

Prediction Assignment Writeup

Andre Jordaan

4/6/2017

Background

Using devices such as Jawbone Up, Nike FuelBand, and Fitbit it is now possible to collect a large amount of data about personal activity relatively inexpensively. These type of devices are part of the quantified self movement - a group of enthusiasts who take measurements about themselves regularly to improve their health, to find patterns in their behavior, or because they are tech geeks. One thing that people regularly do is quantify how much of a particular activity they do, but they rarely quantify how well they do it. In this project, your goal will be to use data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants. They were asked to perform barbell lifts correctly and incorrectly in 5 different ways. 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

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>

The data for this project come from this source: <http://groupware.les.inf.puc-rio.br/har>. If you use the document you create for this class for any purpose please cite them as they have been very generous in allowing their data to be used for this kind of assignment.

What you should submit

The goal of your project is to predict the manner in which they did the exercise. This is the “classe” variable in the training set. You may use any of the other variables to predict with. You should create a report describing how you built your model, how you used cross validation, what you think the expected out of sample error is, and why you made the choices you did. You will also use your prediction model to predict 20 different test cases.

Your submission should consist of a link to a Github repo with your R markdown and compiled HTML file describing your analysis. Please constrain the text of the writeup to < 2000 words and the number of figures to be less than 5. It will make it easier for the graders if you submit a repo with a gh-pages branch so the HTML page can be viewed online (and you always want to make it easy on graders :-). You should also apply your machine learning algorithm to the 20 test cases available in the test data above. Please submit your predictions in appropriate format to the programming assignment for automated grading. See the programming assignment for additional details.

Reproducibility

Due to security concerns with the exchange of R code, your code will not be run during the evaluation by your classmates. Please be sure that if they download the repo, they will be able to view the compiled HTML version of your analysis.

Prepare the data

Load required libraries and global settings

```
# Load the required libraries
library(caret)           # Classification and Regression Training
library(rpart)           # Recursive Partitioning and Regression Trees
library(rpart.plot)      # Plot 'rpart' Models: An Enhanced Version of 'plot.rpart'
library(randomForest)    # Breiman and Cutler's Random Forests for Classification and Regression
library(rattle)          # Graphical User Interface for Data Mining in R
library(doMC);           # Provides a parallel backend for the %dopar% function using
```

Download the data

```
# Create data directory if not exist
if(!file.exists("./data")){dir.create("./data")}

# Download the training data file if it does not exist
if(!file.exists("./data/pml-training.csv")) {
  trainUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
  download.file(trainUrl, destfile="./data/pml-training.csv", method="curl", quiet=FALSE)
}

# Download the test data file if it does not exist
if(!file.exists("./data/pml-testing.csv")) {
  testUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
  download.file(testUrl, destfile="./data/pml-testing.csv", method="curl", quiet=FALSE)
}
```

