

---

# Do (wo)men talk too much in films?

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Introduction

In a movie, there is almost always a main character that the rest of the film is in some way centered around. In some films the main character is not very nice, but most of the time, the audience is expected to sympathize with this character or even look up to, and be inspired by them. This is especially true in children's movies. Since it is often easier to be inspired by someone you can relate to, it is important to make sure that the distribution of main characters in movies at least approximates the distribution of the audience. If the main character is always a man, that makes it more difficult than necessary for girls and women to find someone to be inspired by and vice versa.

This paper investigates how the gender of the lead actor can be predicted using metrics such as the year when the movie was made and the age of said actor. Being able to make such predictions could give an insight into how the movie industry works with respect to the gender of the lead actor. This could provide clues about what areas to investigate further to understand why those connections exist and ultimately, make sure that everyone has a chance to be inspired by someone they can relate to.

## 2 Methods

We have chosen to focus on approaches using logistic regression, k-NN, LDA and QDA to classify the lead actor's gender.

In order to make the methods as comparable as possible, we have used a common set of transformations of the input variables for all tested methods.

To compare between families of models and between which tuning is better we chose to focus on two measures: accuracy (on average, of how often model makes a correct prediction) and AUC (area under ROC curve).

Several methods have a hyper-parameter  $\lambda$  that needs to be tuned. In order to find a value of  $\lambda$  that performs well on the data, cross-validation is used to find the optimal value in a finite set  $\Lambda = \{\lambda_1, \dots, \lambda_k\}$ . Cross-validation works by splitting the data into  $n$  equally sized partitions and training the data separately on the  $n$  choices of  $n - 1$  partitions and testing on the partition that was left out. The test error  $E_{new}$  is estimated by the mean misclassification rate across the partitions. This procedure is repeated for each  $\lambda_j \in \Lambda$  and the value resulting in the lowest estimated test error is chosen.

Since cross-validation is used to estimate the hyper-parameter  $\lambda$ , it cannot be used to estimate the test error of the whole procedure. Instead, the data has to be split into a training set and a testing set. Cross-validation is done on the training set, and the test error is estimated by evaluating performance on the testing set. This can yield significantly different estimates of the test error as only one split into training and testing data is considered. To get a better estimate of the actual test error, a bootstrap procedure is performed.

The full dataset is an iid sample from some unknown distribution so the estimated test error  $\hat{E}_{new}$  is a random variable. By drawing  $B$  independent splits into training and testing data with subsequent fitting and cross-validation, a bootstrap sample of  $\hat{E}_{new}$  is obtained which can be used to obtain a better estimate of  $E_{new}$ . This is very computationally intensive if  $B$  is large.

## 2.1 Input transformations

In the given dataset, there are columns for the total number of words spoken as well as the number of words spoken by the lead, the co-lead etc. This could present a problem if we compare a movie where the lead says 10 out of 100 total words and another movie where the lead says 100 out of 1000 words. Most models would think that the lead speaks more in the second movie and miss the fact that the *proportion* of words spoken by the lead is the same. Hence, we have transformed several input variables to express a proportion instead of absolute numbers. We also believe it might be important to have a dummy variable indicating if the lead or the co-lead is oldest. All transformations are given in Table 1.

Table 1: Transformations of input variables.

Original column	New column	Transformation
Number of words lead	Proportion of words lead	$\frac{\text{Number of words lead}}{\text{Total words}}$
N/A	Proportion of words co-lead	$\frac{\text{Number of words lead} - \text{Difference in words lead and co-lead}}{\text{Total words}}$
Difference in words lead and co-lead	Ratio words co-lead lead	$\frac{\text{Proportion of words co-lead}}{\text{Proportion of words lead}}$
Number words female	Proportion of words female	$\frac{\text{Number words female}}{\text{Total words} - \text{Number of words lead}}$
Number of female actors	Proportion of female actors	$\frac{\text{Number of female actors}}{\text{Number of female actors} + \text{Number of male actors}}$
Number of male actors	Number of actors	$\frac{\text{Number of male actors} + \text{Number of female actors}}{\text{Number of male actors} + \text{Number of female actors}}$
N/A	Older lead	$\begin{cases} 1, \text{Age lead} > \text{Age Co-Lead} \\ 0, \text{else} \end{cases}$

Note that when determining 'Proportion of words female', this should only measure the words spoken by non-lead female actors so we have to subtract the lead's contribution to the total number of words. The column 'Number of male actors' was dropped since all necessary information in this column is contained in 'Proportion of female actors' together with 'Number of actors'. In order to improve regularization and k-NN, all remaining numerical input variables were centered and scaled by their standard deviation. This means that columns with proportions have values in the unit interval  $[0,1]$  and the other numerical variables have values that are of roughly the same magnitude. This scaling was not done for QDA as it is not necessary for that method.

