots, which include the "faceted-by-genre" scatterplot and the "overall" scatterplot + 2D-density plot. Relating to the scatterplots, I thank Nitheesh for his work as our scatterplot pioneer. He authored the scatterplot from our in-person presentation, which drew interest and rich feedback that substantially informed my design.

My role on the team was to connect visualizations when possible, and draft new ones (intermediaries, ideally) when it wasn't. In a similar vein of keeping things connected, I asserted the importance of inconvenient tasks that bore crucially on the consistency of our visualizations. One such example is when I implemented adjustment for inflation, shared that new dataset, and advocated its adoption. Another is when I updated and reran some of the final visualizations to improve the consistency of their color-coding. (I regret that I couldn't figure out how to do this with my Sankey.)

I learned a lot about data visualization in this project, particularly that it's much harder than I realized to visualize with proper intentionality. To give an example: our heatmap of average revenues revealed that some movie genres are particularly successful in the summer. Finding this interesting (and also following some helpful milestone guidance) I sought to drill down on this pattern and visualize it. But somehow, instead of thinking to simply visualize the averages that I was interested in, I went about it in a roundabout way. I fooled myself into thinking it clever to illustrate the genres' shares of count and revenue, then demand the viewer to divide the latter bar by the former. It wasn't until days later, too late, that I realized that I had sidestepped what I truly meant to communicate, that I could have directly plotted the averages, the values that motivated the whole endeavor. But I had been too taken with my own cleverness or something to realize that straightforward option. Intentional visualization had totally eluded me.

It's my impression that this mistake is tied to the creative process in general. That the challenge of getting out of one's own head is intrinsic to the exercise. But luckily I do know that there are helpful tricks for this, particularly one that was mentioned in lecture: hand drawing the visual first, to sketch it out and see if it still makes sense on paper. So in future projects where I use these data visualization skills we've learned, I'm going to work harder to resist the temptation to immediately run ggplot over and over, and instead reach for the pencil.



A.3 Gulbanu Madiyarova

- Worked individually and crafted exploratory visuals during the exploration phase.
- Actively contributed to group discussions, providing my insights.
- Specifically focused on exploring innovative heatmap visualizations.

This class was highly informative, beginning with the basic exploration phase of data and progressing to the creation of exceptional visuals tailored to the audience's needs. I've gained insight into identifying clutter and avoiding being misled by visuals with misleading axes. Additionally, concepts such as data-to-ink ratio and optimal encoding numbers per visual have been particularly helpful and have become something I enjoy in everyday life.

Participating in class discussions allowed me to share various visuals and provide feedback, enabling me to explore different types of visuals available and discuss their strong and weak points, which I found enjoyable.

We devoted significant time to our team, working on the final project. We selected the data we wanted to explore, learned different cleaning tools, and determined the direction we wanted to explore in our project. Collaborating with my teammates on the final project was a highlight, as we exchanged ideas, shared knowledge, and reached common conclusions together.

A.4 Nitheesh Samiappan

Visualization Techniques and Rich Data Displays:

- Created Bar charts which compared genres with total revenue and awards, identifying top-performing genres and those receiving critical acclaim.
- Scatter plots visualized the revenue-budget relationship, grouped by genre, revealing potential correlations and outliers.
- Mosaic plots explored the distribution and associations between genres and other categorical variables like production companies or directors, uncovering potential patterns or dependencies.
- The combination of these techniques provided a multi-layered, comprehensive view, enabling the exploration of financial performance, critical reception, genre trends, and success factors through interactive and dynamic visualizations.

Design Considerations, Interactivity, and Complementary Techniques:

- Utilized effective design elements, including hue colors with reduced opacity for scatter plots, log scales for data point distribution, and the viridis color palette for mosaic plots.



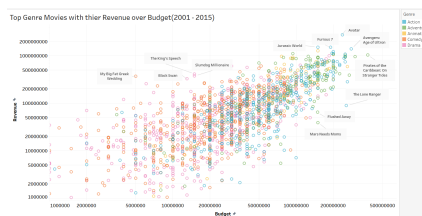
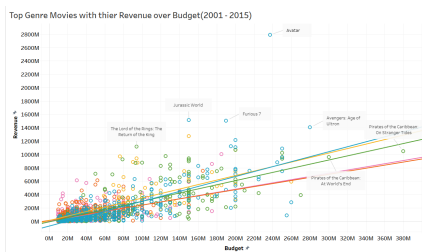
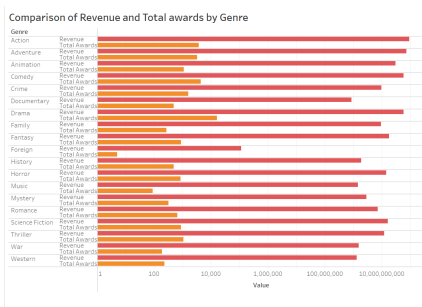
-
- Incorporated interactivity in scatter plots, allowing users to interact with individual data points, view movie titles and details, and display trend lines for each genre.
 - The bar charts, scatter plots, and mosaic plots complemented each other, displaying different aspects of the data and providing a comprehensive view of financial success, critical acclaim, genre trends, and potential success factors.

Conclusions, Further Analysis:

- The visualizations enable conclusions about top-performing genres, genre-specific revenue-budget correlations, and potential genre-company or director associations.
- Statistical techniques like regression analysis, cluster analysis, or time series forecasting could further explore relationships and validate hypotheses generated from the visualizations.

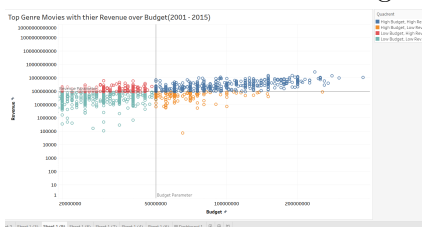
Learning Experience:

- Worked on various visualizations, such as bar charts, mosaic plots, and scatter plots, which provided me with valuable insights on how to enhance the visualizations for better comprehension.
- I learned a great deal about effectively utilizing axis ranges and creating quadrants in scatter plots for improved plot interpretation. As part of a team, I gained a better understanding of my visualizations through constant review and feedback from my teammates, which proved to be an enriching learning experience.
- Additionally, I got to learn about how to use R Studio and Tableau to create proper and neat visualizations for complex and large datasets, as well as different color scales to use for effective visual communication.
- Participating in discussions helped me identify and rectify errors and mistakes, further contributing to my learning journey.



((b)) Scatter plot of top genre movies with revenue and budge.

((c)) Top genre movies with revenue and budget.



((d)) Quadrant Scatter plot.

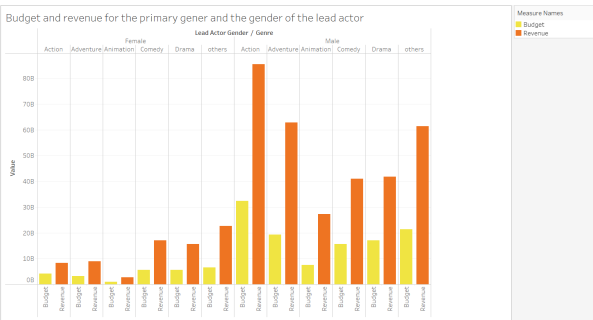
Figure 11: Early drafts



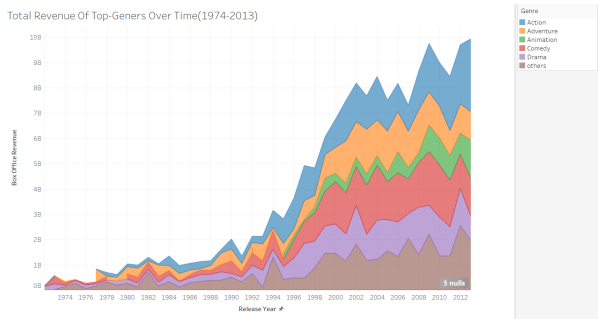
A.5 Rohith Reddy Patlolla

- I produced some creative visualizations with the goal of identifying patterns in the overall revenue of the most popular genres throughout a given period (1974-2013).
- Created side-by-side bar charts that contrasted revenue and budget for the main genre and the lead actor’s gender, providing a detailed insight of the financial aspects and casting decisions in the film business.
- Actively participated in group discussions, offering insights and suggestions.
- Worked actively with the team to make sure that project goals were reached properly and on time, displaying a dedication to the project’s accomplishment.
- By combining all these visualization techniques, the project’s altogether influence was increased through interactive and dynamic visualizations to offer a multi-layered and thorough picture of variables like financial performance, critical reception, genre trends, and success.

This project provided me with a great learning opportunity, allowing me to apply a number of visualization techniques learnt during the course. I focused the significance of all aspects while focusing on encoding techniques and improving visual appeal. Learning how to effectively collaborate with others, including assigning responsibilities and managing group goals, was a big takeaway from this experience. Challenges like making sure that visualization design has clearly brought attention to the significance of intent and clear communication. With teamwork and practical experience using programs like R Studio and Tableau, I improved my ability to produce powerful visuals. Considering all things, this project has improved my understanding of data visualization and prepared me for more research in the area.



((a)) Side-by-side bar chart of budget and revenue.



((b)) Area chart of revenue over time.

Figure 12: Early drafts



A.6 Anwesh Ramesh

- Worked in crafting explanatory and exploratory visuals.
- Crafted Bar charts, Line Charts and Area Charts.
- Actively participated in group reviews and discussions.
- Aided in creating Mosaic Plots for gender proportion for respective genres.
- Created Interactive visuals with calculated fields for highlight actions.

I appreciate the in-depth teachings I received throughout this course, which enabled me to understand and incorporate the various methods, techniques, and the power in visual representation of data. The essential takeaways were towards understanding the audience and tailoring the visualizations accordingly; and significance of storytelling by crafting the narrative.

Additionally, from rich feedback and group reviews I learned a lot about polishing the visualizations-encoding, editing and formatting axes, colour choice, legends, LOD expressions for calculating and transforming fields has been a very challenging and informative industrial experience while working as part of the team.

