

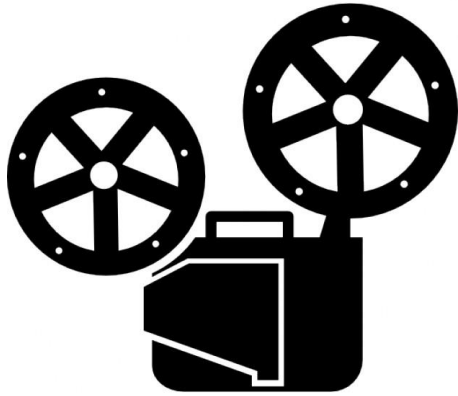
*Using Crowded-source based feature from social
media and Conventional features to predict the
movies popularity*

Paper reproduce and other attempts

*Chong Hu(ch3467)
Wenjie Chen(wc2685)
Haoran Zhang(hz2619)*

Dataset description

Movies Dataset



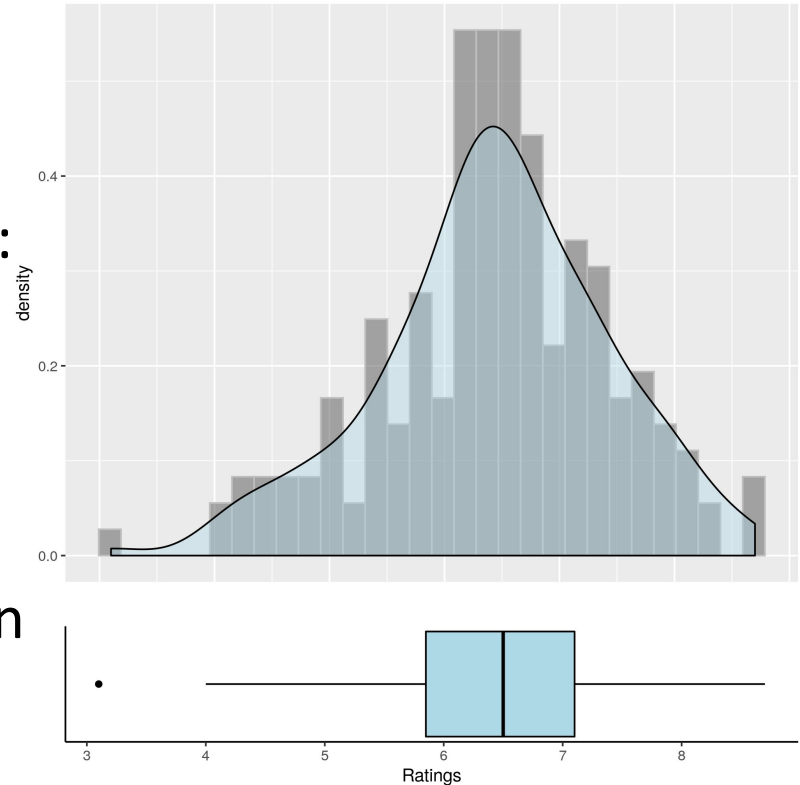
Predicting the success of movies has been of interest to economists and investors as well as predictive analysts. A number of attributes such as **cast**, **genre**, **budget**, **production house**, **PG rating** affect the popularity of a movie. Social media such as Twitter, YouTube etc. are major platforms where people can share their **views about the movies**. In this project, we collect all these features to make a prediction of movies' ratings using machine learning techniques.

Dataset Description

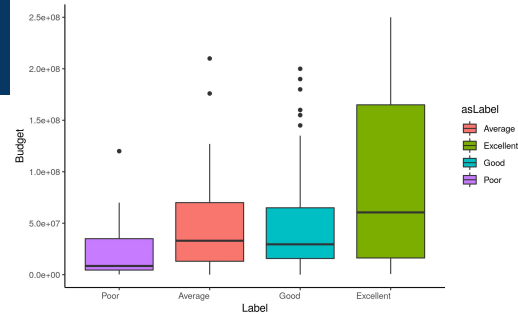
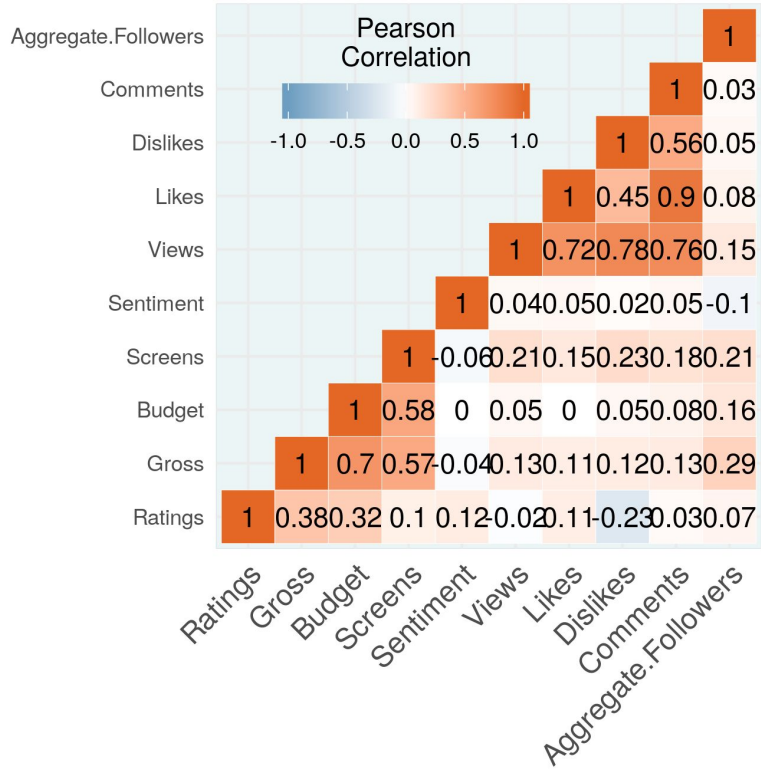
Resource: IMDB, Youtube and Twitter

11 features are divided into two groups:
Conventional Features (5 types)&
Social Media Features (6 types)
To Predict **Rating**

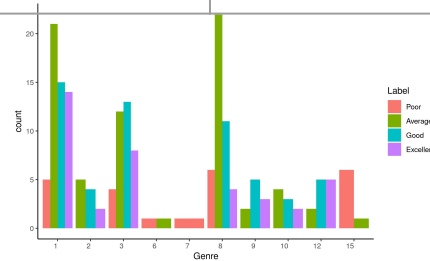
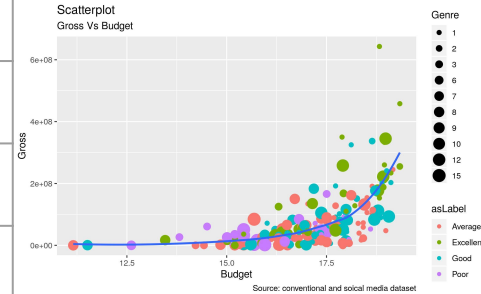
Movie Ratings distribution can be seen
from the graph on the right



Dataset Description



Ratings	Label
0-5.2	Poor
5.2-6.4	Average
6.4-7.2	Good
7.2-10	Excellent



Methods

- ❑ Linear Regression
- ❑ Decision Tree (J48)
- ❑ Random Forest
- ❑ Support Vector Machine
- ❑ Naïve Bayesian
- ❑ LDA, QDA
- ❑ Artificial Neural Network

