# Statistical Machine Learning Mini Project 2022 - Do (wo)men talk too much in films?

**Dap De Bruijckere**
Engineering physics

**Felix Dryselius**
Social techincal systems engineering

**Vidar Lundgren**
Engineering physics

**Irina Zarankina**
Financial mathematics

## Abstract

This project studies gender inequality in films and tries to create a classification model that can predict gender of the main character, given some data. In order to achieve this, multiple classification methods have been tuned and evaluated. In the end, QDA was chosen as the best model to predict gender of lead character.

## 1  Data analysis task

Questions to answer:

- Do men or women dominate speaking roles in Hollywood movies?

- Has gender balance in speaking roles changed over time (i.e. years)?

- Do films in which men do more speaking make a lot more money than films in which women speak more?

### 1.1  Answers

The data comes from the Film dialog data set created by Hanah Anderson and Matt Daniels in 2016 link. The full set contains 2000 entries but in this article a sub-set of 1039 randomly selected entries were used. This set is skewed with almost all features having a positive skewness value. In total skewness ranged from -1.27 to 3.79 and this was calculated using pandas.dataFrame.skew() function (6). The data is also imbalanced, with 785 points having class male and 254 having female. This means that a model that only predicts male will have a misclassification error of circa 24% ($\frac{254}{1039}$).

As visible from 1 gender imbalance is present in Hollywood. Women have both fewer roles and speak fewer words on average per movie for all years in the data set except one. The only exception is during 1958 when the proportion of words spoken by females exceeds those of men. This year might be an outlier though, because the data set is heavily skewed, with the mean of sampled films being in the year 1999.

Gender imbalance seems to improve over time though, and interestingly, movie gross (money earned) does not seem to correlate with male or female dominance.
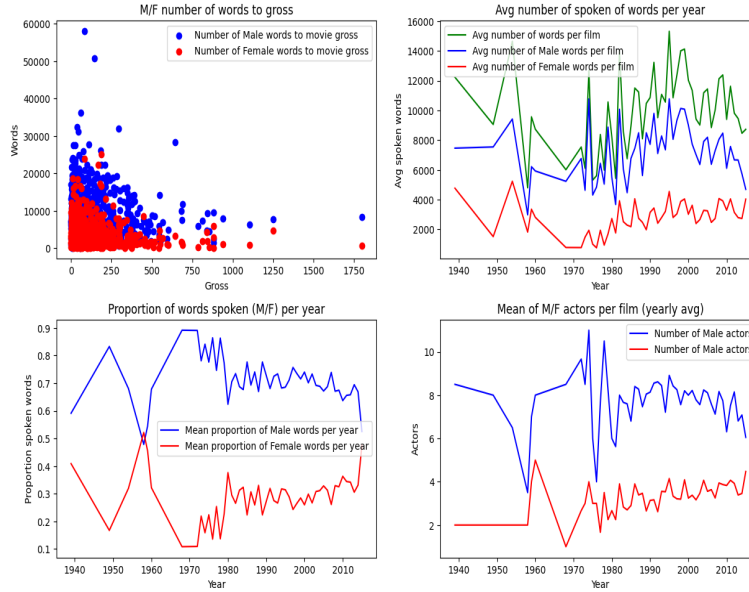
Figure 1: Inequality over time

# 2 Description of methods used

## 2.1 Logistic regression

Logistic regression can be seen as a modification of the linear regression model so it can be applied to a classification problem. This modification is obtained by using the logistic function to model a binary dependent variable based on some input variables. More specifically, let $g(x)$ be a function that approximates the conditional probability of the positive class, $g(x) = \frac{\epsilon^{\theta^T x}}{1+\epsilon^{\theta^T x}}$, where $\theta^T x$ is the linear regression. Thus, the linear regression can be "squeezed" into an interval [0,1] by using the logistic function (1). Then numerical optimisation is applied for learning the parameters $\theta$ of the model.

## 2.2 Discriminant analysis

The discriminant analysis methods: Linear Discriminant Analysis (LDA), and Quadratic Discriminant Analysis (QDA), are based on fitting a gaussian probability density to each class. For LDA, the same covariance matrix is assumed for each class, giving linear decision boundaries. QDA, as the name suggests, gives quadratic decision boundaries instead, due to the omission of the aforementioned assumption of covariance matrices being the same.

## 2.3 K-nearest neighbour

kNN is a distance-based method, that in classification, predicts class of a new data point by taking the majority vote from the k-nearest data points in the training set. It is a non-parametrized method and data must be normalized before use 8.

## 2.4 Tree-based methods

Tree-based methods work by splitting the data with as little error as possible. This was done by using gini index, which divides the data by creating the lowest amount of impurities on both sides. Because of that any redundant or useless data does not affect the accuracy of the method, but affects the runtime. To determine the best depth cross validation was used. Also when finding the minimum sample size, the minimum size that a node needs to be to allow it to be split, cross validation was used. To visualize the tree and see which data was used, graphviz was implemented.
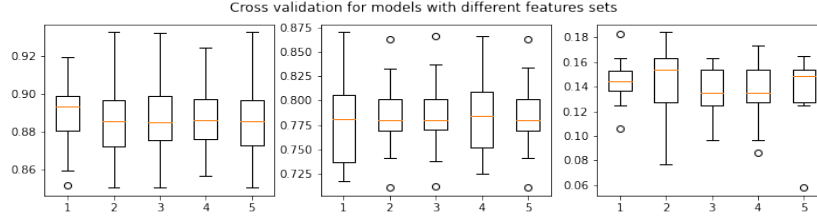
Figure 2: ROC AUC, PC AUC and error of 5 best log regression models

# 3 How the methods were applied to the data

## 3.1 Logistic regression

The base logistic regression model was implemented on the given data using scikit-learn package. The data was scaled using the 'StandardScaler' from the 'sklearn' library, since regularization works better when data is standardized. Evaluation metrics ROC AUC, PC AUC and misclassification error were used while tuning and evaluating performance.

At first, some assumptions were tested manually to tune the parameters. The data was divided into training and validation sets and a random seed was fixed to provide comparability and reproducibility of experiments. It was discovered that a 'liblinear' solver gave the best performance. Using parameter 'class weight' of the logistic regression function didn't improve the results.

After testing different subsets of features as input variables to the model (brute-force approach using library 'itertools'), it became clear that models with 8 and 9 input variables gave the best values of evaluation metrics. Thus, complexity of the models was reduced, which resulted in the improved performance of the models. (More details will follow in feature importance section) Then the 5 models with the best performance (ROC AUC) were compared using 10-fold cross-validation based on evaluation metrics: ROC AUC, PC AUC, misclassification error (the threshold rate was chosen comparing f1 coefficients).

As a result, the model number 3 was chosen (see the Figure 2 above). Then 10-fold cross-validation using randomized search was performed to determine if regularization l1 or l2 is needed. Cross-validation was also used to find regularization coefficient (the interval from 0.0001 to 100 of possible values of parameter C was checked, which is an inverse of the regularization parameter $\lambda$ (1), ($C = 1/\lambda$ ) and threshold rate.

The chosen logistic regression model had the following features:

The list of features: ['Total words', 'Age Co-Lead', 'Number of female actors', 'Difference in words lead and co-lead', 'Number of words lead', 'Age Lead', 'Number words female', 'Number of male actors'] 'C': 18.31792545756584, 'penalty': 'l1', 'solver': 'liblinear', threshold rate 0.435

## 3.2 Discriminant analysis

Discriminant analysis does not have any hyperparameters that can be tuned, which can work as either an advantage or disadvantage. This means the implementation is fairly straight forward, and that the method will often perform well right away, as has proven to be the case in practice. (2)

Using the important features discussed in section 3.1 was the only way the methods were tuned to the problem. As some features were collinear, this greatly improved QDA performance in particular.

## 3.3 K-nearest neighbour

Application of kNN to the data was done in three steps: analysis and normalization of the data, decision of hyper-variables (k) and estimation of E[Accuracy_new], and estimation of other evaluation terms. For results see 4.2.

Normalization of the data *(all features except "Lead" was interpreted as quantitative)* was performed by using sklearn's StandardScaler method (5). This assumes that all features behave like a random
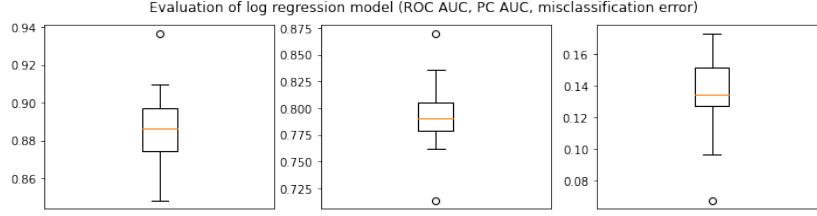
Figure 3: ROC AUC, PC AUC and error of 5 best logistic regression models

variable with normal distribution, and it normalizes each feature with mean 0 and standard deviation 1. Another re-scaler, "MinMaxScaler", was also evaluated but not used, because no significant performance difference was observed.

kNN was implemented using sklearn's KNeighborsClassifier with uniform weights, Euclidian distance, and the 'distance calculation algorithm' set to auto (in this generally low dimensional problem, the "KD Tree algorithm" was probably automatically chosen). This was done to reduce solution complexity and minimize calculation time (4). kNN was also evaluated using distance based weights but this was not implemented due to reduced performance. When plotting the missclassification error against the value of "k" with distance based weights it became clear that kNN only performed well for a specific k-value. This indicated that the model with distance weights would risk being over-fitted in a production environment, see 8.

For the evaluation of kNN, the whole data set was used with no splitting in a test and train data set. This was done after realising that performance improved considerably when using all 1039 entries. For calculations and estimations, cross validation was used instead (1). Cross validation was also used to estimate other performance data than missclassification error by taking the average or sum of the "evaluation term" in each fold.

The value of "k" was set in order to minimize missclassification error and was chosen from a range of 1-80. This range was analysed using cross-validation with 200 folds and the results can be seen in table 4.2. These ranges were capped at these values when the results started converging. Because the data set was imbalanced, evaluation was also performed with different thresholds. One tested the threshold was r=24% for positive class ('Male') and these results are also visible in 4.2. This threshold was chosen because it mimics the statistical chance of randomly selected data point having class ('Female').

### 3.4 Tree-based methods

The tree-based model was implemented using the 'DecisionTreeClassifier' from the 'sklearn' library. This uses the Gini index to partition the data several times until the depth is reached or the minimum amount of samples in a node is too low. The Gini index is calculated as $G = \Sigma_{i=1}^{2} p(i) * (1 - p(i))$, where G is the Gini index and p is the probability that it is a male or a female lead. This is done with all the possible places to split and with all the data. The combination with the lowest Gini index is used. Instead of the Gini index entropy can be used, but it wont make a difference in the accuracy of the tree. To determine both the best depth and the minimum sample a 40 fold cross validation was used 10 times. the range of depth to be evaluated was between 2 and 15 and the minimum sample was between 1 and 40.