Overall, the course has been instrumental in shaping my approach to presenting data effectively in my project. By applying the principles and techniques learned, I was able to create rich, informative, polished visualizations and interactive dashboards that not only showcase my results but also facilitate understanding and decision-making for audience.

B Appendix B: Code

B.1 Python

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import os
6 import ast
7 import json
8 from scipy import stats
9 import requests
10 from IPython.display import HTML, display
11
12 pd.set_option('display.max_columns', None)
```



```
1 PATH = '/media/georgetz/jupyter/Desktop/Classes/Winter 2024/DSC 365 -  
    Data Visualization/Final Project'  
2 os.chdir(PATH)
```

```
1 movies = pd.read_csv(os.path.join(PATH, 'tmdb_5000_movies.csv'))  
2 print(f"Shape: {movies.shape}")  
3 movies.head()
```

```
1 credits = pd.read_csv(os.path.join(PATH, 'tmdb_5000_credits.csv'))  
2 print(f"Shape: {credits.shape}")  
3 credits.head()
```

```
1 df = pd.merge(movies, credits, on='title')  
2 print(f"Shape: {df.shape}")  
3 df.head()
```

```
1 def convert(obj):  
2     L = []  
3     for i in json.loads(obj):  
4         L.append(i["name"])  
5     return L  
6  
7 df["genres"] = df["genres"].apply(convert)  
8 df["keywords"] = df["keywords"].apply(convert)  
9 df['PrimaryGenre'] = df['genres'].apply(lambda x: x[0] if isinstance(x,  
    list) and len(x)>=1 else None)  
10 df.head()
```

```
1 def update_cast(key: str, value: int):  
2     def inner_function(json_string: str):  
3         result = []  
4         for item in ast.literal_eval(json_string):  
5             if item[key] <= value:  
6                 result.append(item['name'].strip())  
7         return result  
8     return inner_function
```

```
1 df['cast'] = df['cast'].apply(update_cast('order', value=5))  
2 df['cast'].head()
```



```
1 def update_crew(key: str, values: List[str]):
2     def inner_function(json_string: str):
3         result = []
4         for item in ast.literal_eval(json_string):
5             if item[key] in values:
6                 result.append(item['name'].strip())
7         return result
8     return inner_function
```

```
1 df['Directors'] = df['crew'].apply(update_crew('job', values=['Director'
2     ]))
3 df['PrimaryDirector'] = df['Directors'].apply(lambda x: x[0] if
4     isinstance(x, list) and len(x)>=1 else None)
5 df[['Directors', 'PrimaryDirector']].head()
```

```
1 df['Producers'] = df['crew'].apply(update_crew('job', values=['Producer'
2     ]))
3 df['PrimaryProducer'] = df['Producers'].apply(lambda x: x[0] if
4     isinstance(x, list) and len(x)>=1 else None)
5 df[['Producers', 'PrimaryProducer']].head()
```

```
1 df['crew'] = df['crew'].apply(update_crew('job', values=['Screenplay', '
2     Producer', 'Editor', 'Writer', 'Director']))
3 df['crew'].head()
```

```
1 def update_production_companies(key: str):
2     def inner_function(json_string: str):
3         result = []
4         for item in ast.literal_eval(json_string):
5             if item[key] <= value:
6                 result.append(item['name'].strip())
7         return result
8     return inner_function
```

```
1 df['production_companies'] = df.production_companies.apply(lambda x: [i[
2     'name'] for i in ast.literal_eval(x)])
3 df['PrimaryProductionCompany'] = df['production_companies'].apply(lambda
4     x: x[0] if len(x) > 0 else None)
5 df[['production_companies', 'PrimaryProductionCompany']].head()
```



```
1 df['spoken_languages'] = df.spoken_languages.apply(lambda x: [i['name']  
    for i in ast.literal_eval(x)] )  
2 df['production_countries'] = df.production_countries.apply(lambda x: [i['  
    'name'] for i in ast.literal_eval(x)] )  
3 df['PrimaryProductionCountry'] = df['production_countries'].apply(lambda  
    x: x[0] if len(x) > 0 else None)  
4 df[['production_countries', 'PrimaryProductionCountry']].head()
```

```
1 for coli in df.columns:  
2     df[coli].replace('[]', np.nan, inplace=True)  
3 df.budget.replace(0, np.nan, inplace=True)  
4 df.revenue.replace(0, np.nan, inplace=True)  
5 df.release_date = pd.to_datetime(df.release_date)
```

```
1 def get_season(dates):  
2     seasons = {'spring': (3, 1), 'summer': (6, 1), 'autumn': (9, 1), '  
    'winter': (12, 1)}  
3  
4     month_to_season = {}  
5     for season, (start_month, _) in seasons.items():  
6         if season == 'winter':  
7             month_to_season[1] = month_to_season[2] = month_to_season  
            [12] = season  
8         elif season == 'spring':  
9             month_to_season[3] = month_to_season[4] = month_to_season[5]  
            = season  
10        elif season == 'summer':  
11            month_to_season[6] = month_to_season[7] = month_to_season[8]  
            = season  
12        elif season == 'autumn':  
13            month_to_season[9] = month_to_season[10] = month_to_season  
            [11] = season  
14  
15  
16        return dates.map(lambda x: month_to_season.get(x.month, 'unknown')  
            if pd.notna(x) else np.nan)  
17  
18 df['release_year'] = df.release_date.dt.year  
19 df['release_month'] = df.release_date.dt.strftime('%B')  
20 df['release_season'] = get_season(df.release_date)  
21 df[['release_date', 'release_year', 'release_month', 'release_season']].  
    head()
```

```
1 df['cast'] = df['cast'].apply(lambda x: ast.literal_eval(x) if x is not  
    np.nan else np.nan)
```



```

2 df['crew'] = df['crew'].apply(lambda x: ast.literal_eval(x) if x is not
    np.nan else np.nan)
3 df['Directors'] = df['Directors'].apply(lambda x: ast.literal_eval(x) if
    x is not np.nan else np.nan)
4 df['Producers'] = df['Producers'].apply(lambda x: ast.literal_eval(x) if
    x is not np.nan else np.nan)
5 df['genres'] = df['genres'].apply(lambda x: ast.literal_eval(x) if x is
    not np.nan else np.nan)
6 df['keywords'] = df['keywords'].apply(lambda x: ast.literal_eval(x) if x
    is not np.nan else np.nan)
7 df['production_companies'] = df['production_companies'].apply(lambda x:
    ast.literal_eval(x) if x is not np.nan else np.nan)
8 df['production_countries'] = df['production_countries'].apply(lambda x:
    ast.literal_eval(x) if x is not np.nan else np.nan)
9 df['spoken_languages'] = df['spoken_languages'].apply(lambda x: ast.
    literal_eval(x) if x is not np.nan else np.nan)

```

```

1 lead_actor_gender = pd.Series(credits.cast.apply(lambda x: x[0]['gender'
    ] if len(x) != 0 else np.nan).map({1: 'Female', 2: 'Male'}), name='
    LeadActorGender')
2 df = pd.merge(df, pd.concat([credits.movie_id, lead_actor_gender], axis
    =1), left_on='id', right_on='movie_id').drop(columns='movie_id_y')
3 df['release_date'] = pd.to_datetime(df['release_date'])
4 df.head()

```

```

1 df['MovieProfited'] = df.profit.apply(lambda x: x > 0)
2 df[['budget', 'revenue', 'profit', 'MovieProfited']].head()

```

```

1 def assign_decade(year):
2     if 1910 <= year < 1960:
3         return '1910-1950'
4     elif 1960 <= year < 2000:
5         return str(year // 10 * 10)
6     elif year >= 2000:
7         return str(year // 10 * 10)
8     else:
9         return np.nan
10
11 df['DecadeBin'] = df['release_date'].dt.year.apply(assign_decade)

```

```

1 def adjust_currency(amount, date):
2     if amount is not np.nan and type(date) != pd._libs.tslibs.nattypes.
    NaTType:

```



```
3         adjusted_revenue = cpi.inflate(amount, date.year)
4         return adjusted_revenue
5     else:
6         return np.nan
7
8 df['revenue_adj'] = df.apply(lambda x: adjust_currency(x['revenue'], x['
9     release_date']), axis=1)
10 df['budget_adj'] = df.apply(lambda x: adjust_currency(x['budget'], x['
    release_date']), axis=1)
11 df['profit_adj'] = df.apply(lambda x: adjust_currency(x['profit'], x['
    release_date']), axis=1)
```

```
1 rating_mapping = {
2     'R': 'R',
3     'PG-13': 'PG-13',
4     'PG': 'PG',
5     'Not Rated': 'Unrated',
6     'G': 'G',
7     'Unrated': 'Unrated',
8     'Approved': 'PG',
9     'Passed': 'PG',
10    'TV-MA': 'TV-MA',
11    'NC-17': 'TV-MA',
12    'TV-14': 'PG-13',
13    'TV-G': 'G',
14    'TV-PG': 'PG',
15    'GP': 'PG',
16    'M/PG': 'PG-13',
17    'X': 'TV-MA',
18    'M': 'PG',
19    '13+': 'PG-13',
20    '18+': 'TV-MA'
21 }
22 df['RatedUpdated'] = df.Rated.map(rating_mapping)
23 df['RatedUpdated'].value_counts(dropna=False)
```

```
1 API_KEY = '#####'
2
3 def get_budget(movie_id, api_key=API_KEY):
4     url = f'https://api.themoviedb.org/3/movie/{movie_id}?api_key={
5         api_key}'
6
7     response = requests.get(url)
8     if response.status_code == 200:
9         data = response.json()
10        budget = data['budget']
11        if budget == 0:
```



```
11         return np.nan
12     return budget
13 else:
14     return np.nan
```

```
1 def get_budget(movie_id, api_key=API_KEY):
2     url = f'https://api.themoviedb.org/3/movie/{movie_id}?api_key={
3         api_key}'
4
5     response = requests.get(url)
6     if response.status_code == 200:
7         data = response.json()
8         budget = data['budget']
9         if budget == 0:
10             return np.nan
11         return budget
12     else:
13         return np.nan
14
15 missing_budget_indices = df[df.budget.isna()].index
16 for i in missing_budget_indices:
17     movieid = df.loc[i, 'id']
18     try:
19         df.loc[i, 'budget'] = get_budget(movieid)
20     except:
21         print(i)
22         break
```

