



ΧΑΡΟΚΟΠΕΙΟ ΠΑΝΕΠΙΣΤΗΜΙΟ  
HAROKOPIO UNIVERSITY

# DATA MINING

Harokopio University of Athens

Kokkalis Konstantinos – it22047

Professor: Varlamis Iraklis



## Contents

<b>A.DATA PRE-PROCESSING</b> .....	3
<b>STEPS</b> .....	3
Data frame Initialization .....	5
Remove redundant columns.....	6
Missing values .....	9
Data formatting.....	13
One-Hot encoding .....	14
Delete duplicates.....	17
Extra functions.....	18
Result .....	19
<b>B1.CLASSIFICATION</b> .....	20
<b>STEPS</b> .....	21
Load and Prepare Data.....	21
Model Training .....	23
<b>B2.PREDICTIONS</b> .....	25
<b>HOW TO RUN THE CODE</b> .....	25



## A.DATA PRE-PROCESSING

### STEPS

1. Data frame Initialization.
2. Remove redundant columns.
3. Handle missing values.
4. Data formatting.
5. One Hot Encoding.
6. Delete duplicates.

```
def executePreprocess(self,type=None):  
    df=initDataframe(self.fileToProcess)  
    df=dropUseless(df,USELESS_COL) # DELETE USELESS COLUMNS  
    df=columnDataFormating(df)  
    df=numericMissingValues(df) # RETRIEVE MISSING VALUES  
    df=stringMissingValues(df) # RETRIEVE MISSING VALUES  
    df=columnDataFormating(df)  
    df=oneHotEncoding(df) # ONE HOT ENCODING  
    df=deleteDuplicate(df) # CHECK FOR DUPLICATE ROWS  
    df=dropUseless(df,['film','year']) # DELETE USELESS COLUMNS
```



### Note:

At the beginning we save in some arrays the numeric columns, the redundant columns, the target strings, the no target strings as long as the script types. The reason we do that is to handle easier decision changes in the future, such as which columns we are not going to use in our classification model.

```
ALL_NUMERIC=['year', 'rotten tomatoes critics', 'metacritic critics', 'average critics', 'rotten tomatoes audience', 'metacritic audience', 'rotten tomatoes vs metacritic deviance', 'average audience', 'audience vs critics deviance', 'opening weekend', 'opening weekend ($million)', 'domestic gross', 'domestic gross ($million)', 'foreign gross ($million)', 'foreign gross', 'worldwide gross', 'worldwide gross ($million)', 'budget ($million)', 'of gross earned abroad', 'budget recovered', 'budget recovered opening weekend', 'imdb rating', 'distributor', 'imdb vs rt disparity']
NO_TARGET_STRINGS=['script type', 'primary genre', 'genre', 'release date (us)'] #except 'film' 'oscar winners', 'oscar detail'
TARGET_STRINGS=['oscar winner', 'oscar detail']
TYPES=['adaptation', 'original', 'based on a true story', 'sequel', 'remake']
USELESS_COL=['rotten tomatoes critics', 'metacritic critics', 'rotten tomatoes audience', 'metacritic audience', 'rotten tomatoes vs metacritic deviance', 'audience vs critics deviance', 'primary genre', 'opening weekend ($million)', 'domestic gross ($million)', 'foreign gross ($million)', 'worldwide gross ($million)', 'worldwide gross', 'budget recovered opening weekend', 'distributor', 'imdb vs rt disparity']
```

**ALL\_NUMERICS** = all the columns that they contain numeric values.

**TARGET\_STRINGS**=all the string columns that have Oscar details.

**NO\_TARGET\_STRINGS**=all the other string columns except from target strings.

**USELESS\_COLUMNS**=all the columns that we do not need for the classification model.



### Data frame Initialization

1. We say the to the data frame to replace '-' and '0' with Not a Number values so we can handle them later.
2. We remove redundant spaces from column names and we transform them to lower case so we can handle them easier. We also do the same for the whole dataset.
3. We update the dataset.
4. We update the ALL\_NUMERIC and the NO\_TRAGET\_STRINGS table according to the USELESS\_COLUMNS array.
5. Return the updated dataset.

```
df=initDataframe(self.fileToProcess)
```

```
def initDataframe(xFile):
    try:
        df=pd.read_excel(xFile, sheet_name = 'Sheet1',na_values=['-','0'])
        df.columns = df.columns.str.lower().str.replace(r'\s+', ' ', regex=True)
        df = df.map(lambda x: x.lower() if isinstance(x, str) else x)
        df.columns = df.columns.str.strip()
        df.to_excel(xFile, index=False)
        for item in USELESS_COL:
            if item in ALL_NUMERIC:
                ALL_NUMERIC.remove(item)
            if item in NO_TRAGET_STRINGS:
                NO_TRAGET_STRINGS.remove(item)
    except:
        raise RuntimeError(f'{colors.RED}A problem occurred while initializing the excel file!
        {colors.END}')
    print(f"{colors.GREEN}INITIALIZATION HAS BEEN SUCCESFULLY CONVERTED!{colors.END}")
    return df
```



