

# Week 4 Assignment: Predicting Exercise Classe from Wearable Devices

## Analysis

### Environment setup

```
library(caret)
library(randomForest)

if (!file.exists('train.csv')) {
  download.file(url = 'https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv',
    destfile = 'train.csv', method = 'curl', quiet = TRUE)
}

if (!file.exists('test.csv')) {
  download.file(url = 'https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv',
    destfile = 'test.csv', method = 'curl', quiet = TRUE)
}

trainRaw <- read.csv('train.csv')
testRaw <- read.csv('test.csv')
```

### Preprocessing

1. First look at the data for each column and remove variables unrelated to exercise (column number and time stamps):

```
str(trainRaw)
```

```
## 'data.frame': 19622 obs. of 160 variables:
## $ X : int 1 2 3 4 5 6 7 8 9 10 ...
## $ user_name : chr "carlitos" "carlitos" "carlitos" "carlitos" ...
## $ raw_timestamp_part_1 : int 1323084231 1323084231 1323084231 1323084232 1323084232 1323084232 1323084232 1323084232 1323084232 ...
## $ raw_timestamp_part_2 : int 788290 808298 820366 120339 196328 304277 368296 440390 484323 484434 ...
## $ cvtd_timestamp : chr "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23" ...
## $ new_window : chr "no" "no" "no" "no" ...
## $ 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 : chr "" "" "" "" "" ...
## $ kurtosis_pitch_belt : chr "" "" "" "" "" ...
## $ kurtosis_yaw_belt : chr "" "" "" "" "" ...
## $ skewness_roll_belt : chr "" "" "" "" "" ...
## $ skewness_roll_belt.1 : chr "" "" "" "" "" ...
## $ skewness_yaw_belt : chr "" "" "" "" "" ...
## $ max_roll_belt : num NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_belt : int NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_belt : chr "" "" "" "" "" ...
## $ min_roll_belt : num NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_belt : int NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_belt : chr "" "" "" "" "" ...
## $ amplitude_roll_belt : num NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_belt : int NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_belt : chr "" "" "" "" "" ...
## $ var_total_accel_belt : num NA NA NA NA NA NA NA NA NA ...
## $ avg_roll_belt : num NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_belt : num NA NA NA NA NA NA NA NA NA ...
## $ var_roll_belt : num NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_belt : num NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_belt : num NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_belt : num NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_belt : num NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_belt : num NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_belt : num 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 0.03 ...
## $ gyros_belt_y : num 0 0 0 0.02 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 ...
## $ 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 ...
## $ avg_roll_arm : num NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_arm : num NA NA NA NA NA NA NA NA NA ...
## $ var_roll_arm : num NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_arm : num NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_arm : num NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_arm : num NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_arm : num NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_arm : num NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_arm : num 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 ...
## $ evros 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 -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 : chr "" "" "" "" "" ...
## $ kurtosis_pitch_arm : chr "" "" "" "" "" ...
## $ kurtosis_yaw_arm : chr "" "" "" "" "" ...
## $ skewness_roll_arm : chr "" "" "" "" "" ...
## $ skewness_pitch_arm : chr "" "" "" "" "" ...
## $ skewness_yaw_arm : chr "" "" "" "" "" ...
## $ max_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_arm : int NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_arm : int NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm : int 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 : chr "" "" "" "" "" ...
## $ kurtosis_pitch_dumbbell : chr "" "" "" "" "" ...
## $ kurtosis_yaw_dumbbell : chr "" "" "" "" "" ...
## $ skewness_roll_dumbbell : chr "" "" "" "" "" ...
## $ skewness_pitch_dumbbell : chr "" "" "" "" "" ...
## $ skewness_yaw_dumbbell : chr "" "" "" "" "" ...
## $ max_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_dumbbell : chr "" "" "" "" "" ...
## $ min_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_dumbbell : chr "" "" "" "" "" ...
## $ amplitude_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]
```

```
train <- trainRaw[, 6:ncol(trainRaw)]
```

