

Prediction Assignment Writeup :Prediction of how was the Exercise Performed

Krishna Mahajan , 22 Dec, 2015

summary

Human activity recognition research has traditionally focused on discriminating between different activities. However, the "how (well)" investigation has only received little attention so far, even though it potentially provides useful information for a large variety of applications, such as sports training [<http://groupware.les.inf.puc-rio.br/har>]

For the prediction of how well individuals performed the assigned exercise six young health participants were asked to perform one set of 10 repetitions of the Unilateral Dumbbell Biceps Curl in five different fashions: exactly according to the specification

Class A: Exactly according to the specification

Class B : Throwing the elbows to the front

Class C : Lifting the dumbbell only halfway

Class D : Lowering the dumbbell only halfway

Class E : Throwing the hips to the front

Data Cleaning

The data for this project come was obtained from [<http://groupware.les.inf.puc-rio.br/har>]. Two data set were available a training set and a test set for which 20 individuals without any classification for the class of exercise was available

#Data Loading

```
pmlTrain<-read.csv("pml-training.csv", header=T, na.strings=c("NA",
"#DIV/0!"))
pmlTest<-read.csv("pml-testing.csv", header=T, na.string=c("NA", "#DIV/0!"))
```

Training data was partitioned and preprocessed using the code described below. In brief, all variables with at least one "NA" were excluded from the analysis. Variables related to time and user information were excluded for a total of 51 variables and 19622 class measurements. Same variables were maintained in the test data set (Validation dataset) to be used for predicting the 20 test cases provided.

```
## NA exclusion for all available variables
noNApmlTrain<-pmlTrain[, apply(pmlTrain, 2, function(x) !any(is.na(x)))]
dim(noNApmlTrain)

## [1] 19622    60
```

```
## variables with user information, time and undefined
cleanpmlTrain<-noNAPmlTrain[,-c(1:8)]
dim(cleanpmlTrain)

## [1] 19622    52

## 20 test cases provided clean info - Validation data set
cleanpmltest<-pmlTest[,names(cleanpmlTrain[, -52])]
dim(cleanpmltest)

## [1] 20 51
```

Data Partitioning and Prediction Process

The cleaned downloaded data set was subset in order to generate a test set independent from the 20 cases provided set. Partitioning was performed to obtain a 75% training set and a 25% test set.

```
#data cleaning
library(caret)
inTrain<-createDataPartition(y=cleanpmlTrain$classe, p=0.75,list=F)
training<-cleanpmlTrain[inTrain,]
test<-cleanpmlTrain[-inTrain,]
#Training and test set dimensions
dim(training)

## [1] 14718    52

dim(test)

## [1] 4904    52
```

Result and conclusion

Random forest trees were generated for the training dataset using cross-validation. Then the generated algorithm was examined under the partitioned training set to examine the accuracy and estimated error of prediction. By using 51 predictors for five classes using cross-validation at a 5-fold an accuracy of 99.2% with a 95% CI [0.989-0.994] was achieved accompanied by a Kappa value of 0.99.

```
library(caret)
set.seed(13333)
fitControl2<-trainControl(method="cv", number=5, allowParallel=T, verbose=T)
rffit<-train(classe~.,data=training, method="rf", trControl=fitControl2,
verbose=F)

## + Fold1: mtry= 2
## - Fold1: mtry= 2
## + Fold1: mtry=26
## - Fold1: mtry=26
## + Fold1: mtry=51
## - Fold1: mtry=51
```

```

## + Fold2: mtry= 2
## - Fold2: mtry= 2
## + Fold2: mtry=26
## - Fold2: mtry=26
## + Fold2: mtry=51
## - Fold2: mtry=51
## + Fold3: mtry= 2
## - Fold3: mtry= 2
## + Fold3: mtry=26
## - Fold3: mtry=26
## + Fold3: mtry=51
## - Fold3: mtry=51
## + Fold4: mtry= 2
## - Fold4: mtry= 2
## + Fold4: mtry=26
## - Fold4: mtry=26
## + Fold4: mtry=51
## - Fold4: mtry=51
## + Fold5: mtry= 2
## - Fold5: mtry= 2
## + Fold5: mtry=26
## - Fold5: mtry=26
## + Fold5: mtry=51
## - Fold5: mtry=51
## Aggregating results
## Selecting tuning parameters
## Fitting mtry = 26 on full training set

predrf<-predict(rffit, newdata=test)
confusionMatrix(predrf, test$classe)

## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1394    7    0    0    0
##           B    0  936    5    0    0
##           C    0    5  850    5    1
##           D    0    1    0  797    2
##           E    1    0    0    2  898
##
## Overall Statistics
##
##               Accuracy : 0.9941
##               95% CI : (0.9915, 0.996)
##       No Information Rate : 0.2845
##       P-Value [Acc > NIR] : < 2.2e-16
##
##               Kappa : 0.9925
## Mcnemar's Test P-Value : NA

```

```
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9993  0.9863  0.9942  0.9913  0.9967
## Specificity      0.9980  0.9987  0.9973  0.9993  0.9993
## Pos Pred Value   0.9950  0.9947  0.9872  0.9963  0.9967
## Neg Pred Value   0.9997  0.9967  0.9988  0.9983  0.9993
## Prevalence       0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate   0.2843  0.1909  0.1733  0.1625  0.1831
## Detection Prevalence 0.2857  0.1919  0.1756  0.1631  0.1837
## Balanced Accuracy 0.9986  0.9925  0.9957  0.9953  0.9980

pred20<-predict(rffit, newdata=cleanpmltest)
# Output for the prediction of the 20 cases provided
pred20

## [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
```

