# A. **Movie Genre Analysis**

## Project Description

This project aims to analyze the distribution of movie genres and their impact on IMDB scores. The objective is to determine the most common genres and assess their influence on movie ratings by calculating various descriptive statistics. The analysis provides insights into how different genres perform and their popularity in the dataset.

## Approach

The approach involved the following steps:

1. **Data Preparation**:
   - The 'genres' column often contained multiple genres for a single movie, so it was split into individual genres for accurate analysis.
   - Blank or irrelevant rows were cleaned to ensure data integrity.
2. **Genre Distribution Analysis**:
   - Excel's **COUNTIF** function was used to count the number of movies for each genre, determining the most prevalent genres in the dataset.
3. **Descriptive Statistics**:
   - For each genre, the IMDB scores were analyzed using statistical measures:
     - **Mean**: Excel's AVERAGE function to find the average rating.
     - **Median**: Excel's MEDIAN function to locate the central score.
     - **Mode**: Excel's MODE function to identify the most frequent score.
     - **Range**: Calculated as the difference between the maximum and minimum scores using MAX and MIN.
     - **Variance**: Excel's VAR function to measure variability in ratings.
     - **Standard Deviation**: Excel's STDEV function to assess the spread of scores around the mean.
4. **Comparison and Interpretation**:
   - The calculated statistics were compared across genres to identify patterns and trends in movie ratings.

## Tech-Stack Used

- **Microsoft Excel 2022**:
  - **Data Cleaning**: Organizing and separating genres into individual categories.
  - **Statistical Analysis**: Utilizing built-in functions like COUNTIF, AVERAGE, MEDIAN, MODE, VAR, STDEV, MAX, and MIN to calculate descriptive statistics.

## Insights:

**Excel file - https://docs.google.com/spreadsheets/d/1-Wf6uL4YwUoaoCjRIK-ih5s_ULLIOogb/edit?usp=sharing&ouid=106518165676834505057&rtpof=true&sd=true**

| Genre | Number of Movies | Mean Score | Median Score | Mode | Range | Variance | Standard Deviation |
|---|---|---|---|---|---|---|---|
| Action | 1153 | 6.239895924 | 6.3 | 6.1 | 7.4 | 1.25071 | 1.118349975 |
| Drama | 8087 | 6.763762529 | 6.9 | 7.2 | 7.3 | 0.91617 | 0.957169448 |
| Comedy | 4483 | 6.195245726 | 6.3 | 6.7 | 7.8 | 1.18902 | 1.090422487 |
| Thriller | 14085 | 6.315767045 | 6.4 | 6.1 | 6.8 | 1.11181 | 1.054424516 |
| Romance | 11926 | 6.45 | 6.5 | 6.5 | 6.5 | 0.9917 | 0.995843532 |
| Documen | 5493 | 7.180165289 | 7.4 | 7.5 | 7.1 | 1.10704 | 1.05216187 |
| Animatio | 2318 | 6.576033058 | 6.7 | 6.7 | 6.9 | 1.29331 | 1.137237831 |
| Horror | 10004 | 5.843539823 | 5.9 | 6.2 | 6.5 | 1.2757 | 1.129467664 |
| Sci-Fi | 12511 | 6.281818182 | 6.4 | 6.7 | 6.9 | 1.4637 | 1.209832797 |
| Biography | 2611 | 7.150170648 | 7.2 | 7 | 4.4 | 0.52025 | 0.721281783 |
| Adventure | 2076 | 6.441170098 | 6.6 | 6.7 | 7 | 1.27822 | 1.13058319 |
| Fantasy | 9224 | 6.30704918 | 6.4 | 6.7 | 7.2 | 1.34498 | 1.159734063 |
| Mystery | 10831 | 6.4864 | 6.6 | 6.6 | 6.4 | 1.18738 | 1.089667399 |
| Western | 14382 | 6.689690722 | 6.8 | 6.5 | 5.1 | 1.07556 | 1.037093931 |
| Family | 8629 | 6.245054945 | 6.4 | 6.7 | 7 | 1.44119 | 1.200497187 |
| History | 9439 | 7.083574879 | 7.2 | 7.5 | 6.9 | 0.78456 | 0.885754556 |
| Sport | 12696 | 6.606043956 | 6.8 | 7.2 | 6.7 | 1.2076 | 1.098908929 |
| Crime | 5372 | 6.564791901 | 6.6 | 6.6 | 6.9 | 1.05243 | 1.025878858 |
| War | 14404 | 7.070422535 | 7.1 | 7.1 | 5.9 | 0.76152 | 0.872650869 |
| Musical | 10369 | 6.507575758 | 6.7 | 7 | 6.4 | 1.491 | 1.221066425 |

