

Problem 1 Part a)

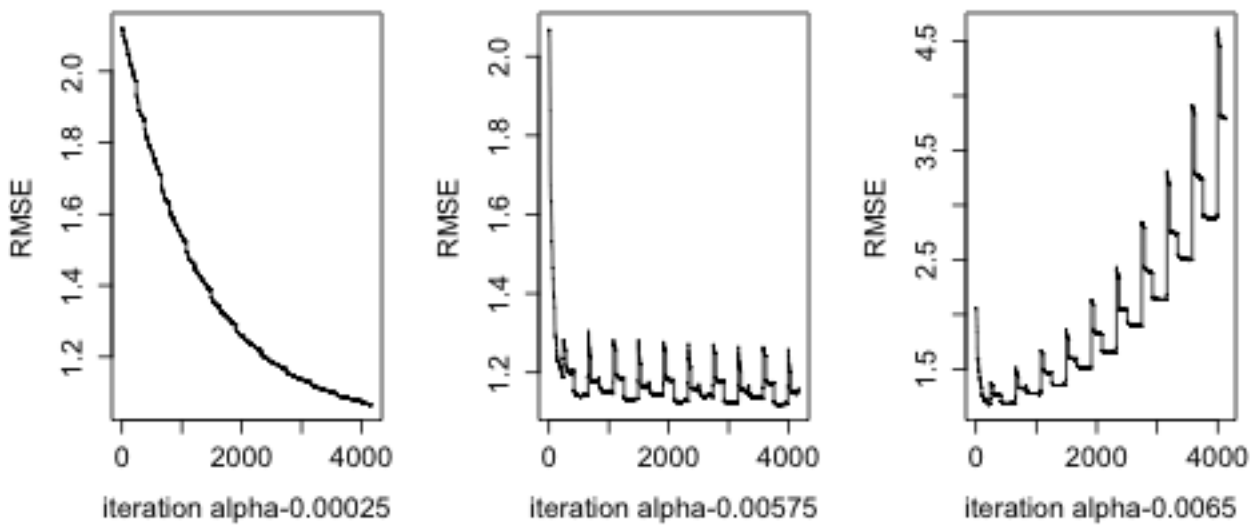


Image indicating the Error values for the various values of Alpha

Training Number	
α	Iteration Number
0.00025	4170
0.00575	4137
0.0065	4133

Error Values	
Model	MSE
<i>MLR</i>	1564
<i>SGD</i> $\alpha = 0.00025$	41.6
<i>SGD</i> $\alpha = 0.00575$	69.2
<i>SGD</i> $\alpha = 0.0065$	69.8

The SGD model has been able to reduce the amount of error by a huge factor. This is because now the coefficients are being adjusted in a manner which reduces error value. The smaller the value of α the smaller the MSE. This is because α determines the step size, and in SGD, the smaller the step size, the better the approximation. However, as the value of α reduces the number of iterations increases, thus increasing the runtime. When the value of the threshold is reduced, from 0.1 to 0.01, the number of iterations needed for the training to complete jumps from the range to 4000 to the range of 40000. The behavior of SGD depends a lot on the choice of gradient function and α . Therefore for different gradient functions, for same value of α we might get different iteration numbers when the solution would converge.

Problem 2 For a n-D Space let's say that the convex hulls intersect and thus there exists a z for two such set of points. The convex hull of such a set of points can be represented as :

$$\sum_i \alpha_i x_i \text{ where } \alpha_i > 0 \text{ and } \sum_i \alpha_i = 1$$

Thus z can be represented as

$$z = \sum_i \alpha_i x_i = \sum_j \gamma_j y_j$$

$$\text{and } \sum_i \alpha_i = \sum_j \gamma_j = 1$$

If these two sets of points are linearly separable, there exists a separating hyper-plane determined by w, w_0 such that

$$w^T x_i + w_0 > 0 \forall i$$

$$w^T y_j + w_0 < 0 \forall j$$

$$\begin{aligned} w^T z + w_0 &= w^T \sum_i \alpha_i x_i + w_0 \\ &= \sum_i \alpha_i w^T x_i + w_0 \sum_i \alpha_i \text{ as } \sum_i \alpha_i = 1 \\ &= \sum_i \alpha_i (w^T x_i + w_0) \end{aligned}$$

Similarly

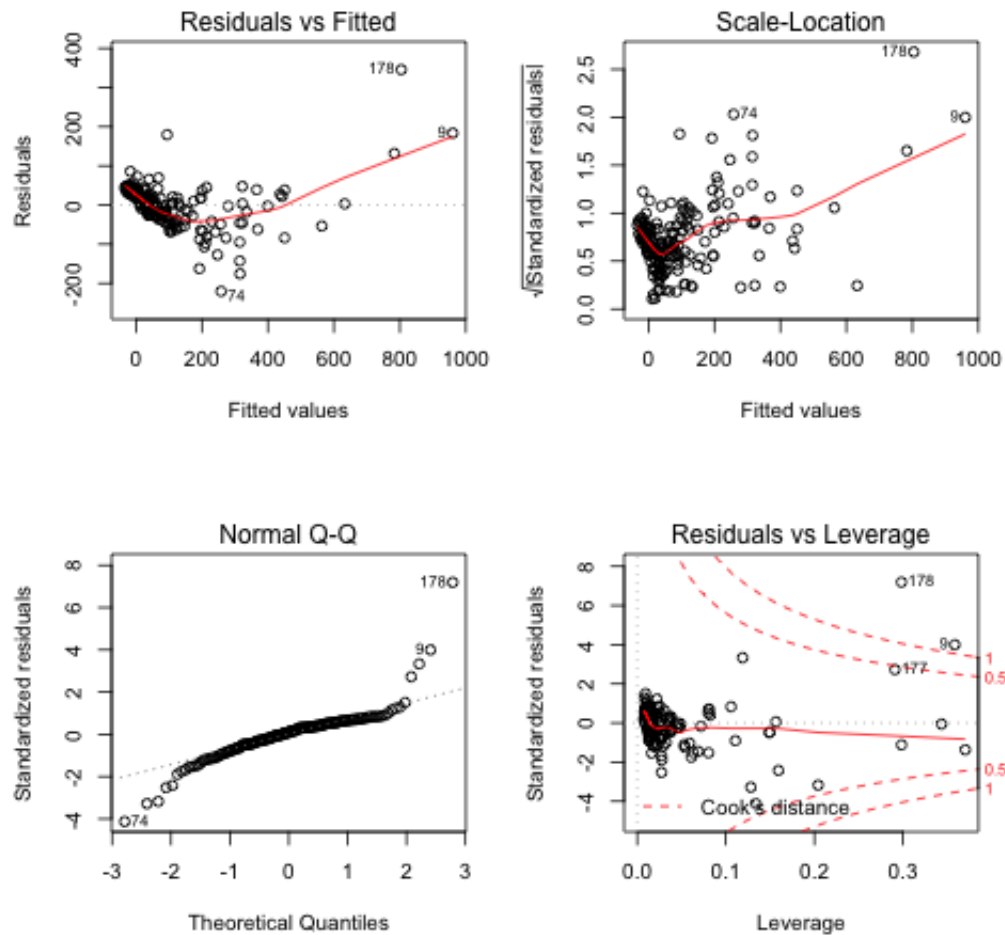
$$\begin{aligned} w^T z + w_0 &= \sum_i \gamma_i (w^T y_i + w_0) \\ \sum_i \alpha_i (w^T x_i + w_0) &= \sum_i \gamma_i (w^T y_i + w_0) \\ \sum_i \alpha_i (w^T x_i + w_0) &= \sum_i \gamma_i (w^T y_i + w_0) \end{aligned}$$