## 4 Performance evaluation

### 4.1 Logistic regression

Results of the model's performance evaluation using 10-fold cross validation are shown in table 1.

### 4.2 K-nearest neighbour evaluation

The results from the model are shown in table 4.2. As is visible, the model tends to over-predict on the positive class (Male) with the FPR being around 73%, however do to the imbalance in the data

4

Table 1: Results from logistic regression evaluation

| Term | Value |
|---|---|
| ROC AUC | 0.887 |
| PC AUC | 0.792 |
| TPR | 0.675 |
| FPR | 0.068 |
| E[Accuracy_new] | 0.867 |
| F1 | 0.712 |

set this does not translate to a high missclassification rate. Note also that changing the threshold does not translate into a similar change of the results, this is somewhat expected given that kNN is a non-parametrized model. For a full table and graphs please see A and B.

Table 2: Results from kNN evaluation with two thresholds

| Term (r=50%) | Values | Term (r=24%) | Values |
|---|---|---|---|
| Optimal_k | 16 | Optimal_k | 5 |
| TPR | 0.971975 | TPR | 0.950318 |
| FPR | 0.728346 | FPR | 0.700787 |
| E[Accuracy_new] | 0.80077 | E[Accuracy_new] | 0.791145 |
| F1 | 0.880554 | F1 | 0.873025 |

Table 3: Discriminant Analysis Confusion Matrices

| LDA | Term | QDA | Term |
|---|---|---|---|
| Accuracy | 0.860 | Accuracy | 0.888 |
| TPR | 0.975 | TPR | 0.868 |
| FPR | 0.497 | FPR | 0.227 |

## 4.3   Discriminant analysis

Because there is no tuning in LDA or QDA, the methods could not be evaluated using cross-validation. Instead, 200 runs were performed with random splits of the test and training data for each run. The accuracy was estimated by the mean average accuracy of the 200 runs. Please see 3 for results.

## 4.4   Tree-based method

Cross validation showed that the best depth was either 8 or 9 with the average being 8.6. It also showed that the minimum sample size varied from 7 to 15,

Table 4: Results from tree-based method evaluation

| Term | Value |
|---|---|
| TPR | 0.862 |
| FPR | 0.399 |
| E[Accuracy_new] | 0.805 |
| missclassification | 0.204 |
| Best depth | 8.6 |
| Best min sample size | 12.1 |

## 4.5   Gini importance

The Gini importance looks at how many times each feature has been used to split a node, divided by all the split in the tree. figure(4). shows that "Number words female", "Number of female actors" and
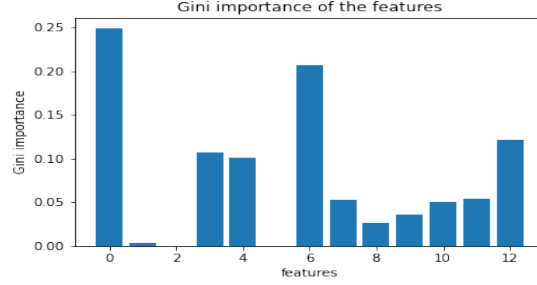
5

Figure 4: Inequality over time

"Age Co-Lead" are the most used and "Total words", "Number of words lead" and "Year" are not used.

# 5   Choice of method

Usually when picking which method to use, one would look at a comparison of different performance metrics for the available methods and weigh in how important each metric is for the problem at hand. For example, if you were working on a machine learning algorithm for identifying disease of some kind, one of the top concerns would be to minimise the false negative rate so not to risk the condition go untreated. For this project, the outcome of false negatives and false positives does not matter in the same way; therefore, our primary metric of choice is accuracy, $\frac{TN+TP}{n}$

In place of unseen data, cross validation (when applicable) has been used to estimate the accuracies in table 5. From the table it is evident that QDA gives the best results and should be used "in production".

Table 5: Accuracies

| Method | Accuracy |
|---|---|
| Logistic Regression | 0.867 |
| LDA | 0.860 |
| QDA | 0.888 |
| Tree-based | 0.804 |
| k-NN | 0.801 |

# 6   Conclusions

It is clear that model performance depends a lot on the features used. With the feature optimization performed, the best performing model is QDA: Quadratic Discriminant Analysis when studying accuracy, TPR and FPR. If a new test set was given we estimate that QDA would have an accuracy of 0.888.

# 7   Feature importance and selection

Based on the data analysis and correlation matrix between features (figure 5), the assumption was made that some variables would not be significant for the model.

Running logistic regression with different number of features showed the following – models with 8 and 9 features had the the highest ROC-AUC, PC-AUC and misclassification error values.

ROC-AUC was chosen as the main metric based on following factors:

- During manual experimentation it was discovered that if one metric improves or degrades, so does the others in the majority of cases.
- It is easier to compare logistic regression models by ROC-AUC and PC-AUC as there is no need to take in account a threshold rate

- The data is imbalanced, but not severely imbalanced

As there could have been several models with close values of ROC-AUC, the decision was made to examine which features subsets gave the best average results (ROC-AUC) on 20 runs. As a result, 5 best subsets of features were selected: 3 with 8 features, 2 with 9 features.

The best subsets of features varied just by 1-2 features. The most important features are those, which were in all these subsets. More specifically, 'Number words female', 'Age Co-Lead''Number of female actors', 'Difference in words lead and co-lead', 'Number of words lead', 'Age Lead', 'Number of male actors'. 'Age Lead' was also a popular feature.

The feature 'Number words female' was included in all subsets of features, which gave the best ROC-AUC metric. Moreover, if we look at the models' coefficients, this factor had the highest coefficient value (approx. 2.0). The feature 'Number words male' was in two of the 5 feature subsets with best performance, and its coefficients were not that heavy, just approximately 0.6.

Such features as Year of release ('Year') and Money made by film ('Gross') were not included in the feature subsets at all. So, these factors give worse prediction than chosen features. They are not important as when we try the logistic regression model with all features, their weight coefficients are very low (<0.1), while 'Number words female' and 'Number words male' have higher coefficients.

The models with one variable didn't perform well, they didn't have true positives or even false positive at all, so they performed like the worst-case classifer, which always predicts the same output class. It has a misclassification error approx. 25 percent according to the proportion of female lead roles to the number of all movies.

## 8   Discussion

The model choice was based on evaluation using an accuracy metric. However, the models differ in various metrics, for example, TPR, FPR, F1. These metrics demonstrate the characteristics of the models that differ from accuracy.

In real life it would be right to choose an evaluation metric or combination of metrics depending on how the results of the prediction would be used.

For example, assume, we have a client, i.e. a person who will use the model. If the main reason of the prediction is to obtain the list of films with women leads, then TPR, precision (the proportion of TP to all predicted positives), and F1 metrics should be taken in account as well as accuracy. This is because the data is imbalanced, and we can have a case when true positives are 100% correct, but we have 13% misclassification rate. This is, in fact, would be the ratio of false positives to data. Thus, if we had 100 films, we would get 13 FP and 25 TP (given the data is imbalanced in proportion 1:3). I.e., every third film from the list with women leads wouldn't belong there.

Maybe it is not important how many films with men leads got in the list with women leads and the client cares just about TPR and accuracy.

Another case would occur if both lists of films with men and women leads should be used by the client and he/she doesn't care how many films with men leads got on the list with women leads. We would be interested in the accuracy and some balance between TPR and FPR in this case.

Probably, if the client wants just one list with films classified by the gender of the lead actor, the most important is to obtain the minimum misclassification rate.

Since the project description does not specify how the model's predictions will be used, accuracy was chosen as a "safe" bet. However, if TPR and FPR are considered, the best model could be log regression as it has low FPR (0.068), good values of TPR (0.675), and the accuracy (0.867).

# References

[1] A. Lindholm, N. Wahlström, F. Lindsten, T. B. Schön, MACHINE LEARNING - *A First Course for Engineers and Scientists*, Draft version: April 30, 2021

[2] 1.2. Linear and Quadratic Discriminant Analysis, Scikit-learn.org. Retrieved February 19, 2022.

[3] Scikit-learn documentation, (KNeighborsClassifier) Scikit-learn.org. Retrieved February 21, 2022.

[4] Scikit-learn documentation, (Nearest Neighbors) Scikit-learn.org. Retrieved March 3, 2022.

[5] Scikit-learn documentation, (StandardScaler) Scikit-learn.org. Retrieved February 21, 2022.

[6] Pandas documentation, (DataFrame.skew) pandas.pydata.org. Retrieved February 21, 2022.

# A  Table Appendix

Table 6: Results from kNN evaluation with two thresholds

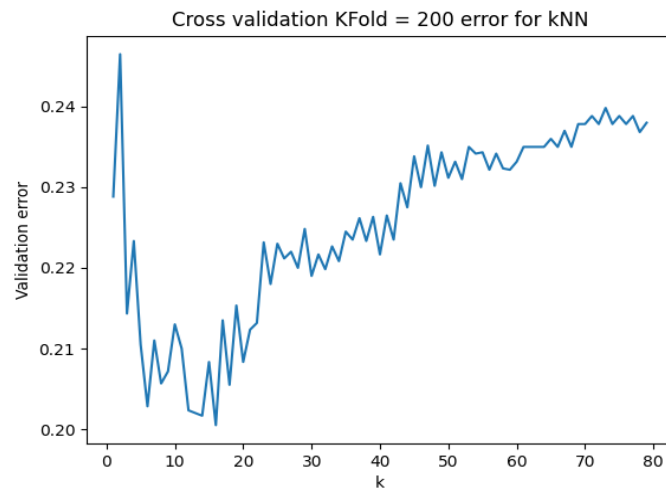| Term (r=50%) | Values | Term (r=24%) | Values |
|---|---|---|---|
| P | 785 | P | 785 |
| N | 254 | N | 254 |
| P_star | 948 | P_star | 924 |
| N_star | 91 | N_star | 115 |
| TN | 69 | TN | 76 |
| FP | 185 | FP | 178 |
| FN | 22 | FN | 39 |
| TP | 763 | TP | 746 |
| Optimal_k | 16 | Optimal_k | 5 |
| TPR | 0.971975 | TPR | 0.950318 |
| FPR | 0.728346 | FPR | 0.700787 |
| E[Accuracy_new] | 0.80077 | E[Accuracy_new] | 0.791145 |
| E[Error_new] | 0.2005 | E[Error_new] | 0.2165 |
| precision | 0.804852 | precision | 0.807359 |
| recall | 0.971975 | recall | 0.950318 |
| F1 | 0.880554 | F1 | 0.873025 |

