# Data Collection and Preprocessing Phase
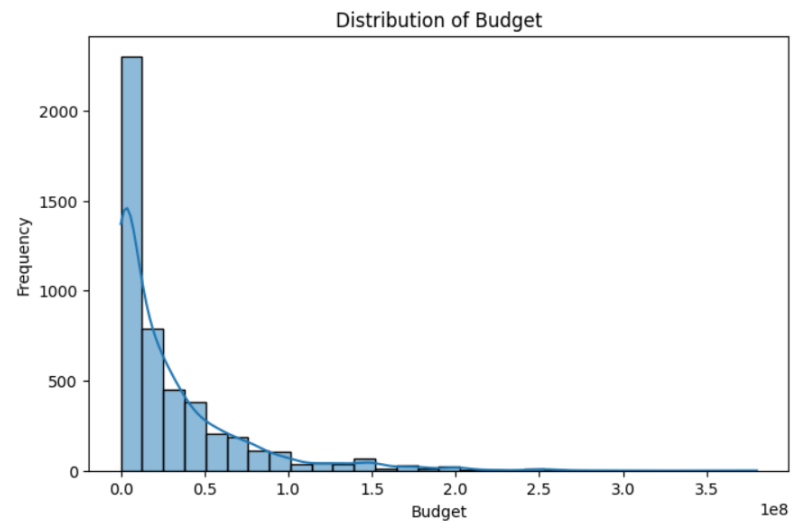
| | |
|---|---|
| Date | 23 September 2024 |
| Team ID | LTVIP2024TMID24986 |
| Project Title | Movie Box Office Gross Prediction using Machine Learning |
| Maximum Marks | 6 Marks |

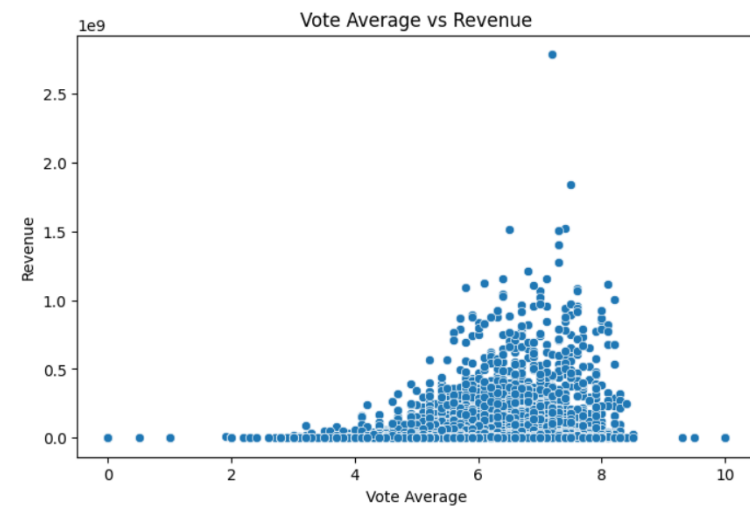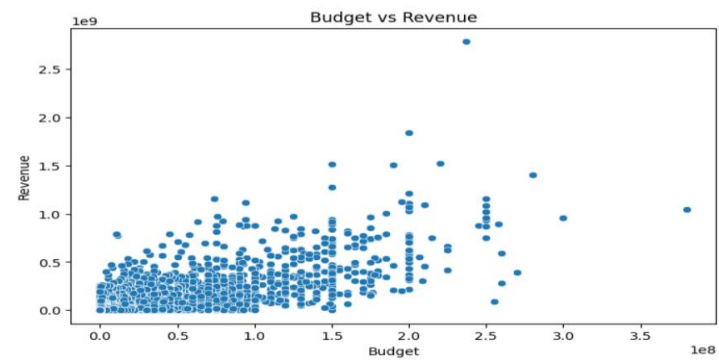**Data Exploration and Preprocessing Report**

