

December 13,2021

## PREDICTIVE MODEL FOR ONLINE NEWS POPULARITY

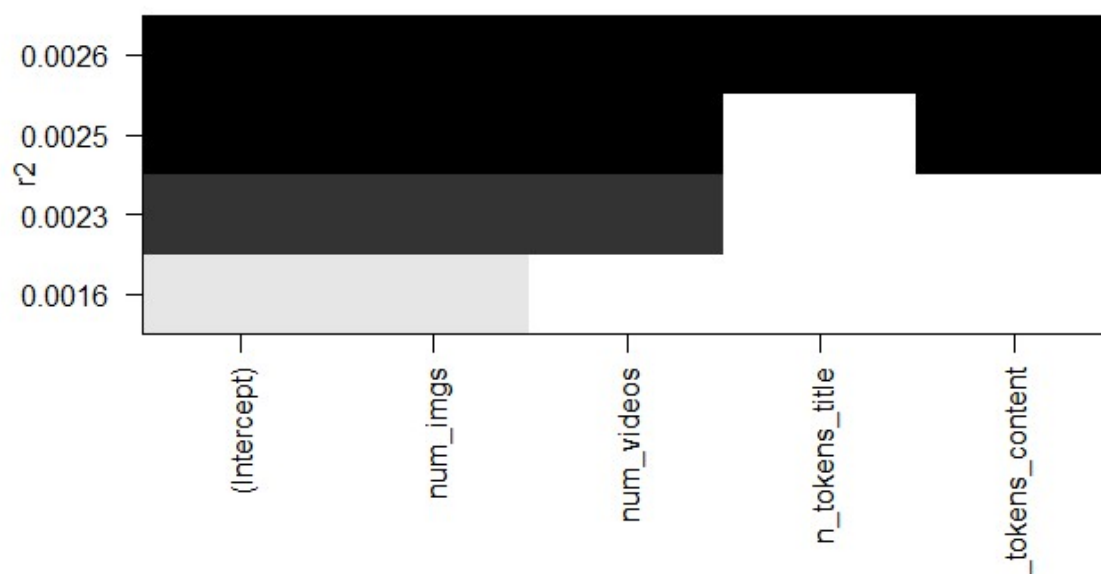
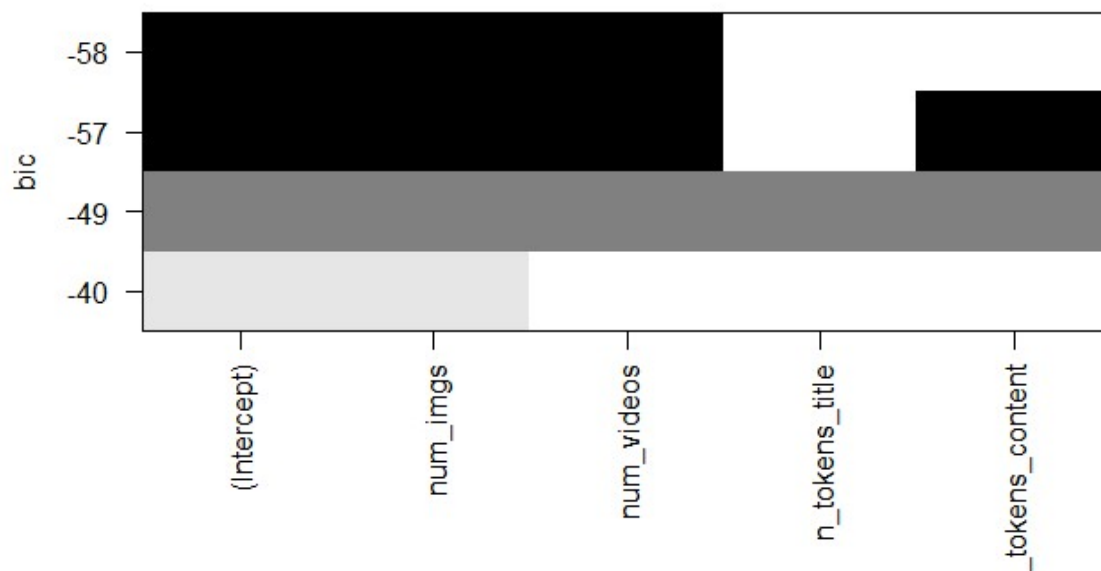
The Online News Popularity dataset was analyzed . It contains 61 characteristics of 39797 articles shared on Mashable.com .Some of the characteristics measured were the number of images , number of videos, title and article length and the number of shares each article received. This data set had 61 variables and 39797 samples. Processing the entire dataset led to memory issues and severe lagging in the code so the data set was reduced and the number of variables was decreased to 5. The number of images , number of videos, title and article length and the number of shares each article received were chosen since those values were easily measurable by the average person . The purpose of the predictive model was to predict the number of article shares using the article characteristics chosen.