# B  Graph Appendix



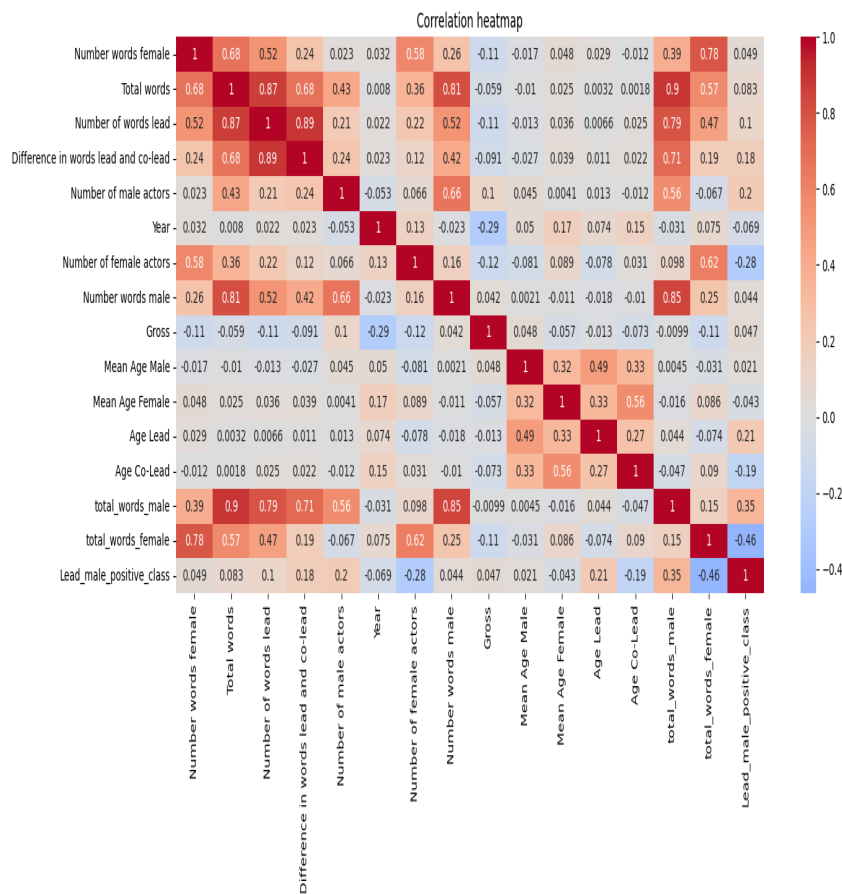Figure 6: kNN, CV missclassification to number of 'k', r=50%

9

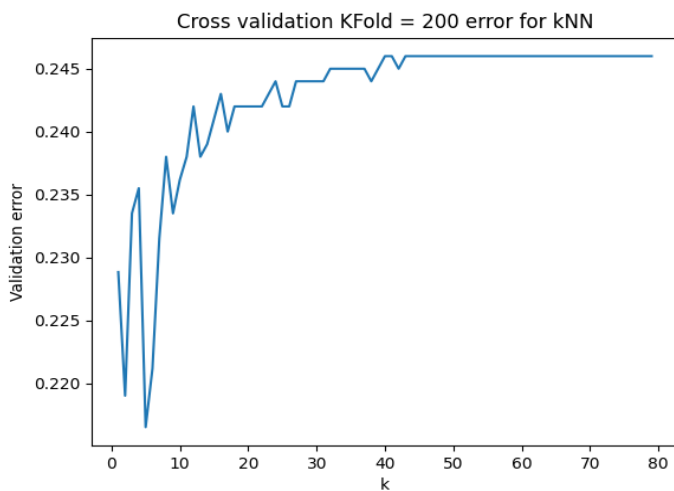Figure 5: Feature correlation heatmap



Figure 7: kNN, CV missclassification to number of 'k', r=24%
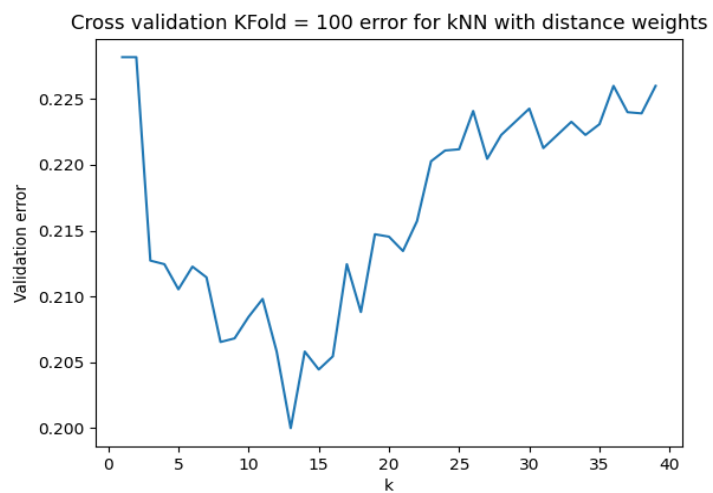
10

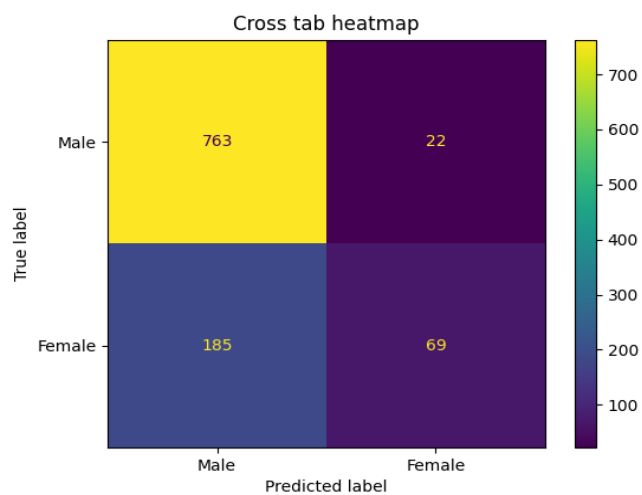Figure 8: kNN, CV missclassification to 'k', distance weights, r=50%
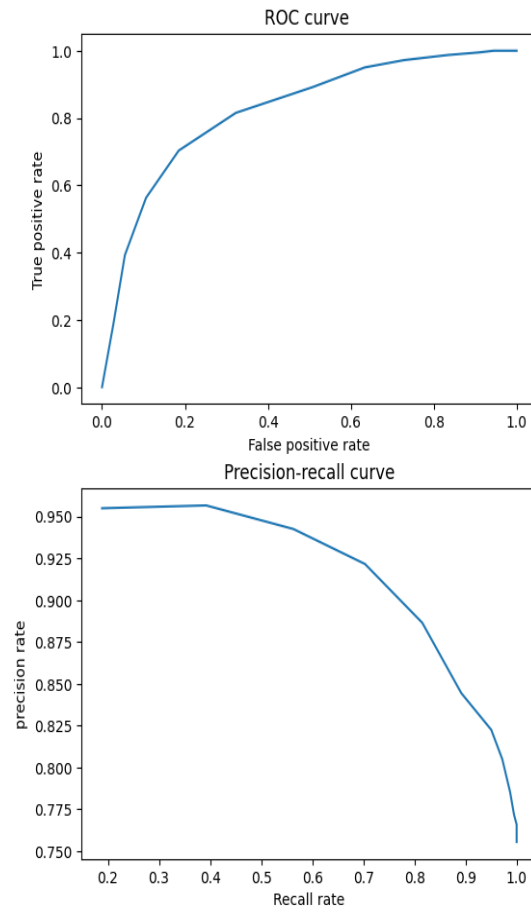


Figure 9: kNN, Crosstab heatmap, r=50%

Figure 10: kNN, ROC and Precision-recall

# C   Code Appendix

## C.1   Data equality plotting

```python
# ----------------- Data analysis ------------------

# This program creates the plots to answers the task:
# "Data analysis task"


import pandas as pd
import matplotlib.pyplot as plt
import os

cwd = os.getcwd()
URI = (cwd+"\\train.csv")
film_data = pd.read_csv(URI, dtype={"Lead":str}).dropna().reset_index(
    drop=True)

# Creating dummy values of the category "Lead"
film_data_dummies = pd.get_dummies(film_data).copy()

# Adding two columns that summerizes the number
# of words M/F with the lead's words if  it is
# the same gender
film_data_dummies = film_data_dummies.assign(total_words_male=lambda
    row: (row['Number words male'] + row['Number of words lead'] * row
    ['Lead_Male']))
film_data_dummies = film_data_dummies.assign(total_words_female=lambda
     row: (row['Number words female'] + row['Number of words lead'] *
    row['Lead_Female']))


# Only saving "Lead_Male" as category and
# renaming it to "Lead"
film_data_dummies.drop("Lead_Female",1, inplace=True)
y = film_data_dummies['Lead_Male'].rename("Lead")

# Creating the features set without the category
x = film_data_dummies.drop(columns=['Lead_Male'])

# ----------------- Scatter plot ------------------
# Creates a scatter plot of movie earnings vs
# total number words M/F The aim is to determine
# if movie earnings correlates to M/F dominace

df = film_data_dummies.copy()

plt.subplot(2,2,1)
ax_male_word_gross = plt.scatter(df["Gross"], df["total_words_male"],
    color="b", label="Number of Male words to movie gross")
ax_female_word_gross = plt.scatter(df["Gross"], df["total_words_female
    "], color="r", label="Number of Female words to movie gross")
plt.xlabel("Gross")
plt.ylabel("Words")
plt.title("M/F number of words to gross")
plt.legend()


# ----------------- Line plot 1/3 ------------------
# Creates a line plot where the yearly mean of
# "Total words", total_words_female", and
# "total_words_male" are plotted against "Year".
# The aim is to discern if eqality has
# improved over time.
```

```python
55
56  df = film_data_dummies.copy()
57
58  #average nr spoken words (M/F) per film per year
59  df_avg_male_words = df.groupby("Year")["total_words_male"].mean()
60  df_avg_female_words = df.groupby("Year")["total_words_female"].mean()
61  df_avg_total_words= df.groupby("Year")["Total words"].mean()
62  plt.legend()
63
64  plt.subplot(2,2,2)
65  df_avg_total_words_plot = df_avg_total_words.plot(kind="line", y="
        Total words", x="Year", color="g", label="Avg number of words per
         film")
66  df_avg_male_words_plot = df_avg_male_words.plot(kind="line", y="
        total_words_male", x="Year", color="b", label="Avg number of Male
        words per film", ax=df_avg_total_words_plot)
67  df_avg_female_words_plot = df_avg_female_words.plot(kind="line", y="
        total_words_female", x="Year", color ="r",label="Avg number of
        Female words per film", xlabel="Year", ylabel="Avg spoken words",
        title="Avg number of spoken of words per year", ax=
        df_avg_male_words_plot)
68  plt.legend()
69
70  # ----------------- Line plot 2/3 -----------------
71  # Creates a line plot where yearly mean of
72  # proportion of "total_words_female" and
73  # "total_words_male" to "Total words" is
74  # plotted against "Year".
75  # The aim is to discern if eqality has
76  # improved over time.
77
78  # Proportion of spoken words (M/F) per film per year
79  def proportion(x, y):
80      x_new = x/(y)
81      return x_new
82
83  df['Proportion female words'] = df.apply(lambda row : proportion(row['
        total_words_female'],row['Total words']), axis = 1)
84  df['Proportion male words'] = df.apply(lambda row : proportion(row['
        total_words_male'], row['Total words']), axis = 1)
85
86  df_male_words_proportion = df.groupby("Year")["Proportion male words"
        ].mean()
87  df_female_words_proportion = df.groupby("Year")["Proportion female
        words" ].mean()
88
89  plt.subplot(2,2,3)
90  df_male_words_proportion_plot = df_male_words_proportion.plot(kind="
        line", y="Proportion male words", x="Year", color="b", label="Mean
         proportion of Male words per year")
91  df_female_words_proportion_plot = df_female_words_proportion.plot(kind
        ="line", y="Proportion female words", x="Year", color ="r", \
92       label="Mean proportion of Female words per year", xlabel="Year",
        ylabel="Proportion spoken words", title="Proportion of words
        spoken (M/F) per year", \
93          ax=df_male_words_proportion_plot)
94  plt.legend()
95
96
97  # ----------------- Line plot 3/3 -----------------
98  # Creates a line plot where yearly mean of
99  # "Number of male actors" and
100 # "Number of female actors" is plotted against "Year".
101 # The aim is to discern if eqality has
102 # improved over time.
```

```
103
104 df_male_actors = df.groupby("Year", as_index=True)["Number of male
        actors"].mean()
105 df_female_actors = df.groupby("Year", as_index=True)["Number of female
         actors"].mean()
106
107 df_male_actors.columns = ["Year", "Number of male actors"]
108 df_female_actors.columns = ["Year", "Number of female actors"]
109
110 plt.subplot(2,2,4)
111 df2 = df_male_actors.plot(kind="line", y="Number of male actors",x="
        Year", color="b", label="Number of Male actors")
112 df3 = df_female_actors.plot(kind="line", y="Number of female actors",x
        ="Year", color="r", label="Number of Male actors", ax=df2, xlabel=
        "Year", ylabel="Actors", title="Mean of M/F actors per film (
        yearly avg)",figsize=(15,10))
113
114
115 plt.legend()
116 plt.show()
```