### Remove redundant columns

```
df=dropUseless(df,USELESS_COL) # DELETE USELESS COLUMNS
```

```
def dropUseless(file,uselessColumns):  
    try:  
        for item in uselessColumns:  
            if item not in file.columns:  
                uselessColumns.remove(item)  
                continue  
        if len(uselessColumns)!=0:  
            file.drop(uselessColumns,axis=1,inplace=True) # deleting useless data from the excel  
    except:  
        raise RuntimeError(f'{colors.RED}A problem occured while dropping useless columns.{colors.END}')  
    print(f"{colors.GREEN}USELESS COLUMNS HAS BEEN SUCCESFULLY DELETED!{colors.END}")  
    return file
```

```
USELESS_COL=['rotten tomatoes critics','metacritic critics','rotten  
tomatoes audience','metacritic audience','rotten tomatoes vs metacritic  
deviance','audience vs critics deviance','primary genre','opening  
weekend ($million)','domestic gross ($million)','foreign gross  
($million)','worldwide gross ($million)','worldwide gross','budget  
recovered opening weekend','distributor','imdb vs rt disparity']
```

*Why each of them has been deleted?*

**distributor, imdb vs rt disparity, primary genre:** We dropped these columns because we had no, or a few data comparing to the size of the tuples, which means they would not give any information and they would not help the training model.

**opening weekend (\$million), domestic gross (\$million), foreign gross (\$million), worldwide gross (\$million):** We dropped these columns because we have already the information of each, in other columns with a bigger precision (opening weekend, domestic gross, foreign gross).

**rotten tomatoes vs metacritic deviance, audience vs critics deviance:** We dropped these columns because they do not provide any useful information.



**worldwide gross:** We dropped this column because worldwide gross is a result of domestic and foreign gross columns that we already kept, so there is no point to keep those two features and at the same time the sum of those two features.

**rotten tomatoes critics, metacritic critics, rotten tomatoes audience, metacritic audience:**

Here we checked the correlation of these attributes and the averages columns to decide if we are going to keep the averages or each one individually. The result was to keep the averages and drop the individuals. The decision was made according to these results:

	rotten tomatoes critics	metacritic critics	average critics
rotten tomatoes critics	1.000000	0.941663	0.991123
metacritic critics	0.941663	1.000000	0.977929
average critics	0.991123	0.977929	1.000000

We can notice in the previous picture that all the correlations are very close to 1 and they also very similar between them. That means that we can keep only one of them. In that case we will keep average critics because it is also the mean of the other two features.

	rotten tomatoes audience	metacritic audience	average audience
rotten tomatoes audience	1.000000	0.719496	0.935264
metacritic audience	0.719496	1.000000	0.897723
average audience	0.935264	0.897723	1.000000

We can notice here that all the correlations are very close to 1 and they also very similar between them. That means that we can keep only one of them. In that case we will keep average audience because it is also the mean of the other two features.

**budget recovered opening weekend:** Here we checked again the correlation of this attribute, budget recovered and budget (\$million). The decision was made according to these results:

	budget (\$million)	budget recovered	budget recovered opening weekend
budget (\$million)	1.000000	-0.001887	-0.002231
budget recovered	-0.001887	1.000000	0.998427
budget recovered opening weekend	-0.002231	0.998427	1.000000

We can notice here that all the correlations are very close to 0 (so there is no correlation) and they also very similar between them. We can only notice that budget recovered and budget recovered opening weekend they have correlation close to 1 so one of them could be deleted. So, we chose to keep budget recovered feature.



Code for getting correlation:

```
subset = df[['budget ($million)', 'budget recovered', 'budget recovered opening weekend']]
# subset = df[['rotten tomatoes critics', 'metacritic critics', 'average critics']]
# subset = df[['rotten tomatoes audience', 'metacritic audience', 'average audience']]
getCorrelation(subset)
def getCorrelation(items):
    correlation_matrix = items.corr()
    print(correlation_matrix)
```

### Note:

**film, year:** We delete film and year because the name of the film is a string that will not give helpful information to the classification model, on the other side it will confuse it. Year also does not provide any useful information, but we can use both to retrieve other data. So, we do not drop them yet.

