HOME (/)
WHAT WE DO (/SYSTEMATIC-INVESTING-AND-MACHINE-LEARNING/)
INVESTOR LETTERS (/USING-HISTORY-AND-DATA-IN-SYSTEMATIC-INVESTING)
DATA POSTS (/DATA-POSTS/)
ABOUT
LOGIN (/LOGIN/)

June 14, 2015 (/machine-learning-in-practice/2015/6/12/r-caret-and-parameter-tuning-c50)

# R, caret, and Parameter Tuning C5.0 (/machine-learning-in-practice/2015/6/12/r-caret-and-parameter-tuning-c50)

By John Alberg (https://www.linkedin.com/pub/john-alberg/3/9a7/8b0)

Boosted C5.0 classifiers are known to perform well when stacked up against other classifiers (see, for example, this paper (http://jmlr.org/papers/volume15/delgado14a/delgado14a.pdf)).

The c (http://topepo.github.io/caret/index.html)aret library for the R programming language is an exceptional environment for automatic parameter tuning and training of classifiers. However, caret does not allow for out-of-box tuning of C5.0 tree complexity. This post shows how you can customize caret to do just that.

Caret has built in capabilities for tuning the C5.0 meta parameters trials, model, and winnow. The C5.0 documentation (http://cran.r-project.org/web/packages/C50/C50.pdf)describes these parameters in detail. The following code illustrates the ease of tuning and training a C5.0 classifier with a custom tuning grid:

```
1   library(caret)
2   library(C50)
3   library(mlbench)
4
5   fitControl <- trainControl(method = "repeatedcv",
6     number = 10,
7     repeats = 10, returnResamp="all")
8
9   # Choose the features and classes
10  data(PimaIndiansDiabetes2)
11  x <- PimaIndiansDiabetes2[c("age","glucose","insulin","mass","pedigree","pregnant","pressure","triceps")]
12  y <- PimaIndiansDiabetes2$diabetes
13
14  grid <- expand.grid( .winnow = c(TRUE,FALSE), .trials=c(1,5,10,15,20), .model="tree" )
15
16  mdl<- train(x=x,y=y,tuneGrid=grid,trControl=fitControl,method="C5.0",verbose=FALSE)
17
18  mdl
```
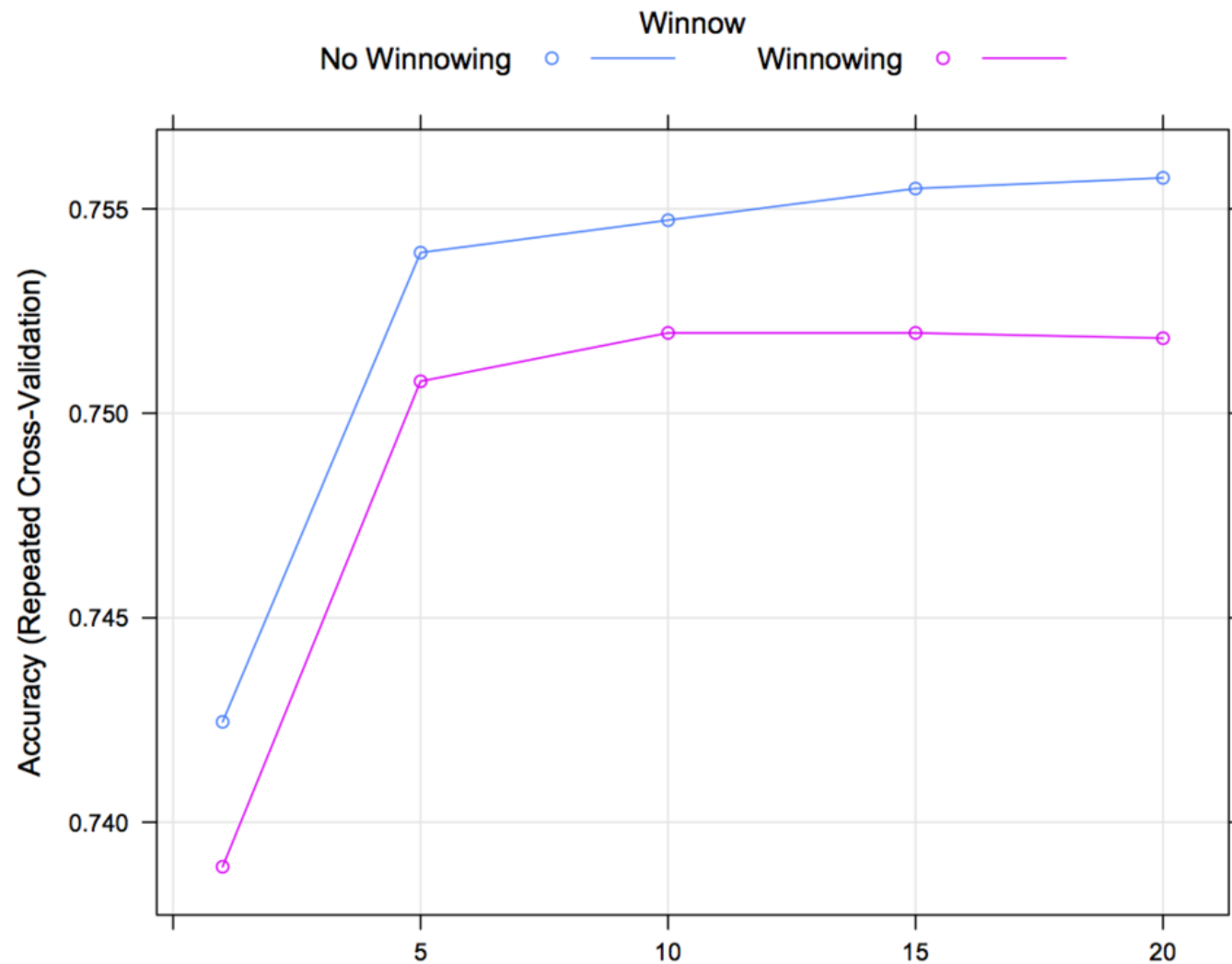
```
19
20    # visualize the resample distributions
21    xyplot(mdl,type = c("g", "p", "smooth"))
```

