

Prediction Assignment Writeup

Shu Yan

February 22, 2015

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. In this project, we will be performing machine learning algorithms to predict the correctness of barbell lifting based on data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants. 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>).

Data preprocessing

We first load the training data set from the given *URL*. It is reasonable to believe that the prediction only depends on the body motions. Therefore we only select predictors coming from the data of the sensors. Meanwhile we should remove predictors that contain any *NAs*. (fortunately, we don't have any after the first cleaning step)

```
urlTrain <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
pmlTrain <- read.csv(urlTrain, na.strings=c("NA", "#DIV/0!"))
feature <- c(grep("belt_x|arm_x|bell_x|belt_y|arm_y|bell_y|belt_z|arm_z|bell_z",
                 names(pmlTrain)), length(pmlTrain))
pmlTrain <- pmlTrain[, feature]
#pmlTrain <- pmlTrain[-which(sapply(pmlTrain, anyNA))]
urlTest <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
pmlTest <- read.csv(urlTest, na.strings=c("NA", "#DIV/0!"))
pmlTest <- pmlTest[-which(sapply(pmlTest, anyNA))]
```

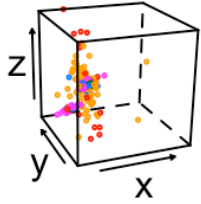
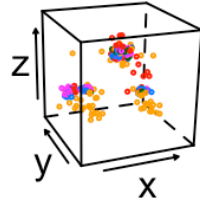
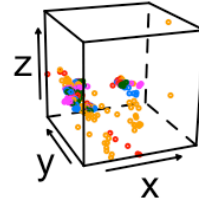
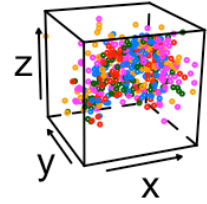
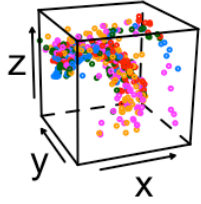
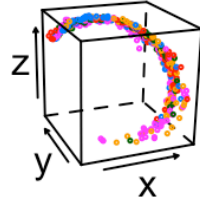
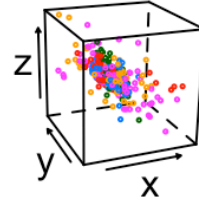
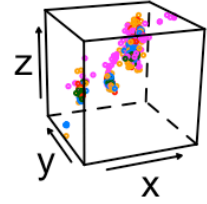
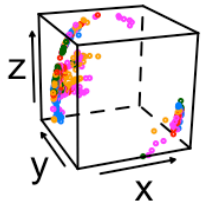
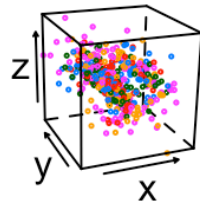
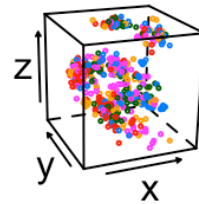
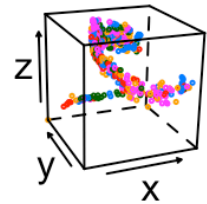
data visualization

Our data now contains 12 features, each of which has 3 spatial components, we can visualize the features in 3D plots. Here we randomly select a small portion of all observations. (Note that the figures are not drawn on same scales.)

```

library(lattice)
vis <- sample(1:19622,500)
i <- 1
has_more = TRUE
par.set <-
  list(axis.line = list(col = "transparent"),
        clip = list(panel = "off"))
while(i < length(pmlTrain)-1) {
  if(i > 33) {
    has_more = FALSE
  }
  j = (i+2)/3
  str <- names(pmlTrain[i])
  print(cloud(pmlTrain[vis,i+2] ~ pmlTrain[vis,i] * pmlTrain[vis,i+1],
    data = pmlTrain[vis,], cex = .2,
    groups = classe,
    xlab = "x", ylab = "y", zlab = "z",
    main = substr(str, 1, nchar(str)-2),
    screen = list(z = 20, x = -70, y = 3),
    par.settings = par.set,
    scales = list(col = "black")),
    split = c((j-1)%4+1,ceiling(j/4),4,3), more = has_more)
  i = i + 3
}

```

gyros_belt**accel_belt****magnet_belt****gyros_arm****accel_arm****magnet_arm****gyros_dumbbell****accel_dumbbell****magnet_dumbbell****gyros_forearm****accel_forearm****magnet_forearm**

Data slicing

The data sets are splitted into training set and testing set as follow

```
library(caret)
```

```
## Loading required package: ggplot2
```

```
inTrain <- createDataPartition(y=pmlTrain$classe, p=.7, list=FALSE)
training <- pmlTrain[inTrain,]
testing  <- pmlTrain[-inTrain,]
dim(training); dim(testing)
```

```
## [1] 13737    37
```

```
## [1] 5885     37
```

Data training

The **rpart** is fast but gives poor prediction. Here we will be using **gbm** method, which is accurate but quite time consuming.

```
library(caret)
modFit <- train(classe~., method="gbm",data=training,verbose=FALSE)
```

```
## Loading required package: gbm
## Loading required package: survival
## Loading required package: splines
##
## Attaching package: 'survival'
##
## The following object is masked from 'package:caret':
##
##      cluster
##
## Loading required package: parallel
## Loaded gbm 2.1
## Loading required package: plyr
```

```
predicted <- predict(modFit,newdata=testing)
confusionMatrix(predicted,testing$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A     B     C     D     E
##           A 1601    79    19    13    12
##           B   15   986    55    13    40
##           C   27    56   935    73    13
##           D   31     5    15   850    30
##           E    0    13     2    15   987
##
## Overall Statistics
##
##           Accuracy : 0.9106
##           95% CI : (0.903, 0.9178)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.8868
##           McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9564    0.8657    0.9113    0.8817    0.9122
## Specificity      0.9708    0.9741    0.9652    0.9835    0.9938
## Pos Pred Value   0.9287    0.8891    0.8469    0.9130    0.9705
## Neg Pred Value   0.9825    0.9680    0.9810    0.9770    0.9805
## Prevalence       0.2845    0.1935    0.1743    0.1638    0.1839
## Detection Rate   0.2720    0.1675    0.1589    0.1444    0.1677
## Detection Prevalence 0.2929    0.1884    0.1876    0.1582    0.1728
## Balanced Accuracy 0.9636    0.9199    0.9383    0.9326    0.9530
```

Predict

Finally we use the test data set for prediction

```
predict(modFit,newdata=pmlTest)
```

```
##  [1] B A B A A E D B A A A C B A E E A B B B
## Levels: A B C D E
```