

Coursera - Practical Machine Learning - Course Assignment

Iñigo Fernández del Amo

1/2/2021

Executive Summary

This report explains a procedure to predict the *classe* of 20 different test cases belonging to the Weight Lifting Exercise data set. Among various ML algorithms trained using k-fold cross-validation ($k = 10$), Random Forests obtained the lowest out-of-sample error rate (0.7). Thus, making it the most accurate algorithm to predict the abovementioned outcomes.

1. Introduction

The Weight Lifting Exercise data set comprises information collected from various sensors (e.g., accelerometers) in wearable devices (e.g., Jawbone Up, Fitbit, etc.) regarding the movements of people when performing Weight Lifting exercises. This data set, and the different variables within it, can be utilized to predict how well wearable devices' users perform certain Weight Lifting exercises.

The Weight Lifting Exercise data set has been presented in “**Qualitative Activity Recognition of Weight Lifting Exercises**” by *Velloso, Bulling, Gellersen, Ugulino and Fuks (2013)*. For more information on it, please visit: http://groupware.les.inf.puc-rio.br/har#weight_lifting_exercises.

For more information on the original code employed to conduct this course assignment, please visit: https://github.com/ifernandezblanc/coursera_dss_practicalmachinelearning.

2. Data preparation

The following are the R libraries utilized for this assignment.

```
library(caret) # For training models and predicting outcomes
library(randomForest) # For applying random forests' algorithms
set.seed(234332) # For reproducibility purposes
```

The training and test data sets can be downloaded directly from the Internet.

```
trainURL <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
training <- read.csv(url(trainURL), na.strings = c("NA", "", "#DIV/0!"))
training$classe <- as.factor(training$classe) # For classe to be treated as factor
dim(training)
```

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

```
testURL <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
testing <- read.csv(url(testURL), na.strings = c("NA", "", "#DIV/0!"))
dim(testing)
```

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

For each data set a total of 160 variables are loaded, with 19622 observations for the training set and 20 observations for the test set. Due to the large number of observations in the training data set, this can be further split for cross-validation purposes (using a 70/30 ratio).

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

A brief exploration to the training data set shows that many of those variables have a large amount of NA values, while some others have very small variances. Besides, the first seven columns are identifiers with no relevance as predictors. All the variables can therefore be removed.

```
allNA <- sapply(trainSet, function(x) mean(is.na(x))) > 0.95 # remove vars with NAs > 95%
trainSet <- trainSet[, allNA == FALSE]
valSet <- valSet[, allNA == FALSE]
nzv <- nearZeroVar(trainSet) # remove variables with near-zero variance
trainSet <- trainSet[, -nzv]
valSet <- valSet[, -nzv]
idVARS <- c(1:7)
trainSet <- trainSet[, -idVARS]
valSet <- valSet[, -idVARS]
dim(trainSet)
```

```
## [1] 13737    52
```

```
dim(valSet)
```

```
## [1] 5885    52
```

These data tidying operations leave a total of 52 variables for *classe* prediction.

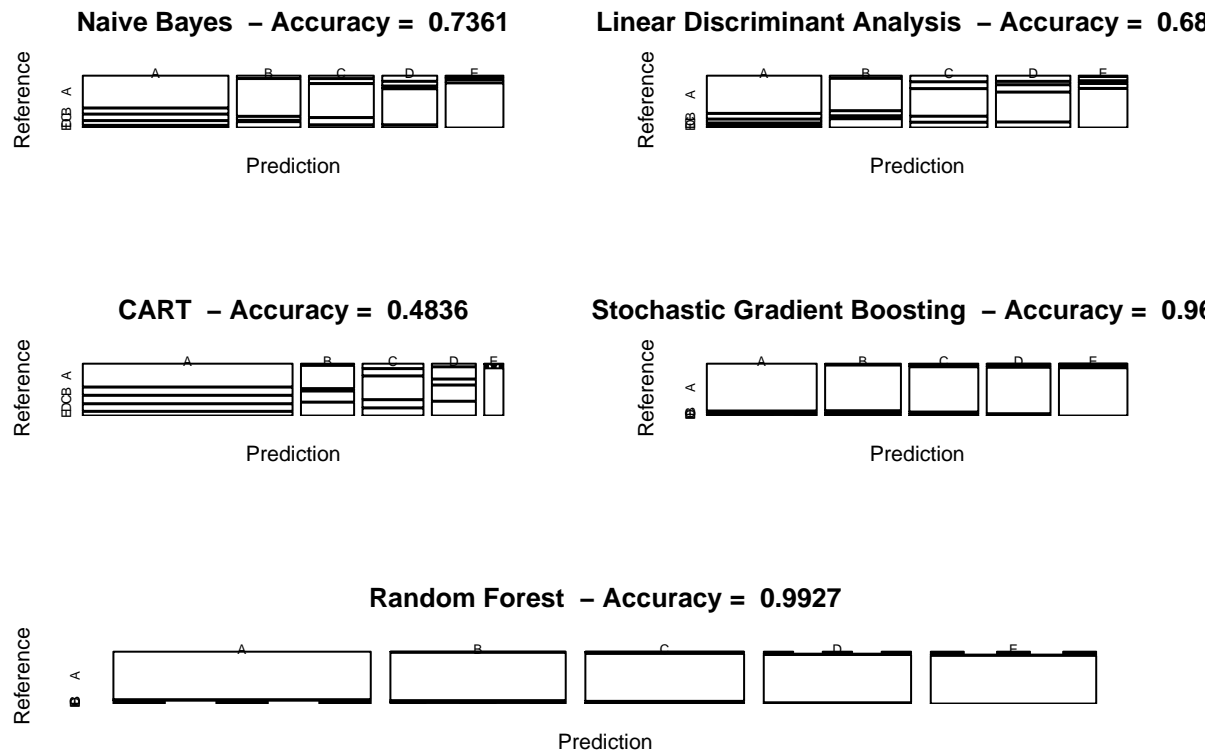
3. Prediction modelling with ML algorithms

