

Machine Learning Writeup, Data Science Specialization

Charles Carter*

November 17, 2014

Contents

1	Introduction	1
2	Instructions	2
2.1	Background	2
2.2	Data	2
2.3	What you should submit	2
2.4	Reproducibility	3
3	User Defined Functions	3
3.1	Loading Packages	3
3.2	Loading Data Files	3
3.3	Delete Empty Columns	4
4	Logic	4
4.1	Calling Helper Functions	4
4.2	Validation of Results	5
4.3	Prediction	6
4.4	Cross Validation and Out of Sample Error	7
5	Results and Conclusion	7

1 Introduction

This paper is the writeup of the course project in Practical Machine Learning, which is part of the Data Science specialization offered by Johns Hopkins University through Coursera. The project requires the predictive analysis of a data set, and the evaluation of various models. This paper consists of an (1) Introduction, the (2) Instructions for completing the project, a description of

*ccarter@troy.edu

the (3) User Defined Functions called by the main logic of the script, the (4) Logic of the script, and the (5) Results and Conclusions.

2 Instructions

2.1 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, your goal will be to use data from accelerometers on the belt, forearm, arm, and dumbbell 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://groupware.les.inf.puc-rio.br/har> (see the section on the Weight Lifting Exercise Dataset).

2.2 Data

The training data for this project are available here: <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>

The test data are available here: <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>

The data for this project come from this source: <http://groupware.les.inf.puc-rio.br/har>. If you use the document you create for this class for any purpose please cite them as they have been very generous in allowing their data to be used for this kind of assignment.

2.3 What you should submit

The goal of your project is to predict the manner in which they did the 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.

1. Your submission should consist of a link to a Github repo with your R markdown and compiled HTML file¹ describing your analysis. Please constrain the text of the writeup to less than 2000 words and the number of figures to be

¹I have not submitted MarkDown or HTML files, but Rnw and PDF files. The reason is that I almost never use HTML output in my job, but frequently generate PDF files, so I wanted to focus on the production of output in PDF format.

less than 5. It will make it easier for the graders if you submit a repo with a gh-pages branch so the HTML page can be viewed online (and you always want to make it easy on graders.)

2. You should also apply your machine learning algorithm to the 20 test cases available in the test data above. Please submit your predictions in appropriate format to the programming assignment for automated grading. See the programming assignment for additional details.

2.4 Reproducibility

Due to security concerns with the exchange of R code, your code will not be run during the evaluation by your classmates. Please be sure that if they download the repo, they will be able to view the compiled HTML version of your analysis.

3 User Defined Functions

3.1 Loading Packages

My analysis uses a number of helper functions. This sections details these functions. The first order of business is to load the required packages. We need `caret` to perform the prediction and training, `knitr` to prepare the PDF, and `ggplot2` for visualization.

```
load.libraries <- function()
{
  library(caret)
  library(knitr)
  library(ggplot2)
  library(randomForest)
}
```

3.2 Loading Data Files

Next, we load and clean the data. Function `load.data()` assumes that the data file is in the current directory. This function loads the data into the R environment.

```
load.data <- function(filename)
{
  file <- read.csv(filename, stringsAsFactors = TRUE, na.strings=c("NA", ""))
  cat(filename, " set loaded, returns\n")
  return(file)
}
```

3.3 Delete Empty Columns

The data files contain many variables with no values, or very sparse values. The function `delete.empty.cols.1` reduces the datasets, both training and testing, by deleting those variables. The numeral *1* appended to the end of the function results from previous, exploratory versions.

```
delete.empty.cols.1 <- function(df)
{
  x <- df[, apply(df, 2, function(x)
    (sum(grepl("[A-Za-z0-9]", x, perl = TRUE)) / length(x)) > 0.5)]
  delete <- grep("X|timestamp|window", names(df), perl = TRUE)
  x <- x[, -delete]
  print(" set cleaned and prepared, returns\n")
  return(x)
}
```

4 Logic

4.1 Calling Helper Functions

After creating the helper functions, I call them with the following listing. `train` consists of the original, raw data. `trainA` consists of the data after cleaning. This results in a model fit for the testing

```
##this is the program logic
print("loading libraries\n")

## [1] "loading libraries\n"

load.libraries()
print("calling train <- load.data.training('pml-training.csv')\n")

## [1] "calling train <- load.data.training('pml-training.csv')\n"

train <- load.data('pml-training.csv')

## pml-training.csv set loaded, returns

print("calling delete.empty.cols.1(), returns trainA\n")

## [1] "calling delete.empty.cols.1(), returns trainA\n"

trainA <- delete.empty.cols.1(train)

## [1] " set cleaned and prepared, returns\n"
```

```

print("creating training and testing sets, returns 'training' and'testing'")

## [1] "creating training and testing sets, returns 'training' and'testing'"

inTrain <- createDataPartition(y = trainA$classe, p = 0.7, list = FALSE)
training <- trainA[inTrain, ]
testing <- trainA[-inTrain, ]
cat("Dimensions of training data are: ", dim(training), "\n and dimensions of the testing da

## Dimensions of training data are: 13737 54
## and dimensions of the testing data are: 5885 54 .

cat("setting seed and calling train() with GLM\n")

## setting seed and calling train() with GLM

set.seed(141017)
suppressWarnings( modelFit <- train(classe~ ., data = training, method = "glm") )

## Error: final tuning parameters could not be determined

prediction <- predict(modelFit, newdata = testing)

