# Peer Graded Assignment: Prediction Assignment Writeup

#### 1. Loading add-on package and set seed

```
set.seed(12345)

library(caret)
```

#### 2. Download rawdata and submit\_data

```
url_train <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-trainin
g.csv"

rawdata <- read.csv(url_train, na.strings = c("", "NA"))

url_submit <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testin
g.csv"

submit_data <- read.csv(url_submit, na.strings = c("", "NA"))</pre>
```

# 3. Cleaning data

We should delete the column that contains NA to avoid the error. In addition, in order to make accurate predictions, columns that is not related exercise must also be deleted. In particular "X", "user\_name", "raw\_timestamp\_part\_1", "raw\_timestamp\_part\_2", "cvtd\_timestamp", "new\_window", "num\_window" are deleted.

```
#Remove NA cols
colname <- colnames(rawdata)[!colSums(is.na(rawdata)) > 0]
colname
   [1] "X"
                                "user name"
                                                        "raw timestamp part 1"
   [4] "raw timestamp part 2" "cvtd timestamp"
                                                        "new window"
   [7] "num window"
                                "roll belt"
                                                        "pitch belt"
## [10] "yaw belt"
                                "total accel belt"
                                                        "gyros belt x"
## [13] "gyros belt y"
                                "gyros belt z"
                                                        "accel belt x"
## [16] "accel belt y"
                                "accel belt z"
                                                        "magnet belt x"
                                "magnet belt z"
## [19] "magnet belt y"
                                                        "roll arm"
## [22] "pitch arm"
                                "yaw arm"
                                                        "total accel arm"
```

```
## [25] "gyros arm x"
                              "gyros arm y"
                                                     "gyros arm z"
## [28] "accel arm x"
                              "accel arm y"
                                                     "accel arm z"
## [31] "magnet arm x"
                              "magnet arm y"
                                                     "magnet arm z"
## [34] "roll dumbbell"
                              "pitch dumbbell"
                                                     "yaw dumbbell"
  [37] "total accel dumbbell" "gyros dumbbell x"
                                                     "gyros dumbbell y"
  [40] "gyros dumbbell z"
                              "accel dumbbell x"
                                                     "accel dumbbell y"
## [43] "accel dumbbell z"
                              "magnet dumbbell x"
                                                     "magnet dumbbell y"
## [46] "magnet dumbbell z"
                              "roll forearm"
                                                     "pitch forearm"
## [49] "yaw forearm"
                              "total accel forearm" "gyros forearm x"
## [52] "gyros forearm y"
                              "gyros forearm z"
                                                     "accel forearm x"
## [55] "accel forearm y"
                              "accel forearm z"
                                                     "magnet forearm x"
## [58] "magnet forearm y"
                              "magnet forearm z"
                                                     "classe"
#Slice data related with exercise
colname <- colname[8: length(colname)]</pre>
df wo NA <- rawdata[colname]</pre>
#Check the colnames of df wo NA is in submit data.
#The last colname is "classe"
is.element(colname, colnames(submit data))
   [1] TRUE TRUE TRUE TRUE TRUE TRUE
                                           TRUE
                                                 TRUE TRUE TRUE
                                                                   TRUE
## [12] TRUE
              TRUE TRUE TRUE TRUE TRUE
                                            TRUE
                                                 TRUE TRUE TRUE
                                                                   TRUE
## [23] TRUE TRUE TRUE TRUE TRUE TRUE
                                                 TRUE TRUE TRUE TRUE
                                            TRUE
## [34] TRUE TRUE TRUE TRUE TRUE TRUE
                                            TRUE
                                                 TRUE TRUE TRUE TRUE
## [45] TRUE TRUE TRUE TRUE TRUE TRUE TRUE FALSE
```

## 4. Split data into random train and test

```
inTrain = createDataPartition(df_wo_NA$classe, p = 3/4)[[1]]
training = df_wo_NA[ inTrain,]
testing = df_wo_NA[-inTrain,]
```

#### 4. Random Forest

It takes a very long time for training, but it has a high accuracy.

```
model_rf <- train(classe ~ ., data = training, method = "rf")
pred_rf <- predict(model_rf, testing)</pre>
```

```
confusionMatrix(testing$classe, pred rf)
## Confusion Matrix and Statistics
##
##
           Reference
## Prediction
               A B C
                            D
                                 Ε
         A 1395
                        0
               9 934
##
##
               0
                    3 848
##
          D
               0
                    0
                        3 801
                                  0
##
               0
                    0
                        1
                            1 899
##
## Overall Statistics
##
##
                Accuracy: 0.9945
##
                  95% CI: (0.992, 0.9964)
##
     No Information Rate: 0.2863
##
     P-Value [Acc > NIR] : < 2.2e-16
##
                   Kappa: 0.993
##
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
                     Class: A Class: B Class: C Class: D Class: E
##
## Sensitivity
                      0.9936 0.9968 0.9883 0.9938 1.0000
## Specificity
                      1.0000 0.9962 0.9983 0.9993 0.9995
## Pos Pred Value
                      1.0000 0.9842 0.9918 0.9963 0.9978
## Neg Pred Value
                      0.9974 0.9992 0.9975 0.9988 1.0000
## Prevalence
                       0.2863 0.1911 0.1750
                                                0.1644 0.1833
                      0.2845
## Detection Rate
                               0.1905 0.1729
                                                0.1633 0.1833
## Detection Prevalence 0.2845 0.1935 0.1743 0.1639 0.1837
## Balanced Accuracy
                       0.9968 0.9965 0.9933
                                                0.9965 0.9998
```

# **5. Liner Discriminant Analysis**

It takes a short time but poor accuracy.

