
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 2.1.

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

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

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

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

30 $Y|X \sim \text{Ber}(p)$, we get $E(Y|X) = p = g^{-1}(X \cdot \beta)$. In other words, $P(Y = 1|X = x) = g^{-1}(x \cdot \beta)$,
31 which we can use to predict Y given data x .

32 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
33 minimizes the mean squared error (MSE) loss function. A potential problem with this approach is
34 that there are no restrictions on the components of β and that can lead to overfitting, especially if n is
35 not much larger than p . To address that issue, one can introduce regularization.

36 In general, regularization is done by adding a penalizing term to the loss function that restricts β
37 in some way. If $L(\beta; x_i, y_i)$ is the loss function before regularization, we instead consider the new
38 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
39 penalizing function and λ is a hyper-parameter that can be tuned. The two most common forms of
40 regularization is LASSO and Ridge regression.

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

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

49 Since cross-validation is used to estimate the hyper-parameter λ , this method cannot be used to
50 estimate the test error of the whole procedure. Instead, the dataset has to be split into a training set
51 and a testing set with a specified fraction of the total data in each set. The whole procedure above is
52 done on the training set and the test error is estimated by evaluating the performance of the model on
53 the testing set. However, this can yield significantly different estimates of the test error since only
54 one split into training and testing data is considered. To get a better estimate of the actual testing
55 error, a bootstrap procedure is performed.

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

62 2.3 k-Nearest Neighbors

63 2.4 LDA and QDA

64 3 Results

65 3.1 Logistic Regression

66 3.1.1 k-Nearest Neighbors

67 3.2 LDA and QDA

68 4 Conclusions

69 5 Feature Importance