$$\text{This implies } \alpha_i = \gamma_j = 0 \forall i, j$$

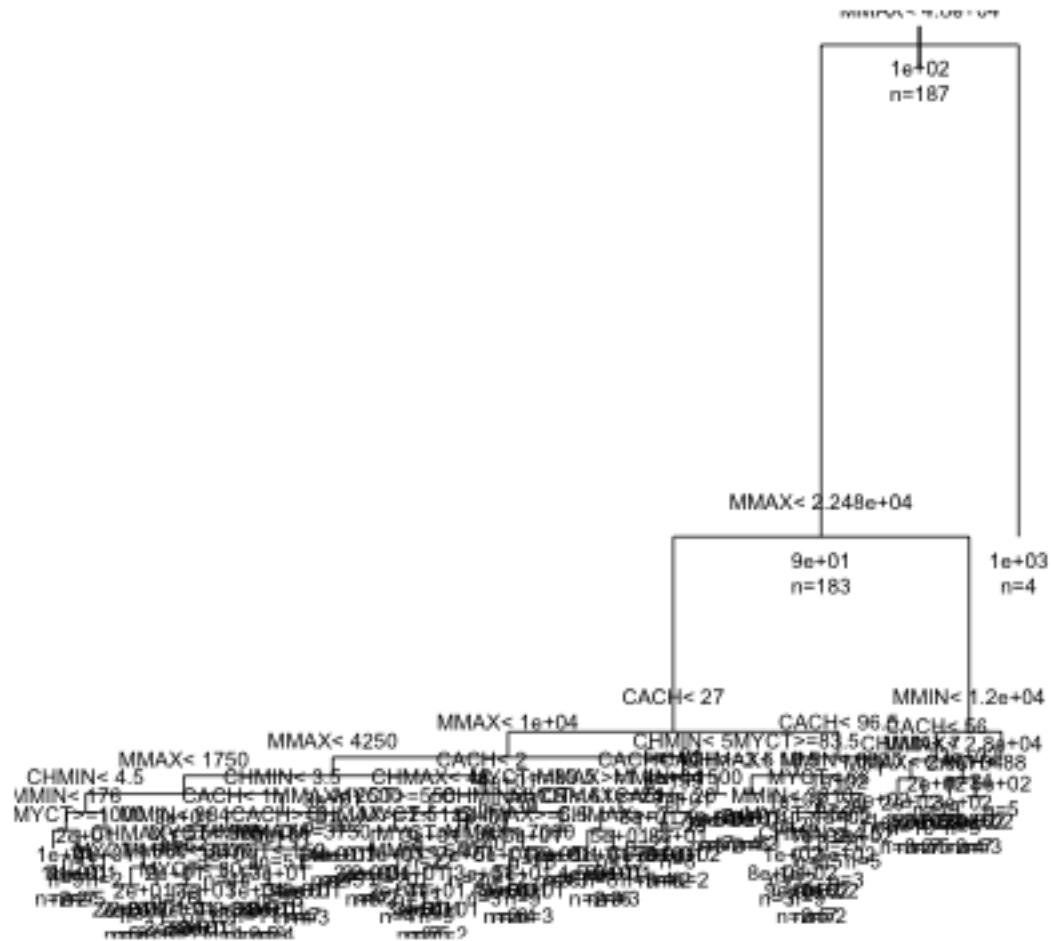
However this contradicts our previous assumption. Thus assuming the two sets of points are linearly separable leads us to a contradiction. Thus if the convex hulls of two separate sets of points intersect \Rightarrow the two sets are not linearly separable.

Problem 3 Part a) A Regression tree is first generated with $cp = 0$. This creates an overfitted model. And then the tree is pruned so that the overfitting can be removed. The pruned tree is used for the computations needed in the rest of the parts.

RT	MSE	MAE	R^2	RMSE
$Train = 1, Test = 1$	2053.3	28.2	0.924	45.3

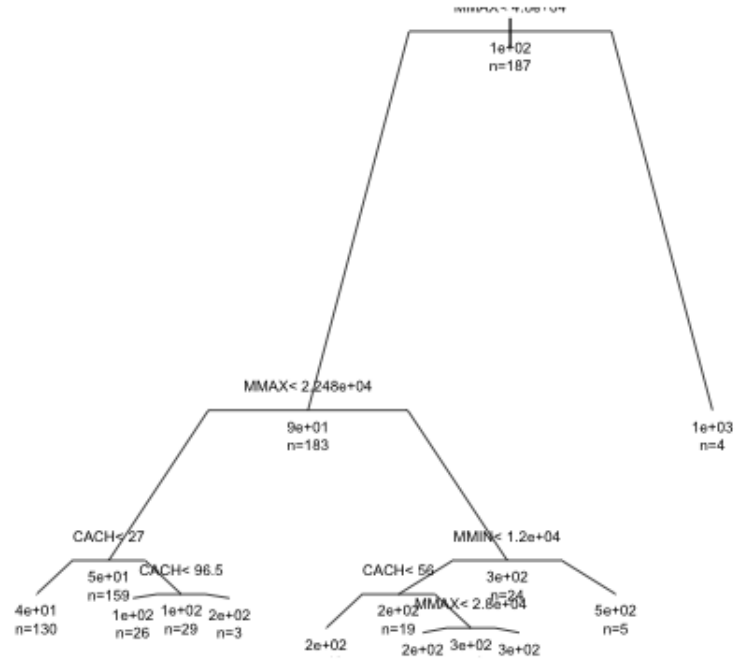


MLR PLOT

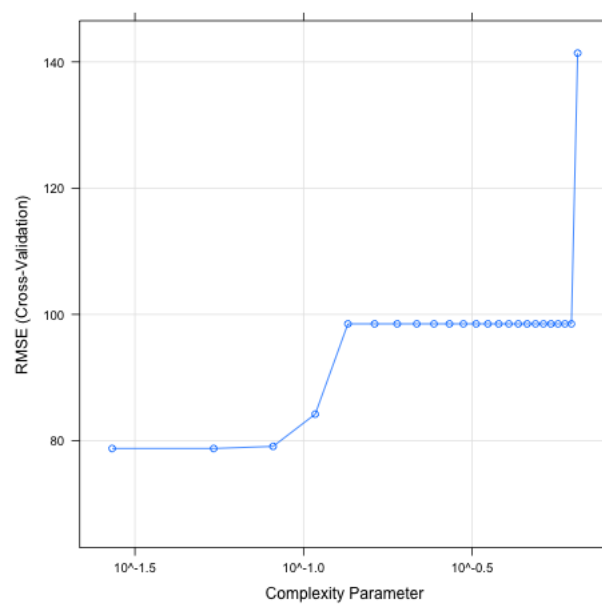


OverFit Regression tree

Part b)



Pruned Regression tree to overcome the Overfitting done in previous step after choosing the value of cp which corresponds to minimum error value.

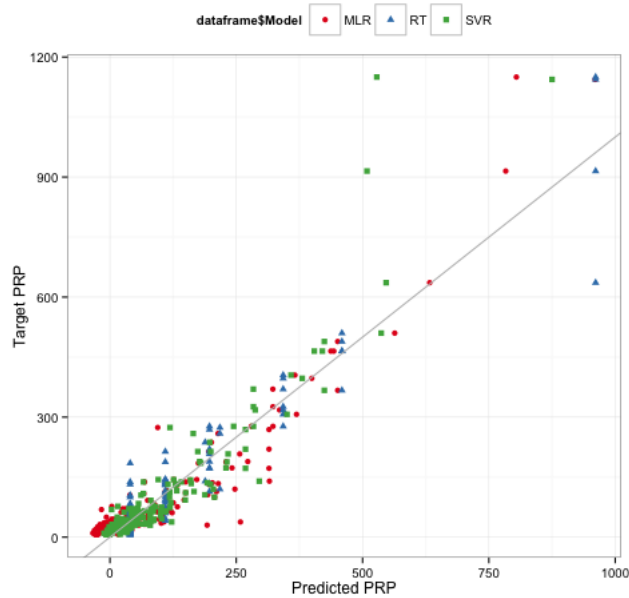


5-fold Cross Validation results for the Regression tree plotted for the various values of complexity parameter

Part c)

Error Values				
SVR	MSE	MAE	R^2	RMSE
$Train = 1, Test = 1$	4640	32.6	0.85	68.1
$Train = 1, Test = test$	3938.7	33.0	0.761	62.8
$Train = 2, Test = test$	321276	1136	0.742	567

Part d)



ggPlot showing the combined plots of the three models

With the addition of outliers the overall error values increased for all the three models. However out of the three models, MLR was affected the most with the addition of outliers which is expected. Because SVR makes use of Huber loss function the increase in the error values for SVR is not huge. Regression trees have the least recorded error values for all three types of test cases out of the three models. Thus, Regression trees had the best performance with given tuning. The choice of complexity of the model, will affect the performance of this model.

Error Values				
MLR	MSE	MAE	R^2	RMSE
$Train = 1, Test = 1$	3166.4	37.9	0.882	56.3
$Train = 1, Test = test$	5475.4	44.1	0.69	74.0
$Train = 2, Test = test$	394011	1247	0.622	628
RT	MSE	MAE	R^2	RMSE
$Train = 1, Test = 1$	2053.3	28.2	0.924	45.3
$Train = 1, Test = test$	3054.8	40.4	0.817	55.3
$Train = 2, Test = test$	401037	1114	0.565	633
SVR	MSE	MAE	R^2	RMSE
$Train = 1, Test = 1$	4640	32.6	0.85	68.1
$Train = 1, Test = test$	3938.7	33.0	0.761	62.8
$Train = 2, Test = test$	321276	1136	0.742	567

Part f)

Error Values				
Random Forest	MSE	MAE	R^2	RMSE
$Train = 1, Test = 1$	793.1	14.1	0.973	28.3
RT	MSE	MAE	R^2	RMSE
$Train = 1, Test = 1$	2053.3	28.2	0.924	45.3

The performance of Random Forests is even more promising than that of the Regression Tress.

Problem 4 $P(C_1|X) = \frac{P(C_1)P(X|C_1)}{P(X)}$

The division line occurs where $P(C_1|X) = P(C_2|X)$

$P(C_1) = 1/5$ and $P(C_2) = 4/5$

$P(X|C_1) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{(X-\mu)(X-\mu)^T}{2\sigma^2}}$ where σ = variance of C_1 which is 1 (because the covariance matrix for C_1 and C_2 is an identity matrix) $\mu = [2 \ 0]$

Similarly $P(X|C_2) = \frac{1}{\sqrt{2\pi}}e^{-\frac{(X-\mu)(X-\mu)^T}{2}}$ where σ = variance of C_2 which is 1 $\mu = [0 \ 2]$

Thus $\frac{1}{5} * \frac{1}{\sqrt{2\pi}}e^{-\frac{(X-\mu_1)(X-\mu_1)^T}{2}} = \frac{4}{5} \frac{1}{\sqrt{2\pi}}e^{-\frac{(X-\mu_2)(X-\mu_2)^T}{2}}$
 $e^{-\frac{(X-\mu_1)(X-\mu_1)^T}{2}} = 4 * e^{-\frac{(X-\mu_2)(X-\mu_2)^T}{2}}$ Taking log of both sides.

