# Peer-graded Assignment: Learning Machine Course Project

*Christian Frei*

*21 June 2017*

# Executive Summary

## 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 dumbell 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 (http://groupware.les.inf.puc-rio.br/har) (see the section on the Weight Lifting Exercise Dataset).

## Data

The training data for this project are available here:

https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv (https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv)

The test data are available here:

https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv (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 (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.

## Goal

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.

# Download and Load Data

In the beginning of this project, data will be downloaded and load into the memory.

```r
dest.subdirectory <- "./data/"
dest.filename     <- c("pml-training.csv", "pml-testing.csv")
dest.filepath     <- paste0(dest.subdirectory, dest.filename)
source.fileURL    <- c("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-tr
aining.csv",

                          "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-te
sting.csv")

if (!file.exists("data")) {
    dir.create("data")
}

for (download in 1:2) {
    if (!file.exists(dest.filepath[download])){
        download.file(source.fileURL[download], dest.filepath[download], method="cur
l")
    }
}

pml.training.csv <- read.csv(dest.filepath[1], header=TRUE, sep=",", na.strings=c(
"NA","#DIV/0!"))
pml.testing.csv <- read.csv(dest.filepath[2], header=TRUE, sep=",", na.strings=c("
NA","#DIV/0!"))
```

```

# Cleaning the data

Next, data have to be explored regarding its structure.

```r
str(pml.training.csv)

str(pml.testing.csv)
```

At the first glance, the first column is just the row number and some predictors have got a lot of NAs or 0 values. Therefore, those columns have to been eliminated. Furthermore, some rows consist of more than 70% NAs. Those rows will be removed.

```
pml.training.csv  <- pml.training.csv[,c(-1)]

nzvCol            <- nearZeroVar(pml.training.csv)
pml.training.csv  <- pml.training.csv[, -nzvCol]

uselessColumns    <- c()

for (i in 1:length(pml.training.csv)) {
    if (sum(is.na(pml.training.csv[ , i])) / nrow(pml.training.csv) >= .70) {
        uselessColumns <- rbind(uselessColumns, i)
    }
}

pml.training.csv  <- pml.training.csv[,-uselessColumns]
```

After cleaning training data, validation data (pml.testing.csv) must be brought into the same shape. Therefore, only those columns will be taken over which are included in the training data (pml.training.csv).

```
usedColumns <- colnames(pml.training.csv)
pml.testing.csv   <- pml.testing.csv[, usedColumns[1:length(usedColumns)-1]]
```

# Data Partitioning

In this step, the training data will be partitioned into training and testing data.

```
inTrain    <- createDataPartition(pml.training.csv$classe, p=0.6, list=FALSE)
training   <- pml.training.csv[inTrain, ]
testing    <- pml.training.csv[-inTrain, ]
```

The function read.csv loaded data of pml-training.csv & pml-testing.csv. Unfortunately, it often identifies different class for those columns, which are available in both of the loaded data. Therefore, all data of pml-testing.csv have to be converted according to the types of pml-training.csv.