Paper detail and reproduce

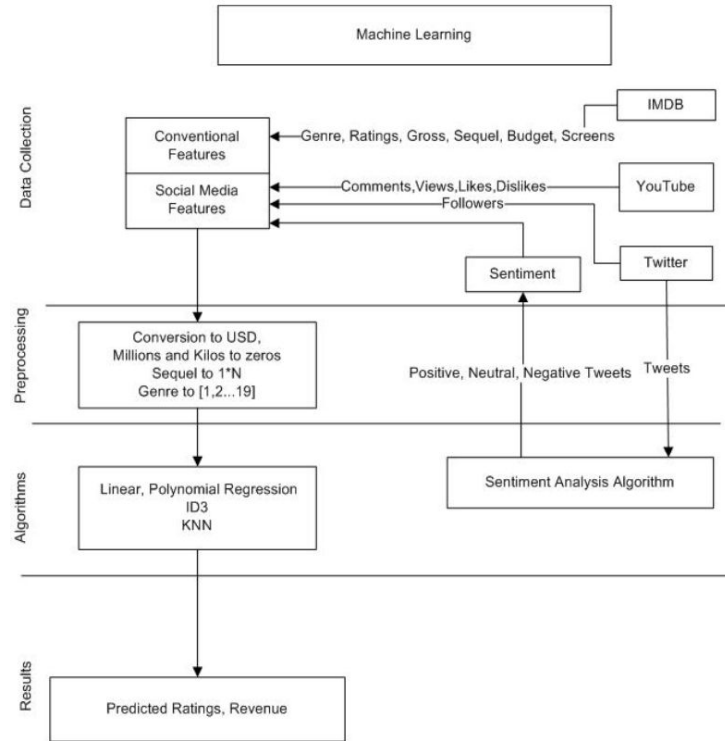
Section Subtitle Section Subtitle

First, values of Ratings are predicted using all other attributes except Gross Income (as gross income is not available before release). As Rating is a continuous numeric attributes, we have performed Linear Regression in order to predict the values.

Second, apply linear regression.

$$Accuracy_1 = \frac{\text{Number of Movies with exact prediction of Rating}}{\text{Total Number of Predictions}}$$

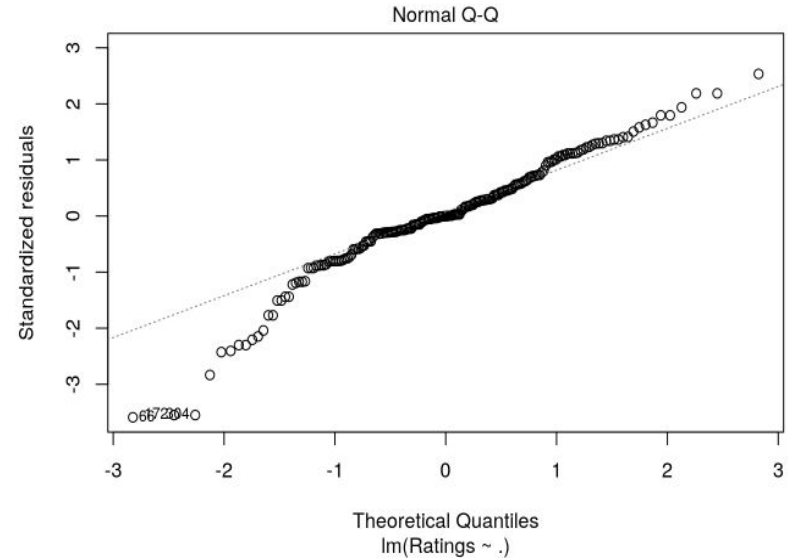
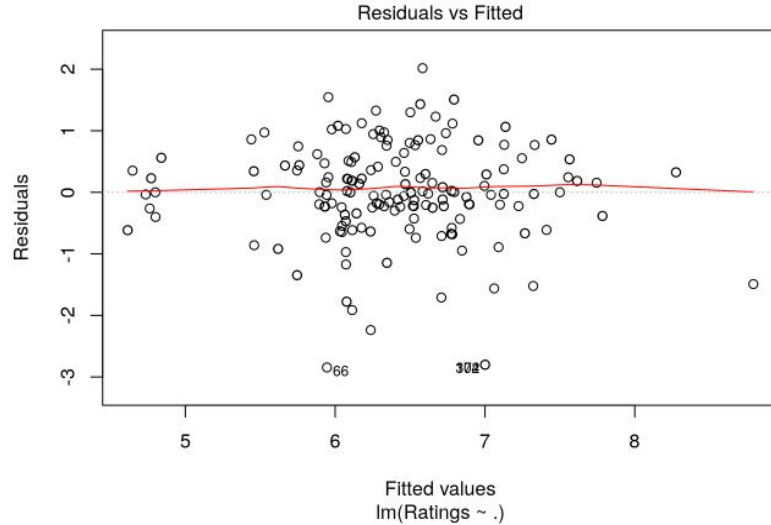
$$Accuracy_2 = \frac{\text{Number of Movies with approx prediction of Rating}}{\text{Total Number of Predictions}}$$



Linear Regression

conventional features

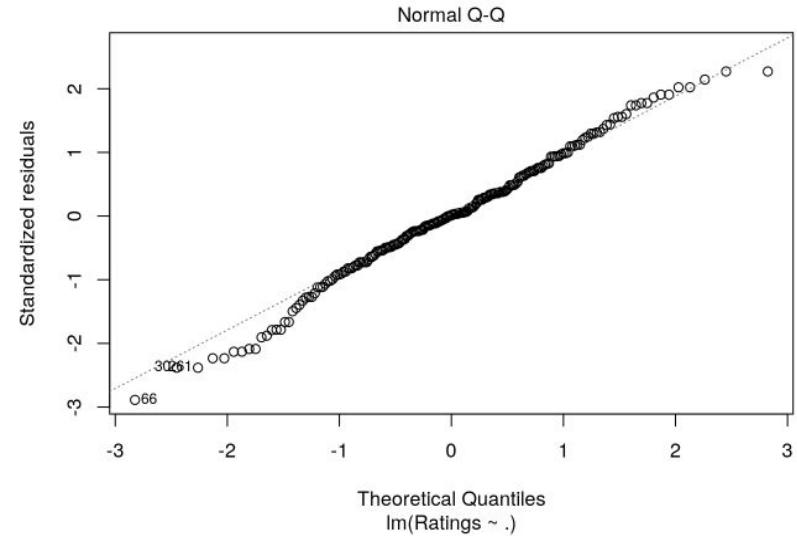
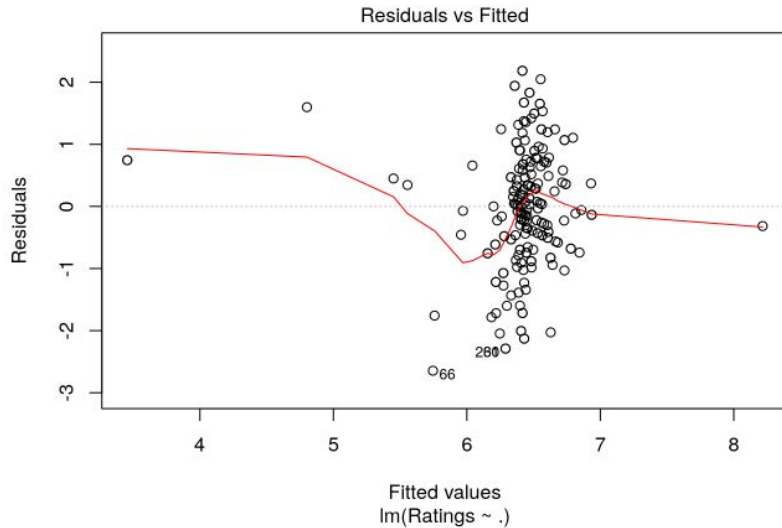
formula = Ratings ~ ., data = conventional_train.df



