

Title: “Practical Machine Learning - Course Project”

Author: “Uma Balakrishnan”

Date: “January 13, 2016”

## Question

6 participants were participated in a barbell lifting in 5 different ways.

Class A: Exactly according to the specification

Class B: Throwing the elbows to the front

Class C: Lifting the dumbbell only halfway

Class D: Lowering the dumbbell only halfway

Class E: Throwing the hips to the front

Class A perform barbell lifts correctly, while Classes B-E perform incorrectly.

By processing data gathered from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants in a machine learning algorithm, the question is can the appropriate activity quality (class A-E) can be predicted on testing data?

## Input Data:

Initialize library

```
library(AppliedPredictiveModeling)
library(caret)
library(randomForest)
```

```
library(rattle)
library(rpart.plot)
library(kernlab)
```

Downloading data from source and reading data. Treating empty values as NA.

```
if(!file.exists('pml-training.csv')){
  download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv",
               destfile='pml-training.csv') }
if(file.exists('pml-training.csv')){
  training <- read.csv('pml-training.csv', na.strings = c("", " ", "NA"), header = TRUE) }
str(training)
```

```
## '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 132
3084232 1323084232 1323084232 1323084232 1323084232 ...
##  $ raw_timestamp_part_2    : int   788290 808298 820366 120339 196328 304277 368296 440390 48
4323 484434 ...
##  $ 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      : Factor w/ 396 levels "-0.016850","-0.021024",...: NA NA NA NA NA
NA NA NA NA NA ...
## $ kurtosis_pitch_belt    : Factor w/ 316 levels "-0.021887","-0.060755",...: NA NA NA NA NA
NA NA NA NA NA ...
## $ kurtosis_yaw_belt      : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_belt     : Factor w/ 394 levels "-0.003095","-0.010002",...: NA NA NA NA NA
NA NA NA NA NA ...
## $ skewness_roll_belt.1   : Factor w/ 337 levels "-0.005928","-0.005960",...: NA NA NA NA NA
NA NA NA NA NA ...
## $ skewness_yaw_belt      : Factor w/ 1 level "#DIV/0!": NA NA NA NA 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           : Factor w/ 67 levels "-0.1","-0.2",...: NA NA NA NA NA NA NA NA N
A 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           : Factor w/ 67 levels "-0.1","-0.2",...: NA NA NA NA NA NA NA NA N
A 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     : Factor w/ 3 levels "#DIV/0!","0.00",...: 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 0.02 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 : Factor w/ 329 levels "-0.02438","-0.04190",...: NA NA NA NA NA N
A NA NA NA NA NA ...
## $ kurtosis_picth_arm : Factor w/ 327 levels "-0.00484","-0.01311",...: NA NA NA NA NA N
A NA NA NA NA NA ...
## $ kurtosis_yaw_arm  : Factor w/ 394 levels "-0.01548","-0.01749",...: NA NA NA NA NA N
A NA NA NA NA NA ...
## $ skewness_roll_arm : Factor w/ 330 levels "-0.00051","-0.00696",...: NA NA NA NA NA N
A NA NA NA NA NA ...
## $ skewness_pitch_arm : Factor w/ 327 levels "-0.00184","-0.01185",...: NA NA NA NA NA N
A NA NA NA NA NA ...
## $ skewness_yaw_arm  : Factor w/ 394 levels "-0.00311","-0.00562",...: NA NA NA NA NA N
A NA NA NA NA NA ...
## $ max_roll_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_picth_arm     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_arm       : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_arm     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_arm       : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm  : int   NA 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 : Factor w/ 397 levels "-0.0035","-0.0073",...: NA NA NA NA NA NA
NA NA NA NA NA ...
## $ kurtosis_picth_dumbbell : Factor w/ 400 levels "-0.0163","-0.0233",...: NA NA NA NA NA NA