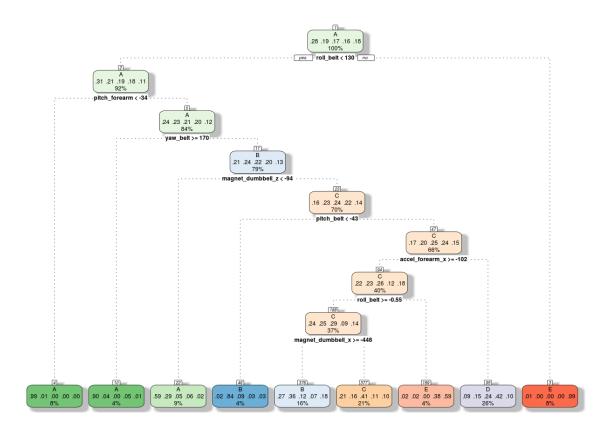
```
model lda <- train(classe ~ ., data = training, method = "lda")</pre>
pred lda <- predict(model lda, testing)</pre>
confusionMatrix(testing$classe, pred lda)
## Confusion Matrix and Statistics
##
          Reference
## Prediction A B C D E
##
         A 1150 29 103 103 10
         B 159 589 125 34 42
##
##
         C 100 71 558 101 25
         D 34 40 89 609 32
##
##
         E 31 151 84 83 552
##
## Overall Statistics
##
##
               Accuracy: 0.7051
##
                  95% CI: (0.6922, 0.7179)
     No Information Rate: 0.3006
##
     P-Value [Acc > NIR] : < 2.2e-16
##
##
##
                  Kappa : 0.6267
## Mcnemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##
                    Class: A Class: B Class: C Class: D Class: E
## Sensitivity
                      0.7802 0.6693 0.5819 0.6548 0.8351
## Specificity
                      0.9286 0.9105 0.9247 0.9509 0.9177
## Pos Pred Value
                      0.8244 0.6207 0.6526 0.7575 0.6127
## Neg Pred Value 0.9077 0.9264 0.9010 0.9217 0.9728
                      0.3006 0.1794 0.1956 0.1896 0.1348
## Prevalence
                  0.2345 0.1201 0.1138 0.1242 0.1126
## Detection Rate
## Detection Prevalence 0.2845 0.1935 0.1743 0.1639 0.1837
```

## 6. Recursive Partitioning and Regression Trees

The results can be confirmed visually, but poor accuracy.

```
model rpart <- train(classe ~ ., data = training, method = "rpart")</pre>
pred_rpart<- predict(model_rpart, testing)</pre>
confusionMatrix(testing$classe, pred rpart)
## Confusion Matrix and Statistics
##
           Reference
## Prediction A B C
         A 859 215 203 109
##
          В 151 434 166 197
         C 27 117 442 269 0
##
##
          D 43 65 131 492 73
          E 14 143 111 133 500
##
##
## Overall Statistics
##
##
                Accuracy: 0.5561
                  95% CI: (0.542, 0.57)
##
     No Information Rate: 0.2447
##
     P-Value [Acc > NIR] : < 2.2e-16
##
##
##
                   Kappa : 0.4442
   Mcnemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
                      Class: A Class: B Class: C Class: D Class: E
##
## Sensitivity
                       0.7852  0.4456  0.41975  0.4100  0.8576
## Specificity
                       0.6158 0.4573 0.51696 0.6119 0.5549
## Pos Pred Value
## Neg Pred Value
                       0.9330 0.8635 0.84910 0.8273 0.9793
```

```
## Prevalence
                          0.2231
                                   0.1986 0.21472
                                                      0.2447
                                                               0.1189
                          0.1752
                                   0.0885 0.09013
## Detection Rate
                                                      0.1003
                                                               0.1020
## Detection Prevalence
                          0.2845
                                   0.1935
                                           0.17435
                                                      0.1639
                                                               0.1837
## Balanced Accuracy
                          0.8223
                                   0.6573 0.65625
                                                      0.6629
                                                               0.8824
library(rattle)
fancyRpartPlot(model rpart$finalModel)
```



#### 7. Submit data with Random Forest

We can use the high accuracy model to submit data. In this report the Random Forest accuracy has the highest value 99.45. We can show the prediction.

```
submit_rf <- predict(model_rf, submit_data)
submit_rf
## [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</pre>
```