Once the download of the required files are complete, the data will be read into data frames using the read.csv function. Display the size and information of the two data frames.

```
#Read the .csv files into corresponding data frames.
trainDF <- read.csv("./data/pml-training.csv", na.strings=c("NA","NaN","#DIV/0!", ""))
testDF <- read.csv("./data/pml-testing.csv", na.strings=c("NA","NaN","#DIV/0!", ""))

# Display dimensions of the data frames
dim(trainDF)
```

```
## [1] 19622 160
```

```
dim(testDF)
```

```
## [1] 20 160
```

```
str(trainDF)
```

```
## 'data.frame': 19622 obs. of 160 variables:
## $ X : int 1 2 3 4 5 6 7 8 9 10 ...
## $ user_name : Factor w/ 6 levels "adelmo","carlitos",...: 2 2 2 2 2 2 2 2 2 2 ...
## $ raw_timestamp_part_1 : int 1323084231 1323084231 1323084231 1323084232 1323084232 1323084232 ...
## $ raw_timestamp_part_2 : int 788290 808298 820366 120339 196328 304277 368296 440390 484323 484...
## $ cvtd_timestamp : Factor w/ 20 levels "02/12/2011 13:32",...: 9 9 9 9 9 9 9 9 9 9 ...
## $ new_window : Factor w/ 2 levels "no","yes": 1 1 1 1 1 1 1 1 1 1 ...
## $ num_window : int 11 11 11 12 12 12 12 12 12 12 ...
```

```

## $ roll_belt      : num  1.41 1.41 1.42 1.48 1.48 1.45 1.42 1.42 1.43 1.45 ...
## $ pitch_belt     : num  8.07 8.07 8.07 8.05 8.07 8.06 8.09 8.13 8.16 8.17 ...
## $ yaw_belt       : num -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 ...
## $ total_accel_belt : int   3 3 3 3 3 3 3 3 3 3 ...
## $ kurtosis_roll_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_belt  : logi  NA NA NA NA NA NA ...
## $ skewness_roll_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_belt.1 : num  NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_belt  : logi  NA NA NA NA NA NA ...
## $ max_roll_belt     : num  NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_belt    : int   NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_belt      : num  NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_belt     : num  NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_belt    : int   NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_belt      : num  NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_belt : int   NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_belt  : num  NA NA NA NA NA NA NA NA NA NA ...
## $ var_total_accel_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ avg_roll_belt     : num  NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_belt  : num  NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_belt     : num  NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_belt    : num  NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_belt  : num  NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_belt    : num  NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_belt      : num  NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_belt   : num  NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_belt      : num  NA NA NA NA NA NA NA NA NA NA ...
## $ gyros_belt_x      : num  0 0.02 0 0.02 0.02 0.02 0.02 0.02 0.02 0.03 ...
## $ gyros_belt_y      : num  0 0 0 0 0.02 0 0 0 0 0 ...
## $ gyros_belt_z      : num -0.02 -0.02 -0.02 -0.03 -0.02 -0.02 -0.02 -0.02 -0.02 0 ...
## $ accel_belt_x      : int   -21 -22 -20 -22 -21 -21 -22 -22 -20 -21 ...
## $ accel_belt_y      : int    4 4 5 3 2 4 3 4 2 4 ...
## $ accel_belt_z      : int   22 22 23 21 24 21 21 21 24 22 ...
## $ magnet_belt_x     : int    -3 -7 -2 -6 -6 0 -4 -2 1 -3 ...
## $ magnet_belt_y     : int  599 608 600 604 600 603 599 603 602 609 ...
## $ magnet_belt_z     : int -313 -311 -305 -310 -302 -312 -311 -313 -312 -308 ...
## $ roll_arm          : num -128 -128 -128 -128 -128 -128 -128 -128 -128 -128 ...
## $ pitch_arm         : num  22.5 22.5 22.5 22.1 22.1 22 21.9 21.8 21.7 21.6 ...
## $ yaw_arm           : num -161 -161 -161 -161 -161 -161 -161 -161 -161 -161 ...
## $ total_accel_arm   : int   34 34 34 34 34 34 34 34 34 34 ...
## $ var_accel_arm     : num  NA NA NA NA NA NA NA NA NA NA ...
## $ avg_roll_arm      : num  NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_arm   : num  NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_arm      : num  NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_arm     : num  NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_arm  : num  NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_arm     : num  NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_arm       : num  NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_arm    : num  NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_arm       : num  NA NA NA NA NA NA NA NA NA NA ...
## $ gyros_arm_x       : num  0 0.02 0.02 0.02 0 0.02 0 0.02 0.02 0.02 ...
## $ gyros_arm_y       : num  0 -0.02 -0.02 -0.03 -0.03 -0.03 -0.03 -0.02 -0.03 -0.03 ...

