
Do (wo)men talk too much in films?

Authors: Jannes van Poppelen, Linus Falk, Steven Bijl, Niklas Kostrzewa

Abstract

In this work, we train different machine learning models using a data set of 2000 Hollywood movies in order to accurately predict the sex of the lead role of a movie, given a list of features corresponding to each movie. Even though the training data is heavily biased towards a male lead, all trained models end up outperforming a naive classifier that merely predicts a male lead. The generative QDA model was found to give the best performance after a 10-fold cross-validation, as out of all models it has the best mean prediction accuracy of 90% and a near-vanishing variance. Since the selected input features for this model were only briefly touched, a more thorough optimization for the input features could possibly yield an even more accurate model.

1 Introduction

A study of 2000 movies has recently shown that white men tend to dominate movie roles [1]. Attempting to confirm or deny this finding, multiple machine learning models will be trained using a relatively small data set of Hollywood movies, with the goal to predict the gender of the main actor in Hollywood screenplays. An additional focus will be put on analyzing gender bias for the lead- and co-lead actors. The provided data set has a wide variety of Hollywood movies, ranging over multiple decades and genres. Each entry consists of a total of 14 attributes. The 13 input parameters are the number of actors per gender, gross profit, the total number of words spoken, the number of words spoken per gender, the number of words spoken by lead, the difference in the number of words by lead and co-lead, lead- and co-lead age, and the mean age of male and female characters. The last data entry is the gender of the lead actor, which is the output that should be predicted.

2 Data Analysis

In order to get a better understanding of the gender bias within the data we can ask ourselves *which sex dominates speaking roles in Hollywood movies*. By summing and comparing the number of male- and female actors for all movies in the data set we can clearly see that there is male dominance in speaking roles, as is shown in figure 1. While this is a nice static observation, it makes more sense to ask the question in a dynamic context, by studying how this dominance has changed over time. To do this we start by calculating the percentage of female actors in each film year by year. Then, using NumPy's polyfit function, a straight line is fitted to the data. The resulting line shows a positive trend, indicating that the percentage of female actors per movie is increasing, but is still not close to 50 %. Next, we can examine economic differences between male- and female-dominated movies, where the most straightforward would be to analyze *if male-dominated movies generate more gross than female-dominated movies*. By sorting out all movies by the gender with the most words spoken and calculating the mean gross we can see that indeed male-dominated movies generate more gross, which is seen in the lower plot in figure 1. In order to get a bigger look at the data set and locate outliers, a histogram and pairwise scatter plot of the normalized features are made. These can also indicate obvious correlations between the features. The plot is found in the appendix in figure A.3, which does confirm the existence of some outliers. Since all features are numerical, there is no need for 'dummy' variables. Most features do not follow a normal distribution, which makes it difficult to exclude the outliers without transforming the data. Another way to analyze the

data would be to compute the covariance matrix, by calculating the covariance between all features $\text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y))$ [2], however, since the majority of the features does not follow a normal distribution, this would again be of limited use.

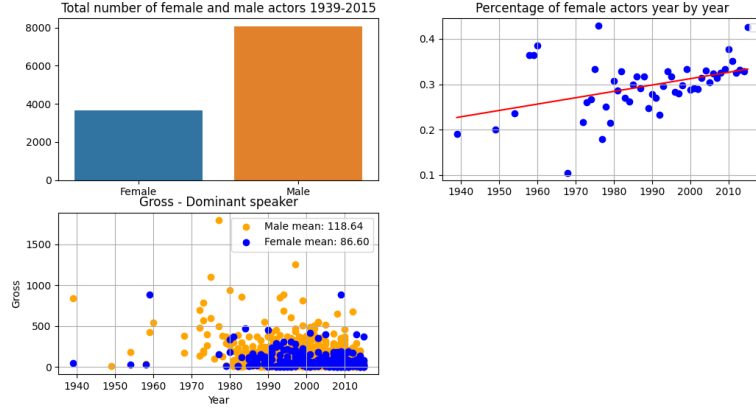


Figure 1: Various plots analyzing the gender bias in Hollywood movies.

3 Implementation of methods

There is an abundance of machine learning methods used for classification problems. This analysis, however, restricts itself to 5 different methods, these being logistic regression, discriminant analysis, k -nearest neighbour model, tree-based methods, and deep learning. The analysis of most of these methods relies heavily on the concept of cross-validation. This is a statistical technique used to test and train these models, by splitting the data set into training- and validation sets, whence the final error is obtained upon averaging the previously obtained errors.

