

machine learning(732A99) lab2

Anubhav Dikshit(anudi287)

10 December 2018

Contents

Assignment 2	3
2.1 Import the data to R and divide into training/validation/test as 50/25/25: use data partitioning code specified in Lecture 1e.	3
2.2 Fit a decision tree to the training data by using the following measures of impurity: a. Deviance b. Gini index and report the misclassification rates for the training and test data. Choose the measure providing the better results for the following steps.	3
3. Use training and validation sets to choose the optimal tree depth. Present the graphs of the dependence of deviances for the training and the validation data on the number of leaves. Report the optimal tree, report it's depth and the variables used by the tree. Interpret the information provided by the tree structure. Estimate the misclassification rate for the test data.	6
4. Use training data to perform classification using Naïve Bayes and report the confusion matrices and misclassification rates for the training and for the test data. Compare the results with those from step 3.	9
5. Use the optimal tree and the Naïve Bayes model to classify the test data by using the following principle: where $\text{prob}(Y \text{'good'})=A$, where $A=0.05, 0.10, \dots, 0.95$	13
6. Repeat Naïve Bayes classification as it was in step 4 but use the following loss matrix (good loss 1, bad loss 10) and report the confusion matrix for the training and test data. Compare the results with the results from step 4 and discuss how the rates has changed and why.	14
Assignment 3	16
1. Reorder your data with respect to the increase of MET and plot EX versus MET. Discuss what kind of model can be appropriate here. Use the reordered data in steps 2-5.	16
2. Use package tree and fit a regression tree model with target EX and feature MET in which the number of the leaves is selected by cross-validation, use the entire data set and set minimum number of observations in a leaf equal to 8 (setting minsize in tree.control). Report the selected tree. Plot the original and the fitted data and histogram of residuals. Comment on the distribution of the residuals and the quality of the fit.	17
3. Compute and plot the 95% confidence bands for the regression tree model from step 2 (fit a regression tree with the same settings and the same number of leaves as in step 2 to the resampled data) by using a non-parametric bootstrap. Comment whether the band is smooth or bumpy and try to explain why. Consider the width of the confidence band and comment whether results of the regression model in step 2 seem to be reliable.	22
4. Compute and plot the 95% confidence and prediction bands the regression tree model from step 2 (fit a regression tree with the same settings and the same number of leaves as in step 2 to the resampled data) by using a parametric bootstrap, assume Normal distribution with mean as labels in the tree leaves, while variance is residual variance. Consider the width of the confidence band and comment whether results of the regression model in step 2 seem to be reliable. Does it look like only 5% of data are outside the prediction band? Should it be?	24
5. Consider the histogram of residuals from step 2 and suggest what kind of bootstrap is actually more appropriate here.	24
Assignment 4	25

1. Conduct a standard PCA by using the feature space and provide a plot explaining how much variation is explained by each feature. Does the plot show how many PC should be extracted? Select the minimal number of components explaining at least 99% of the total variance. Provide also a plot of the scores in the coordinates (PC1, PC2). Are there unusual diesel fuels according to this plot?	25
2. Make trace plots of the loadings of the components selected in step 1. Is there any principle component that is explained by mainly a few original features?	27
Appendix	31

Loading The Libraries

```
if (!require("pacman")) install.packages("pacman")
pacman::p_load(xlsx, ggplot2, MASS, tidyr, dplyr, reshape2, gridExtra,
               tree, caret, e1071, pROC, boot, factoextra)

set.seed(12345)
options("jtools-digits" = 2, scipen = 999)
```

Assignment 2

2.1 Import the data to R and divide into training/validation/test as 50/25/25: use data partitioning code specified in Lecture 1e.

```
set.seed(12345)
credit_data <- read.xlsx("creditscoring.xls", sheetName = "credit")
credit_data$good_bad <- as.factor(credit_data$good_bad)

n=NROW(credit_data)
set.seed(12345)
id=sample(1:n, floor(n*0.4))
train=credit_data[id,]

id1=setdiff(1:n, id)
set.seed(12345)
id2=sample(id1, floor(n*0.3))
valid=credit_data[id2,]

id3=setdiff(id1,id2)
test=credit_data[id3,]
```

2.2 Fit a decision tree to the training data by using the following measures of impurity: a. Deviance b. Gini index and report the misclassification rates for the training and test data. Choose the measure providing the better results for the following steps.

```
set.seed(12345)

# Create a decision tree model
credit_tree_deviance <- tree(good_bad~., data=train, split = c("deviance"))
credit_tree_gini <- tree(good_bad~., data=train, split = c("gini"))

# Visualize the decision tree with rpart.plot
summary(credit_tree_deviance)

##
## Classification tree:
## tree(formula = good_bad ~ ., data = train, split = c("deviance"))
```

```

## Variables actually used in tree construction:
## [1] "duration" "history" "marital" "existcr" "amount" "purpose"
## [7] "savings" "resident" "age" "other"
## Number of terminal nodes: 22
## Residual mean deviance: 0.7423 = 277.6 / 374
## Misclassification error rate: 0.1869 = 74 / 396

summary(credit_tree_gini)

##
## Classification tree:
## tree(formula = good_bad ~ ., data = train, split = c("gini"))
## Variables actually used in tree construction:
## [1] "foreign" "coapp" "depends" "telephon" "existcr" "savings"
## [7] "history" "property" "amount" "marital" "duration" "resident"
## [13] "job" "installp" "purpose" "employed" "housing"
## Number of terminal nodes: 53
## Residual mean deviance: 0.9468 = 324.7 / 343
## Misclassification error rate: 0.2247 = 89 / 396

# predicting on the test dataset to get the misclassification rate.
predict_tree_deviance <- predict(credit_tree_deviance, newdata = test, type = "class")
predict_tree_gini <- predict(credit_tree_gini, newdata = test, type = "class")

conf_tree_deviance <- table(test$good_bad, predict_tree_deviance)
names(dimnames(conf_tree_deviance)) <- c("Actual Test", "P2redicted Test")
caret::confusionMatrix(conf_tree_deviance)

## Confusion Matrix and Statistics
##
##               P2redicted Test
## Actual Test bad good
##      bad   49   43
##      good   56  152
##
##               Accuracy : 0.67
##               95% CI : (0.6136, 0.723)
##      No Information Rate : 0.65
##      P-Value [Acc > NIR] : 0.2539
##
##               Kappa : 0.2534
##  McNemar's Test P-Value : 0.2278
##
##               Sensitivity : 0.4667
##               Specificity : 0.7795
##      Pos Pred Value : 0.5326
##      Neg Pred Value : 0.7308
##      Prevalence : 0.3500
##      Detection Rate : 0.1633
##      Detection Prevalence : 0.3067
##      Balanced Accuracy : 0.6231
##
##      'Positive' Class : bad
##

```

```
conf_tree_gini <- table(test$good_bad, predict_tree_gini)
names(dimnames(conf_tree_gini)) <- c("Actual Test", "Predicted Test")
caret::confusionMatrix(conf_tree_gini)
```

```
## Confusion Matrix and Statistics
##
##           Predicted Test
## Actual Test bad good
##      bad    25    67
##      good   42   166
##
##              Accuracy : 0.6367
##              95% CI : (0.5794, 0.6912)
##      No Information Rate : 0.7767
##      P-Value [Acc > NIR] : 1.00000
##
##              Kappa : 0.0755
##  Mcnemar's Test P-Value : 0.02152
##
##      Sensitivity : 0.37313
##      Specificity : 0.71245
##      Pos Pred Value : 0.27174
##      Neg Pred Value : 0.79808
##      Prevalence : 0.22333
##      Detection Rate : 0.08333
##      Detection Prevalence : 0.30667
##      Balanced Accuracy : 0.54279
##
##      'Positive' Class : bad
##
```

Analysis: On the Training dataset model with ‘deviance’ had a misclassification rate of 18.7% while the model with ‘gini’ split had the misclassification rate of 22.8%.

For the test dataset we see that the model with ‘deviance’ type of split has a accuracy of 67% or misclassification rate of 33%, we see that to predict ‘good’ the accuracy is 73.08% but for predicting bad its just 53.26%. Thus our model is heavily biased towards predicting cases as ‘good’.

For the test dataset we see that the model with ‘gini’ type of split has a accuracy of 62.7% or misclassification rate of 37.3%, we also see that to predict ‘good’ the accuracy is 78.8% but for predicting bad its just 26%. Thus our model is heavily biased towards predicting cases as ‘good’ even more than the model which uses ‘deviance’ to split variable.

Both our models would lead to many bad loan applicants to be given loans which is never a good thing, however among the model the one using ‘deviance’ mode for split is better by 27%.

Thus we will select model using ‘deviance’ for further model building.