**script type, genre, Oscar detail:** We will one-hot encode them and some of them so we will delete the main feature of one-hot encoding after this process.





### Missing values

```
df=numericMissingValues(df) # RETRIEVE MISSING VALUES  
df=stringMissingValues(df) # RETRIEVE MISSING VALUES
```

```
def numericMissingValues(file): # replacing ',' and missing values with  
the mean of the year it belongs to  
    try:  
        for item in ALL_NUMERIC:  
            if item not in file.columns:  
                ALL_NUMERIC.remove(item)  
                continue  
        if len(ALL_NUMERIC)==0:  
            print(f'{colors.GREEN}NO NUMERIC COLUMNS TO BE  
PROCESSED.SUCCESFULLY COMPLETED{colors.END}')  
            return file  
        # file=externalIMDb(file) # IMDb  
        for element in ALL_NUMERIC:  
            for index, row in file.iterrows():  
                year = row['year']  
                mean = file[file['year'] == year][element].mean()  
                file[element] = pd.to_numeric(file[element], errors='coerce')  
                if pd.isna(row[element]):  
                    if pd.isna(mean):  
                        file[element]=file[element].ffill()  
                        continue  
                    file.loc[index, element] = mean  
            file[element].bfill(inplace=True)  
    except:  
        print(element,index,row,mean)  
        raise ValueError(f'{colors.RED}A problem occured while processing  
numeric missing values{colors.END}')  
    print(f'{colors.GREEN}NUMERIC MISSING VALUES HAS BEEN SUCCESFULLY  
RESTORED!{colors.END}')  
    return file
```

For numeric missing values we get **external knowledge** only for imdb rating feature. For all the others we take the mean value of each numeric feature according to year the film released.



(e.g Missing value at foreign gross. Year=2007, we complete the missing value with the foreign gross mean of the films released in 2007).

In case we cannot retrieve external knowledge or calculate the mean for a missing value we are doing a ffill() and if this is also not possible we do bfill().

```
def stringMissingValues(file):
    try:
        for item in NO_TRAGET_STRINGS:
            if item not in file.columns:
                NO_TRAGET_STRINGS.remove(item)
                print(f"{item} column removed from the array because it doesn't exist in the dataset")
                continue
            if len(ALL_NUMERIC)==0:
                print(f'{colors.GREEN}NO STRING COLUMNS TO BE PROCESSED.SUCCESFULLY COMPLETED{colors.END}')
                return file
        file=externalGenre(file)
        for index, row in file.iterrows():
            if 'oscar winners' in file.columns:
                if (pd.isna(row['oscar winners'])):
                    file.loc[index, 'oscar winners'] = 0
                else:
                    file.loc[index, 'oscar winners'] = 1
        for j in NO_TRAGET_STRINGS:
            if file[j].isnull().any():
                file[j]=file[j].ffill()
            else:
                file[j]=file[j].bfill()
    except:
        raise RuntimeError(f'{colors.RED}A problem occured while processing string missing values{colors.END}')
        print(f'{colors.GREEN}OTHER MISSING VALUES HAS BEEN SUCCESFULLY RESTORED!{colors.END}')
    return file
```

In string missing values we get **external knowledge** only for genre.

If any other value is missing we are doing a ffill() and if this is also not possible we do bfill().



### External knowledge

In this part of external knowledge we just use a python library of IMDb and we get the genre. If we cannot get the genre then we complete the value as NaN.

```
def externalGenre(file):  
    ia=IMDb()  
    try:  
        for index,row in file.iterrows():  
            if pd.isna(row['genre']):  
                movies = ia.search_movie(row['film'])  
                if movies:  
                    movie = ia.get_movie(movies[0].movieID)  
                    genre=",".join(movie['genres']).lower()  
                else:  
                    print("Movie not found.")  
                    genre=np.nan  
                file.at[index, 'genre'] = genre  
    except:  
        raise RuntimeError(f'{colors.RED}A problem occurred while  
receiving external knowledge in string missing values{colors.END}')  
    return file
```