```
valdiation  <- rbind(training[1, 1:ncol(training)-1], pml.testing.csv)
valdiation  <- valdiation[2:nrow(valdiation), ]
row.names(valdiation) <- 1:nrow(valdiation)
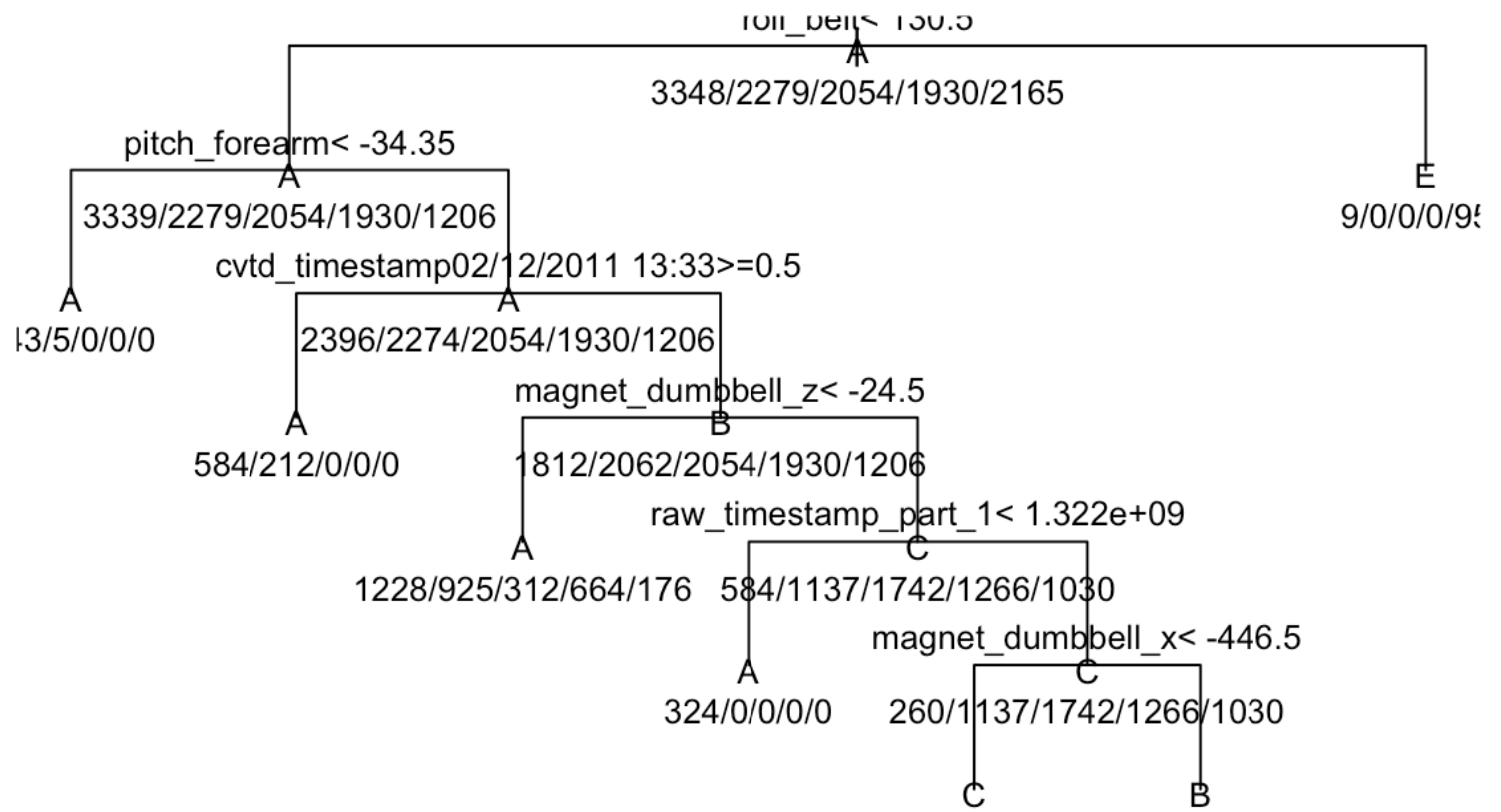```

# Prediction Model 1: Decision Tree

The first prediction model which will be calculated, is the decision tree.

```
set.seed(12345)

model.rpart <- train(classe ~ ., data=training, method="rpart")

# fancyRpartPlot(model.rpart1)
plot(model.rpart$finalModel, uniform=TRUE, main="Classification Tree")
text(model.rpart$finalModel, use.n=TRUE, all=TRUE, cex=.8)
```