Linear Regression

social media

formula = Ratings ~ ., data = socialmedia_train.df



Linear Regression

convention

soft accuracy: 78.57%

classification accuracy: 35.71%

MSE: 0.7018

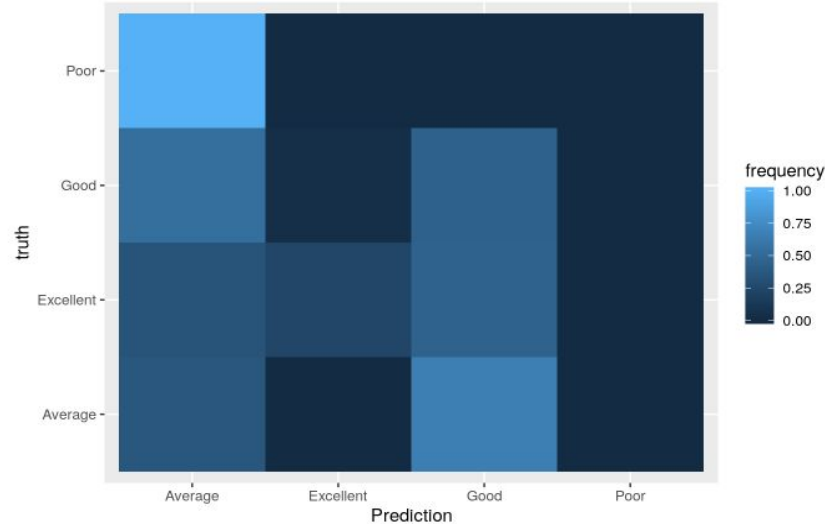
social media

soft accuracy: 75.0%

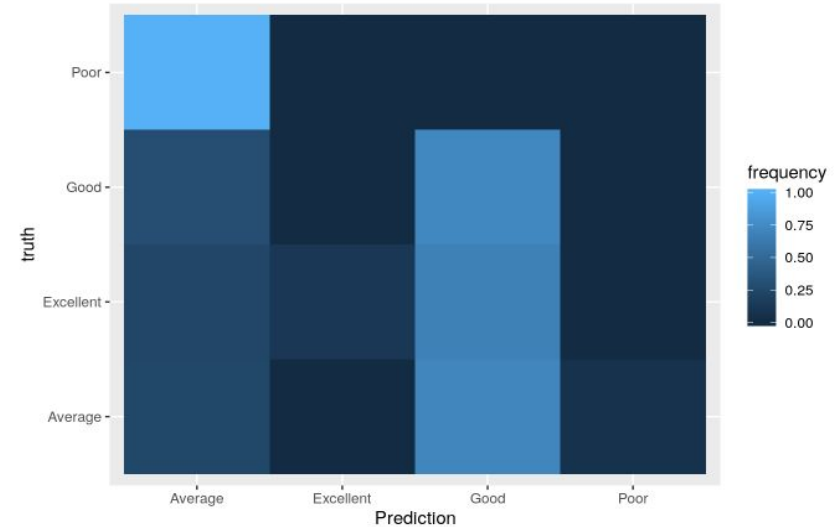
classification accuracy: 44.64%

MSE: 0.7167

Confusion Matrix of linear regression for conventional features



Confusion Matrix of linear regression for social media features



Generalised Linear Regression

convention

use Gamma (link = "log") as link function

social media

soft accuracy: 76.78%

classification accuracy: 34.71%

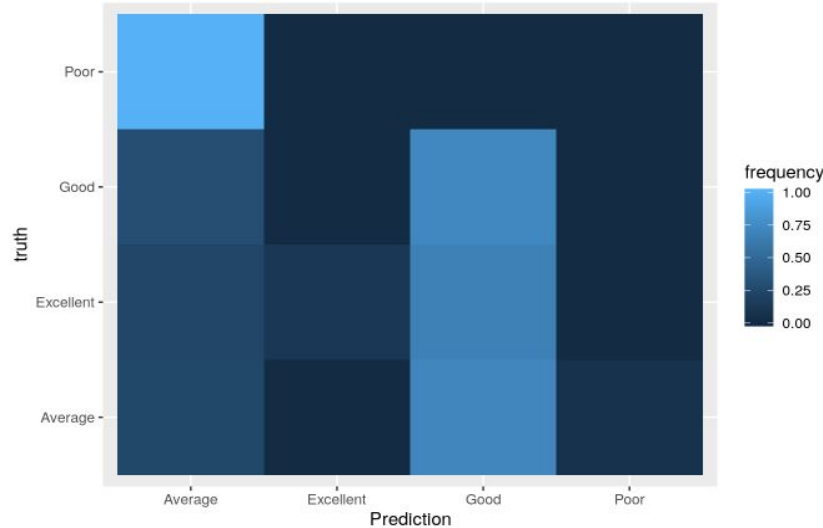
MSE: 0.7018

soft accuracy: 73.21%

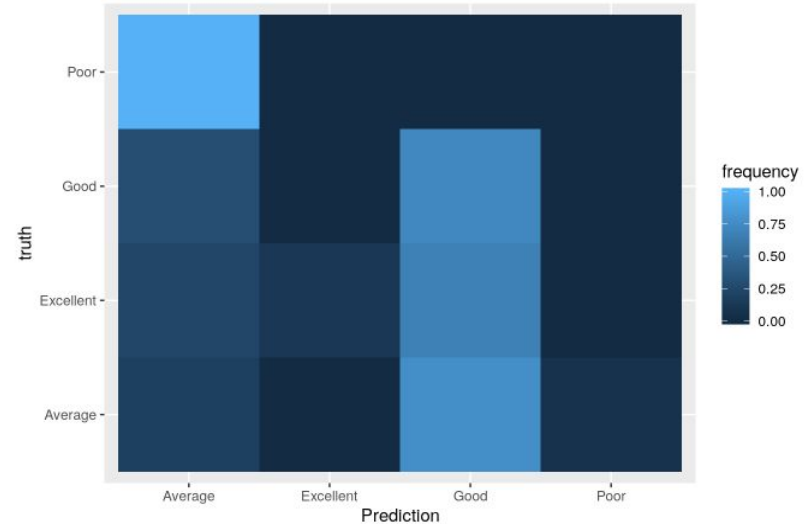
classification accuracy: 42.85%

MSE: 0.7157

Confusion Matrix of linear regression for social media features



