

```
1 from google.colab import drive
2 drive.mount('/content/drive')

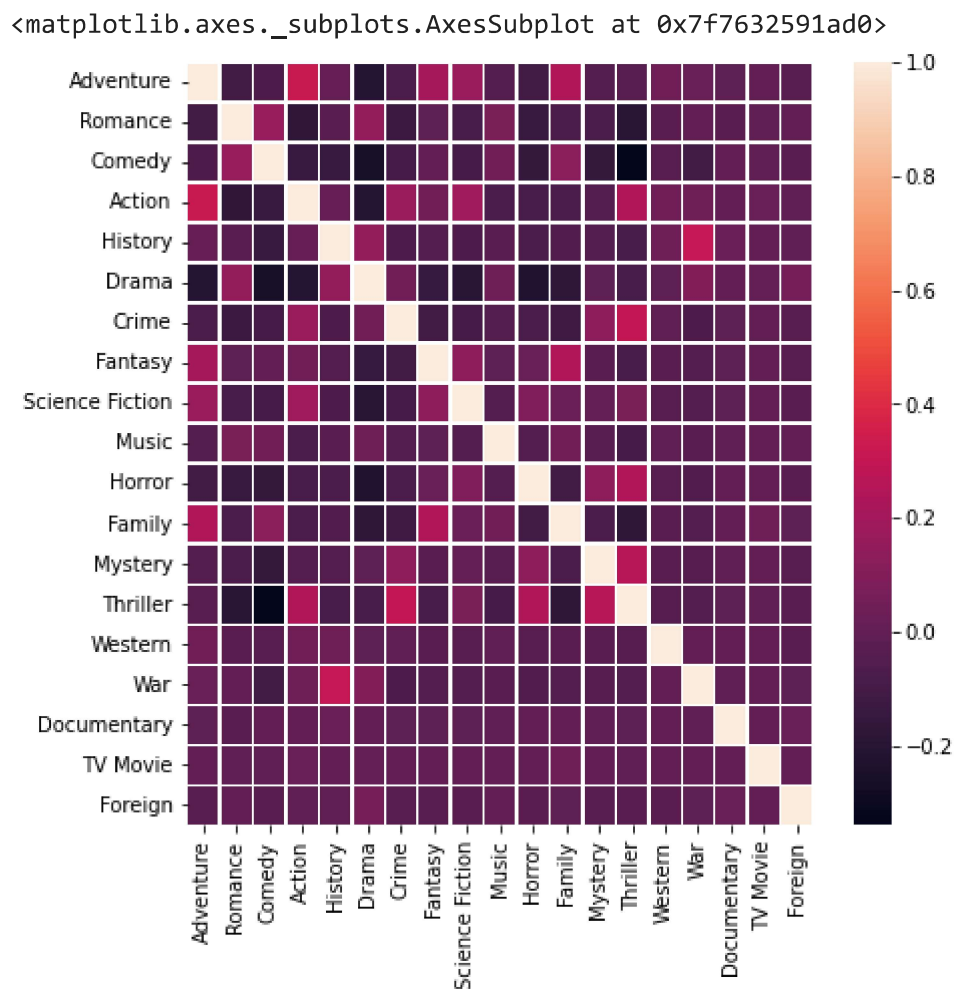
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
```

```
1 # Pandas is used to convert our csv to workable dataframe.
2 import pandas as pd
3
4 # Seabord is an advanced visualising tool, used here , for Heat maps.
5 import seaborn as sns
6
7 # Matplotlib is used to plot data into graphs.
8 import matplotlib.pyplot as plt
9
10 # Using Sklearn to model and retrieve evaluation metrics.
11 from sklearn.model_selection import train_test_split
12 from sklearn.linear_model import LogisticRegression
13 from sklearn.metrics import accuracy_score, confusion_matrix
14 from sklearn.metrics import precision_score, f1_score, roc_curve
15 from sklearn.metrics import auc, recall_score, roc_auc_score
16 from sklearn.linear_model import LogisticRegression
17 from sklearn.neighbors import KNeighborsClassifier
18 from sklearn import tree
19 from sklearn.ensemble import RandomForestClassifier
```

```
1 # Importing the cleaned Dataset into a pandas variable called df
2 df = pd.read_csv("/content/drive/MyDrive/IDMP PROJECT/Final Dataset/Final Dataset.csv")
```

```
1 # Accessing only the genres related columns from the dataframe
2 genres = df[df.columns[126:146]]
```

```
1 # Using the matplotlib library to retireve a subplot of size 7x7
2 fig, ax = plt.subplots(figsize=(7,7))
3 # Using the Seaborn library to plot the heatmap on our subplot
4 sns.heatmap(genres.corr(), linewidths=.5, ax=ax)
```



```
1 # Converting the Bool column "adult" into 1s and 0s
2 df["adult"].replace(True,'1', inplace=True)
3 df['adult'].replace(False,'0', inplace=True)
```

```
1 # Retriving the correlation values of only the target variable
2 # corr() gives the correlation values of all variables against each other.
3 # Since df is a pandas dataframe, we use column name(hit/not)
4 # to retrieve the correlation vlaues relative to that column
5
6 cor = df.corr()
7 cor_tar = cor["hit/not"]
8
9 #Using a threshold value of 0.04 we choose the columns only that
10 # have correlation greater than threshold
11
```