# Classification Tree

roll_belt< 130.5

3348/2279/2054/1930/2165

pitch_forearm< -34.35

3339/2279/2054/1930/1206

E

9/0/0/0/95

A

43/5/0/0/0

cvtd_timestamp02/12/2011 13:33>=0.5

2396/2274/2054/1930/1206

A

584/212/0/0/0

magnet_dumbbell_z< -24.5

1812/2062/2054/1930/1206

B

A

1228/925/312/664/176

raw_timestamp_part_1< 1.322e+09

584/1137/1742/1266/1030

C

A

324/0/0/0/0

magnet_dumbbell_x< -446.5

260/1137/1742/1266/1030

C

C

B

```
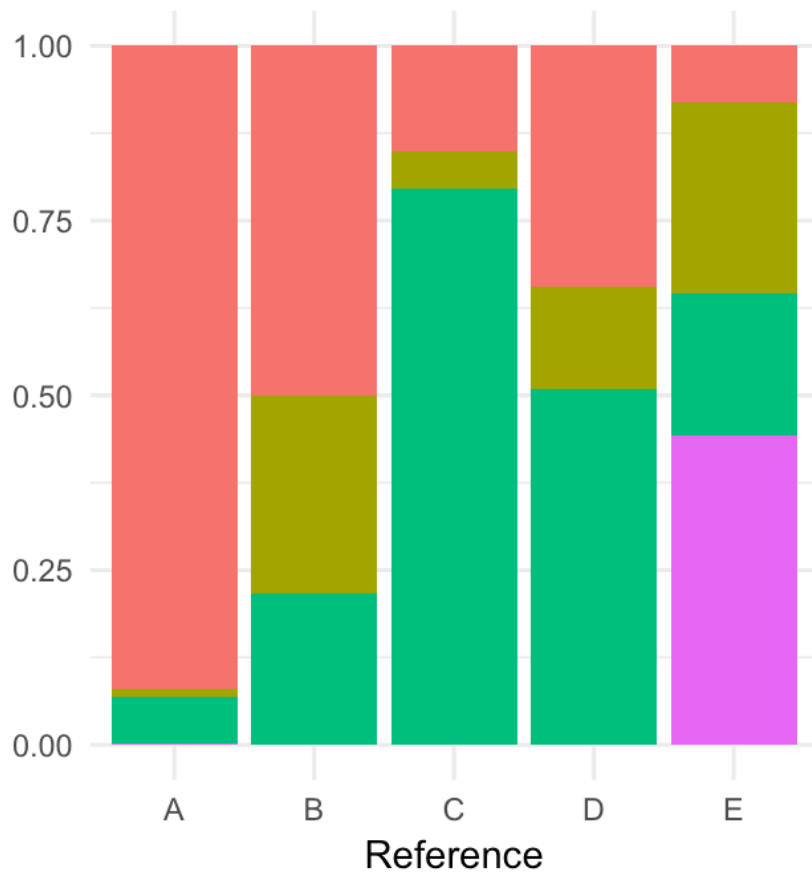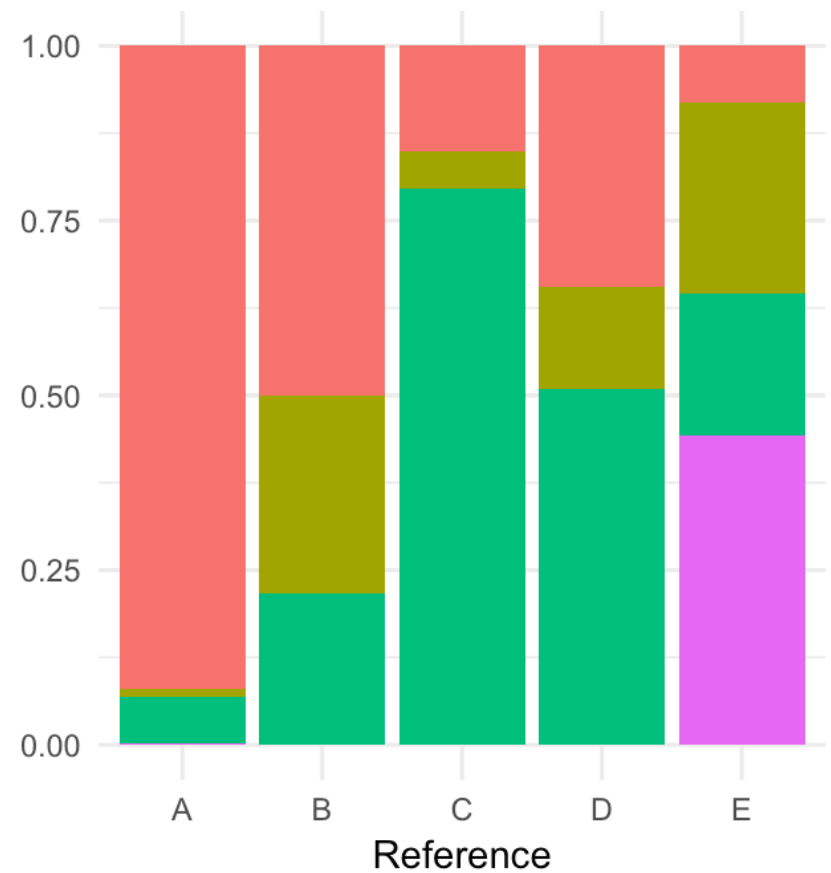prediction.training.rpart <- predict(model.rpart, newdata = training, method="class")
cm.training.rpart <- confusionMatrix(prediction.training.rpart, training$classe)

prediction.testing.rpart <- predict(model.rpart, newdata = testing, method="class")
cm.testing.rpart <- confusionMatrix(prediction.testing.rpart, testing$classe)
```

