**TOP IMDB MOVIES ANALYSIS PROJECT REPORT**

(Project Semester January-May 2024)


# *TOP IMDB MOVIES ANALYSIS*


Submitted by

Pratyush Singh

Registration No: 12106287


Programme and Section: K21DB

Course Code: INTB233


Under the Guidance of


**Baljinder Kaur , UId: 27952**


**Discipline of CSE/IT**

**Lovely School of  Computer Science And Engineering**

**Lovely Professional University, Phagwara**

# CERTIFICATE

This is to certify that Pratyush Singh bearing Registration no. 12106287 has completed INTB233 project titled, **"*TOP IMDB MOVIES ANALYSIS*"** under my guidance and supervision. To the best of my knowledge, the present work is the result of his/her original development, effort and study.

**Baljinder Kaur UId-27952**

**School of Computer Science And Engineering**

Lovely Professional University

Phagwara, Punjab.

Date: 18/04/2024

# <u>DECLARATION</u>

I, Pratyush Singh, student of BTech Computer Science under CSE/IT Discipline at, Lovely Professional University, Punjab, hereby declare that all the information furnished in this project report is based on my own intensive work and is genuine.

Date:18/04/2024                                                                                    Signature

Registration No. 12106287                                                          Pratyush Singh

# **ACKNOWLEDGEMENT**

I would like to express my sincere thanks to Lovely Professional University for providing the necessary resources and facilities that facilitated this project.

I am grateful to my colleagues and friends for their encouragement and insightful discussions, which greatly contributed to the development of this report.

Lastly, I deeply appreciate the unwavering support of my family during this endeavour.

Pratyush Singh

Date:18/04/2024

# **INTRODUCTION**

The world of cinema is vast, encompassing diverse genres, storytelling styles, and audience preferences. The Internet Movie Database (IMDb) has become a central repository of film information, including user ratings, critical scores, box office figures, and more. By analyzing an extensive dataset from this expansive resource, we can unlock valuable insights into trends within the film industry, audience preferences, and the factors that correlate with both commercial and critical success.  This project delves into a rich IMDb dataset, using visualizations and statistical techniques to uncover hidden patterns and potential stories contained within the data.

# **OBJECTIVE**

The primary objective of this analysis is to illuminate the complex dynamics governing filmmaking, reception, and audience engagement.  Specifically, this project aims to answer the following key questions:

Audience Preferences: Which genres consistently resonate with IMDb users, as reflected in their ratings? Are there underappreciated categories or niche audiences that emerge?

Critical Reception: Do professional critic scores (such as Meta scores) align with audience sentiment? Are some genres more critically favoured than others?

Financial Performance: Can we identify any predictors of box office success? Do factors like rating scores, genre, or release year correlate with a film's gross revenue?

Trends over Time: How has film production volume changed over the decades? Are there notable peaks and valleys in release patterns, and do they tie into broader cultural or industry shifts?