Listing 1: Code for equality analysis

## C.2 Discriminant Analysis

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

import sklearn.preprocessing as skl_pre
import sklearn.linear_model as skl_lm
import sklearn.discriminant_analysis as skl_da
import sklearn.neighbors as skl_nb


url = 'train.csv'
indata = pd.read_csv(url, na_values='?', dtype={'ID': str}).dropna().
    reset_index()

allparams = ['Number words female', 'Total words', 'Number of words
    lead', 'Difference in words lead and co-lead', 'Number of male
    actors', 'Year', 'Number of female actors', 'Number words male', '
    Gross', 'Mean Age Male', 'Mean Age Female', 'Age Lead', 'Age Co-
    Lead']
optparams = ['Number words female', 'Total words', 'Number of words
    lead', 'Difference in words lead and co-lead', 'Number of male
    actors', 'Number of female actors', 'Age Lead', 'Age Co-Lead']

N = 200          # Number of random data sets tested
seeds = []       # For np.random

for i in range(0, N):
    seeds.append(int(np.random.random()*1000))   # Random seed 0-999

LDA_Accuracy = []
QDA_Accuracy = []

for s in seeds:

    np.random.seed(s)
    trainI = np.random.choice(indata.shape[0], size=300, replace=False
    )
    trainIndex = indata.index.isin(trainI)
    train = indata.iloc[trainIndex]
    test = indata.iloc[~trainIndex]

    X_train = train[optparams]
    Y_train = train['Lead']
    X_test = test[optparams]
    Y_test = test['Lead']


    # ---------  LDA  ---------

    model = skl_da.LinearDiscriminantAnalysis()
    model.fit(X_train, Y_train)

    predict_prob_L = model.predict_proba(X_test)

    prediction_L = np.empty(len(X_test), dtype=object)
    prediction_L = np.where(predict_prob_L[:, 0]>=0.5, 'Female', 'Male
    ')

    # Accuracy
    LDA_Accuracy.append(np.mean(prediction_L == Y_test))


    # ---------  QDA  ---------
```

```
55
56    model = skl_da.QuadraticDiscriminantAnalysis()
57    model.fit(X_train, Y_train)
58
59    predict_prob_Q = model.predict_proba(X_test)
60
61    prediction_Q = np.empty(len(X_test), dtype=object)
62    prediction_Q = np.where(predict_prob_Q[:, 0]>=0.5, 'Female', 'Male
      ')
63
64    # Accuracy
65    QDA_Accuracy.append(np.mean(prediction_Q == Y_test))
66
67 # Accuracy results and sample confusion matrix:
68
69 print(f"LDA Accuracy: {np.mean(LDA_Accuracy):.3f}")
70 print(f"QDA Accuracy: {np.mean(QDA_Accuracy):.3f} \n")
71
72 # Confusion Matrix LDA
73 print('LDA Confusion Matrix:')
74 print(pd.crosstab(prediction_L, Y_test), '\n')
75
76 # Confusion Matrix QDA
77 print('QDA Confusion Matrix:')
78 print(pd.crosstab(prediction_Q, Y_test))
```
Listing 2: Code for LDA & QDA

## C.3    Decision-tree based method

```python
#!/usr/bin/env python
# coding: utf-8

# In[1]:


import pandas as pd
import numpy as np
import matplotlib
import matplotlib.pyplot as plt

from sklearn import tree
from sklearn.ensemble import BaggingClassifier, RandomForestClassifier
from sklearn.model_selection import cross_val_score
from sklearn.metrics import roc_auc_score
import graphviz


# In[19]:


Oscar = pd.read_csv('train.csv')
np.random.seed(1)
Oscar_index=np.random.choice(Oscar.shape[0],size=800,replace=False)
train_Oscar=Oscar.iloc[Oscar_index]
test_Oscar=Oscar.drop(Oscar_index)
x_train_Oscar=train_Oscar.drop(columns=['Lead'])
y_train_Oscar=train_Oscar['Lead']
y_test_Oscar=test_Oscar['Lead']


# In[3]:


def create_right_train(x_train_Oscar,test_Oscar,name_of_new_x):
    #removes the parameters that are not used in the new training and
    test set
    new_x_train_Oscar=x_train_Oscar[name_of_new_x]
    new_x_test_Oscar=test_Oscar[name_of_new_x]
    return new_x_train_Oscar,new_x_test_Oscar


# In[4]:


new_one_train,new_one_test=create_right_train(x_train_Oscar,test_Oscar
    ,["Number words female","Total words"])
new_one_test


# In[5]:


def create_model_samples(x_train_Oscar,y_train_Oscar,depth,min_samples
    ):
    #Create the decision-tree with specified sample rate
    model=tree.DecisionTreeClassifier(max_depth=depth,
    min_samples_split=min_samples)
    model.fit(X=x_train_Oscar,y=y_train_Oscar)
    return model


# In[6]:
```

```
60
61
62  def create_model(x_train_Oscar,y_train_Oscar,depth):
63      #Create the decision-tree
64      model=tree.DecisionTreeClassifier(max_depth=depth)
65      model.fit(X=x_train_Oscar,y=y_train_Oscar)
66      return model
67
68
69  # In[7]:
70
71
72  def make_graph(model,x_train_Oscar):
73      #Make a graph of the decision-tree
74      dot_data=tree.export_graphviz(model,out_file=None,feature_names=
        x_train_Oscar.columns,
75                                    class_names=model.classes_,filled=True,
        rounded=True,
76                                    leaves_parallel=True,proportion=True)
77      graph=graphviz.Source(dot_data)
78      return graph
79
80
81  # In[8]:
82
83
84  def test_graph(test_Oscar,y_test_Oscar,model):
85      #Test the decision-tree on the test set
86      #x_test_Oscar=test_Oscar.drop(columns=['Lead'])
87      #y_test_Oscar=test_Oscar['Lead']
88
89      y_predict=model.predict(test_Oscar)
90      true_male=0
91      false_male=0
92      true_female=0
93      false_female=0
94      y_tests=y_test_Oscar.tolist()
95      for nr in range(0,len(y_predict)):
96          if y_predict[nr]==y_tests[nr]:
97              if y_predict[nr]=="Male":
98                  true_male+=1
99              else:
100                 true_female+=1
101         else:
102             if y_predict[nr]=="Female":
103                 false_male+=1
104             else:
105                 false_female+=1
106
107     err=np.mean(y_predict != y_test_Oscar)
108     #print('Error rate for tree: '+ str(err))
109     #print('Accuracy rate is %.2f' % np.mean(y_predict==y_test_Oscar))
110     #print("TM "+str(true_male)+" FM: "+str(false_male)+" TF: "+ str(
        true_female)+" FF: "+str(false_female))
111     return err,true_male,false_male,true_female,false_female
112
113
114 # In[9]:
115
116
117 def cross_val(new_x_train,depths,min_sample):
118
119     cv=40
120     scoring='accuracy'
121     cross_scores=[]
```

19