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

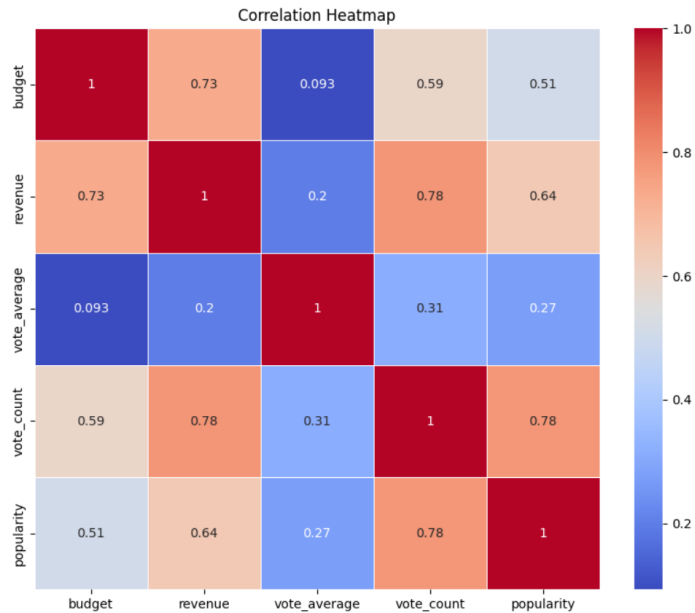| Section | Description |
|---|---|
| Data Overview | Dimension:<br>4083rows × 23columns<br>Descriptive statistics:<br><br>|  | budget | id | popularity | revenue | runtime | vote_average | vote_count |<br>|---|---|---|---|---|---|---|---|<br>| count | 4.803000e+03 | 4803.000000 | 4803.000000 | 4.803000e+03 | 4801.000000 | 4803.000000 | 4803.000000 |<br>| mean | 2.904504e+07 | 57165.484281 | 21.492301 | 8.226064e+07 | 106.875859 | 6.092172 | 690.217989 |<br>| std | 4.072239e+07 | 88694.614033 | 31.816650 | 1.628571e+08 | 22.611935 | 1.194612 | 1234.585891 |<br>| min | 0.000000e+00 | 5.000000 | 0.000000 | 0.000000e+00 | 0.000000 | 0.000000 | 0.000000 |<br>| 25% | 7.900000e+05 | 9014.500000 | 4.668070 | 0.000000e+00 | 94.000000 | 5.600000 | 54.000000 |<br>| 50% | 1.500000e+07 | 14629.000000 | 12.921594 | 1.917000e+07 | 103.000000 | 6.200000 | 235.000000 |<br>| 75% | 4.000000e+07 | 58610.500000 | 28.313505 | 9.291719e+07 | 118.000000 | 6.800000 | 737.000000 |<br>| max | 3.800000e+08 | 459488.000000 | 875.581305 | 2.787965e+09 | 338.000000 | 10.000000 | 13752.000000 | |

| Univariate Analysis |  Distribution of Budget |
| --- | --- |
| Bivariate Analysis |  Budget vs Revenue  Vote Average vs Revenue |

| | |
|---|---|
| Multivariate Analysis |  |
| Outliers and Anomalies | - |

## Data Preprocessing Code Screenshots

| | |
|---|---|
| Loading Data | ```python
credits=pd.read_csv("/content/tmdb_5000_credits.csv")
movies_df=pd.read_csv("/content/tmdb_5000_movies.csv")

credits.head()
```
 |
| Handling Missing Data | ```python
from sklearn.preprocessing import LabelEncoder
from collections import Counter as c
cat=['director','genres']
for i in movies_box[cat]:
  print("LABEL ENCODING OF:",i)
  LE = LabelEncoder()
  print(c(movies_box[i]))
  movies_box[i] = LE.fit_transform(movies_box[i])
  print(c(movies_box[i]))
``` |

| Data Transformation | ```python
movies['log_revenue'] = np.log1p(movies['revenue'])
movies['log_budget'] = np.log1p(movies['budget'])

movies_box = movies.drop(['homepage','id','keywords','original_language','original_title',
                          'overview','production_countries','release_date','spoken_languages',
                          'status','tagline','title_x','title_y','cast',
                          'log_revenue','log_budget'],axis = 1)

movies_box=movies_box.drop(['production_companies'],axis=1)
``` + Code  + Text ```python
movies_box.isnull().sum()
``` |
|---|---|
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | - |