

Practical Machine Learning Prediction Project

Jatin Kumar

08 August 2018

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).

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.

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 file describing your analysis. Please constrain the text of the writeup to < 2000 words and the number of figures to be 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.

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.

Prepare the datasets

Read the training data into a data table.

```
require(data.table)
url <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
D <- fread(url)
```

Read the testing data into a data table.

```
url <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
DTest <- fread(url)
```

Which variables in the test dataset have zero NAs? Use this tip: finding columns with all missing values in R.

Belt, arm, dumbbell, and forearm variables that do not have any missing values in the test dataset will be **predictor candidates**.

```
isAnyMissing <- sapply(DTest, function(x) any(is.na(x) | x == ""))
isPredictor <- !isAnyMissing & grepl("belt|^(fore)arm|dumbbell|forearm", names(isAnyMissing))
predCandidates <- names(isAnyMissing)[isPredictor]
predCandidates
```

```
## [1] "roll_belt"           "pitch_belt"         "yaw_belt"
## [4] "total_accel_belt"    "gyros_belt_x"       "gyros_belt_y"
## [7] "gyros_belt_z"        "accel_belt_x"       "accel_belt_y"
## [10] "accel_belt_z"        "magnet_belt_x"      "magnet_belt_y"
## [13] "magnet_belt_z"       "roll_arm"           "pitch_arm"
## [16] "yaw_arm"             "total_accel_arm"    "gyros_arm_x"
## [19] "gyros_arm_y"         "gyros_arm_z"        "accel_arm_x"
## [22] "accel_arm_y"         "accel_arm_z"        "magnet_arm_x"
## [25] "magnet_arm_y"        "magnet_arm_z"       "roll_dumbbell"
## [28] "pitch_dumbbell"      "yaw_dumbbell"       "total_accel_dumbbell"
## [31] "gyros_dumbbell_x"    "gyros_dumbbell_y"   "gyros_dumbbell_z"
## [34] "accel_dumbbell_x"    "accel_dumbbell_y"   "accel_dumbbell_z"
## [37] "magnet_dumbbell_x"   "magnet_dumbbell_y"   "magnet_dumbbell_z"
## [40] "roll_forearm"        "pitch_forearm"      "yaw_forearm"
## [43] "total_accel_forearm" "gyros_forearm_x"    "gyros_forearm_y"
## [46] "gyros_forearm_z"     "accel_forearm_x"     "accel_forearm_y"
## [49] "accel_forearm_z"     "magnet_forearm_x"    "magnet_forearm_y"
## [52] "magnet_forearm_z"
```

Subset the primary dataset to include only the **predictor candidates** and the outcome variable, **classe**.

```
varToInclude <- c("classe", predCandidates)
D <- D[, varToInclude, with=FALSE]
dim(D)
```

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

```
names(D)
```

```
## [1] "classe"           "roll_belt"         "pitch_belt"
## [4] "yaw_belt"         "total_accel_belt"  "gyros_belt_x"
## [7] "gyros_belt_y"     "gyros_belt_z"     "accel_belt_x"
## [10] "accel_belt_y"     "accel_belt_z"     "magnet_belt_x"
## [13] "magnet_belt_y"    "magnet_belt_z"    "roll_arm"
## [16] "pitch_arm"        "yaw_arm"          "total_accel_arm"
## [19] "gyros_arm_x"      "gyros_arm_y"      "gyros_arm_z"
```

```
## [22] "accel_arm_x"      "accel_arm_y"      "accel_arm_z"
## [25] "magnet_arm_x"     "magnet_arm_y"     "magnet_arm_z"
## [28] "roll_dumbbell"    "pitch_dumbbell"   "yaw_dumbbell"
## [31] "total_accel_dumbbell" "gyros_dumbbell_x" "gyros_dumbbell_y"
## [34] "gyros_dumbbell_z" "accel_dumbbell_x" "accel_dumbbell_y"
## [37] "accel_dumbbell_z" "magnet_dumbbell_x" "magnet_dumbbell_y"
## [40] "magnet_dumbbell_z" "roll_forearm"     "pitch_forearm"
## [43] "yaw_forearm"      "total_accel_forearm" "gyros_forearm_x"
## [46] "gyros_forearm_y"  "gyros_forearm_z"  "accel_forearm_x"
## [49] "accel_forearm_y"  "accel_forearm_z"  "magnet_forearm_x"
## [52] "magnet_forearm_y" "magnet_forearm_z"
```