3. Use training and validation sets to choose the optimal tree depth. Present the graphs of the dependence of deviances for the training and the validation data on the number of leaves. Report the optimal tree, report its depth and the variables used by the tree. Interpret the information provided by the tree structure. Estimate the misclassification rate for the test data.

```
set.seed(12345)

credit_tree <- tree(good_bad~., data=train, split = c("deviance"))

credit_tree_purned_train <- prune.tree(credit_tree, method = c("deviance"))
credit_tree_purned_valid <- prune.tree(credit_tree, newdata = valid ,method = c("deviance"))

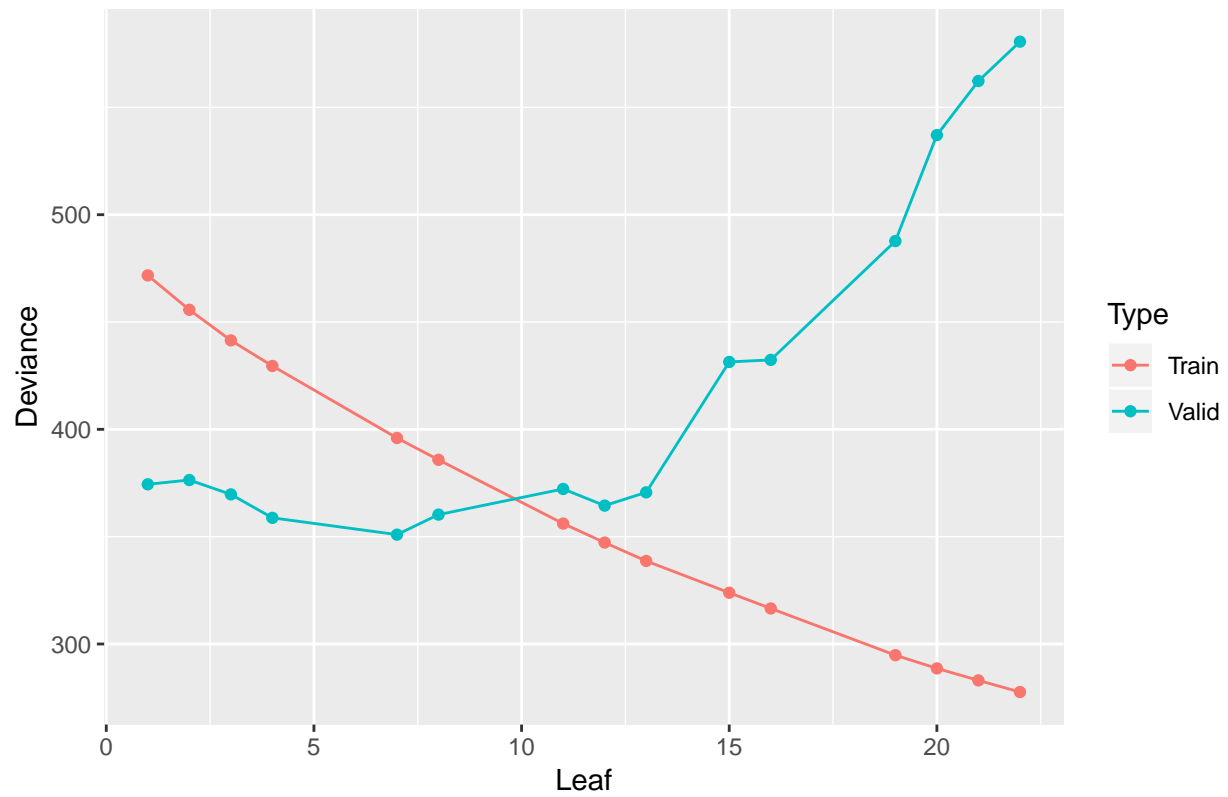
result_train <- cbind(credit_tree_purned_train$size,
                      credit_tree_purned_train$dev, "Train")
result_valid <- cbind(credit_tree_purned_valid$size,
                     credit_tree_purned_valid$dev, "Valid")

result <- as.data.frame(rbind(result_valid, result_train))
colnames(result) <- c("Leaf", "Deviance", "Type")

result$Leaf <- as.numeric(as.character(result$Leaf))
result$Deviance <- as.numeric(as.character(result$Deviance))

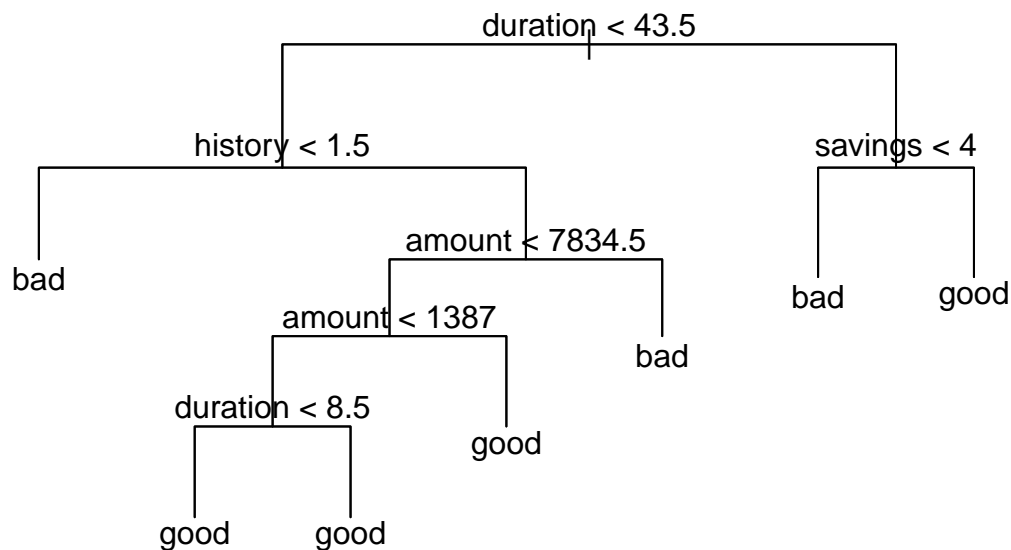
# plot of deviance vs. number of leafs
ggplot(data = result, aes(x = Leaf, y = Deviance, colour = Type)) +
  geom_point() + geom_line() +
  ggtitle("Plot of Deviance vs. Tree Depth (Shows Deviance least at 7)")
```

Plot of Deviance vs. Tree Depth (Shows Deviance least at 7)



```
# prune the tree to the required depth
credit_tree_sniped <- prune.tree(credit_tree, best=7)

plot(credit_tree_sniped)
text(credit_tree_sniped)
```



misclassification rate for best pruned tree

```
result_prune_test <- predict(credit_tree_sniped, newdata = test, type = "class")
```

```
conf_prune_tree_test <- table(test$good_bad, result_prune_test)
```

```
names(dimnames(conf_prune_tree_test)) <- c("Actual Test", "Predicted Test")
```

```
caret::confusionMatrix(conf_prune_tree_test)
```

Confusion Matrix and Statistics

##

Predicted Test

Actual Test bad good

bad 32 60

good 19 189

##

Accuracy : 0.7367

95% CI : (0.683, 0.7856)

No Information Rate : 0.83

P-Value [Acc > NIR] : 1

##

Kappa : 0.2929

McNemar's Test P-Value : 0.000006784

##

Sensitivity : 0.6275

Specificity : 0.7590

Pos Pred Value : 0.3478


```
##          Neg Pred Value : 0.9087
##          Prevalence : 0.1700
##          Detection Rate : 0.1067
##    Detection Prevalence : 0.3067
##          Balanced Accuracy : 0.6932
##
##          'Positive' Class : bad
##
```

Analysis: Choosing optimal depth tree we get that '7' as the best depth. The variables used in the best tree are- duration, history, savings, amount.

From the tree structure we can see that the following variables are best to split on, 'duration' < 43.5 then 'savings' < 4, 'history' < 1.5 and finally 'amount'.

The accuracy on the model trained on 'train' dataset is (77% and misclassification 23%) and on the 'test' dataset accuracy is 73.67%, thus the misclassification rate is 26.33%. We see that model predicts 'good' applicants very well (accuracy of 90%) while it classifies 'bad' applicant way badly (accuracy is 34.78%).

Thus this model would be very bad for the business and would likely to run the business bankrupt.

4. Use training data to perform classification using Naïve Bayes and report the confusion matrices and misclassification rates for the training and for the test data. Compare the results with those from step 3.

```
#Fitting the Naive Bayes model
credit_naive_model = naiveBayes(good_bad ~., data=train)
credit_naive_model
```

```
##
## Naive Bayes Classifier for Discrete Predictors
##
## Call:
## naiveBayes.default(x = X, y = Y, laplace = laplace)
##
## A-priori probabilities:
## Y
##    bad    good
## 0.285 0.715
##
## Conditional probabilities:
##      duration
## Y      [,1]    [,2]
## bad 24.85965 14.95946
## good 19.13287 11.11076
##
##      history
## Y      [,1]    [,2]
## bad  2.263158 1.113494
## good 2.744755 1.033419
##
##      purpose
## Y      [,1]    [,2]
## bad  2.883929 2.859300
## good 2.718310 2.409833
```

```

##
##      amount
## Y      [,1]      [,2]
## bad  3661.202 3504.882
## good 2995.524 2431.313
##
##      savings
## Y      [,1]      [,2]
## bad   1.736842 1.343965
## good   2.314685 1.664572
##
##      employed
## Y      [,1]      [,2]
## bad   3.236842 1.250151
## good   3.500000 1.147843
##
##      installp
## Y      [,1]      [,2]
## bad   3.175439 0.9798466
## good   2.898601 1.1271131
##
##      marital
## Y      [,1]      [,2]
## bad   2.596491 0.7252622
## good   2.716783 0.6650668
##
##      coapp
## Y      [,1]      [,2]
## bad   1.105263 0.4072011
## good   1.171329 0.5383695
##
##      resident
## Y      [,1]      [,2]
## bad   2.824561 1.074611
## good   2.807692 1.121439
##
##      property
## Y      [,1]      [,2]
## bad   2.456140 1.073744
## good   2.300699 1.069743
##
##      age
## Y      [,1]      [,2]
## bad   33.91228 10.97589
## good   36.51399 11.21050
##
##      other
## Y      [,1]      [,2]
## bad   2.500000 0.8227947
## good   2.751748 0.6362558
##
##      housing
## Y      [,1]      [,2]
## bad   1.938596 0.5991895

```

```
## good 1.965035 0.5014028
##
## existcr
## Y      [,1]      [,2]
## bad  1.368421 0.5988656
## good 1.458042 0.5893693
##
## job
## Y      [,1]      [,2]
## bad  2.824561 0.7194592
## good 2.849650 0.6503305
##
## depends
## Y      [,1]      [,2]
## bad  1.140351 0.3488843
## good 1.164336 0.3712295
##
## telephon
## Y      [,1]      [,2]
## bad  1.315789 0.4668818
## good 1.388112 0.4881745
##
## foreign
## Y      [,1]      [,2]
## bad  1.017544 0.1318659
## good 1.045455 0.2086640
```

```
#Prediction on the dataset
```

```
predict_naive_train = predict(credit_naive_model, newdata=train, type = "class")
predict_naive_test = predict(credit_naive_model, newdata=test, type = "class")

conf_naive_train <- table(train$good_bad, predict_naive_train)
names(dimnames(conf_naive_train)) <- c("Actual Train", "Predicted Train")
caret::confusionMatrix(conf_naive_train)
```

```
## Confusion Matrix and Statistics
```

```
##
##              Predicted Train
## Actual Train bad good
##      bad    62    52
##      good   55   231
##
##              Accuracy : 0.7325
##              95% CI : (0.6863, 0.7753)
##      No Information Rate : 0.7075
##      P-Value [Acc > NIR] : 0.1480
##
##              Kappa : 0.3488
##      McNemar's Test P-Value : 0.8467
##
##              Sensitivity : 0.5299
##              Specificity : 0.8163
##      Pos Pred Value : 0.5439
##      Neg Pred Value : 0.8077
##              Prevalence : 0.2925
```

```
##          Detection Rate : 0.1550
##    Detection Prevalence : 0.2850
##          Balanced Accuracy : 0.6731
##
##          'Positive' Class : bad
##

conf_naive_test <- table(test$good_bad, predict_naive_test)
names(dimnames(conf_naive_test)) <- c("Actual Test", "Predicted Test")
caret::confusionMatrix(conf_naive_test)

## Confusion Matrix and Statistics
##
##          Predicted Test
## Actual Test bad good
##          bad    50   42
##          good   45  163
##
##          Accuracy : 0.71
##          95% CI : (0.6551, 0.7607)
##    No Information Rate : 0.6833
##    P-Value [Acc > NIR] : 0.1762
##
##          Kappa : 0.3242
##  Mcnemar's Test P-Value : 0.8302
##
##          Sensitivity : 0.5263
##          Specificity : 0.7951
##          Pos Pred Value : 0.5435
##          Neg Pred Value : 0.7837
##          Prevalence : 0.3167
##          Detection Rate : 0.1667
##    Detection Prevalence : 0.3067
##          Balanced Accuracy : 0.6607
##
##          'Positive' Class : bad
##
```

Analysis:

For the train dataset using NaiveBayes method we get accuracy 73% or misclassification of 27%, here we also notice that the accuracy of class 'bad' is 54% while for class 'good' is 80%, thus the model is more balanced in predicting, thus its still biased in predict one class over the other.

For the test dataset using NaiveBayes method we get accuracy 71% or misclassification of 29%, here we also notice that the accuracy of class 'bad' is 54% while for class 'good' is 78%, thus the model is almost the same compared to train.

Compared to step3, we see that for the 'train' dataset the optimal tree has accuracy of 77% while it is 73% on the 'test' dataset. For the NaiveBayes model, accuracy on the 'train' dataset is 73% and while it is 71% on the 'test' dataset.

Accuracy is only part of the story what we see is better here is that this model classifies 'bad' customers better better for both train and test dataset than decision tree (54% for both train and test for naive compared 38% train, 34% test for decision tree).

Thus the model is better to be used for the business than the one in the step3, the risk of providing loans to bad applicant is lesser than the previous model but its still not good enough!

5. Use the optimal tree and the Naïve Bayes model to classify the test data by using the following principle: where $\text{prob}(Y|\text{'good'})=A$, where $A=0.05, 0.10, \dots, 0.95$.