## 2.2 Logistic Regression

Logistic regression is a *general linear model* (GLM), i.e. the relationship between the data  $X \in \mathcal{X} \subseteq \mathbb{R}^p$  and the outcome  $Y$  is on the form

$$E(Y|X = x) = g^{-1}(x \cdot \beta) \quad (1)$$

where  $\beta \in \mathbb{R}^p$  and  $g$  is the link function. In the case of logistic regression,  $Y|(X = x) \sim \text{Ber}(p(x))$  and the canonical link function is the logit link  $g(x) = \log\left(\frac{x}{1-x}\right)$  with  $g^{-1}(x) = \frac{\exp(x)}{1+\exp(x)}$ . Since  $Y|(X = x) \sim \text{Ber}(p(x))$ , we get  $E(Y|X = x) = p(x) = g^{-1}(x \cdot \beta)$ . In other words,  $P(Y = 1|X = x) = g^{-1}(x \cdot \beta)$ , which we can use to predict  $Y$  given data  $x$ .

To do the regression, we find  $\hat{\beta} \in \arg \min_{\beta} \sum_{i=1}^n (y_i - \hat{y}(x_i; \beta))^2$  where  $\hat{y}(x; \beta) = g^{-1}(x \cdot \beta)$ . This is the MLE estimator of  $\beta$ . This minimizes the mean squared error (MSE) loss function.

69 A potential problem with this approach is that there are no restrictions on the components  
70 of  $\beta$  and that can lead to overfitting, especially if  $n$  is not much larger than  $p$ . To address  
71 that issue, one can introduce regularization.



72 In general, regularization is done by adding a penalizing term to the loss function that restricts  
73  $\beta$  in some way. If  $L(\beta; x_i, y_i)$  is the loss function before regularization, we instead consider  
74 the new loss function  $L(\beta; x_i, y_i) + \lambda R(\beta)$  and find  $\hat{\beta}_{reg} \in \arg \min_{\beta} (L(\beta; x_i, y_i) + \lambda R(\beta))$ .  
75  $R$  is some penalizing function and  $\lambda$  is a hyper-parameter that can be tuned. The two most  
76 common forms of regularization is LASSO and Ridge regression.

77 LASSO regression uses  $L_1$ -regularization, meaning that  $R_{LASSO}(\beta) = \|\beta\|_1 = \sum_{i=1}^p |\beta_i|$   
78 while Ridge regression uses  $L_2$ -regularization,  $R_{Ridge}(\beta) = \|\beta\|_2^2 = \sum_{i=1}^p \beta_i^2$ .

## 79 2.3 k-Nearest Neighbors

80 The  $k$ -nearest neighbors ( $k$ -NN) method is based on the simple principle of finding the  $k$   
81 closest neighboring points with respect to the input data  $X \in \mathcal{X} \subseteq \mathbb{R}^p$ . In the case of  
82 classification the outcome  $Y$  is then determined by a majority vote among the  $k$  nearest data  
83 points. The method is based on the idea that if a test data point is close to some training  
84 data point then the prediction should be that they have the same outcome  $Y$ .

85 The algorithm for  $k$ -NN can be implemented in a simple manner with a brute force algorithm  
86 for measuring the distance from the test data point  $\mathbf{x}_\star$  to each training data point  $\mathbf{x}_i$ , where  
87  $i = 1, \dots, n$  using some distance function  $d(x, y)$ . It is standard to use the Minkowski distance  
88 for a certain order  $p$ , depending on the problem, which is given by

$$d(\mathbf{x}, \mathbf{y}) = \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}, \text{ where } \mathbf{x} = (x_1, \dots, x_p), \mathbf{y} = (y_1, \dots, y_p) \in \mathbb{R}^p. \quad (2)$$

89 Note that  $p = 1$  gives the Manhattan distance, and  $p = 2$  gives the Euclidean distance. The  
90 brute force algorithm for  $k$ -NN is given by [Lindholm et al., 2021]

91

- 
- 92 1. Calculate the distance  $d(\mathbf{x}_i, \mathbf{x}_\star)$  for each  $i = 1, \dots, n$
  - 93 2. Set  $\mathcal{N}_\star = \{\mathbf{x}_i : \text{Where } \mathbf{x}_i \text{ is one of the } k \text{ nearest points}\}$
  - 94 3. Return  $\hat{y}(\mathbf{x}_\star) = \text{MajorityVote}\{y_j : j \in \mathcal{N}_\star\}$
- 

95

96 A problem with the brute force algorithm is that for each point we need to calculate the  
97 distance to every other point, which is computationally demanding for larger datasets.