Make classe into a factor.

```
D <- D[, classe := factor(D[, classe])]
D[, .N, classe]
```

```
##   classe    N
## 1:     A 5580
## 2:     B 3797
## 3:     C 3422
## 4:     D 3216
## 5:     E 3607
```

Split the dataset into a 60% training and 40% probing dataset.

```
require(caret)
seed <- as.numeric(as.Date("2014-10-26"))
set.seed(seed)
inTrain <- createDataPartition(D$classe, p=0.6)
DTrain <- D[inTrain[[1]]]
DProbe <- D[-inTrain[[1]]]
```

Preprocess the prediction variables by centering and scaling.

```
X <- DTrain[, predCandidates, with=FALSE]
preProc <- preProcess(X)
preProc
```

```
## Created from 11776 samples and 52 variables
##
## Pre-processing:
##   - centered (52)
##   - ignored (0)
##   - scaled (52)
```

```
XCS <- predict(preProc, X)
DTrainCS <- data.table(data.frame(classe = DTrain[, classe], XCS))
```

Apply the centering and scaling to the probing dataset.

```
X <- DProbe[, predCandidates, with=FALSE]
XCS <- predict(preProc, X)
DProbeCS <- data.table(data.frame(classe = DProbe[, classe], XCS))
```

Check for near zero variance.

```
nzv <- nearZeroVar(DTrainCS, saveMetrics=TRUE)
if (any(nzv$nzv)) nzv else message("No variables with near zero variance")
```

```
## No variables with near zero variance
```

Examine groups of prediction variables.

```
histGroup <- function (data, regex) {  
  col <- grep(regex, names(data))  
  col <- c(col, which(names(data) == "classe"))  
  require(reshape2)  
  n <- nrow(data)  
  DMelted <- melt(data[, col, with=FALSE][, rownum := seq(1, n)], id.vars=c("rownum", "classe"))  
  require(ggplot2)  
  ggplot(DMelted, aes(x=classe, y=value)) +  
    geom_violin(aes(color=classe, fill=classe), alpha=1/2) +  
  #   geom_jitter(aes(color=classe, fill=classe), alpha=1/10) +  
  #   geom_smooth(aes(group=1), method="gam", color="black", alpha=1/2, size=2) +  
  facet_wrap(~ variable, scale="free_y") +  
  scale_color_brewer(palette="Spectral") +  
  scale_fill_brewer(palette="Spectral") +  
  labs(x="", y="") +  
  theme(legend.position="none")  
}  
histGroup(DTrainCS, "belt")
```

```
## Loading required package: reshape2
```