```

```

NA NA NA NA ...
## $ kurtosis_yaw_dumbbell : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell : Factor w/ 400 levels "-0.0082","-0.0096",...: NA NA NA NA NA NA
NA NA NA NA ...
## $ skewness_pitch_dumbbell : Factor w/ 401 levels "-0.0053","-0.0084",...: NA NA NA NA NA NA
NA NA NA NA ...
## $ skewness_yaw_dumbbell : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA ...
## $ max_roll_dumbbell : num NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_dumbbell : Factor w/ 72 levels "-0.1","-0.2",...: NA NA NA NA NA NA NA NA N
A NA ...
## $ min_roll_dumbbell : num NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_dumbbell : Factor w/ 72 levels "-0.1","-0.2",...: NA NA NA NA NA NA NA NA N
A NA ...
## $ amplitude_roll_dumbbell : num NA NA NA NA NA NA NA NA NA ...
## [list output truncated]

```

```

if(!file.exists('pml-testing.csv')){
  download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv",
               destfile='pml-testing.csv') }
if(file.exists('pml-testing.csv')){
  testing <- read.csv('pml-testing.csv', na.strings = c("", " ", "NA"), header = TRUE) }
str(testing)

```

```

## '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 5 2 3
...
## $ raw_timestamp_part_1 : int  1323095002 1322673067 1322673075 1322832789 1322489635 132

```

```

2673149 1322673128 1322673076 1323084240 1322837822 ...
## $ raw_timestamp_part_2 : int 868349 778725 342967 560311 814776 510661 766645 54671 916
313 384285 ...
## $ 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_pitch_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_pitch_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_picth_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_picth_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_picth_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 ...
## $ max_roll_dumbbell : logi NA NA NA NA NA NA ...
## $ max_pitch_dumbbell : logi NA NA NA NA NA NA ...
## $ max_yaw_dumbbell : logi NA NA NA NA NA NA ...
## $ min_roll_dumbbell : logi NA NA NA NA NA NA ...
## $ min_pitch_dumbbell : logi NA NA NA NA NA NA ...
## $ min_yaw_dumbbell : logi NA NA NA NA NA NA ...
## $ amplitude_roll_dumbbell : logi NA NA NA NA NA NA ...
## [list output truncated]
```

# Features

Finding the columns not matching between training and testing sets

```
ind <- which(is.na(pmatch(names(training), names(testing))))
names(training)[ind]
```

```
## [1] "classe"
```

```
names(testing)[ind]
```

```
## [1] "problem_id"
```

Perform machine learning algorithm to column “classe” in training set to predict testing data.

For machine learning algorithm to predict effectively, drop the columns which has lot of NA's and also drop first 7 columns which consist of personal information regarding participants.

```
na_count <- sapply(training, function(y) sum(length(which(is.na(y)))))
na_count <- data.frame(na_count)
nonzeroind <- which(na_count == 0)
nonzeroind <- nonzeroind[8:length(nonzeroind)]
```

Following columns are considered in the training set for prediction algorithm.

```
training <- as.data.frame(training[,nonzeroind], drop = FALSE)
names(training)
```

```
## [1] "roll_belt" "pitch_belt" "yaw_belt"
## [4] "total_accel_belt" "gyros_belt_x" "gyros_belt_y"
## [7] "gyros_belt_z" "accel_belt_x" "accel_belt_y"
## [10] "accel_belt_z" "magnet_belt_x" "magnet_belt_y"
## [13] "magnet_belt_z" "roll_arm" "pitch_arm"
## [16] "yaw_arm" "total_accel_arm" "gyros_arm_x"
## [19] "gyros_arm_y" "gyros_arm_z" "accel_arm_x"
## [22] "accel_arm_y" "accel_arm_z" "magnet_arm_x"
## [25] "magnet_arm_y" "magnet_arm_z" "roll_dumbbell"
## [28] "pitch_dumbbell" "yaw_dumbbell" "total_accel_dumbbell"
## [31] "gyros_dumbbell_x" "gyros_dumbbell_y" "gyros_dumbbell_z"
## [34] "accel_dumbbell_x" "accel_dumbbell_y" "accel_dumbbell_z"
## [37] "magnet_dumbbell_x" "magnet_dumbbell_y" "magnet_dumbbell_z"
## [40] "roll_forearm" "pitch_forearm" "yaw_forearm"
## [43] "total_accel_forearm" "gyros_forearm_x" "gyros_forearm_y"
## [46] "gyros_forearm_z" "accel_forearm_x" "accel_forearm_y"
## [49] "accel_forearm_z" "magnet_forearm_x" "magnet_forearm_y"
## [52] "magnet_forearm_z" "classe"
```