3.1 Logistic regression

Logistic regression is a classification model, while the name itself could suggest something else. Here we present a short description of the fundamentals behind the model, and how to implement and tune it for this specific problem. Assuming the reader has a basic understanding of linear regression we can start from there and build our logistic regression model: $z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p$. To fit our model to a prediction of probability $p(y = 1|x)$ we use the logistic function also known as the Sigmoid function, which is defined as $f(z) = \frac{e^z}{1+e^z} \in [0, 1]$. Given the probability of one class in the case of a classification problem with 2 classes, we have the probability for the second class $p(y = -1|x) = 1 - f(z)$. By choosing 1 and -1 as the labels we can simplify our expressions such that we are left with $f(z) = \frac{e^z}{1+e^z} = 1 - f(z) = \frac{e^{\theta^T x}}{1+e^{\theta^T x}}$. Next, by using the training data, the goal is to find an optimal set of parameters $\hat{\theta}$ for the model. Using the maximum likelihood approach we find $\hat{\theta} = \arg \max p(\mathbf{y} | \mathbf{X}; \theta) = \arg \max \sum_{i=1}^n \ln p(y_i | \mathbf{x}_i; \theta)$, which can be turned into a minimization problem by using the negative log-likelihood as the cost function. Since the labels are chosen in a clever way, we end up with the cost function: $J(\theta) = \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y_i \theta^T x_i})$. Then, the optimal parameters for the logistic regression model is found by solving the minimization problem

$$\hat{\theta} = \arg \min \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y_i \theta^T x_i}) \quad (1)$$

However, this modification to the regression model is not "perfect", as there exists no closed-form solution. To find the parameters $\hat{\theta}$, the minimization problem has to be solved numerically [2]. In the scikit-learn library, different numerical solvers for the logistic regression model can be chosen, all of which have their own pros and cons. This hyperparameter is one that can be tuned in logistic

70 regression. Others are the type of regularization and its strength.

71

72 Regularization is a method to decrease the flexibility, or in other words, to avoid overfitting the
73 training data. To achieve this, a penalty term is added to the loss function. L1 regularization penalizes
74 the parameter weights in proportion to the sum of the absolute value of the weights. This penalty can
75 drive the weight for an irrelevant feature parameter to 0. An other option is L2 regularization, which
76 penalizes the parameter weights in proportion to the sum of squares of the weights. The strength of
77 the penalty is regulated with the coefficient: "Regularization strength" or C in the scikit-learn library.

78 Tuning the model further can be done by optimizing the features that are used. Since the features
79 in the data set are all related to the problem and there is no simple way to exclude any of them One
80 option is therefore to test all possible combinations of features by "brute forcing".

81 We start tuning the model by testing different solvers using the default settings for the other hy-
82 perparameters. First, a model is trained and tested using 10-fold cross-validation, after which the
83 solver and set of features is picked out. Since some of the solvers have convergence problems with
84 non-normalized data, the standardScaler and MinMaxScaler from the scikit-learn library are used to
85 scale the input data. Even though all features are not normally distributed, it proved empirically better
86 to use standardScaler rather than MinMaxScaler. The test resulted in the following solver: 'lbfgs'
87 and set of features: 'Number words female', 'Number of words lead', 'Difference in words lead and
88 co-lead', 'Number of male actors', 'Number of female actors', 'Number words male', 'Mean Age
89 Male', 'Age Lead', 'Age Co-Lead'.

90 Next, the remaining optimal set hyperparameters to the logistic regression model, these being the
91 type of regularization and its strength, are found by again testing all possible combinations and seeing
92 which combinations perform best in a 10-fold cross-validation. A useful package from scikit-learn
93 library is GridSearchCV, which allows us to set up lists of regularization methods and strengths and
94 test all of them, after which it returns the best combination. The input for the grid search is Solver:
95 'lbfgs', penalty: ['l2'], C: np.logspace(-3,3,7), which returned the results below.

96 • **Model:** LogisticRegression(solver='lbfgs', penalty='l2', C=1.0)

97 • **Mean accuracy:** 87.5%

98 Comparing the result to a naive classifier that only predicts a male lead, which has an accuracy of
99 around 75%, we can see an improvement of about 12%.