```
1 def get_revenue(movie_id, api_key=API_KEY):
2     url = f'https://api.themoviedb.org/3/movie/{movie_id}?api_key={
3         api_key}'
4
5     response = requests.get(url)
6     if response.status_code == 200:
7         data = response.json()
8         revenue = data['revenue']
9         if revenue == 0:
10             return np.nan
11         return revenue
12     else:
13         return np.nan
14
15 missing_revenue_indices = df[df.revenue.isna()].index
16 for i in missing_revenue_indices:
17     movieid = df.loc[i, 'id']
18     try:
19         df.loc[i, 'revenue'] = get_revenue(movieid)
```




```
19 except:
20     print(i)
21     break
```

```
1 def get_producers(movie_id, api_key=API_KEY):
2     url = f"https://api.themoviedb.org/3/movie/{movie_id}/credits?
3         api_key={api_key}"
4     response = requests.get(url)
5     if response.status_code == 200:
6         data = response.json()
7         producers = [member['name'] for member in data['crew'] if member
8             ['job'] == 'Producer']
9         if len(producers) == 0:
10             return np.nan
11         return str(producers)
12     else:
13         return np.nan
14
15 missing_producer_indices = df[df.Producers.isna()].index
16 for i in missing_producer_indices:
17     movieid = df.loc[i, 'id']
18     try:
19         df.loc[i, 'Producers'] = get_producers(movieid)
20     except:
21         print(i)
22         print(movieid)
23         break
```

```
1 def get_tagline(movie_id, api_key=API_KEY):
2     url = f'https://api.themoviedb.org/3/movie/{movie_id}?api_key={
3         api_key}'
4
5     response = requests.get(url)
6     if response.status_code == 200:
7         data = response.json()
8         tagline = data['tagline']
9         if tagline is not np.nan:
10             if len(tagline) == 0:
11                 return np.nan
12             return tagline
13     else:
14         return np.nan
15
16 missing_tagline_indices = df[df.tagline.isna()].index
17 for i in missing_tagline_indices:
18     movieid = df.loc[i, 'id']
19     try:
```



```
19         df.loc[i, 'tagline'] = get_tagline(movieid)
20     except:
21         print(i)
22         break
```

```
1 def get_production_companies(movie_id, api_key=API_KEY):
2     # Accessing the main movie details endpoint
3     url = f"https://api.themoviedb.org/3/movie/{movie_id}?api_key={
4         api_key}"
5     response = requests.get(url)
6
7     if response.status_code == 200:
8         data = response.json()
9         # Extracting production companies
10        production_companies = [company['name'] for company in data.get(
11            'production_companies', [])]
12
13        if not production_companies:
14            return np.nan
15        return str(production_companies)
16    else:
17        return np.nan
18
19 missing_production_companies_indices = df[df.production_companies.isna()
20 ].index
21 for i in missing_production_companies_indices:
22     movieid = df.loc[i, 'id']
23     try:
24         df.loc[i, 'production_companies'] = get_production_companies(
25             movieid)
26     except:
27         print(i)
28         break
```

```
1 def get_production_countries(movie_id, api_key=API_KEY):
2     # Accessing the main movie details endpoint
3     url = f"https://api.themoviedb.org/3/movie/{movie_id}?api_key={
4         api_key}"
5     response = requests.get(url)
6
7     if response.status_code == 200:
8         data = response.json()
9         # Extracting production companies
10        production_countries = [company['name'] for company in data.get(
11            'production_countries', [])]
12
13        if not production_countries:
```



```
12         return np.nan
13         return str(production_countries)
14     else:
15         return np.nan
16
17 missing_production_countries_indices = df[df.production_countries.isna()
18 ].index
19 for i in missing_production_countries_indices:
20     movieid = df.loc[i, 'id']
21     try:
22         df.loc[i, 'production_countries'] = get_production_countries(
23             movieid)
24     except:
25         print(i)
26         break
```

```
1 headers = headers = {
2     "accept": "application/json",
3     "Authorization": "Bearer eyJhbGciOiJIUzI1NiJ9.
4     eyJhdWQiOiJlYzYxOWZiMzgZMDcxMmUzNzY3ZTc1ODI4OThhODU5MiIsInN1YiI6IjYwOTQ0Mm
5     .zhJcElpbNKItZyIiRAWcJ9QvCXA59reAmEg2GZ9jQmqE"
6 }
7
8 def get_num_services(movie_id, headers=headers):
9     url = f"https://api.themoviedb.org/3/movie/{movie_id}/watch/
10     providers"
11     response = requests.get(url, headers=headers)
12     if response.status_code == 200:
13         data = response.json()['results']
14         stream_services = set()
15         rent_services = set()
16         buy_services = set()
17         for countryi in data.keys():
18             stream_services.update([i['provider_id'] for i in data[
19                 countryi].get('flatrate', [])])
20             rent_services.update([i['provider_id'] for i in data[
21                 countryi].get('rent', [])])
22             buy_services.update([i['provider_id'] for i in data[countryi
23                 ].get('buy', [])])
24
25         return len(list(stream_services)), len(list(rent_services)), len
26             (list(buy_services))
27     else:
28         return np.nan, np.nan, np.nan
29
30 num_stream_services, num_rent_services, num_buy_services = list(), list
31     (), list()
32 for idi in df.id:
33     results = get_num_services(idi)
```



```

25     num_stream_services.append(results[0])
26     num_rent_services.append(results[1])
27     num_buy_services.append(results[2])
28
29 df['NoStreamProviders'] = num_stream_services
30 df['NoRentProviders'] = num_rent_services
31 df['NoBuyProviders'] = num_buy_services

```

```

1 headers = headers = {
2     "accept": "application/json",
3     "Authorization": "Bearer eyJhbGciOiJIUzI1NiJ9.
        eyJhdWQiOiJlYzYxOWZiMzgzMDCxMmUzNzY3ZTc1ODI4OThhODU5MiIsInN1YiI6IjYwOTQOMm
        .zhJcEIpbNKItZyiRAWcJ9QvCXA59reAmEg2GZ9jQmqE"
4 }
5 def get_num_providers(movie_id, headers=headers):
6     url = f"https://api.themoviedb.org/3/movie/{movie_id}/watch/
        providers"
7     response = requests.get(url, headers=headers)
8     num_providers = set()
9     if response.status_code == 200:
10         data = response.json()['results']
11         num_countries = list(data.keys())
12         for key in num_countries:
13             for typei in ['buy', 'rent', 'flatrate']:
14                 ids = [i['provider_id'] for i in data[key].get(typei,
15                     [])]
16                 num_providers.update(ids)
17             return len(num_providers), len(num_countries)
18         else:
19             return np.nan, np.nan
20
21 num_global_providers = list()
22 num_countries = list()
23 for idi in df.id:
24     result = get_num_providers(idi)
25     num_global_providers.append(result[0])
26     num_countries.append(result[1])
27
28 df['NoTotalProviders'] = num_global_providers
29 df['NoCountriesAvailable'] = num_countries

```

```

1 df.isnull().sum().sort_values(ascending=False)
2
3 tmp = df.groupby('PrimaryGenre')['revenue'].mean().sort_values()
4 tmp = tmp / 1000000
5 fig, ax = plt.subplots(figsize=(8,4))

```



```

6 ax.barh(tmp.index, tmp.values)
7 ax.set_xlabel('Revenue (in millions of dollars)', fontsize=12)
8 ax.set_ylabel('Genre', fontsize=12)
9 ax.set_title('Average Revenue by Movie Genre', fontsize=18)
10 fig.savefig('mean_revenue_by_genre.png')
11 plt.show()

```

```

1 numeric_features = df.select_dtypes(exclude='object').drop(columns=['id'
    , 'movie_id', 'release_date']).columns.to_list()
2 print(len(numeric_features))
3 print(numeric_features)
4
5 fig, ax = plt.subplots(3, 2, figsize=(8,8))
6 for coli, axi in zip(numeric_features, ax.ravel()):
7     skew = stats.skew(df[coli], nan_policy='omit')
8     sns.histplot(df[coli], kde=True, ax=axi, label=f"Skew= {np.round(
        skew, 2)}")
9     axi.set_title(coli.capitalize(), fontsize=12)
10    axi.set_xlabel('')
11    axi.legend()
12
13 fig.suptitle('Distribution of Numerical Features', fontsize=18)
14 fig.tight_layout()
15 fig.savefig('numeric_distributions.png')
16 plt.show()

```

```

1 df['HTML_Poster'] = df['Poster'].apply(lambda x: f'')
2 df['HTML_Poster_100px'] = df['Poster'].apply(lambda x: f'')
3 HTML(df[['HTML_Poster_100px', 'title', 'budget', 'revenue']].head().
    to_html(escape=False))

```

```

1 tmp = df.PrimaryGenre.value_counts().sort_values()
2
3 fig, ax = plt.subplots(figsize=(8,4))
4 ax.barh(tmp.index, tmp.values)
5 ax.set_title('Count of Movie Genres', fontsize=16)
6 ax.set_xlabel('Count', fontsize=12)
7 ax.set_ylabel('Genre', fontsize=12)
8 fig.savefig('count_movie_genres.png')
9 plt.show()

```

```

1 df['LeadActor'] = df.cast.apply(lambda x: x[0] if x is not np.nan else
    np.nan)

```



```
2
3 tmp = df.LeadActor.value_counts().sort_values()[-15:]
4 fig, ax = plt.subplots(figsize=(8,4))
5 ax.barh(tmp.index, tmp.values)
6 ax.set_title('Most Common Lead Actors in a Movie', fontsize=16)
7 ax.set_xlabel('Count', fontsize=12)
8 ax.set_ylabel('Actor', fontsize=12)
9 fig.savefig('count_actors.png')
10 plt.show()
```

```
1 tmp = df.PrimaryDirector.value_counts().sort_values()[-15:]
2 fig, ax = plt.subplots(figsize=(8,4))
3 ax.barh(tmp.index, tmp.values)
4 ax.set_title('Most Common Directors for a Movie', fontsize=16)
5 ax.set_xlabel('Count', fontsize=12)
6 ax.set_ylabel('Director', fontsize=12)
7 fig.savefig('count_directors.png')
8 plt.show()
```