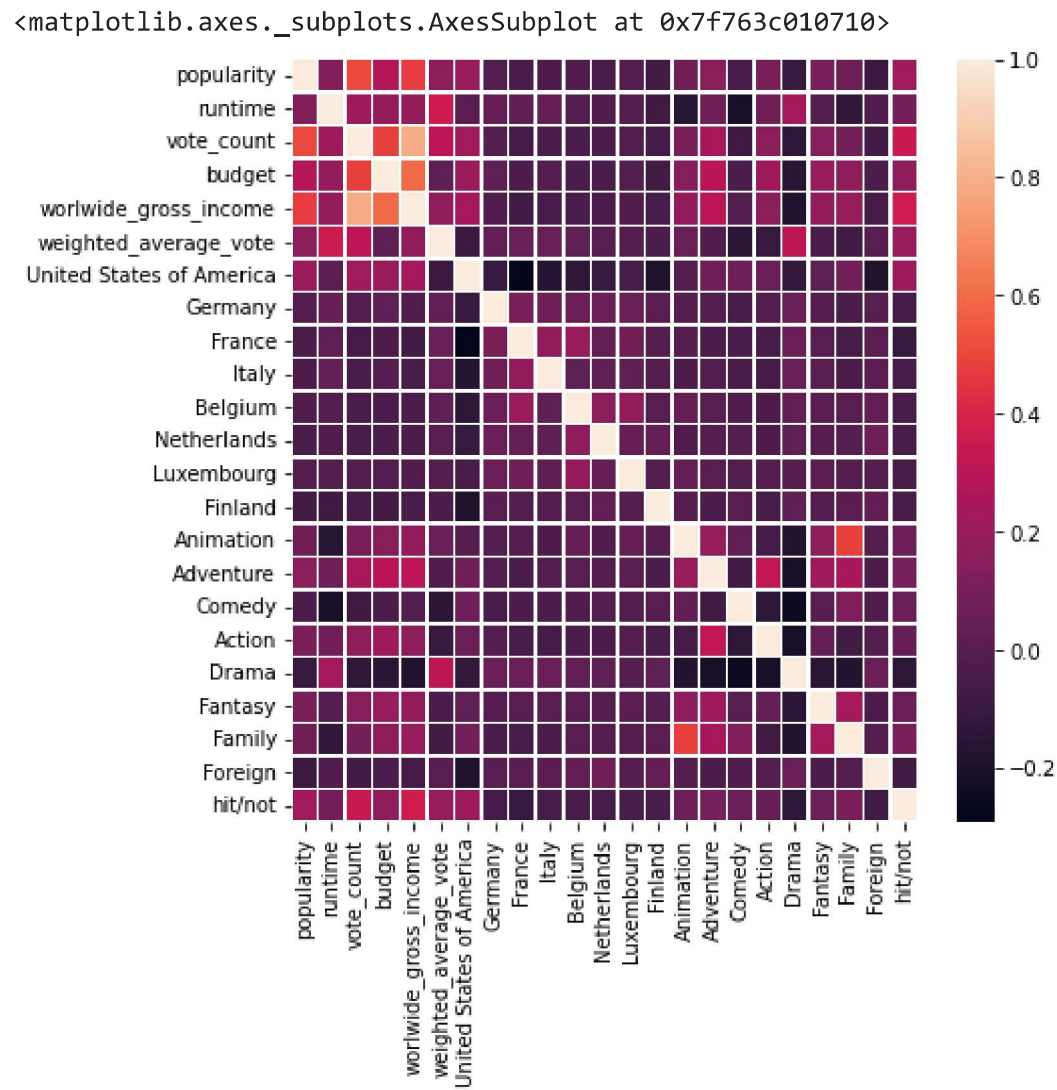
```
12 threshold = 0.04
13 rel_cor = cor_tar[abs(cor_tar) > threshold ]

1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8992 entries, 0 to 8991
Columns: 147 entries, adult to hit/not
dtypes: float64(4), int64(135), object(8)
memory usage: 10.1+ MB

1 # rel_cols only contain those columns that have greater correlation with our
2 # target variable.
3
4 rel_cols = ["popularity","runtime","vote_count","budget","worlwide_gross_income",
5             "weighted_average_vote","United States of America","Germany","France",
6             "Italy","Belgium","Netherlands","Luxembourg","Finland","Animation",
7             "Adventure","Comedy","Action","Drama","Fantasy","Family","Foreign",
8             "hit/not"]
9 # Filtering out the columns that are irrelevant to our model
10 non_rel = []
11 for col in df.columns:
12     if col not in rel_cols:
13         non_rel.append(col)
14
15 # using drop() to drop those columns from our dataframe
16 df = df.drop(non_rel, axis=1)
```

```
1 # creating a Heat map for relevant Columns
2 fig, ax = plt.subplots(figsize=(7,7))
3 sns.heatmap(df.corr(), linewidths=.5, ax=ax)
```



```
1 ## Splitting data into Test and Training set
2
3 X = df.iloc[:, :-1:]
4 y = df.iloc[:, -1]
5
6 x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
7                                                     random_state=0)

1 # Evaluation Metrics
2 # Using Sklearns to calculate metrics
3 # We are calculating accuracy, precision, F1-score, AUC and Specificity
4 def eval_metrics(y_test, y_pred):
5     accuracy = accuracy_score(y_test, y_pred)
6     cm = confusion_matrix(y_test, y_pred)
7     precision = precision_score(y_test, y_pred)
8     recall = recall_score(y_test, y_pred)
9     f1 = f1_score(y_test, y_pred)
10    fpr, tpr, thresholds = roc_curve(y_test, y_pred, pos_label=1)
```

```

11 A = auc(fpr, tpr)
12 roc = roc_auc_score(y_test, y_pred)
13 tn, fp, fn, tp = cm.ravel()
14 specificity = tn / (tn+fp)
15 return accuracy*100, precision, recall, f1, A, specificity
16