In this part of external knowledge we just use a python library of IMDb and we get the imdb rating. If we cannot get the rating then we complete the value as NaN and they will be replaced later with the mean, ffill() or bfill()

```
def externalIMDb(file):
    ia = IMDb()
    i=2
    for movieTitle in file['film']:
        movies = ia.search_movie(f"{movieTitle}")
        try:
            if movies:
                movie = ia.get_movie(movies[0].movieID)
                print(f"Title: {movie['title']}")
                year=movie['year']
                rating=movie['rating']
            else:
                print("Movie not found.")
                continue
            print(i, " ", rating)
            if not rating:
                file.at[i-2, 'imdb rating'] = np.nan
                continue
            else:
                file.at[i-2, 'imdb rating'] = rating
                file.at[i-2, 'year']=year
        except:
            file.at[i-2, 'imdb rating'] = np.nan
            continue
        i=i+1
    print(f'{colors.GREEN}EXTERNAL KNOWLEDGE \'IMDb\' HAS BEEN
    SUCCESSFULLY ADDED{colors.END}')
    # file.to_excel("clone.xlsx")
    return file
```



### Note:

this process has already been executed and we copied the data to another file (**moviesUpdated.xlsx**), because the process to fetch data for almost 1400 rows takes a lot of time.

### Data formatting

```
df=columnDataFormating(df)
```

After receiving all the new values and we have already dropped the most of the columns we do not need, we have to format all the data to be the same so the model can work correctly.

```
def columnDataFormating(file):
    try:
        # Convert each value of the df that is ending with the '%' to
        decimal (divide by 100)
        file = file.map(lambda val: float(val.rstrip('%')) / 100 if
isinstance(val, str) and val.endswith('%') else val)
    except:
        raise RuntimeError(f'{colors.RED}A problem occurred while converting
percentages to decimal.{colors.END}')
    try:
        file[ALL_NUMERIC] = file[ALL_NUMERIC].replace(',', ' ',
regex=True).apply(pd.to_numeric, errors='coerce')
        file['budget ($million)'] = file['budget ($million)'] * 1000000
        # file[NO_TARGET_STRINGS] = file[NO_TARGET_STRINGS].replace(',', ' ',
regex=True)
        file['genre'] = file['genre'].str.replace(',', ' ',
).str.replace('.', ' ').str.replace('\s+', ' ',
regex=True).str.strip()
        # file['oscar detail'] = file['oscar
detail'].str.extract(r'([^\(]+)')
        file['oscar detail'] = file['oscar detail'].str.split('(',
n=1).str[0].str.strip()
    except:
        raise RuntimeError(f'{colors.RED}A problem occurred while replacing
characters.{colors.END}')
    print(f"{colors.GREEN}ALL COLUMNS HAS BEEN SUCCESSFULLY
FORMATED!{colors.END}")
    return file
```



### One-Hot encoding

```
df=oneHotEncoding(df) # ONE HOT ENCODING
```

We do have nominal features that we need to represent them with a way to help the model to work better. Those are script type, oscar detail, release date (us), genre

```
# SCRIPT TYPE
if 'script type' in file.columns:
    try:
        one_hot_encoded = pd.get_dummies(file['script type'].apply(lambda
x: next((t for t in TYPES if str(x).startswith(t)), None))).astype(int)
        file = pd.concat([file, one_hot_encoded], axis=1)
        file=dropUseless(file,['script type'])
    except:
        raise RuntimeError(f'{colors.RED}A problem occurred while one-hot
encoding script-type{colors.END}')
```

```
# OSCAR DETAILS
if 'oscar detail' in file.columns:
    try:
        one_hot_encoded = file['oscar detail'].str.get_dummies(',
').astype(int)
        file = pd.concat([file, one_hot_encoded], axis=1)
        file=dropUseless(file,['oscar detail'])
    except:
        raise RuntimeError(f'{colors.RED}A problem occurred while one-hot
encoding oscar details{colors.END}')
```



In the previous ones we did directly normal one-hot encoding but genre and release date (us) need a special treatment.

