# PML Course Project

*DePoint*

*Sunday, June 14, 2015*

## Intro

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 dumbell of 6 participants. They were asked to perform barbell lifts correctly and incorrectly in 5 different ways. These are captured in the (outcome) variable called "classe"

## Setup the Environment

```
library(caret);library(rattle);library(dplyr)
```

```
## Loading required package: lattice
## Loading required package: ggplot2
## Rattle: A free graphical interface for data mining with R.
## Version 3.4.1 Copyright (c) 2006-2014 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.
##
## Attaching package: 'dplyr'
##
## The following object is masked from 'package:stats':
##
##     filter
##
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

## get the data and bring it into memory only

```
trainUrl <-"http://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
testUrl <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
training <- read.csv(url(trainUrl))
testing <- read.csv(url(testUrl))
```

# clean the data

lets examine missing values, near zero variables, redundant variables, etc

```
nsv <- nearZeroVar(training,saveMetrics=TRUE)
nsvDropList <- names(training) %in% as.list(row.names(subset(nsv, nzv=='TRUE',select=nz
v)))
slimTraining <- training[!nsvDropList]
#crunch through again to remove all missing
slimTraining <- slimTraining[, colSums(is.na(slimTraining)) == 0]
```

After cycling through some modeling, these variables cause some issues. lets remove

```
slimTraining <- slimTraining[, !grepl("^X", names(slimTraining))]
#slimTraining <- select(slimTraining,
#                   -c(raw_timestamp_part_1,raw_timestamp_part_2,cvtd_timestamp))
dim(slimTraining);str(slimTraining)
```

```
## [1] 19622     58
```

```
## 'data.frame':     19622 obs. of  58 variables:
##  $ 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 1
323084232 1323084232 1323084232 1323084232 1323084232 ...
##  $ raw_timestamp_part_2: int  788290 808298 820366 120339 196328 304277 368296 440390
484323 484434 ...
##  $ cvtd_timestamp     : Factor w/ 20 levels "02/12/2011 13:32",..: 9 9 9 9 9 9 9 9 9
9 ...
##  $ 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 -9
4.4 ...
##  $ total_accel_belt   : int  3 3 3 3 3 3 3 3 3 3 ...
##  $ 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 ...
##  $ 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 ...
##  $ 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 ...
##  $ total_accel_dumbbell: int  37 37 37 37 37 37 37 37 37 37 ...
##  $ gyros_dumbbell_x    : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ gyros_dumbbell_y    : num  -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02
## -0.02 ...
##  $ gyros_dumbbell_z    : num  0 0 0 -0.02 0 0 0 0 0 0 ...
##  $ accel_dumbbell_x    : int  -234 -233 -232 -232 -233 -234 -232 -234 -232 -235 ...
##  $ accel_dumbbell_y    : int  47 47 46 48 48 48 47 46 47 48 ...
##  $ accel_dumbbell_z    : int  -271 -269 -270 -269 -270 -269 -270 -272 -269 -270 ...
##  $ magnet_dumbbell_x   : int  -559 -555 -561 -552 -554 -558 -551 -555 -549 -558 ...
##  $ magnet_dumbbell_y   : int  293 296 298 303 292 294 295 300 292 291 ...
##  $ magnet_dumbbell_z   : num  -65 -64 -63 -60 -68 -66 -70 -74 -65 -69 ...
##  $ roll_forearm        : num  28.4 28.3 28.3 28.1 28 27.9 27.9 27.8 27.7 27.7 ...
##  $ pitch_forearm       : num  -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.8 -63.8 -6
## 3.8 ...
##  $ yaw_forearm         : num  -153 -153 -152 -152 -152 -152 -152 -152 -152 -152 ...
##  $ total_accel_forearm : int  36 36 36 36 36 36 36 36 36 36 ...
##  $ gyros_forearm_x     : num  0.03 0.02 0.03 0.02 0.02 0.02 0.02 0.02 0.03 0.02 ...
##  $ gyros_forearm_y     : num  0 0 -0.02 -0.02 0 -0.02 0 -0.02 0 0 ...
##  $ gyros_forearm_z     : num  -0.02 -0.02 0 0 -0.02 -0.03 -0.02 0 -0.02 -0.02 ...
##  $ accel_forearm_x     : int  192 192 196 189 189 193 195 193 193 190 ...
##  $ accel_forearm_y     : int  203 203 204 206 206 203 205 205 204 205 ...
##  $ accel_forearm_z     : int  -215 -216 -213 -214 -214 -215 -215 -213 -214 -215 ...
##  $ magnet_forearm_x    : int  -17 -18 -18 -16 -17 -9 -18 -9 -16 -22 ...
##  $ magnet_forearm_y    : num  654 661 658 658 655 660 659 660 653 656 ...
##  $ magnet_forearm_z    : num  476 473 469 469 473 478 470 474 476 473 ...
##  $ classe              : Factor w/ 5 levels "A","B","C","D",..: 1 1 1 1 1 1 1 1 1 1
## ...
```

# Partition the data

Create a training and validation data set out of the Slim Training data set