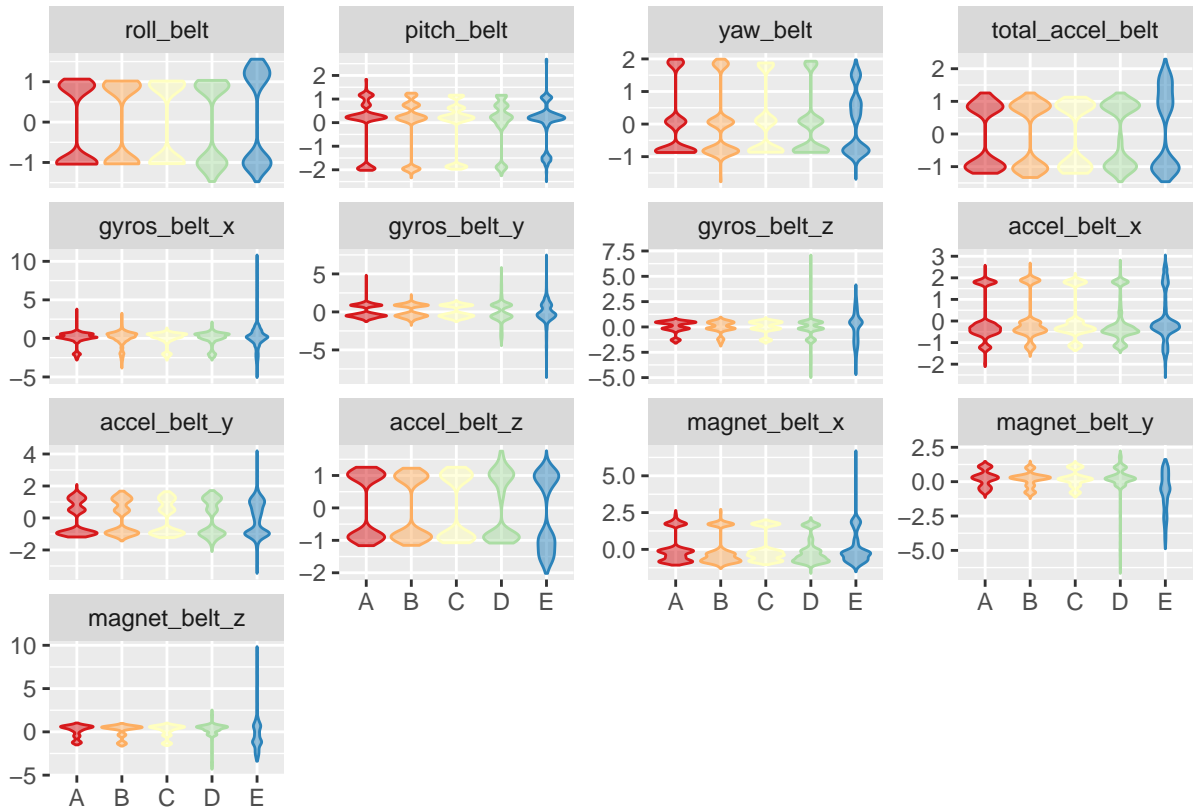
```
##
```

```
## Attaching package: 'reshape2'
```