As far as the genre is concerned there were a lot genres with wrong spelling mistakes or the same genre with another word , for example 'thriller' and 'thrill'. All of those had to be handled before one-hot encoding.

```
# GENRE
if 'genre' in file.columns:
    try:
        file['genre']=file['genre'].str.lower()
        genres=set()
        correctGenre=[]
        for item in file['genre']:
            for word in item.split():
                # correctedWords=[]
                correctedWord=str(Word(word).correct().lower())
                # print(correctedWord)
                correctGenre.append(correctedWord)
        genres.update(correctGenre)
        # print("c:",genres)
        words_to_remove = set()
        for word1 in genres:
            for word2 in genres:
                if word1 != word2 and len(word1) >3 and len(word2) > 3:
                    common_substrings = set([word1[i:i+5] for i in
range(len(word1)-4) if word1[i:i+5] in word2])
                    if common_substrings:
                        shorter_word = word1 if len(word1) < len(word2)
else word2
                        words_to_remove.add(shorter_word)
        genres.difference_update(words_to_remove)
        # print(genres)
        file['genre'] = file['genre'].apply(lambda cell: ' '.join(
            [next((word_set_word for word_set_word in genres if
word_set_word[:3] == word[:3]), word) for word in cell.split()]
        ))
        print('Set:::',genres)
        for genre in genres:
            file[genre] = file['genre'].apply(lambda x: 1 if genre in
x.split() else 0)
        file=dropUseless(file,['genre'])
    except:
```



```
raise RuntimeError(f'{colors.RED}A problem occurred while one-hot encoding genres{colors.END}')
```

As far as the date is concerned before one hot encoding we played with the format and we kept only the released month instead of the whole day, month and year, because we believed all of this information will not give us any advantage.

```
# DATE
if 'release date (us)' in file.columns:
    try:
        file['release date (us)'] = pd.to_datetime(file['release date (us)'], format='mixed')
        # Extract month and day
        file['release date (us)'] = file['release date (us)'].dt.strftime('%m').astype(int)
        monthMapping = {
            1: 'january',
            2: 'february',
            3: 'march',
            4: 'april',
            5: 'may',
            6: 'june',
            7: 'july',
            8: 'august',
            9: 'september',
            10: 'october',
            11: 'november',
            12: 'december'
        }

        # Nominalize the 'release date (us)' column
        file['release date (us)'] = file['release date (us)'].map(monthMapping)
        one_hot_encoded = pd.get_dummies(file['release date (us)'], prefix='').astype(int)
        file = pd.concat([file, one_hot_encoded], axis=1)
        file = dropUseless(file, ['release date (us)'])
    except:
        raise RuntimeError(f'{colors.RED}A problem occurred while one-hot encoding dates{colors.END}')
```

### Note:

Of course, after one-hot encoding we drop each feature that got one-hot encoded.





### Delete duplicates

```
df=deleteDuplicate(df) # CHECK FOR DUPLICATE ROWS
```

after all the process we have to check if we have duplicate data.

```
def deleteDuplicate(file):  
    try:  
        if (file.duplicated().sum() != 0) or (not  
file[file.duplicated(subset=['film'])].empty):  
            print(f'The dataset contains  
{(file.duplicated(subset=["film"])).sum()} duplicate films that need to  
be removed.')  
            print(f'The dataset contains {file.duplicated().sum()} duplicate  
rows that need to be removed.')  
            file.drop_duplicates(inplace=True)  
            file = file.drop_duplicates(subset=['film'], keep='first')  
    except:  
        raise RuntimeError(f'{colors.RED}A problem occurred while deleting  
duplicates{colors.END}')  
    print(f"{colors.GREEN}DUPLICATE ROWS HAVE BEEN SUCCESSFULLY  
DELETED!{colors.END}")  
    return file
```

we check the whole subset for duplicates, but also, we check the subset film for duplicate film names. If a film has a duplicate name, then we keep the data of the film we first found in our dataset.



### Extra functions

We also created some functions for the preprocessing such as scaling and normalization because we do not know yet what classification model we are going to use and what are its prerequisites for it to work.

```
def scaling(file):
    for item in ALL_NUMERIC:
        if item not in ALL_NUMERIC:
            ALL_NUMERIC.remove(item)
            continue
    if len(ALL_NUMERIC) != 0:
        try:
            scaler = StandardScaler()
            file[ALL_NUMERIC] = scaler.fit_transform(file[ALL_NUMERIC])
        except:
            raise ValueError(f'{colors.RED}A problem occurred while scaling values{colors.END}')
    print(f'{colors.GREEN}SCALING HAS BEEN SUCCESSFULLY COMPLETED!{colors.END}')
    return file

def normalization(file):
    for item in ALL_NUMERIC:
        if item not in ALL_NUMERIC:
            ALL_NUMERIC.remove(item)
            continue
    if len(ALL_NUMERIC) != 0:
        try:
            columns_to_normalize = file[ALL_NUMERIC]
            scaler = MinMaxScaler()
            normalized_columns = pd.DataFrame(scaler.fit_transform(columns_to_normalize), columns=ALL_NUMERIC)
            file[ALL_NUMERIC] = normalized_columns
        except:
            raise ValueError(f'{colors.RED}A problem occurred while normalising values{colors.END}')
    print(f'{colors.GREEN}NORMALISING HAS BEEN SUCCESSFULLY COMPLETED!{colors.END}')
    return file
```