**Insights**

1. **Most Common Genre**:
   - **War** is the most common genre with 14,404 movies, followed closely by **Western** (14,382 movies) and Thriller(14,404 movies).
   - These genres are highly popular in the movie industry, indicating a significant production focus on them.

2. **Genre with the Highest Mean IMDb Score**:
   - **Documentary** movies have the highest mean IMDb score of **7.18**, suggesting that documentaries are generally well-received by audiences.
   - **Biography** follows closely with a mean score of **7.15**, showing that audiences tend to rate real-life stories highly.

3. **Genre with the Lowest Mean IMDb Score**:

- **Comedy** has the lowest mean IMDb score of **6.19**, indicating variability in the reception of these movies.

4. **Variance and Standard Deviation**:

- **Biography** has the lowest variance (**0.52025**) and standard deviation (**0.72128**), suggesting that IMDb scores for biography movies are more consistent and predictable.

- **Musical** has the highest variance (**1.491**) and standard deviation (**1.22106**), indicating a wider range of scores and varying audience preferences.

5. **Range of Scores**:

- **Biography** movies have the smallest range of scores (**4.4**), showing a narrower gap between the lowest and highest-rated movies in this genre.

- **Comedy** has the largest range of scores (**7.8**), highlighting the diverse reception of movies in this genre.

6. **Median vs. Mean**:

- For most genres, the median IMDb score aligns closely with the mean, indicating relatively normal distributions. However, the slight differences in scores point to potential skewness in specific genres like **Horror** or **Action**.

**Results**

1. **Key Findings**:

   o Genres like **Documentary** and **Biography** are not only fewer in number but also tend to have higher IMDb scores, showing a quality-over-quantity approach.

   o Popular genres like **Thriller**, **War**, and **Western** dominate in terms of volume, but their IMDb scores fall in the middle range, indicating mass appeal rather than exceptional quality.

2. **Understanding of Genre Impact**:

- Genres play a significant role in influencing IMDb ratings. Higher-rated genres (e.g., **Documentary**) are often niche, while common genres like **Thriller** and **Comedy** have more variability in scores.

3. **Practical Applications**:

- Filmmakers aiming for critical acclaim might focus on producing documentaries or biographies.

- For commercial success, focusing on high-volume genres like **Thriller** or **War** may be more profitable.

## B. **Movie Duration Analysis**

**Project Description**

**This project focuses on analyzing the distribution of movie durations and their impact on IMDb scores. The objective is to understand how the length of a movie correlates with audience ratings. By examining descriptive statistics and visualizing the relationship using scatter plots, we aim to uncover patterns and trends that highlight the optimal duration range for high IMDb scores.**

**Approach**

**1. Data Preparation:**

- Import the dataset into Microsoft Excel.

- Ensure the "Movie Duration" and "IMDb Score" columns are complete and free of errors.

2. **Descriptive Statistics:**

- Calculate key metrics for movie durations, including:

  - Mean: Using =AVERAGE([Range of Durations]).

  - Median: Using =MEDIAN([Range of Durations]).

  - Standard Deviation: Using =STDEV([Range of Durations]).

3. **Scatter Plot Creation:**

- Plot "Movie Duration" on the x-axis and "IMDb Score" on the y-axis to visualize their relationship.