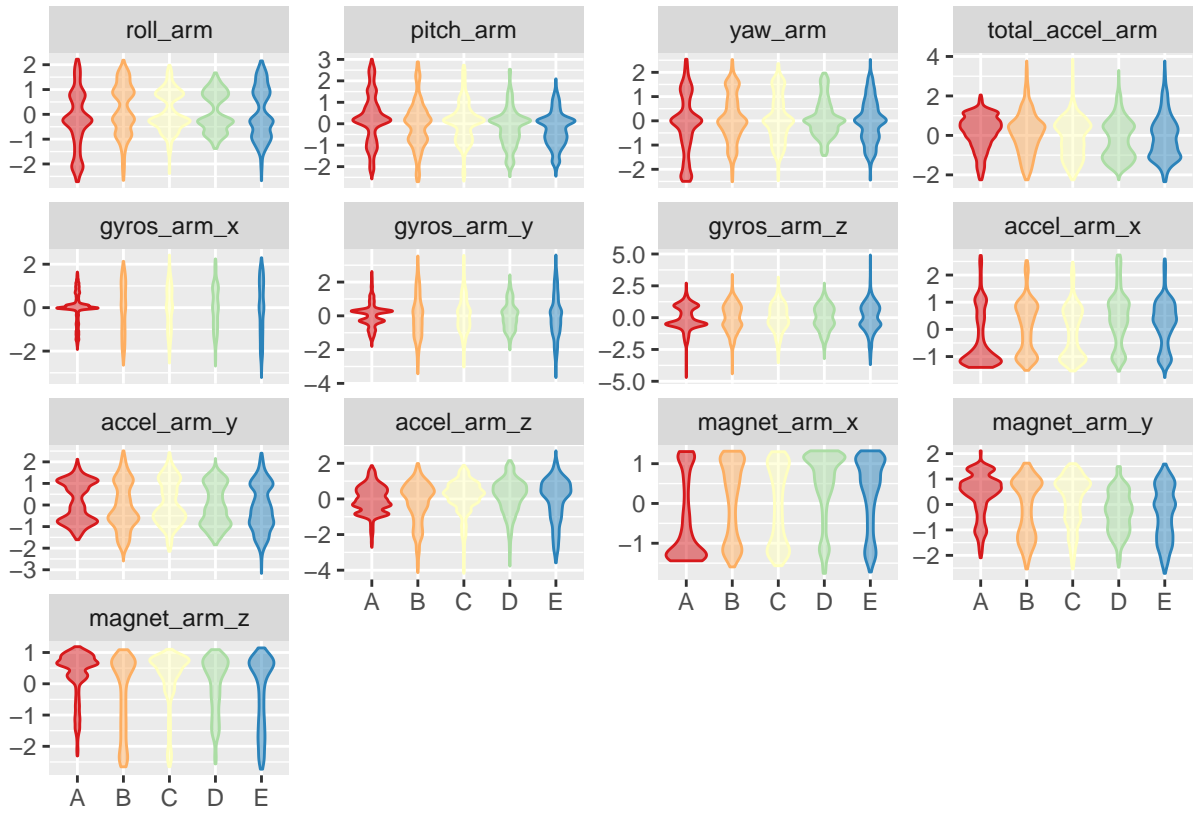
```
## The following objects are masked from 'package:data.table':
```

```
##
```