```

```
## $ gyros_arm_z : num -0.02 -0.02 -0.02 0.02 0 0 0 0 -0.02 -0.02 ...
## $ accel_arm_x : int -288 -290 -289 -289 -289 -289 -289 -289 -288 -288 ...
## $ accel_arm_y : int 109 110 110 111 111 111 111 111 109 110 ...
## $ accel_arm_z : int -123 -125 -126 -123 -123 -122 -125 -124 -122 -124 ...
## $ magnet_arm_x : int -368 -369 -368 -372 -374 -369 -373 -372 -369 -376 ...
## $ magnet_arm_y : int 337 337 344 344 337 342 336 338 341 334 ...
## $ magnet_arm_z : int 516 513 513 512 506 513 509 510 518 516 ...
## $ kurtosis_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_arm : int NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_arm : int NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm : int NA NA NA NA NA NA NA NA NA NA ...
## $ roll_dumbbell : num 13.1 13.1 12.9 13.4 13.4 ...
## $ pitch_dumbbell : num -70.5 -70.6 -70.3 -70.4 -70.4 ...
## $ yaw_dumbbell : num -84.9 -84.7 -85.1 -84.9 -84.9 ...
## $ kurtosis_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_dumbbell : logi NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_dumbbell : logi NA NA NA NA NA NA ...
## $ max_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]
```

```
str(testDF)
```

```
## 'data.frame': 20 obs. of 160 variables:
## $ X : int 1 2 3 4 5 6 7 8 9 10 ...
## $ user_name : Factor w/ 6 levels "adelmo","carlitos",...: 6 5 5 1 4 5 5 2 3 ...
## $ raw_timestamp_part_1 : int 1323095002 1322673067 1322673075 1322832789 1322489635 1322673149 ...
## $ raw_timestamp_part_2 : int 868349 778725 342967 560311 814776 510661 766645 54671 916313 3842 ...
## $ cvtd_timestamp : Factor w/ 11 levels "02/12/2011 13:33",...: 5 10 10 1 6 11 11 10 3 2 ...
## $ new_window : Factor w/ 1 level "no": 1 1 1 1 1 1 1 1 1 1 ...
## $ num_window : int 74 431 439 194 235 504 485 440 323 664 ...
## $ roll_belt : num 123 1.02 0.87 125 1.35 -5.92 1.2 0.43 0.93 114 ...
## $ pitch_belt : num 27 4.87 1.82 -41.6 3.33 1.59 4.44 4.15 6.72 22.4 ...
## $ yaw_belt : num -4.75 -88.9 -88.5 162 -88.6 -87.7 -87.3 -88.5 -93.7 -13.1 ...
## $ total_accel_belt : int 20 4 5 17 3 4 4 4 4 18 ...
## $ kurtosis_roll_belt : logi NA NA NA NA NA NA ...
```

```

## $ kurtosis_picth_belt      : logi  NA NA NA NA NA NA ...
## $ kurtosis_yaw_belt       : logi  NA NA NA NA NA NA ...
## $ skewness_roll_belt      : logi  NA NA NA NA NA NA ...
## $ skewness_roll_belt.1    : logi  NA NA NA NA NA NA ...
## $ skewness_yaw_belt       : logi  NA NA NA NA NA NA ...
## $ max_roll_belt           : logi  NA NA NA NA NA NA ...
## $ max_picth_belt          : logi  NA NA NA NA NA NA ...
## $ max_yaw_belt            : logi  NA NA NA NA NA NA ...
## $ min_roll_belt           : logi  NA NA NA NA NA NA ...
## $ min_pitch_belt          : logi  NA NA NA NA NA NA ...
## $ min_yaw_belt            : logi  NA NA NA NA NA NA ...
## $ amplitude_roll_belt     : logi  NA NA NA NA NA NA ...
## $ amplitude_pitch_belt    : logi  NA NA NA NA NA NA ...
