# Prediction Assignment Writeup

Julia Maier

2025-12-08

## 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, our 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. More information is available from the website here: http://web.archive.org/web/20161224072740/http:/groupware.les.inf.puc-rio.br/har (web.archive.org) (see the section on the Weight Lifting Exercise Dataset).

**Reference**:

Velloso, E.; Bulling, A.; Gellersen, H.; Ugulino, W.; Fuks, H. Qualitative Activity Recognition of Weight Lifting Exercises. Proceedings of 4th International Conference in Cooperation with SIGCHI (Augmented Human '13) . Stuttgart, Germany: ACM SIGCHI, 2013.

## Goal of this project

The goal of this project is to predict the manner in which participants did the predefined 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.

## Preparation of dataset

In a first step, data are loaded from the website. They are then processed to create clean datasets.

```r
# Create datasets

## Download and save datasets

url_train <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
url_test  <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"

if(!file.exists("pml-training.csv")) {
  download.file(url_train, "pml-training.csv", method = "curl")
}

if(!file.exists("pml-testing.csv")) {
  download.file(url_test, "pml-testing.csv", method = "curl")
}


## Load TRAINING DATA

# 1) Read training data
raw <- readLines("pml-training.csv", warn = FALSE)

# 2) Remove Backslash-escaped quotes (\"), and all " Signs
    raw <- gsub('\\\\\"', '', raw, perl = TRUE)
    clean <- gsub('"', '', raw, fixed = TRUE)

# 3) Save clean data set
    writeLines(clean, "pml-training_noquotes.csv", useBytes = TRUE)

    data_train <- read.csv("pml-training_noquotes.csv", na.strings = c("NA", "#DIV/0!", ""),
stringsAsFactors = FALSE)


## Load TEST DATA

# 1) Read testing data
raw <- readLines("pml-testing.csv", warn = FALSE)

# 2) Remove Backslash-escaped quotes (\"), and all " Signs
    raw <- gsub('\\\\\"', '', raw, perl = TRUE)
    clean <- gsub('"', '', raw, fixed = TRUE)

# 3) Save clean data set
    writeLines(clean, "pml-testing_noquotes.csv", useBytes = TRUE)

    data_test <- read.csv("pml-testing_noquotes.csv", na.strings = c("NA", "#DIV/0!", ""), st
ringsAsFactors = FALSE)

dim(data_train)
```

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

```r
dim(data_test)
```

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

```r
 ## Remove NA > 95%

  na_thresh <- 0.95
  valid_cols <- colMeans(is.na(data_train)) < na_thresh
  data_train <- data_train[, valid_cols]

 ## Remove irrelevant variables
  drop_cols <- c("X", "user_name", "cvtd_timestamp", "raw_timestamp_part_1", "raw_timestamp_p
art_2", "new_window", "num_window")

  data_train <- data_train[, !(names(data_train) %in% drop_cols)]


 ## Define Factor

  data_train$classe <- factor(data_train$classe)


 ## Data overview

  str(data_test)
```

```
## 'data.frame':    20 obs. of  160 variables:
##  $ X                    : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ user_name            : chr  "pedro" "jeremy" "jeremy" "adelmo" ...
##  $ raw_timestamp_part_1 : int  1323095002 1322673067 1322673075 1322832789 1322489635 1
322673149 1322673128 1322673076 1323084240 1322837822 ...
##  $ raw_timestamp_part_2 : int  868349 778725 342967 560311 814776 510661 766645 54671 9
16313 384285 ...
##  $ cvtd_timestamp       : chr  "05/12/2011 14:23" "30/11/2011 17:11" "30/11/2011 17:11"
"02/12/2011 13:33" ...
##  $ new_window           : chr  "no" "no" "no" "no" ...
##  $ 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 ...
##  $ kurtosis_picth_arm       : logi  NA NA NA NA NA NA ...
##  $ kurtosis_yaw_arm         : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_arm        : logi  NA NA NA NA NA NA ...
##  $ skewness_pitch_arm       : logi  NA NA NA NA NA NA ...
##  $ skewness_yaw_arm         : logi  NA NA NA NA NA NA ...
##  $ max_roll_arm             : logi  NA NA NA NA NA NA ...
##  $ max_picth_arm            : logi  NA NA NA NA NA NA ...
##  $ max_yaw_arm              : logi  NA NA NA NA NA NA ...
##  $ min_roll_arm             : logi  NA NA NA NA NA NA ...
##  $ min_pitch_arm            : logi  NA NA NA NA NA NA ...
##  $ min_yaw_arm              : logi  NA NA NA NA NA NA ...
##  $ amplitude_roll_arm       : logi  NA NA NA NA NA NA ...
##  $ amplitude_pitch_arm      : logi  NA NA NA NA NA NA ...
##  $ amplitude_yaw_arm        : logi  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 ...
##  $ kurtosis_picth_dumbbell  : logi  NA NA NA NA NA NA ...
##  $ kurtosis_yaw_dumbbell    : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_dumbbell   : logi  NA NA NA NA NA NA ...
##  $ skewness_pitch_dumbbell  : logi  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_picth_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]
```

```
str(data_train)
```

```
## 'data.frame':    19622 obs. of  53 variables:
##  $ 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 ...
##  $ 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 -63.8
## ...
##  $ 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 ...
```

