

Introduction to Machine Learning

Lab 2 Block 2

Rasmus Holm

2016-11-30

Contents

Assignment 1a	2
Assignment 1b	3
Assignment 2a	4
1	4
2	4
3	4
Assignment 2b	5
Assignment 3a	6
1	6
2	6
Assignment 4a	7
Appendix	9
Code for Assignment 1a	9
Code for Assignment 1b	9
Code for Assignment 2a	9
Code for Assignment 2b	9
Code for Assignment 3a	9
Code for Assignment 4a	9

Assignment 1a

Assignment 1b

Assignment 2a

1

2

3

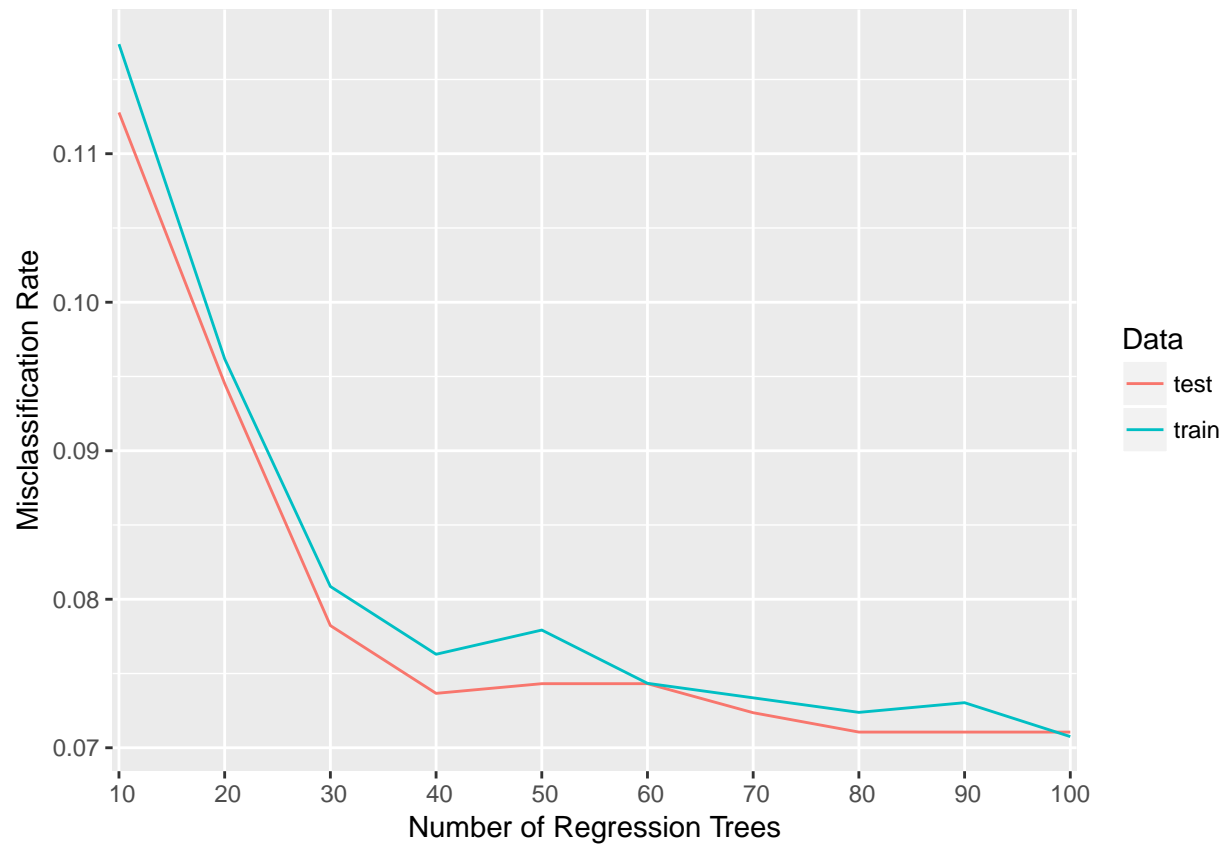
Assignment 2b

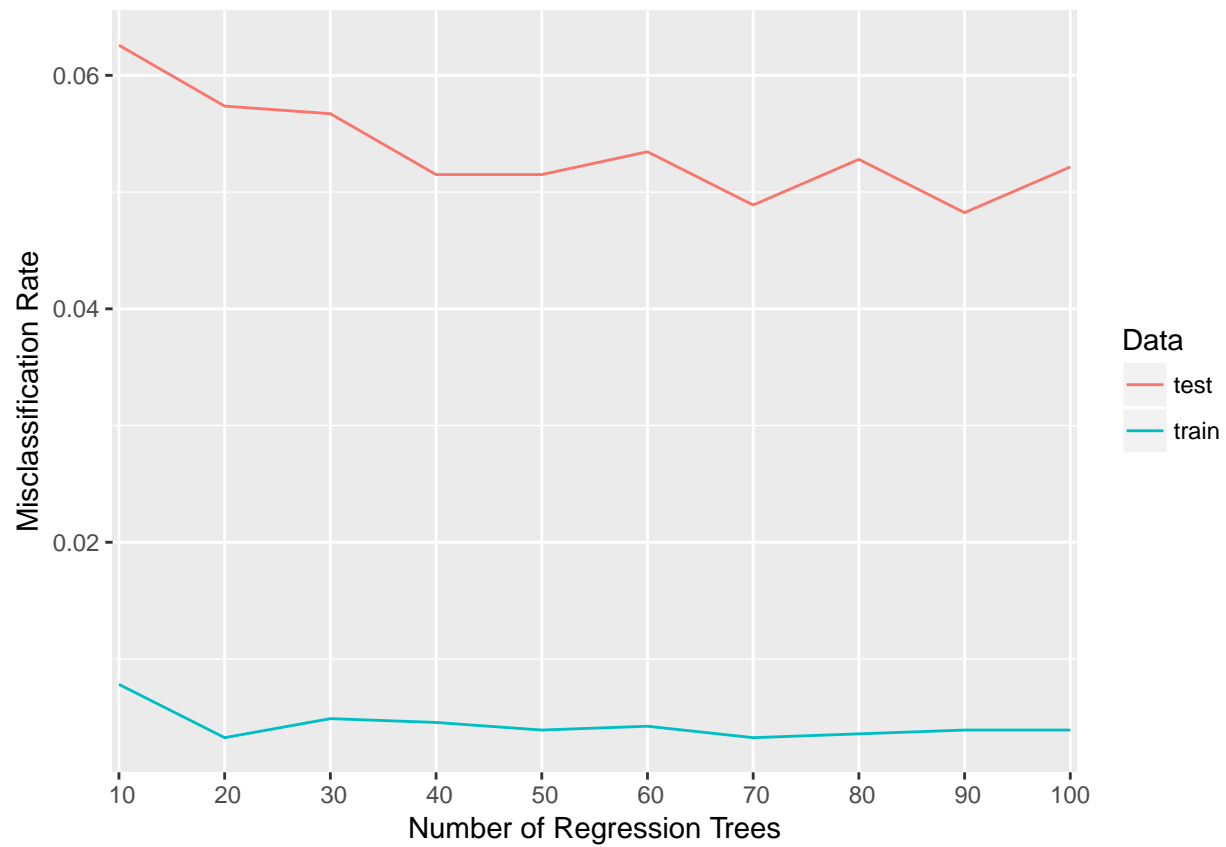
Assignment 3a

1

2

Assignment 4a





Appendix

Code for Assignment 1a

Code for Assignment 1b

Code for Assignment 2a

Code for Assignment 2b

Code for Assignment 3a

Code for Assignment 4a

```
library(mboost)
library(randomForest)
library(ggplot2)
library(reshape2)

data <- read.csv2("../data/spambase.csv")
data$Spam <- as.factor(data$Spam)

set.seed(1234567890)
train_idx <- sample(nrow(data), floor(nrow(data) * (2 / 3)))
train <- data[train_idx,]
test <- data[-train_idx,]
tree_counts <- seq(10, 100, by=10)
test_errors <- rep(0, length(tree_counts))
train_errors <- rep(0, length(tree_counts))

for (i in 1:length(tree_counts)) {
  fit <- blackboost(Spam ~ ., data=train, family=AdaExp(),
                    control=boost_control(mstop=tree_counts[i]))
  test_error <- 1 - (sum(predict(fit, test, type="class") == test$Spam) / nrow(test))