### Result

```
INITIALIZATION HAS BEEN SUCCESSFULLY CONVERTED!
USELESS COLUMNS HAS BEEN SUCCESSFULLY DELETED!
ALL COLUMNS HAS BEEN SUCCESSFULLY FORMATED!
NUMERIC MISSING VALUES HAS BEEN SUCCESSFULLY RESTORED!
OTHER MISSING VALUES HAS BEEN SUCCESSFULLY RESTORED!
ALL COLUMNS HAS BEEN SUCCESSFULLY FORMATED!
USELESS COLUMNS HAS BEEN SUCCESSFULLY DELETED!
USELESS COLUMNS HAS BEEN SUCCESSFULLY DELETED!
USELESS COLUMNS HAS BEEN SUCCESSFULLY DELETED!
USELESS COLUMNS HAS BEEN SUCCESSFULLY DELETED!
ONE HOT ENCODING HAS BEEN SUCCESSFULLY COMPLETED!
The dataset contains 6 duplicate films that need to be removed.
The dataset contains 0 duplicate rows that need to be removed.
DUPLICATE ROWS HAVE BEEN SUCCESSFULLY DELETED!
USELESS COLUMNS HAS BEEN SUCCESSFULLY DELETED!
-----PRE-PROCESSING-FINISHED-----

DataFrame does not contain NaN values.
  average critics  average audience  opening weekend  domestic gross  ...  _may  _november  _october  _september
0             56.0             80.0       70885301       210614939  ...    0         0         0         0
1             82.0             80.0       14035033       53606916  ...    0         0         0         1
2             52.0             61.0       15951902       39568996  ...    0         0         1         0
3             55.0             78.0        3824988       24343673  ...    0         0         1         0
4             22.0             38.0       10059425       41797066  ...    0         0         0         0

[5 rows x 83 columns]
PS C:\Users\kokka\GitHub projects\Oscars---Data-Mining> |
```

### Important Note:

The features that have been selected to be used they are not the final as long as the classification model might not give us the best results. So modifications to USELESS\_COL might happen.



## B1.CLASSIFICATION

Before we start explaining about the classification process is important to mention that the features we use now are different. That is because with the previous features we could not receive very good results so the model was overfitting or it was not giving any results. The way we pre process the data is exactly the same. The only thing that changes is the columns we use.

New Dropped Columns:

```
USELESS_COL=['id','imdb vs rt disparity','oscar detail','distribu-  
tor','primary genre','domestic gross ($million)','foreign gross ($mil-  
lion)','worldwide gross ($million)','worldwide gross','opening weekend  
($million)','metacritic audience','genre','script type','release date  
(us)','opening weekend','average critics','foreign gross','rotten toma-  
toes audience']
```

The way we decided which columns to drop since we were not getting good results was running the model and checking at the same time the importance of each feature. Features with zero or very low importance were deleted.

```
<bound method NDFrame.head of                                     Feature  Importance  
27          metacritic critics          0.195014  
30          of gross earned abroad      0.124744  
12      budget recovered opening weekend  0.116819  
45          worldwide gross            0.081943  
11          budget recovered            0.077720  
6          average audience             0.073446  
22          foreign gross ($million)     0.051833  
37          rotten tomatoes audience     0.048502  
42          thrilled                    0.044097  
26          metacritic audience          0.042145  
16          domestic gross              0.038523  
9          biography                    0.030103  
38          rotten tomatoes critics      0.027516  
33          original                    0.017639  
10          budget ($million)            0.016786  
17          domestic gross ($million)    0.013170  
39      rotten tomatoes vs metacritic deviance  0.000000  
29          mystery                     0.000000  
8          based on a true story         0.000000  
13          comedy                      0.000000  
14          crime                       0.000000  
15          documentary                  0.000000  
18          drama                       0.000000  
19          family                      0.000000  
20          fantasy                      0.000000  
21          foreign gross                0.000000  
1          action                       0.000000  
24          horror                      0.000000  
46          worldwide gross ($million)   0.000000>
```