```
1 genre_summary = df.groupby('PrimaryGenre').agg({'budget': 'mean', '
    revenue': 'mean', 'profit': 'mean'}).reset_index()
2 genre_summary[['budget', 'revenue', 'profit']] = genre_summary[['budget',
    , 'revenue', 'profit']] / 1000000
3 genre_summary = genre_summary[~((genre_summary.PrimaryGenre == 'Foreign
    ') | (genre_summary.PrimaryGenre == 'TV Movie'))]
4
5 # Plotting
6 fig, ax = plt.subplots(figsize=(8,4))
7 sns.barplot(x='profit', y='PrimaryGenre', data=genre_summary, label="
    Profit", color="#145a32", zorder=1)
8 sns.barplot(x='revenue', y='PrimaryGenre', data=genre_summary, label="
    Revenue", color="#2e86c1", alpha=0.6, zorder=2)
9 sns.barplot(x='budget', y='PrimaryGenre', data=genre_summary, label="
    Budget", color="#4a235a", alpha=0.6, zorder=3)
10
11 ax.legend(ncol=3, loc="lower right", frameon=True)
12 ax.set_xlabel('Amount (in millions of dollars)', fontsize=12)
13 ax.set_ylabel('Genre', fontsize=12)
14 ax.set_title('Comparison of Budget, Revenue, and Profit by Genre',
    fontsize=16)
15 fig.savefig('genres_budget_revenue_profit.png')
16 plt.show()
```

```
1 indices = df[(df.release_year >= 2000 )& (df.release_year < 2017)].
    groupby('release_year')['revenue'].idxmax()
```



```
2 HTML(df.loc[indices, ['release_year', 'HTML_Poster_100px', 'title']].  
    set_index('release_year').to_html(escape=False))
```

```
1 fig, ax = plt.subplots(figsize=(15, 8)) # Increase figure size  
2 bar_width = 0.5 # Decrease bar width  
3 ax.bar(tmp.release_year, tmp.revenue, width=bar_width, color='skyblue')  
    # Adjust bar color  
4  
5 # Rotate x-ticks for better readability  
6 plt.xticks(rotation=45, ha='right')  
7  
8 # Add horizontal gridlines  
9 ax.yaxis.grid(True, linestyle='--', which='major', color='grey', alpha  
    =0.5)  
10  
11 # Set y-axis label  
12 ax.set_ylabel('Revenue (in millions of dollars)')  
13  
14 # Set title  
15 ax.set_title('Yearly Movie Revenue with Posters')  
16  
17 # Adjust the image placement and size  
18 for i, url in enumerate(tmp.Poster):  
19     img = get_image_from_url(url)  
20     imagebox = OffsetImage(img, zoom=0.2) # Adjust zoom as needed  
21     ab = AnnotationBbox(imagebox, (tmp.release_year[i], tmp.revenue[i]),  
        frameon=False) # Offset_value to adjust the image placement  
        above the bar  
22     ax.add_artist(ab)  
23 ax.set_ylim(0, ax.get_ylim()[1]+200)  
24 # Show the plot  
25 plt.tight_layout() # Adjust subplot params for better fit  
26 plt.show()
```

```
1 from matplotlib import patheffects  
2 fig, ax = plt.subplots(figsize=(8, 15)) # Increase figure size  
3 bar_height = 0.8 # This is now the height of the horizontal bars  
4  
5 # Create a horizontal bar chart  
6 ax.barh(tmp.release_year, tmp.revenue, height=bar_height, color='skyblue'  
    ,')  
7  
8 # Rotate y-ticks for better readability  
9 plt.yticks(rotation=45, ha='right')  
10  
11 # Add vertical gridlines
```



```

12 ax.xaxis.grid(True, linestyle='--', which='major', color='grey', alpha
    =0.5)
13
14 # Set x-axis label (since this is now the axis for revenue)
15 ax.set_xlabel('Revenue (in millions of dollars)')
16
17 # Set title
18 ax.set_title('Highest Revenue Movie by Year', fontsize=16)
19
20 # Adjust the image placement and size for horizontal bar chart
21 for i, url in enumerate(tmp.Poster):
22     img = get_image_from_url(url)
23     imagebox = OffsetImage(img, zoom=0.15) # Adjust zoom as needed
24     # Note: The y coordinate (tmp.release_year[i]) now comes before the
        x coordinate (tmp.revenue[i])
25     ab = AnnotationBbox(imagebox, (2700, tmp.release_year[i]), frameon=
        False, xycoords='data', pad=0.2)
26     ax.add_artist(ab)
27
28 text_path_effect = [patheffects.withStroke(linewidth=4, foreground='
    black'')]
29
30 for year, title in zip(tmp.release_year, tmp.title):
31     text = ax.text(20, year, title, color='white', fontsize=12,
32                   ha='left', va='center') # ha and va for horizontal
        and vertical alignment
33     text.set_path_effects(text_path_effect)
34
35
36 # Adjust the limits to fit images
37 ax.set_xlim(0, ax.get_xlim()[1]+200) # You may need to adjust this
    based on your data
38 ax.set_yticks(np.arange(2000, 2017))
39 # Show the plot
40 plt.tight_layout() # Adjust subplot params for better fit
41 fig.savefig('top_movie_by_revenue_by_year.png')
42 plt.show()

```

```

1 tmp = df[~df.PrimaryGenre.isin(['Foreign', 'Music', 'TV Movie'])].
    groupby(['release_month', 'PrimaryGenre'])['profit'].mean().unstack()
    .loc[
2     ['January', 'February', 'March', 'April', 'May', 'June', 'July', '
        August', 'September', 'October', 'November', 'December']].T /
        1000000
3 fig, ax = plt.subplots(figsize=(8,6))
4 sns.heatmap(tmp, cmap=sns.diverging_palette(220, 20, as_cmap=True), ax=
    ax)
5 ax.set_xlabel('Month of release', fontsize=12)

```




```
6 ax.set_ylabel('Genre', fontsize=12)
7 ax.set_title('Average Profit by Genre and Month', fontsize=16)
8 ax.collections[0].colorbar.set_label('Profit (in millions)')
9 plt.show()
```

```
1 tmp = df[~df.PrimaryGenre.isin(['Foreign', 'Music', 'TV Movie'])].
    groupby(['release_month', 'PrimaryGenre'])['BoxOfficeRevenue'].mean()
    .unstack().loc[
2     ['January', 'February', 'March', 'April', 'May', 'June', 'July', '
        August', 'September', 'October', 'November', 'December']].T /
        1000000
3 fig, ax = plt.subplots(figsize=(8,6))
4 sns.heatmap(tmp, cmap=sns.cubehelix_palette(start=2, rot=0, dark=0,
    light=.95, reverse=False, as_cmap=True), ax=ax)
5 ax.set_xlabel('Month of release', fontsize=12)
6 ax.set_ylabel('Genre', fontsize=12)
7 ax.set_title('Average Box Office Revenue by Genre and Month', fontsize
    =16)
8 ax.collections[0].colorbar.set_label('Revenue (in millions)')
9 fig.tight_layout()
10 fig.savefig('avg_box_office_revenue_heatmap.png')
11 plt.show()
```

```
1 import networkx as nx
2
3 # Create an empty graph
4 G = nx.Graph()
5
6 # Add nodes and edges
7 for index, row in df.iterrows():
8     movie = row['title']
9     actors = row['cast']
10
11     if movie is not np.nan and actors is not np.nan:
12         for actor in actors:
13             G.add_node(actor, type='actor')
14             G.add_node(movie, type='movie')
15             G.add_edge(actor, movie)
16
17 # # Optionally, differentiate node colors by type
18 # colors = ['red' if G.nodes[node]['type'] == 'actor' else 'blue' for
    node in G]
19
20 # # Draw the graph
21 # plt.figure(figsize=(8, 8)) # Adjusted for better visibility
22 # nx.draw(G, with_labels=False, node_color=colors, node_size=20,
    font_size=8, alpha=0.7) # with_labels set to False for clarity
```



```
23 # plt.title("Movie-Actor Network")
24 # plt.show()
25
26 nx.write_graphml(G, 'actors_movies.graphml')
```

```
1 G = nx.Graph()
2
3 for actors in df['cast']:
4     if actors is not np.nan:
5         for i in range(len(actors)):
6             for j in range(i + 1, len(actors)):
7                 actor_pair = (actors[i], actors[j])
8                 if G.has_edge(*actor_pair):
9                     G[actor_pair[0]][actor_pair[1]]['weight'] += 1
10                else:
11                    G.add_edge(actor_pair[0], actor_pair[1], weight=1)
12
13 nx.write_graphml(G, 'actor_pairs.graphml')
14
15 sorted_edges = sorted(G.edges(data=True), key=lambda x: x[2]['weight'],
16                        reverse=True)
17
18 N = 5 # for example, to get the top 5 pairs
19 for edge in sorted_edges[:N]:
20     print(f"{edge[0]} and {edge[1]} appeared together in {edge[2]['weight']} movies.")
```

```
1 def plot_test_results(male, female, title, savefig=None):
2     fig = plt.figure(figsize=(10,6))
3     ax1 = plt.subplot2grid((3,2), (0,0), rowspan=2)
4     ax2 = plt.subplot2grid((3,2), (0,1), rowspan=2)
5     ax3 = plt.subplot2grid((3,2), (2,0), colspan=2)
6
7     sns.kdeplot(male, common_norm=False, fill=True, ax=ax1, color=sns.
8                 color_palette()[0], label='Male')
9     sns.kdeplot(female, common_norm=False, fill=True, ax=ax1, color=sns.
10                color_palette()[1], label='Female')
11
12     medians = pd.Series([male.median(), female.median()], index=['Male',
13                          'Female'])
14     errors = pd.Series([male.sem(), female.sem()], index=['Male', 'Female'])
15     medians.plot(kind='bar', yerr=errors, capsize=4, ax=ax2, color=sns.
16                  color_palette()[:2])
17
18     u_statistic, p_value = stats.mannwhitneyu(male.dropna(), female.
19        dropna(), alternative='two-sided')
```