Corresponding columns in the testing set.

```
testing <- as.data.frame(testing[,nonzeroind], drop = FALSE)
names(testing)
```

```
## [1] "roll_belt"      "pitch_belt"      "yaw_belt"
## [4] "total_accel_belt" "gyros_belt_x"    "gyros_belt_y"
## [7] "gyros_belt_z"    "accel_belt_x"    "accel_belt_y"
## [10] "accel_belt_z"    "magnet_belt_x"   "magnet_belt_y"
## [13] "magnet_belt_z"   "roll_arm"        "pitch_arm"
## [16] "yaw_arm"         "total_accel_arm" "gyros_arm_x"
## [19] "gyros_arm_y"     "gyros_arm_z"     "accel_arm_x"
## [22] "accel_arm_y"     "accel_arm_z"     "magnet_arm_x"
## [25] "magnet_arm_y"    "magnet_arm_z"    "roll_dumbbell"
## [28] "pitch_dumbbell"  "yaw_dumbbell"    "total_accel_dumbbell"
## [31] "gyros_dumbbell_x" "gyros_dumbbell_y" "gyros_dumbbell_z"
## [34] "accel_dumbbell_x" "accel_dumbbell_y" "accel_dumbbell_z"
## [37] "magnet_dumbbell_x" "magnet_dumbbell_y" "magnet_dumbbell_z"
## [40] "roll_forearm"    "pitch_forearm"   "yaw_forearm"
## [43] "total_accel_forearm" "gyros_forearm_x" "gyros_forearm_y"
## [46] "gyros_forearm_z"  "accel_forearm_x" "accel_forearm_y"
## [49] "accel_forearm_z"  "magnet_forearm_x" "magnet_forearm_y"
## [52] "magnet_forearm_z" "problem_id"
```

Check for covariates that have virtually no variability

```
nsv <- nearZeroVar(training, saveMetrics = TRUE)
nsv
```

```
##          freqRatio percentUnique zeroVar  nsv
## roll_belt      1.101904      6.7781062  FALSE FALSE
```

## pitch_belt	1.036082	9.3772296	FALSE	FALSE
## yaw_belt	1.058480	9.9734991	FALSE	FALSE
## total_accel_belt	1.063160	0.1477933	FALSE	FALSE
## gyros_belt_x	1.058651	0.7134849	FALSE	FALSE
## gyros_belt_y	1.144000	0.3516461	FALSE	FALSE
## gyros_belt_z	1.066214	0.8612782	FALSE	FALSE
## accel_belt_x	1.055412	0.8357966	FALSE	FALSE
## accel_belt_y	1.113725	0.7287738	FALSE	FALSE
## accel_belt_z	1.078767	1.5237998	FALSE	FALSE
## magnet_belt_x	1.090141	1.6664968	FALSE	FALSE
## magnet_belt_y	1.099688	1.5187035	FALSE	FALSE
## magnet_belt_z	1.006369	2.3290184	FALSE	FALSE
## roll_arm	52.338462	13.5256345	FALSE	FALSE
## pitch_arm	87.256410	15.7323412	FALSE	FALSE
## yaw_arm	33.029126	14.6570176	FALSE	FALSE
## total_accel_arm	1.024526	0.3363572	FALSE	FALSE
## gyros_arm_x	1.015504	3.2769341	FALSE	FALSE
## gyros_arm_y	1.454369	1.9162165	FALSE	FALSE
## gyros_arm_z	1.110687	1.2638875	FALSE	FALSE
## accel_arm_x	1.017341	3.9598410	FALSE	FALSE
## accel_arm_y	1.140187	2.7367241	FALSE	FALSE
## accel_arm_z	1.128000	4.0362858	FALSE	FALSE
## magnet_arm_x	1.000000	6.8239731	FALSE	FALSE
## magnet_arm_y	1.056818	4.4439914	FALSE	FALSE
## magnet_arm_z	1.036364	6.4468454	FALSE	FALSE
## roll_dumbbell	1.022388	84.2065029	FALSE	FALSE
## pitch_dumbbell	2.277372	81.7449801	FALSE	FALSE
## yaw_dumbbell	1.132231	83.4828254	FALSE	FALSE
## total_accel_dumbbell	1.072634	0.2191418	FALSE	FALSE
## gyros_dumbbell_x	1.003268	1.2282132	FALSE	FALSE
## gyros_dumbbell_y	1.264957	1.4167771	FALSE	FALSE