```
122     cross_std=[]
123     cross_mean=[]
124     accuracy=[]
125     for depth in depths:
126         model=create_model_samples(new_x_train,y_train_Oscar,depth,
    min_sample)
127         #err=test_graph(new_x_test,y_test_Oscar,model)
128         cross_score=cross_val_score(model,new_x_train,y_train_Oscar,cv
    =cv,scoring=scoring)
129         cross_scores.append(cross_score)
130         cross_mean.append(cross_score.mean())
131         cross_std.append(cross_score.std())
132         accuracy.append(model.fit(new_x_train,y_train_Oscar).score(
    new_x_train,y_train_Oscar))
133     cross_mean=np.array(cross_mean)
134     cross_std= np.array(cross_std)
135     accuracy=np.array(accuracy)
136     return cross_mean,cross_std,accuracy
137
138
139 # In[10]:
140
141
142 def cross_val_samples(new_x_train,depths,min_samples):
143
144     cv=40
145     scoring='accuracy'
146     cross_scores=[]
147     cross_std=[]
148     cross_mean=[]
149     accuracy=[]
150     for sample in min_samples:
151         model=create_model_samples(new_x_train,y_train_Oscar,depth,
    sample)
152         #err=test_graph(new_x_test,y_test_Oscar,model)
153         cross_score=cross_val_score(model,new_x_train,y_train_Oscar,cv
    =cv,scoring=scoring)
154         cross_scores.append(cross_score)
155         cross_mean.append(cross_score.mean())
156         cross_std.append(cross_score.std())
157         accuracy.append(model.fit(new_x_train,y_train_Oscar).score(
    new_x_train,y_train_Oscar))
158     cross_mean=np.array(cross_mean)
159     cross_std= np.array(cross_std)
160     accuracy=np.array(accuracy)
161     return cross_mean,cross_std,accuracy
162
163
164 # In[37]:
165
166
167 #Show the gini importance of the
168 names_of_char=["Number words female","Total words","Number of words
    lead","Difference in words lead and co-lead",
169             "Number of male actors","Year","Number of female actors"
    ,
170             "Number words male","Gross","Mean Age Male","Mean Age
    Female","Age Lead","Age Co-Lead"]
171 new_x_train,new_x_test=create_right_train(x_train_Oscar,test_Oscar,
    names_of_char)
172 depth=9
173 model=create_model_samples(new_x_train,y_train_Oscar,depth,15)
174 for l in range(1,101):
175     gini+=model.feature_importances_
176 gini=gini/100
```

```
177  plt.bar([i for i in range(len(gini))],gini)
178  plt.xlabel("features")
179  plt.ylabel("Gini importance")
180  plt.title("Gini importance of the features")
181  plt.savefig("Gini_importance.png")
182  #plt.bar(names_of_char,gini)
183  print(str(names_of_char[0])+", "+str(names_of_char[6])+ " and "+str(
         names_of_char[12]))
184  print(str(names_of_char[1])+", "+str(names_of_char[2])+ " and "+str(
         names_of_char[5]))
185  print(sum(gini))
186
187
188  # In[17]:
189
190
191  #Find the accuracy, misclassification and true/false male/female rate
192  names_of_char=["Number words female","Total words","Number of words
         lead","Difference in words lead and co-lead",
193              "Number of male actors","Year","Number of female actors"
         ,
194              "Number words male","Gross","Mean Age Male","Mean Age
         Female","Age Lead","Age Co-Lead"]
195  new_x_train,new_x_test=create_right_train(x_train_Oscar,test_Oscar,
         names_of_char)
196  depth=9
197  err_all=[]
198  TMA=[]
199  FMA=[]
200  TFA=[]
201  FFA=[]
202  for n in range(1,101):
203      model=create_model_samples(new_x_train,y_train_Oscar,depth,15)
204      err,TM,FM,TF,FF=test_graph(new_x_test,y_test_Oscar,model)
205      err_all.append(err)
206      TMA.append(TM)
207      FMA.append(FM)
208      TFA.append(TF)
209      FFA.append(FF)
210  TMS=sum(TMA)/len(TMA)
211  FMS=sum(FMA)/len(FMA)
212  TFS=sum(TFA)/len(TFA)
213  FFS=sum(FFA)/len(FFA)
214  print("Accuracy :"+str(1-sum(err_all)/len(err_all)))
215  print("False male: "+str(FMS/(FMS+TFS)))
216  print("True male: "+str(TMS/(TMS+FFS)))
217  print("True female: "+ str(TFS/(TFS+FMS)))
218  print("False Female: "+ str(FFS/(FFS+TMS)))
219  print("Misclassification : "+ str((FFS+FMS)/(FFS+FMS+TFS+TMS)))
220
221
222  # In[12]:
223
224
225  #Find the accuracy of the decision-tree
226  names_list=[0,1,2,3,4,5,6,7,8,9,10,11,12]
227  #top 0
228  #second row 6
229  #third row 3 4
230  #fourth row 12 6 0 3
231  #fifth row 11  10
232  #names_list=[0,3,4,6,12]
233  #names_list=[0,6,7,10]
234  names_of_char=["Number words female","Total words","Number of words
         lead","Difference in words lead and co-lead",
```

```
235                "Number of male actors","Year","Number of female actors"
      ,
236                "Number words male","Gross","Mean Age Male","Mean Age
      Female","Age Lead","Age Co-Lead"]
237 name_of_new_x=[]
238 depth=9
239 min_sample=12
240 accuracies=[]
241 for nr in names_list:
242     name_of_new_x.append(names_of_char[nr])
243 new_x_train,new_x_test=create_right_train(x_train_Oscar,test_Oscar,
      name_of_new_x)
244 #cross_mean,cross_std,accuracy=cross_val(new_x_train,depth,min_sample)
245 for n in range(1,51):
246     model=create_model_samples(new_x_train,y_train_Oscar,depth,
      min_sample)
247     err=test_graph(new_x_test,y_test_Oscar,model)
248 #print(err)
249     accuracy=model.fit(new_x_train,y_train_Oscar).score(new_x_test,
      y_test_Oscar)
250     accuracies.append(accuracy)
251 print(max(accuracies))
252 print(sum(accuracies)/50)
253 print(min(accuracies))
254
255
256 # In[13]:
257
258
259 #Plot the model
260 dot_data=tree.export_graphviz(model,out_file=None,feature_names=
      new_x_train.columns,
261                             class_names=model.classes_,filled=True,
      rounded=True,
262                             leaves_parallel=True,proportion=True)
263 graph=graphviz.Source(dot_data)
264 graph
265
266
267 # In[14]:
268
269
270 #plot the
271 names_list=[0,1,2,3,4,5,6,7,8,9,10,11,12] #0.8225 #0.81125
272 names_list_2=[0,3,4,6,12] #0.81875
273 #names_list_2=[1,2,5,7,8,9,10,11] #0.775
274 #names_list_2=[0,3,4,6,10,12] #0.81500
275 #names_list_2=[0,3,4,6,11,12] #0.815
276 #names_list_2=[0,3,4,6,7,10,11,12] #0.82875
277 #names_list_2=[0,3,4,6,7,10] #0.8024
278 #names_list_2=[0,3,4,6,12] #0.821
279 names_of_char=["Number words female","Total words","Number of words
      lead","Difference in words lead and co-lead",
280                "Number of male actors","Year","Number of female actors"
      ,
281                "Number words male","Gross","Mean Age Male","Mean Age
      Female","Age Lead","Age Co-Lead"]
282 name_of_new_x_1=[]
283 name_of_new_x_2=[]
284 depths=[2,3,4,5,6,7,8,9,10,11,12,13,14,15]
285 min_samples=12
286 for nr in names_list:
287     name_of_new_x_1.append(names_of_char[nr])
288     if nr in names_list_2:
289         name_of_new_x_2.append(names_of_char[nr])
```

```
290 new_x_train_1 , new_x_test_1 = create_right_train ( x_train_Oscar , test_Oscar
        , name_of_new_x_1 )
291 new_x_train_2 , new_x_test_2 = create_right_train ( x_train_Oscar , test_Oscar
        , name_of_new_x_2 )
292 cross_mean_1 , cross_std_1 , accuracy_1 = cross_val ( new_x_train_1 , depths ,
        min_samples )
293 cross_mean_2 , cross_std_2 , accuracy_2 = cross_val ( new_x_train_2 , depths ,
        min_samples )
294 fig , ax = plt . subplots (1 ,1 , figsize =(15 ,5) )
295 ax . plot ( depths , cross_mean_1 , color ='red ')
296 ax . plot ( depths , cross_mean_2 )
297 ax . fill_between ( depths , cross_mean_1 -2* cross_std_1 , cross_mean_1 +2*
        cross_std_1 , alpha =0.2 , color ='red ')
298 id_max_1 = cross_mean_1 . argmax ()
299 id_max_2 = cross_mean_2 . argmax ()
300 best_depth_1 = depths [ id_max_1 ]
301 best_depth_2 = depths [ id_max_2 ]
302 best_score_1 = cross_mean_1 [ id_max_1 ]
303 best_score_2 = cross_mean_2 [ id_max_2 ]
304 print (" NR 1 the depth : "+ str ( best_depth_1 )+" with score : "+ str (
        best_score_1 ))
305 print (" NR 2 the depth : "+ str ( best_depth_2 )+" with score : "+ str (
        best_score_2 ))
306
307
308 # In [15]:
309
310
311 # Determine depth
312 names_of_char =[" Number words female "," Total words "," Number of words
        lead "," Difference in words lead and co - lead ",
313                " Number of male actors "," Year "," Number of female actors "
        ,
314                " Number words male "," Gross "," Mean Age Male "," Mean Age
        Female "," Age Lead "," Age Co - Lead "]
315 depths =[2 ,3 ,4 ,5 ,6 ,7 ,8 ,9 ,10 ,11 ,12 ,13 ,14 ,15]
316 min_samples =9
317 times =10
318 best_depths =[]
319 new_x_train , new_x_test = create_right_train ( x_train_Oscar , test_Oscar ,
        names_of_char )
320 for n in range (1 , times +1) :
321     cross_mean_1 , cross_std_1 , accuracy_1 = cross_val ( new_x_train , depths ,
        min_samples )
322     id_max_1 = cross_mean_1 . argmax ()
323     best_depth_1 = depths [ id_max_1 ]
324     best_depths . append ( best_depth_1 )
325 print (" Max depth "+ str ( max ( best_depths ))+" Min depth "+ str ( min (
        best_depths )))
326 print (" Mean "+ str ( sum ( best_depths )/10) )
327
328
329 # In [16]:
330
331
332 # Determine min samples
333 name_of_new_x =[" Number words female "," Total words "," Number of words
        lead "," Difference in words lead and co - lead ",
334                " Number of male actors "," Year "," Number of female actors "
        ,
335                " Number words male "," Gross "," Mean Age Male "," Mean Age
        Female "," Age Lead "," Age Co - Lead "]
336 depths =9
337 min_samples = list ( range (2 ,40) )
338 times =10
```

```
339 best_samples=[]
340 new_x_train,new_x_test=create_right_train(x_train_Oscar,test_Oscar,
        name_of_new_x)
341 for n in range(1,times+1):
342     print(n)
343     cross_mean,cross_std,accuracy=cross_val_samples(new_x_train,depths
        ,min_samples)
344     id_max=cross_mean.argmax()
345     best_sample=min_samples[id_max]
346     best_score_1=cross_mean[id_max]
347     best_samples.append(best_sample)
348
349 print("Max sample size: "+str(max(best_samples))+" Min sample size: "+
        str(min(best_samples)))
350 print(sum(best_samples)/10)
351
352
353 # In[ ]:
```

Listing 3: Code for tree-based methods

### C.4 Logistic regression and feature selection