100 3.2 Discriminant analysis: LDA, QDA

101 The Gaussian Mixture Model (GMM) is a generative model which attempts to model $p(\mathbf{x}, y)$, where
102 \mathbf{x} is assumed to be numerical- and y a categorical variable. Here, we will look at two GMMs:
103 Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA). These methods
104 are considered when working with fully labeled data. Unlike the *discriminative* models, the GMM
105 describes the joint distribution of both inputs and outputs. $p(\mathbf{x}, y)$ states that, given input parameters
106 \mathbf{x} , one finds a categorical output y with a certain probability. The GMM makes use of the following
107 factorization of the joint probability density function: $p(\mathbf{x}, y) = p(\mathbf{x}|y)p(y)$. Since y is categorical,
108 the categorical distribution can be written as $p(y = m) = \pi_m$, where the number of classifications
109 ranges from 1 to m . Additionally, it is assumed that all probabilities follow a Gaussian distribution
110 $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mu_y, \Sigma_y)$, where μ_y is the mean vector of the parameter and Σ_y the covariance matrix.
111 The GMM model can be trained using a supervised and semi-supervised approach. The goal of
112 GMM is to determine the parameters $\theta = \{\mu_m, \Sigma_m, \pi_m\}_{m=1}^M$, which can be achieved using the
113 log-likelihood function, $\hat{\theta} = \arg \max_{\theta} \ln p(\mathcal{T}|\theta)$, where $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ is the training data. This
114 yields the following solutions: $\hat{\pi} = \frac{n_m}{n}$, $\hat{\mu}_m = \frac{1}{n_m} \sum_i \mathbf{x}_i$, and $\hat{\Sigma}_m = \frac{1}{n_m} \sum_i (\mathbf{x}_i - \hat{\mu}_m)(\mathbf{x}_i - \hat{\mu}_m)^T$,
115 where n_m is the total number of data point in class m , and n the total number of data points [2].

116 The next step is to use the trained model to make predictions y_* , i.e., we want to find $p(y|\mathbf{x})$. Using
117 Bayes' rule we get

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} = \frac{p(\mathbf{x}, y)}{\sum_{j=1}^M p(\mathbf{x}, y = j)} \quad (2)$$

118 The final category is selected by using the largest probability, given the input \mathbf{x}_*

$$\hat{y} = \underset{m}{\operatorname{argmax}} \ln p(y = m | \mathbf{x}_*), \quad (3)$$

119 which is known as the decision boundary. The equation above can be extended to the Quadratic
120 Discriminant Analysis (QDA) method by taking the logarithm. Putting this all together, we get the
121 final equation

$$\hat{y}_* = \underset{m}{\operatorname{argmax}} \{ \ln \hat{\pi}_m + \ln \mathcal{N}(\mathbf{X}_* | \hat{\mu}_m, \hat{\Sigma}_m) \} \quad (4)$$

122 Taking the logarithm of the Gaussian distribution then gives rise to a quadratic decision boundary.

123 3.2.1 Training and evaluating LDA and QDA models

124 In order to complete this task, the scikit-learn library for supervised machine learning was used.
125 This is a tool that makes it easy to implement the training and optimization of LDA and QDA
126 methods. This library has functions for LDA and QDA, which creates the models and optimizes
127 hyperparameters depending on the data that is worked with. The main hyperparameter for the LDA
128 methods is the solver. There are three different solvers: "svd" (singular value decomposition), "lsqr"
129 (least squares solution), and "eigen" (eigenvalue decomposition). When choosing lsqr or eigen, there
130 is an additional parameter that can be tuned, which is the shrinkage parameter. The hyperparameter
131 that was tuned for QDA was the regularization parameter, which ranges from 0 to 1. This parameter
132 tries to combat overfitting and collinearity. Here it was found that the standard svd solver gives the
133 highest accuracy for the LDA method, while a regularization parameter of 0.50 was found for the
134 QDA method.

135 The next step is feature selection. This was optimized by training the models while withholding one
136 feature from the dataset. Essentially, 13 different data sets were constructed, which in turn created 13
137 different models. This was tested using the svd solver for the LDA method and with regularization
138 parameter of 0 and 0.50 for the QDA method. Using the LDA method, it was found that the highest
139 misclassification accuracy was found when omitting the data from "Mean Age Female". The highest
140 misclassification accuracy found with the LDA method was 13.5%. For the QDA it was found that
141 the highest misclassification accuracy was obtained when using a regularization parameter of 0.5 and
142 discarding the data from "Year". The highest misclassification accuracy obtained was 10.1%.