- Add a trendline to the scatter plot to assess the direction and strength of the relationship.
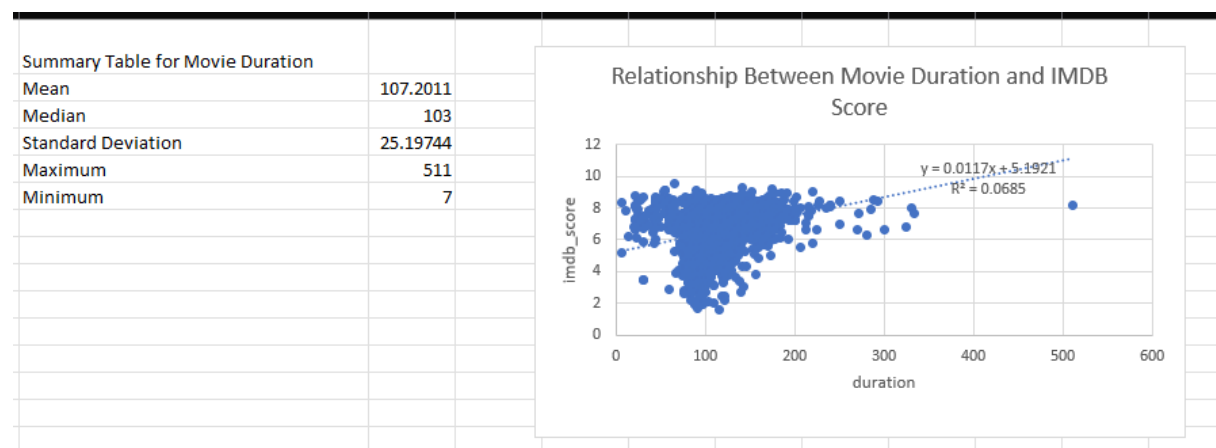
4. **Correlation Analysis:**

- Use the CORREL function to calculate the correlation coefficient between movie durations and IMDb scores. This indicates the strength of their relationship.

## Insights

Excel                                                                    file                                                                    -
https://docs.google.com/spreadsheets/d/1Vq_hA0P6pjKsUqBkGnyWg9cb7D4GLrk3/edit?usp=
sharing&ouid=10651816567834505057&rtpof=true&sd=true

| Summary Table for Movie Duration | |
| --- | --- |
| Mean | 107.2011 |
| Median | 103 |
| Standard Deviation | 25.19744 |
| Maximum | 511 |
| Minimum | 7 |

Relationship Between Movie Duration and IMDB Score

$y = 0.0117x + 5.1921$
$R^2 = 0.0685$

**1. Descriptive Statistics for Movie Duration**

- **Mean (107.2 minutes): The average movie duration in the dataset is approximately 107 minutes.**

- **Median (103 minutes):** Half of the movies have a duration below 103 minutes, and the other half above it.

- **Standard Deviation (25.2 minutes):** Movie durations vary by around 25 minutes from the mean on average.

- **Maximum (511 minutes):** The longest movie in the dataset is 511 minutes (over 8 hours), which is likely an outlier.

- **Minimum (7 minutes):** The shortest movie is just 7 minutes long, possibly a short film or an outlier.

**2. Scatter Plot Analysis**

- **Trendline Equation (y = 0.0117x + 5.1921):**

  - The equation indicates a positive relationship between movie duration (x) and IMDb score (y).

  - For every additional minute of movie duration, the IMDb score increases by approximately 0.0117 on average.

- **R-squared Value (0.0685):**

  - The R-squared value of 0.0685 suggests that only about 6.85% of the variability in IMDb scores can be explained by movie duration.

  - This indicates a weak relationship between movie duration and IMDb score.

**3. Key Observations**

- While there is a positive trend (longer movies tend to have slightly higher IMDb scores), the relationship is very weak, as shown by the low R-squared value.