98 There are however more computationally efficient algorithms to find the  $k$ -NN compared to  
99 the brute force search such as ball-tree and  $k$ -d tree which are not explained in detail here.  
100 All three of these algorithms were tested and no significant difference in the results were  
101 noted thus the choice of algorithm was set to "auto" which chooses the best suited algorithm  
102 for a given problem, further described in the Scikit-learn documentation [skl, a].

103 For our problem we consider the Minkowski distance and let  $p$  and  $k$  be hyper-parameters  
104 which are to be tuned. This is done in an analogous manner to the case of finding the hyper-  
105 parameter  $\lambda$  in the regularization problem with logistic regression previously considered.

106 An alternative approach is to use weighted  $k$ -NN. The idea is to let the distance of the  
107 training data point to the test data point influence the strength of the vote. We compared  
108 *uniform weights* (standard  $k$ -NN, where all weights equal 1) and *distance weights* where the  
109 weights are given by

$$\frac{1}{d(\mathbf{x}_\star, \mathbf{x}_i)}, \quad (3)$$

110 for each of the  $k$ -nearest neighbors. This reinforces the idea that proximity of test data  
111 points to training data points ought be a good predictor [skl, b].

## 112 2.4 LDA and QDA

113 For classification we construct a discriminative classifier from a generative model based on  
 114 Bayes' theorem for the classes  $m = 1, 2, \dots, M$  by

$$p(y = m | \mathbf{x}) = \frac{p(\mathbf{x} | y = m)p(y = m)}{\sum_{i=1}^M p(\mathbf{x} | y = i)p(y = i)} \quad (4)$$

115 where  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ . We estimate the *uninformative prior probability* as  $\hat{p}(y = m) = \frac{n_m}{n}$   
 116 where  $n_m = \sum_{i=1}^n \mathbb{1}\{y_i = m\}$  and assume that  $p(\mathbf{x} | y = m)$  is a normal density with expected  
 117 value  $\mu_m$  and covariance matrix  $\Sigma_m$ . The assumption that distinguishes LDA and QDA is  
 118 that for LDA we assume that  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_M$  but for QDA we make no such assumption,  
 119 that is, we allow for the covariance matrices to differ. A consequence is that LDA is a special  
 120 case of QDA, hence QDA is a model of higher complexity.

121 The estimates for the normal distribution parameters for each class is given by


$$\hat{\mu}_m = \frac{1}{n_m} \sum_{i:y_i=m} \mathbf{x}_i \text{ and } \hat{\Sigma}_m = \frac{1}{n_m - 1} \sum_{i:y_i=m} (\mathbf{x}_i - \hat{\mu}_m)(\mathbf{x}_i - \hat{\mu}_m)^T \quad (5)$$

122 derived from maximum likelihood estimation and adjusting  $\hat{\Sigma}_m$  to make it unbiased. The  
 123 *pooled covariance estimate* (weighted average of the covariance matrix estimates within each  
 124 class) is given by

$$\hat{\Sigma} = \frac{\sum_{m=1}^M (n_m - 1) \hat{\Sigma}_m}{\sum_{m=1}^M (n_m - 1)} = \frac{1}{n - M} \sum_{m=1}^M \sum_{i:y_i=m} (\mathbf{x}_i - \hat{\mu}_m)(\mathbf{x}_i - \hat{\mu}_m)^T. \quad (6)$$

125 With these estimators we may express the discriminant analysis classifier as

$$\hat{p}(y = m | \mathbf{x}) = \frac{n_m \mathcal{N}(\mathbf{x} | \hat{\mu}_m, \hat{\Sigma})}{\sum_{i=1}^M n_i \mathcal{N}(\mathbf{x} | \hat{\mu}_m, \hat{\Sigma})} \quad (7)$$

126 where  $\mathcal{N}(\mathbf{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{M/2} |\Sigma|^{1/2}} \exp[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)]$  is the density for the normal  
 127 distribution with mean  $\mu$  and covariance matrix  $\Sigma$ . 

## 128 3 Results

129 When comparing different models, it is important to have a baseline, or a null model to  
 130 compare against. In this case, an obvious null model is the constant model that always  
 131 predicts the same outcome regardless of input. The best null model is the one with highest  
 132 accuracy, i.e. the constant model that predicts the most frequently occurring outcome. The  
 133 model that always predicts a male lead has an accuracy of 0.756 and is thus chosen as the  
 134 baseline.