## Gradient Boosting Model

As a first step, a Gradient Boosting Model was built and trained.

```
set.seed(123)

ctrl <- trainControl(method = "cv", number = 3)

gbm_cv <- train(
  classe ~ .,
  data = data_train,
  method = "gbm",
  trControl = ctrl,
  verbose = FALSE
)

gbm_cv
```

```
## Stochastic Gradient Boosting
##
## 19622 samples
##    52 predictor
##     5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (3 fold)
## Summary of sample sizes: 13083, 13080, 13081
## Resampling results across tuning parameters:
##
##   interaction.depth  n.trees  Accuracy   Kappa
##   1                   50      0.7498725  0.6828383
##   1                  100      0.8195895  0.7716751
##   1                  150      0.8535822  0.8147405
##   2                   50      0.8568947  0.8187154
##   2                  100      0.9062777  0.8814089
##   2                  150      0.9336451  0.9160353
##   3                   50      0.8930277  0.8645596
##   3                  100      0.9415956  0.9260986
##   3                  150      0.9606050  0.9501609
##
## Tuning parameter 'shrinkage' was held constant at a value of 0.1
##
## Tuning parameter 'n.minobsinnode' was held constant at a value of 10
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were n.trees = 150, interaction.depth =
##   3, shrinkage = 0.1 and n.minobsinnode = 10.
```

THe model has been cross-validated using 3 fold cross-validation. The resulting accuracy was .961 with Kappa = .950.

# Random Forest Model

In a second step, a Random Forest Model is trained. It is then tested in the test data set.

```
set.seed(1122)

## Cross-validation

  ctrl <- trainControl(method = "cv",
                    number = 5,
                    verboseIter = TRUE,
                    allowParallel = TRUE)

  rf_cv <- train(classe ~ .,
              data = data_train,
              method = "rf",
              trControl =  trainControl(method = "cv", number = 3))

  rf_cv
```

```
## Random Forest
##
## 19622 samples
##    52 predictor
##     5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (3 fold)
## Summary of sample sizes: 13082, 13080, 13082
## Resampling results across tuning parameters:
##
##   mtry  Accuracy   Kappa
##    2    0.9923045  0.9902658
##   27    0.9920496  0.9899429
##   52    0.9865968  0.9830434
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 2.
```