143 3.3 k-nearest neighbour

144 From a conceptual point of view, the k -nearest neighbour (k -NN) algorithm is arguably the most
145 straightforward one. The core idea of the algorithm is to predict similar outputs for those inputs
146 that are in close proximity to each other in parameter space. The k -NN algorithm belongs to the
147 class of distance-based machine learning methods, and, despite its simplicity, can be applied to
148 both regression and classification problems. In more detail, suppose we have a set of training data
149 $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and an input \mathbf{x}_* whose output we want to predict. Then, one starts by finding the
150 k closest data points $\{\mathbf{x}_i\}_{i=1}^k$ to \mathbf{x}_* , which gives rise to the name k -NN, by computing the distance
151 $\|\mathbf{x}_* - \mathbf{x}_i\|$ between them. Note that the metric need not necessarily be induced from the ℓ^2 norm. In
152 fact, every norm induces a metric, which can be used to compute the distance between the input
153 and data points, whose physical interpretation heavily depends on the type of metric that is used.
154 Norms that are frequently used are the ℓ^1 -, ℓ^2 - and the more general ℓ^p norm, whose metrics are more
155 commonly known as the Manhattan-, Euclidean-, and Minkowski distance, respectively. Moreover,
156 the inclusion of different norms to compute the distances extends the regular k -NN algorithm to
157 the kernel k -NN algorithm, where the similarity of data points can be specified according to the
158 used kernel [3]. Next, once the k -NNs are found, a prediction is made depending on the type of
159 problem (regression/classification). For classification with y classes, the prediction then takes the
160 form of $y_* = \max_y \sum_{(\mathbf{x}_i, y_i) \in \mathcal{T}} \omega_i \mathbb{1}_{\{y=y_i\}}$, where $\omega_i \in [0, 1]$ are so-called weights and $\mathbb{1}$ denotes
161 the indicator function. For unweighted classification, all weights are equal to 1, but in the case where
162 one would like to further decrease the importance of the data points furthest away from \mathbf{x}_* , their
163 respective weights can be scaled down.

165 The k -NN algorithm for classification is quite rigid, as it only leaves room for the hyperparameter
166 k to be varied, as well as deciding whether or not to classify using weights. While the latter is

merely done by setting a flag to yes or no, the former is a bit more flexible and needs to be properly investigated. Furthermore, it is important to note that the k -NN model behaves differently for each k depending on which features are used in the input. This investigation concerns therefore both the input features and k , simultaneously. While there exist roughly 8200 possible combinations of input features, here, without loss of generality, we restrict ourselves only to combinations of 10, 11, 12, and 13 input features, where a minimum 10 features is set to avoid underfitting the model.

Pipelines from the Scikit-learn python library are used to determine the best combination of input features and k , by doing a 10-fold cross-validation for all of the aforementioned combinations of features, while simultaneously varying k to include 1 to 50 nearest neighbours. Here, every input is normalized so that scaling does not play a role, which is a known issue for the k -NN algorithm. The results for various norms and the inclusion/exclusion of weighted classification are presented in table 1. For visual clarity and brevity, the exact input features are omitted. Three of the models reach the

Model	Accuracy	k
Unweighted ℓ^2	0.864	14
Weighted ℓ^2	0.883	12
Unweighted ℓ^1	0.883	8
Weighted ℓ^1	0.883	11
Unweighted ℓ^3	0.874	8

Table 1: Cross-validation accuracy, the optimal value of k , and optimal number input features for different models.

same maximum accuracy. A bigger data set could have allowed for tiny differences in performance, which allows a single optimal model to be found. However, this is not an option, so the weighted ℓ^2 model has been selected as the best model due to its larger optimal value of k , as this allows for more flexibility than the other models. For this model, the accuracy per k is plotted in figure A.1.

Evident from the figure is that the number of nearest neighbours reaches its optimum very quickly. Moreover, the accuracy quickly drops due to overfitting, as is indicated by the long "tail" for large values of k . Finally, the optimal k -NN model for the given data uses the input features 'Number words female', 'Total words', 'Difference in words lead and co-lead', 'Number of male actors', 'Year', 'Number of female actors', 'Number words male', 'Mean Age Male', 'Age Lead', 'Age Co-Lead', together with $k = 12$.

3.4 Tree-based methods: classification trees, random forests, bagging

The basic idea behind tree methods is to split the input space into regions. The splitting is done according to rules, which essentially turns it into a binary tree. At each internal node, a rule $x_j < s$ with input parameter x_j and threshold s decides which branch will be chosen for the given input x . The leaf nodes then decide which class is chosen. Each leaf node corresponds to one region. The deeper the classification tree, the finer the partition [2]. For this project, three methods will be explained and tested. Classification trees are the most basic form. Bagging is an algorithm to improve the variance of the classification tree, and Random Forest uses a random subset of multiple trees to even further decrease the variance of the classification tree.