## $ amplitude_yaw_belt      : logi  NA NA NA NA NA NA ...
## $ var_total_accel_belt    : logi  NA NA NA NA NA NA ...
## $ avg_roll_belt           : logi  NA NA NA NA NA NA ...
## $ stddev_roll_belt        : logi  NA NA NA NA NA NA ...
## $ var_roll_belt           : logi  NA NA NA NA NA NA ...
## $ avg_pitch_belt          : logi  NA NA NA NA NA NA ...
## $ stddev_pitch_belt       : logi  NA NA NA NA NA NA ...
## $ var_pitch_belt          : logi  NA NA NA NA NA NA ...
## $ avg_yaw_belt            : logi  NA NA NA NA NA NA ...
## $ stddev_yaw_belt         : logi  NA NA NA NA NA NA ...
## $ var_yaw_belt            : logi  NA NA NA NA NA NA ...
## $ gyros_belt_x            : num  -0.5 -0.06 0.05 0.11 0.03 0.1 -0.06 -0.18 0.1 0.14 ...
## $ gyros_belt_y            : num  -0.02 -0.02 0.02 0.11 0.02 0.05 0 -0.02 0 0.11 ...
## $ gyros_belt_z            : num  -0.46 -0.07 0.03 -0.16 0 -0.13 0 -0.03 -0.02 -0.16 ...
## $ accel_belt_x            : int   -38 -13 1 46 -8 -11 -14 -10 -15 -25 ...
## $ accel_belt_y            : int    69 11 -1 45 4 -16 2 -2 1 63 ...
## $ accel_belt_z            : int  -179 39 49 -156 27 38 35 42 32 -158 ...
## $ magnet_belt_x           : int   -13 43 29 169 33 31 50 39 -6 10 ...
## $ magnet_belt_y           : int   581 636 631 608 566 638 622 635 600 601 ...
## $ magnet_belt_z           : int  -382 -309 -312 -304 -418 -291 -315 -305 -302 -330 ...
## $ roll_arm                 : num   40.7 0 0 -109 76.1 0 0 0 -137 -82.4 ...
## $ pitch_arm                : num  -27.8 0 0 55 2.76 0 0 0 11.2 -63.8 ...
## $ yaw_arm                  : num   178 0 0 -142 102 0 0 0 -167 -75.3 ...
## $ total_accel_arm          : int    10 38 44 25 29 14 15 22 34 32 ...
## $ var_accel_arm            : logi  NA NA NA NA NA NA ...
## $ avg_roll_arm             : logi  NA NA NA NA NA NA ...
## $ stddev_roll_arm          : logi  NA NA NA NA NA NA ...
## $ var_roll_arm             : logi  NA NA NA NA NA NA ...
## $ avg_pitch_arm            : logi  NA NA NA NA NA NA ...
## $ stddev_pitch_arm         : logi  NA NA NA NA NA NA ...
## $ var_pitch_arm            : logi  NA NA NA NA NA NA ...
## $ avg_yaw_arm              : logi  NA NA NA NA NA NA ...
## $ stddev_yaw_arm           : logi  NA NA NA NA NA NA ...
## $ var_yaw_arm              : logi  NA NA NA NA NA NA ...
## $ gyros_arm_x              : num  -1.65 -1.17 2.1 0.22 -1.96 0.02 2.36 -3.71 0.03 0.26 ...
## $ gyros_arm_y              : num   0.48 0.85 -1.36 -0.51 0.79 0.05 -1.01 1.85 -0.02 -0.5 ...
## $ gyros_arm_z              : num  -0.18 -0.43 1.13 0.92 -0.54 -0.07 0.89 -0.69 -0.02 0.79 ...
## $ accel_arm_x              : int    16 -290 -341 -238 -197 -26 99 -98 -287 -301 ...
## $ accel_arm_y              : int    38 215 245 -57 200 130 79 175 111 -42 ...
## $ accel_arm_z              : int    93 -90 -87 6 -30 -19 -67 -78 -122 -80 ...
## $ magnet_arm_x             : int  -326 -325 -264 -173 -170 396 702 535 -367 -420 ...