  train_error <- 1 - (sum(predict(fit, train, type="class") == train$Spam) / nrow(train))
  test_errors[i] <- test_error
  train_errors[i] <- train_error
}

plot_data <- data.frame(Trees=tree_counts, test=test_errors, train=train_errors)
plot_data <- melt(plot_data, id="Trees", value.name="Error", variable.name="Data")

ggplot(plot_data) +
  xlab("Number of Regression Trees") +
  ylab("Misclassification Rate") +
  geom_line(aes(x=Trees, y=Error, color=Data)) +
  scale_x_discrete(limits=tree_counts)
test_errors <- rep(0, length(tree_counts))
train_errors <- rep(0, length(tree_counts))

for (i in 1:length(tree_counts)) {
  fit <- randomForest(Spam ~ ., data=train, ntree=tree_counts[i])
  test_error <- 1 - (sum(predict(fit, test, type="class") == test$Spam) / nrow(test))
```

```

train_error <- 1 - (sum(predict(fit, train, type="class") == train$Spam) / nrow(train))
test_errors[i] <- test_error
train_errors[i] <- train_error
}
plot_data <- data.frame(Trees=tree_counts, test=test_errors, train=train_errors)
plot_data <- melt(plot_data, id="Trees", value.name="Error", variable.name="Data")

ggplot(plot_data) +
  xlab("Number of Regression Trees") +
  ylab("Misclassification Rate") +
  geom_line(aes(x=Trees, y=Error, color=Data)) +
  scale_x_discrete(limits=tree_counts)

```