

# Practical Machine Learning

Robert Deng

7/11/2017

```
library(caret)
library(randomForest)
library(rattle)
```

```
## Warning: Failed to load RGtk2 dynamic library, attempting to install it.
```

```
library(rpart)
library(rpart.plot)
```

## Background

Wearable technology (Fitbit, Nike FuelBand, and Jawbone Up) has made it possible to collect data about personal activity. People are measuring their own volume of activity, but rarely are they measuring quality. Training data was collected by participants to perform barbell lifts correctly and incorrectly in 5 different ways. Using data from accelerometers, the research question of interest is:

*Can you learn and properly classify correct and incorrect workout form?*

More information is available from the website here: <http://groupware.les.inf.puc-rio.br/har> (see the section on the Weight Lifting Exercise Dataset).

Data Agenda:

-Load Data -Clean Data -Learn -Predict

The training data for this project are available here: <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>

The test data are available here: <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>

Most of the NA / #DIV/0 cleaning can be done within the read.csv function

```
training <- read.csv('https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv',
                    na.strings=c('#DIV/0!', '', 'NA') , stringsAsFactors = T)

testing <- read.csv('https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv',
                   na.strings=c('#DIV/0!', '', 'NA') ,stringsAsFactors = T)
```

Cross validation used to create a standard 60% / 40% split on the training data set. The testing data doesn't have a classe column to compare results with and is used in the final model evaluation phase.

```
inTrain <- createDataPartition(training$classe, p=0.6, list=FALSE)
myTraining <- training[inTrain,]
myTesting <- training[-inTrain,]
dim(myTraining); dim(myTesting)
```

```
## [1] 11776 160
```

```
## [1] 7846 160
```

## Data cleaning

Near zero variance removes the columns with very little variance and columns with more than 70% NAs are removed. This is to keep the columns with more meaningful data. Also the X, id column 1 is removed.

```
#Near Zero Variance
nzv <- nearZeroVar(myTraining, saveMetrics=TRUE)
myTraining <- myTraining[nzv$nzv==FALSE]
myTesting <- myTesting[nzv$nzv==FALSE]
myTraining <- myTraining[-1]
myTesting <- myTesting[-1]
testing <- testing[-1]

#More than 70% NAs
myTraining <- myTraining[, -which(colMeans(is.na(myTraining)) > 0.7)]
myTesting <- myTesting[, -which(colMeans(is.na(myTesting)) > 0.7)]
```

The myTraining and final testing datasets need the same data classes and features coerced together for the RandomForest model. I have a hard time figuring out why, but the RF model works on an appended testing set with 1 row from myTraining. Something about the familiarity with the training data.

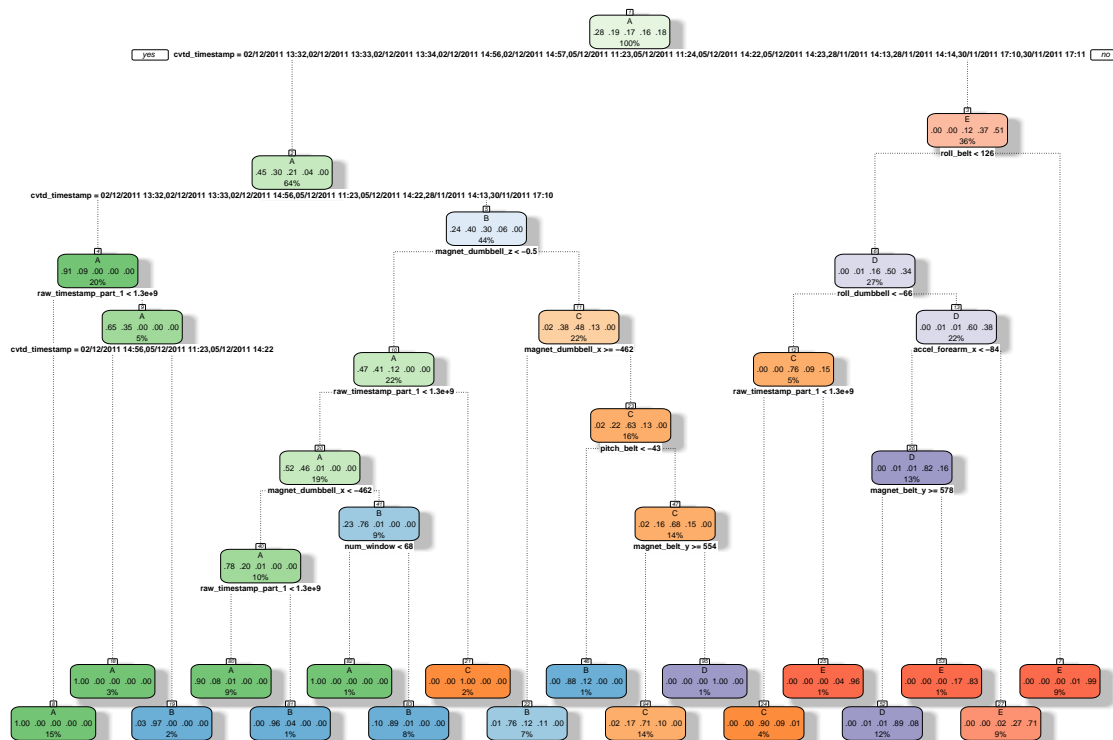
```
#Match columns between myTraining and testing
col_match <- colnames(myTraining)[-58]
testing2 <- testing[col_match]

#Coerce the same datatypes for both testing and myTraining datasets for randomForest
for (i in 1:length(testing2) ) {
  for(j in 1:length(myTraining)) {
    if( length( grep(names(myTraining[i]), names(testing2)[j]) ) ==1) {
      class(testing2[j]) <- class(myTraining[i])
    }
  }
}

#Rbind 1 from of myTraining to testing2
testing2 <- rbind(myTraining[2, -58] , testing2)
```