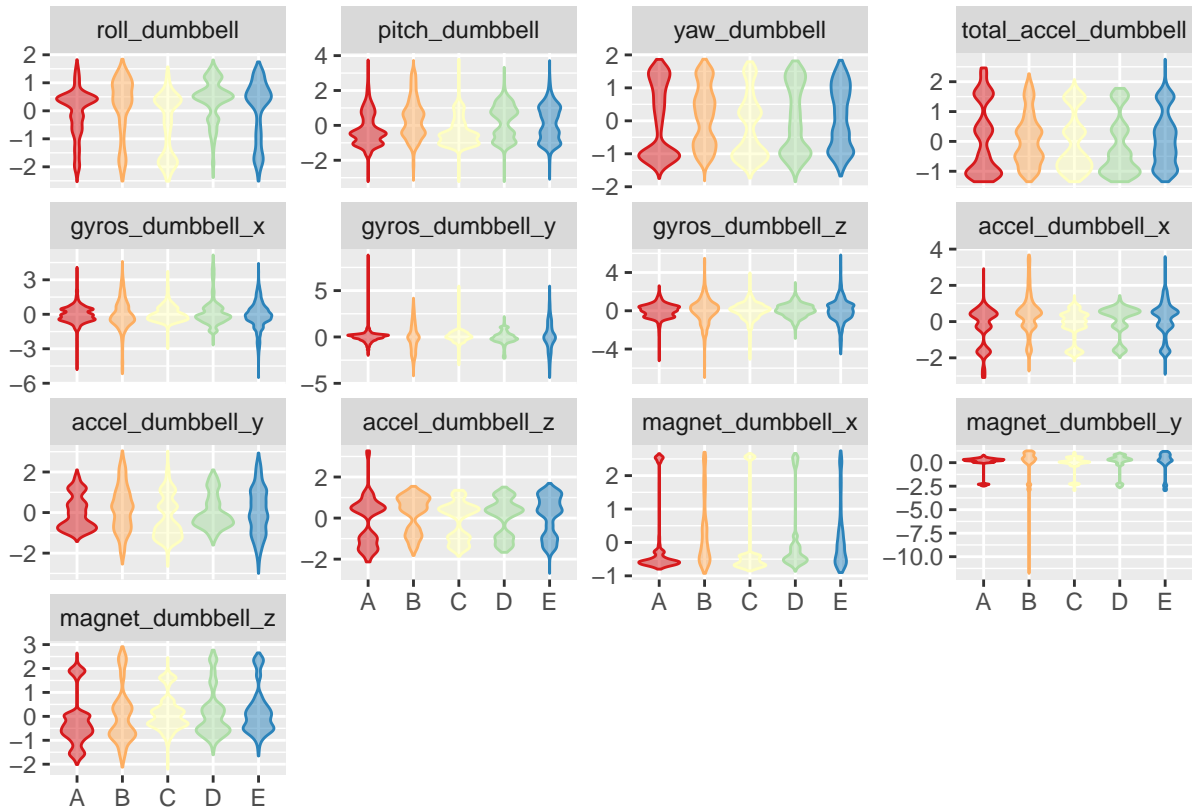
```
##      dcast, melt
```



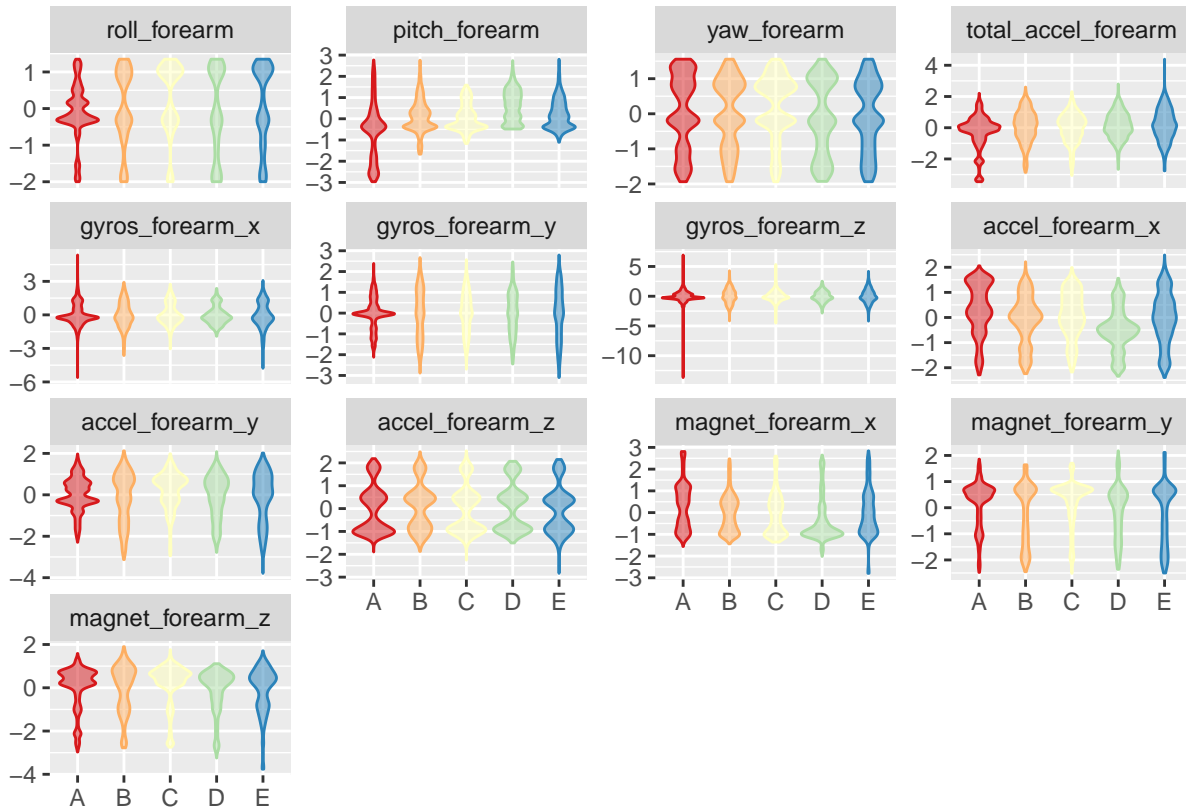
```
histGroup(DTrainCS, "[^(fore)]arm")
```



```
histGroup(DTrainCS, "dumbbell")
```



```
histGroup(DTrainCS, "forearm")
```