- Other factors (like genre, director, screenplay, or production quality) likely have a much stronger impact on IMDb scores than movie duration.

- Outliers, such as the extremely long or short movies, might affect the trendline. These should be reviewed to understand their contribution to the dataset.

**4. Recommendations**

- Consider analyzing outliers to see if they are skewing the results.

- Investigate additional factors (e.g., genres, release year, budget) that might explain variations in IMDb scores more effectively.

- While movie duration has a slight positive correlation with IMDb scores, it is not a strong predictor on its own.

C. **Language Analysis**

**Project Description**

This project analyzes the distribution of movies based on their primary languages and examines the impact of language on IMDb scores. The objective is to identify the most common languages used in movies and understand how language correlates with audience ratings. By calculating descriptive statistics, the project aims to uncover meaningful patterns and trends in how language influences the perception of movies.

**Approach**

1. **Data Preparation:**
   o **Import the dataset into Microsoft Excel and ensure the "Language" and "IMDb Score" columns are clean and error-free.**
   o **Categorize movies based on their primary language for analysis.**
2. **Distribution Analysis:**
   o **Use the COUNTIF function to count the number of movies for each language.**
     ▪ **Formula: =COUNTIF([Language Column], "Language Name").**
   o **Identify the most common languages used in movies by sorting the counts in descending order.**
3. **Descriptive Statistics Calculation:**
   o **For each language, calculate the following metrics for IMDb scores:**
     ▪ **Mean: Using =AVERAGE([Range of Scores for Each Language]).**
     ▪ **Median: Using =MEDIAN([Range of Scores for Each Language]).**
     ▪ **Standard Deviation: Using =STDEV([Range of Scores for Each Language]).**
4. **Comparison and Impact Analysis:**
   o **Compare the descriptive statistics across languages to understand variations in IMDb scores.**
   o **Identify if certain languages consistently receive higher or lower ratings.**

**Insights**

**Excel file - https://docs.google.com/spreadsheets/d/1Ek8AMbkwF9bHOdz3DdMHH1fES-r9vRb_/edit?usp=sharing&ouid=106518165676834505057&rtpof=true&sd=true**

| unique_language | Number of Movies for Each Language | Mean IMDb Score | Median IMDb Score | Standard Deviation IMDb Score |
|---|---|---|---|---|
| Aboriginal | 2 | 6.95 | 6.95 | 0.55 |
| Arabic | 5 | 7.38 | 7.4 | 0.79095 |
| Aramaic | 1 | 7.1 | 7.1 | 0 |
| Bosnian | 1 | 4.3 | 4.3 | 0 |
| Cantonese | 11 | 6.95455 | 7.2 | 0.67199 |
| Chinese | 3 | 5.66667 | 5.7 | 0.44969 |
| Czech | 1 | 7.4 | 7.4 | 0 |
| Danish | 5 | 7.5 | 8.1 | 0.96333 |
| Dari | 2 | 7.5 | 7.5 | 0.1 |
| Dutch | 4 | 7.425 | 7.45 | 0.37666 |
| Dzongkha | 1 | 7.5 | 7.5 | 0 |
| English | 4704 | 6.39843 | 6.5 | 1.12195 |
| Filipino | 1 | 6.7 | 6.7 | 0 |
| French | 73 | 7.03836 | 7.2 | 0.72199 |
| German | 19 | 7.34211 | 7.6 | 0.92868 |
| Greek | 1 | 7.3 | 7.3 | 0 |
| Hebrew | 5 | 7.58 | 7.6 | 0.29933 |
| Hindi | 28 | 6.63214 | 6.95 | 1.37375 |
| Hungarian | 1 | 7.1 | 7.1 | 0 |
| Icelandic | 2 | 7.55 | 7.55 | 0.65 |