```

15     if p_value < 0.05:
16         result = 'Statistically Significant Result'
17     else:
18         result = 'Non-statistically Significant Result'
19
20
21     ax3.set_xticks([])
22     ax3.set_yticks([])
23     ax3.text(0.02,0.7, 'Mann-Whitney U Test Results', fontdict={'fontsize':16, 'fontweight': 'semibold'})
24     ax3.text(0.02, 0.5, f"U-statistic: {u_statistic}, P-value: {p_value}", fontdict={'fontsize':14})
25     ax3.text(0.02, 0.3, result, fontdict={'fontsize':14})
26     rect = Rectangle((0, 0), 1, 1, linewidth=3, edgecolor='black',
27                     facecolor='none', linestyle='-')
28     ax3.add_patch(rect)
29     legend_handles = [Patch(facecolor=sns.color_palette()[0], label='Male'),
30                       Patch(facecolor=sns.color_palette()[1], label='Female')]
31
32     ax3.legend(handles=legend_handles, loc='lower right', title='Gender')
33
34     ax1.set_xlabel('Profit (in millions)', fontsize=12)
35     ax1.set_ylabel('Density', fontsize=12)
36     ax2.set_xlabel('Gender', fontsize=12)
37     ax2.set_ylabel('Median Profit (in millions)', fontsize=12)
38     ax1.set_title('Density Plot', fontsize=18)
39     ax2.set_title('Median Difference Plot', fontsize=18)
40     fig.suptitle(title, fontsize=18, fontweight='semibold')
41     fig.tight_layout()
42     if savefig is not None:
43         fig.savefig(savefig)
44     plt.show()
45
46 plot_test_results(df[df.LeadActorGender == 'Male']['profit']/1000000, df
47                  [df.LeadActorGender == 'Female']['profit']/1000000, 'Significance of
48                  Lead Actor Gender on Profit', savefig='test_actor_gender.png') #
49 plot_test_results(df[df.PrimaryDirectorGender == 'Male']['profit']
50                  /1000000, df[df.PrimaryDirectorGender == 'Female']['profit']
51                  /1000000, 'Significance of Director Gender on Profit', savefig='
52                  test_director_gender.png') #

```

```

1 tmp = df.groupby('PrimaryProductionCompany')['revenue'].mean().
2   sort_values(ascending=False)[:20][::-1]
3 fig, ax = plt.subplots(figsize=(8,4))
4 ax.barh(tmp.index, tmp.values)

```



```
4 ax.set_title('Top 20 Production Companies by Mean Revenue', fontsize=16)
5 ax.set_xlabel('')
```

```
1 tmp = df.groupby('PrimaryProductionCompany')['revenue'].mean().
    sort_values(ascending=False)[:20][::-1] / 1000000
2 fig, ax = plt.subplots(figsize=(8,4))
3 ax.barh(tmp.index, tmp.values)
4 ax.set_title('Top 20 Production Companies by Mean Revenue', fontsize=16)
5 ax.set_xlabel('')
```

B.2 R

```
1 library(dplyr)
2 library(lubridate)
3 library(ggplot2)
4 library(scales)
5
6 # For inflation adjustments
7 library(priceR)
8
9 # For my Sankey
10 library(highcharter)
11 library(htmlwidgets)
12 library(webshot)
13
14 # For pivot_longer()
15 library(tidyverse)
16
17 # For linear correlation annotation on scatterplots
18 library(ggpubr)
19
20 # A couple of my libraries conflict in a way that forced me to specify
21 # "dplyr::" for aggregation functions.
```

```
1 # Loading and preprocessing data.
2 movie_df_full <- read.csv("movies_complete.csv")
3 head(movie_df_full)
```

```
1 ### Creating a subset that primarily contains the attributes I'm
    interested in
2 # First, selecting specific attributes
```



```
3 mov_df <- movie_df_full %>% select(title, release_date, revenue, budget,
  popularity, vote_count, vote_average, runtime, original_language,
  PrimaryGenre, PrimaryProductionCompany, PrimaryProductionCountry )
```

```
1 ### Creating *release month* and *release year* from existing attribute
  *release date*
2 # From the "release_date" attribute, I create several attributes:
3 # release_month, release_year, and release_mon_yr .
4
5 mov_df <- mov_df %>%
6   mutate(release_date = as.POSIXct(strptime(release_date, '%Y-%m-%d')),
7         release_month = month(release_date),
8         release_year = year(release_date))
9
10 head(mov_df)
11
12 mov_df$release_mon_yr = format(as.Date(mov_df$release_date), "%m-%Y")
13
14 # head(mov_df, 10)
15 mov_df %>%
16   count(release_year)
```

```
1 ### Removing data older than 1996 and newer than 2015
2 # 2015 is the most recent year of complete data in this set.
3
4 mov_df <- mov_df %>%
5   filter(release_year <= 2015 & release_year >= 2001)
6
7 mov_df %>%
8   count(release_year)
```

```
1 ### Removing some records with invalid *revenue* and *budget* values
2 # The dataset contains values of zero for revenue and budget that strike
  us as
3 # invalid. Removing tuples that show those values.
4
5 nrow(mov_df)
6
7 mov_df <- mov_df %>%
8   filter(revenue != 0 & budget != 0)
9
10 nrow(mov_df)
```



```
1 ### Defining the function for inflation adjustment
2 country <- "US"
3 inflation_dataframe <- retrieve_inflation_data(country)
4 countries_dataframe <- show_countries()
5
6 # 2022 is the most recent year that this function works for.
7 adjust_for_inflation(100, 2016, country, 2022,
8                       inflation_dataframe = inflation_dataframe,
9                       countries_dataframe = countries_dataframe)
10
11
12 # The cpiLookup function from Tutorial 4 doesn't work with our dataset,
13   which shows repeated release years. The below function cannot be used
14   with mutate() on such a dataset.
15
16 # cpiLookup <- function(y) {
17 #   return(CPITable %>%
18 #     filter(Year %in% y) %>%
19 #     select(Value) %>%
20 #     .$Value ) }
```

```
1 # Creating inflation-adjusted attributes
2
3 mov_df <- mov_df %>%
4   mutate(adjustedRev = adjust_for_inflation(revenue, release_year,
5                                             country, 2022,
6                                             inflation_dataframe = inflation_dataframe,
7                                             countries_dataframe = countries_dataframe))
8
9
10 mov_df <- mov_df %>%
11   mutate(adjustedBudget = adjust_for_inflation(budget, release_year,
12                                                country, 2022,
13                                                inflation_dataframe = inflation_dataframe,
14                                                countries_dataframe = countries_dataframe))
15
16 mov_df %>%
17   select(release_year, budget, adjustedBudget)
```

```
1 # After adjusting for inflation, which movies have the top revenues
2 mov_df %>%
3   select(title, adjustedRev, release_year) %>%
4   arrange(-adjustedRev)
```



```
1 ### Creating an attribute for the quartile of inflation-adjusted revenue
2 breaks_adjRev_qrts <- quantile(mov_df$adjustedRev, c(0.75))
3 breaks_adjRev_qrts
4
5 # Binning accordingly (but approximating the break values with even
6   numbers that are easier for viewers to read.)
7 mov_df <- mov_df %>%
8   mutate(adjRev_quartile = cut(adjustedRev,
9     breaks = c(0, 209953708, 10000000000),
10    labels = c("Lower 3 quartiles of revenue"
11      ,
12      "Top revenue quartile"),
13    include.lowest=TRUE))
```

```
1 ### Creating a *release season* attribute
2 # Defining a function getSeason (found on StackOverflow) and use it to
3   create a release_season attribute from release_date.
4
5 getSeason <- function(input.date){
6   numeric.date <- 100*month(input.date)+day(input.date)
7   ## input Seasons upper limits in the form MMDD in the "break =" option
8   :
9   cuts <- base::cut(numeric.date, breaks = c(0,319,0620,0921,1220,1231))
10  # rename the resulting groups (could've been done within cut(...levels
11    =) if "Winter" wasn't double
12  levels(cuts) <- c("Winter","Spring","Summer","Fall","Winter")
13  return(cuts)
14 }
15
16 mov_df <- mov_df %>%
17   mutate(release_season = getSeason(release_date))
```

```
1 # Creating attribute that tracks whether the season of release was
2   warmer or
3   cooler. Used to analyze how some genres thrive in warmer months.
4
5 mov_df <- mov_df %>%
6   mutate(warmSzn_movie = ifelse(release_season %in% c("Spring", "Summer"
7     ),
8     paste("Warmer Season (", PrimaryGenre, "
9       )"),
10     paste("Cooler Season (", PrimaryGenre, "
11       )"))))
12
13 head(mov_df, 10)
```



```
11 # Writing this preprocessed data to file
12 write.csv(mov_df, "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul MS Work
    \\DSC 465\\Homework\\HW_4\\mov_df_processed.csv", row.names = FALSE)
```

```
1 # Prepping data for Sankey, specifically
2 # Including only film records that show a PrimaryGenre that is top-
    selling
3
4 # Identifying the Genres that earned the most revenue in this period
5 mov_df %>%
6   dplyr::group_by(PrimaryGenre) %>%
7   dplyr::summarise(total_rev = sum(adjustedRev)) %>%
8   arrange(-total_rev) %>%
9   head(10)
```

```
1 # Improved Seasonal Sankey
2 # Narrowing dataset to include only records associated with a top-5-
    earning genre.
3 summer_movie_sankey_df <- mov_df %>%
4   filter(PrimaryGenre %in% c("Family", "Science Fiction", "Adventure", "
    Animation"))
```

```
1 # TESTING WARMER COOLER SEASON
2 summer_movie_sankey_df <- summer_movie_sankey_df %>%
3   select(PrimaryGenre, warmSzn_movie, adjRev_quartile)
4
5
6 summer_movie_Sankey <- hchart(data_to_sankey(summer_movie_sankey_df), "
    sankey", name = "Relating season of release, genre, and adj-revenue
    quartile")
7
8 summer_movie_Sankey <- summer_movie_Sankey %>%
9   hc_title(text = "Seasonal blockbusting rate is distinct from seasonal
    average") %>%
10  hc_subtitle(text = "Though these genres might all register higher
    average revenues during the warmer months, their proportions of
    blockbusters in that time don't necessarily improve.")
11
12 summer_movie_Sankey
```