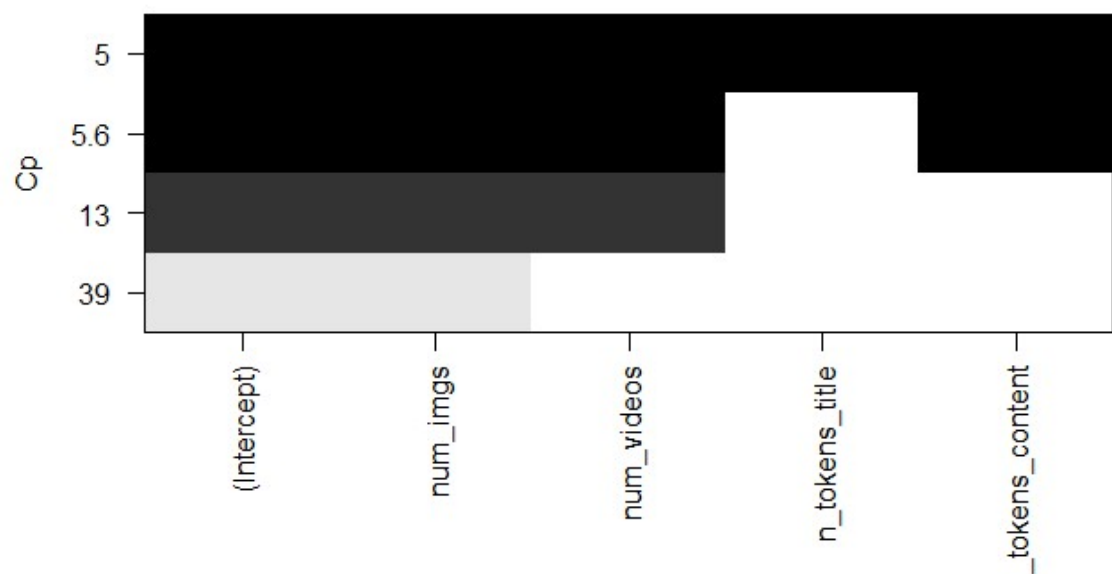
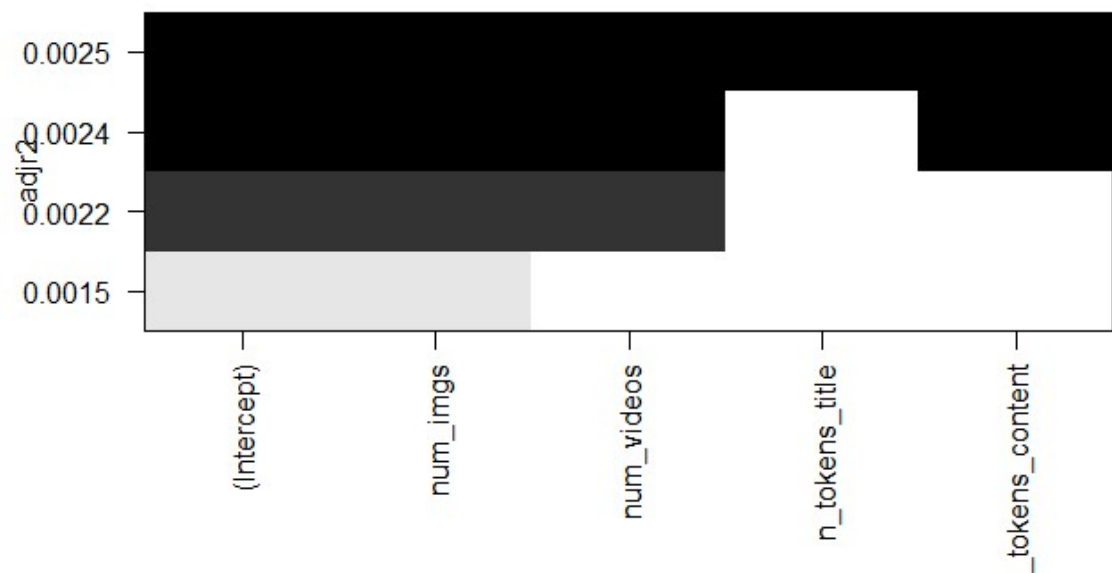
The distribution of these variables were plotted ,summary tables were made and a correlation matrix was also created. This matrix showed that all of the variables had a low correlation and similar correlation with shares so the four variables chosen were not worse predictors than the other variables. To test the best prediction model , best subset selection and forward and backward stepwise selection was implemented. 10-fold Cross Validation was also performed on each selection technique. Prior to this, cross validation was also performed on models that contained only one variable and a model with all variables .

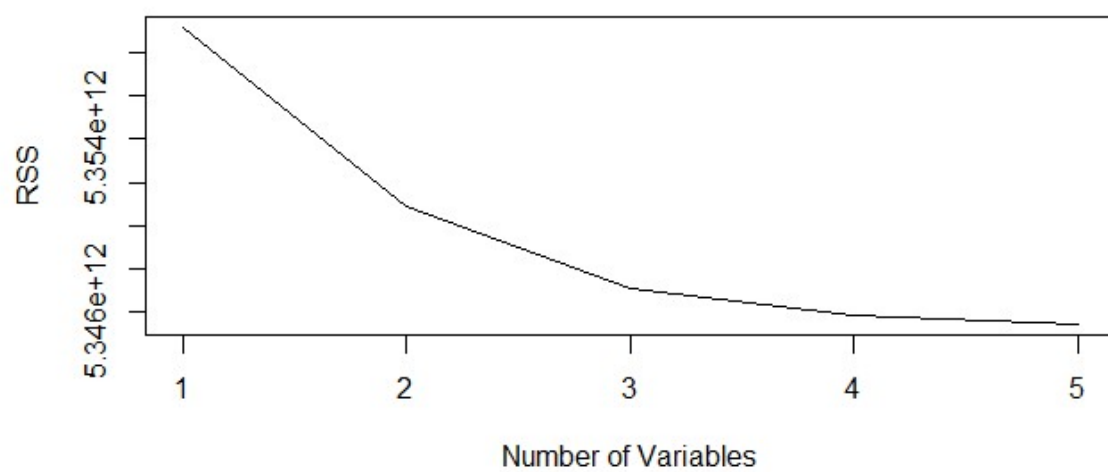
To choose the best models, linear regression was performed on each variable individually then on a model that included all four variables. From this analysis, it was predicted that the best model was a four variable model since this model had the lowest MSE. However, the testing was