The model has been cross-validated using 3 fold cross-validation. The resulting accuracy was .993 (N = 13081) with Kappa = .991.

```
# Select variables
test_use <- data_test[, names(data_train)[names(data_train) != "classe"]]


# use test data set

rf_pred1 <- predict(gbm_cv, test_use)
rf_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
```

```
rf_pred2 <- predict(rf_cv, test_use)
rf_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
```
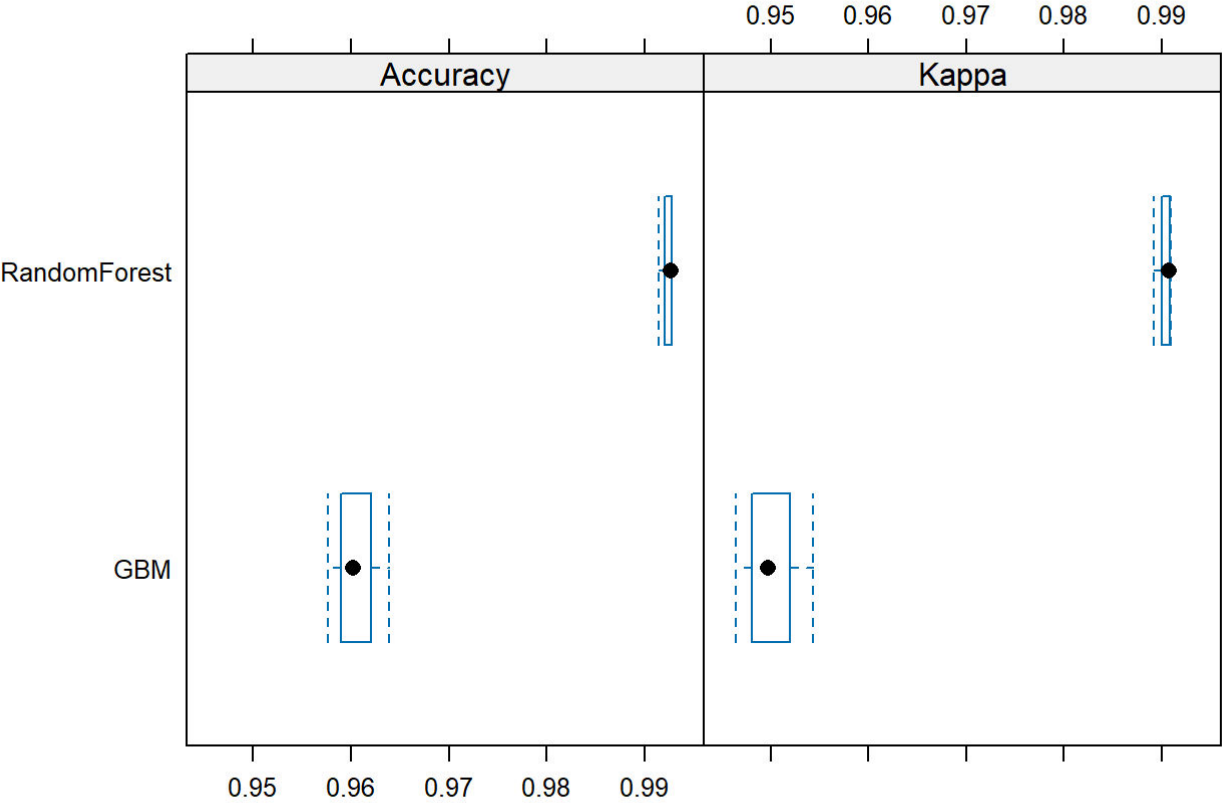
# Comparison of models

In a last step, both models are compared as regards their prediction accuracy.

```
results <- resamples(list(RandomForest = rf_cv, GBM = gbm_cv))
summary(results)
```

```
##
## Call:
## summary.resamples(object = results)
##
## Models: RandomForest, GBM
## Number of resamples: 3
##
## Accuracy
##                   Min.   1st Qu.    Median      Mean   3rd Qu.      Max. NA's
## RandomForest 0.9914373 0.9920489 0.9926606 0.9923045 0.9927381 0.9928157    0
## GBM          0.9576388 0.9589448 0.9602507 0.9606050 0.9620881 0.9639254    0
##
## Kappa
##                   Min.   1st Qu.    Median      Mean   3rd Qu.      Max. NA's
## RandomForest 0.9891709 0.9899429 0.9907148 0.9902658 0.9908132 0.9909117    0
## GBM          0.9464185 0.9480545 0.9496904 0.9501609 0.9520321 0.9543738    0
```

```
bwplot(results, main = "Model Comparison (CV)")
```

# Model Comparison (CV)



Considering the direct comparison of both models, Random Forest allows for a better prediction than Gradient Boosting.