## gyros_dumbbell_z	1.060100	1.0498420	FALSE	FALSE
## accel_dumbbell_x	1.018018	2.1659362	FALSE	FALSE
## accel_dumbbell_y	1.053061	2.3748853	FALSE	FALSE
## accel_dumbbell_z	1.133333	2.0894914	FALSE	FALSE
## magnet_dumbbell_x	1.098266	5.7486495	FALSE	FALSE
## magnet_dumbbell_y	1.197740	4.3012945	FALSE	FALSE
## magnet_dumbbell_z	1.020833	3.4451126	FALSE	FALSE
## roll_forearm	11.589286	11.0895933	FALSE	FALSE
## pitch_forearm	65.983051	14.8557741	FALSE	FALSE
## yaw_forearm	15.322835	10.1467740	FALSE	FALSE
## total_accel_forearm	1.128928	0.3567424	FALSE	FALSE
## gyros_forearm_x	1.059273	1.5187035	FALSE	FALSE
## gyros_forearm_y	1.036554	3.7763735	FALSE	FALSE
## gyros_forearm_z	1.122917	1.5645704	FALSE	FALSE
## accel_forearm_x	1.126437	4.0464784	FALSE	FALSE
## accel_forearm_y	1.059406	5.1116094	FALSE	FALSE
## accel_forearm_z	1.006250	2.9558659	FALSE	FALSE
## magnet_forearm_x	1.012346	7.7667924	FALSE	FALSE
## magnet_forearm_y	1.246914	9.5403119	FALSE	FALSE
## magnet_forearm_z	1.000000	8.5771073	FALSE	FALSE
## classe	1.469581	0.0254816	FALSE	FALSE

Non zero variance is FALSE for all columns considered in the training set. Hence there is no need to eliminate any covariates due to lack of variability.

## Machine Learning Algorithm

Training set has 19,622 observations (large in size). It will be time consuming to perform algorithm on large data set.

**Step 1.** For classe variable, divide the given training set into 4 almost equal parts (part1, part2, part3, part4).

**Step 2.** Split each part into a training (60%) (part1.train, part2.train, part3.train, part4.train) and testing set (40%)

(part1.test, part2.test, part3.test, part4.test).

**Step 3.** Construct regression model using train on each training set (part1.train, part2.train, part3.train, part4.train) by defining method and cross validation.

**Step 4.** Predict respective testing set (part1.test, part2.test, part3.test, part4.test) and calculate the accuracy of prediction.

**Step 5.** Alter the method used for regression analysis, based on the accuracy for predicting testing sets.

**Step 6.** After achieving desired accuracy, predict the given testing set (20 Quiz prediction) using all four regression models and compare the results of prediction and accuracy.

### Step 1.

```
set.seed(2121)
partition <- createDataPartition(y = training$classe, p = 0.25, list = FALSE)
part1 <- training[partition,]
rest<- training[-partition,]
set.seed(2121)
partition <- createDataPartition(y = rest$classe, p = 0.33, list = FALSE)
part2 <- rest[partition,]
rest <- rest[-partition,]
set.seed(2121)
partition <- createDataPartition(y = rest$classe, p = 0.5, list = FALSE)
part3 <- rest[partition,]
part4 <- rest[-partition,]
```