A boosting algorithm was also run to confirm and be able to compare predictions. Data is not shown but the boosting approach presented less accuracy (96%) (Data not shown). However, when the predictions for the 20 test cases were compared match was same for both ran algoritms.

```
fitControl2<-trainControl(method="cv", number=5, allowParallel=T, verbose=T)
gmbfit<-train(classe~.,data=training, method="gbm", trControl=fitControl2,
verbose=F)

## + Fold1: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## - Fold1: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## + Fold1: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## - Fold1: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## + Fold1: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## - Fold1: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## + Fold2: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## - Fold2: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## + Fold2: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## - Fold2: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## + Fold2: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
```

```

## - Fold2: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## + Fold3: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## - Fold3: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## + Fold3: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## - Fold3: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## + Fold3: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## - Fold3: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## + Fold4: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## - Fold4: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## + Fold4: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## - Fold4: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## + Fold4: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## - Fold4: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## + Fold5: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## - Fold5: shrinkage=0.1, interaction.depth=1, n.minobsinnode=10,
n.trees=150
## + Fold5: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## - Fold5: shrinkage=0.1, interaction.depth=2, n.minobsinnode=10,
n.trees=150
## + Fold5: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## - Fold5: shrinkage=0.1, interaction.depth=3, n.minobsinnode=10,
n.trees=150
## Aggregating results
## Selecting tuning parameters
## Fitting n.trees = 150, interaction.depth = 3, shrinkage = 0.1,
n.minobsinnode = 10 on full training set

gmbfit$finalModel

## A gradient boosted model with multinomial loss function.
## 150 iterations were performed.
## There were 51 predictors of which 43 had non-zero influence.

class(gmbfit)

```

```
## [1] "train"          "train.formula"

predgmb<-predict(gmbfit, newdata=test)
confusionMatrix(predgmb, test$classe)

## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 1374    27     0     0     0
##      B   14   888    21     9    10
##      C    4    30   824    20    11
##      D    2     1     8   765    10
##      E    1     3     2    10   870
##
## Overall Statistics
##
##              Accuracy : 0.9627
##              95% CI : (0.957, 0.9678)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9528
##  Mcnemar's Test P-Value : 0.0001673
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9849   0.9357   0.9637   0.9515   0.9656
## Specificity          0.9923   0.9863   0.9839   0.9949   0.9960
## Pos Pred Value       0.9807   0.9427   0.9269   0.9733   0.9819
## Neg Pred Value       0.9940   0.9846   0.9923   0.9905   0.9923
## Prevalence           0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate       0.2802   0.1811   0.1680   0.1560   0.1774
## Detection Prevalence 0.2857   0.1921   0.1813   0.1603   0.1807
## Balanced Accuracy    0.9886   0.9610   0.9738   0.9732   0.9808

predtrain<-predict(gmbfit, newdata=training)
confusionMatrix(predtrain, training$classe)

## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 4150    75     0     0     0
##      B   25 2713    48    13    15
##      C    6   56 2489    69    24
##      D    4     0    26 2316    19
##      E    0     4     4    14 2648
##
## Overall Statistics
```

```
##
##           Accuracy : 0.9727
##           95% CI : (0.9699, 0.9753)
##       No Information Rate : 0.2843
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9654
##  McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9916  0.9526  0.9696  0.9602  0.9786
## Specificity      0.9929  0.9915  0.9872  0.9960  0.9982
## Pos Pred Value   0.9822  0.9641  0.9414  0.9793  0.9918
## Neg Pred Value   0.9967  0.9887  0.9935  0.9922  0.9952
## Prevalence       0.2843  0.1935  0.1744  0.1639  0.1839
## Detection Rate   0.2820  0.1843  0.1691  0.1574  0.1799
## Detection Prevalence 0.2871  0.1912  0.1796  0.1607  0.1814
## Balanced Accuracy 0.9923  0.9720  0.9784  0.9781  0.9884

predtrain<-predict(gmbfit, newdata=training)
confusionMatrix(predtrain, training$classe)

## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 4150   75    0    0    0
##           B   25 2713   48   13   15
##           C    6   56 2489   69   24
##           D    4    0   26 2316   19
##           E    0    4    4   14 2648
##
## Overall Statistics
##
##           Accuracy : 0.9727
##           95% CI : (0.9699, 0.9753)
##       No Information Rate : 0.2843
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9654
##  McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9916  0.9526  0.9696  0.9602  0.9786
## Specificity      0.9929  0.9915  0.9872  0.9960  0.9982
## Pos Pred Value   0.9822  0.9641  0.9414  0.9793  0.9918
```

## Neg Pred Value	0.9967	0.9887	0.9935	0.9922	0.9952
## Prevalence	0.2843	0.1935	0.1744	0.1639	0.1839
## Detection Rate	0.2820	0.1843	0.1691	0.1574	0.1799
## Detection Prevalence	0.2871	0.1912	0.1796	0.1607	0.1814
## Balanced Accuracy	0.9923	0.9720	0.9784	0.9781	0.9884

Once, the predictions were obtained for the 20 test cases provided, the below shown script was used to obtain single text files to be uploaded to the courses web site to comply with the submission assignment. 20 out of 20 hits also confirmed the accuracy of the obtained models.

```
getwd()

## [1] "C:/Users/Krish/Google Drive/Coursera/backup/Practical Machine
Learning/Project"

pml_write_files = function(x){
  n = length(x)
  for(i in 1:n){
    filename = paste0("problem_id_",i,".txt")

write.table(x[i],file=filename,quote=FALSE,row.names=FALSE,col.names=FALSE)
  }
}
pml_write_files(pred20)
```