### 135 3.1 Logistic Regression

136 For all logistic regression models fitted, the set of regularization parameters,  $\Lambda$ , consisted of  
 137 10 logarithmically spaced values between  $10^{-4}$  and  $10^4$ . This was the default value in the  
 138 methods from scikit learn and having more densely packed values did not affect the model  
 139 performance in any appreciable way. The number of folds used in cross-validation was also  
 140 10, no improvement was observed by increasing this value.

141 In Tables 2 and 3, the accuracy and AUC are estimated using the mean of 100 bootstrap  
 142 samples in the case of LASSO regression and 400 in the case of Ridge regression. The reason  
 143 for having different sample sizes is that computing the LASSO regression is much more  
 144 computationally demanding.

145 We see that the regularization does not affect the model performance much. LASSO and  
 146 Ridge regularization perform almost identically and yield at best around 0.003 extra accuracy  
 147 but considering that the different splits of the data yielded estimated test errors in a range  
 148 from 0.8 to 0.98, we cannot reject that regularization does not matter in this case.

Table 2: Accuracy and AUC for logistic regression models. 70% training data.

Input	Regularization	Accuracy	AUC
Before transformations	None	0.870	0.878
	LASSO	0.871	0.880
	Ridge	0.871	0.880
After transformations	None	0.893	0.920
	LASSO	0.895	0.921
	Ridge	0.894	0.921

Table 3: Accuracy and AUC for logistic regression models. 90% training data.

Input	Regularization	Accuracy	AUC
Before transformations	None	0.876	0.878
	LASSO	0.875	0.883
	Ridge	0.871	0.880
After transformations	None	0.895	0.924
	LASSO	0.897	0.924
	Ridge	0.898	0.923

### 149 3.2 k-Nearest Neighbors

150 When hyper-tuning  $k$ -NN the set of  $p$ -values and  $k$ -values were given by  $\{1, 1.25, 1.5, \dots, 4\}$   
 151 and  $\{1, 2, 3, \dots, 25\}$  respectively. The number of folds used in cross-validation was again set to  
 152 10. We found that  $p = 2$  (Euclidean distance) and  $k = 4$  performed best for our transformed  
 153 data set. Using these parameters the  $k$ -NN algorithm was tested and performance was  
 154 measured with the mean of 100 bootstraps samples, the size of the sample was chosen with  
 155 regards to  $k$ -NN being computationally demanding. The results are summarized in tables 4  
 and 5.

Table 4: Accuracy and AUC using k-NN with 70% training data.

Input	Weighted k-NN	Accuracy	AUC
Before transformations	Uniform	0.745	0.675
	Distance	0.780	0.688
After transformations	Uniform	0.864	0.883
	Distance	0.872	0.888

156

Table 5: Accuracy and AUC using k-NN with 90% training data.

Input	Weighted k-NN	Accuracy	AUC
Before transformations	Uniform	0.750	0.678
	Distance	0.783	0.693
After transformations	Uniform	0.875	0.891
	Distance	0.882	0.901

157 It is obvious that  $k$ -NN is drastically improved by transforming the data, before transfor-  
 158 mations the model performance was even outperformed by, or just slightly better than the  
 159 best null model. Weighted  $k$ -NN with the distance weight seemed to perform better than  
 160 the uniform weight both before and after the transformation, the impact seems to be an  
 161 increase of 0.007-0.008 in accuracy after transformation and an increase of almost 0.03 in  
 162 accuracy before transformations.

### 3.3 LDA and QDA

The given dataset was bootstrapped 400 times and both DA models were tested before and after input transformations and dropping 'Year' and 'Gross'. From the Tables 6 and 7

Table 6: Accuracy and AUC for discriminant analysis models using bootstrap. 70% training data.

Input	Model	Accuracy	AUC
Before transformations	LDA	0.856	0.870
	QDA	0.818	0.849
After transformations	LDA	0.900	0.917
	QDA	0.945	0.984

Table 7: Accuracy and AUC for discriminant analysis models using bootstrap. 90% training data.