```

```
## $ magnet_arm_y      : int  385 447 474 257 275 176 15 215 335 294 ...
## $ magnet_arm_z      : int  481 434 413 633 617 516 217 385 520 493 ...
## $ kurtosis_roll_arm : logi  NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_arm : logi  NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_arm   : logi  NA NA NA NA NA NA NA ...
## $ skewness_roll_arm  : logi  NA NA NA NA NA NA NA ...
## $ skewness_pitch_arm : logi  NA NA NA NA NA NA NA ...
## $ skewness_yaw_arm   : logi  NA NA NA NA NA NA NA ...
## $ max_roll_arm       : logi  NA NA NA NA NA NA NA ...
## $ max_pitch_arm      : logi  NA NA NA NA NA NA NA ...
## $ max_yaw_arm        : logi  NA NA NA NA NA NA NA ...
## $ min_roll_arm       : logi  NA NA NA NA NA NA NA ...
## $ min_pitch_arm      : logi  NA NA NA NA NA NA NA ...
## $ min_yaw_arm        : logi  NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : logi  NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : logi  NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm  : logi  NA NA NA NA NA NA NA ...
## $ roll_dumbbell      : num  -17.7 54.5 57.1 43.1 -101.4 ...
## $ pitch_dumbbell     : num   25 -53.7 -51.4 -30 -53.4 ...
## $ yaw_dumbbell       : num  126.2 -75.5 -75.2 -103.3 -14.2 ...
## $ kurtosis_roll_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ skewness_pitch_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ skewness_yaw_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ max_roll_dumbbell   : logi  NA NA NA NA NA NA NA ...
## $ max_pitch_dumbbell  : logi  NA NA NA NA NA NA NA ...
## $ max_yaw_dumbbell    : logi  NA NA NA NA NA NA NA ...
## $ min_roll_dumbbell   : logi  NA NA NA NA NA NA NA ...
## $ min_pitch_dumbbell  : logi  NA NA NA NA NA NA NA ...
## $ min_yaw_dumbbell    : logi  NA NA NA NA NA NA NA ...
## $ amplitude_roll_dumbbell : logi  NA NA NA NA NA NA NA ...
## [list output truncated]
```