Compute the TPR and FPR values for the two models and plot the corresponding ROC curves. Conclusion?

```
set.seed(12345)

credit_tree <- tree(good_bad~., data=train, split = c("deviance"))
credit_naive_model = naiveBayes(good_bad ~., data=train)

# prune the tree to the required depth
credit_tree_sniped <- prune.tree(credit_tree, best=7)

# predicting class, getting probability
predict_prune_test_prob <- predict(credit_tree_sniped, newdata = test)
predict_naive_test_prob <- predict(credit_naive_model, newdata=test, type = "raw")

# data mugging
probability_data_naive <- as.data.frame(cbind(predict_naive_test_prob,
                                              as.character(test$good_bad), "naivebayes"))
probability_data_tree <- as.data.frame(cbind(predict_prune_test_prob,
                                              as.character(test$good_bad), "tree"))

probability_data_combined <- rbind(probability_data_tree, probability_data_naive)
colnames(probability_data_combined) <- c("prob_bad", "prob_good",
                                         "actual_test_class", "model")

# final dataset
probability_data_combined$prob_good <- as.numeric(as.character(probability_data_combined$prob_good))

# changing the threshold and printing the probability

tree_list <- NULL
naive_list <- NULL
final <- NULL
for(threshold in seq(from = 0.05, to = 0.95, by = 0.05)){
  probability_data_combined$predicted_class <- ifelse(probability_data_combined$prob_good > threshold,

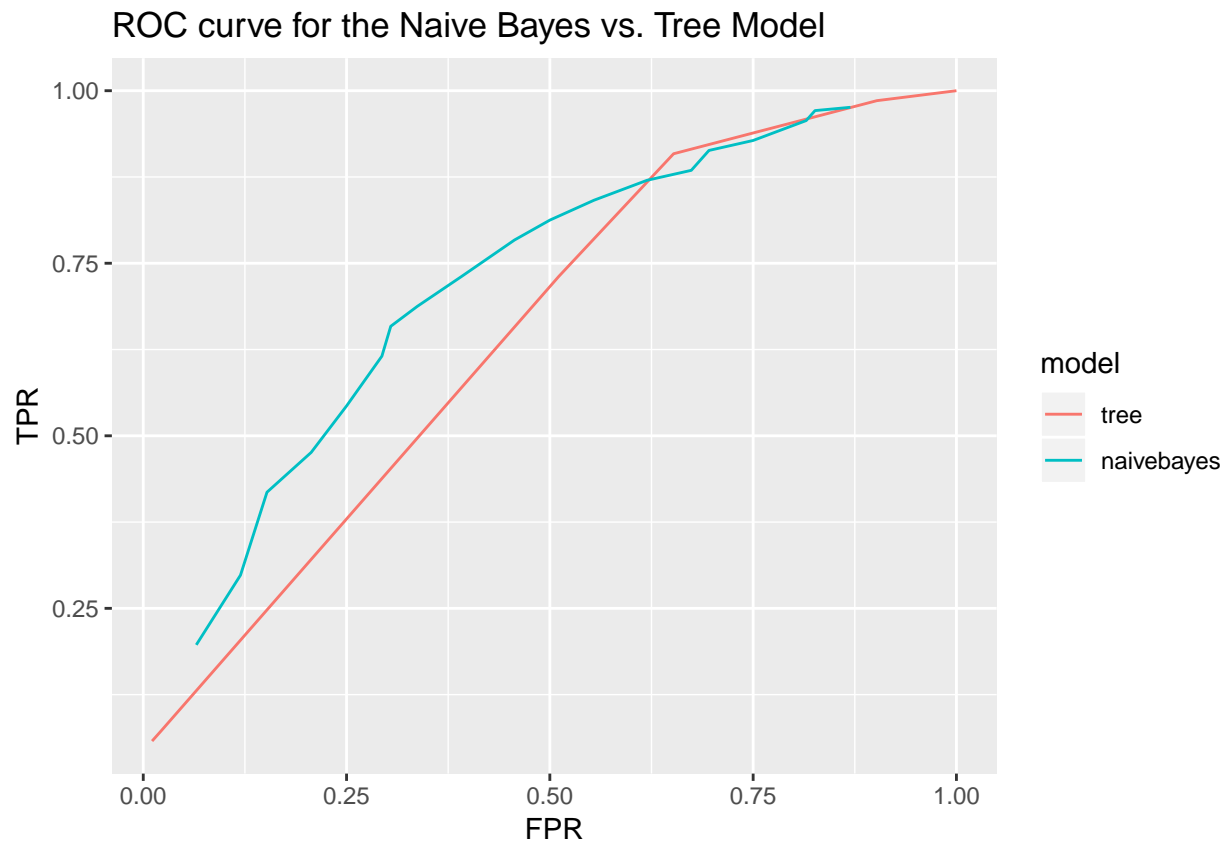
  df2 <- probability_data_combined[,c("model", "actual_test_class", "predicted_class")]
  df2$threshold <- threshold
  df2$match <- ifelse(df2$actual_test_class == df2$predicted_class, 1, 0)

  final <- rbind(df2, final)
}

# Creating the FRP and TRP for each model and threshold
final$temp <- 1
final_summary <- final %>%
group_by(model, threshold) %>%
summarise(total_positive = sum(temp[actual_test_class == "good"]),
          total_negative = sum(temp[actual_test_class == "bad"]),
          correct_positive = sum(temp[actual_test_class == "good" & predicted_class == "good"]),
          false_positive = sum(temp[actual_test_class == "bad" & predicted_class == "good"])) %>%
```

```
mutate(TPR = correct_positive/total_positive, FPR = false_positive/total_negative) %>%
select(model, threshold, TPR, FPR)

ggplot(data = final_summary, aes(x = FPR, y=TPR)) + geom_line(aes(colour = model)) +
  ggtitle("ROC curve for the Naive Bayes vs. Tree Model")
```



Analysis: We find that 'naivebayes' model is better than 'tree' model for across varying threshold values.

6. Repeat Naïve Bayes classification as it was in step 4 but use the following loss matrix (good loss 1, bad loss 10) and report the confusion matrix for the training and test data. Compare the results with the results from step 4 and discuss how the rates has changed and why.

```
set.seed(12345)

credit_naive_model = naiveBayes(good_bad ~., data=train)

# predicting class, getting probability
predict_naive_train_prob <- predict(credit_naive_model, newdata=train, type = "raw")
predict_naive_test_prob <- predict(credit_naive_model, newdata=test, type = "raw")

train <- cbind(predict_naive_train_prob, train)
test <- cbind(predict_naive_test_prob, test)

# class based on the loss matrix
```

```

train$predicted_class <- ifelse(train$good > 10*train$bad, "good", "bad")
test$predicted_class <- ifelse(test$good > 10*test$bad, "good", "bad")

# confusion matrix
conf_naive_train <- table(train$good_bad, train$predicted_class)
names(dimnames(conf_naive_train)) <- c("Actual Train", "Predicted Train")
caret::confusionMatrix(conf_naive_train)

```

```

## Confusion Matrix and Statistics
##
##               Predicted Train
## Actual Train bad good
##      bad  105    9
##      good 202   84
##
##               Accuracy : 0.4725
##               95% CI : (0.4227, 0.5227)
##      No Information Rate : 0.7675
##      P-Value [Acc > NIR] : 1
##
##               Kappa : 0.1423
##  Mcnemar's Test P-Value : <0.0000000000000002
##
##      Sensitivity : 0.3420
##      Specificity : 0.9032
##      Pos Pred Value : 0.9211
##      Neg Pred Value : 0.2937
##      Prevalence : 0.7675
##      Detection Rate : 0.2625
##      Detection Prevalence : 0.2850
##      Balanced Accuracy : 0.6226
##
##      'Positive' Class : bad
##

```

```

conf_naive_test <- table(test$good_bad, test$predicted_class)
names(dimnames(conf_naive_test)) <- c("Actual Test", "Predicted Test")
caret::confusionMatrix(conf_naive_test)

```

```

## Confusion Matrix and Statistics
##
##               Predicted Test
## Actual Test bad good
##      bad   81   11
##      good 147   61
##
##               Accuracy : 0.4733
##               95% CI : (0.4157, 0.5315)
##      No Information Rate : 0.76
##      P-Value [Acc > NIR] : 1
##
##               Kappa : 0.123
##  Mcnemar's Test P-Value : <0.0000000000000002
##

```

```
##           Sensitivity : 0.3553
##           Specificity : 0.8472
##           Pos Pred Value : 0.8804
##           Neg Pred Value : 0.2933
##           Prevalence : 0.7600
##           Detection Rate : 0.2700
##           Detection Prevalence : 0.3067
##           Balanced Accuracy : 0.6012
##
##           'Positive' Class : bad
##
```

Assignment 3

1. Reorder your data with respect to the increase of MET and plot EX versus MET. Discuss what kind of model can be appropriate here. Use the reordered data in steps 2-5.

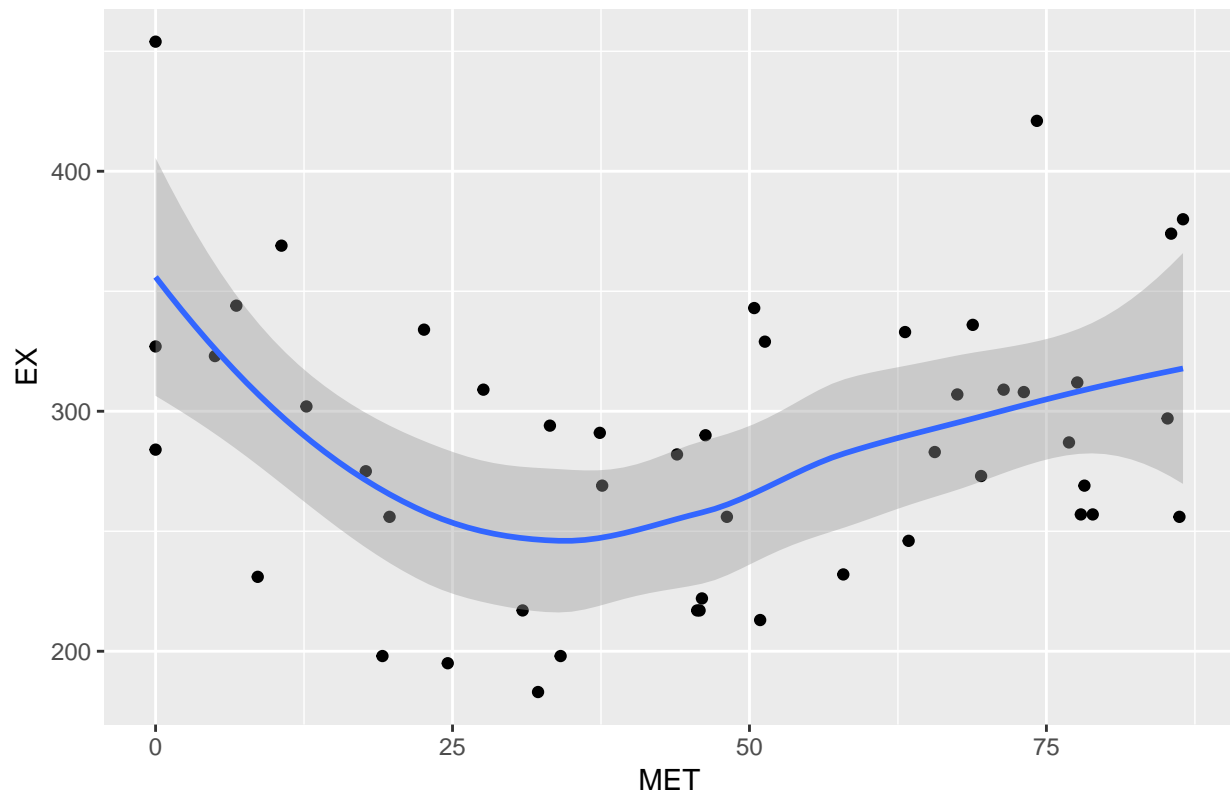
```
rm(list=ls())

set.seed(12345)
state_data <- read.csv2("state.csv")

state_data <- state_data %>% arrange(MET)

ggplot(data = state_data, aes(x=MET, y = EX)) +
  geom_point() +
  geom_smooth(method = 'loess') +
  ggtitle("Plot of MET vs. EX")
```