2. Split the data into 70% training and 30% testing set :

```
set.seed(12345)
inTrain <- createDataPartition(y = train$classe, p = 0.7, list = F)
training <- train[inTrain, ]
testing <- train[-inTrain, ]
```

3. Remove the variables with a lot of similarities :

```
nzv <- nearZeroVar(train, saveMetrics = T)
keepFeat <- row.names(nzv[nzv$nzv == FALSE, ])
training <- training[, keepFeat]
```

4. Remove the variables with all NAs :

```
training <- training[, colSums(is.na(training)) == 0]
dim(training)
```

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

This is a rather stringent cutoff but there is still >50 features after removal !

## Model training

1. Set up 5-fold cross validation for training :

```
modCtl <- trainControl(method = 'cv', number = 5)
```

2. Fit a model with random forests :

```
set.seed(12345)
modRf <- train(classe ~., data = training, method = 'rf', trControl = modCtl)
```

• Read the summary of the model built with random forests :

```
modRf$finalModel
```

```
##
## Call:
## randomForest(x = x, y = y, mtry = param$mtry)
##              Type of random forest: classification
##              Number of trees: 500
## No. of variables tried at each split: 27
##
## OOB estimate of error rate: 0.23%
## Confusion matrix:
##      A   B   C   D   E class.error
## A 3904   1   0   0   1 0.0005120328
## B   6 2648   3   1   0 0.0037622272
## C   0   6 2388   2   0 0.0033388982
## D   0   0   8 2244   0 0.0035523979
## E   0   0   0   3 2522 0.0011881188
```

- Predict with the validation set and check the confusion matrix and accuracy :

```
predRf <- predict(modRf, newdata = testing)
tst = factor(testing$classe)
confusionMatrix(predRf, tst)$overall[1]
```

```
## Accuracy
## 0.9991504
```

```
confusionMatrix(predRf, tst)$table
```

```
##           Reference
## Prediction  A   B   C   D   E
##      A 1674   1   0   0   0
##      B    0 1138   1   0   0
##      C    0   0 1025   2   0
##      D    0   0   0 962   1
##      E    0   0   0   0 1081
```

The accuracy is **~99.6%** under **5-fold** cross validation

3. Fit a model with gradient boosting method :

```
modGbm <- train(classe ~., data = training, method = 'gbm', trControl = modCtl, verbose = F)
```

- Read the summary of the model built with gbm :

```
modGbm$finalModel
```

```
## A gradient boosted model with multinomial loss function.
## 150 iterations were performed.
## There were 53 predictors of which 53 had non-zero influence.
```

- Predict with the validation set and check the confusion matrix and accuracy :

```
predGbm <- predict(modGbm, newdata = testing)
tst = factor(testing$classe)
confusionMatrix(predRf, tst)$overall[1]
```

```
## Accuracy
## 0.9991504
```

```
confusionMatrix(predRf, tst)$table
```

```
##           Reference
## Prediction  A   B   C   D   E
##      A 1674   1   0   0   0
##      B    0 1138   1   0   0
##      C    0   0 1025   2   0
##      D    0   0   0 962   1
##      E    0   0   0   0 1081
```

The accuracy is **~98.8%** under **5-fold** cross validation

## Quiz

Since random forests gives the highest accuracy under the validation set, this model will be selected and used for prediction in the test set :

```
predRfTest <- predict(modRf, newdata = testRaw)
predRfTest
```

```
## [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
```

The gbm model can also be used for prediction and the results can be compared to above :

```
predGbmTest <- predict(modGbm, newdata = testRaw)
table(predRfTest, predGbmTest)
```

```
##           predGbmTest
## predRfTest A B C D E
##      A 7 0 0 0 0
##      B 0 8 0 0 0
##      C 0 0 1 0 0
##      D 0 0 0 1 0
##      E 0 0 0 0 3
```

The two models produce the same results, as shown in the confusion matrix !