

---

# 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 (who plays the main character), can be predicted using metrics such as the year when the movie was made and the age of the lead 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 ROC/AUC.

### 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 since 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. For that reason 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.

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	$\frac{\text{Number of male actors} + \text{Number of female actors}}{\text{Number of female actors} + \text{Number of male actors}}$
N/A	Older lead	$\begin{cases} 1, \text{Age lead} > \text{Age Co-Lead} \\ 0, \text{else} \end{cases}$

Table 1: Transformations of input variables.

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.

## 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 minimizes the mean squared error (MSE) loss function. A potential problem with this approach is that there are no restrictions on the components of  $\beta$  and that can lead to overfitting, especially if  $n$  is not much larger than  $p$ . To address that issue, one can introduce regularization.

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

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

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  $\lambda$ , this method cannot be used to estimate the test error of the whole procedure. Instead, the dataset has to be split into a training set and a testing set with a specified fraction of the total data in each set. The whole procedure above is done on the training set and the test error is estimated by evaluating the performance of the model on the testing set. However, this can yield significantly different estimates of the test error since only one split into training and testing data is considered. To get a better estimate of the actual testing error, a bootstrap procedure is performed.

Since the full dataset is an iid sample from some unknown distribution, the estimated test error  $\hat{E}_{new}$  is a random variable. By repeating the whole procedure  $B$  times (i.e.  $B$  independent splits into training and testing data and subsequent fitting and cross-validation), a bootstrap sample of  $\hat{E}_{new}$  is obtained which can be used to estimate the distribution (or at least properties thereof) of  $\hat{E}_{new}$ . This is very computationally intensive but gives a much clearer view of the variability of the test error compared to just computing it for one split.

## 2.3 k-Nearest Neighbors

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

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

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

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

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

A problem with the brute force algorithm is that for each point we need to calculate the distance to every other point, which is computationally demanding for larger datasets. There are however algorithms based on the same principle as brute force search such as the ball-tree and k-d tree that are less computationally demanding which were used but since they are mostly improvements for calculation speed we do not explain these here. All three of these algorithms were tested and no significant difference in the results were noted thus the choice of algorithm was set to "auto" which chooses the best suited algorithm for a given problem, further described in the Scikit-learn documentation. [2]

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

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

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

for each of the  $k$ -nearest neighbors. This reinforces the idea that proximity of test data points to training data points ought be a good predictor. [3]

## 2.4 LDA and QDA

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

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

We estimate the *uninformative prior probability* as  $\hat{p}(y = m) = \frac{n_m}{n}$  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 value  $\mu_m$  and covariance matrix  $\Sigma_m$ . The assumption that distinguishes LDA and QDA is that for LDA we assume that  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_M$  but for QDA we make no such assumption, that is, we allow for the covariance matrices to differ. A consequence is that LDA is a special case of QDA, hence QDA is a model of higher complexity.

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

$$\hat{\mu}_m = \frac{1}{n_m} \sum_{i:y_i=m} \mathbf{x}_i, \quad (5)$$

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

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

$$\hat{\Sigma} = \frac{\sum_{m=1}^M (n_m - 1) \hat{\Sigma}_m}{\sum_{m=1}^M (n_m - 1)} = \frac{1}{n - M} \sum_{m=1}^M \sum_{i:y_i=m} (\mathbf{x}_i - \hat{\mu}_m)(\mathbf{x}_i - \hat{\mu}_m)^T. \quad (7)$$

With these estimators we may express the discriminant analysis classifier as

$$\hat{p}(y = m | \mathbf{x}) = \frac{n_m \mathcal{N}(\mathbf{x} | \hat{\mu}_m, \hat{\Sigma})}{\sum_{i=1}^M n_i \mathcal{N}(\mathbf{x} | \hat{\mu}_m, \hat{\Sigma})} \quad (8)$$

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 distribution with mean  $\mu$  and covariance matrix  $\Sigma$ .

## 3 Results

### 3.1 Logistic Regression

When comparing different models, it is important to have a baseline, or a null model to compare against. In this case, an obvious null model is the constant model that always

135 predicts the same outcome regardless of input. The best null model is the one with highest  
 136 accuracy, i.e. the constant model that predicts the most frequently occurring outcome. The  
 137 model that always predicts a male lead has an accuracy of 0.756 and is thus chosen as the  
 138 baseline.

139 For all logistic regression models fitted, the set of regularization parameters,  $\Lambda$ , consisted of  
 140 10 logarithmically spaced values between  $10^{-4}$  and  $10^4$ . This was the default value in the  
 141 methods from scikit learn and having more densely packed values did not affect the model  
 142 performance in any appreciable way. The number of folds used in cross-validation was also  
 143 10, no improvement was observed by increasing this value.

144 The model performance was measured by accuracy (1 - misclassification rate) and AUC  
 145 (area under ROC curve). In Tables 2 and 3, the accuracy and AUC are estimated using  
 146 the mean of 100 bootstrap samples in the case of LASSO regression and 400 in the case of  
 147 Ridge regression. The reason for having different sample sizes is that computing the LASSO  
 148 regression is much more computationally demanding.