Plot of MET vs. EX



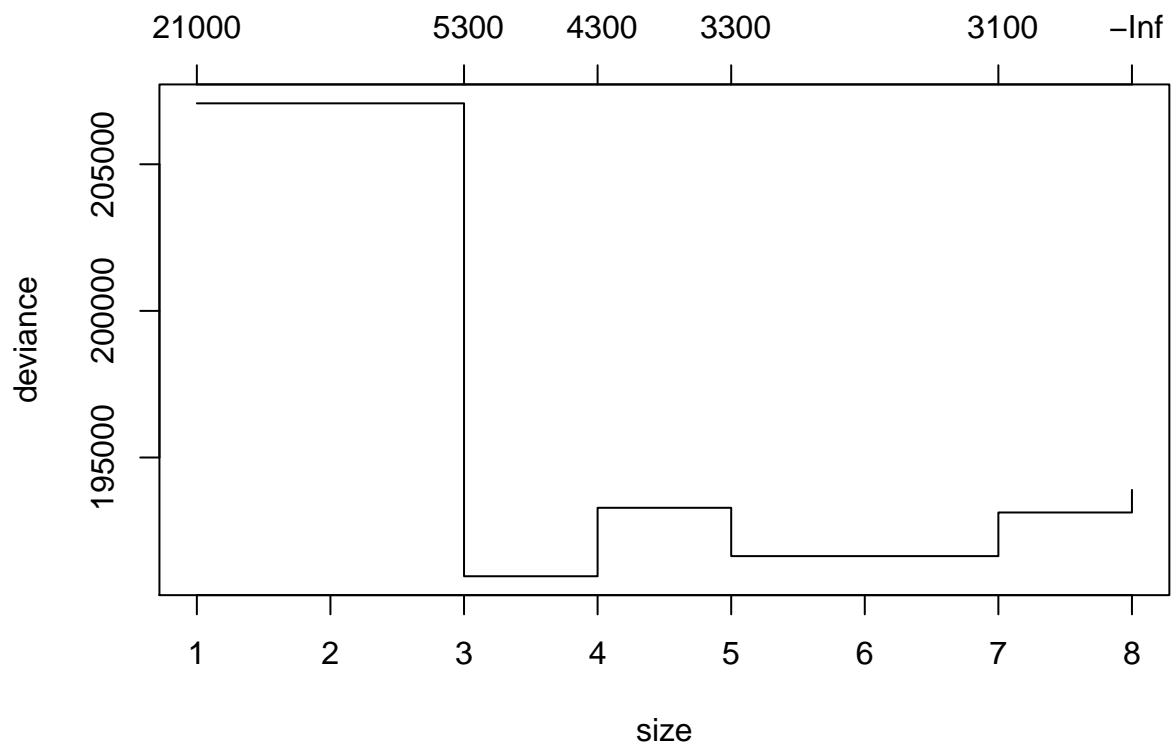
Analysis: As evident from the graph the best model, linear regression will not be a good fit, even the trend is non linear. Piece wise linear model(spline) might be a good one, thus regression per group/cluster will be a good approach.

2. Use package tree and fit a regression tree model with target EX and feature MET in which the number of the leaves is selected by cross-validation, use the entire data set and set minimum number of observations in a leaf equal to 8 (setting minsize in tree.control). Report the selected tree. Plot the original and the fitted data and histogram of residuals. Comment on the distribution of the residuals and the quality of the fit.

```
set.seed(12345)

state_tree_regression <- tree(data = state_data, EX~MET,
                             control = tree.control(nobs=NROW(state_data),
                                                      minsize = 8))

state_cv_tree <- cv.tree(state_tree_regression, FUN = prune.tree)
plot(state_cv_tree)
```



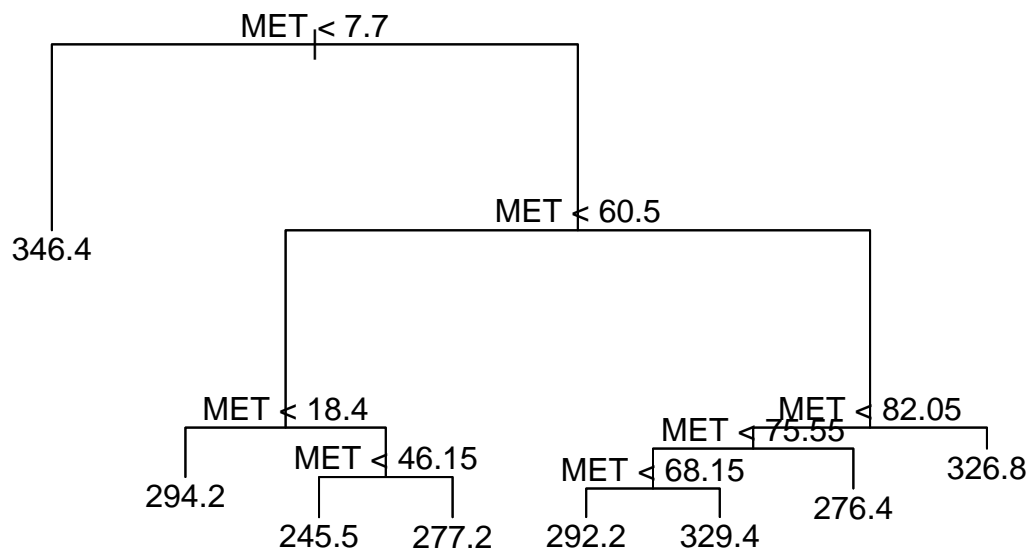
The best size is either 3 or 4

purging the tree for leaf size of 3

```
state_cv_tree_purned <- prune.tree(state_tree_regression, k = 3)
```

```
plot(state_cv_tree_purned, main="Pruned Tree for the given dataset")
```

```
text(state_cv_tree_purned)
```

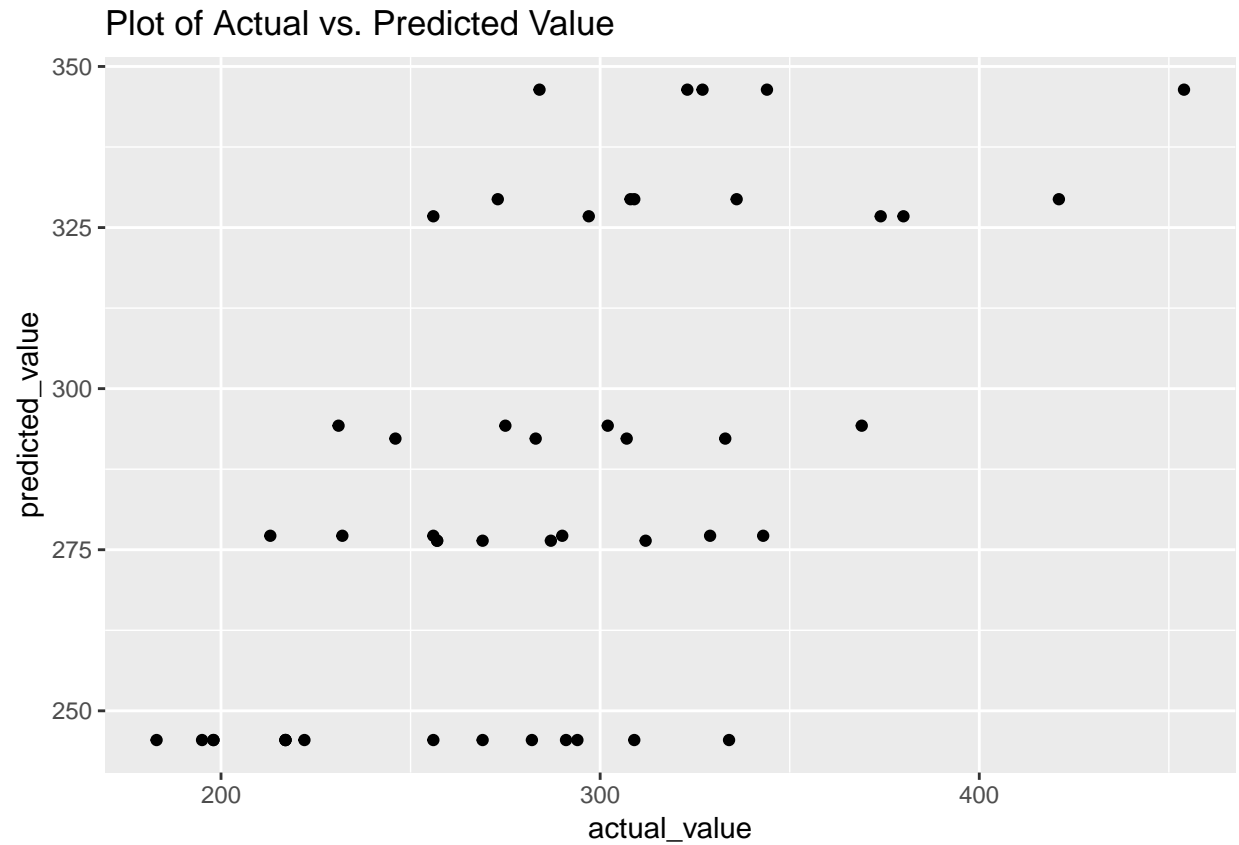


```

# Original vs. Fitted values
compare_data <- predict(state_cv_tree_purned, newdata = state_data)
compare_data <- cbind(compare_data, state_data$EX)
compare_data <- as.data.frame(compare_data)
colnames(compare_data) <- c("predicted_value", "actual_value")
compare_data$residual <- compare_data$actual_value - compare_data$predicted_value

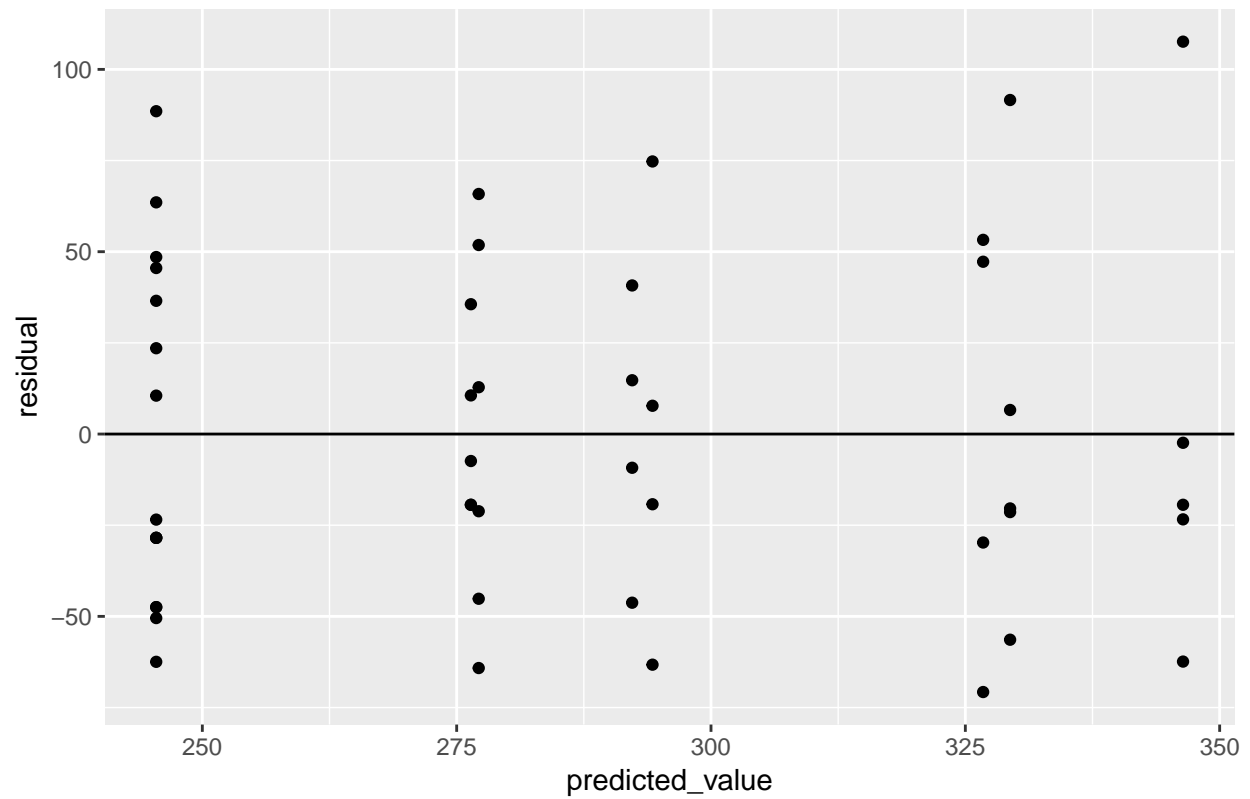
# plots
ggplot(compare_data, aes(x = actual_value, y = predicted_value)) +
  geom_point() +
  ggtitle("Plot of Actual vs. Predicted Value")

