

Practical Machine Learning Course Project

Daniel Blackburn

July 8, 2018

Introduction

This project is being carried out in completion of the “Practical Machine Learning” Coursera course.

A dataset of measurement data has been provided by the course. The dataset is comprised of measurements of acceleration made by individuals who are carrying out one of five classes of physical activity. According to this project’s instructions, the measurements are made using devices worn on the belt, forearm, arm, and a dumbbell.

Additional information the dataset is available here: <http://groupware.les.inf.puc-rio.br/har>

My task is to create a model that can predict the which class of activity is being done.

Download dataset

```
if(!file.exists("training.csv")) {  
  download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv",  
                "training.csv")  
}  
  
if(!file.exists("test.csv")) {  
  download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv",  
                "test.csv")  
}
```

Prepare data for model training, cross validation, and testing.

```
set.seed(2738255)  
  
df <- read.csv("training.csv", stringsAsFactors=TRUE)  
# Remove the column named "X", which represents the observation number and is not relevant to predicting outcome  
df <- df[,which(names(df) != "X")]  
# Let's ignore the data dimensions that have near zero variance. This eliminates rarely varying  
# parameters, which are presumably not very useful for prediction. In this were a more thorough  
# study, the effect of removing these parameters would be determined.  
df <- df[,-nearZeroVar(df)]  
  
# Let's partition the training data into training and cross-validation sets. The  
# cross-validation set is a hold-out set that allows us to measure the accuracy of the model.  
inTrain = createDataPartition(df$classe, p = 0.8)[[1]]  
training = df[ inTrain,]  
cvng = df[-inTrain,]  
  
testing <- read.csv("test.csv")  
testing <- testing[, which(names(df) != "X")]
```

Assess data

Let's inspect the data to see if there's any missing values.

```
print("Dimensions of training data:")  
  
## [1] "Dimensions of training data:"  
paste(dim(df))  
  
## [1] "19622" "99"  
print("Dimensions of training data, removing observations with NA values:")  
  
## [1] "Dimensions of training data, removing observations with NA values:"  
print(dim(df[rowSums(is.na(df)) == 0, ]))  
  
## [1] 406  99  
print("Dimensions of training data, removing data dimensions with NA values:")  
  
## [1] "Dimensions of training data, removing data dimensions with NA values:"  
print(dim(df[, colSums(is.na(df)) == 0]))  
  
## [1] 19622     58
```

Given that only 406 observations contain zero NAs, it seems that NAs were included by design.

Build model

Now that the data is partitioned into distinct sets for training, cross validating, and testing, we can build our model using the training set. The following three models were the ones I wished to try:

1. The package rpart allows a model to be built that contains NA values. Since the dataset has many NAs, this simplifies the problem of how to use the provided data.
2. Random forest with observations removed if NA is present.
3. Random forest with dimensions removed if NA is present. I am assuming that the testing set and training set both have the same NA values. If the testing set were to have NAs, we could try imputing the missing values. (In this assignment, this approach was not necessary.)

I ruled out 2. because so few training examples apply to it. Furthermore, it could not be applied to any test examples that contain NA values.

Build rpart model

```
# Method 1  
rpartModel <- rpart(classe ~ ., data=training, method="class")  
confusionMatrix(predict(rpartModel, cvng, type="class"),  
                 cvng$classe)  
  
## Confusion Matrix and Statistics  
##  
##          Reference  
## Prediction   A     B     C     D     E  
##           A 1076    27     3     1     0
```

```

##      B   28   616    32    27     0
##      C   12   110   633   117    24
##      D     0     6    10   431    43
##      E     0     0     6    67   654
##
## Overall Statistics
##
##          Accuracy : 0.8692
## 95% CI : (0.8583, 0.8796)
## No Information Rate : 0.2845
## P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.8346
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9642   0.8116   0.9254   0.6703   0.9071
## Specificity       0.9890   0.9725   0.9188   0.9820   0.9772
## Pos Pred Value    0.9720   0.8762   0.7065   0.8796   0.8996
## Neg Pred Value    0.9858   0.9556   0.9832   0.9382   0.9790
## Prevalence        0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate    0.2743   0.1570   0.1614   0.1099   0.1667
## Detection Prevalence 0.2822   0.1792   0.2284   0.1249   0.1853
## Balanced Accuracy  0.9766   0.8920   0.9221   0.8262   0.9421

```

At 95% confidence level, the accuracy is in the range of 85.8-88.0%. Let's compare the rpart model's performance with that of a random forest model.

Build random forest model

```

# Method 3: Remove columns with missing values and do train with random forest method

rfTraining <- training[, colSums(is.na(df)) == 0]
rfCving <- cving[, colSums(is.na(df)) == 0]

rfModel <- train(classe ~ ., data=rfTraining, method="rf")
confusionMatrix(predict(rfModel, rfCving), rfCving$classe)

## Confusion Matrix and Statistics
##
##          Reference
## Prediction   A     B     C     D     E
##           A 1116    0     0     0     0
##           B    0   759    1     0     0
##           C    0     0   683    0     0
##           D    0     0     0   643    1
##           E    0     0     0     0   720
##
## Overall Statistics
##
##          Accuracy : 0.9995

```

```

##                               95% CI : (0.9982, 0.9999)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##                           Kappa : 0.9994
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                                Class: A Class: B Class: C Class: D Class: E
## Sensitivity                 1.0000   1.0000   0.9985   1.0000   0.9986
## Specificity                  1.0000   0.9997   1.0000   0.9997   1.0000
## Pos Pred Value                1.0000   0.9987   1.0000   0.9984   1.0000
## Neg Pred Value                1.0000   1.0000   0.9997   1.0000   0.9997
## Prevalence                     0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate                 0.2845   0.1935   0.1741   0.1639   0.1835
## Detection Prevalence          0.2845   0.1937   0.1741   0.1642   0.1835
## Balanced Accuracy              1.0000   0.9998   0.9993   0.9998   0.9993

```

The proof is in the confusion matrix of the cross validation sample. The overall accuracy is determined to be greater than 99.8% at the 95% confidence level.

Because the random forest model has exceptional accuracy and seems equally suited to identifying all classes, it's the final model selected for this project.