Confusion Matrix
Decision Tree (Training)
Accuracy = 0.5363

Confusion Matrix
Decision Tree (Testing)
Accuracy = 0.5382

# Prediction Model 2: Random Forest

The second prediction model which will be calculated, is the random forest.

```
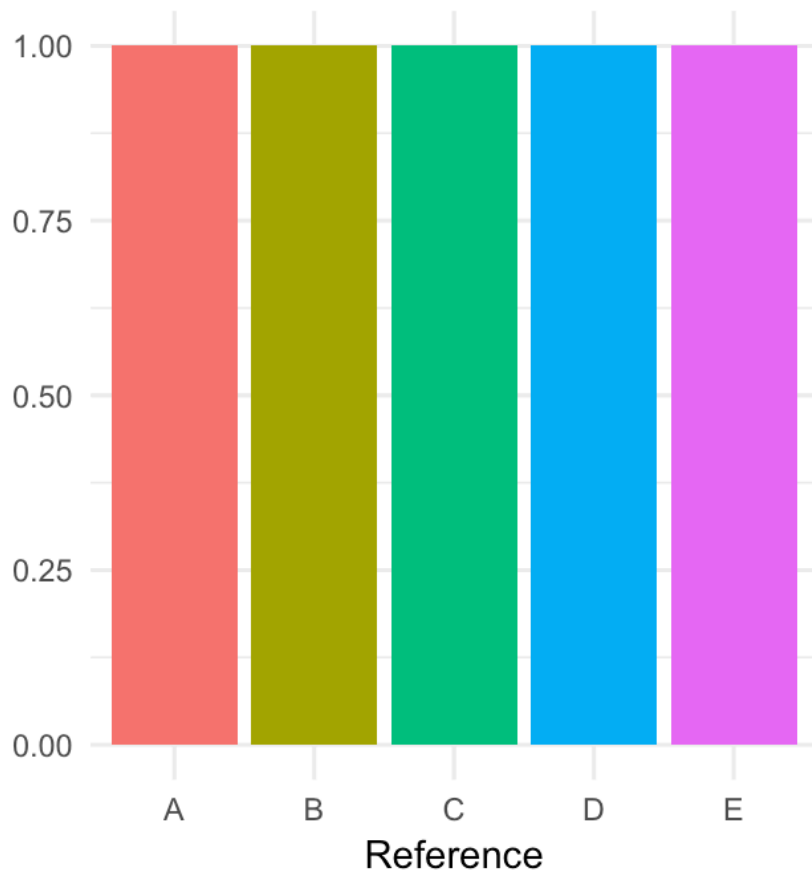set.seed(12345)

# model.rf <- train(classe ~ ., data=training, method="rf")
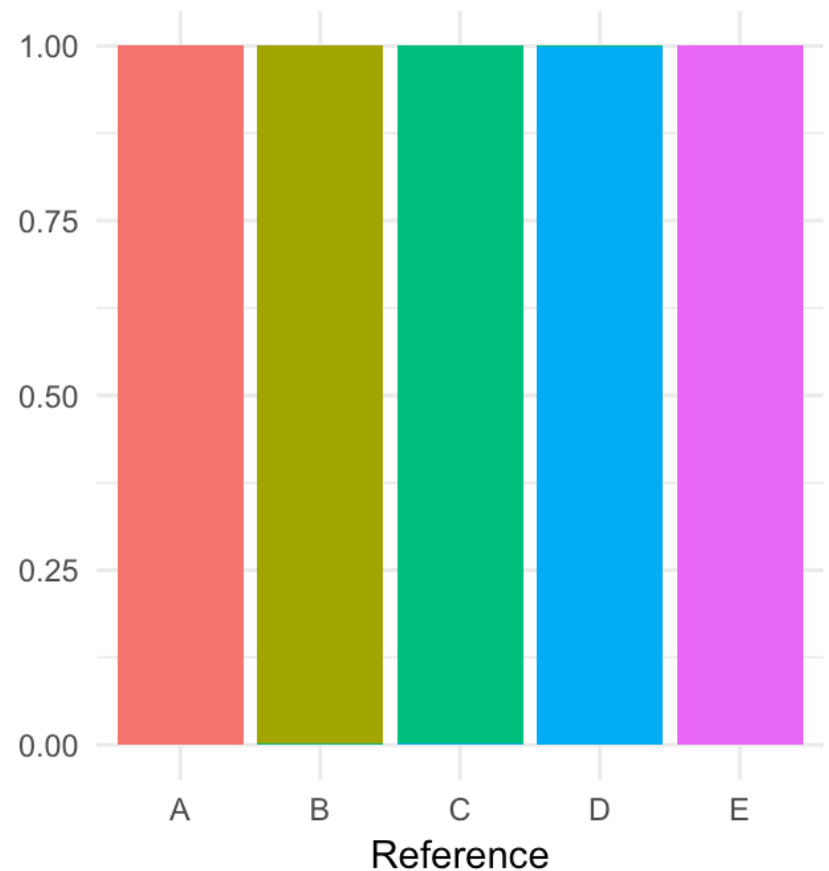model.rf <- randomForest(classe ~ ., data=training)

prediction.training.rf <- predict(model.rf, newdata = training)
cm.training.rf <- confusionMatrix(prediction.training.rf, training$classe)

prediction.testing.rf <- predict(model.rf, newdata = testing)
cm.testing.rf <- confusionMatrix(prediction.testing.rf, testing$classe)
```