Confusion Matrix of linear regression for social media features



Generalised Additive Model

use 'Gamma (link = "inverse")' as link function

convention

soft accuracy: 73.21%

classification accuracy: 39.28%

MSE: 0.7079

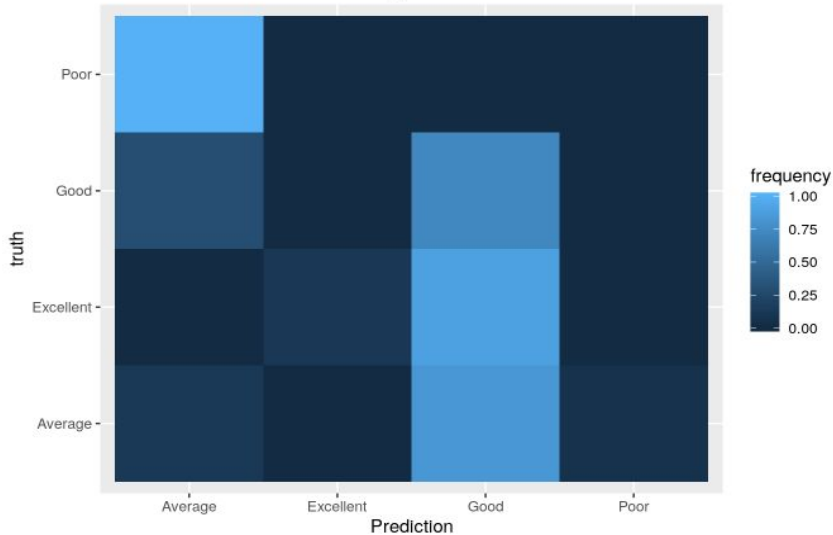
social media

soft accuracy: 70.56%

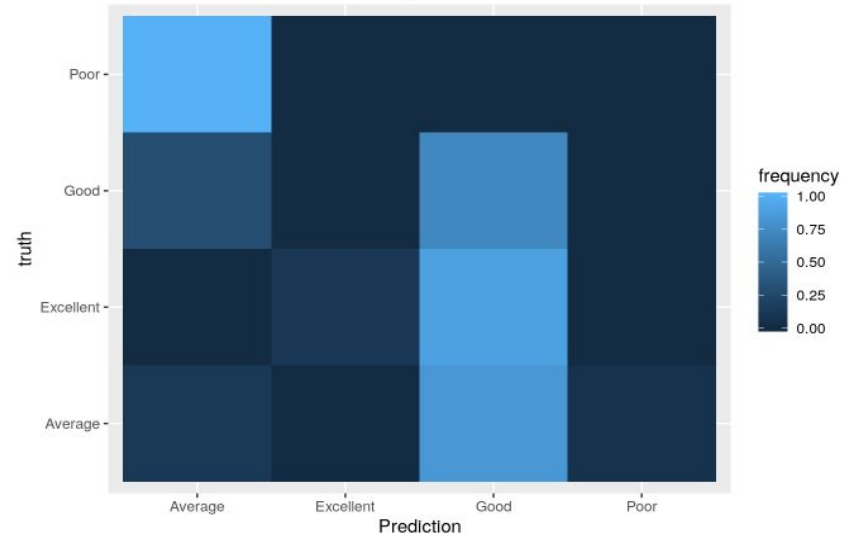
classification accuracy: 41.07%

MSE: 0.7321

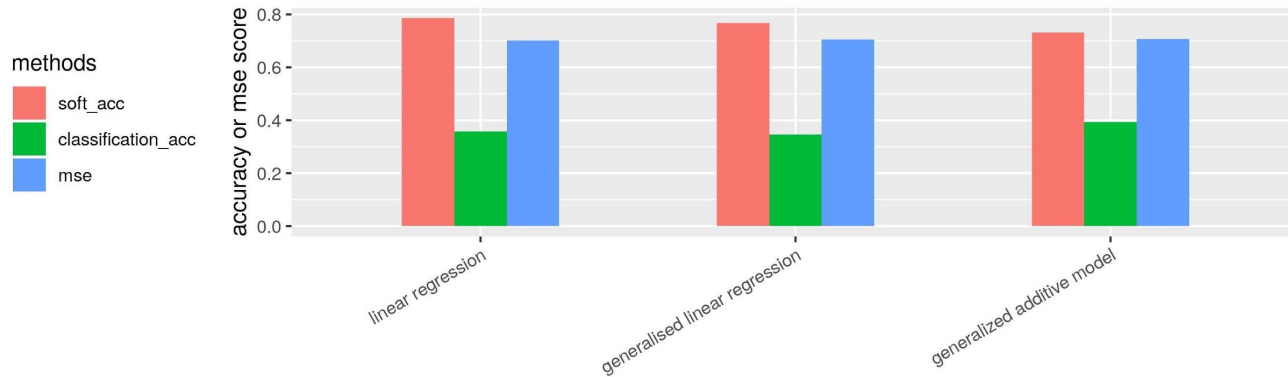
Confusion Matrix of linear regression for social media features



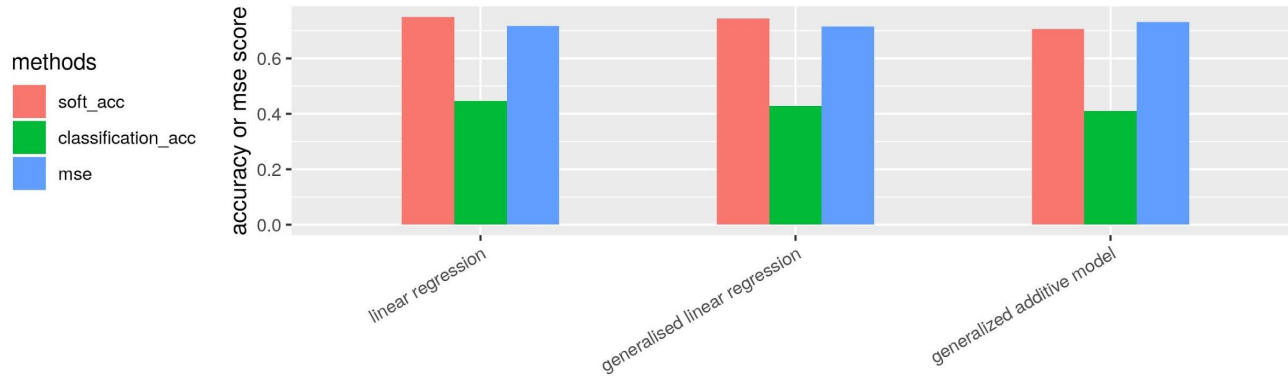
Confusion Matrix of linear regression for social media features



Linear Regression



different methods evaluations on conventional featurea



different methods evaluations on social media feature

$$\text{Accuracy} = \frac{\text{Number of Movies with correct prediction of BandofRating}}{\text{Total Number of Predictions}}$$