Now the fun part. First train for the rpart model using myTraining, then predict using myTesting.

```
set.seed(98765)
model1 <- rpart(classe ~ ., data=myTraining, method="class")
fancyRpartPlot(model1)
```



Rattle 2017-Jul-11 16:36:52 robo

Out of sample error is ~12% with most of the misclassifications on D.

```
predResults1 <- predict(model1, myTesting, type = "class")
confusionMatrix(predResults1, myTesting$classe)
```

## Confusion Matrix and Statistics

```
##
##           Reference
## Prediction   A    B    C    D    E
##           A 2150   58    7    1    0
##           B   64 1284   86   69    0
##           C   18  168 1250  141    4
##           D    0    8   16  868   83
##           E    0    0    9  207 1355
```

## Overall Statistics

```
##
##           Accuracy : 0.8803
##           95% CI : (0.8729, 0.8874)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
```

```
##
##           Kappa : 0.8486
```

```
## McNemar's Test P-Value : NA
```

## Statistics by Class:

```
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity   0.9633   0.8458   0.9137   0.6750   0.9397
```

```
## Specificity          0.9882  0.9654  0.9489  0.9837  0.9663
## Pos Pred Value      0.9702  0.8543  0.7906  0.8903  0.8625
## Neg Pred Value      0.9854  0.9631  0.9812  0.9392  0.9861
## Prevalence          0.2845  0.1935  0.1744  0.1639  0.1838
## Detection Rate      0.2740  0.1637  0.1593  0.1106  0.1727
## Detection Prevalence 0.2824  0.1916  0.2015  0.1243  0.2002
## Balanced Accuracy    0.9758  0.9056  0.9313  0.8293  0.9530
```

To validate, RandomForest yielded a 0.2% error rate. This outperformed the rpart model but could be slightly overfitting. RandomForest does a better job selecting the most important features via gini index and is more resilient to outliers.

```
model2 <- randomForest(classe ~. , data=myTraining, trControl = trainControl(method = "cv", 5))
predResults2 <- predict(model2, myTesting, type = "class")
confusionMatrix(predResults2, myTesting$classe)
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 2232     1     0     0     0
##      B     0 1517     0     0     0
##      C     0     0 1365     4     0
##      D     0     0     3 1282     2
##      E     0     0     0     0 1440
##
## Overall Statistics
##
##              Accuracy : 0.9987
##              95% CI : (0.9977, 0.9994)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9984
##      McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          1.0000  0.9993  0.9978  0.9969  0.9986
## Specificity          0.9998  1.0000  0.9994  0.9992  1.0000
## Pos Pred Value       0.9996  1.0000  0.9971  0.9961  1.0000
## Neg Pred Value       1.0000  0.9998  0.9995  0.9994  0.9997
## Prevalence           0.2845  0.1935  0.1744  0.1639  0.1838
## Detection Rate       0.2845  0.1933  0.1740  0.1634  0.1835
## Detection Prevalence 0.2846  0.1933  0.1745  0.1640  0.1835
## Balanced Accuracy    0.9999  0.9997  0.9986  0.9981  0.9993
```

Now for the final testing rounds:

```
predFinal1 <- predict(model1, testing2, type = "class")
predFinal1
```

```
##  2  1 21  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
##  A  B  A  C  A  A  E  D  C  A  A  B  C  B  A  E  E  A  B  B  B
## Levels: A B C D E
```

```
#Ignore the 31 (stitched from myTraining) in the model
predFinal2 <- predict(model2, testing2, type = "class")
predFinal2 <- predFinal2[-4]
predFinal2
```

```
##  2  1 21  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
##  A  B  A  A  A  E  D  B  A  A  B  C  B  A  E  E  A  B  B  B
## Levels: A B C D E
```

Write out results:

```
pml_write_files = function(x) {
  n = length(x)
  for (i in 1:n) {
    filename = paste0("problem_id_", i, ".txt")
    write.table(x[i], file = filename, quote = FALSE, row.names = FALSE,
               col.names = FALSE)
  }
}
pml_write_files(predFinal2)
```