Input	Model	Accuracy	AUC
Before transformations	LDA	0.866	0.877
	QDA	0.840	0.869
After transformations	LDA	0.900	0.918
	QDA	0.947	0.984

we conclude that QDA seems to be more apt for this problem for the final inputs. It is unclear which method performs better for the original inputs. For QDA on the original inputs, the variance in accuracy is high which can be explained by the fact that the original inputs are close to being colinear resulting in inaccurate matrix inversion of  $\Sigma$ . For example, the standard deviation of accuracy and AUC of the QDA classifier, on the original inputs, with 400 bootstrapped datasets and 90% training data were 0.076 and 0.081 respectively, compared to 0.021 and 0.019 with the final inputs, ceteris paribus. Fortunately, the final inputs are not colinear.

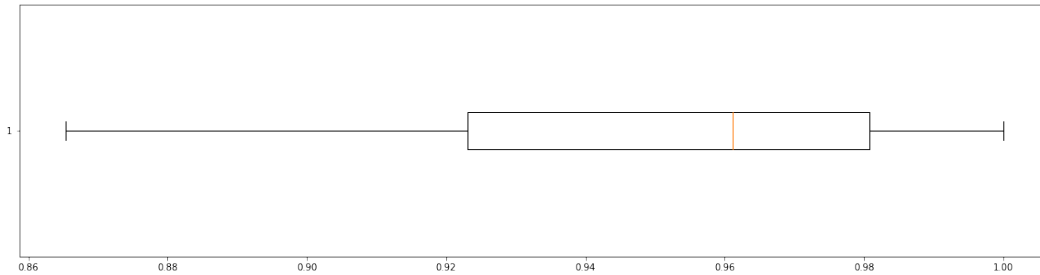


Figure 1: Accuracy estimation using cross validation with 20 folds.

Cross validation was carried out to estimate the accuracy using 20 folds resulting in an estimated accuracy of 0.949 using 70% training data. As can be seen in Figure 1, there is a noticeable variance in accuracy for the different training sets suggesting that there might be outliers in the data that the model has problems accounting for. Increasing the number of folds also increases the minimum accuracy that can be found in the corresponding box-plot, as to be expected. Also using cross validation to compare the effects of input transformations we see a 0.090 increase in accuracy using the final inputs compared to the original inputs. Adding the variables 'Years' and 'Gross' back into the inputs we see a small decrease in accuracy hence they are left out.



## 4 Conclusions

Based on the results we draw the conclusion that  $k$ -NN is the worst performing model of the ones considered. Nevertheless, it is worth noting the drastic effect that data transformations has on the  $k$ -NN model. Before transformations the model performed at the level of the baseline model, after the transformations we saw a accuracy of 0.882 and a AUC of 0.901 which is found in Table 5. The logistic regression model slightly outperformed  $k$ -NN. Looking at the best performing setup we see that logistic regression had an accuracy of 0.898 and AUC 0.923 which is seen in Table 3. The LDA model had an accuracy of 0.903 which is only a 0.005 increase from the logistic regression model. Taking variance into account, we draw the conclusion that we can not say with any confidence which model performs better. The QDA model performs best among all considered models by a margin with an accuracy of 0.945 and 0.984 AUC according to the bootstrap test and accuracy 0.949 according to the cross validation.

We answer the following questions assuming that the given data set is representative of Hollywood movies in general.

1. **Do men or women dominate speaking roles in Hollywood?** In the given data we see that in 75.6% of movies there is a male in the lead. Further, the mean of 'Proportion of words female' is 34,6% indicating that Hollywood movies consist of almost 65% male speaking.
2. **Has gender balance in speaking roles changed over time (i.e. years)?** None of the models considered suffered performance loss from removing the input variable 'Year'. However, the correlation in the given data set between 'Year' and all of 'Numbers of words female', 'Number of female actors' and 'Mean Age female' is positive and not trivially small (0.032, 0.134 and 0.170 respectively). We conclude that this indicates a slight change in the gender balance over time but that the model does not necessarily use this relationship to make predictions.
3. **Do films in which men do more speaking make a lot more money than films in which women speak more?** None of the models considered suffered performance loss from removing the input variable 'Gross' and the QDA model even improved a bit indicating that 'Gross' may be a bad input variable for this model. This might be because over the years, the change in 'Gross' for men and women may have changed at different rates and because of inflation; this makes 'Gross' a potentially bad input variable for future predictions.

All of the models fitted in this project focus on prediction instead of inference, hence it is not possible to use them to answer the above questions directly.