Classification Tree: To find the best partition, which means minimizing the training error, a greedy recursive binary splitting is used. It splits regions recursively in two half-spaces until a stopping criterion is met. The class prediction in each region can be expressed by $p(y = m|x_*) \approx \sum_{l=1}^L \hat{\pi}_{lm} \mathbb{I}\{x_* \in R_l\}$, with the proportion of data points in region l of class m described by $\hat{\pi}_{lm} = \frac{1}{n_l} \sum_{i: x_i \in R_l} \mathbb{I}\{y_i = m\}$. At each internal node, the input variable and threshold with a minimal error are chosen. Common error measurements are Log Loss, Entropy, and Gini index [2].

Bagging: Bootstrap Aggregation or Bagging tries to lower the variance without raising the bias. To get a low bias, a tree has to grow deep, which can result in overfitting, yielding a high variance. Multiple slightly different classification trees are trained with different training sets in order to avoid this. Bootstrapping is a technique to create multiple datasets of size n from one separate dataset of size n . To obtain multiple data sets from one training set, one samples n times from the training set

211 and repeat this B times to get B data sets. After that, one can learn B models whereafter averaging
 212 over them results in a low variance.

213 **Random Forest:** Since the individual models used in Bagging are correlated, the variance reduction
 214 is limited. The Random Forest approach tries to fix this issue by using only a random subset of
 215 inputs as splitting variables at each internal node. This results in more different individual trees since
 216 dominant input variables are not picked all the time. Computationally this is also better, as it allows
 217 for calculations to be done in parallel, on top of using smaller subsets of variables.

218 **Implementation** The scikit-learn library provides implementations to all the methods, as well as a
 219 gridsearch implementation. The feature selection was done per the 'brute force' method, for each
 220 of the three methods, individually. All possible combinations of input features were used to fit the
 221 model and the accuracy was tested using 10-fold cross-validation. To tune the hyperparameter of
 222 each method, the gridsearch is performed with each method and different parameters. The result-
 223 ing accuracy is used to determine the best parameters. The following final models are calculated:
 224 **Classification Tree** with input features ['Number words female', 'Number of words lead', 'Differ-
 225 ence in words lead and co-lead', 'Number of male actors', 'Number of female actors', 'Number
 226 words male', 'Age Lead'] and parameter (max_depth=6, min_samples_split=3, min_samples_leaf=5,
 227 max_features='log2', splitter='best', criterion='gini'), **Random Forest** with ['Number words female',
 228 'Number of words lead', 'Difference in words lead and co-lead', 'Number of male actors', 'Num-
 229 ber of female actors', 'Number words male', 'Age Co-Lead'] and (n_estimators=5, max_depth=9,
 230 min_samples_split=5, min_samples_leaf=2, max_features=None, criterion='entropy') and **Bagging**
 231 with ['Number words female', 'Number of words lead', 'Difference in words lead and co-lead',
 232 'Number of male actors', 'Number of female actors', 'Number words male', 'Mean Age Female'].
 233 Since Bagging performed better without hyperparameter tuning, no gridsearch was done, and instead
 234 default parameters were used. The final accuracies are for the basic tree is 0.837, while Random
 235 Forest achieves an accuracy of 0.877 and Bagging reaches 0.861.

236 3.5 Deep learning

237 Deep Learning (DL) can be considered a subset of Machine Learning. DL deals with hierarchical
 238 machine learning models and can describe more complex relations between input and output. A Deep
 239 Neural Net (DNN) consists of three regions: the input layer, the hidden layer, and the output layer.
 240 In this case, the input layer and output layer are the same as with the statistical machine learning
 241 methods. It is in the hidden layers where DL differentiates itself from standard machine learning.
 242 The hidden layer can be understood as many independent layers of linear regression models with
 243 non-linear activation functions.

244 There are many types of neural networks, here a fully connected neural network (FCN) is used. The
 245 general structure of a FCN looks as follows:

$$y_{jk}(x) = \sum_{i=1}^N \sigma(W_{jk}x_i + b_j). \quad (5)$$

246 Here the final output y_{jk} is given by the matrix multiplication of N hidden layers. Each layer consists
 247 of the weight matrix W and the bias b , and the number of nodes in each layer is represented by
 248 x_i . The σ denotes the activation function of each layer. It is standard to use the rectified linear
 249 activation function (ReLU) for each hidden layer, and a separate activation function for the final layer,
 250 depending on how the output looks (i.e. Sigmoid activation function in the case of a binary output).
 251 As can be seen in the equation above, the number of layers and nodes in each layer gives rise to a
 252 large sequence of matrix multiplication in order for us to find the output. Thus, computation time and
 253 training of the DNN heavily depend on the size of your network [2].