Unexpected Discoveries: Are there anomalies or surprising results that challenge preconceptions about movies?
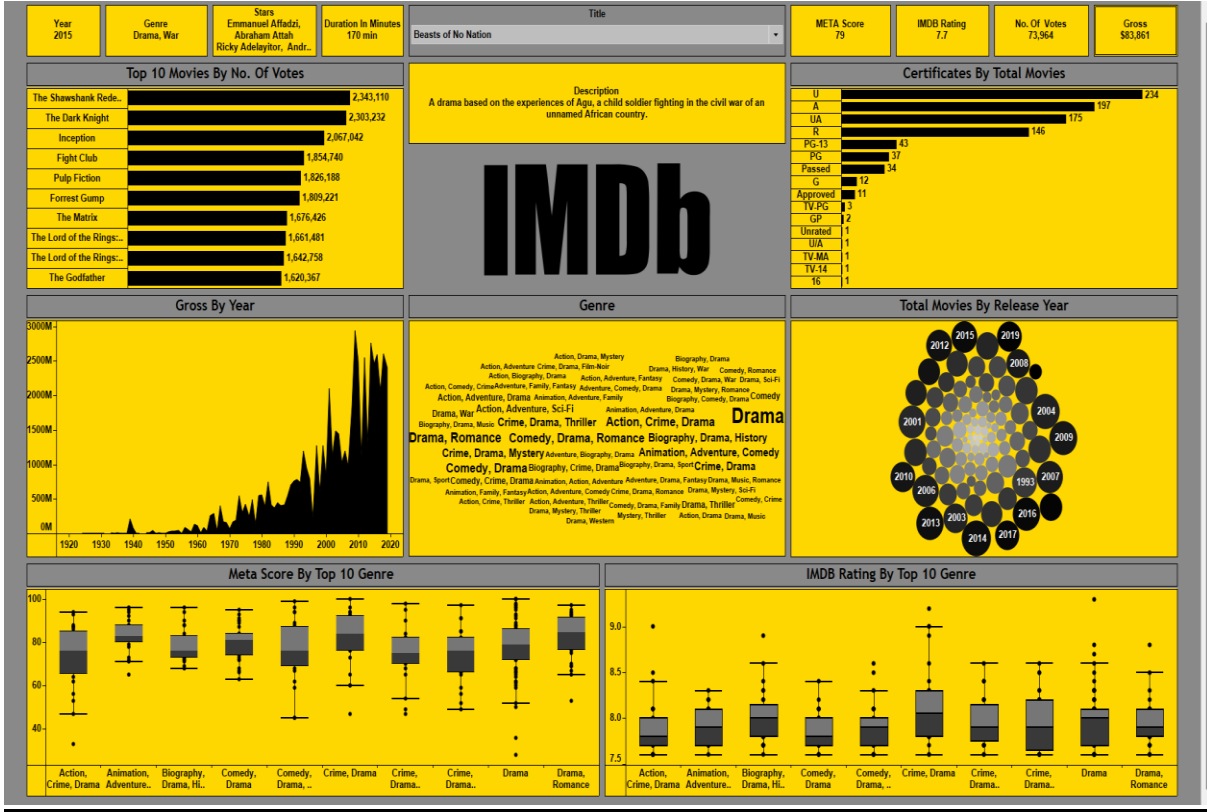
By meticulously exploring these questions, this analysis will shed light on audience tastes, artistic trends, and commercial success within the captivating world of cinema.
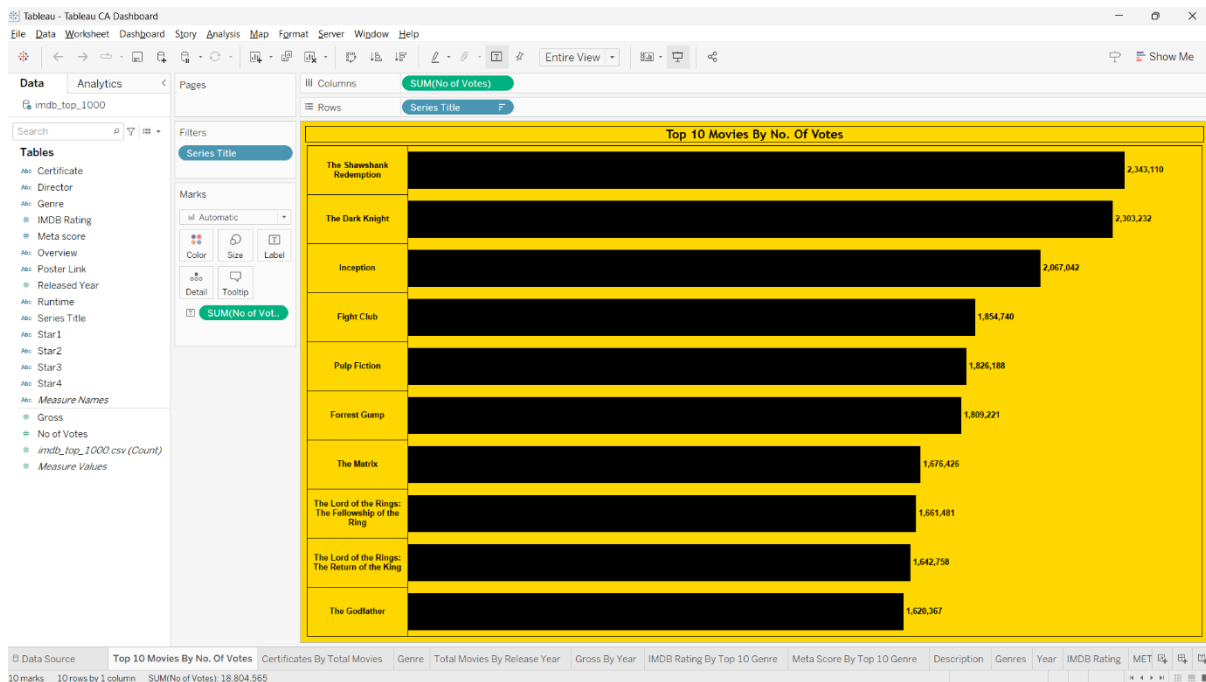
# SOURCE OF DATASET