## 5 Feature Importance

In order to determine which of the features 'Year', 'Gross' or 'Number of female actors' is most important, we fitted models excluding one or more of the features. Using a logistic regression model with LASSO (L1) regularization, we found that the model performance in terms of prediction accuracy was completely unaffected by removing either or both of the features 'Year' and 'Gross'. On the other hand, any model that excluded the variable 'Proportion of words female' (which contains all information about how much male and female actors speak), had its performance reduced drastically. As seen in Table 3, the model including all features had an accuracy of 0.9. This dropped to 0.8 when removing the 'Proportion of words female'. Comparing this to the null accuracy of 0.756, we conclude that the 'Proportion of words female' is a very important feature in the model. Interestingly, this 'Proportion of words female' seems very important for all models, but especially so for QDA; this suggests that QDA may depend more heavily on this input for predictions. The same results hold for both QDA and  $k$ -NN as well, the models are completely unaffected by removing either 'Year' or 'Gross' (or both), while suffering 0.08 accuracy loss for  $k$ -NN and 0.12 loss for QDA. Further, not only accuracy is affected but AUC is affected in the same way.

## 235 Appendix A: Code

236 Transformations of input variables and various common functions.

```
237 import numpy as np
238 import pandas as pd
239 import sklearn.preprocessing as skl_pre
240
241 rawData = pd.read_csv('train.csv')
242
243 cols_to_norm = [
244     'Total words',
245     'Year',
246     'Gross',
247     'Mean Age Male',
248     'Mean Age Female',
249     'Age Lead',
250     'Age Co-Lead',
251     'Number of actors'
252 ]
253
254 def pre_process(raw_data, cols_to_norm):
255     data = raw_data.copy()
256
257     data['Lead'] = pd.get_dummies(data['Lead'])
258     data['Number of words co-lead'] = data['Number of words lead'] -
259         data['Difference in words lead and co-lead']
260     data['Proportion of words lead'] = data['Number of words lead']/data
261         ['Total words']
262     data['Proportion of words co-lead'] = data['Number of words co-lead']
263         /data['Total words']
264     data['Ratio words co-lead lead'] = data['Number of words co-lead']/
265         data['Number of words lead']
266     data['Proportion of words female'] = data['Number words female']/((
267         data['Total words'] - data['Number of words lead']))
268     data['Number of actors'] = data['Number of male actors'] + data['
269         Number of female actors']
270     data['Proportion of female actors'] = data['Number of female actors']
271         /data['Number of actors']
272     data['Older lead'] = data['Age Lead'] < data['Age Co-Lead']
273     data['Older lead'] = pd.get_dummies(data['Older lead'])
274
275     scaler = skl_pre.StandardScaler()
276     data[cols_to_norm] = scaler.fit_transform(data[cols_to_norm])
277
278     return data
279
280 data = pre_process(rawData, cols_to_norm)
281
282 def fit_and_test(classifier, train, test, features, target, suppress_output
283     = False):
284     classifier.fit(train[features], train[target])
285     if not suppress_output:
286         skl_met.plot_roc_curve(classifier, test[features], test[
287             target])
288         print('accuracy: ' + str(classifier.score(test[features],
289             test[target])))
290         print(' auc: ' + str(skl_met.roc_auc_score(test[target],
291             classifier.predict_proba(test[features])[:,1])) + '\n')
```



```

292         print(skl_met.classification_report(test[target], classifier.
293             predict(test[features])))
294     return classifier
295
296 rawFeatures = [
297     'Year',
298     'Number words female',
299     'Total words',
300     'Number of words lead',
301     'Difference in words lead and co-lead',
302     'Number of male actors',
303     'Number of female actors',
304     'Number words male',
305     'Gross',
306     'Mean Age Male',
307     'Mean Age Female',
308     'Age Lead',
309     'Age Co-Lead'
310 ]
311
312 featureSet1 = [
313     'Year',
314     'Gross',
315     'Number of actors',
316     'Proportion of female actors',
317     'Mean Age Male',
318     'Mean Age Female',
319     'Age Lead',
320     'Age Co-Lead',
321     'Total words',
322     'Proportion of words lead',
323     'Proportion of words co-lead',
324     'Ratio words co-lead lead',
325     'Proportion of words female',
326     'Older lead'
327 ]
328
329
330 print('Null accuracy: ' + str(max([np.mean(data[target]), 1 - np.mean(data[
331     target]))]))
332
333 Logistic Regression
334
335 trainRatio = config['Train Ratio'][0]
336 seed = config['Random Seed'][0]
337 train, test = skl_ms.train_test_split(data, train_size=trainRatio)
338
339 features = featureSet1.copy()
340 #features.remove('Proportion of words female')
341 #features.remove('Year')
342 #features.remove('Gross')
343 #features = ['Proportion of words lead']
344 target = 'Lead'
345
346 # No regularization
347
348 B = 100
349 accuracies = []
350 aucs = []