Confusion Matrix
Random Forest (Training)
Accuracy = 1

Confusion Matrix
Random Forest (Testing)
Accuracy = 0.9992

# Prediction Model 3: Generalized Boosted Regression

The last prediction model which will be calculated, is the random forest.

```
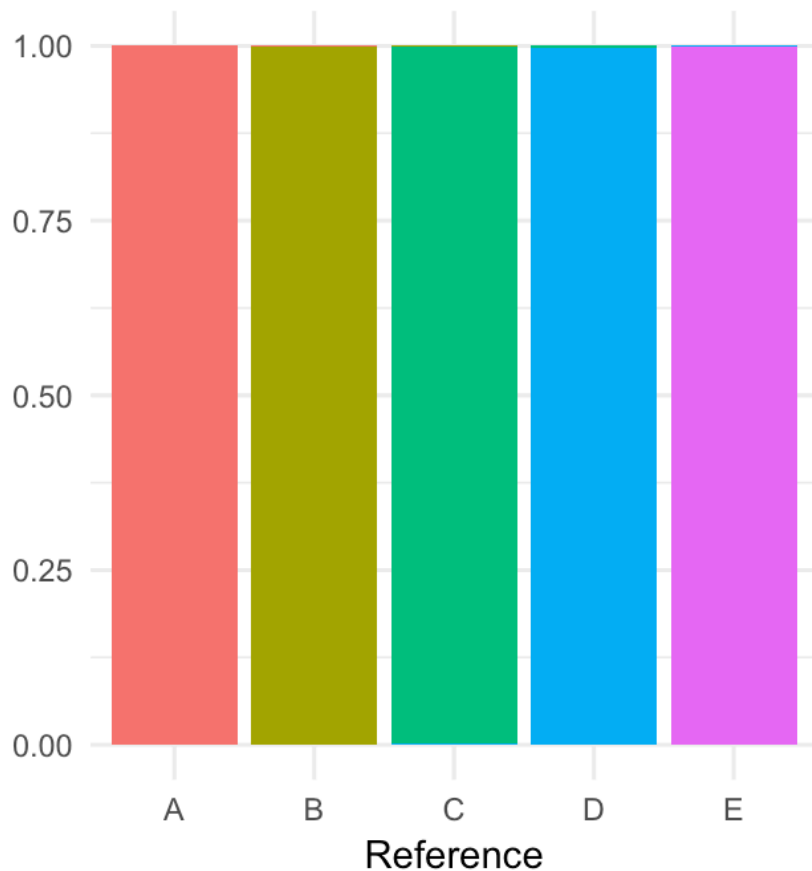set.seed(12345)
fitControl <- trainControl(method = "repeatedcv",
                           number = 5,
                           repeats = 1)

model.gbm <- train(classe ~ ., data=training, method = "gbm",
                 trControl = fitControl,
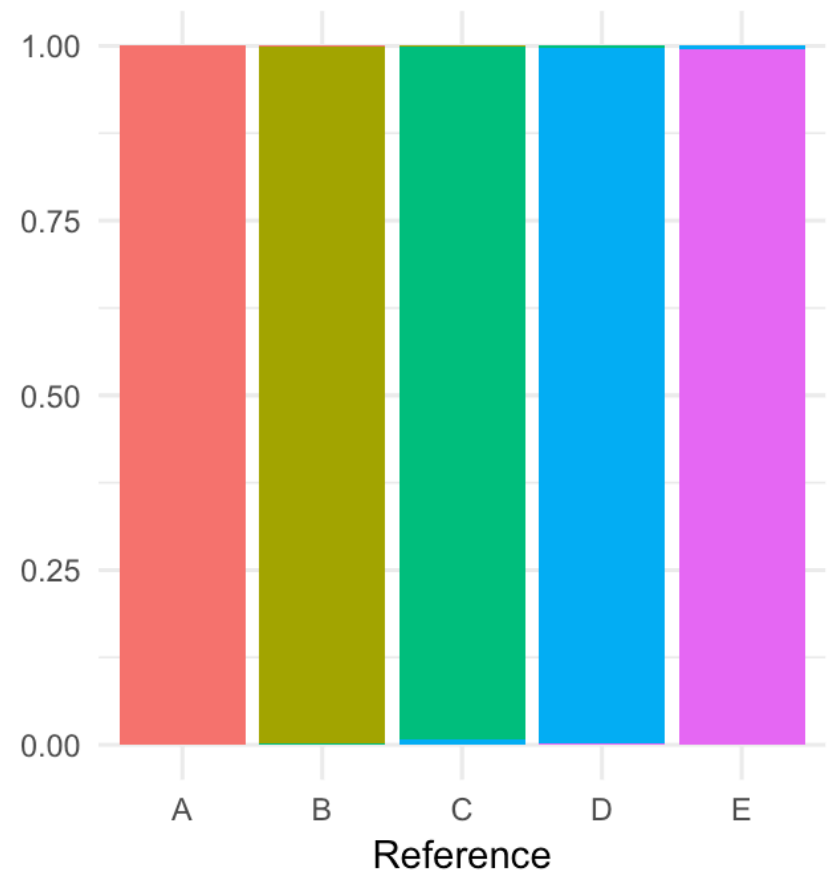                 verbose = FALSE)


prediction.training.gbm <- predict(model.gbm, newdata = training)
cm.training.gbm <- confusionMatrix(prediction.training.gbm, training$classe)

prediction.testing.gbm <- predict(model.gbm, newdata = testing)
cm.testing.gbm <- confusionMatrix(prediction.testing.gbm, testing$classe)
```