$$\text{Info}(D) = \text{Entropy}(D) = - \sum_j p(j|D) \log p(j|D)$$

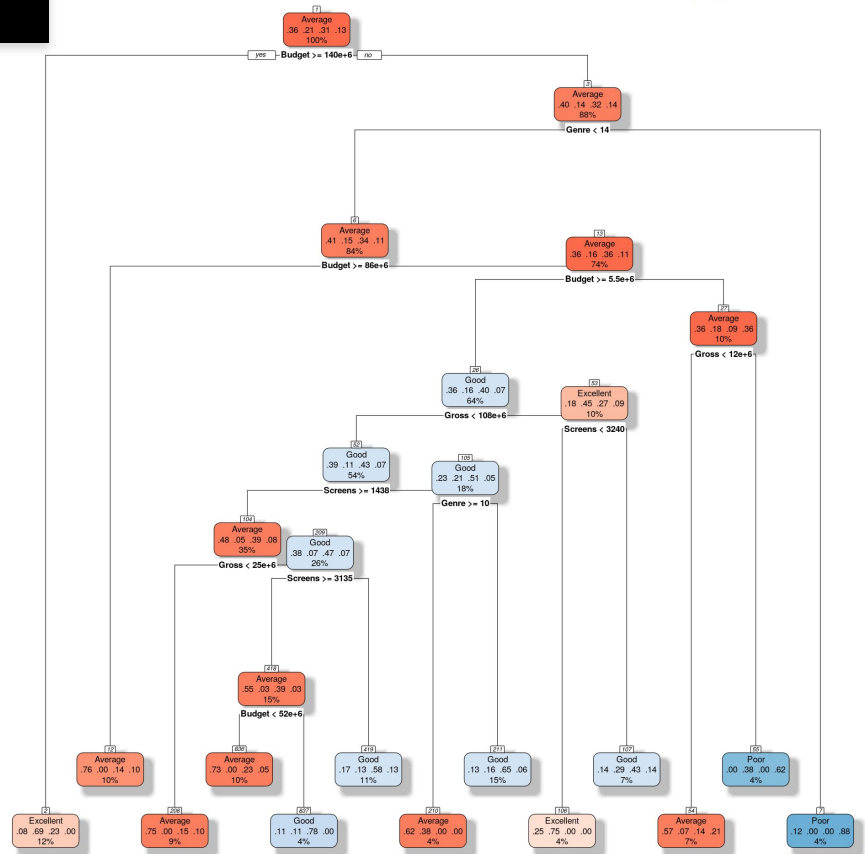
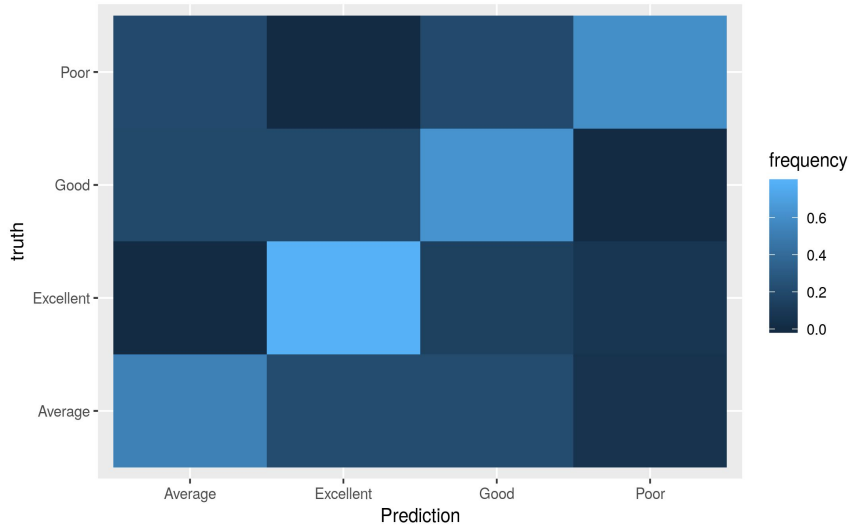
$$\text{Info}_A(D) = \sum_{i=1}^v \frac{n_i}{n} \text{Info}(D_i)$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

J48 Tree

conventional
Accuracy: 62.96%

Confusion Matrix of decision tree J48 for conventional features

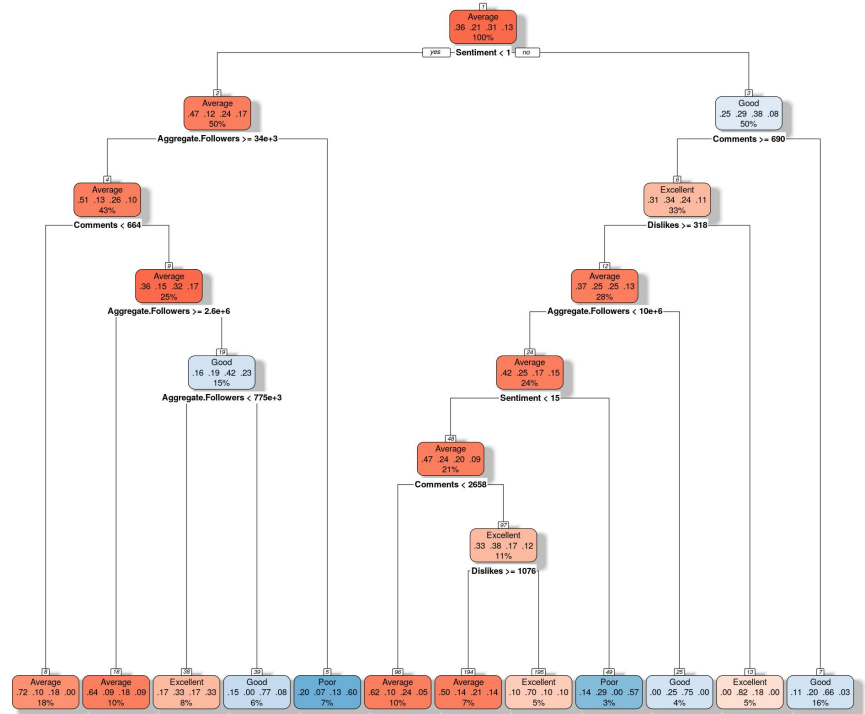
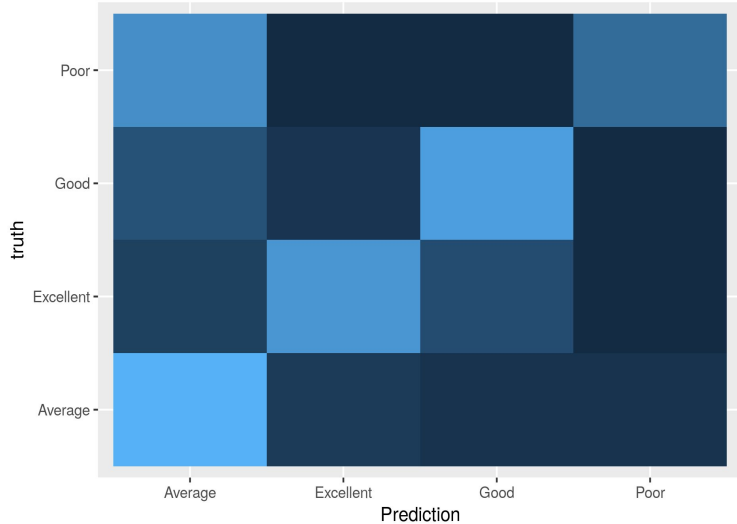