```

```

350 for i in range(B):
351     train, test = skl_ms.train_test_split(data, train_size=trainRatio)
352     logReg = fit_and_test(skl_lm.LogisticRegression(penalty='none',
353         solver='newton-cg'), train, test, features, target,
354         suppress_output=True)
355     accuracies.append(logReg.score(test[features], test[target]))
356     aucs.append(skl_met.roc_auc_score(test[target], logReg.predict_proba
357         (test[features])[:,1]))
358
359 # LASSO
360
361 B = 100
362 accuracies = []
363 aucs = []
364 for i in range(B):
365     train, test = skl_ms.train_test_split(data, train_size=trainRatio)
366     logRegLasso = fit_and_test(skl_lm.LogisticRegressionCV(Cs=10, cv=10,
367         penalty='l1', solver='liblinear', n_jobs=10), train, test,
368         features, target, suppress_output=True)
369     accuracies.append(logRegLasso.score(test[features], test[target]))
370     aucs.append(skl_met.roc_auc_score(test[target], logRegLasso.
371         predict_proba(test[features])[:,1]))
372
373 # Ridge
374
375 B = 400
376 accuracies = []
377 aucs = []
378 for i in range(B):
379     train, test = skl_ms.train_test_split(data, train_size=trainRatio)
380     logRegLasso = fit_and_test(skl_lm.LogisticRegressionCV(Cs=10, cv=10,
381         penalty='l2', solver='liblinear', n_jobs=10), train, test,
382         features, target, suppress_output=True)
383     accuracies.append(logRegLasso.score(test[features], test[target]))
384     aucs.append(skl_met.roc_auc_score(test[target], logRegLasso.
385         predict_proba(test[features])[:,1]))
386
387 k-NN
388 import sklearn.neighbors as skl_nb
389 X = data[featureSet1]
390 y = 'Lead'
391
392 ## Using gridsearch to find the optimal parameter
393 knn2 = skl_nb.KNeighborsClassifier()
394 param_grid = {'n_neighbors':np.arange(1,25),
395     'p':np.linspace(1,4,13)}
396
397 knn_gscv = GridSearchCV(knn2, param_grid, cv =10)
398
399 knn_gscv.fit(X,y)
400
401 ## k-NN Distance weight (swap 'distance' to 'uniform' for standard distance)
402     after transformations (swap featureSet1 to rawFeatures for before
403     transformations)
404
405 B = 100
406 accuracies = []
407 aucs = []

```

```

408 knn = skl_nb.KNeighborsClassifier(n_neighbors = 4, p=2, weights='distance')
409 features = featureSet1
410 for i in range(B):
411     train, test = train_test_split(data, test_size=0.3)
412     KNN_gscv = fit_and_test(knn,train, test, features, 'Lead' )
413     accuracies.append(KNN_gscv.score(test[features], test['Lead']))
414     aucs.append(skl_met.roc_auc_score(test['Lead'], KNN_gscv.predict_proba(
415         test[features][:,1]))
416
417 ## k-NN Uniform weight (swap 'distance' to 'uniform' for standard distance)
418 after transformations (swap featureSet1 to rawFeatures for before
419 transformations)
420
421
422 B = 100
423 accuracies = []
424 aucs = []
425 knn = skl_nb.KNeighborsClassifier(n_neighbors = 4, p=2, weights='uniform')
426 features = featureSet1
427 for i in range(B):
428     train, test = train_test_split(data, test_size=0.3)
429     KNN_gscv = fit_and_test(knn,train, test, features, 'Lead' )
430     accuracies.append(KNN_gscv.score(test[features], test['Lead']))
431     aucs.append(skl_met.roc_auc_score(test['Lead'], KNN_gscv.predict_proba(
432         test[features][:,1]))
433
434
435 import pandas as pd
436 import numpy as np
437 import matplotlib.pyplot as plt
438 import sklearn.preprocessing as skl_pre
439 import sklearn.linear_model as skl_lm
440 import sklearn.discriminant_analysis as skl_da
441 from IPython.core.pylabtools import figsize
442 import sklearn.model_selection as skl_ms
443 import sklearn.metrics as skl_met
444 import itertools
445 import math
446
447 url = 'train.csv'
448 dataset = pd.read_csv(url, na_values='?', dtype={'ID': str}).dropna().
449     reset_index()
450 dataset
451
452 ## Accuracy estimation using cross validation
453 ##
454 #X = dataset[dataset.columns[1:14]]
455 X = pre_process(dataset).drop(columns=['Lead'])
456 X=X[{'
457     'Year',
458     'Gross',
459     'Number of actors',
460     'Proportion of female actors',
461     'Mean Age Male',
462     'Mean Age Female',
463     'Age Lead',
464     'Age Co-Lead',
465     'Total words',