Then, typing the following

```
> plot(mdl)
```

from the R console will give you this nice chart showing the classifiers performance across the tuning parameters.

# Boosting Iterations

Alas, there is an important C5.0 tuning parameter that is not baked into caret. This parameter is minCases. The minCases parameter specifies the minimum number of cases (training examples) that must be put in at least two of the splits. Essentially it controls the depth of the trees created by C5.0 (depth cannot be controlled directly) and hence it is intimately connected with the resulting tree complexity. The purpose of tuning meta parameters is to find the optimal trade-off between model complexity and the training set size and so minCases is an important parameter to tune. That said, tuning minCases is problematic under cross-validation because the number of cases in the training folds are different than the number of cases in the entire dataset so the optimal value of minCases found in cross validation will not be equal to the true optimum for the entire data set (which the final model will be trained on). To overcome this obstacle we can define minCases as a proportion of the data set size and tune the proportional parameter instead. If we define minCases as

```
minCases <- length(y)/splits
```

then as "splits" increases, so will the depth and complexity of the resulting trees. The code below customizes the standard caret functions to allow for the tuning of "splits" along with the other C5.0 meta parameters.

```
1    library(caret)
2    library(C50)
3    library(mlbench)
4
5    C5CustomSort <- function(x) {
6
7      x$model <- factor(as.character(x$model), levels = c("rules","tree"))
8      x[order(x$trials, x$model, x$splits, !x$winnow),]
9
10   }
11
12   C5CustomLoop <- function (grid)
13   {
14       loop <- ddply(grid, c("model", "winnow","splits"), function(x) c(trials = max(x$trials)))
15       submodels <- vector(mode = "list", length = nrow(loop))
16       for (i in seq(along = loop$trials)) {
17           index <- which(grid$model == loop$model[i] & grid$winnow ==
18               loop$winnow[i] & grid$splits == loop$splits[i])
19           trials <- grid[index, "trials"]
20           submodels[[i]] <- data.frame(trials = trials[trials !=
21               loop$trials[i]])
22       }
23       list(loop = loop, submodels = submodels)
24   }
25
26   C5CustomGrid <- function(x, y, len = NULL) {