## STEPS

1. Load Data
2. Prepare Data
3. Train Model
4. Predict

```
if __name__ == '__main__':  
    argOne, argTwo, argThree = handleArgs()  
  
    trainDataset, predictDataset = preprocess(argThree)  
    trainTarget, trainData = separateData(trainDataset) #seperate train from target data of the  
    predictData = fixTestFile(trainData, predictDataset)  
    scaledTrainData, scaledPredictData = scaleData(trainData, predictData)  
    model, X_train, X_valid, y_train, y_valid = doTraining(scaledTrainData, trainTarget, argOne)  
    y_pred = model.predict(X_valid)  
    if sys.argv[2] == 'stats':  
        printStats(model, y_valid, y_pred, scaledTrainData, trainTarget)  
    if sys.argv[1] != 'knn' and sys.argv[1] != 'lr':  
        featureImportances(model, trainData)  
    doPredictions(model, scaledPredictData, predictData) # Predict in the model
```

### Load and Prepare Data

Before the training of the model and before making the predictions we have to be sure that there are no missing values in both training and test datasets as long as the two datasets must have exactly the same features. We also have to separate features from target in the training dataset and scale the data using MinMax scaler.

So first we preprocess both datasets

```
def preprocess(command='prepro'):  
    if command == 'prepro':  
        dp = DataPreprocessor(TRAIN_PATH, TRAIN_PATH_PROCESSED)  
        df = DataPreprocessor(PREDICT_PATH, PREDICT_PATH_PROCESSED)  
        trainDataset = dp.executePreprocess()  
        predictDataset = df.executePreprocess(predict=True) #options: pre-  
dict=True/False  
    else:  
        # SAVING TIME-----  
        trainDataset = pd.read_excel(TRAIN_PATH_PROCESSED, sheet_name =  
'Sheet1')  
        predictDataset = pd.read_excel(PREDICT_PATH_PROCESSED, sheet_name =  
'Sheet1')  
        print('all files has been succesfully preprocessed')  
        return trainDataset, predictDataset
```



and right after we make sure that both datasets have the same columns

```
def fixTestFile(train_file, test_file):  
    diff_train_columns = set(train_file.columns) - set(test_file.columns)  
    diff_test_columns = set(test_file.columns) - set(train_file.columns)  
  
    for test_column in diff_test_columns:  
        for train_column in diff_train_columns:  
            # Check if the columns share at least four letters  
            if len(set(test_column).intersection(train_column)) >= 4:  
                # Replace the column name in test_file  
                test_file = test_file.rename(columns={test_column:  
train_column})  
                break # Stop searching for similar columns once a match  
is found  
  
    # Add missing columns to test_file with values set to 0  
    for train_column in diff_train_columns:  
        test_file[train_column] = 0  
    test_file = test_file.loc[:, ~test_file.columns.duplicated(keep='first')]  
    test_file = test_file.drop(columns= set(test_file.columns) -  
set(train_file.columns))  
  
    test_file.to_excel(PREDICT_PATH_PROCESSED)  
    # print(train_file.shape, test_file.shape)  
    return test_file.sort_index(axis=1)
```

separating features from target and scaling

```
def seperateData(ds):  
    if 'oscar winners' not in ds.columns:  
        raise ValueError('Oscar winners not in the dataset')  
    target=ds['oscar winners']  
    data=ds.drop(columns='oscar winners')  
    data=data.sort_index(axis=1)  
    return (target, data)
```

```
def scaleData(trainData, predictData):  
    scaler=MinMaxScaler()  
    scaledTrainData=scaler.fit_transform(trainData)  
    scaledPredictData=scaler.transform(predictData)  
    return scaledTrainData, scaledPredictData
```



### Model Training

Now that the data are ready we can proceed to train the model by splitting our train dataset in train and test data.

```
def doTraining(scaledTrainData,trainTarget,modelName='knn'):  
    X_train, X_valid, y_train, y_valid =  
train_test_split(scaledTrainData,trainTarget, test_size=0.25,ran-  
dom_state=42)  
    if modelName=='rf':  
        model = RandomForestClassifier(random_state=42)  
    elif modelName=='lr':  
        model = LogisticRegression(max_iter=1500, random_state=42)  
    elif modelName=='dtt':  
        model = DecisionTreeClassifier(random_state=42)  
    elif modelName=='knn':  
        model=KNeighborsClassifier(n_neighbors=3)  
    model.fit(X_train, y_train)  
    return model,X_train,X_valid, y_train, y_valid
```

after the splitting we try different models to see which one works better for our dataset

