

# Predicting using the PML dataset

## 1. loading the training and test sets

```
library(caret)
```

```
## Loading required package: lattice
```

```
## Loading required package: ggplot2
```

```
data <- read.csv("pml-training.csv")
finaltest <- read.csv("pml-testing.csv")
```

```
names(data)
```

```
##      [1] "X"                                "user_name"
##      [3] "raw_timestamp_part_1"           "raw_timestamp_part_2"
##      [5] "cvtd_timestamp"                "new_window"
##      [7] "num_window"                    "roll_belt"
##      [9] "pitch_belt"                    "yaw_belt"
##     [11] "total_accel_belt"              "kurtosis_roll_belt"
##     [13] "kurtosis_pitch_belt"           "kurtosis_yaw_belt"
##     [15] "skewness_roll_belt"            "skewness_roll_belt.1"
##     [17] "skewness_yaw_belt"             "max_roll_belt"
##     [19] "max_pitch_belt"                "max_yaw_belt"
##     [21] "min_roll_belt"                 "min_pitch_belt"
##     [23] "min_yaw_belt"                  "amplitude_roll_belt"
##     [25] "amplitude_pitch_belt"          "amplitude_yaw_belt"
##     [27] "var_total_accel_belt"          "avg_roll_belt"
##     [29] "stddev_roll_belt"              "var_roll_belt"
##     [31] "avg_pitch_belt"                "stddev_pitch_belt"
##     [33] "var_pitch_belt"                "avg_yaw_belt"
##     [35] "stddev_yaw_belt"               "var_yaw_belt"
##     [37] "gyros_belt_x"                  "gyros_belt_y"
##     [39] "gyros_belt_z"                  "accel_belt_x"
##     [41] "accel_belt_y"                  "accel_belt_z"
##     [43] "magnet_belt_x"                 "magnet_belt_y"
##     [45] "magnet_belt_z"                 "roll_arm"
##     [47] "pitch_arm"                     "yaw_arm"
##     [49] "total_accel_arm"               "var_accel_arm"
##     [51] "avg_roll_arm"                  "stddev_roll_arm"
##     [53] "var_roll_arm"                  "avg_pitch_arm"
##     [55] "stddev_pitch_arm"              "var_pitch_arm"
##     [57] "avg_yaw_arm"                   "stddev_yaw_arm"
##     [59] "var_yaw_arm"                   "gyros_arm_x"
##     [61] "gyros_arm_y"                   "gyros_arm_z"
##     [63] "accel_arm_x"                   "accel_arm_y"
##     [65] "accel_arm_z"                   "magnet_arm_x"
##     [67] "magnet_arm_y"                   "magnet_arm_z"
##     [69] "kurtosis_roll_arm"             "kurtosis_pitch_arm"
##     [71] "kurtosis_yaw_arm"              "skewness_roll_arm"
##     [73] "skewness_pitch_arm"            "skewness_yaw_arm"
```

```
## [75] "max_roll_arm"          "max_picth_arm"
## [77] "max_yaw_arm"           "min_roll_arm"
## [79] "min_pitch_arm"         "min_yaw_arm"
## [81] "amplitude_roll_arm"    "amplitude_pitch_arm"
## [83] "amplitude_yaw_arm"     "roll_dumbbell"
## [85] "pitch_dumbbell"        "yaw_dumbbell"
## [87] "kurtosis_roll_dumbbell" "kurtosis_picth_dumbbell"
## [89] "kurtosis_yaw_dumbbell" "skewness_roll_dumbbell"
## [91] "skewness_pitch_dumbbell" "skewness_yaw_dumbbell"
## [93] "max_roll_dumbbell"     "max_picth_dumbbell"
## [95] "max_yaw_dumbbell"      "min_roll_dumbbell"
## [97] "min_pitch_dumbbell"    "min_yaw_dumbbell"
## [99] "amplitude_roll_dumbbell" "amplitude_pitch_dumbbell"
## [101] "amplitude_yaw_dumbbell" "total_accel_dumbbell"
## [103] "var_accel_dumbbell"    "avg_roll_dumbbell"
## [105] "stddev_roll_dumbbell"  "var_roll_dumbbell"
## [107] "avg_pitch_dumbbell"    "stddev_pitch_dumbbell"
## [109] "var_pitch_dumbbell"    "avg_yaw_dumbbell"
## [111] "stddev_yaw_dumbbell"   "var_yaw_dumbbell"
## [113] "gyros_dumbbell_x"      "gyros_dumbbell_y"
## [115] "gyros_dumbbell_z"      "accel_dumbbell_x"
## [117] "accel_dumbbell_y"      "accel_dumbbell_z"
## [119] "magnet_dumbbell_x"     "magnet_dumbbell_y"
## [121] "magnet_dumbbell_z"     "roll_forearm"
## [123] "pitch_forearm"         "yaw_forearm"
## [125] "kurtosis_roll_forearm" "kurtosis_picth_forearm"
## [127] "kurtosis_yaw_forearm"  "skewness_roll_forearm"
## [129] "skewness_pitch_forearm" "skewness_yaw_forearm"
## [131] "max_roll_forearm"      "max_picth_forearm"
## [133] "max_yaw_forearm"       "min_roll_forearm"
## [135] "min_pitch_forearm"     "min_yaw_forearm"
## [137] "amplitude_roll_forearm" "amplitude_pitch_forearm"
## [139] "amplitude_yaw_forearm"  "total_accel_forearm"
## [141] "var_accel_forearm"     "avg_roll_forearm"
## [143] "stddev_roll_forearm"   "var_roll_forearm"
## [145] "avg_pitch_forearm"     "stddev_pitch_forearm"
## [147] "var_pitch_forearm"     "avg_yaw_forearm"
## [149] "stddev_yaw_forearm"    "var_yaw_forearm"
## [151] "gyros_forearm_x"       "gyros_forearm_y"
## [153] "gyros_forearm_z"       "accel_forearm_x"
## [155] "accel_forearm_y"       "accel_forearm_z"
## [157] "magnet_forearm_x"      "magnet_forearm_y"
## [159] "magnet_forearm_z"      "classe"
```