not sufficient to conclude that the four variable model was the best choice. Forward and backward selection were also performed to choose the best prediction model for shares . Forward and Backward selection also determined that the four variable model was the best model to predict shares. The  $R^2$ , adjusted  $R^2$ ,  $C_p$ , BIC and RSS values were the best for the four variable model. The four variable model chosen was :Shares= 2761.5554047 + 65.9358207 num\_imgs + 80.4838476 num\_videos + 44.2298988 n\_tokens\_title + -0.4141284 n\_tokens\_content.

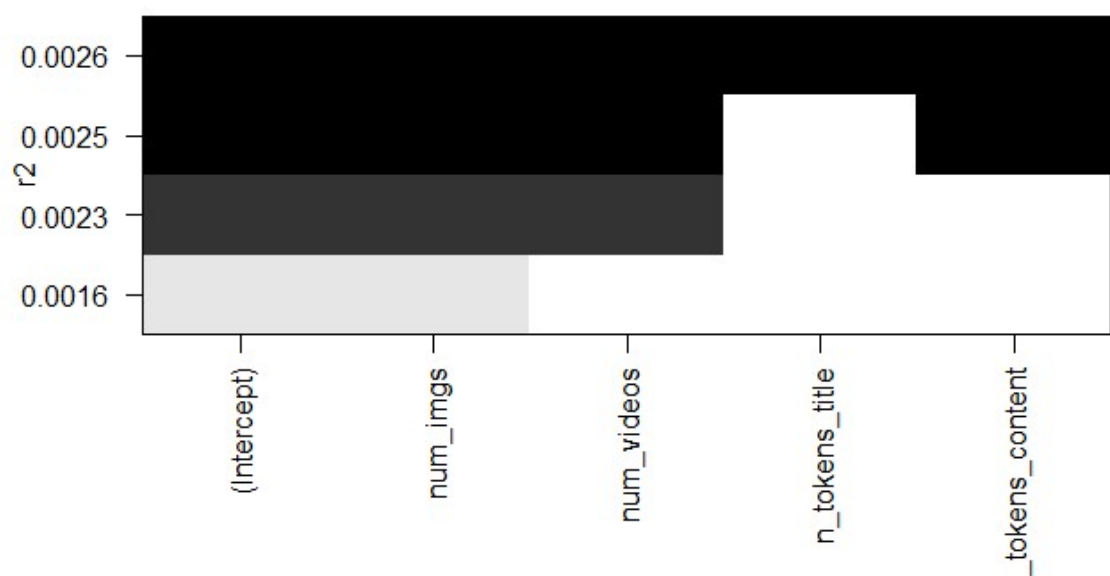
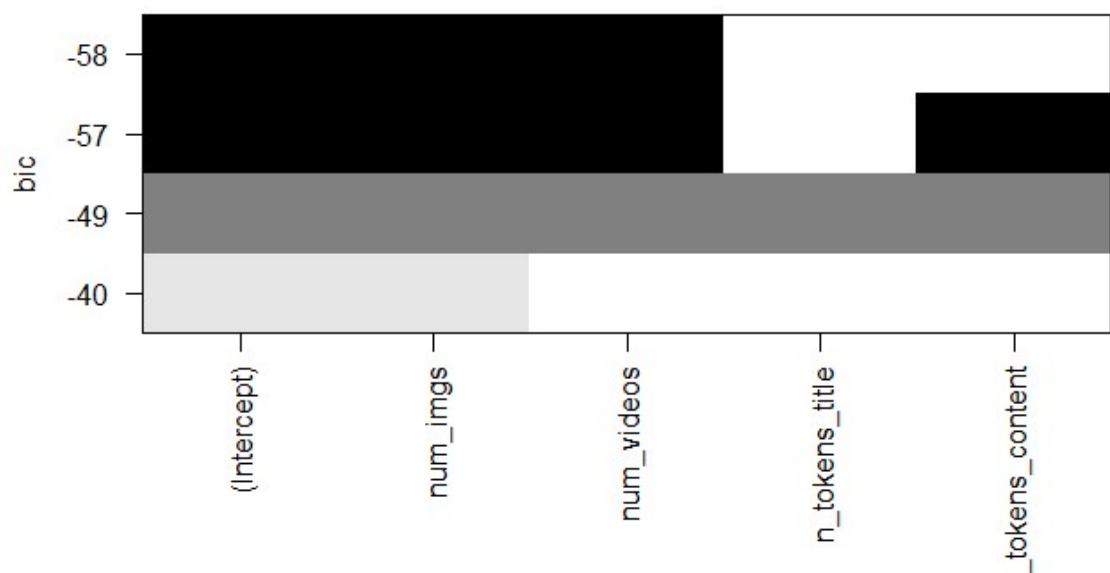
## BACKWARD