### Step 2.

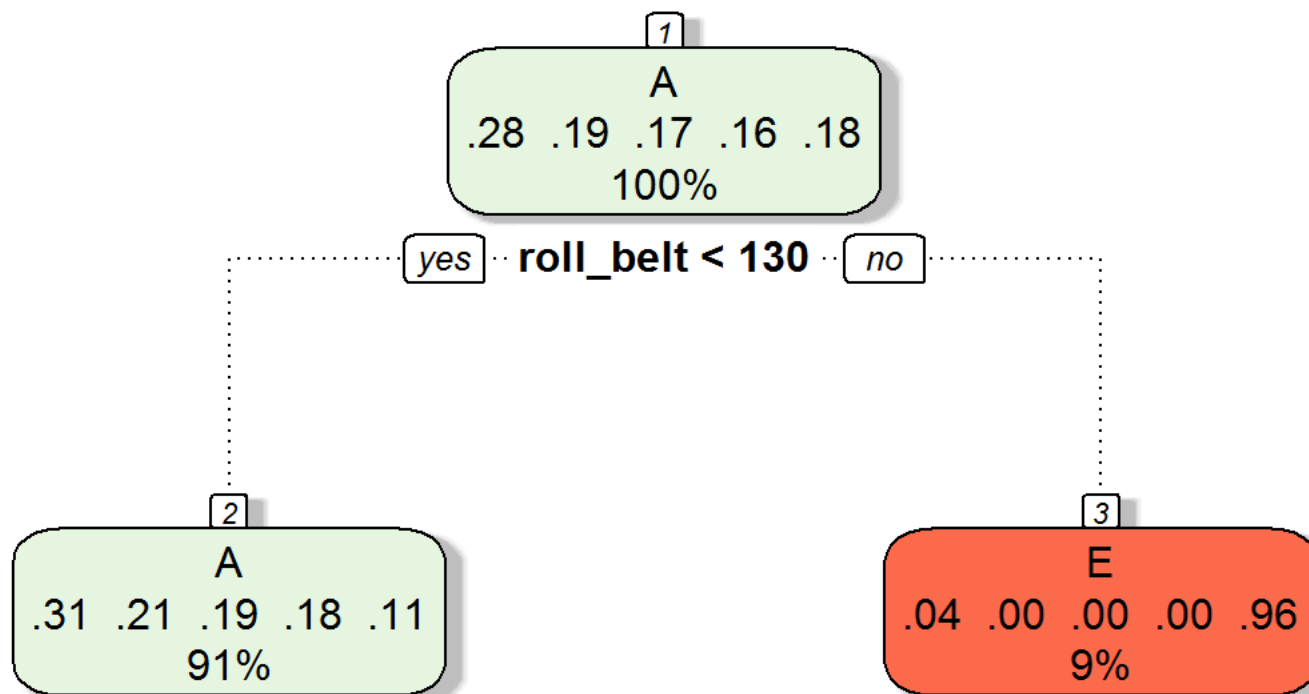
```
set.seed(2121)
inTrain <- createDataPartition(y = part1$classe, p = 0.6, list = FALSE)
part1.train <- part1[inTrain,]
part1.test <- part1[-inTrain,]
set.seed(2121)
inTrain <- createDataPartition(y = part2$classe, p = 0.6, list = FALSE)
```

```
part2.train <- part2[inTrain,]  
part2.test <- part2[-inTrain,]  
set.seed(2121)  
inTrain <- createDataPartition(y = part3$classe, p = 0.6, list = FALSE)  
part3.train <- part3[inTrain,]  
part3.test <- part3[-inTrain,]  
set.seed(2121)  
inTrain <- createDataPartition(y = part4$classe, p = 0.6, list = FALSE)  
part4.train <- part4[inTrain,]  
part4.test <- part4[-inTrain,]
```

### Step 3.

```
set.seed(2121)  
part1.modFit <- train(part1.train$classe ~., method = "rpart", data = part1.train,  
                      trControl = trainControl(method = "cv", number = 3))  
fancyRpartPlot(part1.modFit$finalModel)
```





Rattle 2016-Jan-14 01:49:10 Hariharan

```
part1.pred <- predict(part1.modFit, part1.test)
confusionMatrix(part1.pred, part1.test$classe)$overall['Accuracy']
```

```
## Accuracy
## 0.3686894
```

Accuracy is just 37% using method = "rpart" and cross validation. Accuracy is pretty low. We can try another method called randomForest, method = "rf".