## 2. splitting data to training and testing sets

```
inTrain <- createDataPartition(data$classe, p = 0.7, list = FALSE)
training <- data[inTrain,]
testing <- data[-inTrain,]

dim(training)
```

```
## [1] 13737 160
```

```
dim(testing)
```

```
## [1] 5885 160
```

### 3. excluding variables with variance near zero

```
nzv <- nearZeroVar(training)
training <- training[,-nzv]
testing <- testing[,-nzv]
```

### 4. excluding variables that are mostly NA

```
allNA <- sapply(training, function(x) mean(is.na(x)) > 0.95)
training <- training[,allNA == FALSE]
testing <- testing[,allNA == FALSE]
dim(training)
```

```
## [1] 13737 59
```

```
dim(testing)
```

```
## [1] 5885 59
```

### 5. excluding id variables (from 1 to 5)

```
training <- training[,-(1:5)]
testing <- testing[,-(1:5)]
names(training)
```

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

```
dim(training)
```

```
## [1] 13737    54
```

## 5. training a random forest model

```
set.seed(12345)
controlRF <- trainControl(method = "cv", number = 3, verboseIter = FALSE)
RFmodel <- train(classe ~ ., method = "rf", data = training, trControl = controlRF)
```

## 6. checking out-of-sample accuracy

```
preds <- predict(RFmodel, testing)
confusionMatrix(preds, testing$classe)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction    A    B    C    D    E
##           A 1673     1     0     0     0
##           B     1 1137     3     0     0
##           C     0     1 1023     2     0
##           D     0     0     0  962     2
##           E     0     0     0     0 1080
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##           Accuracy : 0.9983
##           95% CI : (0.9969, 0.9992)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##           Kappa : 0.9979
```

```
##
```

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

```
##
```

```
## Statistics by Class:
```

```
##
```

```
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9994  0.9982  0.9971  0.9979  0.9982
## Specificity      0.9998  0.9992  0.9994  0.9996  1.0000
## Pos Pred Value   0.9994  0.9965  0.9971  0.9979  1.0000
## Neg Pred Value   0.9998  0.9996  0.9994  0.9996  0.9996
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2843  0.1932  0.1738  0.1635  0.1835
## Detection Prevalence 0.2845  0.1939  0.1743  0.1638  0.1835
## Balanced Accuracy 0.9996  0.9987  0.9982  0.9988  0.9991
```

## 7. predicting the 20 test samples

```
testpred <- predict(RFmodel, finaltest)
testDF <- data.frame(id = seq(length(testpred)), class = testpred )
testDF
```

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