```python
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4
5  import sklearn.linear_model as skl_lm
6  import sklearn.discriminant_analysis as skl_da
7  import sklearn.neighbors as skl_nb
8  import sklearn.preprocessing as skl_pre
9
10 from google.colab import files
11 uploaded = files.upload()
12
13 import io
14 data = pd.read_csv(io.BytesIO(uploaded['train.csv']), na_values = '?')
      .dropna().reset_index(drop=True)
15 # Dataset is now stored in a Pandas Dataframe
16
17 data['Lead'] = np.where(data['Lead'] == 'Female', 1, 0)
18
19 # Split the data randomly into a training set and a test set of
      approximately similar size.
20
21 np.random.seed(2)
22 trainI = np.random.choice(data.shape[0], size= 700, replace=False)
23 trainIndex = data.index.isin(trainI)
24
25 train = data.iloc[trainIndex]
26 test = data.iloc[~trainIndex]
27
28 # Try different feature sets (with data scaling)
29
30 import itertools
31 def findsubsets(s, n):
32     return list(itertools.combinations(s, n))
33
34 from sklearn.metrics import roc_auc_score
35 from sklearn.metrics import auc
36 from sklearn.metrics import precision_recall_curve
37
38 min_error = data.shape[0]
39 error_a = 0
40 features_max = []
41 coef =0
42 conf_matrix =0
43 chosen_model =0
44 prob = 0
45
46 pc_auc = 0
47 roc_auc = 0
48 f1 = 0
49 roc_auc_a = 0
50 pc_auc_a = 0
51 pc_auc_max = 0
52 roc_auc_max = 0
53 f1_a = 0
54
55 s = {'Number words female', 'Total words', 'Number of words lead', '
      Difference in words lead and co-lead', 'Number of male actors', '
      Year', 'Number of female actors', 'Number words male', 'Gross', '
      Mean Age Male', 'Mean Age Female','Age Lead', 'Age Co-Lead'}
56 n = 13
57
58 # define model
```

```python
59  model = skl_lm.LogisticRegression(solver ='liblinear')
60
61  subsets = findsubsets(s, n)
62  for s in subsets:
63    input_variables = list(s)
64
65    x_train = train[input_variables]
66    y_train = train['Lead']
67    x_test = test[input_variables]
68    y_test = test['Lead']
69
70    scaler = skl_pre.StandardScaler().fit(x_train)
71    model.fit(scaler.transform(x_train), y_train )
72    predict_prob = model.predict_proba(scaler.transform(x_test))
73    prediction = np.where(predict_prob[:,1] > 0.5, 1, 0)
74
75    # misclassification error
76    error = np.mean(prediction != y_test)
77
78    # calculate roc auc
79    roc_auc = roc_auc_score(y_test, predict_prob[:,1])
80
81    # calculate the precision-recall auc
82    precision, recall, _ = precision_recall_curve(y_test, predict_prob
       [:,1])
83    pc_auc = auc(recall, precision)
84    f1 = f1 = 2*precision*recall/(precision+recall)
85
86    # chose the model with the max ROC AUC
87    if roc_auc > roc_auc_max:
88      roc_auc_max = roc_auc
89      error_a = error
90      features_max = input_variables
91      chosen_model = model
92      conf_matrix = pd.crosstab(prediction, y_test)
93      prob = predict_prob
94      pc_auc_a = pc_auc
95      f1_a = f1
96
97  print(features_max)
98  print('Misclassification error',error_a)
99  print(f'ROC AUC = {roc_auc_max}')
100 print('PC AUC =', pc_auc_a)
101 print('F1 =', np.max(f1_a))
102 print(chosen_model.coef_)
103 print(chosen_model.classes_)
104 print(conf_matrix)
105
106 # Find best features subsets (with 8,9,10 features)
107 X = data
108 y = data['Lead']
109
110 import sklearn.model_selection as skl_ms
111 import sklearn.preprocessing as skl_pre
112 n_runs = 20
113 features_lsts_with_big_roc_auc =[]
114 s = {'Number words female', 'Total words', 'Number of words lead', '
       Difference in words lead and co-lead', 'Number of male actors', '
       Year', 'Number of female actors', 'Number words male', 'Gross', '
       Mean Age Male', 'Mean Age Female','Age Lead', 'Age Co-Lead'}
115
116 subsets_1 = findsubsets(s, 8)
117 subsets_2 = findsubsets(s, 9)
118 subsets_3 = findsubsets(s, 10)
119
```

```python
120 subsets = subsets_1 + subsets_2 + subsets_3
121 roc_auc_n = np.zeros((n_runs, len(subsets)))
122 pc_auc = 0
123 error = 0
124 pc_auc = 0
125 f1 = 0
126
127 for i in range(n_runs):
128   X_train, X_val, y_train, y_val = skl_ms.train_test_split(X, y,
      test_size = 0.3)
129
130   for j, s in enumerate(subsets):
131     input_variables = list(s)
132
133     x_train = X_train[input_variables]
134     x_test = X_val[input_variables]
135     y_test = y_val
136
137     model = skl_lm.LogisticRegression(solver ='liblinear')
138     scaler = skl_pre.StandardScaler().fit(x_train)
139
140     model.fit(scaler.transform(x_train), y_train)
141     predict_prob = model.predict_proba(scaler.transform(x_test))
142     prediction = np.where(predict_prob[:,1] > 0.5, 1, 0)
143
144     roc_auc = roc_auc_score(y_test, predict_prob[:,1])
145     roc_auc_n[i,j] = roc_auc
146
147 roc_auc_avg = np.mean(roc_auc_n, axis = 0)
148
149 for idx ,el in enumerate(roc_auc_avg):
150   if el > np.max(roc_auc_avg)*(1-0.001):
151     features_lsts_with_big_roc_auc.append(list(subsets[idx]))
152
153 for el in features_lsts_with_big_roc_auc:
154   print(el)
155   print('Number of features',len(el))
156
157 # find the decision threshold rate
158 def find_r(prob, y_test):
159   recall = []
160   precision = []
161   f1_max = 0
162   r_f1 = 0
163
164   positive_class = 1
165   negative_class = 0
166
167   P = np.sum(y_test == positive_class)
168   prediction = np.empty(len(x_test), dtype ='object')
169   tr_val = np.linspace(0.00 , 1, num =101)
170
171   for r in tr_val:
172     prediction = np.where(prob[:,1] > r, positive_class,
      negative_class)
173     P_star = np.sum(prediction == positive_class)
174     tr_pos = np.sum((prediction == y_test)&(prediction ==
      positive_class))
175
176     rec = tr_pos/P
177     prec = tr_pos/P_star
178
179     f1 = 2*prec*rec/(prec+rec)
180     if f1_max < f1:
181       f1_max =f1
```

```
182        r_f1 = r
183
184    return r_f1
185
186 # cross-validation of the threshold rate
187 def calculate_average_r(features, X, y):
188 # features -list of features
189 # X - data with all features
190 # y - column 'Lead'
191    n_runs = 10
192
193    tr_rates_n = np.zeros((n_runs))
194    for i in range(n_runs):
195
196        X_train, X_val, y_train, y_val = skl_ms.train_test_split(X, y,
       test_size = 0.3)
197        model = skl_lm.LogisticRegression(solver ='liblinear')
198        scaler = skl_pre.StandardScaler().fit(X_train[features])
199        model.fit(scaler.transform(X_train[features]), y_train)
200        predict_prob = model.predict_proba(scaler.transform(X_val[features
       ]))
201
202        tr_rates_n[i] = find_r(predict_prob, y_val)
203
204    return np.average(tr_rates_n)
205
206 # cross-validation of models with different feature sets with k-fold
207 features_list = features_lsts_with_big_roc_auc
208
209 n_runs = 10
210
211 pc_auc_n =np.zeros((n_runs, len(features_list)))
212 roc_auc_n = np.zeros((n_runs, len(features_list)))
213 misclassification_n = np.zeros((n_runs, len(features_list)))
214 tr_rate_models = []
215
216 X = data.drop(columns = 'Lead')
217 y = data['Lead']
218
219 n_fold = 10
220 cv = skl_ms.KFold(n_splits = n_fold, random_state = 1, shuffle = True)
221
222 for i, (train_index, val_index) in enumerate(cv.split(X)):
223    X_train, X_val = X.iloc[train_index], X.iloc[val_index]
224    y_train, y_val = y.iloc[train_index], y.iloc[val_index]
225
226    for j,features in enumerate(features_list):
227        model = skl_lm.LogisticRegression(solver ='liblinear')
228        scaler = skl_pre.StandardScaler().fit(X_train[features])
229        model.fit(scaler.transform(X_train[features]), y_train)
230        predict_prob = model.predict_proba(scaler.transform(X_val[features
       ]))
231
232        precision, recall, _ = precision_recall_curve(y_val, predict_prob
       [:,1])
233        pc_auc_n[i,j] = auc(recall, precision)
234        roc_auc_n[i,j] = roc_auc_score(y_val, predict_prob[:,1])
235
236        tr_rate = calculate_average_r(features, X, y)
237        prediction = np.where(predict_prob[:,1] > tr_rate, 1, 0)
238        misclassification_n[i,j] = (np.mean(prediction != y_val))
239        tr_rate_models.append(tr_rate)
240
241 fig, axs = plt.subplots(1, 3, figsize=(12, 3), sharey=False)
242 axs[0].boxplot(roc_auc_n)
```

```
243  axs[1].boxplot(pc_auc_n)
244  axs[2].boxplot(misclassification_n)
245
246  fig.suptitle('Cross validation for models with different features sets
        ')
247  plt.show()
248
249  # Model (logistic regression) validation: find regularization type and
        rate
250  from scipy.stats import loguniform
251  from sklearn.linear_model import LogisticRegression
252  from sklearn.model_selection import RepeatedStratifiedKFold
253  from sklearn.model_selection import RandomizedSearchCV
254
255  def validate_model(X,y):
256    model = LogisticRegression()
257    cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state
        =1)
258    space = dict()
259    space['solver'] = ['liblinear']
260    space['penalty'] = ['none', 'l1', 'l2', 'elasticnet']
261    space['C'] = loguniform(1e-5, 100)
262
263    search = RandomizedSearchCV(model, space, n_iter=500, scoring='
        roc_auc', n_jobs=-1, cv=cv, random_state=1)
264    res = search.fit(X, y)
265
266    print('Score: %s' % res.best_score_)
267    print('Hyperparameters: %s' % res.best_params_)
268    return res.best_params_
269
270  # insert input variables (list of features) to validate the model
271  features = ['Total words', 'Age Co-Lead', 'Number of female actors', '
        Difference in words lead and co-lead', 'Number of words lead', '
        Age Lead', 'Number words female', 'Number of male actors']
272  validate_model(X[features], y)
273
274  # Find threshold rate for the best model (result = 0.435)
275  features = ['Total words', 'Age Co-Lead', 'Number of female actors', '
        Difference in words lead and co-lead', 'Number of words lead', '
        Age Lead', 'Number words female', 'Number of male actors']
276  # parameters {'C': 18.31792545756584, 'penalty': 'l1', 'solver': '
        liblinear'}
277
278  n_runs = 10
279  tr_rates_n = np.zeros((n_runs))
280  X = data[features]
281  y = data['Lead']
282
283  for i in range(n_runs):
284
285    X_train, X_val, y_train, y_val = skl_ms.train_test_split(X, y,
        test_size = 0.3)
286
287    model = skl_lm.LogisticRegression(solver ='liblinear', C =
        18.31792545756584, penalty ='l1')
288    scaler = skl_pre.StandardScaler().fit(X_train)
289    model.fit(scaler.transform(X_train), y_train)
290    predict_prob = model.predict_proba(scaler.transform(X_val))
291    tr_rates_n[i] = find_r(predict_prob, y_val)
292
293  r_res = np.average(tr_rates_n)
294  print(r_res)
295
296  # Evaluate performance of logistic regression model
```