Train a prediction model

Using random forest, the out of sample error should be small. The error will be estimated using the 40% probing sample. I would be quite happy with an error estimate of 3% or less.

Set up the parallel clusters.

```
require(parallel)
```

```
## Loading required package: parallel
```

```
require(doParallel)
```

```
## Loading required package: doParallel
```

```
## Loading required package: foreach
```

```
## Loading required package: iterators
```

```
cl <- makeCluster(detectCores() - 1)
registerDoParallel(cl)
```

Set the control parameters.

```
ctrl <- trainControl(classProbs=TRUE,
                      savePredictions=TRUE,
                      allowParallel=TRUE)
```

Fit model over the tuning parameters.


```
method <- "rf"
system.time(trainingModel <- train(classe ~ ., data=DTrainCS, method=method))
```

```
##      user  system elapsed
## 41.315   0.381 786.109
```

Stop the clusters.

```
stopCluster(cl)
```

Evaluate the model on the training dataset

```
trainingModel
```

```
## Random Forest
##
## 11776 samples
##    52 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 11776, 11776, 11776, 11776, 11776, 11776, ...
## Resampling results across tuning parameters:
##
##  mtry  Accuracy  Kappa
##    2    0.9855896 0.9817724
##   27    0.9873282 0.9839718
##   52    0.9771134 0.9710496
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.
```

```
hat <- predict(trainingModel, DTrainCS)
confusionMatrix(hat, DTrain[, classe])
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 3348    0    0    0    0
##      B    0 2279    0    0    0
##      C    0    0 2054    0    0
##      D    0    0    0 1930    0
##      E    0    0    0    0 2165
##
## Overall Statistics
##
##              Accuracy : 1
##              95% CI : (0.9997, 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
```

Evaluate the model on the probing dataset

```
hat <- predict(trainingModel, DProbeCS)
confusionMatrix(hat, DProbeCS[, classe])
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 2229   16    0    0    0
##           B    2 1498    7    1    2
##           C    0    4 1352   18    7
##           D    0    0    9 1265    8
##           E    1    0    0    2 1425
##
## Overall Statistics
##
##           Accuracy : 0.9902
##           95% CI : (0.9877, 0.9922)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9876
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9987  0.9868  0.9883  0.9837  0.9882
## Specificity      0.9971  0.9981  0.9955  0.9974  0.9995
## Pos Pred Value   0.9929  0.9921  0.9790  0.9867  0.9979
## Neg Pred Value   0.9995  0.9968  0.9975  0.9968  0.9974
## Prevalence       0.2845  0.1935  0.1744  0.1639  0.1838
## Detection Rate   0.2841  0.1909  0.1723  0.1612  0.1816
## Detection Prevalence 0.2861  0.1925  0.1760  0.1634  0.1820
## Balanced Accuracy 0.9979  0.9925  0.9919  0.9905  0.9939
```