```

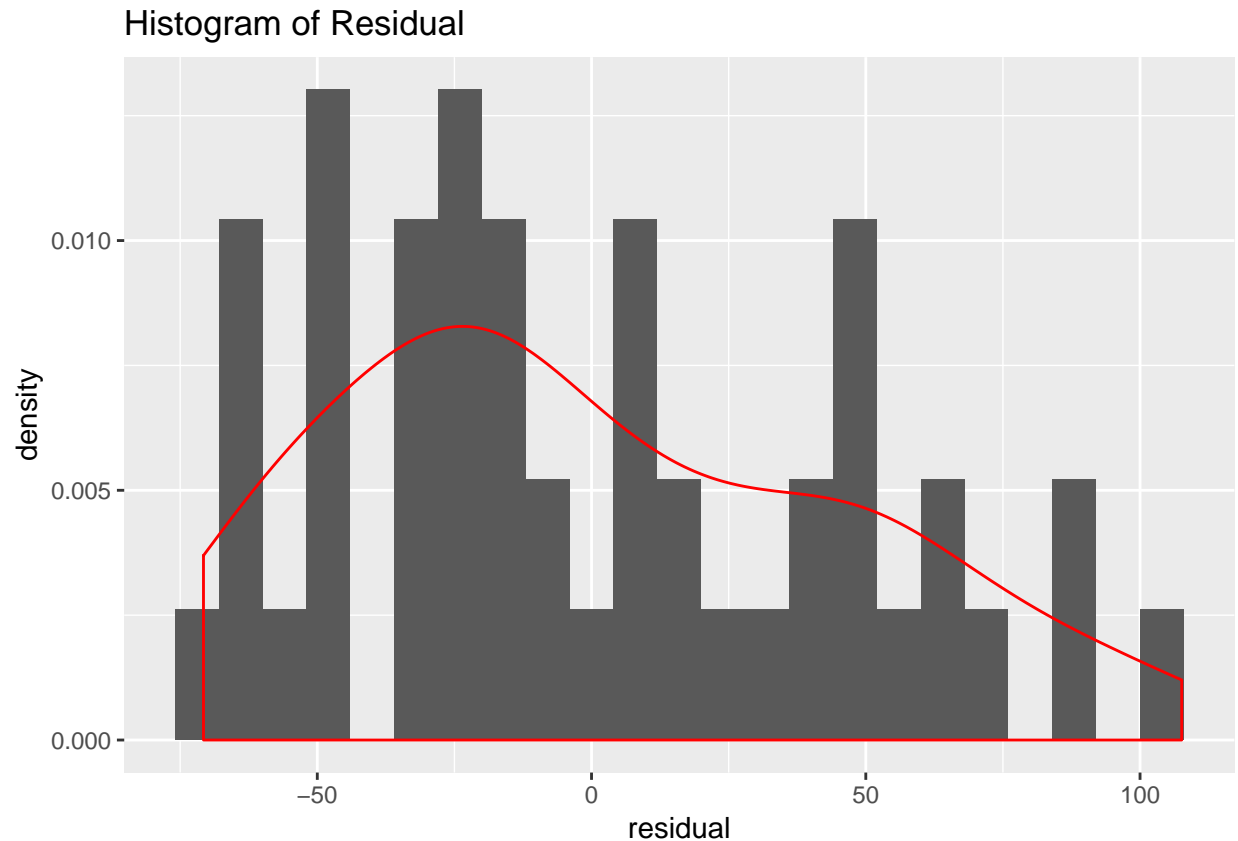


```
ggplot(compare_data, aes(x = predicted_value, y = residual)) +  
  geom_point() + geom_abline(slope=0, intercept=0) +  
  ggtitle("Plot of Predicted Value vs. Residual")
```

Plot of Predicted Value vs. Residual



```
ggplot(data = compare_data, aes(x = residual)) +  
  geom_histogram(aes(y = ..density..), binwidth = 8) +  
  geom_density(colour = "red") +  
  ggtitle("Histogram of Residual")
```



Analysis:

The predicted vs. Actual provides us the insight that for lower values of variable, our model is over predicting (actual value ~200) while the predicted value is ~250. While for larger values (~400) our model is under predicting (~330). At around the mean value (~300) the predicted values are both under and over predicted thus no bias. Thus for values that are away from the mean our model is biased towards over/under predicted while values close to mean our model is not biased.

From the plot of Predicted vs. Residual values we can see that error appears random, neither is large bias/concentration of the error towards any value except at lower/higher values (more points on one side of the line)

From the histogram we can see that the histogram has higher values on the left of zero, and a longer tail on the right. Thus from the above three points we can see that there is scope of improvement in the model especially in the extreme values of the predicted values.

3. Compute and plot the 95% confidence bands for the regression tree model from step 2 (fit a regression tree with the same settings and the same number of leaves as in step 2 to the resampled data) by using a non-parametric bootstrap. Comment whether the band is smooth or bumpy and try to explain why. Consider the width of the confidence band and comment whether results of the regression model in step 2 seem to be reliable.

```
set.seed(12345)
```

```

# function to obtain R-Squared from the data
rsq <- function(data, indices) {
  d <- data[indices,] # allows boot to select sample

  model <- tree(data = d,
                EX~MET,
                control = tree.control(nobs=NROW(data), minsize = 8))

  model_purned <- prune.tree(model, k = 3)

  compare_data <- predict(model_purned, newdata = d)

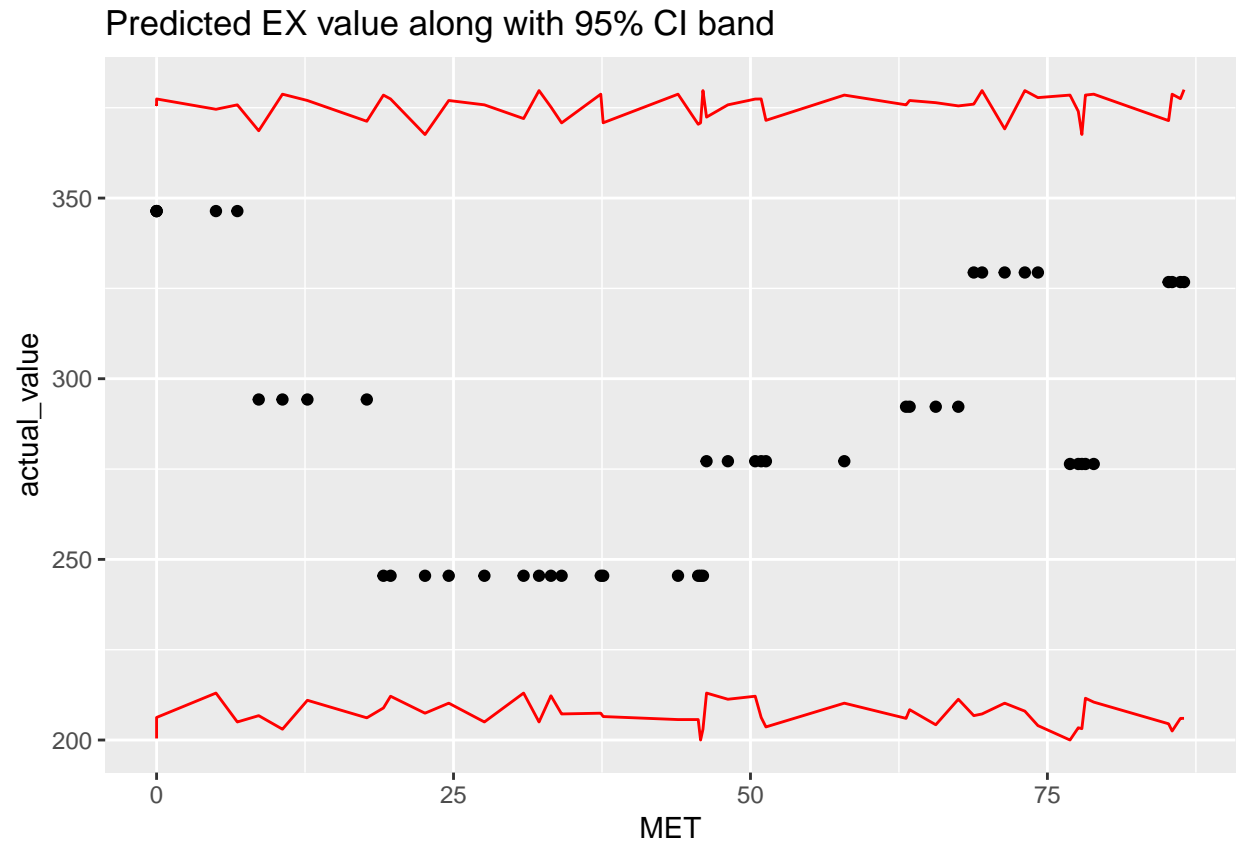
  return(compare_data)
}
# bootstrapping with 1000 replications
results <- boot(data=state_data, statistic=rsq, R=1000)

e=envelope(results) #compute confidence bands

compare_data_non_para_boot <- compare_data
compare_data_non_para_boot$MET <- state_data$MET
compare_data_non_para_boot$lower_bound <- e$point[2,]
compare_data_non_para_boot$upper_bound <- e$point[1,]

ggplot(data=compare_data_non_para_boot, aes(x = MET, y = actual_value)) +
  geom_point(aes(y=predicted_value)) +
  geom_line(aes(y=upper_bound), colour="red") + # first layer
  geom_line(aes(y=lower_bound), colour="red") +
  ggtitle("Predicted EX value along with 95% CI band")

```



Analysis:

The confidence bands certainly appear to be bumpy and not smooth, the confidence bands are bumpy because the predicted values shows large fluctuations. From the width of the confidence band we can assume that our model is not a good one. Ideally we want the confidence band to be narrow thus a wider band suggests we need further tuning to the model.