$-\frac{(X-\mu_1)(X-\mu_1)^T}{2} = 2\log 2 + -\frac{(X-\mu_2)(X-\mu_2)^T}{2}$
 $-(X-\mu_1)(X-\mu_1)^T = 4\log 2 + -(X-\mu_2)(X-\mu_2)^T$

Replacing the value of X by x,y and values of μ_1 and μ_2

$(x-2)^2 + y^2 = -4\log 2 + x^2 + (y-2)^2$
 $= x^2 - 4x + 4 + y^2 = -4\log 2 + x^2 + y^2 - 4y + 4$
 $-x = -\log 2 - y$

$\log 2 = x - y$

Part b)

Now because the cost function has changed , the division line will occurs where $2P(C_2|X) - P(C_1|X) = 0$ This is because the total cost for misclassification is : $2k * P(C_2|X) + k * P(C_1|X)$.

We need to minimize this value. Thus it is ok to misclassify a C_2 value as C_1 unless the penalty of misclassifying is not twice the cost.

Using the values given above and the above equation we get.

$$\begin{aligned} \frac{1}{5} * \frac{1}{\sqrt{2\pi}} e^{-\frac{(X-\mu_1)(X-\mu_1)^T}{2}} &= 2 * \frac{4}{5} \frac{1}{\sqrt{2\pi}} e^{-\frac{(X-\mu_2)(X-\mu_2)^T}{2}} \\ e^{-\frac{(X-\mu_1)(X-\mu_1)^T}{2}} &= 8 * e^{-\frac{(X-\mu_2)(X-\mu_2)^T}{2}} \text{ Taking log of both sides.} \\ -\frac{(X-\mu_1)(X-\mu_1)^T}{2} &= \log 8 - \frac{(X-\mu_2)(X-\mu_2)^T}{2} \\ (X-\mu_1)(X-\mu_1)^T &= -2\log 8 + (X-\mu_2)(X-\mu_2)^T \\ \text{Replacing the value of X by x,y and values of } \mu_1 \text{ and } \mu_2 & \\ (x-2)^2 + y^2 &= -2\log 8 + x^2 + (y-2)^2 \\ x^2 - 4x + 4 + y^2 &= -2\log 8 + x^2 + y^2 - 4y + 4 \\ -2x &= -\log 8 - 2y \\ \log 8 &= 2x - 2y \end{aligned}$$

Problem 5 a)

$$\begin{aligned} P(B = \text{good}, F = \text{empty}, G = \text{empty}, S = \text{yes}) &= P(B = \text{good}) * P(F = \text{empty}) * P(G = \text{empty} | B = \text{good}, F = \text{empty}) * P(S = \text{yes} | B = \text{good}, F = \text{empty}) \\ &= 0.9 * 0.2 * 0.8 * 0.2 = 0.00576 \end{aligned}$$

b)

$$\begin{aligned} P(B = \text{bad}, F = \text{empty}, G = \text{notempty}, S = \text{no}) &= P(B = \text{bad}) * P(F = \text{empty}) * P(G = \text{notempty} | B = \text{bad}, F = \text{empty}) * P(S = \text{no} | B = \text{bad}, F = \text{empty}) \\ &= 0.1 * 0.2 * 0.1 * 1.0 = 0.002 \end{aligned}$$

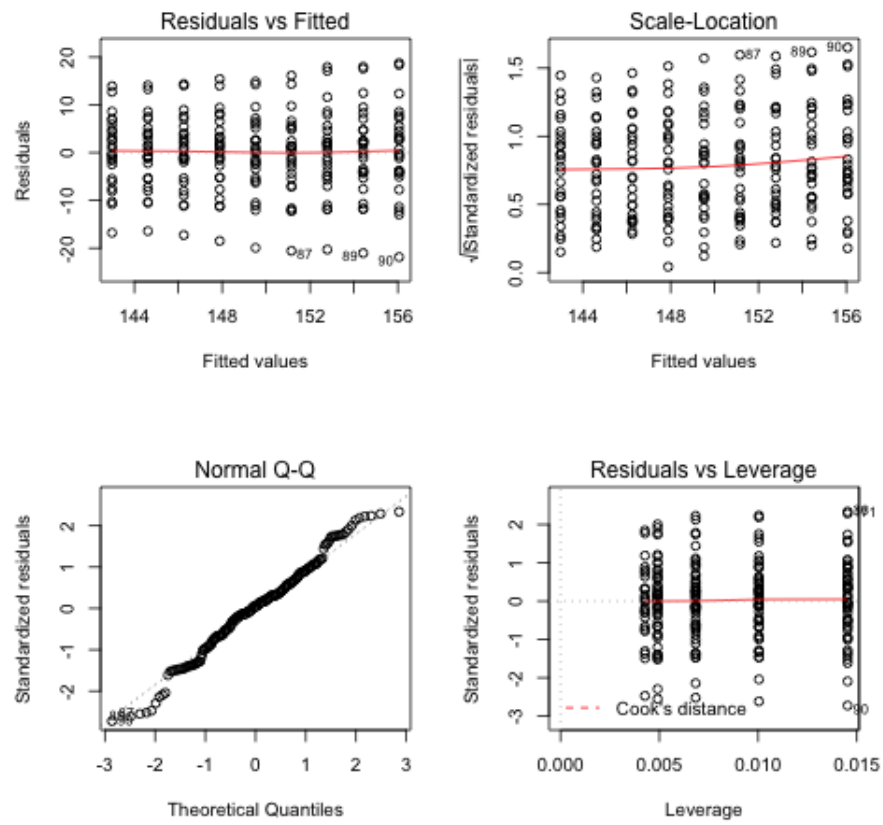
c) Given that the battery is bad, compute the probability that the car will start. =

$$\begin{aligned} P(S = \text{yes} | B = \text{Bad}) &= \sum_{g,f} P(S = \text{yes}, B = \text{Bad}, G = g, F = f) / P(B = \text{Bad}) \\ &= (P(S = \text{yes}, B = \text{bad}, F = \text{empty}, G = \text{empty}) + P(S = \text{yes}, B = \text{bad}, F = \text{empty}, G = \text{notempty}) + P(S = \text{yes}, B = \text{bad}, F = \text{notempty}, G = \text{empty}) + P(S = \text{yes}, B = \text{bad}, F = \text{notempty}, G = \text{notempty})) / P(B = \text{Bad}) \end{aligned}$$