```python
297 from sklearn.metrics import f1_score
298 features = ['Total words', 'Age Co-Lead', 'Number of female actors', '
        Difference in words lead and co-lead', 'Number of words lead', '
        Age Lead', 'Number words female', 'Number of male actors']
299 X = data[features]
300 y = data['Lead']
301 pc_auc_n2 =np.zeros((n_runs))
302 roc_auc_n2 = np.zeros((n_runs))
303 misclassification_n2 = np.zeros((n_runs))
304 true_positive_rate = []
305 false_positive_rate = []
306 f1_rate = []
307 tr_pos = 0
308 fal_pos = 0
309
310 positive_class = 1
311 negative_class = 0
312
313 n_fold = 10
314 cv = skl_ms.KFold(n_splits = n_fold, random_state = 1, shuffle = True)
315
316 for i, (train_index, val_index) in enumerate(cv.split(X)):
317   X_train, X_val = X.iloc[train_index], X.iloc[val_index]
318   y_train, y_val = y.iloc[train_index], y.iloc[val_index]
319
320   P = np.sum(y_val == positive_class)
321   N = np.sum(y_val == negative_class)
322
323   model = skl_lm.LogisticRegression(solver ='liblinear', C =
        18.31792545756584, penalty ='l1')
324
325   scaler = skl_pre.StandardScaler().fit(X_train)
326   model.fit(scaler.transform(X_train[features]), y_train)
327   predict_prob = model.predict_proba(scaler.transform(X_val))
328   # calculate the metrics
329   precision, recall, _ = precision_recall_curve(y_val, predict_prob
        [:,1])
330
331   pc_auc_n2[i] = auc(recall, precision)
332   roc_auc_n2[i] = roc_auc_score(y_val, predict_prob[:,1])
333
334   tr_rate = 0.435
335   prediction = np.where(predict_prob[:,1] > tr_rate, 1, 0)
336
337   misclassification_n2[i] = (np.mean(prediction != y_val))
338   tr_pos = np.sum((prediction == y_val)&(prediction == positive_class)
        )
339   fal_pos = np.sum((prediction != y_val)&(prediction == positive_class
        ))
340   true_positive_rate.append(tr_pos/P)
341   false_positive_rate.append(fal_pos/N)
342
343   f1 = f1_score(y_val, prediction, average='binary')
344   f1_rate.append(f1)
345
346 print('ROC AUC', np.average(roc_auc_n2))
347 print('PC AUC', np.average(pc_auc_n2))
348 print('Error', np.average(misclassification_n2))
349 print('TPR', np.average(true_positive_rate))
350 print('FPR', np.average(false_positive_rate))
351 print('f1 coef', np.average(f1_rate))
352
353 fig, axs = plt.subplots(1, 3, figsize=(12, 3), sharey=False)
354 axs[0].boxplot(roc_auc_n2)
355 axs[1].boxplot(pc_auc_n2)
```

```
356 axs[2].boxplot(misclassification_n2)
357
358 axs[0].set_xticks(range(1))
359 axs[1].set_xticks(range(1))
360 axs[2].set_xticks(range(1))
361
362 fig.suptitle('Evaluation of log regression model (ROC AUC, PC AUC,
        misclassification error)')
363 plt.show()
```

Listing 4: Code for logistic regression and feature selection

## C.5   kNN evaluation code

```python
# Imports
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import sklearn.neighbors as skl_nb
import sklearn.model_selection as skl_ms
import sklearn.metrics as met
import os
from tabulate import tabulate


def main():
    # Load data which has already been normalized
    # using StandardScaler()
    cwd = os.getcwd()
    #URI = (cwd+"\\train.csv")
    URI = (cwd+"\\train_standardscaler.csv")
    #URI = (cwd+"\\train_minmaxscaler.csv")
    film_data = pd.read_csv(URI, dtype={"Lead":str}).dropna().
    reset_index(drop=True)

    x = film_data.drop(columns=['Lead']) #, 'Year', 'Mean Age Male', '
    Mean Age Female', 'Number of male actors', 'Number of female
    actors', 'Age Co-Lead'])
    y = film_data['Lead']
    positive_class = "Male"
    negative_class = "Female"



    # Setting values to k-range, number of folds, and threshold
    k_max = 80
    n_folds = 200
    threshold = 0.24

    # Call of function "evaluate_with_kfold",
    # please see function def for more info
    missclassification_k_error, evaluation_terms, plotting_terms =
    evaluate_with_kfold(x, y,positive_class, negative_class, threshold
    , n_folds, k_max)

    # Finds what k-value minimizes missclassification
    # and find missclassification at that point
    min_error_k = find_best_k(missclassification_k_error)
    new_error_estimate = np.min(missclassification_k_error)
    evaluation_terms["Optimal_k"] = min_error_k
    evaluation_terms["E[Error_new]"] = new_error_estimate

    # Creates list for latex
    data = list(evaluation_terms.items())
    data.insert(0, ["Term","Value"])
    data_np = np.array(data)
    print(tabulate(data_np, headers=("firstrow"), tablefmt="latex"))


    # Create the missclassification to 'k' value
    # graph
    plt.plot(np.arange(1,k_max), missclassification_k_error)
    plt.title(f"Cross validation KFold = {n_folds} error for kNN")
    plt.xlabel("k")
    plt.ylabel("Validation error")
    plt.show()
```

```
59      # Create crosstab heatmap
60      cm=np.array([ [evaluation_terms["TP"],evaluation_terms["FN"]],\
61          [evaluation_terms["FP"],evaluation_terms["TN"]] ])
62      disp = met.ConfusionMatrixDisplay(confusion_matrix=cm,
        display_labels=["Male", "Female"])
63      disp.plot()
64      disp.ax_.set_title("Cross tab heatmap")
65      plt.show()
66
67      # Create ROC and Recall-Precision graph
68      figure, axis = plt.subplots(2)
69      figure.set_constrained_layout(True)
70
71      axis[0].plot(plotting_terms["FPR_curve"], plotting_terms["
        TPR_curve"]);
72      axis[0].set_title("ROC curve ")
73      axis[0].set_xlabel("False positive rate")
74      axis[0].set_ylabel("True positive rate")
75
76      axis[1].plot(plotting_terms["recall_curve"], plotting_terms["
        precision_curve"]);
77      axis[1].set_title("Precision-recall curve")
78      axis[1].set_xlabel("Recall rate")
79      axis[1].set_ylabel("precision rate")
80      plt.show()
81
82  def evaluate_with_kfold(x, y, positive_class, negative_class,
        threshold=0.5, n_folds=None, k_max=None):
83      # ---------------- function description -----------------
84      # This function has three returns:
85      # 'missclassification_k_error', 'evaluation_terms', and
86      # 'plotting_terms'.
87      # These are all dictionaries with 'values' estimated
88      # by cross validation and 'key' as the name of the
89      # variable or list of plotting points.
90      #
91      # 'missclassification_k_error' is a np.array where index+1
92      # represents k and value at index represents
93      # missclassification error for that k
94      #
95      # 'evaluation_terms' is a dictionary containing values
96      # of 'evaluation terms', (ex. TP,FP...) where the name
97      # of the value is the key
98      #
99      # 'plotting_terms' is a dictionary containing np.arrays.
100     # Each np.array has elements representing the value of
101     # an 'evalueation term' at a certain 'r' threshold. The
102     # index represents the 'r' value which ranges from 0 to 1
103     # with increments of 0.01.
104
105
106     if n_folds==None:
107         n_folds=x.shape[0]//4
108     if k_max == None:
109         k_max= x.shape[1]*4
110
111     # Defining the splits and setting random_state to have
        reproductable results
112     Kfold_cv = skl_ms.KFold(shuffle=True,n_splits=n_folds,
        random_state=1)
113     k_range = np.arange(1,k_max)
114
115     # Initializing variables for later usage
116     missclassification_k_error = np.zeros(len(k_range))
117     evaluation_terms = {}
```

```python
118    threshold_terms = {}
119    plotting_terms = {}
120
121    # ---------------- 1st cross validation ----------------
122    # Finds the optimal value for 'k' and an estimation of
123    # E[Error_new] using cross validation with uniform weights.
124    # Number of folds and range of k-values tested are
125    # given by n_folds and k_max
126
127    for train_index, val_index in Kfold_cv.split(x):
128        x_train, x_test = x.iloc[train_index], x.iloc[val_index]
129        y_train, y_test = y.iloc[train_index], y.iloc[val_index]
130
131        temp_missclassification_k_error = evaluate_k_kNN(k_range,
    x_train, y_train, x_test, y_test,positive_class, negative_class,"
    uniform", "auto",threshold)
132        missclassification_k_error = np.add(missclassification_k_error
    ,temp_missclassification_k_error)
133
134    missclassification_k_error /= n_folds
135    min_error_k = find_best_k(missclassification_k_error)
136
137    # ---------------- 2nd cross validation ----------------
138    # Finds values of 'evaluation terms' and 'plotting
139    # ranges(terms)'
140
141    for train_index, val_index in Kfold_cv.split(x):
142        x_train, x_test = x.iloc[train_index], x.iloc[val_index]
143        y_train, y_test = y.iloc[train_index], y.iloc[val_index]
144        model = skl_nb.KNeighborsClassifier(n_neighbors=min_error_k,
    weights="uniform",algorithm="auto")
145        model.fit(x_train, y_train)
146
147        temp_evaluation_terms, temp_threshold_terms =
    get_evaluation_terms(model, x_test, y_test, positive_class,
    negative_class, threshold)
148
149        # The sum of each 'evaluation terms' generated
150        # per fold
151        for key in temp_evaluation_terms.keys():
152            if key in evaluation_terms:
153                evaluation_terms[key] += temp_evaluation_terms[key]
154            else:
155                # If evaluation_terms is empty
156                evaluation_terms[key] = temp_evaluation_terms[key]
157
158        # The sum of each 'plotting term' per index, per fold
159        # This works because all plotting_terms have the same
160        # size
161        for key in temp_threshold_terms.keys():
162            if key in threshold_terms:
163                for i, k in enumerate(temp_threshold_terms[key]):
164                    val = threshold_terms[key][i] +
    temp_threshold_terms[key][i]
165                    threshold_terms[key][i]= val
166            else:
167                # If plotting_terms is empty
168                threshold_terms[key] = temp_threshold_terms[key]
169
170
171
172    # Creation of additional evaluation terms and addition to
173    # dictionary "evaluation_terms"
174    evaluation_terms["TPR"] = evaluation_terms["TP"]/evaluation_terms[
    "P"]
```

34