Confusion Matrix
GBM (Training)
Accuracy = 0.9988

Confusion Matrix
GBM (Testing)
Accuracy = 0.9962

# Conclusion

The random forest and generalized boosted regression are best.

Decision Tree:

- Accuracy: 0.5382
- Out of sample error: 0.4618

Random Forest:

- Accuracy: 0.9992
- Out of sample error: 810^{-4}

Generalized Boosted Regression:

- Accuracy: 0.9962
- Out of sample error: 0.0038

Finally, the model with the highest accuracy and lowest out of sample value will be applied on validation data.

```
prediction.valdiation.rf <- predict(model.rf, newdata = valdiation)
prediction.valdiation.rf
```

```
##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
##  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
```

```
pml_write_files = function(x) {
  n = length(x)
  for (i in 1:n) {
    filename = paste0("problem_id_", i, "_", prediction.valdiation.rf[i], ".txt")
    write.table(x[i], file=filename, quote=FALSE,row.names=FALSE, col.names=FALSE)
  }
}


pml_write_files(prediction.valdiation.rf)
```

# Appendix

## Cleaning the data

```
str(pml.training.csv)
```

```
## 'data.frame':    19622 obs. of  58 variables:
##  $ user_name           : Factor w/ 6 levels "adelmo","carlitos",..: 2 2 2 2 2 2
2 2 2 2 ...
##  $ raw_timestamp_part_1: int  1323084231 1323084231 1323084231 1323084232 13230
84232 1323084232 1323084232 1323084232 1323084232 1323084232 ...
##  $ raw_timestamp_part_2: int  788290 808298 820366 120339 196328 304277 368296
440390 484323 484434 ...
##  $ cvtd_timestamp      : Factor w/ 20 levels "02/12/2011 13:32",..: 9 9 9 9 9 9
9 9 9 9 ...
##  $ 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 ...
##  $ 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 ...
##  $ 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 ...
##  $ 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 ...
##  $ total_accel_dumbbell: int  37 37 37 37 37 37 37 37 37 37 ...
##  $ gyros_dumbbell_x   : num  0 0 0 0 0 0 0 0 0 0 ...
##  $ gyros_dumbbell_y   : num  −0.02 −0.02 −0.02 −0.02 −0.02 −0.02 −0.02 −0.02 −
0.02 −0.02 ...
##  $ gyros_dumbbell_z   : num  0 0 0 −0.02 0 0 0 0 0 0 ...
##  $ accel_dumbbell_x   : int  −234 −233 −232 −232 −233 −234 −232 −234 −232 −235
...
##  $ accel_dumbbell_y   : int  47 47 46 48 48 48 47 46 47 48 ...
##  $ accel_dumbbell_z   : int  −271 −269 −270 −269 −270 −269 −270 −272 −269 −270
...
##  $ magnet_dumbbell_x  : int  −559 −555 −561 −552 −554 −558 −551 −555 −549 −558
...
##  $ magnet_dumbbell_y  : int  293 296 298 303 292 294 295 300 292 291 ...
##  $ magnet_dumbbell_z  : num  −65 −64 −63 −60 −68 −66 −70 −74 −65 −69 ...
##  $ roll_forearm       : num  28.4 28.3 28.3 28.1 28 27.9 27.9 27.8 27.7 27.7 .
..
##  $ pitch_forearm      : num  −63.9 −63.9 −63.9 −63.9 −63.9 −63.9 −63.9 −63.8 −
63.8 −63.8 ...
##  $ yaw_forearm        : num  −153 −153 −152 −152 −152 −152 −152 −152 −152 −152
...
##  $ total_accel_forearm : int  36 36 36 36 36 36 36 36 36 36 ...
##  $ gyros_forearm_x    : num  0.03 0.02 0.03 0.02 0.02 0.02 0.02 0.02 0.03 0.02
...
##  $ gyros_forearm_y    : num  0 0 −0.02 −0.02 0 −0.02 0 −0.02 0 0 ...
##  $ gyros_forearm_z    : num  −0.02 −0.02 0 0 −0.02 −0.03 −0.02 0 −0.02 −0.02 .
..
##  $ accel_forearm_x    : int  192 192 196 189 189 193 195 193 193 190 ...
##  $ accel_forearm_y    : int  203 203 204 206 206 203 205 205 204 205 ...
##  $ accel_forearm_z    : int  −215 −216 −213 −214 −214 −215 −215 −213 −214 −215
...
##  $ magnet_forearm_x   : int  −17 −18 −18 −16 −17 −9 −18 −9 −16 −22 ...
##  $ magnet_forearm_y   : num  654 661 658 658 655 660 659 660 653 656 ...
```

```
## $ magnet_forearm_z    : num  476 473 469 469 473 478 470 474 476 473 ...
## $ classe              : Factor w/ 5 levels "A","B","C","D",..: 1 1 1 1 1 1 1 1
1 1 ...
```