J48 Tree

erage
cellent
Good
Poor

social media
Accuracy:68.52%

Confusion Matrix of decision tree J48 for social media features

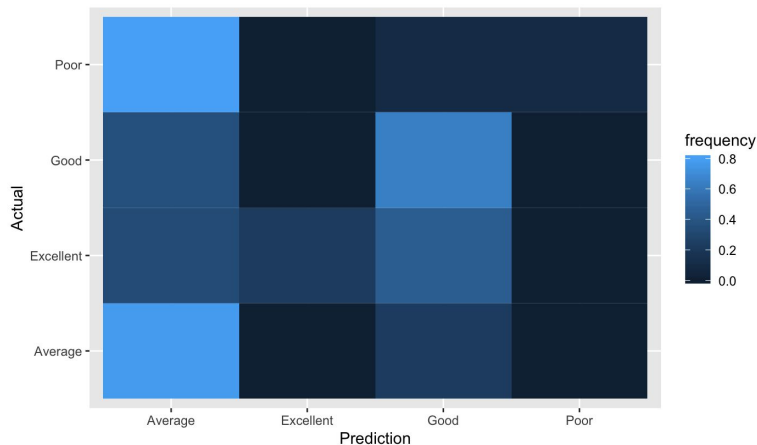


Other Techniques

Support Vector Machine

Accuracy:0.4359

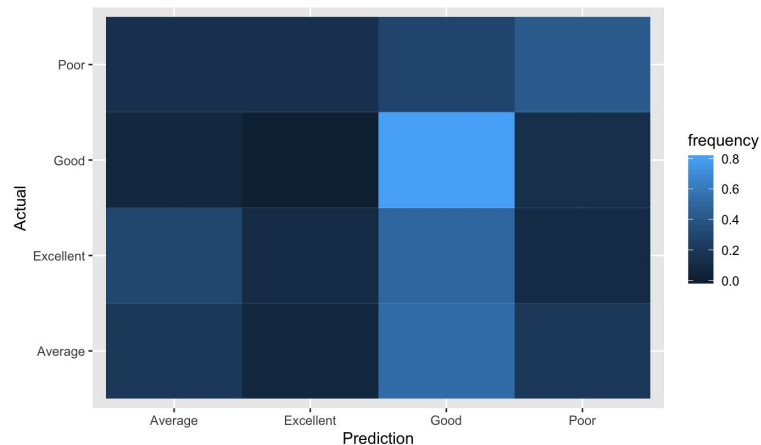
Confusion Matrix of Support Vector Machine Conventional Feature



Kernel: radial

Accuracy:0.3846

Confusion Matrix of NaiveBayesian social media feature



Kernel: polynomial

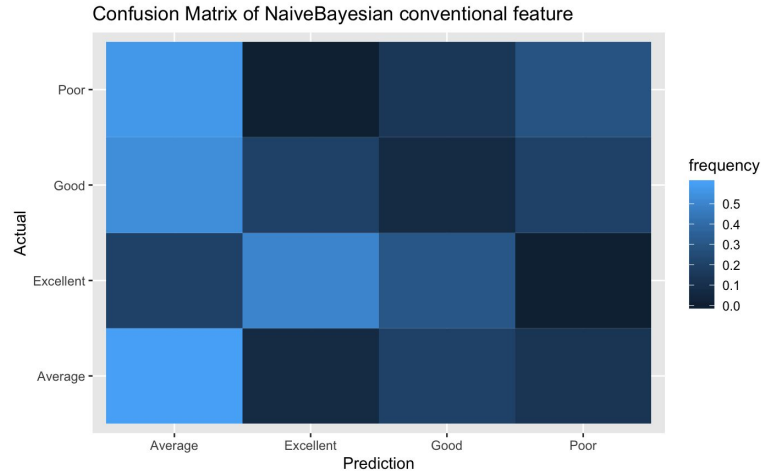
gamma: 1

coefficient: 7

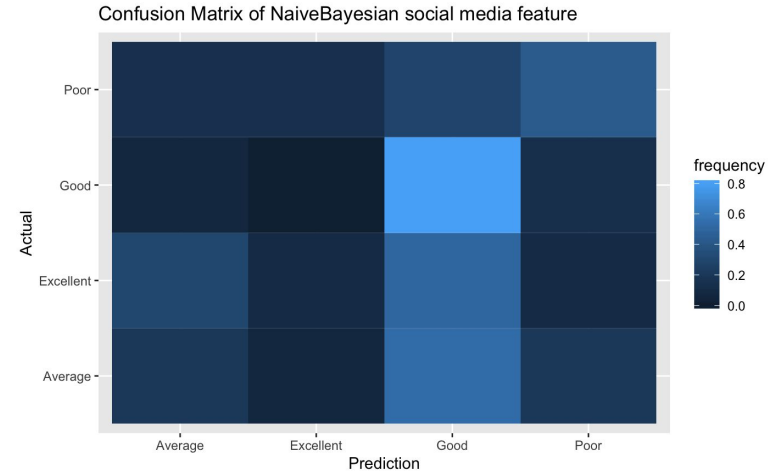
degree: 3

Naive Bayesian

Accuracy:0.3617

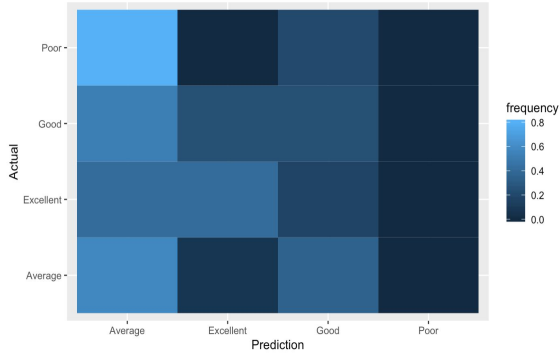


Accuracy:0.4043

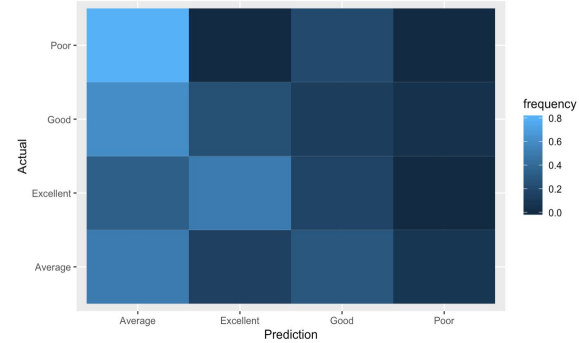


LDA & QDA

Confusion Matrix of LDA Conventional Feature

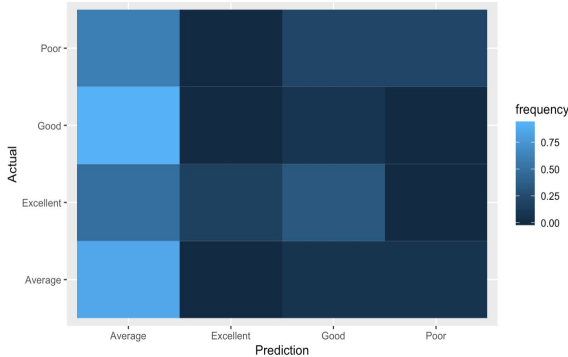


Confusion Matrix of QDA Conventional Feature

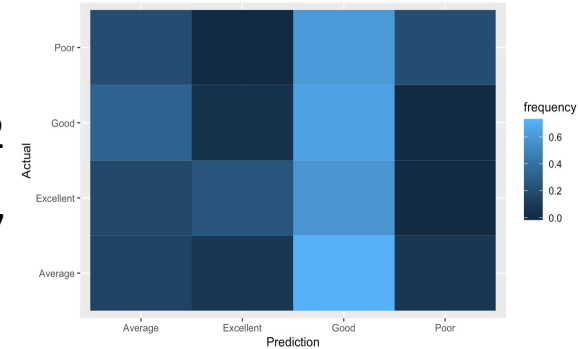


Poor Performance
|
Curse of High Dimension

Confusion Matrix of LDA Social Media Feature



Confusion Matrix of QDA Social Media Feature



	LDA	QDA
Conventional	0.3035	0.3392
Social Media	0.3928	0.2857

Artificial Neural Network

Architecture :