Due to the number of variables included in the study, there are several methods that can be applied to obtain predictions on the outcome *classe*: (a) Naive Bayes, (b) Linear discriminant analysis, (c) Decision trees, (d) Random forests and (e) Stochastic Gradient Boosting (generalized boosted models).

The results of training these algorithms with the *trainSet* and evaluating them with the *valSet* are shown below. The coding employed to run those is show in **Appendix A**. The results of models' training and evaluation are fully displayed in **Appendix B** and **Appendix C**, respectively. It is worthy to note that k-fold cross-validation ($k = 10$) was employed to train these models.

```
# List ML algorithms, train them with trainSet and evaluate tthem with valSet
mlMethods <- c("nb", "lda", "rpart", "gbm", "rf")
mlModels <- lapply(mlMethods, mlTraining)
mlResults <- lapply(mlModels, mlPredicting)

# Visualize confusion matrices and accuracies
layout(matrix(c(1,3,5,2,4,5), nrow = 3, ncol = 2))
mlPlots <- lapply(mlResults, mlPlotting)
```



```
# Visualize in-sample and out-of-sample errors
mlErrors <- data.frame(do.call(rbind,(lapply(mlResults, mlEvaluating))))
mlErrors
```

	Algorithm	In_Sample_Error_Rate
## Accuracy	Naive Bayes	0.2595
## Accuracy1	Linear Discriminant Analysis	0.3081
## Accuracy2	CART	0.4901
## Accuracy3	Stochastic Gradient Boosting	0.0399
## Accuracy4	Random Forest	0.0074
##	Out_Of_Sample_Error_Rate	
## Accuracy		0.2639
## Accuracy1		0.3176
## Accuracy2		0.5164
## Accuracy3		0.0399
## Accuracy4		0.0073

These results show the in-sample and out-of-sample error rates obtained by the algorithms for both, the training set (*trainSet*) and the cross-validation set (*valSet*). Based on them, it is possible to say that Random Forests is the algorithm obtaining the lowest out-of-sample error rate (~0.7%). This is closely followed by Stochastic Gradient Boosting (~4%), while Naive Bayes (~26%), Linear Discriminant Analysis (~32%) and Decision Trees - CART (~52%) obtain less accurate results. Therefore, it seems reasonable to use Random Forests to predict the *classe* outcome on the *testing* data set.

4. Data prediction

Based on previous results, the Random Forests' model is utilized to predict the 20-quiz results (*testing* data set).

```
predict(mlModels[[5]], newdata = testing)
```

```
## [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
```

5. Conclusions

This report shows the application of different ML algorithms to predict the *classe* outcome on the Weight Lifting Exercise data set. Due to the large number of NA values in some variables, and the low variances in others, only 52 variables out of 160 were used to predict the outcome. Besides, different ML algorithms (Naive Bayes, LDA, CART, GBM and RF) were trained (using 10-fold cross-validation) and evaluated with cross-validation data sets (*valSet*) to compare them according to their out-of-sample error rates ($1 - Accuracy$). Random Forests (RF) was the algorithm with the lowest out-of-sample error (0.7) and so, it was utilized to predict the results on the *testing* data set.