Another way to calculate this can be

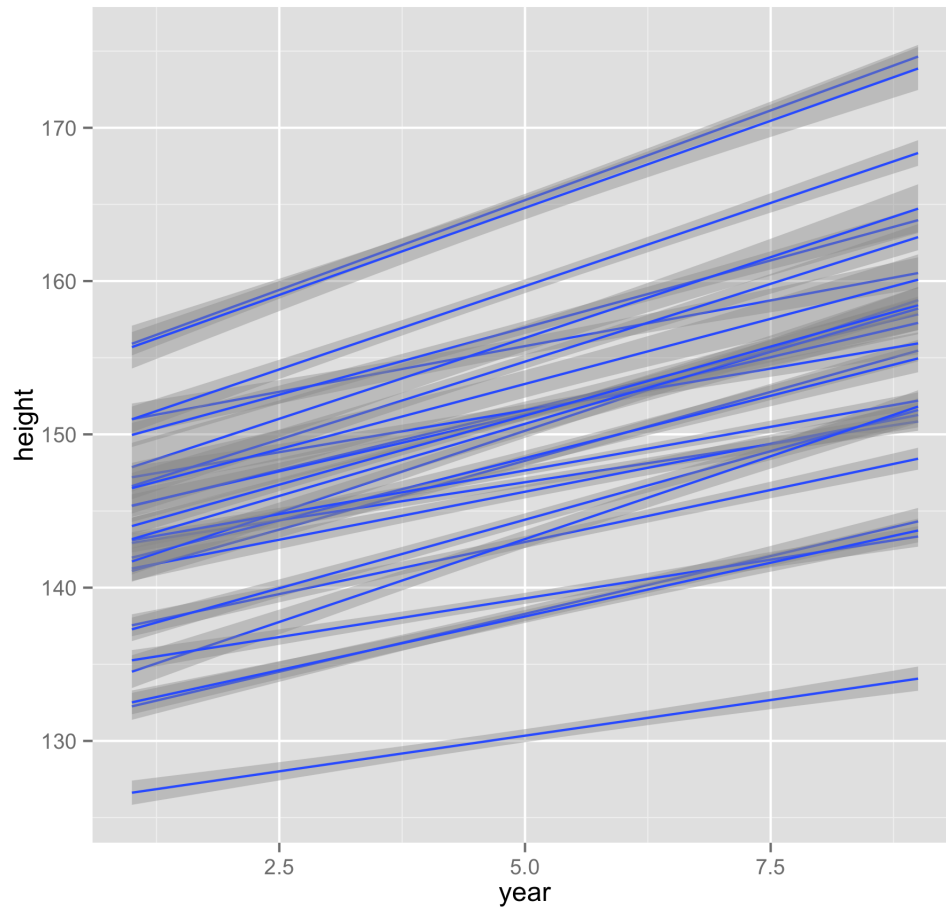
$$\begin{aligned} P(S = \text{yes} | B = \text{Bad}) &= P(F = \text{empty}) * P(S = \text{yes} | B = \text{bad}, F = \text{empty}) + P(F = \text{notempty}) * P(S = \text{yes} | B = \text{bad}, F = \text{notempty}) \\ &= (0.2 * 0 + 0.8 * 0.1) = 0.08 \end{aligned}$$

Problem 6 Part a)

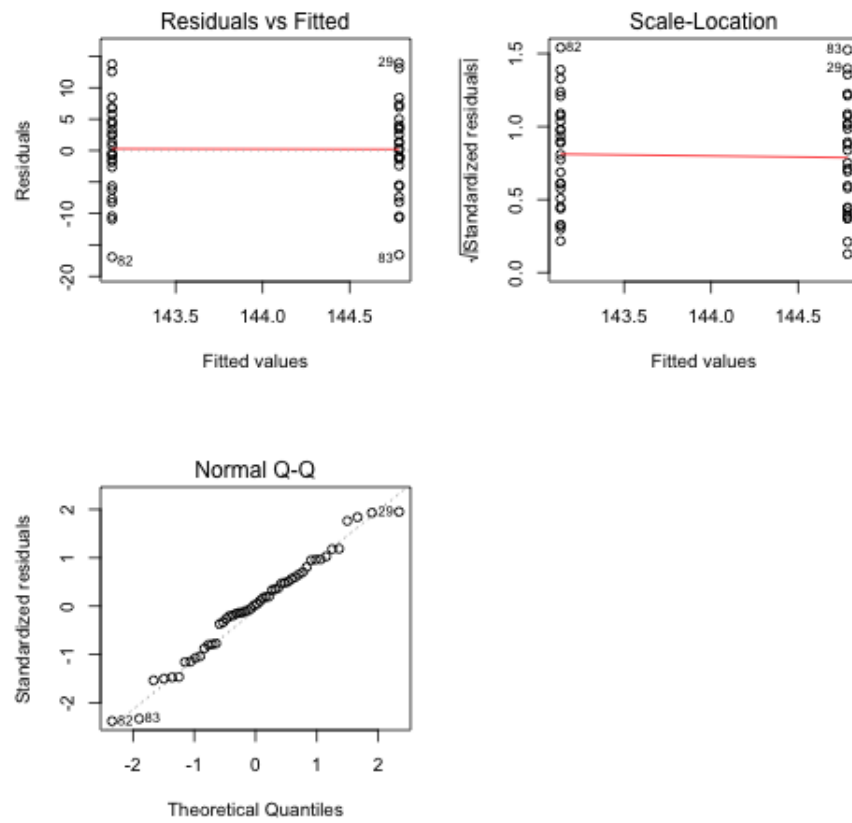


MLR where the model is ignoring the variable id

Part b) This image is able to capture all the multiple relations which can be used to define the given height and year relation. This model is grouping the values of the variable id.