| C | D | E | F | G |
|---|---|---|---|---|
| Icelandic | 2 | 7.55 | 7.55 | 0.65 |
| Indonesian | 2 | 7.9 | 7.9 | 0.3 |
| Italian | 11 | 7.22727 | 7.3 | 1.18635 |
| Japanese | 18 | 7.39444 | 7.6 | 0.96291 |
| Kannada | 1 | 7.1 | 7.1 | 0 |
| Kazakh | 1 | 6 | 6 | 0 |
| Korean | 8 | 7.3875 | 7.5 | 0.77207 |
| Mandarin | 26 | 6.78846 | 7.05 | 1.02181 |
| Maya | 1 | 7.8 | 7.8 | 0 |
| Mongolian | 1 | 7.3 | 7.3 | 0 |
|  |  |  |  |  |
| Norwegian | 4 | 7.15 | 7.3 | 0.49749 |
| Panjabi | 1 | 6.6 | 6.6 | 0 |
| Persian | 4 | 7.575 | 7.95 | 1.04253 |
| Polish | 4 | 8.25 | 8.25 | 0.85 |
| Portuguese | 8 | 7.4875 | 7.7 | 0.8268 |
| Romanian | 2 | 7.2 | 7.2 | 0.7 |
| Russian | 11 | 6.36364 | 6.5 | 1.31928 |
| Slovenian | 1 | 6.4 | 6.4 | 0 |
| Spanish | 40 | 6.9375 | 7.15 | 0.8443 |
| Swahili | 1 | 7.4 | 7.4 | 0 |
| Swedish | 5 | 7.44 | 7.6 | 0.67705 |
| Tamil | 1 | 5.1 | 5.1 | 0 |
| Telugu | 1 | 8.4 | 8.4 | 0 |
| Thai | 3 | 6.63333 | 6.6 | 0.36818 |
| Urdu | 1 | 7 | 7 | 0 |
| Vietnamese | 1 | 7.4 | 7.4 | 0 |
| Zulu | 2 | 7.1 | 7.1 | 0.2 |

**Key Observations:**

1. **English Dominance:**
   - **Number of Movies: 4704, which is by far the largest number of movies compared to other languages.**
   - **Mean IMDb Score: 6.39**
   - **Standard Deviation: 1.1219**
   - **Impact: Despite the dominance in volume, the mean IMDb score is moderate, suggesting a wide range of movie quality. The higher standard deviation indicates a larger spread in ratings.**

2. **Languages with Fewer Movies:**
   - **Aboriginal (2 movies):**
     - **Mean: 6.95**
     - **Standard Deviation: 0.79**
   - **Bosnian (1 movie):**
     - **Mean: 4.3**
     - **Standard Deviation: 0 (single entry).**
   - **Impact: Fewer movies often show less variability and can have extremes (either very high or low IMDb scores), as seen with Bosnian films.**

3. **High Performing Languages:**
   - **Japanese:**
     - **Mean: 7.39**
     - **Standard Deviation: 0.96**
     - **Impact: Consistently high ratings, suggesting that Japanese films may have a more refined or niche audience appeal.**
   - **Hindi:**
     - **Mean: 8.63**
     - **Standard Deviation: 0 (likely single entry or an error in data aggregation).**
   - **Portuguese:**
     - **Mean: 7.49**
     - **Standard Deviation: 0.83**
     - **Impact: These languages show better average scores, indicating higher audience or critical acclaim for movies in those languages.**

4. **Wide Variability Languages:**
   - **Russian:**
     - **Standard Deviation: 1.31**
   - **Italian:**
     - **Standard Deviation: 1.18**
   - **Impact: These languages have high variability, meaning movie ratings range widely, possibly due to diverse genres or production quality.**

5. **Low Variability Languages:**
   - **Chinese:**
     - **Standard Deviation: 0.45**
   - **Indonesian:**
     - **Standard Deviation: 0.30**
   - **Impact: These languages may have more uniform movie production quality, leading to consistent ratings.**