```
1 # Saving the sankey HTML and png files
2 htmlwidgets::saveWidget(widget = summer_movie_Sankey, file = 'C:\\Users
    \\maxru\\OneDrive\\Documents\\DePaul MS Work\\DSC 465\\Homework\\HW_
    4\\summer_blockbuster_sankey.html')
```




```
3
4
5 webshot::install_phantomjs()
6
7 webshot::webshot(url = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul MS
  Work\\DSC 465\\Homework\\HW_4\\summer_blockbuster_sankey.html",
8               file = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul MS
  Work\\DSC 465\\Homework\\HW_4\\summer_blockbuster_
  sankey.png" , delay = 5)
```

```
1 ## SUMMER MOVIE SHARE CHART
2
3 # Summer genre count shares
4 mov_summer_df <- mov_df
5
6 mov_summer_df %>%
7   filter(release_season %in% c("Spring", "Summer"))
8
9 mov_summer_df <- mov_summer_df %>%
10  mutate(summerGenre = ifelse(PrimaryGenre %in% c("Adventure", "Family",
11                                                  "Science Fiction", "Animation
12                                                  ", "Thriller"),
13                                PrimaryGenre ,
14                                "Other"))
15
16 mov_summer_df$summerGenre <- factor(mov_summer_df$summerGenre,
17                                   levels = c("Adventure", "
18                                   Animation",
19                                   "Science Fiction", "
20                                   Family",
21                                   "Thriller", "Other"))
22
23
24 # summer_counts_df
25
26 movie_count <- nrow(mov_summer_df)
27
28 summer_counts_df <- summer_counts_df %>%
29   mutate(count_share = n / movie_count)
30
31 summer_counts_df
```

```
1 # Summer genre revenue shares
2 summer_revs_df <- mov_summer_df %>%
```



```
3 group_by(summerGenre) %>%
4 summarise(tot_adjRev = sum(adjustedRev)) %>%
5 arrange(-tot_adjRev)
6
7 summer_adjRev_total = sum(summer_revs_df$tot_adjRev)
8
9 summer_revs_df <- summer_revs_df %>%
10 mutate(rev_share = tot_adjRev / summer_adjRev_total)
11
12 summer_revs_df
```

```
1 # Summer genre popularity shares
2 summer_pops_df <- mov_summer_df %>%
3 group_by(summerGenre) %>%
4 summarise(tot_pop = sum(popularity)) %>%
5 arrange(-tot_pop)
6
7 summer_pop_total_all = sum(summer_pops_df$tot_pop)
8
9 summer_pops_df <- summer_pops_df %>%
10 mutate(pop_share = tot_pop / summer_pop_total_all)
11
12 summer_pops_df
```

```
1 # joining these aggregations, message = FALSE
2 summer_mov_proport_plot_df <- summer_counts_df %>%
3 inner_join(summer_revs_df) %>%
4 inner_join(summer_pops_df) %>%
5 select(summerGenre, count_share, rev_share, pop_share)
6
7 summer_mov_proport_plot_df
8
9 summer_mov_proport_plot_df <- summer_mov_proport_plot_df %>%
10 pivot_longer(!summerGenre, names_to = "Film_Measure", values_to = "
    Share")
11
12 summer_mov_proport_plot_df <- summer_mov_proport_plot_df %>%
13 mutate(Film_Measure = fct_recode(as.factor(Film_Measure),
14                                   "Sample Count" = "count_share",
15                                   "Revenue" = "rev_share",
16                                   "Popularity" = "pop_share"))
17
18
19
20
21 levels(summer_mov_proport_plot_df$Film_Measure) <- c("Sample Count", "
    Revenue", "Popularity")
```



```
22 levels(summer_mov_proport_plot_df$Film_Measure)
23
24
25 summer_mov_proport_plot_df <- summer_mov_proport_plot_df %>%
26   mutate(Share_rnd = round(Share, 3))
27
28 summer_mov_proport_plot_df
```

```
1 # FINAL VIZ
2 proport_plot_palette <- c("Action"="#94aeca", "Adventure"="#f7ba7e",
3   "Comedy"="#cfaec6", "Drama"="#ed999a", "Animation"="#
4     acd3d0",
5     "Other"="#d5cfcd", "Thriller"="#5a5858", "Science
6     Fiction"="#B9fdbbe", "Family"="#Cfabfb")
7
8 summer_mov_proport_plot_df %>%
9   filter(summerGenre != "Other") %>%
10  ggplot(aes(x = summerGenre, y = Share, fill = summerGenre)) +
11    geom_bar(stat='identity') +
12    facet_wrap(~Film_Measure) +
13    geom_text(aes(label = Share_rnd),
14      color = "white", size = 4.5,
15      vjust = 1, fontface = "bold") +
16    scale_y_continuous(labels = scales::percent) +
17    scale_fill_manual(name = "Genre", values = proport_plot_palette) +
18    guides(fill = guide_legend(ncol = 5)) +
19    ggtitle("Some Genres Thrive in Warmer Weather") +
20    theme_minimal() +
21    theme(plot.title = element_text(hjust = 0.5, size = 20,
22      vjust = 2, face = "bold"),
23      plot.subtitle = element_text(hjust = 0.5),
24      legend.position="bottom",
25      legend.box = "horizontal",
26      legend.title = element_text(size = 16),
27      legend.text = element_text(size = 12),
28      axis.title.x=element_blank(),
29      axis.text.x=element_blank(),
30      axis.ticks.x=element_blank(),
31      strip.text = element_text(size = 16),
32      panel.grid.major.x = element_blank()) +
33    labs(subtitle = "Shares by genre of various totals from Spring and
34      Summer releases, years 2001 to 2015.")
```

```
1 # Making the (previous) Sankey
2
3 # Narrowing dataset to include only records associated with a top-5-
4   earning genre.
```



```
4 mov_df_top_gens <- mov_df %>%
5   filter(PrimaryGenre %in% c("Action", "Adventure", "Drama", "Comedy", "
6     Animation"))
7 # TESTING WARMER COOLER SEASON
8 season_genre_sankey_df <- mov_df_top_gens %>%
9   select(PrimaryGenre, warmSzn_movie, adjRev_quartile)
10
11
12 season_genre_Sankey <- hchart(data_to_sankey(season_genre_sankey_df), "
13   sankey", name = "Relating season of release, genre, and adj-revenue
14   quartile")
15
16 season_genre_Sankey <- season_genre_Sankey %>%
17   hc_title(text = "Season of Release, Genre, and Revenue Quartile") %>%
18   hc_subtitle(text = "Examining how release season, genre, and revenue
19     quartile relate in terms of distribution. Revenues were inflation-
20     adjusted to their 2022 equivalents.")
21
22 season_genre_Sankey
```

```
1 season_genre_sankey_df <- mov_df_top_gens %>%
2   select(release_season, PrimaryGenre, adjRev_quartile)
3
4
5 season_genre_Sankey <- hchart(data_to_sankey(season_genre_sankey_df), "
6   sankey", name = "Relating season of release, genre, and adj-revenue
7   quartile")
8
9 season_genre_Sankey <- season_genre_Sankey %>%
10   hc_title(text = "Season of Release, Genre, and Revenue Quartile") %>%
11   hc_subtitle(text = "Examining how release season, genre, and revenue
12     quartile relate in terms of distribution. Revenues were inflation-
13     adjusted to their 2022 equivalents.")
14
15 season_genre_Sankey
```

```
1 # Saving the sankey HTML and png files
2 htmlwidgets::saveWidget(widget = season_genre_Sankey, file = 'C:\\Users
3   \\maxru\\OneDrive\\Documents\\DePaul MS Work\\DSC 465\\Homework\\HW_
4   4\\season-genre-adjRevQ_sankey_all_years.html')
5
6
7 # webshot::install_phantomjs()
8
9 webshot::webshot(url = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul MS
10   Work\\DSC 465\\Homework\\HW_4\\season-genre-adjRevQ_sankey_all_
```



```
years.html",
file = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul
MS Work\\DSC 465\\Homework\\HW_4\\season-genre-
adjRevQ_sankey_all_years.png" , delay = 5)
```

```
1 # Original Proportionality-by-Genre Bar Chart (of the 5 top-grossing
  genres)
2
3 # proportionality by genre
4 mov_df2 <- mov_df
5
6 # head(mov_df2)
7
8 mov_df2 <- mov_df2 %>%
9   mutate(bigGenre = ifelse(PrimaryGenre %in% c("Action", "Adventure", "
    Drama",
10                                     "Comedy", "Animation"),
11                             PrimaryGenre,
12                             "Other"))
13
14 mov_df2_gen_cnts <- mov_df2 %>%
15   count(bigGenre) %>%
16   arrange(-n)
17
18 movie_count = nrow(mov_df2)
19
20 mov_df2_gen_cnts <- mov_df2_gen_cnts %>%
21   mutate(count_share = n / movie_count)
22
23 mov_df2_gen_cnts
```

```
1 # Revenue shares by genre
2 mov_df2_rev_sums <- mov_df2 %>%
3   group_by(bigGenre) %>%
4   summarise(tot_adjRev = sum(adjustedRev)) %>%
5   arrange(-tot_adjRev)
6
7 adjRev_total = sum(mov_df2_rev_sums$tot_adjRev)
8
9 mov_df2_rev_sums <- mov_df2_rev_sums %>%
10   mutate(rev_share = tot_adjRev / adjRev_total)
11
12 mov_df2_rev_sums
```

```
1 # Popularity shares by genre
```



```
2 mov_df2_pop_sums <- mov_df2 %>%
3   group_by(bigGenre) %>%
4   summarise(tot_pop = sum(popularity)) %>%
5   arrange(-tot_pop)
6
7 pop_total_all = sum(mov_df2_pop_sums$tot_pop)
8
9 mov_df2_pop_sums <- mov_df2_pop_sums %>%
10  mutate(pop_share = tot_pop / pop_total_all)
11
12 mov_df2_pop_sums
```

