

# Prediction Assignment Writeup

*Victor*

*February 14, 2016*

## Practical Machine Learning Assignment Writeup

### Abstract

In this assignment, I built a predictive model to determine whether a particular form of exercise (barbell lifting) is performed correctly, using accelerometer data.

### Data Retrieval

The dataset is downloaded from the Internet:

```
if (! file.exists('./pml-training.csv')) {  
  download.file('http://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv',  
    destfile = './pml-training.csv')  
}  
if (! file.exists('./pml-testing.csv')) {  
  download.file('http://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv',  
    destfile = './pml-testing.csv')  
}
```

Loading CSV data:

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

### Exploratory Analysis

The training data has 19622 observations and 160 features:

```
dim(training)
```

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

Inspection of the data set indicates that many of the 159 predictors are missing in most of the observations:

```
sum(complete.cases(training))
```

```
## [1] 406
```

I've divided test into 80%/20% for training/validation set. The I cleaned the data set by extracting necessary features:

```
library(caret)

## Loading required package: lattice

## Loading required package: ggplot2

trainset <- createDataPartition(training$classe, p = 0.8, list = FALSE)
Training <- training[trainset, ]
Validation <- training[-trainset, ]
nzvcol <- nearZeroVar(Training)
Training <- Training[, -nzvcol]
cntlength <- sapply(Training, function(x) {
  sum(!(is.na(x) | x == ""))
})
nullcol <- names(cntlength[cntlength < 0.6 * length(Training$classe)])
descriptcol <- c("X", "user_name", "raw_timestamp_part_1", "raw_timestamp_part_2",
  "cvtd_timestamp", "new_window", "num_window")
excludecols <- c(descriptcol, nullcol)
Training <- Training[, !names(Training) %in% excludecols]
```

## Predictive Model

For my initial attempt at building a predictive model I chose the random forest algorithm. Random forests have several nice theoretical properties:

1. They deal naturally with non-linearity, and assuming linearity in this case would be imprudent.
2. There's no parameter selection involved. While random forest may overfit a given data set, just as any other machine learning algorithm, it has been shown by Breiman that classifier variance does not grow with the number of trees used (unlike with Adaboosted decision trees, for example). Therefore, it's always better to use more trees, memory and computational power allowing.
3. The algorithm allows for good in-training estimates of variable importance and generalization error, which largely eliminates the need for a separate validation stage, though obtaining a proper generalization error estimate on a testing set would still be prudent.
4. The algorithm is generally robust to outliers and correlated covariates, which seems like a nice property to have when there are known interactions between variables and no data on presence of outliers in the data set.

Given that the problem at hand is a high-dimensional classification problem with number of observations much exceeding the number of predictors, random forest seems like a sound choice.

Let's train a classifier without cross validation. CV tests on such big datasets cause increase in training time significantly and actually unnecessary because of great results without cross validation.

```
set.seed(555)
library(randomForest)

## randomForest 4.6-12

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

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

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

rfModel <- randomForest(classe ~ ., data = Training, importance = TRUE, ntrees = 5)
```

### Training set accuracy (In-Sample)

I expect error rate to be less than 10%. Let's see our in-sample results:

```
ptraining <- predict(rfModel, Training)
print(confusionMatrix(ptraining, Training$classe))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 4464    0    0    0    0
##           B    0 3038    0    0    0
##           C    0    0 2738    0    0
##           D    0    0    0 2573    0
##           E    0    0    0    0 2886
##
## Overall Statistics
##
##           Accuracy : 1
##           95% CI : (0.9998, 1)
##           No Information Rate : 0.2843
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 1
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity         1.0000   1.0000   1.0000   1.0000   1.0000
## Specificity         1.0000   1.0000   1.0000   1.0000   1.0000
## Pos Pred Value      1.0000   1.0000   1.0000   1.0000   1.0000
## Neg Pred Value      1.0000   1.0000   1.0000   1.0000   1.0000
## Prevalence          0.2843   0.1935   0.1744   0.1639   0.1838
## Detection Rate      0.2843   0.1935   0.1744   0.1639   0.1838
## Detection Prevalence 0.2843   0.1935   0.1744   0.1639   0.1838
## Balanced Accuracy    1.0000   1.0000   1.0000   1.0000   1.0000
```

### Validation set accuracy (Out-of-Sample)

Let's see our out-of-sample results:

```
pvalidation <- predict(rfModel, Validation)
print(confusionMatrix(pvalidation, Validation$classe))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1115    6    0    0    0
##           B    0  752    3    0    0
##           C    1    1  681    4    0
##           D    0    0    0  639    0
##           E    0    0    0    0  721
##
## Overall Statistics
##
##           Accuracy : 0.9962
##           95% CI : (0.9937, 0.9979)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9952
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9991  0.9908  0.9956  0.9938  1.0000
## Specificity      0.9979  0.9991  0.9981  1.0000  1.0000
## Pos Pred Value   0.9946  0.9960  0.9913  1.0000  1.0000
## Neg Pred Value   0.9996  0.9978  0.9991  0.9988  1.0000
## Prevalence       0.2845  0.1935  0.1744  0.1639  0.1838
## Detection Rate   0.2842  0.1917  0.1736  0.1629  0.1838
## Detection Prevalence 0.2858  0.1925  0.1751  0.1629  0.1838
## Balanced Accuracy 0.9985  0.9949  0.9969  0.9969  1.0000
```

## Test Set Prediction

Prediction of the test set:

```
pctest <- predict(rfModel, testing)
print(confusionMatrix(pctest, testing$classe))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 5579    6    0    0    0
##           B    0 3790    3    0    0
##           C    1    1 3419    4    0
##           D    0    0    0 3212    0
##           E    0    0    0    0 3607
##
## Overall Statistics
```

```

##
##           Accuracy : 0.9992
##           95% CI   : (0.9987, 0.9996)
##    No Information Rate : 0.2844
##    P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.999
##  McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9998  0.9982  0.9991  0.9988  1.0000
## Specificity      0.9996  0.9998  0.9996  1.0000  1.0000
## Pos Pred Value   0.9989  0.9992  0.9982  1.0000  1.0000
## Neg Pred Value   0.9999  0.9996  0.9998  0.9998  1.0000
## Prevalence       0.2844  0.1935  0.1744  0.1639  0.1838
## Detection Rate   0.2843  0.1932  0.1742  0.1637  0.1838
## Detection Prevalence 0.2846  0.1933  0.1745  0.1637  0.1838
## Balanced Accuracy 0.9997  0.9990  0.9994  0.9994  1.0000

```

I obtained the error rate less than 0.001% on the test set. That's very good result.

## Conclusion

Given that the model obtained using the initial approach appears to be highly successful, further exploration of the matter does not seem to be necessary.