

QUIZ3

Davy Meesemaeker

19/1/2018

Practical machine learning - Quiz 3

Question 1 :

For this quiz we will be using several R packages. R package versions change over time, the right answers have been checked using the following versions of the packages.

AppliedPredictiveModeling: v1.1.6 caret: v6.0.47 ElemStatLearn: v2012.04-0 pgmm: v1.1 rpart: v4.1.8

If you aren't using these versions of the packages, your answers may not exactly match the right answer, but hopefully should be close.

Load the cell segmentation data from the AppliedPredictiveModeling package using the commands:

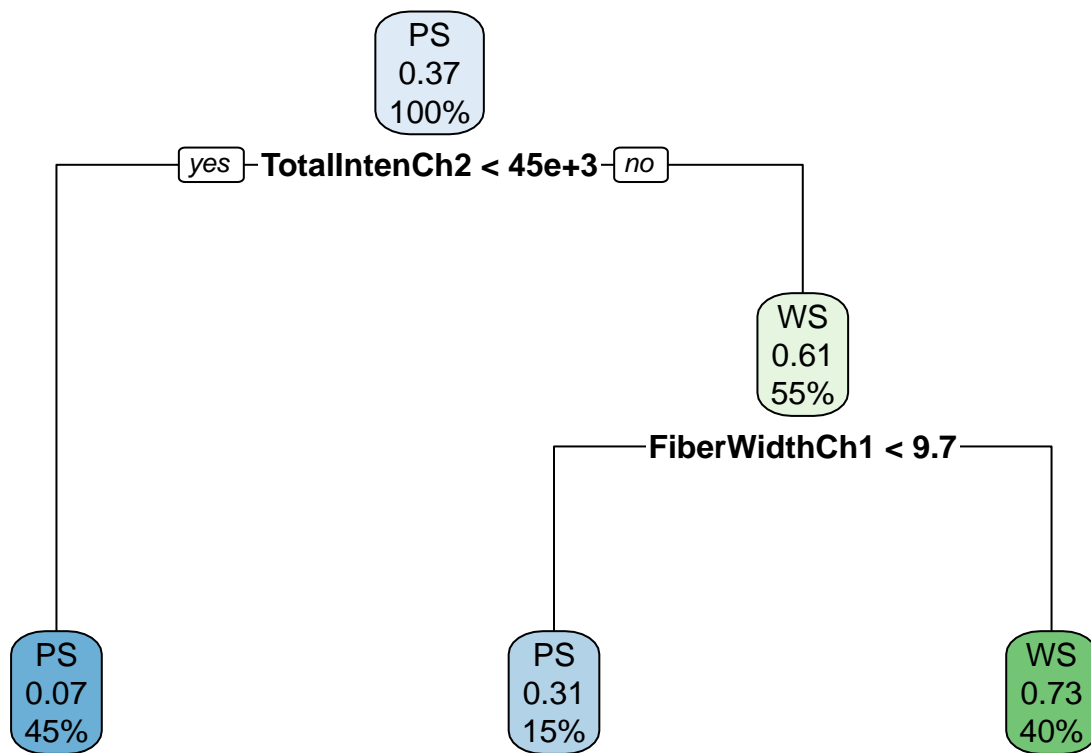
```
library(AppliedPredictiveModeling)
data(segmentationOriginal)
library(caret)
```

1. Subset the data to a training set and testing set based on the Case variable in the data set.
2. Set the seed to 125 and fit a CART model to predict Class with the rpart method using all predictor variables and default caret settings.
3. In the final model what would be the final model prediction for cases with the following variable values:
 - a. TotalIntench2 = 23,000; FiberWidthCh1 = 10; PerimStatusCh1=2
 - b. TotalIntench2 = 50,000; FiberWidthCh1 = 10;VarIntenCh4 = 100
 - c. TotalIntench2 = 57,000; FiberWidthCh1 = 8;VarIntenCh4 = 100
 - d. FiberWidthCh1 = 8;VarIntenCh4 = 100; PerimStatusCh1=2

TIP: Plot the resulting tree and to use the plot to answer this question.

Answer :

```
library(dplyr)
trainS0 <- filter(segmentationOriginal, Case == "Train")
testS0 <- filter(segmentationOriginal, Case == "Test")
set.seed(125)
modfit <- train(Class ~ ., method = "rpart", data = trainS0)
library(rpart.plot)
rpart.plot(modfit$finalModel)
```



Question 2 :

Load the olive oil data using the commands:

```
library(pgmm)
data(olive)
olive = olive[,-1]
```

(NOTE: If you have trouble installing the pgmm package, you can download the `-code-olive-/code-` dataset here: `olive_data.zip`. After unzipping the archive, you can load the file using the `-code-load()-/code-` function in R.)

These data contain information on 572 different Italian olive oils from multiple regions in Italy. Fit a classification tree where `Area` is the outcome variable. Then predict the value of `area` for the following data frame using the `tree` command with all defaults

```
newdata = as.data.frame(t(colMeans(olive)))
```

What is the resulting prediction? Is the resulting prediction strange? Why or why not?

Answer :

```
modfit2 <- train(Area ~., method = "rpart", data = olive)
```

```
## Warning in nominalTrainWorkflow(x = x, y = y, wts = weights, info = 
## trainInfo, : There were missing values in resampled performance measures.
```

```
pred <- predict(modfit2, newdata = as.data.frame(t(colMeans(olive))))
pred
```

```
##          1
## 2.783282
```

Question 4 :

Load the South Africa Heart Disease Data and create training and test sets with the following code:

```
library(ElemStatLearn)
data(SAheart)
set.seed(8484)
train = sample(1:dim(SAheart)[1], size=dim(SAheart)[1]/2, replace=F)
trainSA = SAheart[train,]
testSA = SAheart[-train,]
```

Then set the seed to 13234 and fit a logistic regression model (method="glm", be sure to specify family="binomial") with Coronary Heart Disease (chd) as the outcome and age at onset, current alcohol consumption, obesity levels, cumulative tobacco, type-A behavior, and low density lipoprotein cholesterol as predictors. Calculate the misclassification rate for your model using this function and a prediction on the "response" scale:

```
missClass = function(values, prediction){sum(((prediction > 0.5)*1) != values)/length(values)}
```

What is the misclassification rate on the training set? What is the misclassification rate on the test set?

Answer :

```
set.seed(13234)
modfit3 <- train(chd ~ age + alcohol + obesity + tobacco + typea + ldl, method = "glm", family = "binomial")
missClass(trainSA$chd, predict(modfit3, trainSA))
```

```
## [1] 0.2727273
```

```
missClass(testSA$chd, predict(modfit3, testSA))
```

```
## [1] 0.3116883
```

Question 5 :

Load the vowel.train and vowel.test data sets:

```
library(ElemStatLearn)
data(vowel.train)
data(vowel.test)
```

Set the variable y to be a factor variable in both the training and test set. Then set the seed to 33833. Fit a random forest predictor relating the factor variable y to the remaining variables. Read about variable importance in random forests here: http://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm#ooberr The caret package uses by default the Gini importance.

Calculate the variable importance using the varImp function in the caret package. What is the order of variable importance?

[NOTE: Use randomForest() specifically, not caret, as there's been some issues reported with that approach. 11/6/2016]

Answer :

```
vowel.test$y <- as.factor(vowel.test$y)
vowel.train$y <- as.factor(vowel.train$y)
library(randomForest)
set.seed(33833)
modfit4 <- randomForest(y~. , data = vowel.train)
order(varImp(modfit4), decreasing = TRUE)
```

```
## [1] 2 1 5 6 8 4 9 3 7 10
```