```
1 # adjBudget shares by genre
2 mov_df2_budg_sums <- mov_df2 %>%
3   group_by(bigGenre) %>%
4   summarise(tot_adjBudg = sum(adjustedBudget)) %>%
5   arrange(-tot_adjBudg)
6
7 adjBudg_total_all = sum(mov_df2_budg_sums$tot_adjBudg)
8
9 mov_df2_budg_sums <- mov_df2_budg_sums %>%
10  mutate(budg_share = tot_adjBudg / adjBudg_total_all)
11
12 mov_df2_budg_sums
```

```
1 # joining these aggregations, message = FALSE
2 mov_df_proportional_plot_df <- mov_df2_gen_cnts %>%
3   inner_join(mov_df2_rev_sums) %>%
4   inner_join(mov_df2_pop_sums) %>%
5   inner_join(mov_df2_budg_sums) %>%
6   select(bigGenre, count_share, rev_share, pop_share, budg_share)
```

```
1 # pivoting-longer these joined tables
2 mov_df_proportional_plot_df <- mov_df_proportional_plot_df %>%
3   pivot_longer(!bigGenre, names_to = "Film_Measure", values_to = "Share")
4
5 mov_df_proportional_plot_df
```

```
1 ## Visualizing shares
2
3 # Rough sketch of genre proportions across the various measures
4 mov_df_proportional_plot_df %>%
5   ggplot(aes(x = bigGenre, y = Share)) + geom_bar(stat='identity') +
6     facet_wrap(~Film_Measure)
```



```
1 # Linear correlations between Revenue and Popularity, by Genre
2
3 # Final Viz: Linear corrs by genre
4 corrsByGenre_PopVsRev <- mov_df2 %>%
5   ggplot(aes(x = adjustedRev, y = popularity, color = bigGenre)) +
6   geom_point(alpha = 0.1) +
7   geom_smooth(method = lm, se = FALSE, color = "black", size = 0.2) +
8   stat_cor(p.accuracy = 0.001, r.accuracy = 0.01,
9           inherit.aes = FALSE, aes(x = adjustedRev, y = popularity)) +
10  facet_wrap(~bigGenre) +
11  scale_x_continuous(breaks = 500000000 * (0:4), labels = c("0", "0.5",
12    "1", "1.5", "2"),
13    limits = c(0, 2000000000)) +
14  scale_y_continuous(limits = c(0, 320)) +
15  guides(color = FALSE) +
16  labs(x = "Revenue ($billions)",
17       y = "Popularity (TMDB Score)",
18       title = "Correlations between Revenue and Popularity by Genre")
19 corrsByGenre_PopVsRev
```

```
1 # RECOLORING: Linear corrs by genre
2 # Important Animation genre has been too lightly colored in this scheme.
3 # But
4 # I've already recolored all the other visualizations... :(
5
6 corrsByGenre_PopVsRev <- mov_df2 %>%
7   ggplot(aes(x = adjustedRev, y = popularity, color = bigGenre)) +
8   geom_point(alpha = 0.18) +
9   geom_smooth(method = lm, se = FALSE, color = "black", size = 0.2) +
10  stat_cor(p.accuracy = 0.001, r.accuracy = 0.01,
11          inherit.aes = FALSE, aes(x = adjustedRev, y = popularity)) +
12  facet_wrap(~bigGenre) +
13  scale_color_manual(name = "Genre", values = top5_palette) +
14  scale_x_continuous(breaks = 500000000 * (0:4), labels = c("0", "0.5",
15    "1", "1.5", "2"),
16    limits = c(0, 2000000000)) +
17  scale_y_continuous(limits = c(0, 320)) +
18  guides(color = FALSE) +
19  labs(x = "Revenue ($billions)",
20       y = "Popularity (TMDB Score)",
21       title = "Correlations between Revenue and Popularity by Genre")
22 corrsByGenre_PopVsRev
```



```
1 # Popularity versus Revenue overall, and its density
2
3 ## This is a two-part visualization
4
5 Part 1 - Popularity vs Revenue overall
6 pop_v_adjRev <- mov_df2 %>%
7   ggplot(aes(x = adjustedRev, y = popularity)) +
8   geom_point(alpha = 0.2) +
9   annotate("rect", xmin = 0, xmax = 300000000, ymin = 0, ymax = 125,
10          color = "cyan", fill = NA) +
11   geom_smooth(method = lm, se = FALSE, color = "red", size = 0.2) +
12   stat_cor(p.accuracy = 0.001, r.accuracy = 0.01,
13          inherit.aes = FALSE, aes(x = adjustedRev, y = popularity),
14          color = "red") +
15   scale_x_continuous(labels = label_comma()) +
16   labs(title = "Popularity vs Revenue",
17        x = "Revenue ($)",
18        y = "Popularity (TMDB Score)")
19 pop_v_adjRev
20
21 # The light blue box indicates what area will be featured in the 2D-
22 # density plot
23 # that follows.
```

```
1 # Part 2 - "Zooming in" with a 2D Density Plot
2 pop_v_adjRev_zoomedUp <- mov_df2 %>%
3   ggplot(aes(x = adjustedRev, y = popularity)) +
4   stat_density2d() +
5   geom_point(alpha = 0.08) +
6   annotate("rect", xmin = 0, xmax = 300000000, ymin = 0, ymax = 125,
7          color = "cyan", fill = NA) +
8   # scale_x_continuous(labels = label_comma())
9   scale_x_continuous(limits = c(0, 300000000),
10          labels = label_comma()) +
11   scale_y_continuous(limits = c(0, 125)) +
12   labs(title = "Zooming in on Popularity vs Revenue - A Density Plot",
13        x = "Revenue ($)",
14        y = "Popularity (TMDB Score)")
15
16 pop_v_adjRev_zoomedUp
```

```
1 # Final Viz: Popularity versus Revenue overall, and its density
2 comboPlot_popVsRev <- ggarrange(pop_v_adjRev, pop_v_adjRev_zoomedUp,
3   nrow = 2)
```




```
4 comboPlot_popVsRev
```

```
1 # Stacked Area Chart Redo
2 # Creating a dedicated subset of data
3 # Creating a df specific to this viz
4
5 stackedArea_df <- mov_df2
6
7 # Setting levels to my Genre field, so that the desired order will be
8   reflected
9 # by the area chart.
10
11 # mov_df2_rev_sums
12
13 stackedArea_df$bigGenre <- factor(stackedArea_df$bigGenre,
14                                   levels = c("Action", "Adventure", "
15                                               Comedy",
16                                               "Drama", "Animation", "
17                                               Other"))
18
19 # Cutting df so that the oldest year is 2001
20
21 stackedArea_df <- stackedArea_df %>%
22   filter(release_year >= 2001)
```

```
1 stackedArea_df2 <- stackedArea_df %>%
2   group_by(release_year, bigGenre) %>%
3   summarise(adjRevSum = sum(adjustedRev)) %>%
4   select(release_year, bigGenre, adjRevSum) %>%
5   arrange(release_year)
```

```
1 # Recoloring and reordering area chart
2 top5_palette <- c("Action"="#94aeca", "Adventure"="#f7ba7e",
3                   "Comedy"="#cfaec6", "Drama"="#ed999a", "Animation"="#
4                   acd3d0",
5                   "Other"="#d5cfcd")
6
7 plot_stackedArea_revByGen <- stackedArea_df2 %>%
8   ggplot(aes(x = release_year, y = adjRevSum, fill = bigGenre)) +
9   geom_area() +
10   scale_fill_manual(values = top5_palette, name = "Genre") +
11   labs(title = "Total Revenue of Top-Grossing Genres (2001-2015)",
12        y = "Revenue ($ billion)",
13        x = "Year") +
14   scale_y_continuous(breaks = 10000000000 * (0:3), labels = c("0", "10",
15                                                                "20", "30")) +
```



```

14 scale_x_continuous(breaks = seq(1, 2015, by=2)) +
15 theme_minimal() +
16 theme(panel.grid.major.x = element_blank(),
17        panel.grid.minor.x = element_blank())
18
19 plot_stackedArea_revByGen

```

```

1 # Previous attempts and scratch visualizations
2
3 # Final visualization of genre proportions across various measures
4 mov_df_proportional_plot_df$bigGenre <-
5   factor(mov_df_proportional_plot_df$bigGenre,
6          levels = c('Other', 'Action', 'Adventure', 'Comedy', 'Drama', '
7                     Animation'))
8
9 mov_df_proportional_plot_df <- mov_df_proportional_plot_df %>%
10   mutate(Share_rnd = round(Share, 2))
11
12 mov_df_proportional_plot_df %>%
13   filter(Film_Measure != "budg_share") %>%
14   ggplot(aes(fill = bigGenre, x = bigGenre, y = Share)) +
15   geom_bar(position="dodge", stat='identity') +
16   geom_text(aes(label = Share_rnd),
17             color = "white", size = 4.5,
18             vjust = 1.5, fontface = "bold") +
19   facet_wrap(~factor(Film_Measure,
20                      levels = c('count_share', 'rev_share', 'pop_share')
21                                ,
22                                labels = c('Sample Count', 'Revenue', 'Popularity')
23                                )) +
24   scale_y_continuous(labels = scales::percent,
25                      limits = c(0, 0.3)) +
26   scale_x_discrete(guide = guide_axis(n.dodge = 3)) +
27   scale_fill_brewer(name = "Genre", palette = "Set2") +
28   guides(fill = guide_legend(ncol = 6)) +
29   ggtitle("Genre Proportions of Several Measures") +
30   theme_minimal() +
31   theme(plot.title = element_text(hjust = 0.5, size = 20,
32                                   vjust = 2, face = "bold"),
33         legend.position="bottom",
34         legend.box = "horizontal",
35         legend.title = element_text(size = 18),
36         legend.text = element_text(size = 13),
37         axis.title.x=element_blank(),
38         axis.text.x=element_blank(),
39         axis.ticks.x=element_blank(),
40         strip.text = element_text(size = 16),
41         panel.grid.major.x = element_blank())

```



