

Random Forest

Calvin Seto

January 14, 2016

Random Forest

mtry = number of variables randomly sampled as candidates at each split Defaults: Classification \sqrt{p}
Regression $p/3$ where p = number of variables

```
library(caret)
```

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

```
setwd("~/Dropbox/jhudatascience/8_Practical_Machine_Learning/CourseProject")
```

```
# setwd("C:/Users/Calvin/Calvinsbiz/Dropbox/jhudatascience/8_Practical_Machine_Learning/CourseProject")
```

```
pmlTrain1 <- read.csv("data/pml-training.csv", stringsAsFactors = FALSE, na.strings = c("#DIV/0!", "", "N"))
```

```
pmlTrain1MissingCounts <- sapply(pmlTrain1, function(x)sum(is.na(x)))  
pmlTrain1Complete <- pmlTrain1MissingCounts[pmlTrain1MissingCounts==0]  
pmlTrain2 <- pmlTrain1[,names(pmlTrain1Complete)]
```

```
inTrain <- createDataPartition(y=pmlTrain2$classe,  
                               p=0.75, list=FALSE)
```

```
training <- pmlTrain2[inTrain,]  
testing <- pmlTrain2[-inTrain,]
```

```
predictors <- training[,8:59]  
outcome <- as.factor(training[,60])
```

```
# configure parallel
```

```
library(parallel)  
library(doParallel)
```

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

```
cluster <- makeCluster(detectCores() - 1) # convention to leave 1 core for OS  
registerDoParallel(cluster)
```

```
# seed  
set.seed(168)
```

```
# default is bootstrap  
fitRFControl <- trainControl(method="cv",
```

```
                                number=10
)
```

```
# fitRFGrid <- expand.grid(mtry=
# )
```

```
"Start Time "; Sys.time()
```

```
## [1] "Start Time "
```

```
## [1] "2016-01-14 13:09:02 EST"
```

```
fitRF <- train(x=predictors,
              y=outcome,
              data=training,
              method="rf",
              trControl=fitRFControl
)
```

```
## Loading required package: randomForest
## randomForest 4.6-12
## Type rfNews() to see new features/changes/bug fixes.
```

```
"End Time "; Sys.time()
```

```
## [1] "End Time "
```

```
## [1] "2016-01-14 13:17:48 EST"
```

```
stopCluster(cluster)
```

```
# show model summary
fitRF
```

```
## Random Forest
##
## 14718 samples
##    52 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 13246, 13246, 13247, 13247, 13246, 13246, ...
## Resampling results across tuning parameters:
##
##  mtry  Accuracy   Kappa     Accuracy SD   Kappa SD
##    2    0.9931378 0.9913188 0.002036074   0.002577022
##   27    0.9929336 0.9910610 0.001475788   0.001867019
##   52    0.9843728 0.9802307 0.001811996   0.002295538
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 2.
```

```
fitRF$resample
```

```
##      Accuracy      Kappa Resample
## 1  0.9966033 0.9957040   Fold02
## 2  0.9945652 0.9931251   Fold01
## 3  0.9925221 0.9905423   Fold03
## 4  0.9918478 0.9896869   Fold06
## 5  0.9925272 0.9905448   Fold05
## 6  0.9918423 0.9896785   Fold04
## 7  0.9898167 0.9871126   Fold07
## 8  0.9932019 0.9914016   Fold10
## 9  0.9959239 0.9948441   Fold09
## 10 0.9925272 0.9905478   Fold08
```

```
confusionMatrix.train(fitRF)
```

```
## Cross-Validated (10 fold) Confusion Matrix
##
## (entries are percentages of table totals)
##
##           Reference
## Prediction   A    B    C    D    E
##           A 28.4  0.1  0.0  0.0  0.0
##           B  0.0 19.2  0.2  0.0  0.0
##           C  0.0  0.1 17.3  0.3  0.0
##           D  0.0  0.0  0.0 16.1  0.0
##           E  0.0  0.0  0.0  0.0 18.3
```

```
# Make predictions and make table
```

```
pred <- predict(fitRF,testing)
testing$predRight <- pred==testing$classe
table(pred,testing$classe)
```

```
##
## pred   A    B    C    D    E
##   A 1392    2    0    0    0
##   B   3  942   11    0    0
##   C   0    5  843   22    2
##   D   0    0    1  782    4
##   E   0    0    0    0  895
```