4. Compute and plot the 95% confidence and prediction bands the regression tree model from step 2 (fit a regression tree with the same settings and the same number of leaves as in step 2 to the resampled data) by using a parametric bootstrap, assume Normal distribution with mean as labels in the tree leaves, while variance is residual variance. Consider the width of the confidence band and comment whether results of the regression model in step 2 seem to be reliable. Does it look like only 5% of data are outside the prediction band? Should it be?

```
set.seed(12345)
```

5. Consider the histogram of residuals from step 2 and suggest what kind of bootstrap is actually more appropriate here.

```
set.seed(12345)
```


Assignment 4

```
rm(list=ls())
NIR_data <- read.csv2("NIRSpectra.csv")
```

1. Conduct a standard PCA by using the feature space and provide a plot explaining how much variation is explained by each feature. Does the plot show how many PC should be extracted? Select the minimal number of components explaining at least 99% of the total variance. Provide also a plot of the scores in the coordinates (PC1, PC2). Are there unusual diesel fuels according to this plot?

```
set.seed(12345)

pca_data = select(NIR_data, -c(Viscosity))
pca_result = prcomp(pca_data)

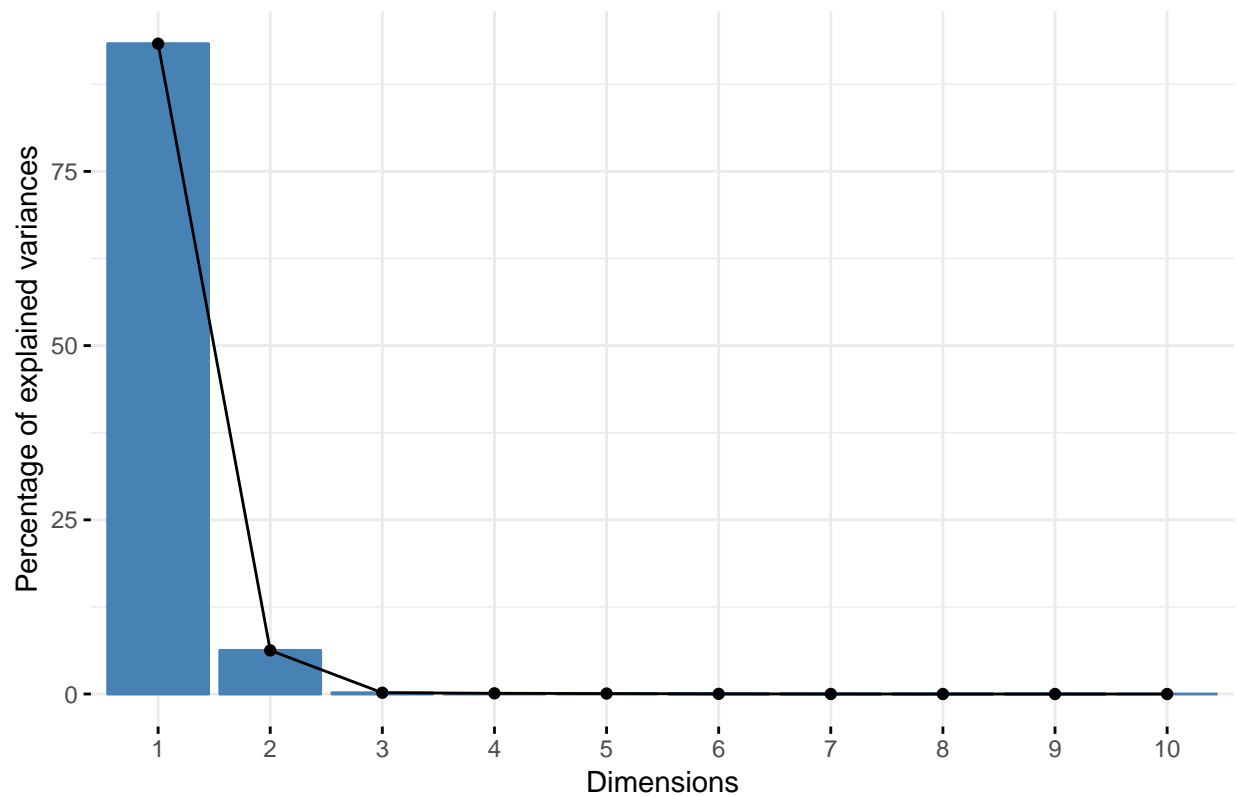
contribution <- summary(pca_result)$importance
knitr::kable(contribution[,1:5],
              caption = "Contribution of PCA axis towards variance explanation")
```

Table 1: Contribution of PCA axis towards variance explanation

	PC1	PC2	PC3	PC4	PC5
Standard deviation	0.122062	0.0316205	0.0054353	0.0040107	0.0033031
Proportion of Variance	0.933320	0.0626300	0.0018500	0.0010100	0.0006800
Cumulative Proportion	0.933320	0.9959600	0.9978100	0.9988200	0.9995000

```
# plots PCA components and the eigen vectors
factoextra::fviz_eig(pca_result)
```

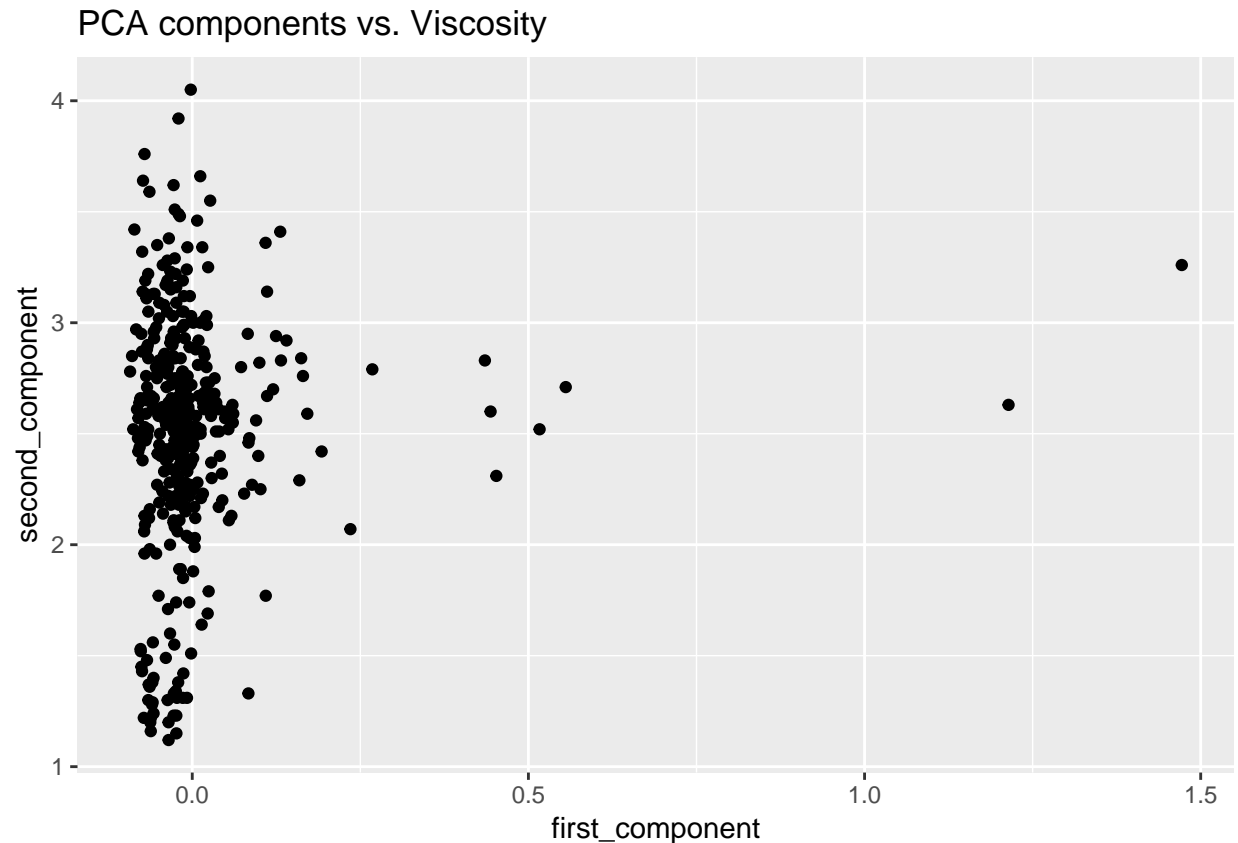
Scree plot



```
# pca components and the viscosity
pca_result_data = cbind(first_component = pca_result$x[,1],
                        second_component = pca_result$x[,2],
                        Viscosity = NIR_data$Viscosity)

pca_result_data = as.data.frame(pca_result_data)

# plotting the data variation and the viscosity
ggplot(data = pca_result_data, aes(x = first_component, y = second_component)) +
  geom_point(aes(y = Viscosity)) + ggtitle("PCA components vs. Viscosity")
```



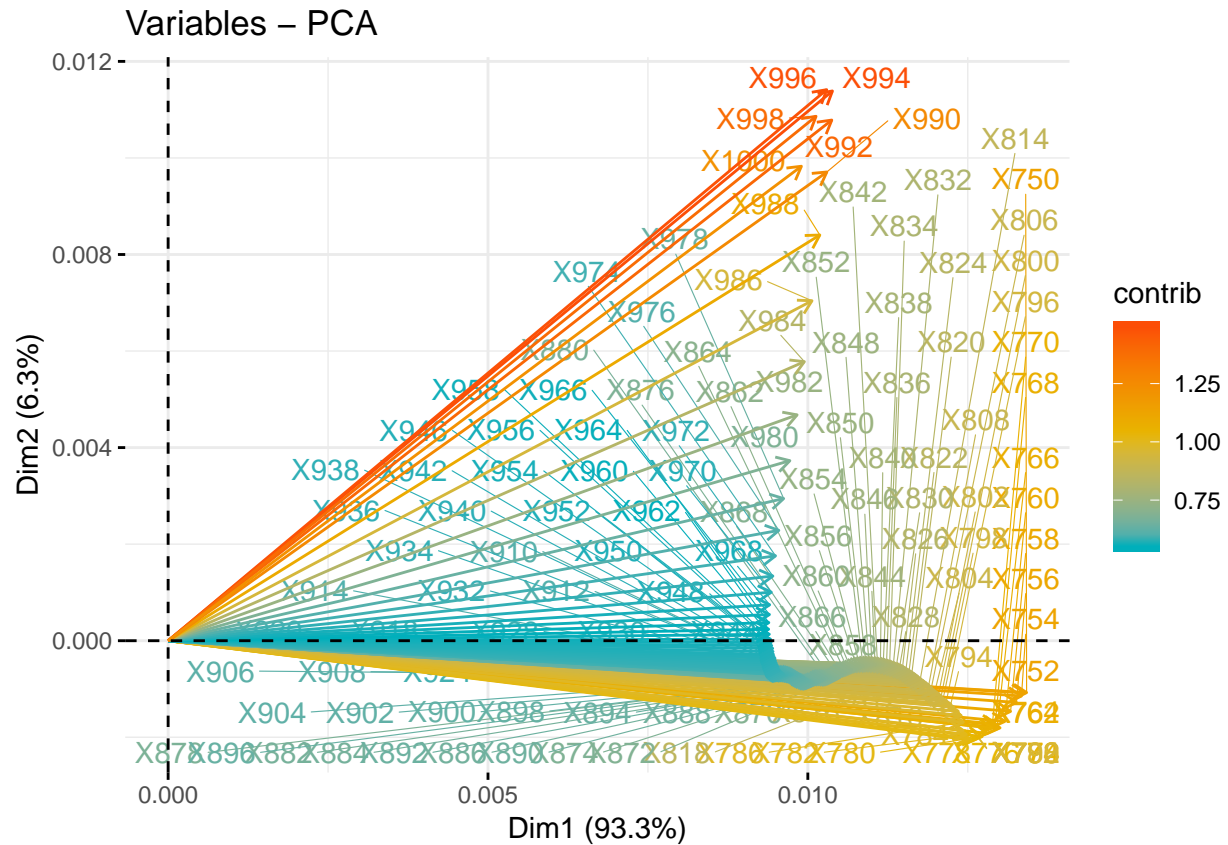
Analysis:

From the plot of PCA component vs. Viscosity, we can see that second component of PCA is needed, this is due to the fact the most of the data is vertically spread, removing this dimension would make it impossible to differentiate the types of diesel.