Summary of the data

- Training data frame has 19622 rows (observations)
- Training data frame has 160 columns (variables)
- 152 out of 160 are sensor readings for 4 sensors
- Sensor identified by: 'belt', 'arm', 'dumbbell' or 'forearm'
- "classe" variable in the training set is the outcome to predict
- Columns 1 - 7 will be removed as the data are not sensor readings

Clean the data

- Remove columns that contain NA missing values.
- Remove columns that do not contribute much to the accelerometer measurements

```
# Remove columns not containing the accelerometer measurements. Keep prediction variable classe.
# Remove columns 1 - 7
reqcol = grep(pattern = "_belt|_arm|_dumbbell|_forearm", names(trainDF))
```

```
trainDF <- trainDF[, c(reqcol,160)]
trainDF <- trainDF[, -c(1:7)]
notreqcol = which(colSums(is.na(trainDF)) > 19000)
trainDF = trainDF[, -notreqcol]
```

```
# Show table dimensions and show table classes
dim(trainDF)
```

```
## [1] 19622    49
```

```
table(sapply(trainDF[1,], class))
```

```
##
##  factor integer numeric
##      1      24      24
```

Training now consists of 19622 observations and 49 variables.

Testing data now consists of 20 observations and 49 variables.

Slice the data

We are now able to split the training set into clean 70% training and 30 validation sets.

```
# Set the seed for reproducibility
set.seed(760924)
```

```
inTrain <- createDataPartition(trainDF$classe, p=0.70, list=F)
training <- trainDF[inTrain, ]
testing <- trainDF[-inTrain, ]
```

```
# Check the dimensions for each set
dim(training)
```

```
## [1] 13737    49
```

```
dim(testing)
```

```
## [1] 5885    49
```

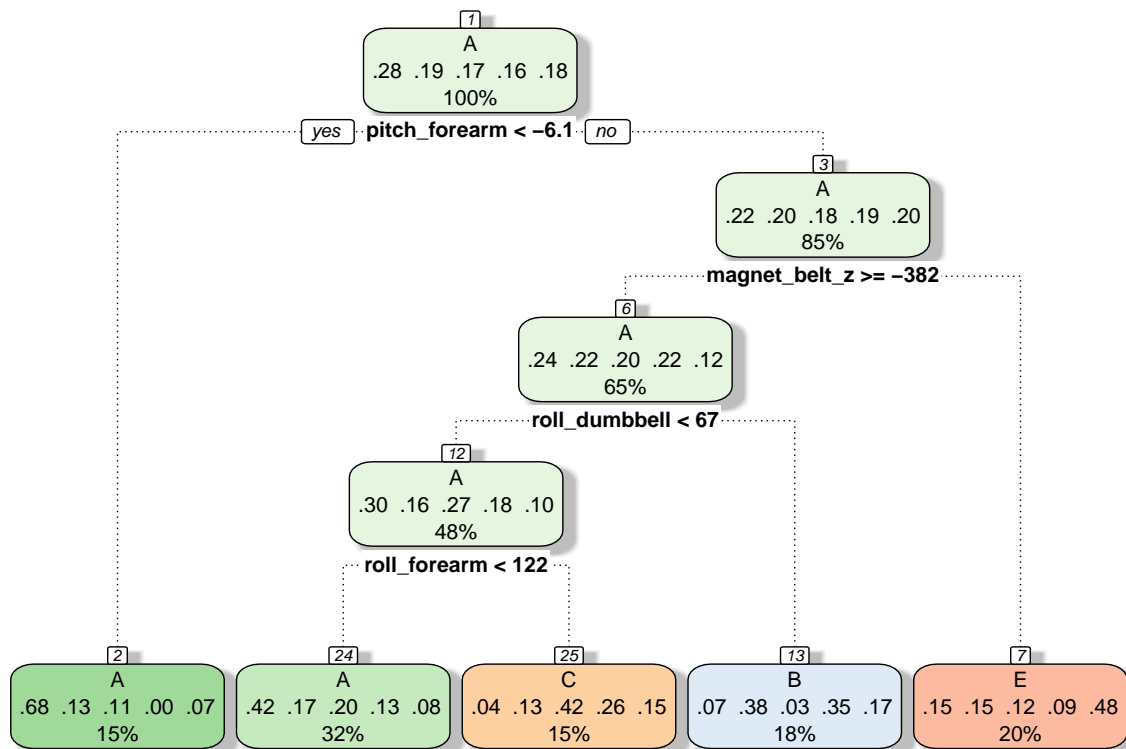
Data Modeling Selection

Model: Decision Tree

We fit a predictive model for activity recognition using Random Forest algorithm because it automatically selects important variables and is robust to correlated covariates & outliers in general.

```
# Fit and train the model
modFitregTree <- rpart(classe ~ ., data=training, minbucket = 2000)
```

```
# Plot the tree
fancyRpartPlot(modFitregTree)
```



Rattle 2017-Apr-09 19:38:17 andre

Show predictors

```

predmodFitregTree <- predict(modFitregTree, testing, type = "class")
confusionMatrix(predmodFitregTree, testing$classe)

```

Confusion Matrix and Statistics

##

Reference

Prediction	A	B	C	D	E
A	1403	469	461	258	206
B	64	395	33	369	180
C	23	89	381	227	117
D	0	0	0	0	0
E	184	186	151	110	579

##

Overall Statistics

##

Accuracy : 0.4686

95% CI : (0.4558, 0.4815)

No Information Rate : 0.2845

P-Value [Acc > NIR] : < 2.2e-16

##

Kappa : 0.3081

McNemar's Test P-Value : < 2.2e-16

##

Statistics by Class:

##

Class: A Class: B Class: C Class: D Class: E

Sensitivity 0.8381 0.34680 0.37135 0.0000 0.53512


```
## Specificity          0.6690  0.86389  0.90615  1.0000  0.86862
## Pos Pred Value      0.5016  0.37944  0.45520    NaN  0.47851
## Neg Pred Value      0.9122  0.84641  0.87223  0.8362  0.89241
## Prevalence          0.2845  0.19354  0.17434  0.1638  0.18386
## Detection Rate       0.2384  0.06712  0.06474  0.0000  0.09839
## Detection Prevalence 0.4753  0.17689  0.14223  0.0000  0.20561
## Balanced Accuracy    0.7535  0.60534  0.63875  0.5000  0.70187
```

Accuracy : 0.4686 making this not the ideal model to use.

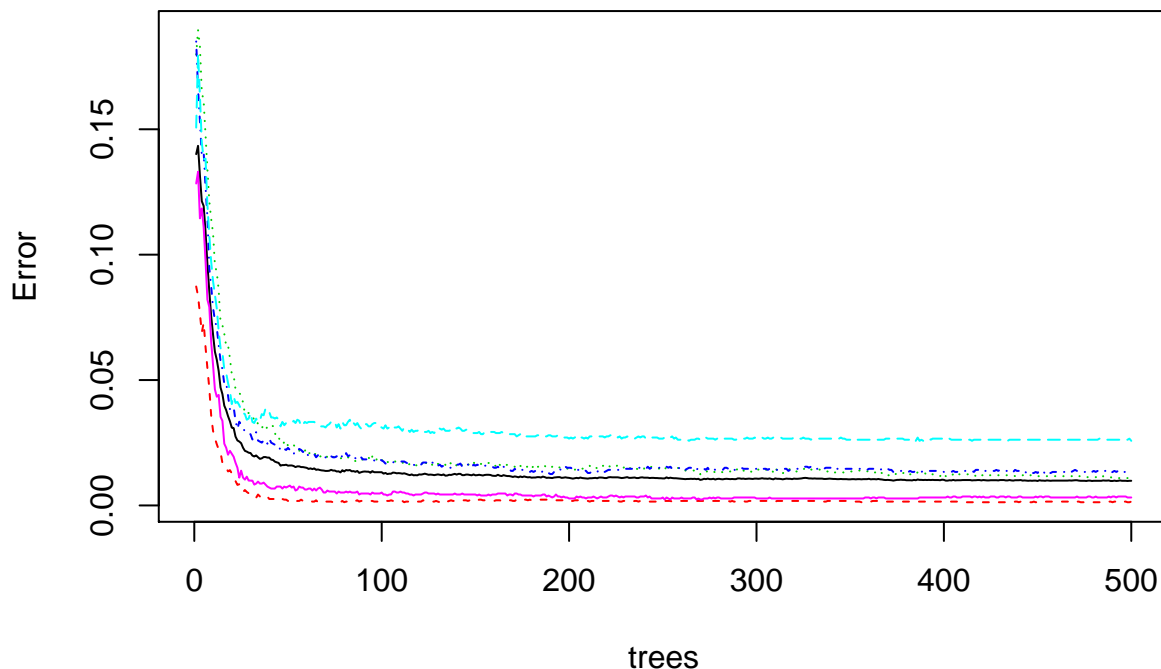
Model: Random Forest

```
#registerDoMC(cores = 4) #Leveraging Multi-Core for Parallelization

# Fit and train the model
modFitrandomForest <- randomForest(classe ~. , data=training, ntree = 500)

# Plot the tree
plot(modFitrandomForest)
```

modFitrandomForest



```
# Show predictors
predmodFitrandomForest = predict(modFitrandomForest, newdata = testing)
confusionMatrix(predmodFitrandomForest, testing$classe)
```