```
set.seed(8675309)
trainflag <- createDataPartition(slimTraining$classe, p=0.70, list=FALSE)
modelingData <- slimTraining[trainflag, ]
validationData <- slimTraining[-trainflag, ]
```
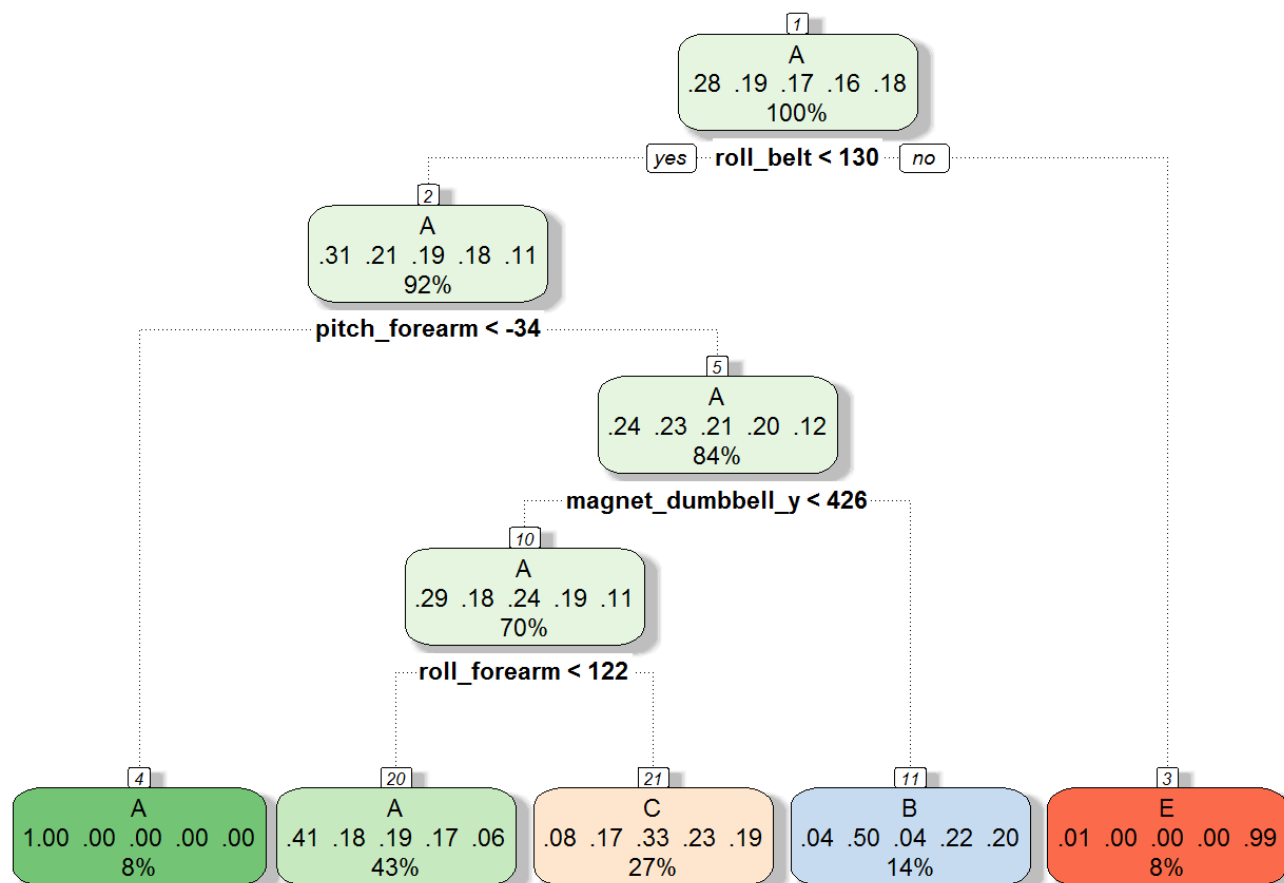
# Modeling the Data

we wish to predict the manner in which people exercise, using classe as the dependent variable