#### RANDOM FOREST

```
PS C:\Users\kokka\GitHub projects\Oscars---Data-Mining> python3 main.py rf stats noprepro  
all files has been succesfully preprocessed  
Confusion Matrix:  
[[325   5]  
 [ 11   8]]  
Classification Report:  
              precision    recall  f1-score   support  
  
     0       0.97         0.98         0.98         330  
     1       0.62         0.42         0.50          19  
  
   accuracy          0.95         0.95         0.95         349  
  macro avg          0.79         0.70         0.74         349  
 weighted avg          0.95         0.95         0.95         349  
  
Accuracy: 0.9541547277936963  
Cross validation: 0.9684767025089606
```

```
#Oscar winners:  5
```



### DECISION TREE CLASSIFIER

```
PS C:\Users\kokka\GitHub projects\Oscars---Data-Mining> python3 main.py dtc stats noprepro
all files has been succesfully preprocessed
Confusion Matrix:
[[324  6]
 [ 5 14]]
Classification Report:
              precision    recall  f1-score   support

     0       0.98      0.98      0.98        330
     1       0.70      0.74      0.72         19

   accuracy          0.97
  macro avg       0.84      0.86      0.85
 weighted avg     0.97      0.97      0.97
```

```
#Oscar winners: 154
```

### K-NEAREST NEIGHBOURS

```
PS C:\Users\kokka\GitHub projects\Oscars---Data-Mining> python3 main.py knn stats noprepro
all files has been succesfully preprocessed
Confusion Matrix:
[[328  2]
 [ 11  8]]
Classification Report:
              precision    recall  f1-score   support

     0       0.97      0.99      0.98        330
     1       0.80      0.42      0.55         19

   accuracy          0.96
  macro avg       0.88      0.71      0.77
 weighted avg     0.96      0.96      0.96

Accuracy: 0.9627507163323782
Cross validation: 0.9613082437275985
#Oscar winners: 7
PS C:\Users\kokka\GitHub projects\Oscars---Data-Mining>
```

### LOGISTIC REGRESSION

Linear Regression and other algorithms like SVC are not suitable for our dataset as long as they are ill defined and they don't return any useful results.

```
Confusion Matrix:
[[330  0]
 [ 19  0]]
Classification Report:
              precision    recall  f1-score   support

     0       0.95      1.00      0.97        330
     1       0.00      0.00      0.00         19

   accuracy          0.95
  macro avg       0.47      0.50      0.49
 weighted avg     0.89      0.95      0.92
```

As we can see the True Negative is 0 which means 0 oscar winners and that is why it is ill-defined, so not suitable.





## B2.PREDICTIONS

As we can see in the models above the best results recall and f-score is coming from Decision Tree classifier. But at the same time it gives too many oscar winners which could mean overfitting. For our predictions **the model we use is KNN** because it is the second better from DTC with recall=0,42 f-score=0,72, Accuracy=96 and Cross-Validation accuracy=96.

```
PS C:\Users\kokka\GitHub projects\Oscars---Data-Mining> python3 main.py knn stats noprepro
all files has been succesfully preprocessed
Confusion Matrix:
[[328  2]
 [ 11  8]]
Classification Report:
              precision    recall  f1-score   support

     0       0.97       0.99       0.98        330
     1       0.80       0.42       0.55         19

   accuracy       0.88       0.71       0.96        349
  macro avg       0.88       0.71       0.77        349
weighted avg       0.96       0.96       0.96        349

Accuracy: 0.9627507163323782
Cross validation: 0.9613082437275985
#Oscar winners: 7
PS C:\Users\kokka\GitHub projects\Oscars---Data-Mining>
```

The predictions are located in 'Data/predictions.csv'

## HOW TO RUN THE CODE

```
python3 main.py <lr/knn/rf/dtc> <stats/nostats> <prepro/noprepro>
```

example:

```
python3 main.py rf nostats prepro
```

To run the final and correct model use this command:

```
python3 main.py knn stats prepro
```



To run data preprocessing separately

Command:

```
python3 dataPreprocessing.py
```

it will run this:

```
if __name__=='__main__':  
    dp=DataPreprocessor("./moviesUpdated.xlsx","final.xlsx")  
    dataset=dp.executePreprocess()  
    if dataset.isna().any().any():  
        print("DataFrame contains NaN values.")  
    else:  
        print("DataFrame does not contain NaN values.")  
    print(dataset.head())
```

### NOTE:

It is recommended to use **moviesUpdated.xlsx** because it is already filled with the imdb rating missing values.

In the file **movies.xlsx** imdb rating missing values are missing in all the rows so it will take a lot of time to fill all the missing values .

**moviesUpdated.xlsx** is a copy of **movies.xlsx** with the only difference that **moviesUpdated.xlsx** has also the IMDb ratings ready.