```
175      evaluation_terms["FPR"] = evaluation_terms["FP"]/evaluation_terms[
         "N"]
176      evaluation_terms["accuracy"] = (evaluation_terms["TP"]+
         evaluation_terms["TN"])/(evaluation_terms["N"]+evaluation_terms["P
         "])
177      evaluation_terms["precision"] = evaluation_terms["TP"]/
         evaluation_terms["P_star"]
178      evaluation_terms["recall"] = evaluation_terms["TP"]/(
         evaluation_terms["TP"]+evaluation_terms["FN"])
179      evaluation_terms["F1"] = 2*(evaluation_terms["precision"]*
         evaluation_terms["TPR"])/(evaluation_terms["precision"]+
         evaluation_terms["TPR"])
180
181
182      # Creation of specific plotting curves
183      plotting_terms["FPR_curve"] = threshold_terms["FP_threshold"]/
         evaluation_terms["N"]
184      plotting_terms["TPR_curve"] = threshold_terms["TP_threshold"]/
         evaluation_terms["P"]
185      plotting_terms["recall_curve"] = threshold_terms["TP_threshold"]/(
         threshold_terms["FN_threshold"]+threshold_terms["TP_threshold"])
186      plotting_terms["precision_curve"] = threshold_terms["TP_threshold"
         ]/(threshold_terms["FP_threshold"]+threshold_terms["TP_threshold"
         ])
187
188
189      return missclassification_k_error, evaluation_terms,
         plotting_terms
190
191
192
193  def evaluate_k_kNN(k_range,x_train, y_train, x_test, y_test,
         positive_class, negative_class, weight_type=None, algorithm_type=
         None, threshold=0.5):
194      # ---------------- function description ----------------
195      # This function returns the missclassification error
196      # of 'k' in range (1-'k_range').
197      # It returns an np.array where index is the value
198      # of ('k'-1) and the value the missclassification error.
199
200      if weight_type == None:
201          weight_type = "uniform"
202
203      if algorithm_type==None:
204          algorithm_type = "auto"
205
206
207      missclassification_k_error = np.zeros(len(k_range))
208      for index, k in enumerate(k_range):
209          model = skl_nb.KNeighborsClassifier(n_neighbors=k,weights=
         weight_type,algorithm=algorithm_type)
210          model.fit(x_train, y_train)
211          missclassification_k_error[index] +=
         get_mean_missclassification(model,x_test,y_test, threshold,
         positive_class, negative_class)
212
213      return missclassification_k_error
214
215
216  def get_mean_missclassification(model,x_test,y_test, threshold,
         positive_class, negative_class):
217      # ---------------- function description ----------------
218      # This function returns the mean missclassification error
219      # from a 'model' evaluated on 'y_test' and with
220      # 'threshold'
```

```python
221
222        # set threshold
223        positive_class_index = np.argwhere(model.classes_== positive_class
       ).squeeze()
224        prediction = np.where(model.predict_proba(x_test)[:,
       positive_class_index] > threshold, positive_class, negative_class)
225
226        # calc missclasification error
227        mean_missclassification =  np.mean(prediction != y_test)
228        return mean_missclassification
229
230
231 def get_evaluation_terms(model, x_test, y_test, positive_class,
       negative_class,threshold):
232        # ---------------- function description ----------------
233        # This function returns the evaluation and plotting terms
234        # given a model, test set and class labels.
235
236        positive_class_index = np.argwhere(model.classes_== positive_class
       ).squeeze()
237
238        # Setting based on 'threshold'
239        prediction = np.where(model.predict_proba(x_test)[:,
       positive_class_index] > threshold, positive_class, negative_class)
240        prediction = model.predict(x_test)
241        predict_prob = model.predict_proba(x_test)
242        P = np.sum(y_test == positive_class) #the same as TP+FN
243        N = np.sum(y_test == negative_class) #the same as TN+FP
244
245        # All variables with *_threshold are lists that contain
246        # plotting values.
247        # Index represents the value of the threshold 'r' in
248        # range 0-1 with 0.01 as increments.
249        FP_threshold = np.zeros(101)
250        TP_threshold = np.zeros(101)
251        FN_threshold = np.zeros(101)
252        TN_threshold = np.zeros(101)
253        threshold_range = np.linspace(0,1,101)
254
255        i=0
256        for r in threshold_range:
257            prediction_curve = np.where(predict_prob[:,
       positive_class_index]> r, positive_class, negative_class)
258            FP_threshold[i] = np.sum((prediction_curve == positive_class)
       & (y_test == negative_class))
259            TP_threshold[i] = np.sum((prediction_curve == positive_class)
       & (y_test == positive_class))
260            FN_threshold[i] = np.sum((prediction_curve == negative_class)
       & (y_test == positive_class))
261            TN_threshold[i] = np.sum((prediction_curve == negative_class)
       & (y_test == negative_class))
262            i +=1
263
264        FP = np.sum((prediction == positive_class) & (y_test ==
       negative_class))
265        TP = np.sum((prediction == positive_class) & (y_test ==
       positive_class))
266        FN = np.sum((prediction == negative_class) & (y_test ==
       positive_class))
267        TN = np.sum((prediction == negative_class) & (y_test ==
       negative_class))
268
269
270        P_star = np.sum(prediction == positive_class) #the same as TP+FP
271        N_star = np.sum(prediction == negative_class) #the same as TN+F
```

```python
272
273     evaluation_terms = {"P":P, "N":N,"P_star":P_star, "N_star":N_star,
         "TN":TN, "FP":FP, "FN":FN, "TP":TP}
274     plotting_terms = {"FP_threshold":FP_threshold,"TP_threshold":
        TP_threshold,"FN_threshold":FN_threshold,"TN_threshold":
        TN_threshold}
275
276     return evaluation_terms, plotting_terms
277
278
279 def find_best_k( missclassification_k_error):
280     # ----------------- function description -----------------
281     # This function returns the index+1 of the minimum valued
282     # element in the missclassification_k_error list.
283     # This equates to the 'k' value for that point
284     #
285     min_error = np.min(missclassification_k_error)
286     min_error_k = [i for i, x in enumerate(missclassification_k_error)
         if x == min_error] [0]+1
287
288     return min_error_k
289
290
291
292 # Functions for creating scaled data sets
293
294 def generate_standard_scaled_datafile(film_data):
295     import sklearn.preprocessing as skl_pre
296     #CREATE NEW DATA WITH STANDARDSCALING
297     x = film_data.drop(columns=['Lead'])
298     y = film_data['Lead']
299
300     standard_scaler = skl_pre.StandardScaler(with_mean=True,with_std=
        True)
301
302     # StandardScaler: mean=0, variance=1
303     scaled_film_data_array = standard_scaler.fit_transform(x)
304
305     x_scaled = pd.DataFrame(scaled_film_data_array, columns = list(x.
        columns))
306
307     film_data_scaled = x_scaled.join(y)
308
309     film_data_scaled.to_csv('train_standardscaler.csv',index=False)
310
311     print("New standardScaled data saved to file: '
        train_standardscaler.csv'")
312
313 def generate_MinMax_scaled_datafile(film_data):
314     import sklearn.preprocessing as skl_pre
315     #CREATE NEW DATA WITH MINMAXSCALER
316     x = film_data.drop(columns=['Lead'])
317     y = film_data['Lead']
318
319     minmax_scaler = skl_pre.MinMaxScaler(feature_range=(0,1))
320
321     # MinMax scaler: min=0, max=1
322     scaled_film_data_array = minmax_scaler.fit_transform(x)
323
324     x_scaled = pd.DataFrame(scaled_film_data_array, columns = list(x.
        columns))
325
326     film_data_scaled = x_scaled.join(y)
327
328     film_data_scaled.to_csv('train_minmaxscaler.csv', index=False)
```

```
329
330      print("New MinMax data saved to file: 'train_minmaxscaler.csv'")
331
332
333
334 if __name__=="__main__":
335      main()
```

Listing 5: Code for evaluating the performance of kNN

## C.6 Quadratic Discriminant Analysis Predictions

```
1  import pandas as pd
2  import numpy as np
3
4  import sklearn.discriminant_analysis as skl_da
5  import csv
6
7  url_train = 'train.csv'
8  train = pd.read_csv(url_train, na_values='?', dtype={'ID': str}).
       dropna().reset_index()
9  url_test = 'test.csv'
10 test = pd.read_csv(url_test, na_values='?', dtype={'ID': str}).dropna
       ().reset_index()
11
12 allparams = ['Number words female', 'Total words', 'Number of words
       lead', 'Difference in words lead and co-lead', 'Number of male
       actors', 'Year', 'Number of female actors', 'Number words male', '
       Gross', 'Mean Age Male', 'Mean Age Female', 'Age Lead', 'Age Co-
       Lead']
13 optparams = ['Number words female', 'Total words', 'Number of words
       lead', 'Difference in words lead and co-lead', 'Number of male
       actors', 'Number of female actors', 'Age Lead', 'Age Co-Lead']
14
15
16 X_train = train[optparams]
17 Y_train = train['Lead']
18 X_test = test[optparams]
19 #Y_test = test['Lead']
20
21
22 # ---------  QDA  ---------
23
24 model = skl_da.QuadraticDiscriminantAnalysis()
25 model.fit(X_train, Y_train)
26
27 predict_prob_Q = model.predict_proba(X_test)
28
29 prediction_Q = np.empty(len(X_test), dtype=object)
30 prediction_Q = np.where(predict_prob_Q[:, 0]>=0.5, '1', '0')
31
32 np.savetxt("predictions.csv", prediction_Q, newline=',', fmt='%s')
```

Listing 6: Code for the predictions, QDA