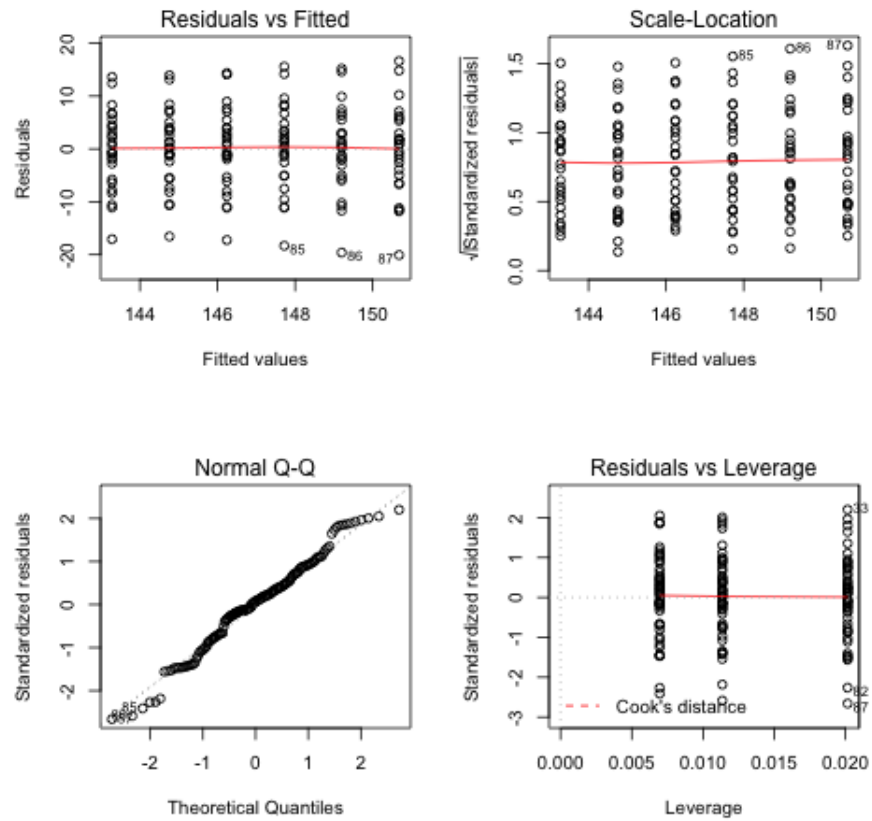


Part c)



Error Values	
Model	MSE
<i>Global</i>	68.7
<i>Local</i>	328
<i>MLM</i>	4.23

Part d)



Error Values	
Model	MSE
Global	80.1
Local	68.7
MLM	2.72

Part e) Looking at these values we can say that MLM can perform better only when there some amount of data for the various classes. When the Amount of training data is less, MLM are not ale to perform well. The Error value for the local model indicates similar detail that when the training value is small the models like MLM are unable to predict correctly and become unreliable. However, they are able to provide very satisfactory results with increase in a small amount of training data.

R - Code

Problem 1

```
1 rm(list = ls())
2 setwd("/Users/bansal/Google Drive/Spring2014/DataMining/HW3/")
3 library(MASS)
4 library(hydroGOF)
5 library(ggplot2)
6 test = read.csv("~/Google Drive/Spring2014/DataMining/HW3/HW3Data/forestfire-test.csv",
7               ,header=TRUE)
8 train = read.csv("~/Google Drive/Spring2014/DataMining/HW3/HW3Data/forestfire-train.csv",
9                 ,header=TRUE)
10
11 for (j in 1:length(train))
12 {
13   train[,j] = (train[,j] - mean(train[,j]))/sd(train[,j])
14   test[,j] = (test[,j] - mean(train[,j]))/sd(train[,j])
15 }
16
17 # squared error cost function
18 cost <- function(x, y, theta) {
19   return (rmse(as.vector(x * theta) , y))
20 }
21
22 #Gradient Function
23 gradFun <- function(xn, yn, theta) {
24   return (xn * (yn - xn * theta) )
25 }
26
27 # learning rate and iteration limit
28 #alphas <- c(0.00025)
29
30 alphas <- c(0.00025, 0.00575, 0.0065)
31 num.iterations <- 10
32 threshold <- 0.1
33
34 temp <- as.matrix(train)
35 x <- temp[,1:9]
36 x[,9] <- x[,9] - x[,9] + 1
37 y <- temp[,c="area"]
38
39 temp <- as.matrix(test)
40 xtest <- temp[,1:9]
41 xtest[,9] <- xtest[,9] - xtest[,9] + 1
42 ytest <- temp[,c="area"]
43
44 # gradient descent
45 gradient.descent <- function(X, Y, alpha=0.1, num.iterations=1000, threshold=0.1) {
46
47   # initialize coefficients
48   xn = X[1,]
49   yn = Y[1]
50   theta1 <- rep(0.5,9)
51   gradient <- gradFun(xn, yn, theta1)
52   theta2 <- theta1 + alpha * gradient
53   #Initiliaze History
```

```

53 error_trace <- NULL
   test_error_trace <- NULL
55 theta_trace <- theta2
   for (i in 1:num.iterations){
57     for (j in 1:417){
       xn = X[j,]
59       yn = Y[j]
       theta1 <- theta2
61       gradient <- gradFun(xn,yn, theta1)
       theta2 <- theta1 + alpha * gradient
63       cost <- rmse(as.vector(X * theta2) , Y)
       # Traces needed for the graph
65       theta_trace <- c(theta_trace,theta2)
       error_trace <- c(error_trace,cost)
67       test_error_trace <- c(test_error_trace, cost(xtest, ytest, theta2))

69       if(abs(theta2 - theta1) > threshold)
         break
71     }
   }
73   return (list(error= error_trace, theta= theta_trace, testError = test_error_trace))
}
75 # plot data and converging fit
vary.alpha <- lapply(alphas, function(alpha) gradient.descent(x,
77                                     y, alpha=alpha, num.
                                     iterations=num.iterations))

79 par(mfrow = c(2, 3))
   for (j in 1:3) {
81     plot((vary.alpha[[j]])$error, ylab="RMSE", xlab=paste("iteration alpha", alphas[j],
       sep="-"), type="l")
   }
83 #print graphs for the iterations
dev.copy(png, 'Q1plot2.png')
85 dev.off()

87 lm_model<- lm(area ~ FFMC + DMC + DC + ISI + temp + RH + wind + rain + 1, data=train)
summary(lm_model)
89
91 png("Q1plot_lm_all.png")
layout(matrix(1:4, ncol = 2))
plot(lm_model)
93 layout(1)
dev.off()
95 pred_lm <- predict(lm_model, newdata = test)
lm.mse <- regr.eval(test$area, pred_lm,
97                 stats = c("mse"),
                 train.y = NULL)
99 lm.mse

101 mse_alpha1<- (vary.alpha[[1]])$testError[4170]
mse_alpha1
103 mse_alpha2<- (vary.alpha[[2]])$testError[4134]
mse_alpha2
105 mse_alpha3<- (vary.alpha[[3]])$testError[4133]
mse_alpha3