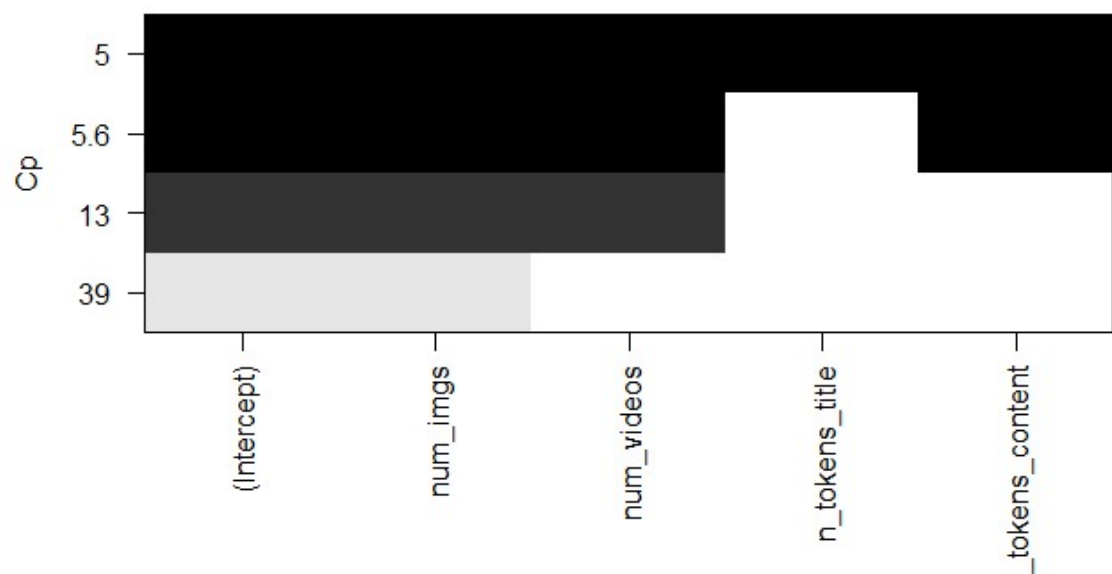
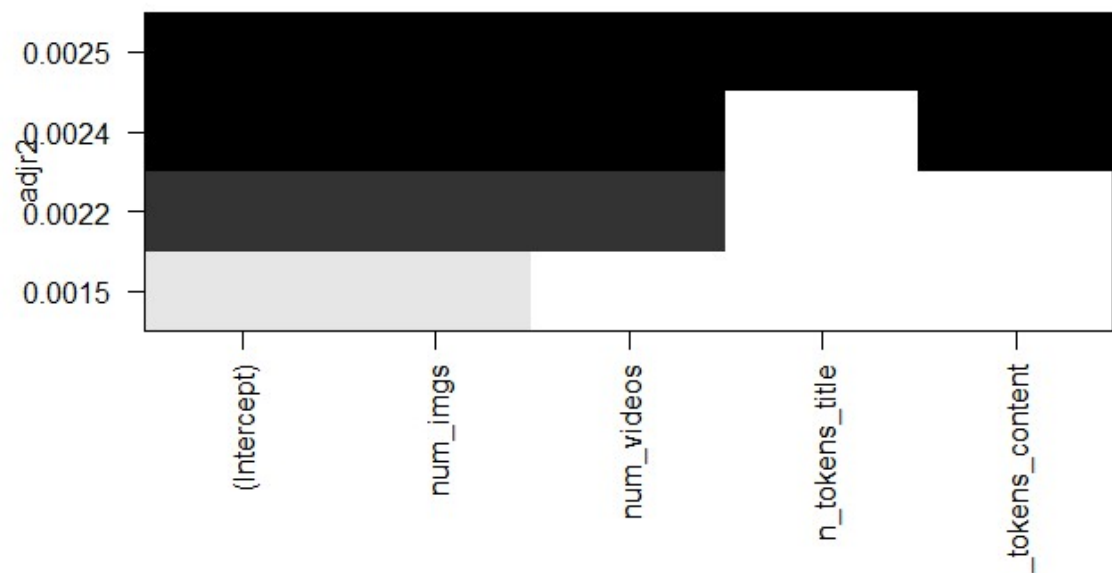