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

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

149 We see that the regularization does not affect the model performance much. LASSO and  
 150 Ridge regularization perform almost identically and yield at best around 0.3% extra accuracy  
 151 but considering that the different splits of the data yielded estimated test errors in a range  
 152 from 0.8 to 0.98, we cannot reject that regularization does not matter in this case.

### 153 3.2 k-Nearest Neighbors

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

161 It is obvious that k-NN is drastically improved by transforming the data, before transforma-  
 162 tions the model performance was even outperformed by, or just slightly better than the best  
 163 null model. Weighted k-NN with the *distance* weight seemed to perform better than the  
 164 *uniform* weight both before and after the transformation, the impact seems to be 0.7 – 0.8%  
 165 extra accuracy after transformation and almost 3% extra accuracy before transformations.

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

Table 4: Accuracy and AUC for k-NN models. 70% 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

Table 5: Accuracy and AUC for k-NN models. 90% training data.

### 166 3.3 LDA and QDA

167 The given dataset was bootstrapped 400 times and both DA models were tested before and after input transformations. From the Tables 6 and 7 we can draw the conclusion that QDA

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

Input	DA 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

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

168  
169 seems to be more apt for this problem after transformations. It is unclear which method  
170 performs better before the transformation. For QDA the variance in accuracy is high which  
171 can be explained by the fact that the original inputs are close to being colinear making  $\Sigma$   
172 close to being singular which in turn results in an inaccurate matrix inversion. For example,  
173 the standard deviation of accuracy and AUC of the QDA classifier, before transformations,  
174 on 400 bootstrapped datasets with 90% training data were 0.076 and 0.081 respectively  
175 compared to 0.021 and 0.019 after transformations ceteris paribus.

176 Cross validation was carried out to estimate the accuracy using 150 folds resulting in an  
177 estimated accuracy of 0.948 using 70% training data. As can be seen in table 8 there is a  
178 noticeable variance in accuracy for the different training sets suggesting that there might be  
179 outliers in the data that the model has problems accounting for. Increasing the number of  
180 folds also increases the minimum accuracy that can be found in the corresponding box-plot,  
181 as to be expected. Also using cross validation to compare the effects of input transformations  
182 we see a 0.090 increase in accuracy using the transformed inputs compared to the original

183 inputs. Adding the variables 'Years' and 'Gross' back into the inputs we see a decrease in  
 184 accuracy of 0.007 hence they are left out.

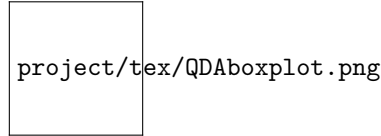


Table 8: Accuracy estimation using cross validation with 150 folds.

## 185 4 Conclusions

186 It seems like  $k$ -NN was the worst performing of the models. Worth noting is the drastic  
 187 effect that data transformations has on the  $k$ -NN model, performing at the same level as  
 188 the baseline model before any transformations. The best results for the  $k$ -NN were a mean  
 189 accuracy score of 0.882 and a mean AUC Score of 0.901 which is found in table 5.

190 Logistic regression slightly outperformed  $k$ -NN. Looking at the best performing setup we see  
 191 that logistic regression had a mean accuracy score of 0.898 and a mean AUC score of 0.923  
 192 which is seen in table 3.

193 LDA had a mean accuracy score of 0.903 which is only a 0.5% increase from the logistic  
 194 regression case, considering that both the results from logistic regression and LDA had some  
 195 variance we cannot reject that LDA and logistic regression performance is similar.

196 QDA however had the best performance among all the models performing at best with a  
 197 mean accuracy score of 0.942 and a mean AUC score of 0.977(!). Almost a 4% increase  
 198 in accuracy compared to the second best method. One problem with the QDA model was  
 199 however the outliers seen in the boxplot of 8. Similar outliers were found in the case of  $k$ -NN  
 200 and logistic regression these could probably be explained by the occurrence of a "bad split"  
 201 where for example many movies deviating from the average movie end up in the test portion  
 202 resulting in a bad performance. However the model performed well enough overall such that  
 203 QDA was chosen for production.

### 204 1. Do men or women dominate speaking roles in Hollywood?

205 Based on the data it seems like males are in the lead. In 75.6% of the cases the  
 206 lead is male which was found by looking at the baseline model. Also the mean for  
 207 "Proportion of words female" was calculated to 34,6% indicating that Hollywood  
 208 movies consist of almost 65% male speaking.

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

210 Since none of the models we used seem to be effected by removing the year variable  
 211 we cannot say that this is the case. However looking at the correlation in the data  
 212 we see that there is a positive correlation between "Years" and each of the variables  
 213 that show female activity in movies (i.e. "Number of words female", "Number of  
 214 female actors" and "Mean Age Female"), indicating that there may be some changes  
 215 in the gender balance over time.