```

## Logistic Regression

```

1 # Using Logistic Regression model from Sklearn to model our data.
2 # Since it is a basic model, the we used mse to calculate error.
3
4 # loading the LogisticRegression module in to lr
5 lr = LogisticRegression()
6
7 # Fitting the model to our train set
8 lr.fit(x_train, y_train)
9
10 # Using the trained model to predict out X_Test
11 pred = lr.predict(x_test)
12
13 # Calculating our evaluation metrics using sklearn build-in package.
14 acc_lr,precision_lr, rec_lr, f1_lr, areaUnderCurve_lr, spec_lr= eval_metrics(y_test, pred)
15
16 # Printing out the Evaluation Metrics
17 print(f"The metrics for Logistic regression are: \n Accuracy: {acc_lr}% ",
18       f"\nPrecision: {precision_lr} \nRecall: {rec_lr} \nF1 score: {f1_lr} ",
19       f"\nArea Under Curve: {areaUnderCurve_lr} \nSpecificity: {spec_lr}")
20

```

```

The metrics for Logistic regression are:
Accuracy: 100.0%
Precision: 1.0
Recall: 1.0
F1 score: 1.0
Area Under Curve: 1.0
Specificity: 1.0

```

## KNN:

```

1 # Using KNeighborsClassifier model from Sklearn to model our data.
2 # Since it is a basic model, the we used mse to calculate error.
3
4 # loading the KNeighborsClassifier module in to knn
5 knn = KNeighborsClassifier()
6
7 # Fitting the model to our train set
8 knn.fit(x_train, y_train)
9
10 # Using the trained model to predict out X_Test
11 y_pred_knn = knn.predict(x_test)
12
13 # Calculating our evaluation metrics using sklearn build-in package.
14 acc_knn, precision_knn, rec_knn, f1_knn, areaUnderCurve_knn, spec_knn = eval_metrics(y_test, y_pred_knn)
15
16 # Printing out the Evaluation Metrics
17 print(f"The metrics for KNN are: \n Accuracy: {acc_knn}% \nPrecision: {precision_knn}"+
18       f" \nRecall: {rec_knn} \nF1 score: {f1_knn} \nArea Under Curve: {areaUnderCurve_knn}"+
19       f" \nSpecificity: {spec_knn}")
20

```

```

The metrics for KNN are:
Accuracy: 99.27737632017788%
Precision: 0.9913606911447084
Recall: 0.9945828819068255
F1 score: 0.9929691725256895
Area Under Curve: 0.9927252309077507
Specificity: 0.9908675799086758

```

## Decision Tree:

```

1 # Using KNeighborsClassifier model from Sklearn to model our data.
2 # For Decision Tree we take Cross-Entropy Loss as loss function.
3 # We set that minimum split samples to be 5
4 # minimum number of leaf nodes is 6
5 # and we set the max_features to auto to facilitate the model to random skip a few features.
6 # Random State is set as 50.
7
8 # loading the KNeighborsClassifier module in to dt with hyperparameters
9 dt = tree.DecisionTreeClassifier(criterion='entropy', min_samples_split=5,