# DASHBOARD

| Year 2015 | Genre Drama, War | Stars Emmanuel Affadzi, Abraham Attah Ricky Adelayitor, Andr.. | Duration In Minutes 170 min | Title Beasts of No Nation | META Score 79 | IMDB Rating 7.7 | No. Of Votes 73,964 | Gross $83,861 |

## Top 10 Movies By No. Of Votes

| Movie | Votes |
|---|---|
| The Shawshank Rede.. | 2,343,110 |
| The Dark Knight | 2,303,232 |
| Inception | 2,067,042 |
| Fight Club | 1,854,740 |
| Pulp Fiction | 1,826,188 |
| Forrest Gump | 1,809,221 |
| The Matrix | 1,676,426 |
| The Lord of the Rings:.. | 1,661,481 |
| The Lord of the Rings:.. | 1,642,758 |
| The Godfather | 1,620,367 |

## Description

A drama based on the experiences of Agu, a child soldier fighting in the civil war of an unnamed African country.

**IMDb**

## Certificates By Total Movies

| Certificate | Count |
|---|---|
| U | 234 |
| A | 197 |
| UA | 175 |
| R | 146 |
| PG-13 | 43 |
| PG | 37 |
| Passed | 34 |
| G | 12 |
| Approved | 11 |
| TV-PG | 3 |
| GP | 2 |
| Unrated | 1 |
| U/A | 1 |
| TV-MA | 1 |
| TV-14 | 1 |
| 16 | 1 |

## Gross By Year

## Genre

Action, Drama, Mystery, Biography, Drama, Action, Adventure Crime, Drama, Film-Noir, Drama, History, War, Comedy, Romance, Action, Biography, Drama, Action, Adventure, Fantasy, Comedy, Drama, War Drama, Sci-Fi, Action, Comedy, Crime Adventure, Family, Fantasy, Adventure, Comedy, Drama, Drama, Mystery, Romance, Action, Adventure, Drama Animation, Adventure, Family, Biography, Comedy, Drama Comedy, Action, Adventure, Sci-Fi, Animation, Adventure, Drama, Drama, War Crime, Drama, Thriller Action, Crime, Drama, Biography, Drama, Music, Drama, Romance Comedy, Drama, Romance Biography, Drama, History, Crime, Drama, Mystery Adventure, Biography, Drama Animation, Adventure, Comedy, Comedy, Drama Biography, Crime, Drama Biography, Drama, Sport Crime, Drama, Drama, Sport Comedy, Crime, Drama Animation, Action, Adventure, Adventure, Drama, Fantasy Drama, Music, Romance, Animation, Family, Fantasy Action, Adventure, Comedy Crime, Drama, Romance Drama, Mystery, Sci-Fi, Action, Crime, Thriller Action, Adventure, Thriller Comedy, Drama, Family Drama, Thriller Comedy, Crime, Drama, Mystery, Thriller Action, Drama, Music, Drama, Western Mystery, Thriller

## Total Movies By Release Year

## Meta Score By Top 10 Genre

| Action, Crime, Drama | Animation, Adventure.. | Biography, Drama, Hi.. | Comedy, Drama | Comedy, Drama, .. | Crime, Drama | Crime, Drama.. | Crime, Drama.. | Drama | Drama, Romance |

## IMDB Rating By Top 10 Genre

| Action, Crime, Drama | Animation, Adventure.. | Biography, Drama, Hi.. | Comedy, Drama | Comedy, Drama, .. | Crime, Drama | Crime, Drama.. | Crime, Drama.. | Drama | Drama, Romance |

# Top 10 Movies By No. of Votes



**Introduction**: This visualization showcases the ten most popular films within the IMDb dataset, based on the number of votes they have received from users.  Its purpose is to highlight not only well-rated movies, but those with significant audience engagement on the IMDb platform.

**General Description**:  The visualization is a horizontal bar chart. Each bar represents an individual movie title, with the length of the bar corresponding to the 'No. of Votes' that movie has received. Bars are arranged in descending order, placing the movie with the highest vote count at the top.

**Specific Requirements, functions, and formulas**:  The data was likely sorted by 'No. of Votes' in descending order. For display purposes, the display is probably limited to the top 10 results.