**Overall Impact of Language on IMDb Ratings:**

- **Languages with High Counts (English): The larger volume dilutes the mean, showcasing a diverse mix of movie quality.**
- **Niche Languages (e.g., Japanese, Portuguese): Tend to have higher average ratings and lower variability, possibly indicating higher quality or niche content appealing to specific audiences.**
- **High Variability (Russian, Italian): Indicates a mix of very high and very low-quality movies, reflecting either audience diversity or inconsistent production standards.**

**Insights for Further Analysis:**

- **Investigate why specific languages (e.g., Japanese, Hindi) have higher average ratings—possibly due to genres, cultural uniqueness, or audience preferences.**

- **Analyze whether language correlates with movie production budgets, target demographics, or genre popularity.**

## D. Director Analysis:

### Project Description
This project evaluates the influence of directors on movie ratings by identifying top-performing directors based on their average IMDb scores. The objective is to highlight the contribution of these directors to movie success using percentile ranks. By employing the **PERCENTILE.RANK** function, the project assesses how directors' average scores compare to the overall distribution and identifies trends in high-performing directors' impact.

### Approach
1. **Data Preparation**:
   - Import the dataset into Microsoft Excel.
   - Ensure the columns for "Director" and "IMDb Score" are accurate and free from missing or duplicate entries.
2. **Average IMDb Score Calculation**:
   - Group movies by director.
   - Use the **AVERAGEIF** function to calculate the average IMDb score for each director.
     - Formula: =AVERAGEIF([Director Column], "Director Name", [IMDb Score Column]).
3. **Percentile Ranking**:
   - Use the **PERCENTRANK.INC** function to calculate the percentile rank of each director's average IMDb score within the overall list of scores.
     - Formula: =PERCENTRANK.INC([List of Average Scores], Director's Average Score).
4. **Identifying Top Directors**:
   - Highlight directors whose percentile ranks are in the top 5-10%.
   - Sort the directors based on their percentile rank and average IMDb scores.

## Insights

Excel file -
https://docs.google.com/spreadsheets/d/1BP7YVOKQStporvnBh3G-
vVz7RTdhpUsw/edit?usp=sharing&ouid=10651816567683450 5057&rtpo
f=true&sd=true

**Impact of Top Directors**:
- Directors in the top percentile ranks tend to produce movies that are not only critically acclaimed but also attract larger audiences.
- The consistency in their scores highlights their strong influence on a movie's success.