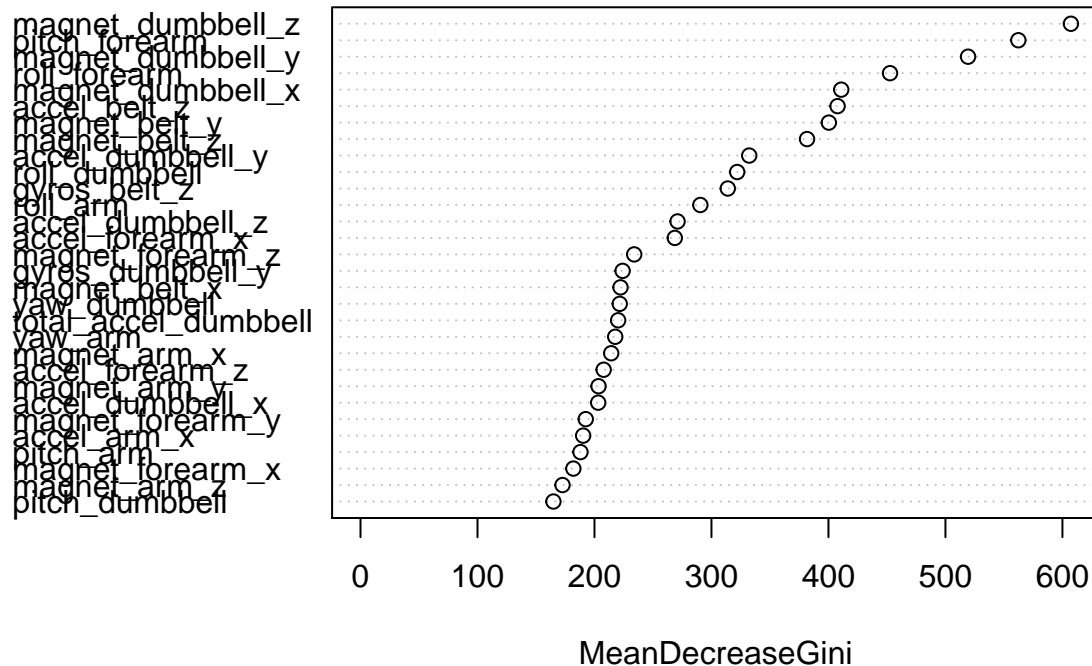
Confusion Matrix and Statistics

```
##
##          Reference
## Prediction    A    B    C    D    E
##          A 1672    7    0    0    0
##          B    1 1131    9    0    0
##          C    0    1 1017   25    0
```

```
##           D      0      0      0 935      3
##           E      1      0      0   4 1079
##
## Overall Statistics
##
##           Accuracy : 0.9913
##           95% CI : (0.9886, 0.9935)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.989
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9988  0.9930  0.9912  0.9699  0.9972
## Specificity      0.9983  0.9979  0.9946  0.9994  0.9990
## Pos Pred Value   0.9958  0.9912  0.9751  0.9968  0.9954
## Neg Pred Value   0.9995  0.9983  0.9981  0.9941  0.9994
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2841  0.1922  0.1728  0.1589  0.1833
## Detection Prevalence 0.2853  0.1939  0.1772  0.1594  0.1842
## Balanced Accuracy 0.9986  0.9954  0.9929  0.9847  0.9981
```

```
#Dotchart of variable importance as measured by a Random Forest
varImpPlot(modFitrandomForest)
```

modFitrandomForest



Random Forests yielded better results, as expected and Cohen's kappa indicator has low out of sample errors!
Accuracy : 0.99 making this model ideal.

Predicting for Test Data Set

```
eval <-testDF[,intersect(names(trainDF),names(testDF))]  
  
predictions<-predict(modFitrandomForest, eval)  
  
pml_write_files = function(x){  
  n = length(x)  
  for(i in 1:n){  
    filename = paste0("answers/problem_id_",i,".txt")  
    write.table(x[i],file=filename,quote=FALSE,row.names=FALSE,col.names=FALSE)  
  }  
}  
  
pml_write_files(predictions)
```