### 216 3. Do films in which men do more speaking make a lot more money than films in which 217 women speak more?

218 Since none of the models seem to be effected by removing the gross variable we  
 219 cannot say that gross is a good indicator for determining whether there is more  
 220 male speaking than female in a given movie. That is films with more male speaking  
 221 cannot be said to make more money than a film with more female speaking.

## 222 5 Feature Importance

223 In order to determine which of the features Year, Gross or Number of female actors is most  
224 important, we fitted models excluding one or more of the features. Using a logistic regression  
225 model with LASSO (L1) regularization, we found that the model performance in terms of  
226 prediction accuracy was completely unaffected by removing either or both of the features  
227 Year and Gross. On the other hand, any model that excluded the variable Proportion of  
228 words female (which contains all information about how much male and female actors speak),  
229 had its performance reduced drastically. As seen in Table 3, the model including all features  
230 had an accuracy of 90%. This dropped to 80% when removing the Proportion of words  
231 female. Comparing this to the null accuracy of 75.6%, we conclude that the Proportion of  
232 words female is a very important feature in the model.

233 The same results hold for both QDA and k-NN as well, the models are completely unaffected  
234 by removing either Year or Gross (or both), while suffering 8% accuracy loss for k-NN and  
235 12% loss for QDA. Further, not only accuracy is affected but AUC is affected in the same  
236 way, showing that Year and Gross are more or less redundant features in the model, while  
237 Proportion of words female is incredibly important for predicting whether the lead is male  
238 or female.



## 239 Appendix A: Code

240 Transformations of input variables and various common functions.

```

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

```

296         print(skl_met.classification_report(test[target], classifier.
297             predict(test[features])))
298     return classifier
299
300 print('Null accuracy: ' + str(max([np.mean(data[target]), 1 - np.mean(data[
301     target]))]))

```

## 302 Logistic Regression

```

303 trainRatio = config['Train Ratio'][0]
304 seed = config['Random Seed'][0]
305 train, test = skl_ms.train_test_split(data, train_size=trainRatio)
306
307 featureSet1 = [
308     'Year',
309     'Gross',
310     'Number of actors',
311     'Proportion of female actors',
312     'Mean Age Male',
313     'Mean Age Female',
314     'Age Lead',
315     'Age Co-Lead',
316     'Total words',
317     'Proportion of words lead',
318     'Proportion of words co-lead',
319     'Ratio words co-lead lead',
320     'Proportion of words female',
321     'Older lead'
322 ]
323
324 features = featureSet1.copy()
325 #features.remove('Proportion of words female')
326 #features.remove('Year')
327 #features.remove('Gross')
328 #features = ['Proportion of words lead']
329 target = 'Lead'
330
331 # No regularization
332
333 B = 100
334 accuracies = []
335 aucs = []
336 for i in range(B):
337     train, test = skl_ms.train_test_split(data, train_size=trainRatio)
338     logReg = fit_and_test(skl_lm.LogisticRegression(penalty='none',
339         solver='newton-cg'), train, test, features, target,
340         suppress_output=True)
341     accuracies.append(logReg.score(test[features], test[target]))
342     aucs.append(skl_met.roc_auc_score(test[target], logReg.predict_proba
343         (test[features])[:,1]))
344
345 # LASSO
346
347 B = 100
348 accuracies = []
349 aucs = []
350 for i in range(B):
351     train, test = skl_ms.train_test_split(data, train_size=trainRatio)

```

```

352     logRegLasso = fit_and_test(sk_lml.LogisticRegressionCV(Cs=10, cv=10,
353                             penalty='l1', solver='liblinear', n_jobs=10), train, test,
354                             features, target, suppress_output=True)
355     accuracies.append(logRegLasso.score(test[features], test[target]))
356     aucs.append(sk_lmet.roc_auc_score(test[target], logRegLasso.
357                                     predict_proba(test[features])[:,1]))
358
359 # Ridge
360
361 B = 400
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(sk_lml.LogisticRegressionCV(Cs=10, cv=10,
367                             penalty='l2', 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(sk_lmet.roc_auc_score(test[target], logRegLasso.
371                                     predict_proba(test[features])[:,1]))

```

## 372 References

- 373 [1] A. Lindholm, N. Wahlström, F. Lindsten, and T. B. Schön, *Supervised Machine Learning*.  
374 Pre-pub:Cambridge university Press, draft version: january 12, 2021 ed., 2021.
- 375 [2] “Scikit-learn, documentation - nearest neighbors.” [https://scikit-learn.org/stable/](https://scikit-learn.org/stable/modules/neighbors.html)  
376 [modules/neighbors.html](https://scikit-learn.org/stable/modules/neighbors.html). Accessed: 2021-02-23.
- 377 [3] “Scikit-learn, documentation - sklearn.neighbors.kneighborsclassifier.” [https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#sklearn-neighbors-kneighborsclassifier)  
378 [KNeighborsClassifier.html#sklearn-neighbors-kneighborsclassifier](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#sklearn-neighbors-kneighborsclassifier). Ac-  
379 cessed: 2021-02-23.