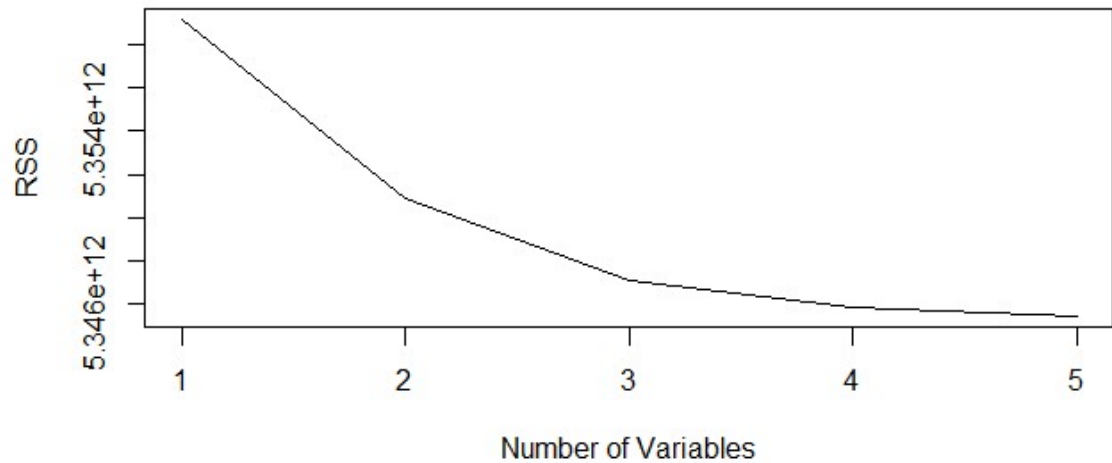




## FORWARD

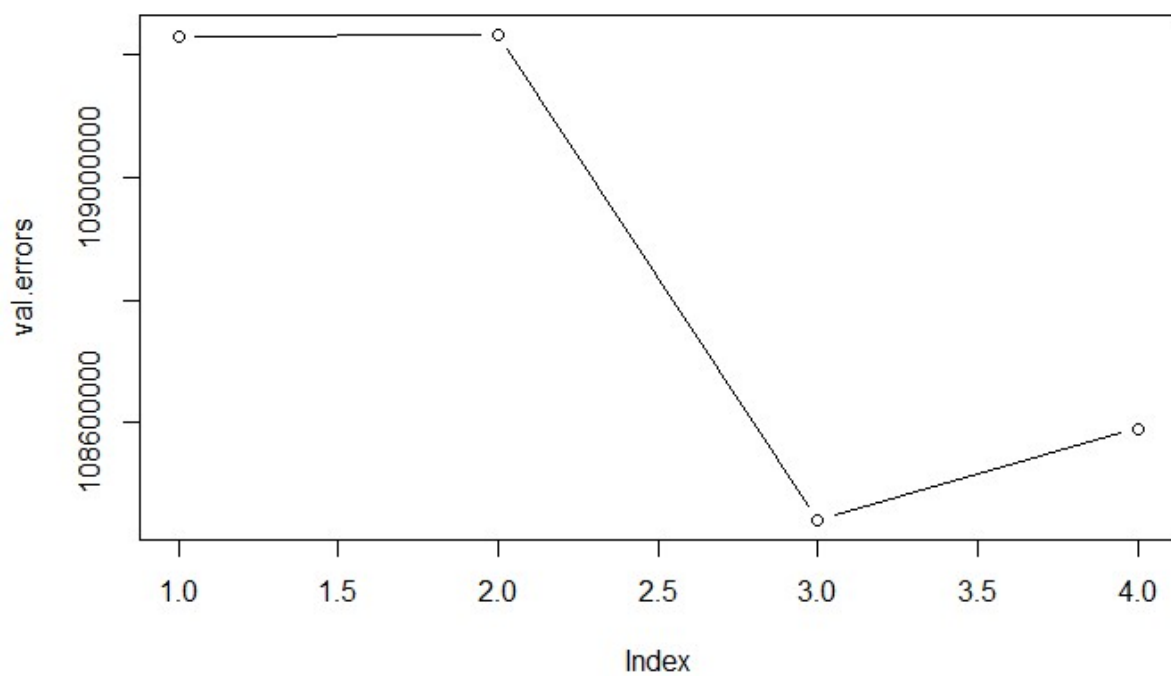
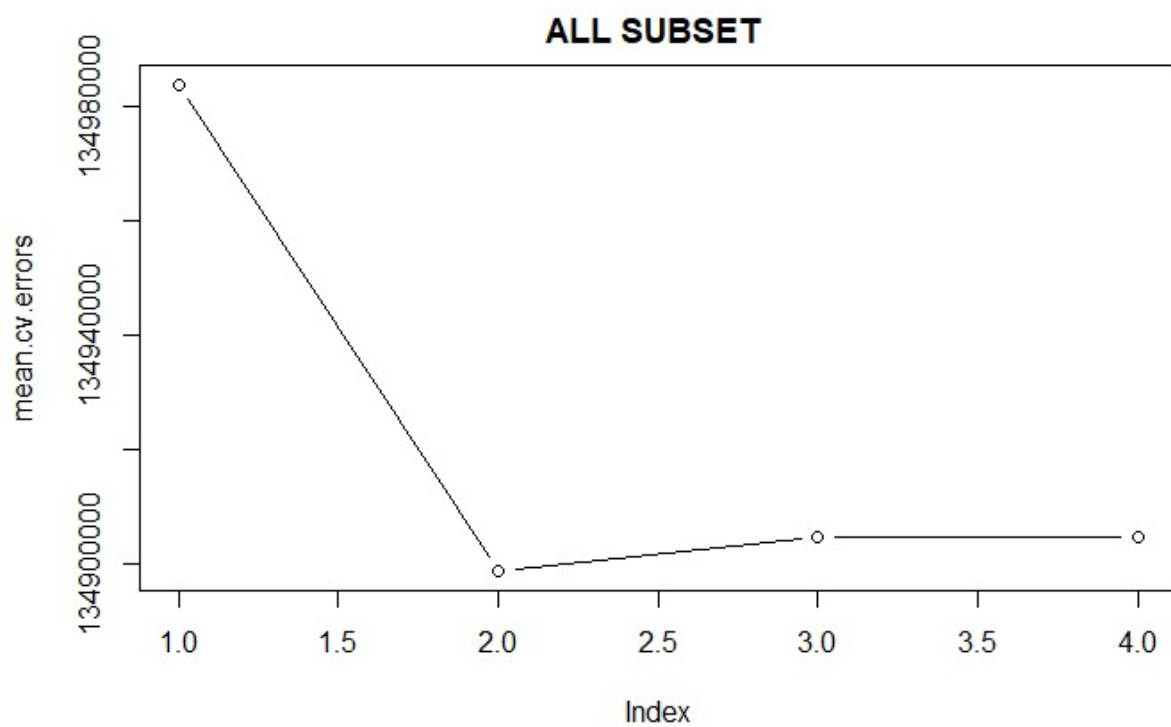