From the plot we can also see that there are evidently some diesel that are outliers (diesel with first_component > 0.5 and second component ~4).

2. Make trace plots of the loadings of the components selected in step 1. Is there any principle component that is explained by mainly a few original features?

```
# showing the components of the PCA components
factoextra::fviz_pca_var(pca_result,
  col.var = "contrib", # Color by contributions to the PC
  gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"),
  repel = TRUE        # Avoid text overlapping
)
```

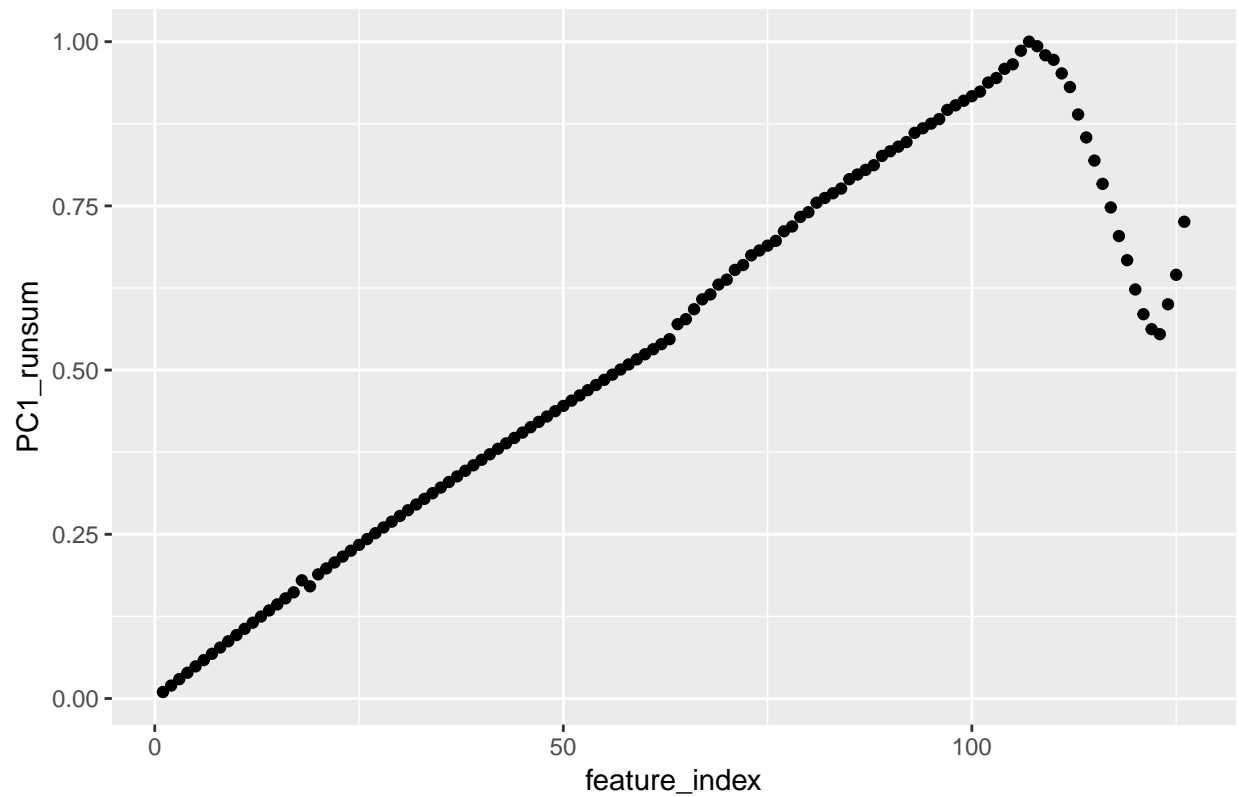


```
# creating extra columns
aload <- abs(pca_result$rotation[,1:2])
components <- sweep(aload, 2, colSums(aload), "/")
components <- as.data.frame(components)
components$feature_name <- rownames(components)
components$feature_index <- 1:nrow(components)

components <- components %>% arrange(-PC1)
components$PC1_runsum <- cumsum(components$PC1)

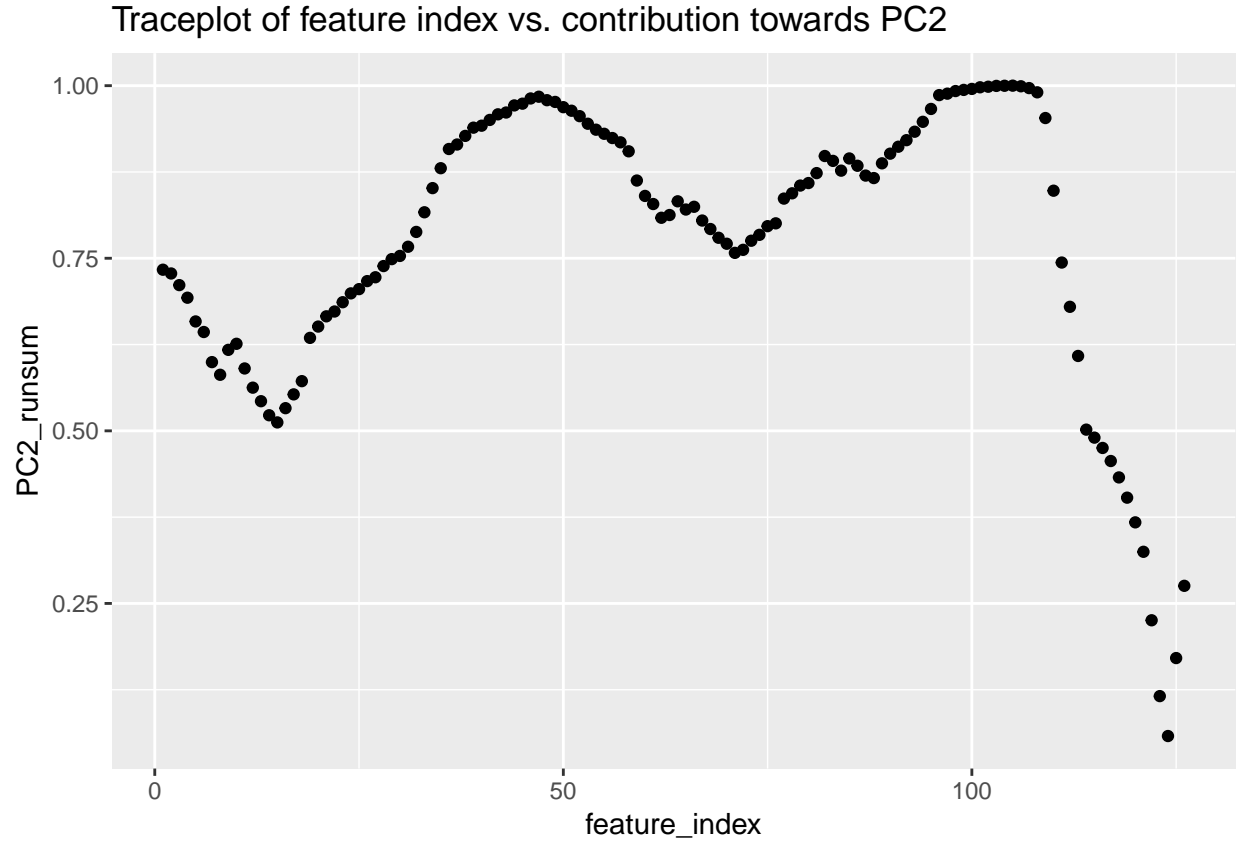
ggplot(data = components, aes(x = feature_index, y = PC1_runsum)) +
  geom_point() +
  ggtitle("Traceplot of feature index vs. contribution towards PC1")
```

Traceplot of feature index vs. contribution towards PC1



```
components <- components %>% arrange(-PC2)
components$PC2_runsum <- cumsum(components$PC2)

ggplot(data = components, aes(x = feature_index, y = PC2_runsum)) +
  geom_point() +
  ggtitle("Traceplot of feature index vs. contribution towards PC2")
```



```
knitr::kable(components[1:10,],
              caption = "Contribution of Features towards the Principle Components")
```

Table 2: Contribution of Features towards the Principle Components

PC1	PC2	feature_name	feature_index	PC1_runsum	PC2_runsum
0.0075609	0.0579423	X996	124	0.6001894	0.0579423
0.0076268	0.0578055	X994	123	0.5546874	0.1157479
0.0074384	0.0551408	X998	125	0.6450879	0.1708887
0.0076192	0.0547439	X992	122	0.5623067	0.2256326
0.0072706	0.0498894	X1000	126	0.7258963	0.2755220
0.0075625	0.0492868	X990	121	0.5850672	0.3248088
0.0074820	0.0426587	X988	120	0.6227157	0.3674674
0.0073922	0.0357468	X986	119	0.6673233	0.4032143
0.0073049	0.0293427	X984	118	0.7040384	0.4325569
0.0072213	0.0237460	X982	117	0.7476261	0.4563030

Analysis:

From the above three plots we see that towards the first principle components axis (93.3% variance accounted) the feature index till 110 are the main contributors(positive contribution), while for the second component(6.3% variance accounted for) we have columns_index 25-45 and 85-100 as the main components.

The corresponding feature name can be accessed by viewing the table used for plot, a sample of the few columns is shown above.