254 Once the network parameters are initialized (weights and biases), DNN can start to optimize its
 255 parameters, which is referred to as training. There is a multitude of ways in which a DNN can
 256 optimize its parameters. This is done by using a loss function $L(\mathbf{x}_i, \mathbf{y}_i, \theta)$ and a cost function $J(\theta)$.
 257 Additionally, one has to choose an optimizer, which is the algorithm for the gradient descent of the
 258 parameters. This often comes in the form of standard packages that are available. They use a variety
 259 of methods to descend toward a global minimum for the parameters. This process is repeated until
 260 one terminates the training when a criterion is fulfilled.

261 The network was tested using multiple layers with multiple different numbers of nodes in each layer.
 262 The best parameter settings were found with three hidden layers, with a number of nodes of 128, 64,
 263 and 32, for each layer, respectively, using the ReLU activation function. A dropout layer showed to
 264 have a negative effect and this option was discarded. The output layer consists of a single neuron
 265 with a Sigmoid activation function. For the optimizer, the "Adam" optimizer was used, as, in the
 266 industry, this optimizer is considered to be the best-performing optimizer. The optimal learning rate
 267 was found to be 0.001. The network was trained with 100 epochs with a batch size of one. The final
 268 accuracy with this model was found to be $87.9 \pm 3.4\%$.

269 3.6 Model selection

270 In order to determine the best model, a 10-fold cross-validation is done for each of the optimized
 271 models, as well as the naive classifier. For the logistic regression and k -NN models, the inputs are
 272 scaled as they were trained with scaled data. The results of the cross-validation are summarized in a
 273 box plot in figure A.2. The performance of each of the models is measured using the misclassification
 274 accuracy, where a lower indicates a better performance. On top of a low misclassification accuracy,
 275 a good model also should be unbiased and have a low variance. It turns out that the QDA model
 276 performs best, as it reaches the lowest misclassification accuracy, whilst also having an almost
 277 vanishing variance. Next, logistic regression and random forest share a similar performance, as they
 278 have a slightly higher misclassification accuracy than QDA, but maintain a low variance. As expected,
 279 random forest performs better than the other tree methods, as it improves them. The Bagging model
 280 has a moderate variance, which is not expected, indicating that something possibly went wrong in
 281 the implementation. The k -NN model has a low variance, however, its misclassification error is
 282 quite poor. Lastly, the LDA and DL methods perform similarly but are not quite as good as the QDA
 283 model.

284 In light of the analysis above, we have decided to put the QDA model 'in production', as its
 285 performance is unmatched by all other models.

286 4 Conclusion

287 In this work, we train different machine learning models using a data set of 2000 Hollywood movies
 288 in order to accurately predict the sex of the lead role of a movie, given a list of features corresponding
 289 to each movie. Here we test and cross-validate five different machine learning models: logistic
 290 regression, linear and quadratic discriminant analysis, k -nearest neighbour, deep learning, and
 291 tree-based methods: classification trees, random forests, and bagging.

292
 293 The data provided was heavily biased towards male actors, both in the number of lead actors and the
 294 total number of female and male speaking roles. This imbalance was also tested when comparing
 295 the percentage of female actors over time. It turns out that, even though a male bias persists, it has
 296 gradually improved over the years, getting closer to a more even and fairer split. There was also
 297 an economic component included in the provided data. When looking at the economic differences
 298 between male- and female-dominated movies, it was found that male-dominated movies generate
 299 more gross than female-dominated movies. All in all, the data provided was heavily biased towards
 300 male leads and that had to be taken into account when constructing the machine learning models.
 301

302 It was found that the generative QDA model gives the best performance after a 10-fold cross-
 303 validation, as out of all models it has the best mean prediction accuracy of 90% and a near-vanishing
 304 variance. Surprisingly, the model performed best when discarding the 'year' feature. Even though
 305 we found a strong correlation between features and the percentage of female actors over the years.
 306 When looking at the hyperparameters of the QDA model, we found they have the highest accuracy
 307 with a regularization parameter of 0.50. Further improvements could be made by brute forcing all
 308 combinations of input features, and by creating new features from old ones.

