

Prediction Assignment

Christian Edelmayer

6 7 2022

Summary

Setup

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.0.5
```

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

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

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

```
library(randomForest)
```

```
## randomForest 4.6-14
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##
```

```
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
##      combine
```

```
## The following object is masked from 'package:ggplot2':
```

```
##
```

```
##      margin
```

Loading & Reading Data

```
trainingRaw <- read.csv("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv")
testingRaw <- read.csv("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv")
```

Data Exploration & Preprocessing

```
dim(trainingRaw)
```

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

```
dim(testingRaw)
```

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

We can see that the training set consists of 19622 observations, while the testing set consists of only 20 observations. Each observation consists of 160 variables.

Let's look at the variables.

```
str(trainingRaw)
```

```
## '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 ...
## $ raw_timestamp_part_2 : int 788290 808298 820366 120339 196328 304277 368296 440390 484323 484...
## $ cvtd_timestamp : chr "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23" "05/12/20...
## $ 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 NA ...
## $ max_pitch_belt : int NA 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 NA ...
## $ min_pitch_belt : int NA 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 NA ...
## $ amplitude_pitch_belt : int NA 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 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 0.03 ...
## $ 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       : 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 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 : 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 NA ...
## $ max_pitch_dumbbell  : num  NA 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 NA ...
## $ min_pitch_dumbbell  : num  NA 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 NA ...
## [list output truncated]
```

We have a lot of variables (160), so we obviously have to reduce that number. Also, a lot of variables seem to be NA, so we decide to simply omit all variables with NA values. Note that this is a very simply imputation method and more advanced methods could be tried instead.

```
training <- trainingRaw[, colSums(is.na(trainingRaw)) == 0]
testing <- testingRaw[, colSums(is.na(testingRaw)) == 0]
```

Now let's take another look.

```
dim(training)
```

```
## [1] 19622 93
```

```
dim(testing)
```

```
## [1] 20 60
```

That's better, but there are still too many variables. Let's remove some more columns.

```
y <- training$classe
```

```
training <- training[, sapply(training, is.numeric)]
testing <- testing[, sapply(testing, is.numeric)]
```

```
training <- training %>% select(-X, -raw_timestamp_part_1, -raw_timestamp_part_2, -num_window)
testing <- testing %>% select(-X, -raw_timestamp_part_1, -raw_timestamp_part_2, -num_window, -problem_id)

training$classe <- y
```

```
dim(training)
```

```
## [1] 19622    53
```

```
dim(testing)
```

```
## [1] 20 52
```

Data Preparation

We now want to split our training data into a train and a validation set.

```
set.seed(42)
```

```
inTrain <- createDataPartition(training$classe, p = .7, list = FALSE)
```

```
train <- training[inTrain,]
```

```
val <- training[-inTrain,]
```

Model Fit

Let's now use a random forest classifier for our training data. The random forest classifier was chosen because in the lectures we were told that they perform very good out of the box for most cases. We use 5 fold cross validation and 100 trees.

```
rf <- train(classe ~ ., data = train, method="rf", trControl = trainControl(method = "cv", 5), ntree=100)
```

```
rf
```

```
## Random Forest
```

```
##
```

```
## 13737 samples
```

```
##    52 predictor
```

```
##    5 classes: 'A', 'B', 'C', 'D', 'E'
```

```
##
```

```
## No pre-processing
```

```
## Resampling: Cross-Validated (5 fold)
```

```
## Summary of sample sizes: 10990, 10990, 10990, 10989, 10989
```

```
## Resampling results across tuning parameters:
```

```
##
```

```
##   mtry  Accuracy  Kappa
```

```
##    2    0.9894447 0.9866460
```

```
##   27    0.9913372 0.9890411
```

```
##   52    0.9835480 0.9791847
```

```
##
```

```
## Accuracy was used to select the optimal model using the largest value.
```

```
## The final value used for the model was mtry = 27.
```

On the test set we get a very high accuracy.

```
valPrediction <- predict(rf, val)
cf <- confusionMatrix(valPrediction, as.factor(val$classe))
cf
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1670    3    0    0    0
##           B    3 1133    6    0    0
##           C    1    3 1020    9    0
##           D    0    0    0  954    5
##           E    0    0    0    1 1077
##
## Overall Statistics
##
##           Accuracy : 0.9947
##           95% CI : (0.9925, 0.9964)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9933
##
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9976  0.9947  0.9942  0.9896  0.9954
## Specificity      0.9993  0.9981  0.9973  0.9990  0.9998
## Pos Pred Value   0.9982  0.9921  0.9874  0.9948  0.9991
## Neg Pred Value   0.9991  0.9987  0.9988  0.9980  0.9990
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2838  0.1925  0.1733  0.1621  0.1830
## Detection Prevalence 0.2843  0.1941  0.1755  0.1630  0.1832
## Balanced Accuracy 0.9984  0.9964  0.9957  0.9943  0.9976
```

And we also get a very high accuracy on the validation set.

Now let's look at the out of sample error.

```
outOfSampleError <- 1 - cf$overall[1]
outOfSampleError
```

```
## Accuracy
## 0.00526763
```

Predict Test Data

Finally, we use our classifier to predict the test data.

```
testPrediction <- predict(rf, testing)
```

```
testPrediction
```

```
## [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
```