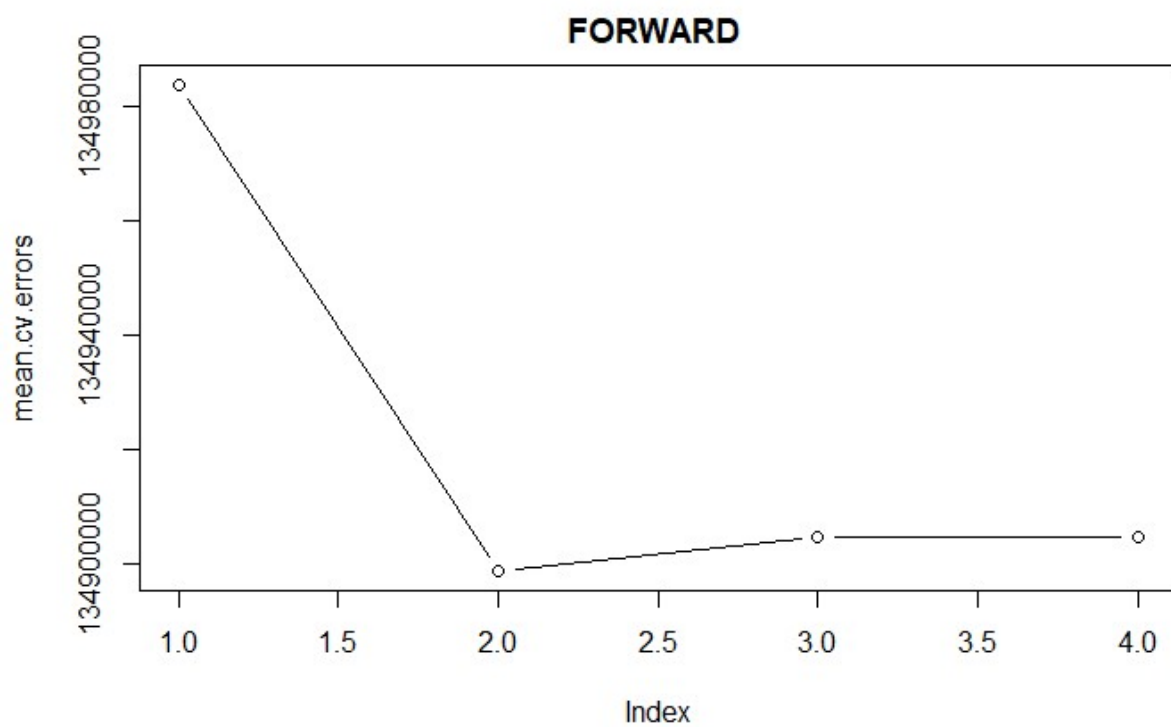


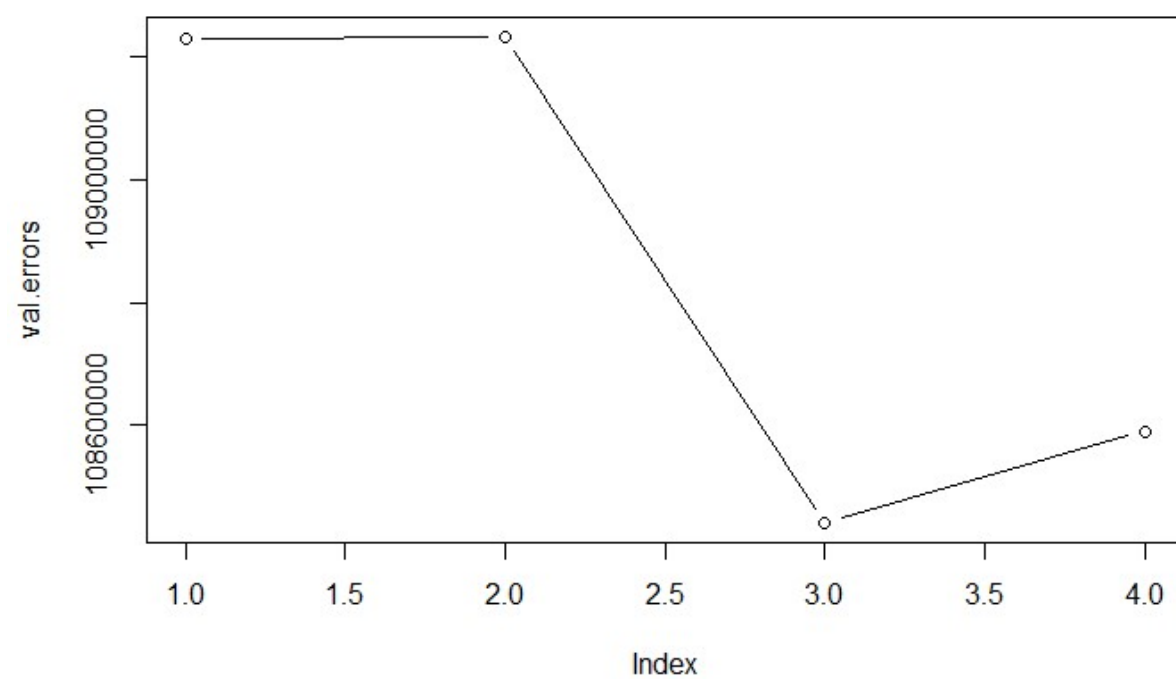


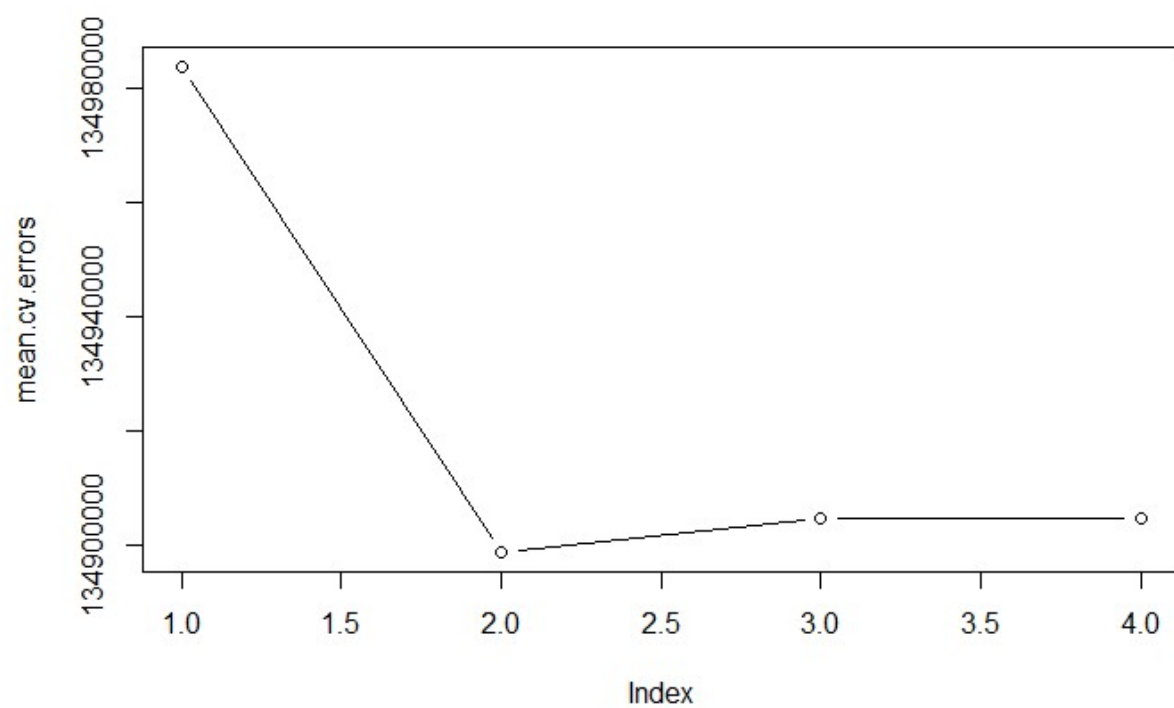
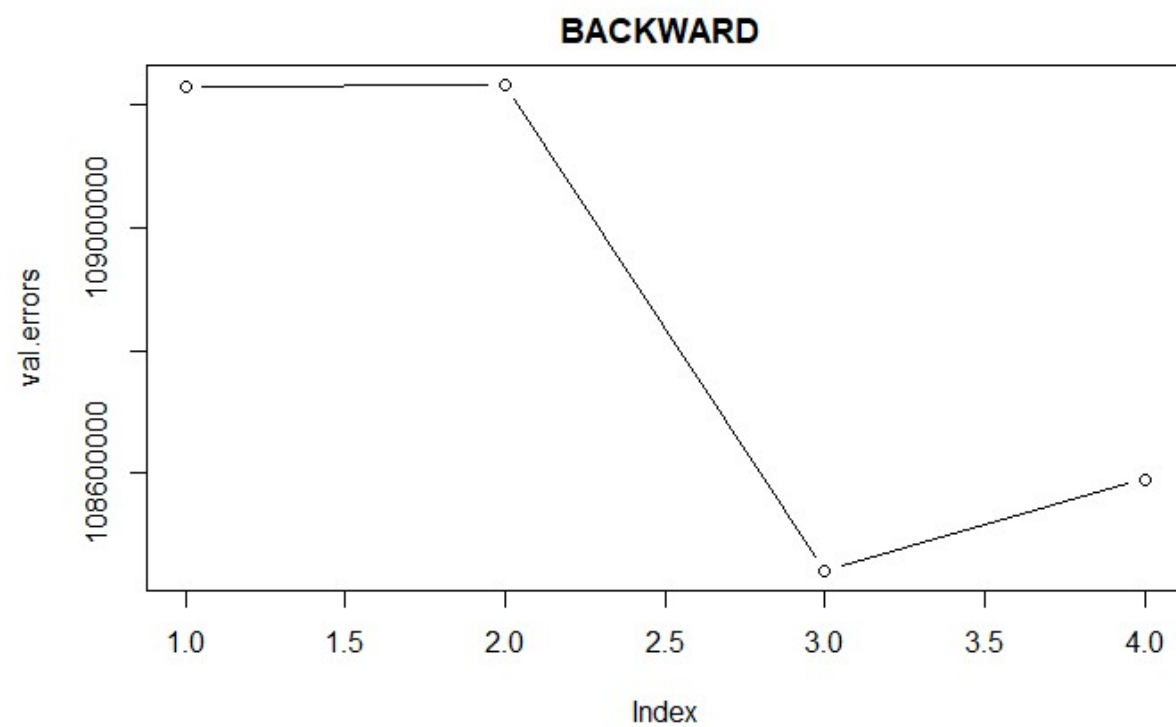
Best subset selection and forward and backward selection were also performed with 10-fold cross-validation . The selection done without cross-validation determined that the four model was the best model . However, when cross-validation was performed the three and two variable model were also deemed as strong models. The three variable model had the lowest MSE and the two variable model had the lowest cross validation error . This model had the lowest cross-validation error in forward, best subset and backward selection :Shares=  $3039.01447 + 57.62953 \text{ num\_imgs} + 75.59873 \text{ num\_videos}$  . This model had the lowest MSE in forward, best subset and backward selection: Shares=  $3268.898882 + 65.517764 \text{ num\_imgs} + 79.859560 \text{ num\_videos} - 0.432109 \text{ n\_tokens\_content}$ .



