
Do (wo)men talk too much in films?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1

2 1 Introduction

3 2 Methods

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

6 In order to make the methods as comparable as possible, we have used a common set of transforma-
7 tions of the input variables for all tested methods.

8 2.1 Input transformations

9 In the given dataset, there are columns for the total number of words spoken as well as the number of
10 words spoken by the lead, the co-lead etc. This could present a problem since if we compare a movie
11 where the lead says 10 out of 100 total words and another movie where the lead says 100 out of 1000
12 words, most models would think that the lead speaks more in the second movie and miss the fact that
13 the *proportion* of words spoken by the lead is the same. For that reason we have transformed several
14 input variables to express a proportion instead of absolute numbers. We also believe it might be
15 important to have a dummy variable indicating if the lead or the co-lead is oldest. All transformations
16 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	$\frac{\text{Number of female actors}}{\text{Number of female actors} + \text{Number of male actors}}$
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.

17 Note that when determining 'Proportion of words female', this should only measure the words spoken
 18 by non-lead female actors so we have to subtract the lead's contribution to the total number of words.

19 The column 'Number of male actors' was dropped since all necessary information in this column is
 20 contained in 'Proportion of female actors'.

21 In order to improve regularization and k-NN, all numerical input variables (after transformation)
 22 where centered and scaled by their standard deviation. This results in a dataset where every numerical
 23 column contains almost all data in the interval $[-3, 3]$ (in the limit, $\approx 99.7\%$ of the data should be in
 24 this interval), with higher density closer to 0.

25 2.2 Logistic Regression

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

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

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

44 2.3 k-Nearest Neighbors

45 2.4 LDA and QDA

46 3 Results

47 3.1 Logistic Regression

48 3.1.1 k-Nearest Neighbors

49 3.2 LDA and QDA

50 4 Conclusions

51 5 Feature Importance