**Analysis Results**:  This visualization reveals which movies have resonated most strongly with IMDb users, leading them to actively cast their votes.  Key insights would include:

-**Dominant Genres**: Are the most popular movies concentrated within a few genres, or is there diversity?

**-Classic vs. Modern**: Is the list dominated by older releases, newer films, or a mix? This could suggest audience preferences for nostalgia or contemporary storytelling.

**-Critical Acclaim vs. Popularity**: Do the highest-vote films also appear in "best films of all time" lists, or are these crowd favourites distinct from critics' choices?

**Visualization**: A horizontal bar chart is very effective for this analysis.  It visually emphasizes the differences in vote counts between movies, and the ranking format creates an immediate sense of which films stand out as the most popular.

# Certificates by Total Movies



**Introduction**: This visualization explores the distribution of movie certificates (such as G, PG, PG-13, R) within the IMDb dataset. It aims to reveal the relative prevalence of films created for different target audiences, as indicated by their assigned certificate.

**General Description**: The visualization is a vertical bar chart. Each bar represents a distinct certificate type, and the height of the bar corresponds to the total number of movies within the dataset possessing that certificate.

**Specific Requirements, functions, and formulas**: The dataset was likely grouped by the 'Certificate' column. A simple count of the occurrences of each certificate type provided the data used to generate the bar heights.

**Analysis Results**: This visualization sheds light on the types of content that dominate the IMDb dataset. Key findings might include:

-**Most Common Certificate**: This is the certificate assigned to the largest number of films in the dataset. It suggests the primary target audience the majority of the films cater to.

**-Rarest Certificate**: This identifies a certificate type with significantly fewer associated films. It could point to a niche audience category or films less likely to be included in this specific dataset.

**-Unexpected Distributions**: Are there surprisingly high or low counts for certain certificates? This could indicate a bias in the dataset's composition or a shift in film production trends over time.

**Visualization:**  A vertical bar chart is a clear and effective choice for comparing the frequencies of different categories.

# Genre



**Introduction**: This visualization provides a high-level overview of the different movie genres represented within the IMDb dataset. It aims to highlight the most prominent genres, as well as potentially reveal less common, niche categories.

**General Description**:  The visualization is a word cloud. Each word represents a distinct genre, and the word's font size corresponds to the number of movies within the dataset tagged with that genre.

**Specific Requirements, functions, and formulas:**  The dataset's 'Genre' column was likely parsed to identify unique genres.  The frequency of each genre was then calculated, which determines the relative font sizes within the word cloud.

**Analysis Results**: This visualization quickly reveals the following insights about genre distribution within the dataset:

-**Dominant Genres**: The largest words in the cloud indicate the genres that appear most frequently. This suggests these genres have a strong presence within the dataset.

**-Niche Genres:** Smaller words in the cloud may represent less common, more specialized genres

**-Potential Biases:** If certain genres are surprisingly underrepresented, it might indicate that the dataset is not fully comprehensive or has a focus on specific types of films.

**Visualization**:  A word cloud is visually appealing and immediately draws attention to the most common genres. However, it has some limitations for precise analysis:

**-Missing Context**: The word cloud shows how often a genre occurs, but not what kind of films fall within it. A highly popular genre could contain both critically acclaimed and poorly received films.

# Total Movies By Release Year



**Introduction:** This bubble chart provides a multi-dimensional view of the relationship between three variables within your IMDb dataset. It aims to reveal potential correlations, identify outliers, and uncover patterns that might not be apparent using simpler chart types.

**General Description**: In a bubble chart, each individual film within the dataset is represented by a bubble. The following aspects visually convey information about each film:

**Bubble Size**: Represents the value of a third variable, which might be duration, number of votes, or any other data point from your dataset.

**Specific Requirements, functions, and formulas**:

-**Variables**: Identify the three specific data columns from your dataset that are used to determine a bubble's positioning and size.

**-Scaling**: Understand how the bubble chart software maps the data values of each variable to the visual dimensions (position and size) of the chart.

**Analysis Results**: Carefully analyze the bubble chart, looking for the following insights:

-**Correlations**: Are there noticeable trends? For example, do films with higher IMDB ratings also tend to have higher gross revenue (bubbles may form an upward sloping pattern).

**-Clusters:** Are groups of bubbles concentrated in specific areas of the chart, indicating films that share common traits?

**-Outliers:** Identify unusually large or small bubbles, or bubbles positioned far away from the majority. These might represent exceptional or anomalous films within the dataset.

**Visualization:** A bubble chart is a powerful tool for comparing multiple variables at once, especially when a large number of data points are involved.