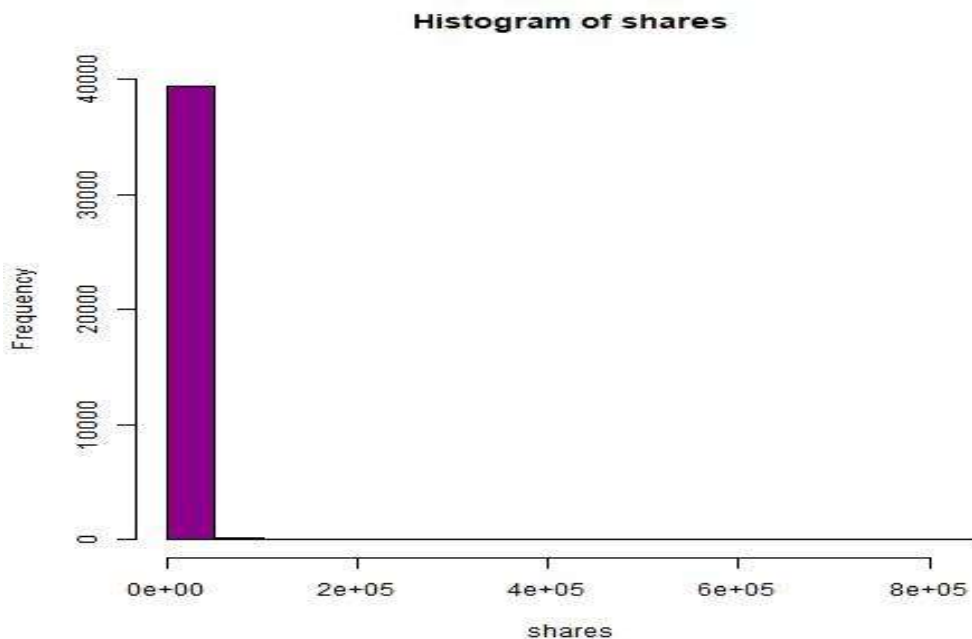








The prediction models created can be a good foundation for determining the characteristics that affect number of article shares. However, a more accurate prediction model for shares would include more variables . Only five of 61 variables were examined . Other variables may have been better predictors for number of shares. However, due to consistent error with processing the large data in the csv file all the variables could not be tested and only four variables were chosen to be tested as possible predictors for shares. In future research, more variables would be tested and compared. Another issue that may have impacted the prediction model was the distribution of shares . The shares distribution was positively skewed. To change the shares distribution to a normal distribution ,the number of shares were logged. However, when using log shares, the MSE value for a model that used all the variables as a predictor was - 0.005016253. This seemed incorrect.



It can be concluded that a one variable model is not a suitable prediction model for shares since the one variable models consistently had higher errors than the models with more variables. However, 10-fold cross-validation, cross validation, forward, best subset and backward stepwise selection sometimes chose different variable models as the best predictor. The three models that were favored were a two variable model, three variable model and a four variable model. This model had the lowest cross-validation error in forward, best subset and backward selection :Shares= 3039.01447 + 57.62953 num\_imgs + 75.59873num\_videos . This model had the lowest MSE in forward, best subset and backward selection: Shares= 3268.898882 + 65.517764 num\_imgs + 79.859560num\_videos -0.432109n\_tokens\_content.The R<sup>2</sup>,adjusted R<sup>2</sup>, Cp,BIC and RSS values were the best for the four variable model,Shares= 2761.5554047 + 65.9358207 num\_imgs + 80.4838476 num\_videos + 44.2298988 n\_tokens\_title -0.4141284 n\_tokens\_content.

**I HAVE ATTACHED MY R CODE,MY R CODE THAT CONSIDERS LOG SHARES  
AND MY KNIT FILE FOR LINE 1 TO 390 OF THE CODE WITH THIS SUBMISSION.**