str(pml.testing.csv)

```
## 'data.frame':    20 obs. of  57 variables:
## $ user_name           : Factor w/ 6 levels "adelmo","carlitos",..: 6 5 5 1 4 5
5 5 2 3 ...
## $ raw_timestamp_part_1: int  1323095002 1322673067 1322673075 1322832789 13224
89635 1322673149 1322673128 1322673076 1323084240 1322837822 ...
## $ raw_timestamp_part_2: int  868349 778725 342967 560311 814776 510661 766645
54671 916313 384285 ...
## $ cvtd_timestamp      : Factor w/ 11 levels "02/12/2011 13:33",..: 5 10 10 1 6
11 11 10 3 2 ...
## $ num_window          : int  74 431 439 194 235 504 485 440 323 664 ...
## $ roll_belt           : num  123 1.02 0.87 125 1.35 -5.92 1.2 0.43 0.93 114 ..
.
## $ pitch_belt          : num  27 4.87 1.82 -41.6 3.33 1.59 4.44 4.15 6.72 22.4
...
## $ yaw_belt            : num  -4.75 -88.9 -88.5 162 -88.6 -87.7 -87.3 -88.5 -93
.7 -13.1 ...
## $ total_accel_belt    : int  20 4 5 17 3 4 4 4 4 18 ...
## $ gyros_belt_x        : num  -0.5 -0.06 0.05 0.11 0.03 0.1 -0.06 -0.18 0.1 0.1
4 ...
## $ gyros_belt_y        : num  -0.02 -0.02 0.02 0.11 0.02 0.05 0 -0.02 0 0.11 ..
.
## $ gyros_belt_z        : num  -0.46 -0.07 0.03 -0.16 0 -0.13 0 -0.03 -0.02 -0.1
6 ...
## $ accel_belt_x        : int  -38 -13 1 46 -8 -11 -14 -10 -15 -25 ...
## $ accel_belt_y        : int  69 11 -1 45 4 -16 2 -2 1 63 ...
## $ accel_belt_z        : int  -179 39 49 -156 27 38 35 42 32 -158 ...
## $ magnet_belt_x       : int  -13 43 29 169 33 31 50 39 -6 10 ...
## $ magnet_belt_y       : int  581 636 631 608 566 638 622 635 600 601 ...
## $ magnet_belt_z       : int  -382 -309 -312 -304 -418 -291 -315 -305 -302 -330
...
## $ roll_arm            : num  40.7 0 0 -109 76.1 0 0 0 -137 -82.4 ...
## $ pitch_arm           : num  -27.8 0 0 55 2.76 0 0 0 11.2 -63.8 ...
## $ yaw_arm             : num  178 0 0 -142 102 0 0 0 -167 -75.3 ...
## $ total_accel_arm     : int  10 38 44 25 29 14 15 22 34 32 ...
## $ gyros_arm_x         : num  -1.65 -1.17 2.1 0.22 -1.96 0.02 2.36 -3.71 0.03 0
.26 ...
## $ gyros_arm_y         : num  0.48 0.85 -1.36 -0.51 0.79 0.05 -1.01 1.85 -0.02
-0.5 ...
## $ gyros_arm_z         : num  -0.18 -0.43 1.13 0.92 -0.54 -0.07 0.89 -0.69 -0.0
2 0.79 ...
## $ accel_arm_x         : int  16 -290 -341 -238 -197 -26 99 -98 -287 -301 ...
## $ accel_arm_y         : int  38 215 245 -57 200 130 79 175 111 -42 ...
## $ accel_arm_z         : int  93 -90 -87 6 -30 -19 -67 -78 -122 -80 ...
## $ magnet_arm_x        : int  -326 -325 -264 -173 -170 396 702 535 -367 -420 ..
.