```

Problem 3

```
#install.packages("randomForest")
2 rm(list = ls(all.names = TRUE))
  setwd("/Users/bansal/Google Drive/Spring2014/DataMining/HW3/")
4 library(MASS)
  library(hydroGOF)
6 library(ggplot2)
  library(rpart)
8 library(caTools)
  library(grid)
10 library(lattice)
  library(DMwR)
12 library(DAAG)
  library(caret)
14 library(randomForest)

16 train_data = read.csv("~/Google Drive/Spring2014/DataMining/HW3/HW3Data/Machine1.csv"
  , header = TRUE)
  test_data = read.csv("~/Google Drive/Spring2014/DataMining/HW3/HW3Data/Machinetest.
    csv", header = TRUE)
18
#regression
20 lm_model<- lm(PRP ~ MYCT + MMIN + MMAX + CACH + CHMIN + CHMAX , data=train_data)
  summary(lm_model)
22
  png("Q3plot_lm_all.png")
24 layout(matrix(1:4, ncol = 2))
  plot(lm_model)
26 layout(1)
  dev.off()
28
#Some additional values which can come handy
30 #lm.rmse<- sqrt(mse(lm_model$fitted.values, train_data$PRP))
  #lm.coef<-coefficients(lm_model) # model coefficients
32 #lm.ci<-confint(lm_model, level=0.95) # CIs for model parameters
  #lm.resd<-residuals(lm_model) # residuals
34 #lm.anova<-anova(lm_model) # anova table
  #lm.cov<-vcov(lm_model) # covariance matrix for model parameters
36 #lm.influ<-influence(lm_model) # regression diagnostics

38 pred_lm<-fitted(lm_model) # predicted values
  pred_lm_test<-predict(lm_model, newdata = test_data )
40 regr.eval(train_data$PRP, pred_lm,
  stats = c("mae", "mse", "rmse"),
  train.y = NULL)
42 cor(pred_lm, train_data$PRP)**2
44
  cv.lm(df=train_data, lm_model, m=5)
46 dev.copy(png, 'Q3plot_lm_cross.png')
  dev.off()
48
#Regression Tree RT
50 rt_model <- rpart(PRP ~ MYCT + MMIN + MMAX + CACH + CHMIN + CHMAX , data = train_data
```

```

,
control = rpart.control(xval = 5, minbucket = 2, cp = 0))
52 #use of cp and max depth to control over and under fitting of the model
plot(rt_model)
54 text(rt_model, use.n=TRUE, all=TRUE, cex=.5)
dev.copy(png, 'Q3plot-rt.png')
56 dev.off()

58 #Train the model using the Cross validation (upto 5)
indx <- createFolds(train_data$PRP, returnTrain = TRUE)
60 ctrl <- trainControl(method = "cv", index = indx, number = 5)
set.seed(187)

62 cartTune <- train(x = train_data, y = train_data$PRP,
64 method = "rpart",
tuneLength = 25,
66 preProc = c("center", "scale"),
trControl = ctrl)

68 cartTune
### Plot the tuning results
70 plot(cartTune, scales = list(x = list(log = 10)))
dev.copy(png, 'Q3plot-rt-crossValidation.png')
72 dev.off()

74 #visualise the Corss validation of the model
plotcp(rt_model)
76 dev.copy(png, 'Q3plot-rt-complexity.png')
dev.off()

78 #prune Chose the tree with ideal value of CP based on the Error values collected.
80 cptable <- as.data.frame(rt_model$cptable)
alpha <- cptable$CP[which.min(cptable$xerror)]
82 rt_pruned <- prune(rt_model, cp=alpha)
plot(rt_pruned, branch = 0.5 )
84 text(rt_pruned, use.n=TRUE, all=TRUE, cex=.5)
dev.copy(png, 'Q3plot-rt-prunedTree.png')
86 dev.off()

88 #Capture the Error Values
#Rsquare
90 rsq.rpart(rt_pruned)
pred_rt <- predict(rt_pruned)
92 cor(pred_rt, train_data$PRP)**2
pred_rt_test<-predict(rt_pruned, newdata = test_data )
94 cor(pred_rt_test, test_data$PRP)**2
#sqrt(mse(pred, train_data$PRP))
96 regr.eval(train_data$PRP, pred_rt,
stats = c("mae", "mse", "rmse"),
98 train.y = NULL)

100
#part c
102 library(kernlab)

104 svr <- ksvm(PRP ~ MYCT + MMIN + MMAX + CACH + CHMIN + CHMAX , data=train_data,
type="eps-svr", kernel="vanilladot", cross=5)

```



```

106 summary(svr)

108 pred_svr <- predict(svr, train_data )
pred_svr_test<-predict(svr, newdata = test_data )

110
112 regr.eval(train_data$PRP, pred_svr,
            stats = c("mae", "mse", "rmse"),
            train.y = NULL)
114 cor(pred_svr, train_data$PRP)**2

116 #Part d
pred <- pred_lm
118 pred <- append(pred, pred_rt)
pred <- append(pred, pred_svr)
120 target <- train_data$PRP
target <- append(target, train_data$PRP)
122 target <- append(target, train_data$PRP)
model <- rep("MLR", dim(train_data)[1])
124 model <- append(model, rep("RT", dim(train_data)[1]))
model <- append(model, rep("SVR", dim(train_data)[1]))

126
dataframe <- data.frame(Pred = pred, Target = target, Model = model)
128 ggplot(data.frame(dataframe), aes(x = dataframe$Pred, y = dataframe$Target, color=
      dataframe$Model, shape=dataframe$Model)) + theme_bw() +
      geom_point() + geom_abline(intercept=0, slope=1, color="grey") +
130 scale_colour_brewer(palette="Set1") + xlab("Predicted PRP") +
      ylab("Target PRP") + theme(legend.position = "top")
132 dev.copy(png, 'Q3plot-part-d.png')
dev.off()

134

136 #part e
regr.eval(test_data$PRP, pred_lm_test,
138         stats = c("mae", "mse", "rmse"),
         train.y = NULL)
140 cor(pred_lm_test, test_data$PRP)**2
regr.eval(test_data$PRP, pred_rt_test,
142         stats = c("mae", "mse", "rmse"),
         train.y = NULL)
144 cor(pred_rt_test, test_data$PRP)**2
regr.eval(test_data$PRP, pred_svr_test,
146         stats = c("mae", "mse", "rmse"),
         train.y = NULL)
148 cor(pred_svr_test, test_data$PRP)**2

150 train_data2 = read.csv("~/Google Drive/Spring2014/DataMining/HW3/HW3Data/Machine2.csv",
      header = TRUE)

152 lm_model<- lm(PRP ~ MYCT + MMIN + MMAX + CACH + CHMIN + CHMAX , data=train_data2)
cv.lm(df=train_data, lm_model, m=5)
154 pred_lm_test<-predict(lm_model, newdata = test_data )
regr.eval(test_data$PRP, pred_lm,
156         stats = c("mae", "mse", "rmse"),
         train.y = NULL)
158 cor(pred_lm_test, test_data$PRP)**2