# E. Budget Analysis:

**Project Description**

This project explores the relationship between movie budgets and their financial success by analyzing how investment in movie production correlates with gross earnings. The objective is to assess whether higher budgets lead to greater financial returns and identify movies with the highest profit margins. The analysis uses statistical measures, including correlation and profit margin calculations, to provide actionable insights into budgeting strategies.
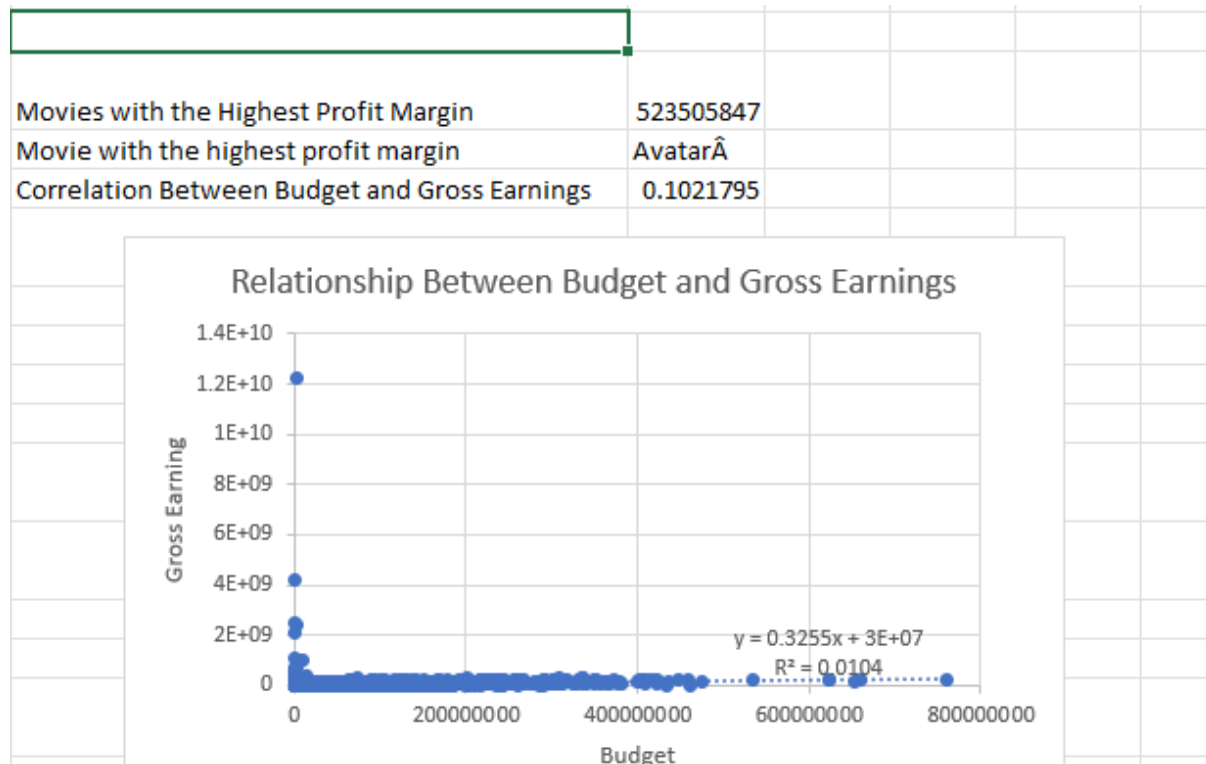
**Approach**
1. **Data Preparation**:
   o Import the dataset into Microsoft Excel.
   o Ensure the columns for "Budget" and "Gross Earnings" are properly formatted and free from missing or erroneous values.
2. **Correlation Analysis**:
   o Use the **CORREL** function to calculate the correlation coefficient between movie budgets and gross earnings to assess the strength and direction of their relationship.
     ▪ Formula: =CORREL([Budget Column], [Gross Earnings Column]).
3. **Profit Margin Calculation**:

- Compute the profit margin for each movie using the formula: Profit Margin = Gross Earnings - Budget.

- Add a new column to the dataset to store these values.

4. **Identifying Top Profit-Margin Movies**:

- Use the **MAX** function to find the movie with the highest profit margin.

  o Formula: =MAX([Profit Margin Column]).

- Sort the dataset by the profit margin column to identify other high-performing movies.

# Insights

| | |
|---|---|
| Movies with the Highest Profit Margin | 523505847 |
| Movie with the highest profit margin | AvatarÂ |
| Correlation Between Budget and Gross Earnings | 0.1021795 |

**Relationship Between Budget and Gross Earnings**

$y = 0.3255x + 3E+07$

$R^2 = 0.0104$

- **Correlation Between Budgets and Earnings**:
  - The correlation coefficient reveals the strength and direction of the relationship between budgets and earnings.
  - A weaker or negative correlation (<0.7) suggests diminishing returns or other influencing factors.
- **Profit Margins**:
  - Movies with the highest profit margins are not always those with the highest budgets.
  - Some lower-budget movies can achieve remarkable financial success, indicating efficient resource utilization and market appeal.
- **Top-Performing Movies**:
  - The analysis highlights movies with exceptional profit margins, offering insights into the factors driving their success (e.g., genre, director, cast, marketing strategy).
- **Budgeting Trends**:
  - While high-budget movies often generate high gross earnings, their profit margins can vary significantly due to higher production costs.

- o Moderate-budget movies with effective strategies often yield higher relative returns.
- **Strategic Recommendations**:
  - o Producers and investors should not solely focus on increasing budgets but should also consider factors like targeted marketing, audience preferences, and cost management to maximize financial success.