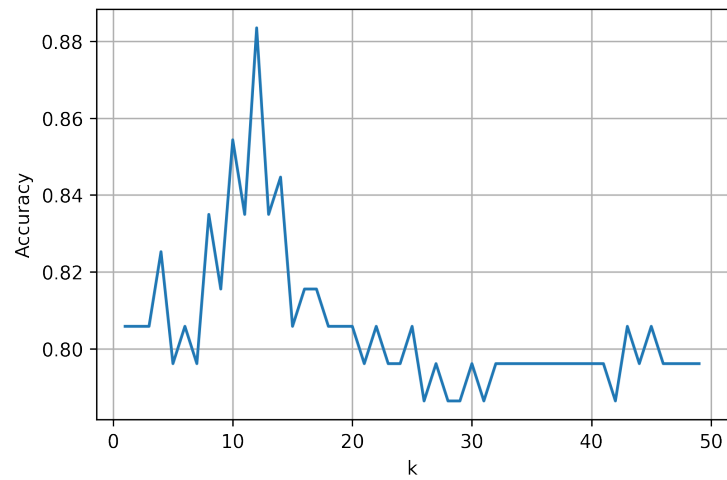


Figure A.1: Accuracy of the k -NN model, which uses the ℓ^2 norm to do a weighted classification, for the number of nearest neighbours k .

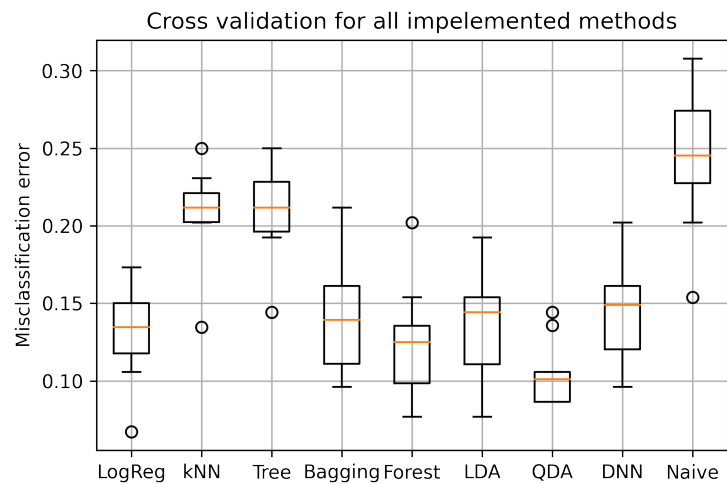


Figure A.2: Box plot of performance of all optimized models.