#### Step 4.

Prediction of given testing set (20 Quiz Prediction) using part1.train

```
set.seed(2121)
part1.modFit <- train(part1.train$classe ~., method = "rf", data = part1.train,
                      trControl = trainControl(method = "cv", number = 3))
part1.pred <- predict(part1.modFit, part1.test)
part1.accuracy <- confusionMatrix(part1.pred, part1.test$classe)$overall['Accuracy']
part1.accuracy
```

```
## Accuracy
## 0.9602244
```

```
pred1 <- predict(part1.modFit, newdata = testing)
pred1
```

```
## [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
```

Prediction of given testing set (20 Quiz Prediction) using part2.train

```
set.seed(2121)
part2.modFit <- train(part2.train$classe ~., method = "rf", data = part2.train,
```

```
trControl = trainControl(method = "cv", number = 3))
part2.pred <- predict(part2.modFit, part2.test)
part2.accuracy <- confusionMatrix(part2.pred, part2.test$classe)$overall['Accuracy']
part2.accuracy
```

```
## Accuracy
## 0.9665121
```

```
pred2 <- predict(part2.modFit, newdata = testing)
pred2
```

```
## [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
```

Prediction of given testing set (20 Quiz Prediction) using part3.train

```
set.seed(2121)
part3.modFit <- train(part3.train$classe ~., method = "rf", data = part3.train,
                      trControl = trainControl(method = "cv", number = 3))
part3.pred <- predict(part3.modFit, part3.test)
part3.accuracy <- confusionMatrix(part3.pred, part3.test$classe)$overall['Accuracy']
part3.accuracy
```

```
## Accuracy
## 0.9563452
```

```
pred3 <- predict(part3.modFit, newdata = testing)
pred3
```

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

Prediction of given testing set (20 Quiz Prediction) using part4.train

```
set.seed(2121)
part4.modFit <- train(part4.train$classe ~., method = "rf", data = part4.train,
                      trControl = trainControl(method = "cv", number = 3))
part4.pred <- predict(part4.modFit, part4.test)
part4.accuracy <- confusionMatrix(part4.pred, part4.test$classe)$overall['Accuracy']
part4.accuracy
```

```
## Accuracy
## 0.960386
```

```
pred4 <- predict(part4.modFit, newdata = testing)
pred4
```

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

Comparing predictions from each part of the training data set.

```
pred <- data.frame(pred1, pred2, pred3, pred4)
names(pred) <- c("Data.part1", "Data.part2", "Data.part3", "Data.part4")
pred
```

##	Data.part1	Data.part2	Data.part3	Data.part4
## 1	B	B	B	B
## 2	A	A	A	A
## 3	B	B	A	A
## 4	A	A	A	A
## 5	A	A	A	A
## 6	E	E	E	E
## 7	D	D	D	D
## 8	B	B	B	B
## 9	A	A	A	A
## 10	A	A	A	A
## 11	B	B	B	B
## 12	C	C	C	C
## 13	B	B	B	B
## 14	A	A	A	A
## 15	E	E	E	E
## 16	E	E	E	E
## 17	A	A	A	A
## 18	B	B	B	B
## 19	B	B	B	A
## 20	B	B	B	B

Calculate accuracy of prediction given by each part of training set and calculate “out of sample error”

```
Accuracy <- data.frame(part1.accuracy, part2.accuracy, part3.accuracy, part4.accuracy)
Accuracy <- round(Accuracy * 100)
OutofSampleError <- 100 - Accuracy
Accuracy <- paste(Accuracy, '%', sep = " ")
```

Accuracy of prediction by each part of training data set is:

```
Accuracy
```

```
## [1] "96%" "97%" "96%" "96%"
```

## Out of Sample Error

Out of sample error for each part of training data set is:

```
OutofSampleError <- paste(OutofSampleError, '%', sep = "")  
OutofSampleError
```

```
## [1] "4%" "3%" "4%" "4%"
```

## Conclusion

Prediction for 20 test data is

```
pred1
```

```
## [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
```