10                                min samples leaf=6, max features='auto',

```

```

11         random_state=50)
12
13 # Fitting the model to our train set
14 dt.fit(x_train, y_train)
15
16 # Using the trained model to predict out X_Test
17 y_pred_dt = dt.predict(x_test)
18
19 # Calculating our evaluation metrics using sklearn build-in package.
20 acc_dt, precision_dt, rec_dt, f1_dt, areaUnderCurve_dt, spec_dt = eval_metrics(y_test, y_pred_dt)
21
22 # Printing out the Evaluation Metrics
23 print(f"The metrics for Decision Tree are: \n Accuracy: {acc_dt}%",
24       f" \nPrecision: {precision_dt} \nRecall: {rec_dt} \nF1 score:",
25       f" {f1_dt} \nArea Under Curve: {areaUnderCurve_dt}",
26       f" \nSpecificity: {spec_dt}")
27
28
29 The metrics for Decision Tree are:
30 Accuracy: 85.38076709282934%
31 Precision: 0.8395061728395061
32 Recall: 0.8840736728060672
33 F1 score: 0.8612137203166227
34 Area Under Curve: 0.8529957405126226
35 Specificity: 0.821917808219178

```

## Random Forest

```

1 # Using RandomForestClassifier model from Sklearn to model our data.
2 # For Random Forest, we take Cross-Entropy Loss as loss function.
3 # We set that minimum split samples to be 5
4 # minimum number of leaf nodes is 6
5 # and we set the max_features to auto to facilitate the model to random skip a few features.
6 # Random State is set as 50.
7
8 # loading the KNeighborsClassifier module in to rf with hyperparameters
9 rf = RandomForestClassifier(n_estimators=50,
10                             criterion='entropy', min_samples_split=5,
11                             min_samples_leaf=6, max_features='auto',
12                             random_state=50)
13
14 # Fitting the model to our train set
15 rf.fit(x_train, y_train)
16
17 # Using the trained model to predict out X_Test
18 y_pred_rf = rf.predict(x_test)
19
20 # Calculating our evaluation metrics using sklearn build-in package.
21 acc_rf, precision_rf, rec_rf, f1_rf, areaUnderCurve_rf, spec_rf = eval_metrics(y_test, y_pred_rf)
22
23 # Printing out the Evaluation Metrics
24 print(f"The metrics for Random Forest are: \n Accuracy: {acc_rf}%" +
25       f" \nPrecision: {precision_rf} \nRecall: {rec_rf} \nF1 score: {f1_rf} " +
26       f" \nArea Under Curve: {areaUnderCurve_rf} \nSpecificity: {spec_rf}")
27
28
29 The metrics for Random Forest are:
30 Accuracy: 95.27515286270149%
31 Precision: 0.9554347826086956
32 Recall: 0.952329360780065
33 F1 score: 0.9538795442213781
34 Area Under Curve: 0.9527628539060143
35 Specificity: 0.9527628539060143

```