## $ magnet_arm_y        : int  385 447 474 257 275 176 15 215 335 294 ...
```

```
##  $ magnet_arm_z         : int   481 434 413 633 617 516 217 385 520 493 ...
##  $ roll_dumbbell        : num   -17.7 54.5 57.1 43.1 -101.4 ...
##  $ pitch_dumbbell       : num   25 -53.7 -51.4 -30 -53.4 ...
##  $ yaw_dumbbell         : num   126.2 -75.5 -75.2 -103.3 -14.2 ...
##  $ total_accel_dumbbell: int   9 31 29 18 4 29 29 29 3 2 ...
##  $ gyros_dumbbell_x     : num   0.64 0.34 0.39 0.1 0.29 -0.59 0.34 0.37 0.03 0.42
## ...
##  $ gyros_dumbbell_y     : num   0.06 0.05 0.14 -0.02 -0.47 0.8 0.16 0.14 -0.21 0.
## 51 ...
##  $ gyros_dumbbell_z     : num   -0.61 -0.71 -0.34 0.05 -0.46 1.1 -0.23 -0.39 -0.2
## 1 -0.03 ...
##  $ accel_dumbbell_x     : int   21 -153 -141 -51 -18 -138 -145 -140 0 -7 ...
##  $ accel_dumbbell_y     : int   -15 155 155 72 -30 166 150 159 25 -20 ...
##  $ accel_dumbbell_z     : int   81 -205 -196 -148 -5 -186 -190 -191 9 7 ...
##  $ magnet_dumbbell_x    : int   523 -502 -506 -576 -424 -543 -484 -515 -519 -531
## ...
##  $ magnet_dumbbell_y    : int   -528 388 349 238 252 262 354 350 348 321 ...
##  $ magnet_dumbbell_z    : int   -56 -36 41 53 312 96 97 53 -32 -164 ...
##  $ roll_forearm         : num   141 109 131 0 -176 150 155 -161 15.5 13.2 ...
##  $ pitch_forearm        : num   49.3 -17.6 -32.6 0 -2.16 1.46 34.5 43.6 -63.5 19.
## 4 ...
##  $ yaw_forearm          : num   156 106 93 0 -47.9 89.7 152 -89.5 -139 -105 ...
##  $ total_accel_forearm  : int   33 39 34 43 24 43 32 47 36 24 ...
##  $ gyros_forearm_x      : num   0.74 1.12 0.18 1.38 -0.75 -0.88 -0.53 0.63 0.03 0
## .02 ...
##  $ gyros_forearm_y      : num   -3.34 -2.78 -0.79 0.69 3.1 4.26 1.8 -0.74 0.02 0.
## 13 ...
##  $ gyros_forearm_z      : num   -0.59 -0.18 0.28 1.8 0.8 1.35 0.75 0.49 -0.02 -0.
## 07 ...
##  $ accel_forearm_x      : int   -110 212 154 -92 131 230 -192 -151 195 -212 ...
##  $ accel_forearm_y      : int   267 297 271 406 -93 322 170 -331 204 98 ...
##  $ accel_forearm_z      : int   -149 -118 -129 -39 172 -144 -175 -282 -217 -7 ...
##  $ magnet_forearm_x     : int   -714 -237 -51 -233 375 -300 -678 -109 0 -403 ...
##  $ magnet_forearm_y     : int   419 791 698 783 -787 800 284 -619 652 723 ...
##  $ magnet_forearm_z     : int   617 873 783 521 91 884 585 -32 469 512 ...
```