Appendices

Appendix A

Functions to train ML algorithms and evaluate and plot confusion matrices

```
# Functions to train, evaluate and plot each ML algorithm sequentially
# Training control parameters
control <- trainControl(  preProcOptions = list(thresh = 0.8),
                          allowParallel=T,
                          savePredictions=T,
                          method = "cv",
                          number = 10)

# Training function for trainSet
mlTraining <- function (mlMethod) {
  modelFit <- train(classe ~ ., data = trainSet,
                   method = mlMethod, trControl = control)
  return(modelFit)
}

# Confusion matrix for valSet
mlPredicting <- function(modelFit) {
  modelPrediction <- predict(modelFit, newdata = valSet)
  modelCM <- confusionMatrix(modelPrediction, valSet$classe)
  modelCM$methodName <- modelFit$modelInfo$label
  modelCM$methodAccuracy <- max(modelFit$results$Accuracy)
  return(modelCM)
}

# Plot on confusion matrix results
mlPlotting <- function(modelResult) {
  mlPlot <- plot(modelResult$table, col = modelResult$byClass,
                main = paste(modelResult$methodName, " - Accuracy = ",
                             round(modelResult$overall["Accuracy"], 4)))
  return(mlPlot)
}

# Calculate errors
mlEvaluating <- function(modelResult) {
  return(data.frame(Algorithm = modelResult$methodName,
                    In_Sample_Error_Rate = round(1 - modelResult$methodAccuracy, 4),
                    Out_Of_Sample_Error_Rate = round(1 - modelResult$overall["Accuracy"], 4)))
}
```

Appendix B

Results of ML algorithms training

mlModels

```
## [[1]]
## Naive Bayes
##
## 13737 samples
##    51 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 12363, 12363, 12364, 12362, 12363, 12363, ...
## Resampling results across tuning parameters:
##
##    usekernel  Accuracy  Kappa
##    FALSE      0.4989386  0.3816917
##    TRUE       0.7404823  0.6683833
##
## Tuning parameter 'fL' was held constant at a value of 0
## Tuning
## parameter 'adjust' was held constant at a value of 1
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were fL = 0, usekernel = TRUE and adjust
## = 1.
##
## [[2]]
## Linear Discriminant Analysis
##
## 13737 samples
##    51 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 12363, 12364, 12363, 12362, 12364, 12362, ...
## Resampling results:
##
##    Accuracy  Kappa
##    0.6919202  0.610044
##
##
## [[3]]
## CART
##
## 13737 samples
##    51 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 12362, 12364, 12363, 12363, 12364, 12365, ...
```

```

## Resampling results across tuning parameters:
##
##      cp          Accuracy    Kappa
##  0.02970196  0.5099395  0.3620551
##  0.03692402  0.4744135  0.3147172
##  0.06410843  0.4275130  0.2373038
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was cp = 0.02970196.
##
## [[4]]
## Stochastic Gradient Boosting
##
## 13737 samples
##    51 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 12364, 12363, 12364, 12363, 12364, 12363, ...
## Resampling results across tuning parameters:
##
##  interaction.depth  n.trees  Accuracy    Kappa
##      1              50      0.7407742  0.6712373
##      1             100      0.8128419  0.7631084
##      1             150      0.8461102  0.8052608
##      2              50      0.8548449  0.8160856
##      2             100      0.9047836  0.8794838
##      2             150      0.9280784  0.9089841
##      3              50      0.8926267  0.8640619
##      3             100      0.9396533  0.9236362
##      3             150      0.9601086  0.9495307
##
## 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.
##
## [[5]]
## Random Forest
##
## 13737 samples
##    51 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 12362, 12365, 12364, 12364, 12364, 12364, ...
## Resampling results across tuning parameters:
##
##      mtry  Accuracy    Kappa
##      2    0.9911187  0.9887642

```

```
## 26 0.9925754 0.9906071
## 51 0.9879161 0.9847124
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 26.
```