27     c5seq <- if(len == 1)  1 else  c(1, 10*((2:min(len, 11)) - 1))
```

```r
    expand.grid(trials = c5seq, splits = c(2,10,20,50), winnow = c(TRUE, FALSE), model = c("tree","rules"))
}

C5CustomFit <- function(x, y, wts, param, lev, last, classProbs, ...) {
  # add the splits parameter to the fit function
  # minCases is a function of splits

  theDots <- list(...)

  splits   <- param$splits
  minCases <- floor( length(y)/splits ) - 1

  if(any(names(theDots) == "control"))
      {
    theDots$control$winnow        <- param$winnow
    theDots$control$minCases      <- minCases
    theDots$control$earlyStopping <- FALSE
  }
  else
    theDots$control <- C5.0Control(winnow = param$winnow, minCases = minCases, earlyStopping=FALSE )

  argList <- list(x = x, y = y, weights = wts, trials = param$trials, rules = param$model == "rules")

  argList <- c(argList, theDots)

  do.call("C5.0.default", argList)

}

GetC5Info <- function() {

  # get the default C5.0 model functions
  c5ModelInfo <- getModelInfo(model = "C5.0", regex = FALSE)[[1]]

  # modify the parameters data frame so that it includes splits
  c5ModelInfo$parameters$parameter <- factor(c5ModelInfo$parameters$parameter,levels=c(levels(c5ModelInfo$parameters$par
  c5ModelInfo$parameters$label <- factor(c5ModelInfo$parameters$label,levels=c(levels(c5ModelInfo$parameters$label),'Spl
  c5ModelInfo$parameters <- rbind(c5ModelInfo$parameters,c('splits','numeric','Splits'))

  # replace the default c5.0 functions with ones that are aware of the splits parameter
  c5ModelInfo$fit  <- C5CustomFit
  c5ModelInfo$loop <- C5CustomLoop
  c5ModelInfo$grid <- C5CustomGrid
  c5ModelInfo$sort <- C5CustomSort

  return (c5ModelInfo)

}

c5info <- GetC5Info()

# Define the structure of cross validation
fitControl <- trainControl(method = "repeatedcv", number = 10,  repeats = 10)
```
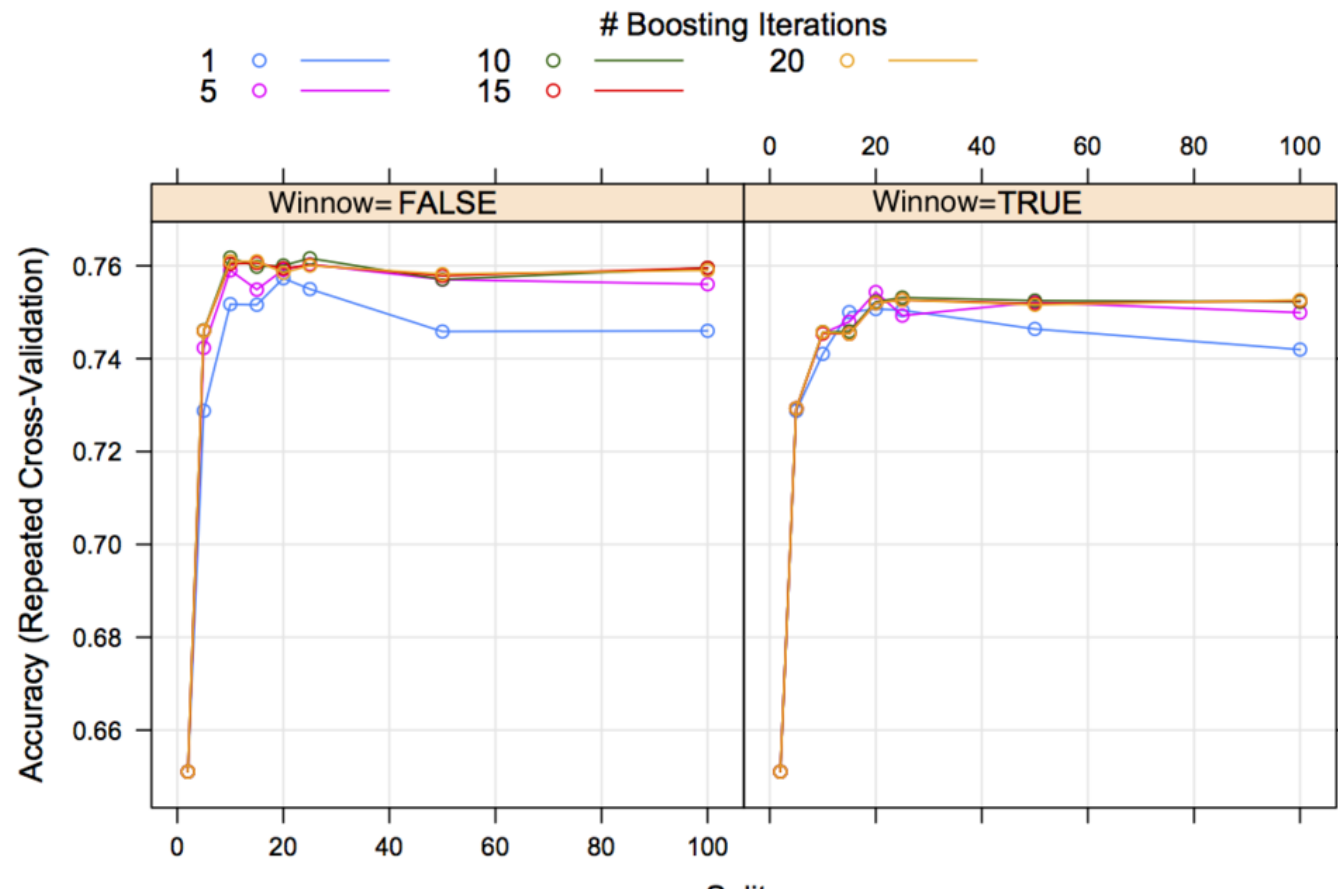
```
81
82    # create a custom cross validation grid
83    grid <- expand.grid( .winnow = c(TRUE,FALSE), .trials=c(1,5,10,15,20), .model=c("tree"), .splits=c(2,5,10,15,20,25,50,10
84
85    # Choose the features and classes
86    data(PimaIndiansDiabetes2)
87    x <- PimaIndiansDiabetes2[c("age","glucose","insulin","mass","pedigree","pregnant","pressure","triceps")]
88    y <- PimaIndiansDiabetes2$diabetes
89
90    # Tune and fit model
91    mdl<- train(x=x,y=y,tuneGrid=grid,trControl=fitControl,method=c5info,verbose=FALSE)
92
93    mdl
```