```

```

160 #Regression Tree RT
    rt_model <- rpart(PRP ~ MYCT + MMIN + MMAX + CACH + CHMIN + CHMAX , data = train_
        data2,
162                     control = rpart.control(xval = 5, minbucket = 2, cp = 0))
    #prune Chose the tree with ideal value of CP based on the Error values collected.
164 cptable <- as.data.frame(rt_model$cptable)
    alpha <- cptable$CP[which.min(cptable$xerror)]
166 rt_pruned <- prune(rt_model, cp=alpha)

168 #Rsquare
    rsq.rpart(rt_pruned)
170 pred_rt_test<-predict(rt_pruned, newdata = test_data )
    regr.eval(test_data$PRP, pred_rt,
172             stats = c("mae", "mse", "rmse"),
                train.y = NULL)
174 cor(pred_rt_test, test_data$PRP)**2

176 #SVR
    svr <- ksvm(PRP ~ MYCT + MMIN + MMAX + CACH + CHMIN + CHMAX , data=train_data2,
178              type="eps-svr", kernel="vanilladot", cross=5)
    pred_svr_test<-predict(svr, newdata = test_data )
180 regr.eval(test_data$PRP, pred_svr,
            stats = c("mae", "mse", "rmse"),
182             train.y = NULL)
    cor(pred_svr_test, test_data$PRP)**2

184 #part f
186 randf <- randomForest(PRP ~ MYCT + MMIN + MMAX + CACH + CHMIN + CHMAX, data=train_
    data, ntree=10 )
    pred_randf <- predict(randf, newdata = train_data)
188 cor(pred_randf, train_data$PRP)**2

190 regr.eval(train_data$PRP, pred_randf,
            stats = c("mae", "mse", "rmse"),
192             train.y = NULL)

```

Problem 6:

```

#install.packages("nlme")
2 rm(list = ls())
    setwd("~/Users/bansal/Google Drive/Spring2014/DataMining/HW3/")
4 library(MASS)
    library(hydroGOF)
6 library(ggplot2)
    library(nlme)
8 oxboys = read.csv("~/Google Drive/Spring2014/DataMining/HW3/HW3Data/oxboys.csv")

10 global <- lm(formula = height ~ year , data = oxboys)
    summary(global)
12 png("Q6-Part-a.png")
    layout(matrix(1:4, ncol = 2))
14 plot(global)
    layout(1)
16 dev.off()

18 ggplot = ggplot(oxboys, aes(x=year, y=height, group=id)) + stat_smooth(method="lm")

```

```

ggsave('Q6-Part-b.png')
20
oxboys.train <- subset(oxboys, year < 3)
22 oxboys.test <- subset(oxboys, year >= 3)

24 global <- lm(formula = height ~ year, data = oxboys.train)
testGolbal <- predict.lm(global, newdata = oxboys.test)
26 mse(testGolbal, oxboys.test$height)
summary(global)
28 png("Q6-Part-c-a.png")
layout(matrix(1:4, ncol = 2))
30 plot(global)
layout(1)
32 dev.off()
list.mse <- NULL
34
for( key in unique(oxboys.train$id)){
36   oxboys.train.ind <- subset(oxboys.train, id==key)
   oxboys.test.ind <- subset(oxboys.test, id==key)
38   # build a linear model with oxboys.train.ind
   mlr <- lm(height ~ year, data = oxboys.train.ind)
40   summary(mlr)
   #plot(mlr)
42
   # predict on oxboys.test.ind, and measure MSE
44   testModel<-predict.lm(mlr, newdata = oxboys.test.ind)
   mlrMSE<-mse(testModel, oxboys.test.ind$height)
46   list.mse <- c(list.mse, mlrMSE)
}
48 sum(list.mse)

50 #lme(height ~ year + id, data = oxboys, random = ~id/year, method = "ML", na.
   action = na.omit)
mlm.obj <- lme(height~year, data=oxboys.train, random=list(id=pdDiag(~year)))
52 testmlm <- predict(mlm.obj, newdata=oxboys.test, level=1)
# calculate MSE
54 mse(testmlm, oxboys.test$height)

56
#PART D
58 oxboys.train <- subset(oxboys, year < 7)
oxboys.test <- subset(oxboys, year >= 7)
60
global <- lm(formula = height ~ year, data = oxboys.train)
62 testGolbal <- predict.lm(global, newdata = oxboys.test)
mse(testGolbal, oxboys.test$height)
64 summary(global)
png("Q6-Part-d-a.png")
66 layout(matrix(1:4, ncol = 2))
plot(global)
68 layout(1)
dev.off()
70 list.mse <- NULL
for( key in unique(oxboys.train$id)){
72   oxboys.train.ind <- subset(oxboys.train, id==key)
   oxboys.test.ind <- subset(oxboys.test, id==key)

```

```

74 # build a linear model with oxboys.train.ind
    mlr <- lm(height ~ year , data = oxboys.train.ind)
76 #summary(mlr)
    # predict on oxboys.test.ind, and measure MSE
78 testModel<-predict.lm(mlr, newdata = oxboys.test.ind)
    mlrMSE<-mse(testModel,oxboys.test.ind$height)
80 list.mse <- c(list.mse,mlrMSE)
    }
82 sum(list.mse)

84 mlm2.obj <- lme(height~year , data=oxboys.train , random=list(id=pdDiag(~year)))
    test2mlm <- predict(mlm2.obj, newdata=oxboys.test , level=1)
86 # calculate MSE
    mse(test2mlm , oxboys.test$height)

```