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

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

1

## 1  Introduction

## 2  Methods

We have chosen to focus on approaches using logistic regression, k-NN and LDA/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.

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

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 where 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

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

where $\beta \in \mathbb{R}^p$ and $g$ is the link function. In the case of logistic regression, $Y|X \sim Ber(p)$ 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

$Y|X \sim Ber(p)$, we get $E(Y|X) = p = g^{-1}(X \cdot \beta)$. 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 $\Lambda = \{\lambda_1, \ldots, \lambda_k\}$. 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 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

### 2.4 LDA and QDA

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

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

3

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

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

### 3.1.1 k-Nearest Neighbors

## 3.2 LDA and QDA

# 4 Conclusions

# 5 Feature Importance