Appendix C

Results of ML algorithms evaluation

mlResults

```
## [[1]]
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  A    B    C    D    E
##           A 1477  236  252  173   77
##           B   48  746   60    5  109
##           C   50   81  682  129   45
##           D   90   65   29  614   38
##           E    9   11    3   43  813
##
## Overall Statistics
##
##           Accuracy : 0.7361
##           95% CI : (0.7246, 0.7473)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.6622
##
## Mcnemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.8823  0.6550  0.6647  0.6369  0.7514
## Specificity      0.8247  0.9532  0.9372  0.9549  0.9863
## Pos Pred Value   0.6668  0.7707  0.6910  0.7344  0.9249
## Neg Pred Value   0.9463  0.9201  0.9298  0.9307  0.9463
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2510  0.1268  0.1159  0.1043  0.1381
## Detection Prevalence 0.3764  0.1645  0.1677  0.1421  0.1494
## Balanced Accuracy 0.8535  0.8041  0.8010  0.7959  0.8688
##
## [[2]]
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  A    B    C    D    E
##           A 1359  170  104   52   57
##           B   49  720   93   44  190
##           C  138  143  654  122  118
##           D  123   56  149  685  119
```

```

##          E      5      50      26      61      598
##
## Overall Statistics
##
##          Accuracy : 0.6824
##          95% CI : (0.6703, 0.6943)
##          No Information Rate : 0.2845
##          P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.5982
##
## McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##          Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.8118  0.6321  0.6374  0.7106  0.5527
## Specificity      0.9090  0.9208  0.8928  0.9092  0.9704
## Pos Pred Value   0.7801  0.6569  0.5566  0.6051  0.8081
## Neg Pred Value   0.9240  0.9125  0.9210  0.9413  0.9059
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2309  0.1223  0.1111  0.1164  0.1016
## Detection Prevalence 0.2960  0.1862  0.1997  0.1924  0.1257
## Balanced Accuracy 0.8604  0.7765  0.7651  0.8099  0.7616
##
## [[3]]
## Confusion Matrix and Statistics
##
##          Reference
## Prediction      A      B      C      D      E
##          A 1528  484  493  440  245
##          B   20  376   23  171  218
##          C   86  122  441  140  138
##          D   37  157   69  213  193
##          E    3    0    0    0  288
##
## Overall Statistics
##
##          Accuracy : 0.4836
##          95% CI : (0.4708, 0.4965)
##          No Information Rate : 0.2845
##          P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.3241
##
## McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##          Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9128  0.33011  0.42982  0.22095  0.26617
## Specificity      0.6053  0.90898  0.89998  0.90734  0.99938
## Pos Pred Value   0.4790  0.46535  0.47573  0.31839  0.98969
## Neg Pred Value   0.9458  0.84971  0.88201  0.85602  0.85806

```



```

## Prevalence          0.2845  0.19354  0.17434  0.16381  0.18386
## Detection Rate      0.2596  0.06389  0.07494  0.03619  0.04894
## Detection Prevalence 0.5421  0.13730  0.15752  0.11368  0.04945
## Balanced Accuracy   0.7591  0.61955  0.66490  0.56415  0.63277
##
## [[4]]
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1640   30    0    1    0
##           B   20 1080   38    6    9
##           C    8   26  973   28   14
##           D    4    2   14  924   26
##           E    2    1    1    5 1033
##
## Overall Statistics
##
##           Accuracy : 0.9601
##           95% CI : (0.9547, 0.9649)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9495
##
## Mcnemar's Test P-Value : 3.736e-08
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9797  0.9482  0.9483  0.9585  0.9547
## Specificity      0.9926  0.9846  0.9844  0.9907  0.9981
## Pos Pred Value   0.9814  0.9367  0.9276  0.9526  0.9914
## Neg Pred Value   0.9919  0.9875  0.9890  0.9919  0.9899
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2787  0.1835  0.1653  0.1570  0.1755
## Detection Prevalence 0.2839  0.1959  0.1782  0.1648  0.1771
## Balanced Accuracy 0.9862  0.9664  0.9664  0.9746  0.9764
##
## [[5]]
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1668    7    0    0    0
##           B    3 1131    6    1    0
##           C    3    1 1016   11    2
##           D    0    0    4  951    4
##           E    0    0    0    1 1076
##
## Overall Statistics
##
##           Accuracy : 0.9927
##           95% CI : (0.9902, 0.9947)

```

```

##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9908
##
##  McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9964   0.9930   0.9903   0.9865   0.9945
## Specificity          0.9983   0.9979   0.9965   0.9984   0.9998
## Pos Pred Value       0.9958   0.9912   0.9835   0.9917   0.9991
## Neg Pred Value       0.9986   0.9983   0.9979   0.9974   0.9988
## Prevalence           0.2845   0.1935   0.1743   0.1638   0.1839
## Detection Rate       0.2834   0.1922   0.1726   0.1616   0.1828
## Detection Prevalence 0.2846   0.1939   0.1755   0.1630   0.1830
## Balanced Accuracy     0.9974   0.9954   0.9934   0.9924   0.9971

```