```
fit1 <- train(classe ~ .,method="rpart",data=modelingData)
```

```
## Loading required package: rpart
```

```
rattle::fancyRpartPlot(fit1$finalModel)
```



Rattle 2015-Jun-14 14:06:54 Matthew

```
predFit1 <- predict(fit1, validationData)
confusionMatrix(validationData$classe, predFit1)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##         A 1500   47  121    0    6
##         B  458  384  297    0    0
##         C  462   49  515    0    0
##         D  421  200  343    0    0
##         E  160  153  274    0  495
##
## Overall Statistics
##
##                Accuracy : 0.4918
##                  95% CI : (0.4789, 0.5046)
##     No Information Rate : 0.5099
##     P-Value [Acc > NIR] : 0.9975
##
##                   Kappa : 0.3365
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            0.4998  0.46098  0.33226       NA  0.98802
## Specificity            0.9397  0.85055  0.88212   0.8362  0.89097
## Pos Pred Value         0.8961  0.33714  0.50195       NA  0.45749
## Neg Pred Value         0.6436  0.90539  0.78699       NA  0.99875
## Prevalence             0.5099  0.14155  0.26338   0.0000  0.08513
## Detection Rate         0.2549  0.06525  0.08751   0.0000  0.08411
## Detection Prevalence   0.2845  0.19354  0.17434   0.1638  0.18386
## Balanced Accuracy      0.7198  0.65577  0.60719       NA  0.93950
```

```
#I had to really mess with this RF to make it work on my $500 computer. frown
fit2 <- train(classe ~ ., data=modelingData, method="rf", trControl=trainControl(metho
d="cv", 5), ntree=250)
```

```
## Loading required package: randomForest
## randomForest 4.6-10
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
##
## The following object is masked from 'package:dplyr':
##
##     combine
```

```
predFit2 <- predict(fit2, validationData)
confusionMatrix(validationData$classe, predFit2)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 1674    0    0    0    0
##          B    0 1139    0    0    0
##          C    0    0 1025    1    0
##          D    0    0    2  961    1
##          E    0    0    0    0 1082
##
## Overall Statistics
##
##                Accuracy : 0.9993
##                  95% CI : (0.9983, 0.9998)
##     No Information Rate : 0.2845
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.9991
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            1.0000   1.0000   0.9981   0.9990   0.9991
## Specificity            1.0000   1.0000   0.9998   0.9994   1.0000
## Pos Pred Value         1.0000   1.0000   0.9990   0.9969   1.0000
## Neg Pred Value         1.0000   1.0000   0.9996   0.9998   0.9998
## Prevalence             0.2845   0.1935   0.1745   0.1635   0.1840
## Detection Rate         0.2845   0.1935   0.1742   0.1633   0.1839
## Detection Prevalence   0.2845   0.1935   0.1743   0.1638   0.1839
## Balanced Accuracy      1.0000   1.0000   0.9989   0.9992   0.9995
```

the accuracy on the fit2 model, a randomforest, was much higher at 0.99. lets choose that one. Create predictions based upon the testing data of 20 observations

```
scoreTest <- predict(fit2, testing)
scoreTest
```

```
##  [1] B A B A A E D B A A B C B A E E A B B B
## Levels: A B C D E
```