```

```

466     'Proportion of words lead',
467     'Proportion of words co-lead',
468     'Proportion of words female',
469     'Older lead'}]
470
471 Y = dataset['Lead']
472 # Split randomized data into training and validation (30% validation)
473 X_train, X_val, Y_train, Y_val = skl_ms.train_test_split(X,Y, test_size
474     =0.1)
475 # List of models
476 models = []
477 #models.append(skl_lm.LogisticRegression(solver='liblinear'))
478 #models.append(skl_da.LinearDiscriminantAnalysis())
479 models.append(skl_da.QuadraticDiscriminantAnalysis())
480
481 n_fold=20
482 accuracy = np.zeros((n_fold, len(models)))
483 model_accuracy = np.zeros(len(models))
484 cv = skl_ms.KFold(n_splits=n_fold, shuffle=True)
485
486 for i, (train_index, val_index) in enumerate(cv.split(X)):
487     X_train, X_val = X.iloc[train_index], X.iloc[val_index]
488     Y_train, Y_val = Y.iloc[train_index], Y.iloc[val_index]
489
490     for m in range(np.shape(models)[0]):
491         model = models[m]
492         model.fit(X_train, Y_train)
493         prediction = model.predict(X_val)
494         accuracy[i,m] = np.mean(prediction == Y_val)
495         model_accuracy[m] += accuracy[i,m]
496 for m in range(np.shape(models)[0]):
497     print('Model accuracy for model ' + str(models[m]) + ' is ')
498     print(model_accuracy[m]/n_fold)
499
500 plot = plt.figure(figsize=(20,5))
501 ax = plot.add_subplot(111)
502 bp = ax.boxplot(accuracy, vert=False)
503 plot.show()
504
505 ## Bootstrapping
506 ##
507 classifier = skl_da.LinearDiscriminantAnalysis()
508 features = [
509     'Number of actors',
510     'Proportion of female actors',
511     'Mean Age Male',
512     'Mean Age Female',
513     'Age Lead',
514     'Age Co-Lead',
515     'Total words',
516     'Proportion of words lead',
517     'Proportion of words co-lead',
518     'Proportion of words female',
519     'Older lead',
520 ]
521 target = 'Lead'
522 X = pre_process(dataset)
523 X=X[{'
524     'Number of actors',

```

```

525     'Proportion of female actors',
526     'Mean Age Male',
527     'Mean Age Female',
528     'Age Lead',
529     'Age Co-Lead',
530     'Total words',
531     'Proportion of words lead',
532     'Proportion of words co-lead',
533     'Proportion of words female',
534     'Older lead', 'Lead'}]
535 B = 400 # number of training sets to sample
536 accuracies = []
537 aucs = []
538 for i in range(B):
539     train, test = skl_ms.train_test_split(X, test_size=0.3)
540     QDA = fit_and_test(classifier, train, test, features, target,
541                        suppress_output=True)
542     accuracies.append(QDA.score(test[features], test[target]))
543     aucs.append(skl_met.roc_auc_score(test[target], QDA.predict_proba(test[
544         features]))[:,1]))
545
546 print('mean accuracy: ' + str(np.mean(accuracies)))
547 print(' mean auc: ' + str(np.mean(aucs)))
548
549 print(' std auc: ' + str(np.std(accuracies)))
550 print(' std auc: ' + str(np.std(aucs)))

```

## 551 References

- 552 Andreas Lindholm, Niklas Wahlström, Fredrik Lindsten, and Thomas B. Schön. *Supervised*  
553 *Machine Learning*. Pre-pub:Cambridge university Press, draft version: january 12, 2021  
554 edition, 2021.
- 555 Scikit-learn, documentation - nearest neighbors. [https://scikit-learn.org/stable/](https://scikit-learn.org/stable/modules/neighbors.html)  
556 [modules/neighbors.html](https://scikit-learn.org/stable/modules/neighbors.html), a. Accessed: 2021-02-23.
- 557 Scikit-learn, documentation - sklearn.neighbors.kneighborsclassifier. [https://scikit-learn.](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#sklearn-neighbors-kneighborsclassifier)  
558 [org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#sklearn-neighbors-kneighborsclassifier)  
559 [sklearn-neighbors-kneighborsclassifier](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#sklearn-neighbors-kneighborsclassifier), b. Accessed: 2021-02-23.