With this code we can generate cross validated results like those in the following chart:

Tagged: machinelearning (/machine-learning-in-practice/?
tag=machinelearning), rstats (/machine-learning-in-practice/?tag=rstats),
R (/machine-learning-in-practice/?tag=R), caret (/machine-learning-in-
practice/?tag=caret), classification (/machine-learning-in-practice/?
tag=classification)

♥ 0 Likes     ≺ Share

COMMENTS (2)                                          Newest First     Subscribe via e-mail

Preview          **POST COMMENT...**

⬭ **Justfor**  3 months ago

Hello,

running the code ,when reaching mdl <- train line, gives 50 warnings as shown below.
Do you have any idea, why and how to fix it?

1: In predict.C5.0(modelFit, newdata, trial = submodels$trials[j]) :
'trials' should be <= 1 for this object. Predictions generated using 1 trials
.. Total of 50 similar lines; sometimes also trials <= 7 or trials <= 9.

⬭ **John Alberg**  3 months ago

C5.0 has an "earlyStopping" parameter that will cause the algorithm to cease adding trials (trees) if the additional trees give poor performance. I updated the gist so that early stopping is set to FALSE. The default is TRUE. Warnings should be gone now. Note, however, that for the purpose of this code the warnings are ok but I changed for cleanliness of execution.

Newer Post
Class distribution plots for machine learning in R and ggplot2 (/machine-learning-in-practice/2015/6/16/class-distribution-plots-for-machine-learning-in-r-and-ggplot2)

**Follow our firm:** **Read our articles in Advisor Perspectives:  Dec '13 (http://advisorperspectives.com/newsletters13/What_Matters_More_When_Investing.php) Jan '14 (http://advisorperspectives.com/newsletters14/Why_Are_There_Timeless_Lessons_That_Do_Not_Get_A '14 (http://www.advisorperspectives.com/newsletters14/Misunderstanding_Buffett.php)**

(https://www.linkedin.com/company/euclidean-technologies)