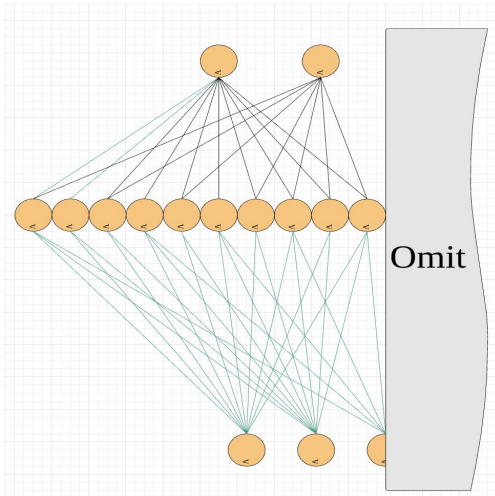
Number of units in the hidden layer = 20,

Initial random weights=0.2,

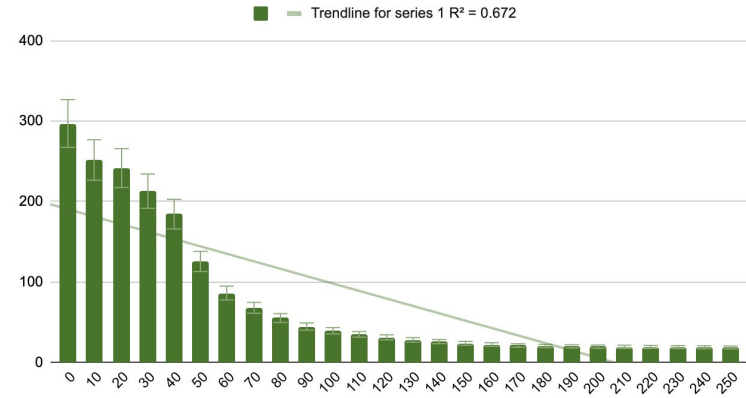
Weight decay=5e-4,

Maximum number of iterations = 250,

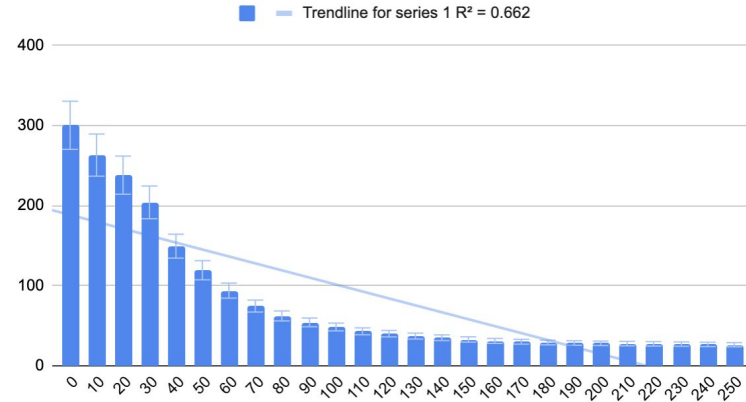
Skip=1.



Train_loss_Conventional



Train_loss_Social

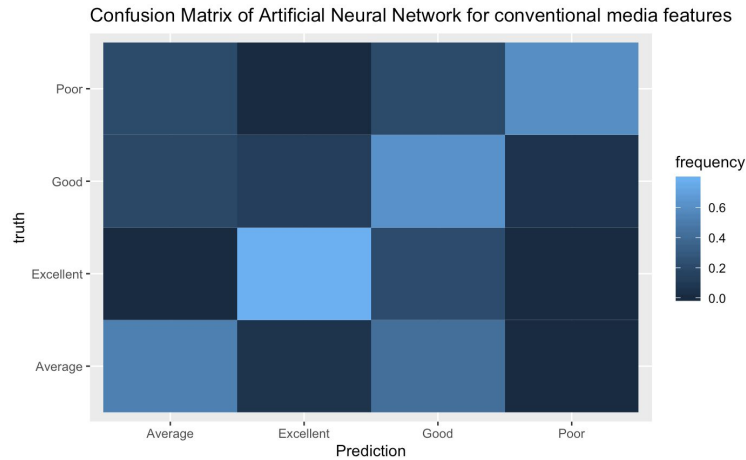


Artificial Neural Network

Test accuracy for conventional media feature and social media feature respectively

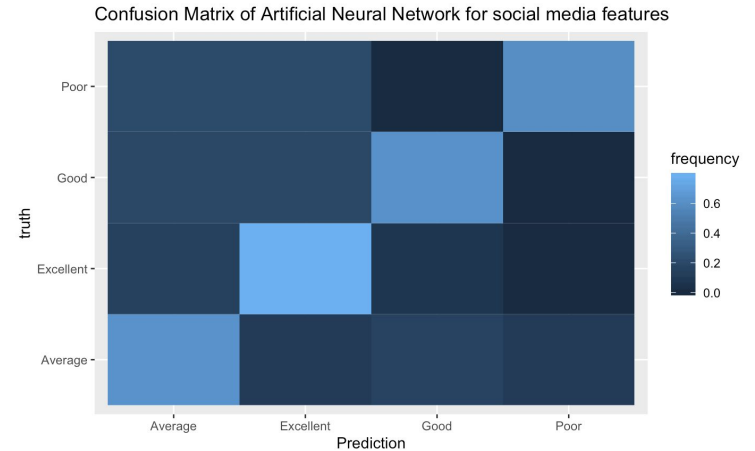
Convention:

accuracy: 66.67%



Social Media

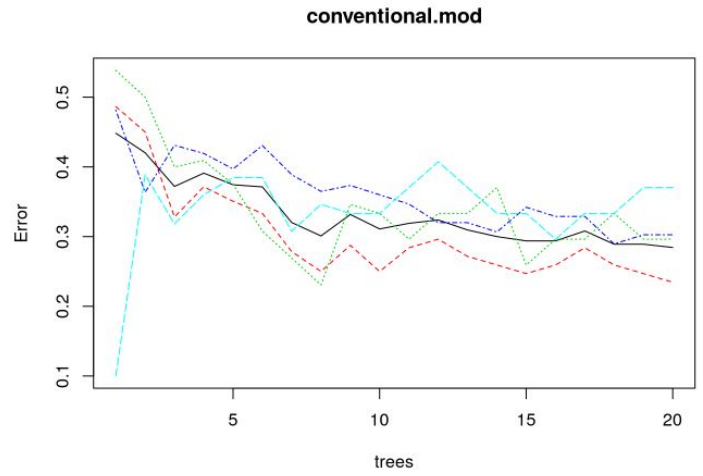
accuracy: 62.96%



Random Forest

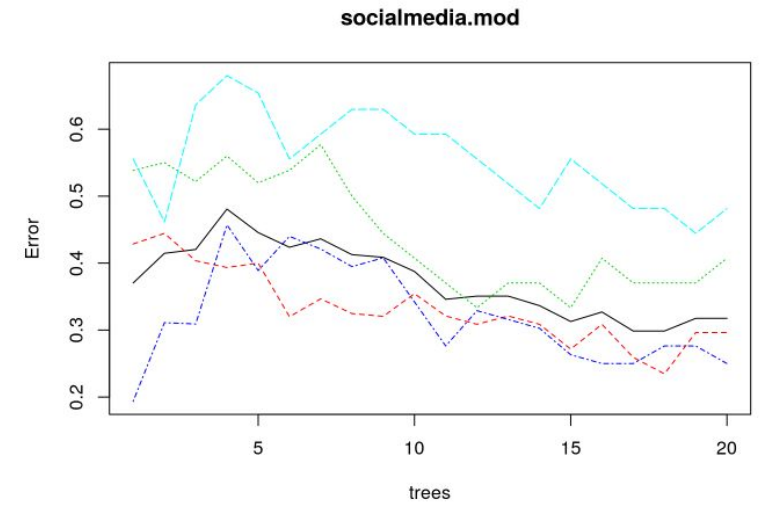
Convention:

training error for 20 subtrees



Social Media

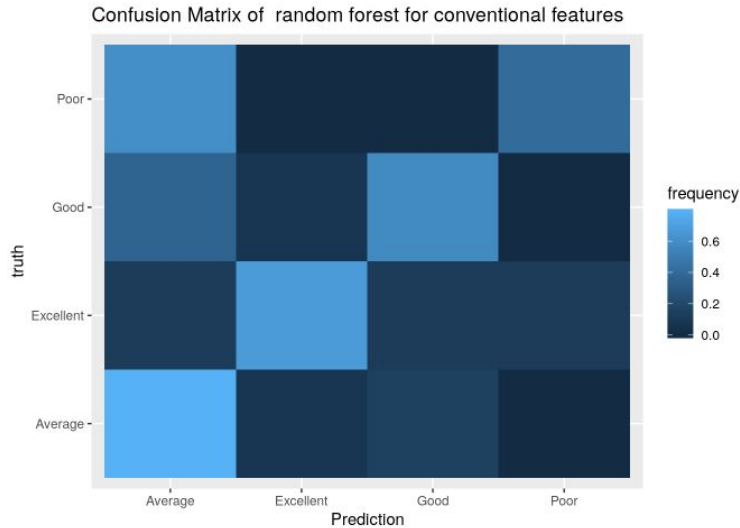
training error for 20 subtrees



Random Forest

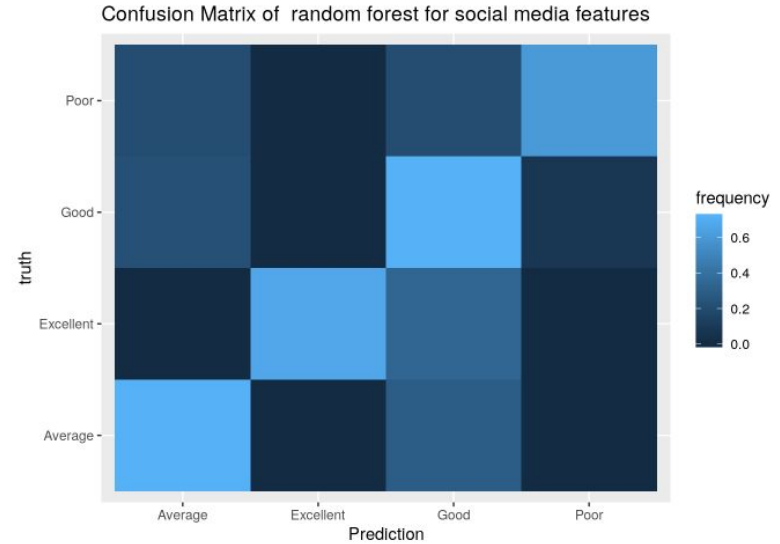
Convention:

accuracy: 62.5%

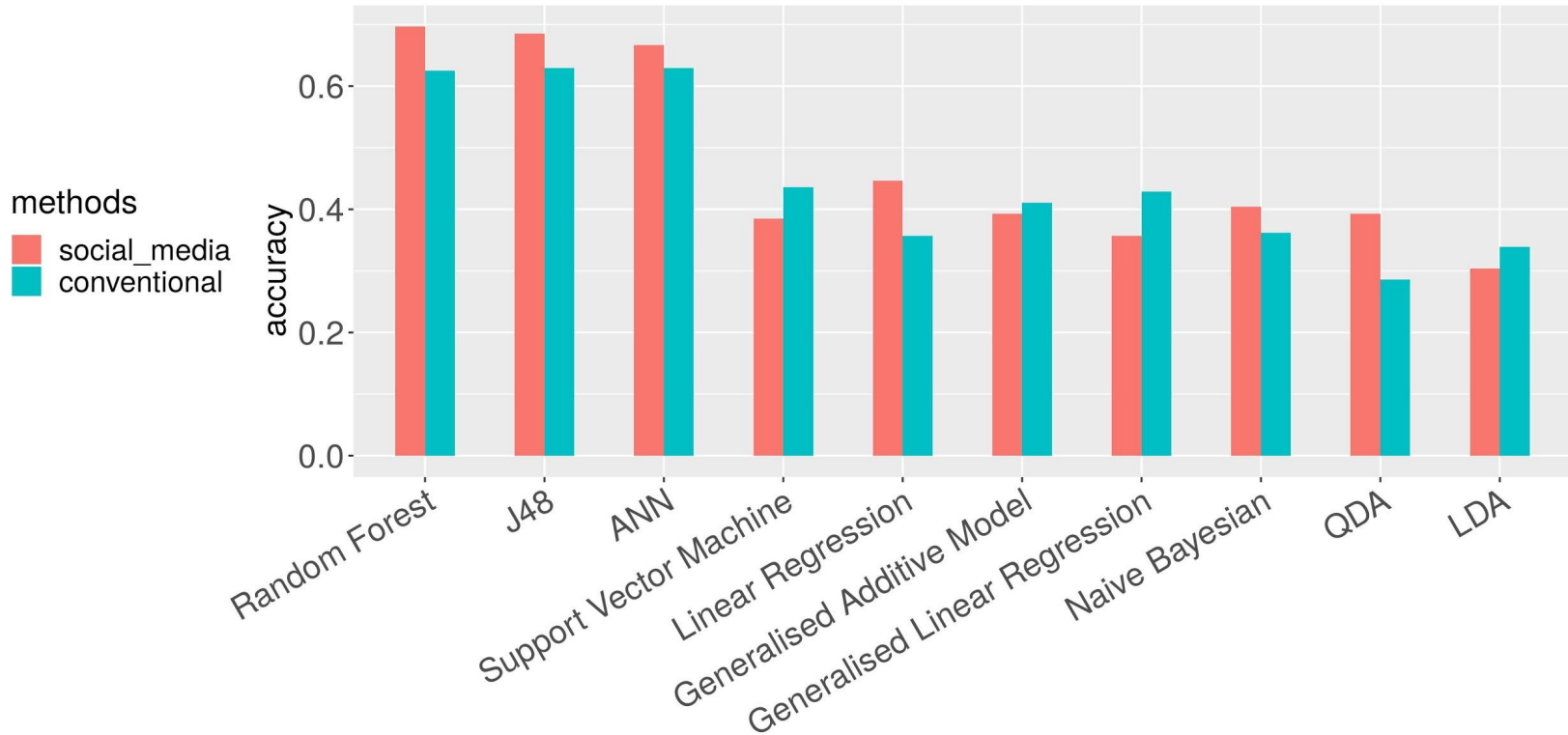


Social Media

accuracy: 69.64%



Conclusion and Discussion



Different methods evaluations on conventional features and social media features

An aerial photograph of New York City, showing Central Park on the left, the Hudson River on the right, and the dense urban landscape in between. The sky is clear and blue.

Thank you

TRANSCENDING DISCIPLINES, TRANSFORMING LIVES