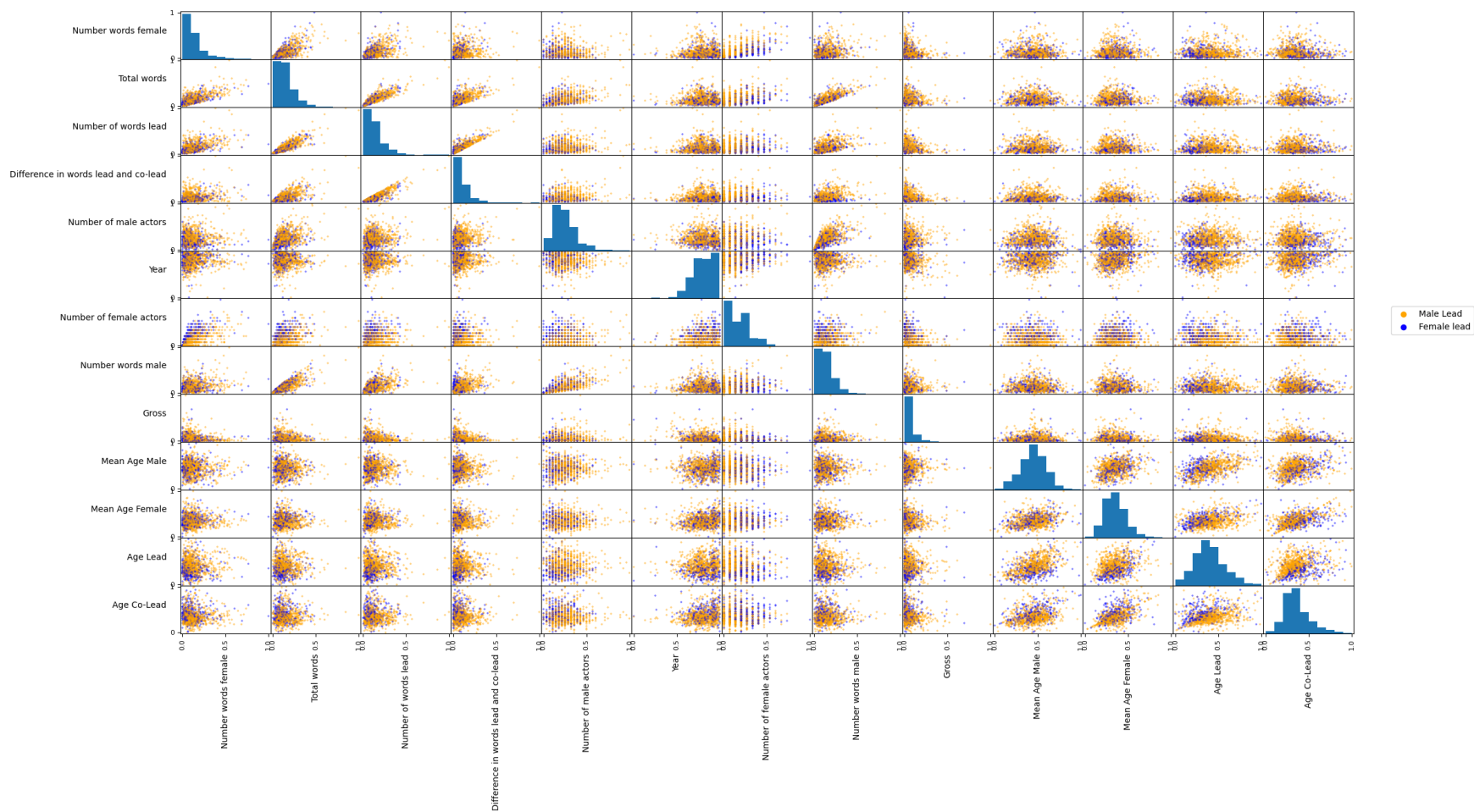


Figure A.3: Histogram and scatter plot.

310 **Supplementary code for the reproduction of results.** A lot of the code is inspired by or taken
311 from the exercise sheets of the Statistical Machine Learning course [4]. These, in turn, use a lot of
312 functions from external python libraries [5–18].

313 References

- 314 [1] The Pudding, “MS Windows NT kernel description.” <https://pudding.cool/2017/03/film-dialogue/>.
- 315
- 316 [2] A. Lindholm, N. Wahlström, F. Lindsten, and T. B. Schön, *Machine Learning - A First Course*
317 *for Engineers and Scientists*. Cambridge University Press, 2022.
- 318 [3] K. Yu, L. Ji, and X. Zhang, “Kernel nearest-neighbor algorithm,” *Neural Processing Letters*,
319 vol. 15, no. 2, pp. 147–156, 2002.
- 320 [4] S. M. L. team, “Statistical machine learning exercise sheets,” *Available through course page on*
321 *Studium*, 2022.
- 322 [5] Contributors to Wikimedia projects, “Covariance - Wikipedia,” Nov. 2022. [Online; accessed 7.
323 Dec. 2022].
- 324 [6] Scikit-learn, “sklearn.ensemble.BaggingClassifier,” Dec. 2022. [Online; accessed 22. Dec.
325 2022].
- 326 [7] Scikit-learn, “sklearn.model_selection.cross_validate,” Dec. 2022. [Online; accessed 22. Dec.
327 2022].
- 328 [8] Scikit-learn, “sklearn.tree.DecisionTreeClassifier,” Dec. 2022. [Online; accessed 22. Dec.
329 2022].
- 330 [9] Scikit-learn, “Scikit-learn: GridSearchCV,” Nov. 2022. [Online; accessed 7. Dec. 2022].
- 331 [10] Scikit-learn, “sklearn.model_selection.KFold,” Dec. 2022. [Online; accessed 22. Dec. 2022].
- 332 [11] Scikit-learn, “Scikit-learn: LogisticRegression,” Nov. 2022. [Online; accessed 7. Dec. 2022].
- 333 [12] Matplotlib, “matplotlib.pyplot.boxplot — Matplotlib 3.6.2 documentation,” Nov. 2022. [Online;
334 accessed 22. Dec. 2022].
- 335 [13] Scikit-learn, “Scikit-learn: MinMaxScaler,” Nov. 2022. [Online; accessed 7. Dec. 2022].
- 336 [14] Numpy, “Numpy: Polyfit,” Nov. 2022. [Online; accessed 7. Dec. 2022].
- 337 [15] Scikit-learn, “sklearn.ensemble.RandomForestClassifier,” Dec. 2022. [Online; accessed 22.
338 Dec. 2022].
- 339 [16] Google, “Google Machine Learning Education: Regularization,” Nov. 2022. [Online; accessed
340 7. Dec. 2022].
- 341 [17] Pandas, “Pandas: Scatter matrix,” Nov. 2022. [Online; accessed 7. Dec. 2022].
- 342 [18] Scikit-learn, “Scikit-learn: StandardScaler,” Nov. 2022. [Online; accessed 7. Dec. 2022].