```
1 # 2D-Density plot specific to genre
2 # I was coding this to be iterated through for each genre, eventually
   arranging
3 # them together as six plots in one image. However, I decided against it
   ,
4 # because I thought it more informative to keep all the axes consistent
   (as
5 # facet_wrap does automatically.)
6
7 the_var = "Action"
8 mov_df2 %>%
9   filter(bigGenre == the_var) %>%
10
11   ggplot(aes(x = adjustedRev, y = popularity)) +
12   stat_density2d(aes(colour=..level..)) +
13   scale_x_continuous(labels = label_comma()) +
14   labs(title = the_var) +
15   theme(plot.title = element_text(hjust = 0.5, size = 14))
```

```
1 # sketching "corr by genre" viz: manually zooming in, away from outliers
2 mov_df2 %>%
3   ggplot(aes(x = adjustedRev, y = popularity, color = bigGenre, alpha =
   0.05)) +
4   geom_point() +
5   facet_wrap(~bigGenre) +
6   scale_x_continuous(limits = c(0, 2000000000)) +
7   scale_y_continuous(limits = c(0, 320))
```

```
1 # sketching "corr by genre" viz: Adding regression lines, final polishes
2 mov_df2 %>%
3   # ggplot(aes(x = adjustedRev, y = popularity, color = bigGenre, alpha
   = 0.05)) +
4   # geom_point() +
5   ggplot(aes(x = adjustedRev, y = popularity, color = bigGenre)) +
6   geom_point(alpha = 0.1) +
7   geom_smooth(method = lm, se = FALSE, color = "black", size = 0.2) +
8   stat_cor(p.accuracy = 0.001, r.accuracy = 0.01,
9             inherit.aes = FALSE, aes(x = adjustedRev, y = popularity)) +
10   facet_wrap(~bigGenre) +
11   scale_x_continuous(breaks = 500000000 * (0:4), labels = c("0", "0.5",
   "1", "1.5", "2"),
12                       limits = c(0, 2000000000)) +
13   scale_y_continuous(limits = c(0, 320)) +
```



```
14 guides(color = FALSE) +
15 labs(x = "Revenue ($billions)",
16      y = "Popularity",
17      title = "Correlations between Revenue and Popularity by Genre")
```

```
1 ## Desperate attempts at faceted density plots
2 Part 1 - Levels only
3 mov_df2 %>%
4   ggplot(aes(x = adjustedRev, y = popularity)) +
5   stat_density2d(aes(colour=..level..)) +
6   facet_wrap(~bigGenre) +
7   scale_x_continuous(limits = c(0, 20000000000)) +
8   scale_y_continuous(limits = c(0, 320))
9
10 Part 2 - Levels, with faint points
11 mov_df2 %>%
12   ggplot(aes(x = adjustedRev, y = adjustedBudget)) +
13   stat_density2d(aes(colour=..level..)) +
14   geom_point(color = "yellow", alpha = 0.02) +
15   scale_x_continuous(labels = label_comma()) +
16   scale_y_continuous(labels = label_comma()) +
17   facet_wrap(~bigGenre)
18
19 Part 3 - Alpha tile-densities, plus points
20 mov_df2 %>%
21   ggplot(aes(x = adjustedRev, y = popularity)) +
22   geom_point(aes(color = bigGenre)) +
23   stat_density2d(aes(alpha=..density..),
24                 geom="tile",
25                 contour=FALSE) +
26   facet_wrap(~bigGenre) +
27   scale_x_continuous(limits = c(0, 20000000000)) +
28   scale_y_continuous(limits = c(0, 320))
```

```
1 ## Early Panel Plot
2
3 # Horizontal panel plot: Popularity vs Revenue by Genre
4 # I found it too wide for the popularity values to be decoded
5 # meaningfully, even
6 # after removing several outliers of adjusted revenue.
7
8 mov_df_top_gens %>%
9   ggplot(aes(x = adjustedRev, y = popularity)) + geom_point() + facet_
10     grid(PrimaryGenre ~ .)
11
12 mov_df_top_gens %>%
```



```
12 filter(adjustedRev < 2000000000) %>%
13 ggplot(aes(x = adjustedRev, y = popularity)) + geom_point() + facet_
    grid(PrimaryGenre ~ .)
```

```
1 # Sankey drafts
2 ## Genre, Revenue Quartile, and Top 10 Production Companies
3 ### Note that this uses raw revenue values (without inflation
4     adjustment.)
5 ### This work was done before inflation-adjustment was implemented.
6
7 # Identifying prod. companies with top revenue sums
8 # Identifying the 10 highest revenue-generating Production Companies
9 top_10_prod_companies <- mov_df_top_gens %>%
10   dplyr::group_by(PrimaryProductionCompany) %>%
11   dplyr::summarise(tot_revenue = sum(revenue)) %>%
12   arrange(-tot_revenue) %>%
13   select(PrimaryProductionCompany) %>%
14   head(10)
15
16 top_10_prod_companies$PrimaryProductionCompany
17
18 # Narrowing dataset to include only records associated with these top
19   -10-earning
20   # production companies.
21 mov_top_gens_and_prodCos <- mov_df_top_gens %>%
22   filter(PrimaryProductionCompany %in% top_10_prod_companies$
23     PrimaryProductionCompany)
```

```
1 # Binning revenue into quartiles - identifying break points
2 genre_prodco_sankey_df <- mov_top_gens_and_prodCos %>%
3   select(PrimaryGenre, revenue, PrimaryProductionCompany)
4
5 breaks_rev_qrts <- quantile(genre_prodco_sankey_df$revenue, c(0.25 *
6   (0:4)))
7 breaks_rev_qrts
```

```
1 # Binning revenue into quartiles - creating "revenueQuartile" attribute
2 # Binning accordingly (somewhat roughly, so that my break values will be
3 # prettier to the eye)
4 genre_prodco_sankey_df <- genre_prodco_sankey_df %>%
5   mutate(revenueQuartile = cut(revenue,
6     breaks = c(0, 41000000, 101000000,
7       217000000, 1850000000),
8     labels = c("$0 - $41,000,000",
```



```
9         "$41,000,000 - $101,000,000",
10         "$101,000,000 - $217,000,000",
11         "$217,000,000 - $1,850,000,000
12         "),
13         include.lowest=TRUE)) %>%
14 select(PrimaryGenre, revenueQuartile, PrimaryProductionCompany)
```

```
1 # Creating Sankey of Genre, Rev Quartile, and Prod Co.
2 # genre_prodco_sankey <- hchart(data_to_sankey(genre_prodco_sankey_df),
3   "sankey", name = "Relating revenue to genre and production company")
4 #
5 # genre_prodco_sankey <- genre_prodco_sankey %>%
6 #   hc_title(text = "Revenue outcomes based on Genre and Production
7   Company") %>%
8 #   hc_subtitle(text = "Examining the distributions of genre and
9   production
10   company for each quartile of revenue earned.")
11 #
12 # genre_prodco_sankey
13
14 sankeySeasons_df <- mov_top_gens_and_prodCos %>%
15   mutate(revenueQuartile = cut(revenue,
16     breaks = c(0, 41000000, 101000000,
17       217000000, 1850000000),
18     labels = c("$0 - $41,000,000",
19       "$41,000,000 - $101,000,000",
20       "$101,000,000 - $217,000,000",
21       "$217,000,000 - $1,850,000,000"
22     ),
23     include.lowest=TRUE)) %>%
24   select(PrimaryGenre, release_season, revenueQuartile)
```

```
1 seasonSankey <- hchart(data_to_sankey(sankeySeasons_df), "sankey", name
2   = "Relating revenue to genre and season of release")
3
4 seasonSankey <- seasonSankey %>%
5   hc_title(text = "Revenue outcomes based on Genre and Release Season")
6   %>%
7   hc_subtitle(text = "Examining the distributions of genre and release
8   season for each quartile of revenue earned.")
9
10 seasonSankey
11
12 # htmlwidgets::saveWidget(widget = seasonSankey, file = 'C:\\Users\\
13   maxru\\OneDrive\\Documents\\DePaul MS Work\\DSC 465\\Homework\\HW_4\\
14   genre-season_sankey.html')
15 # webshot::install_phantomjs()
```



```
11 # webshot::webshot(url = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul
    MS Work\\DSC 465\\Homework\\HW_4\\genre-season_sankey.html",
12 #                   file = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul
    MS Work\\DSC 465\\Homework\\HW_4\\genre-season_sankey.png" , delay =
    5)
13 # # knitr::include_graphics("genre-season_sankey.png")
```

```
1 season_genre_sankey_df <- mov_top_gens_and_prodCos %>%
2   mutate(revenueQuartile = cut(revenue,
3                               breaks = c(0, 41000000, 101000000,
4                                           217000000, 1850000000),
5                               labels = c("$0 - $41,000,000",
6                                           "$41,000,000 - $101,000,000",
7                                           "$101,000,000 - $217,000,000",
8                                           "$217,000,000 - $1,850,000,000
9                                           "),
10          include.lowest=TRUE)) %>%
11 select(release_season, PrimaryGenre, revenueQuartile)
```

```
1 season_genre_Sankey <- hchart(data_to_sankey(season_genre_sankey_df), "
    sankey", name = "Relating revenue to genre and season of release")
2
3 season_genre_Sankey <- season_genre_Sankey %>%
4   hc_title(text = "Season of Release, Genre, and Revenue Quartile") %>%
5   hc_subtitle(text = "Examining how release season, genre, and revenue
    quartile relate in terms of distribution.")
6
7 season_genre_Sankey
8
9 # htmlwidgets::saveWidget(widget = season_genre_Sankey, file = 'C:\\
    Users\\maxru\\OneDrive\\Documents\\DePaul MS Work\\DSC 465\\Homework
    \\HW_4\\season-genre-revQ_sankey.html ')
10 #
11 #
12 # # webshot::install_phantomjs()
13 #
14 # webshot::webshot(url = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul
    MS Work\\DSC 465\\Homework\\HW_4\\season-genre-revQ_sankey.html",
15 #                   file = "C:\\Users\\maxru\\OneDrive\\Documents\\DePaul
    MS Work\\DSC 465\\Homework\\HW_4\\season-genre-revQ_sankey.png" ,
    delay = 5)
16 #
17 # knitr::include_graphics("genre-season_sankey.png")
```