```

4.2 Validation of Results

We validate the results of our prediction by using the test data as the input for `predict()`, and then constructing a confusion matrix with the actual outcomes in the test data and the predicted results. The accuracy of this exercise is 99.5 percent. As shown below, both the positive predictive values and the negative predictive values are 99 percent.

```

prediction <- predict(modelFit, newdata = testing)
cm <- confusionMatrix(prediction, testing$classe)
cm

## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 1673     6    0    0    0
##      B     0 1132     2    0    1
##      C     1     1 1022    10    3
##      D     0     0     2   954    5
##      E     0     0     0     0 1073
##
## Overall Statistics

```

```
##
##           Accuracy : 0.995
##           95% CI : (0.993, 0.996)
##      No Information Rate : 0.284
##      P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.993
##  McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.999   0.994   0.996   0.990   0.992
## Specificity          0.999   0.999   0.997   0.999   1.000
## Pos Pred Value       0.996   0.997   0.986   0.993   1.000
## Neg Pred Value       1.000   0.999   0.999   0.998   0.998
## Prevalence           0.284   0.194   0.174   0.164   0.184
## Detection Rate       0.284   0.192   0.174   0.162   0.182
## Detection Prevalence 0.285   0.193   0.176   0.163   0.182
## Balanced Accuracy     0.999   0.997   0.997   0.994   0.996
```

4.3 Prediction

Finally, we run a prediction on the testing data set and observe the results. `predictions` contains the predictions of the testing set from the model trained in the training set. `result` contains the numerical results of the predictions. We then plot the results, as contained in the testing set, with the predictions, as shown in the plot below.

```
print("calling test <- load.data.training('pml-testing.csv')\n")
## [1] "calling test <- load.data.training('pml-testing.csv')\n"
test <- load.data('pml-testing.csv')
## pml-testing.csv set loaded, returns
print("calling delete.empty.cols.1(), returns testA\n")
## [1] "calling delete.empty.cols.1(), returns testA\n"
testA <- delete.empty.cols.1(test)
## [1] " set cleaned and prepared, returns\n"
outcome.predictions <- predict(modelFit, newdata = outcome1)
```

```
## Error: object 'user_name' not found

outcome.predictions

## [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
```

4.4 Cross Validation and Out of Sample Error

The initial exploratory work was done with the training and testing sets as indicated above. The run time of this code takes approximate 72 seconds on my machine. In order to determine whether the first results I obtained were reasonable, I altered the seed and ran the code ten times (with different seeds). The results did not vary with the example shown above. Therefore, I have omitted these graphs and the code I used.

I also experimented with methods other than `glm`. The function `names(getModelInfo())` lists 169 different methods available with the `train()` function. Other than `glm`, I ran the same code with the methods listed below. Some of the methods ran an inordinately long time — Random Forests ran for about six hours. The plotted results (which was all I checked) were all very similar, so I am comfortable that `glm` works reasonably well as a predictor. I have omitted the R code, the results, and the plots in the interest of not overtaxing the reader. If the reader chooses, he can substitute the lines of code below appropriately and obtain the same results.

```
1 modelFit <- train(classeNum ~ ., data = training, method = "bag")
2 modelFit <- train(classeNum ~ ., data = training, method = "cforest"
  ")
3 modelFit <- train(classeNum ~ ., data = training, method = "lda")
4 modelFit <- train(classeNum ~ ., data = training, method = "dnn")
5 modelFit <- train(classeNum ~ ., data = training, method = "logreg"
  )
6 modelFit <- train(classeNum ~ ., data = training, method = "
  svmLinear")
```

Based on the foregoing, I conclude that the minimum out of sample error is 1 in 4904 records, or approximately 0.02 percent.

5 Results and Conclusion

Using package `caret` enabled a fairly quick and easy way to test various prediction algorithms. Method `glm` precisely predicted the outcomes of the testing data from the training data.

Finally, as conclusive proof of the validity of the results, I uploaded the predictions in Part 2 of the assignment, the 20 test cases contained in the `pml-testing.csv` file, and each of the predictions for the observations were correct.