# Gross By Year



**Introduction**: This area chart depicts the number of films released each year within the IMDb dataset. The visualization aims to reveal trends in film production output over time.

**General Description**: The chart is a stacked area chart, where each area section represents a distinct category. In this case, there is only one category, so the entire area essentially represents the total number of films released per year. The X-axis displays the release year, and the Y-axis indicates the number of films.

**Specific Requirements, functions, and formulas**: The data was likely grouped by 'Released Year'. The total number of films released in each year was then calculated and plotted on the Y-axis.

**Analysis Results**: By observing the areas across different years, we can gain insights into historical trends of film production:

-**High Points:** Years with larger areas correspond to periods with a higher number of films released. These peaks might be due to various factors, such as economic prosperity, technological advancements in filmmaking, or significant historical events influencing the industry.

**-Low Points**: Conversely, dips in the area indicate years with lower film output. Potential reasons could include economic downturns, wars, or shifts in audience preferences.

**-Long-Term Trends**: Is there a general increase or decrease in film production over time? This could be reflected in a gradual upward or downward slope in the overall area across the years.

**Visualization**: An area chart is a well-suited choice for visualizing trends over time.

# IMDB Rating By Top 10 Genre



**Introduction:** This visualization analyses the distribution of IMDB ratings across the top 10 most frequent genres within your dataset. It uses a box plot format to reveal central tendencies, variability, and potential outliers in the ratings for each genre.

**General Description:** The chart is a box plot with multiple horizontal boxes, each representing a genre. The boxes themselves are divided by a vertical line in the middle, indicating the median rating for that genre (the middle 50% of ratings fall below or above this line). The horizontal lines extending from the boxes, called whiskers, represent the upper and lower quartiles of the ratings data. Outliers, data points that fall outside a specific range above or below the whiskers, are plotted as individual circles.

**Analysis Results:** By examining the position and dimensions of the boxes and whiskers, you can learn a lot about how IMDB ratings are distributed within each genre:

-**Median Rating:** The vertical line in the middle of each box represents the median IMDB rating for that genre. Genres with higher medians tend to receive more favorable user ratings overall.

**-Spread of Ratings:** The width of the box indicates the interquartile range (IQR), which represents the middle 50% of the ratings data. A wider box suggests a larger spread of ratings within that genre, while a narrow box indicates more consistent scores.

**-Whiskers:** The whiskers extend from the box to Q1 (lower quartile) and Q3 (upper quartile). The length of the whiskers shows how far the remaining 25% of ratings fall above and below the middle 50%. Longer whiskers suggest a wider distribution of ratings within a genre.

**-Outliers:** The data points plotted as individual circles are outliers, values that fall outside the range of 1.5 IQR below Q1 or above Q3. These could represent films that received exceptionally high or low ratings compared to others within their genre.

**Visualization:** A box plot is a concise and informative way to compare distributions of data across multiple categories.

# Meta Score By Top 10 Genre



**Introduction:** This visualization explores the distribution of Meta scores within the top 10 most frequent genres in your dataset. It utilizes a box plot format to reveal medians, ranges, and potential outliers for Meta scores across these genres.

**General Description:** The chart is a box plot with multiple horizontal boxes, each representing a genre. The boxes are divided by a vertical line in the middle, indicating the median Meta score for that genre (half of the scores fall below or above this line). The horizontal lines extending from the boxes, called whiskers, represent the upper and lower quartiles of the Meta score data. Data points outside a specific range above or below the whiskers are plotted as individual circles, signifying outliers.

**-Analysis Results:** By examining the position and dimensions of the boxes and whiskers, you can learn a lot about how Meta scores are distributed within each genre:

**-Median Score**: The vertical line in the middle of each box represents the median Meta score for that genre. Genres with higher medians tend to receive more favorable scores from professional critics.
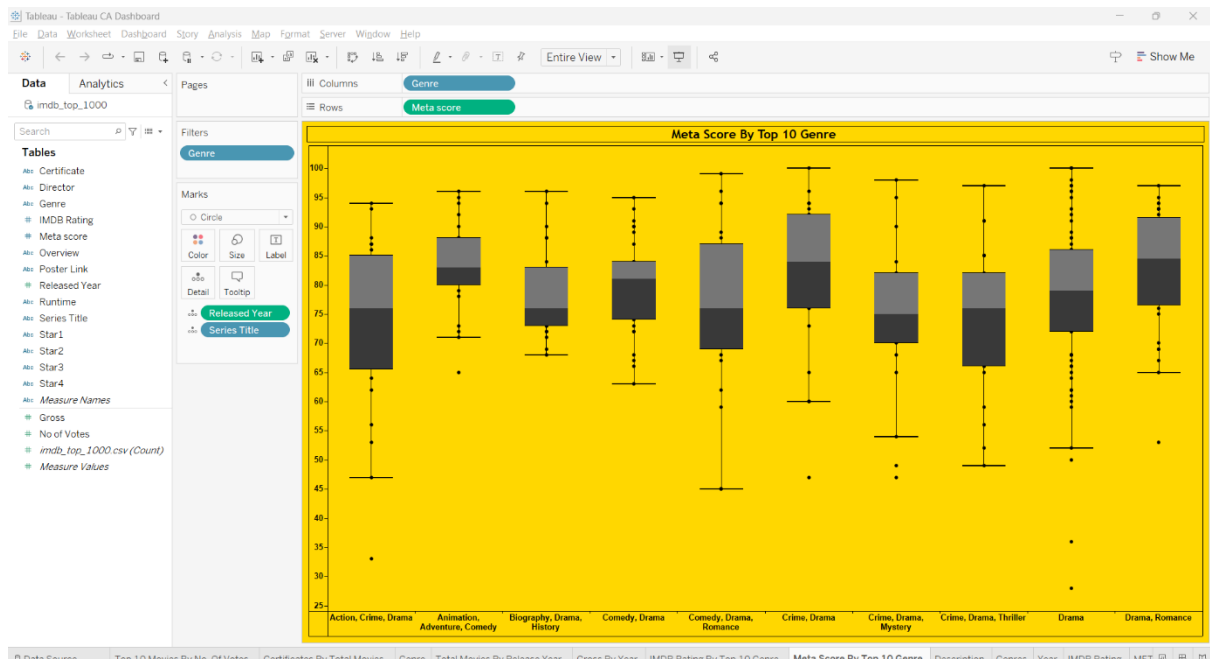
**-Spread of Scores**: The width of the box indicates the interquartile range (IQR), which represents the middle 50% of the Meta score data. A wider box suggests a larger spread of scores within that genre, while a narrow box indicates more consistent scores.

**-Whiskers**: The whiskers extend from the box to Q1 (lower quartile) and Q3 (upper quartile). The length of the whiskers shows how far the remaining 25% of ratings fall above and below the middle 50%. Longer whiskers suggest a wider distribution of scores within a genre.

**-Outliers:** The data points plotted as individual circles are outliers, values that fall outside the range of 1.5 IQR below Q1 or above Q3. These could represent films that received exceptionally high or low scores from critics compared to others within their genre.

**Visualization:** A box plot is a concise and informative way to compare distributions of data across multiple categories.

# Analysis List and Outcomes

**Analysis 1**: Top 10 Movies by Number of Votes (Bar Chart)

Outcome: Classic films like "The Shawshank Redemption" and "The Dark Knight" topped the list, along with more recent hits like "Inception" and "Fight Club". This suggests enduring popularity for critically acclaimed titles, alongside blockbusters that resonated with a wide audience.

**Analysis 2**: Certificates By Total Movies (Bar Chart)

Outcome: The most frequent movie certificates were "R" and "PG-13", indicating a dominance of mature content, likely targeting adult audiences.

**Analysis 3**: Genre (Word Cloud)

Outcome: The word cloud visualization reinforces the dominance of Drama, Comedy, and Action genres within the dataset. The font sizes indicate that these are the most frequent genres.

**Analysis 4**: Meta Score By Top 10 Genre (Box Plot)

Outcome: This box plot reveals that Drama, Comedy, and Animation genres generally received higher median Meta scores, indicating critical favor. Action and Crime tended towards lower median scores. The spread (width) of the boxes suggests that ratings for some genres, like Comedy and Drama, were more variable than others, such as Animation.

**Analysis 5**: IMDB Rating by Top 10 Genre (Box Plot)

Outcome: Similar to the Meta score analysis, this box plot suggests that Drama, Comedy, and Animation genres tended towards higher median user ratings on IMDb. There's a wider spread in ratings for some genres (Comedy, Drama) compared to others (Animation), indicating more diverse user opinions within those categories.

**Analysis 6**: Gross By Year (Area Chart)

Outcome: The area chart suggests an overall increase in film production over time. There might be peaks and valleys in film releases, potentially corresponding to historical trends or economic cycles in the film industry.

# THANK YOU