Appendix

```
knitr::opts_chunk$set(echo = TRUE)
if (!require("pacman")) install.packages("pacman")
pacman::p_load(xlsx, ggplot2, MASS, tidyr, dplyr, reshape2, gridExtra,
               tree, caret, e1071, pROC, boot, factoextra)

set.seed(12345)
options("jtools-digits" = 2, scipen = 999)
set.seed(12345)
credit_data <- read.xlsx("creditscoring.xls", sheetName = "credit")
credit_data$good_bad <- as.factor(credit_data$good_bad)

n=NROW(credit_data)
set.seed(12345)
id=sample(1:n, floor(n*0.4))
train=credit_data[id,]

id1=setdiff(1:n, id)
set.seed(12345)
id2=sample(id1, floor(n*0.3))
valid=credit_data[id2,]

id3=setdiff(id1,id2)
test=credit_data[id3,]
set.seed(12345)

# Create a decision tree model
credit_tree_deviance <- tree(good_bad~., data=train, split = c("deviance"))
credit_tree_gini <- tree(good_bad~., data=train, split = c("gini"))

# Visualize the decision tree with rpart.plot
summary(credit_tree_deviance)
summary(credit_tree_gini)

# predicting on the test dataset to get the misclassification rate.
predict_tree_deviance <- predict(credit_tree_deviance, newdata = test, type = "class")
predict_tree_gini <- predict(credit_tree_gini, newdata = test, type = "class")

conf_tree_deviance <- table(test$good_bad, predict_tree_deviance)
names(dimnames(conf_tree_deviance)) <- c("Actual Test", "P2redicted Test")
caret::confusionMatrix(conf_tree_deviance)

conf_tree_gini <- table(test$good_bad, predict_tree_gini)
names(dimnames(conf_tree_gini)) <- c("Actual Test", "Predicted Test")
caret::confusionMatrix(conf_tree_gini)
set.seed(12345)

credit_tree <- tree(good_bad~., data=train, split = c("deviance"))

credit_tree_purned_train <- prune.tree(credit_tree, method = c("deviance"))
credit_tree_purned_valid <- prune.tree(credit_tree, newdata = valid ,method = c("deviance"))
```

```

result_train <- cbind(credit_tree_purned_train$size,
                     credit_tree_purned_train$dev, "Train")
result_valid <- cbind(credit_tree_purned_valid$size,
                     credit_tree_purned_valid$dev, "Valid")

result <- as.data.frame(rbind(result_valid, result_train))
colnames(result) <- c("Leaf", "Deviance", "Type")

result$Leaf <- as.numeric(as.character(result$Leaf))
result$Deviance <- as.numeric(as.character(result$Deviance))

# plot of deviance vs. number of leafs
ggplot(data = result, aes(x = Leaf, y = Deviance, colour = Type)) +
  geom_point() + geom_line() +
  ggtitle("Plot of Deviance vs. Tree Depth (Shows Deviance least at 7)")

# prune the tree to the required depth
credit_tree_sniped <- prune.tree(credit_tree, best=7)

plot(credit_tree_sniped)
text(credit_tree_sniped)

# misclassification rate for best pruned tree

result_prune_test <- predict(credit_tree_sniped, newdata = test, type = "class")

conf_prune_tree_test <- table(test$good_bad, result_prune_test)
names(dimnames(conf_prune_tree_test)) <- c("Actual Test", "Predicted Test")
caret::confusionMatrix(conf_prune_tree_test)

#Fitting the Naive Bayes model
credit_naive_model = naiveBayes(good_bad ~., data=train)
credit_naive_model

#Prediction on the dataset
predict_naive_train = predict(credit_naive_model, newdata=train, type = "class")
predict_naive_test = predict(credit_naive_model, newdata=test, type = "class")

conf_naive_train <- table(train$good_bad, predict_naive_train)
names(dimnames(conf_naive_train)) <- c("Actual Train", "Predicted Train")
caret::confusionMatrix(conf_naive_train)

conf_naive_test <- table(test$good_bad, predict_naive_test)
names(dimnames(conf_naive_test)) <- c("Actual Test", "Predicted Test")
caret::confusionMatrix(conf_naive_test)
set.seed(12345)

credit_tree <- tree(good_bad~., data=train, split = c("deviance"))
credit_naive_model = naiveBayes(good_bad ~., data=train)

# prune the tree to the required depth
credit_tree_sniped <- prune.tree(credit_tree, best=7)

# predicting class, getting probability

```



```

predict_prune_test_prob <- predict(credit_tree_sniped, newdata = test)
predict_naive_test_prob <- predict(credit_naive_model, newdata=test, type = "raw")

# data mugging
probability_data_naive <- as.data.frame(cbind(predict_naive_test_prob,
                                             as.character(test$good_bad), "naivebayes"))
probability_data_tree <- as.data.frame(cbind(predict_prune_test_prob,
                                             as.character(test$good_bad), "tree"))

probability_data_combined <- rbind(probability_data_tree, probability_data_naive)
colnames(probability_data_combined) <- c("prob_bad", "prob_good",
                                         "actual_test_class", "model")

# final dataset
probability_data_combined$prob_good <- as.numeric(as.character(probability_data_combined$prob_good))

# changing the threshold and printing the probability

tree_list <- NULL
naive_list <- NULL
final <- NULL
for(threshold in seq(from = 0.05, to = 0.95, by = 0.05)){
  probability_data_combined$predicted_class <- ifelse(probability_data_combined$prob_good > threshold,

  df2 <- probability_data_combined[,c("model", "actual_test_class", "predicted_class")]
  df2$threshold <- threshold
  df2$match <- ifelse(df2$actual_test_class == df2$predicted_class, 1, 0)

  final <- rbind(df2, final)
}

# Creating the FRP and TRP for each model and threshold
final$temp <- 1
final_summary <- final %>%
group_by(model, threshold) %>%
summarise(total_positive = sum(temp[actual_test_class == "good"]),
          total_negative = sum(temp[actual_test_class == "bad"]),
          correct_positive = sum(temp[actual_test_class == "good" & predicted_class == "good"]),
          false_positive = sum(temp[actual_test_class == "bad" & predicted_class == "good"])) %>%
  mutate(TPR = correct_positive/total_positive, FPR = false_positive/total_negative) %>%
  select(model, threshold, TPR, FPR)

ggplot(data = final_summary, aes(x = FPR, y=TPR)) + geom_line(aes(colour = model)) +
  ggtitle("ROC curve for the Naive Bayes vs. Tree Model")

set.seed(12345)

credit_naive_model = naiveBayes(good_bad ~., data=train)

# predicting class, getting probability
predict_naive_train_prob <- predict(credit_naive_model, newdata=train, type = "raw")
predict_naive_test_prob <- predict(credit_naive_model, newdata=test, type = "raw")

```

```

train <- cbind(predict_naive_train_prob, train)
test <- cbind(predict_naive_test_prob, test)

# class based on the loss matrix
train$predicted_class <- ifelse(train$good > 10*train$bad, "good", "bad")
test$predicted_class <- ifelse(test$good > 10*test$bad, "good", "bad")

# confusion matrix
conf_naive_train <- table(train$good_bad, train$predicted_class)
names(dimnames(conf_naive_train)) <- c("Actual Train", "Predicted Train")
caret::confusionMatrix(conf_naive_train)

conf_naive_test <- table(test$good_bad, test$predicted_class)
names(dimnames(conf_naive_test)) <- c("Actual Test", "Predicted Test")
caret::confusionMatrix(conf_naive_test)

rm(list=ls())

set.seed(12345)
state_data <- read.csv2("state.csv")

state_data <- state_data %>% arrange(MET)

ggplot(data = state_data, aes(x=MET, y = EX)) +
  geom_point() +
  geom_smooth(method = 'loess') +
  ggtitle("Plot of MET vs. EX")
set.seed(12345)

state_tree_regression <- tree(data = state_data, EX~MET,
                             control = tree.control(nobs=NROW(state_data),
                                                      minsize = 8))

state_cv_tree <- cv.tree(state_tree_regression, FUN = prune.tree)
plot(state_cv_tree)
# The best size is either 3 or 4

# purging the tree for leaf size of 3
state_cv_tree_purned <- prune.tree(state_tree_regression, k = 3)
plot(state_cv_tree_purned, main="Pruned Tree for the given dataset")
text(state_cv_tree_purned)

# Original vs. Fitted values
compare_data <- predict(state_cv_tree_purned, newdata = state_data)
compare_data <- cbind(compare_data, state_data$EX)
compare_data <- as.data.frame(compare_data)
colnames(compare_data) <- c("predicted_value", "actual_value")
compare_data$residual <- compare_data$actual_value - compare_data$predicted_value

# plots
ggplot(compare_data, aes(x = actual_value, y = predicted_value)) +
  geom_point() +

```

```

ggtitle("Plot of Actual vs. Predicted Value")

ggplot(compare_data, aes(x = predicted_value, y = residual)) +
  geom_point() + geom_abline(slope=0, intercept=0) +
  ggtitle("Plot of Predicted Value vs. Residual")

ggplot(data = compare_data, aes(x = residual)) +
  geom_histogram(aes(y = ..density..), binwidth = 8) +
  geom_density(colour = "red") +
  ggtitle("Histogram of Residual")

set.seed(12345)

# function to obtain R-Squared from the data
rsq <- function(data, indices) {
  d <- data[indices,] # allows boot to select sample

  model <- tree(data = d,
                EX~MET,
                control = tree.control(nobs=NROW(data), minsize = 8))

  model_purned <- prune.tree(model, k = 3)

  compare_data <- predict(model_purned, newdata = d)

  return(compare_data)
}

# bootstrapping with 1000 replications
results <- boot(data=state_data, statistic=rsq, R=1000)

e=envelope(results) #compute confidence bands

compare_data_non_para_boot <- compare_data
compare_data_non_para_boot$MET <- state_data$MET
compare_data_non_para_boot$lower_bound <- e$point[2,]
compare_data_non_para_boot$upper_bound <- e$point[1,]

ggplot(data=compare_data_non_para_boot, aes(x = MET, y = actual_value)) +
  geom_point(aes(y=predicted_value)) +
  geom_line(aes(y=upper_bound), colour="red") + # first layer
  geom_line(aes(y=lower_bound), colour="red") +
  ggtitle("Predicted EX value along with 95% CI band")

set.seed(12345)

set.seed(12345)

rm(list=ls())
NIR_data <- read.csv2("NIRSpectra.csv")
set.seed(12345)

```

```

pca_data = select(NIR_data, -c(Viscosity))
pca_result = prcomp(pca_data)

contribution <- summary(pca_result)$importance
knitr::kable(contribution[,1:5],
              caption = "Contribution of PCA axis towards variance explanation")

# plots PCA components and the eigen vectors
factoextra::fviz_eig(pca_result)

# pca components and the viscosity
pca_result_data = cbind(first_component = pca_result$x[,1],
                        second_component = pca_result$x[,2],
                        Viscosity = NIR_data$Viscosity)

pca_result_data = as.data.frame(pca_result_data)

# plotting the data variation and the viscosity
ggplot(data = pca_result_data, aes(x = first_component, y = second_component)) +
  geom_point(aes(y = Viscosity)) + ggtitle("PCA components vs. Viscosity")
# showing the components of the PCA components
factoextra::fviz_pca_var(pca_result,
                        col.var = "contrib", # Color by contributions to the PC
                        gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"),
                        repel = TRUE      # Avoid text overlapping
                        )

# creating extra columns
aload <- abs(pca_result$rotation[,1:2])
components <- sweep(aload, 2, colSums(aload), "/")
components <- as.data.frame(components)
components$feature_name <- rownames(components)
components$feature_index <- 1:nrow(components)

components <- components %>% arrange(-PC1)
components$PC1_runsum <- cumsum(components$PC1)

ggplot(data = components, aes(x = feature_index, y = PC1_runsum)) +
  geom_point() +
  ggtitle("Traceplot of feature index vs. contribution towards PC1")

components <- components %>% arrange(-PC2)
components$PC2_runsum <- cumsum(components$PC2)

ggplot(data = components, aes(x = feature_index, y = PC2_runsum)) +
  geom_point() +
  ggtitle("Traceplot of feature index vs. contribution towards PC2")

knitr::kable(components[1:10,],
              caption = "Contribution of Features towards the Principle Components")

```