Display the final model

```
varImp(trainingModel)

## rf variable importance
##
##    only 20 most important variables shown (out of 52)
##
##              Overall
## roll_belt      100.000
## pitch_forearm   60.142
## yaw_belt        53.838
## pitch_belt      46.489
## roll_forearm    45.165
## magnet_dumbbell_y 43.873
## magnet_dumbbell_z 42.966
## accel_dumbbell_y 21.103
## magnet_dumbbell_x 17.692
## roll_dumbbell   17.641
## accel_forearm_x 17.143
## magnet_forearm_z 14.021
## total_accel_dumbbell 13.993
## accel_dumbbell_z 13.960
## magnet_belt_y    13.765
## magnet_belt_z    13.345
## accel_belt_z     13.344
## yaw_arm          11.625
## gyros_belt_z     11.092
## magnet_belt_x     9.877

trainingModel$finalModel
```

```
##
## Call:
##  randomForest(x = x, y = y, mtry = param$mtry)
##              Type of random forest: classification
##              Number of trees: 500
## No. of variables tried at each split: 27
##
##              OOB estimate of  error rate: 0.86%
## Confusion matrix:
##      A      B      C      D      E class.error
## A 3341      5      2      0      0 0.002090800
## B   17 2253      8      1      0 0.011408513
## C      0   13 2033      8      0 0.010223953
## D      1      1  30 1895      3 0.018134715
## E      0      2      2      8 2153 0.005542725
```

The estimated error rate is less than 1%.

Save training model object for later.

```
save(trainingModel, file="trainingModel.RData")
```

Predict on the test data

Load the training model.

```
load(file="trainingModel.RData", verbose=TRUE)
```

```
## Loading objects:
```

```
##   trainingModel
```

Get predictions and evaluate.

```
DTestCS <- predict(preProc, DTest[, predCandidates, with=FALSE])
```

```
hat <- predict(trainingModel, DTestCS)
```

```
DTest <- cbind(hat, DTest)
```

```
subset(DTest, select=names(DTest)[grep("belt|[^{(fore)}]arm|dumbbell|forearm", names(DTest)), invert=TRUE])
```

```
##      hat V1 user_name raw_timestamp_part_1 raw_timestamp_part_2
## 1:  B  1    pedro          1323095002          868349
## 2:  A  2    jeremy          1322673067          778725
## 3:  B  3    jeremy          1322673075          342967
## 4:  A  4    adelmo          1322832789          560311
## 5:  A  5    eurico          1322489635          814776
## 6:  E  6    jeremy          1322673149          510661
## 7:  D  7    jeremy          1322673128          766645
## 8:  B  8    jeremy          1322673076           54671
## 9:  A  9    carlitos        1323084240          916313
## 10: A 10    charles        1322837822          384285
## 11: B 11    carlitos        1323084277           36553
## 12: C 12    jeremy          1322673101          442731
## 13: B 13    eurico          1322489661          298656
## 14: A 14    jeremy          1322673043          178652
## 15: E 15    jeremy          1322673156          550750
## 16: E 16    eurico          1322489713          706637
## 17: A 17    pedro          1323094971          920315
## 18: B 18    carlitos        1323084285          176314
## 19: B 19    pedro          1323094999          828379
## 20: B 20    eurico          1322489658          106658
##      cvtd_timestamp new_window num_window problem_id
## 1: 05/12/2011 14:23      no         74          1
## 2: 30/11/2011 17:11      no        431          2
## 3: 30/11/2011 17:11      no        439          3
## 4: 02/12/2011 13:33      no        194          4
## 5: 28/11/2011 14:13      no        235          5
## 6: 30/11/2011 17:12      no        504          6
## 7: 30/11/2011 17:12      no        485          7
## 8: 30/11/2011 17:11      no        440          8
## 9: 05/12/2011 11:24      no        323          9
## 10: 02/12/2011 14:57      no        664         10
## 11: 05/12/2011 11:24      no        859         11
## 12: 30/11/2011 17:11      no        461         12
## 13: 28/11/2011 14:14      no        257         13
## 14: 30/11/2011 17:10      no        408         14
```

## 15:	30/11/2011 17:12	no	779	15
## 16:	28/11/2011 14:15	no	302	16
## 17:	05/12/2011 14:22	no	48	17
## 18:	05/12/2011 11:24	no	361	18
## 19:	05/12/2011 14:23	no	72	19
## 20:	28/11/2011 14:14	no	255	20