## A Short Introduction to the caret Package

## Max Kuhn max.kuhn@pfizer.com

November 9, 2016

The caret package (short for classification and regression training) contains functions to streamline the model training process for complex regression and classification problems. The package utilizes a number of R packages but tries not to load them all at package start-up<sup>1</sup>. The package "suggests" field includes 28 packages. caret loads packages as needed and assumes that they are installed. Install caret using

```
> install.packages("caret", dependencies = c("Depends", "Suggests"))
```

to ensure that all the needed packages are installed.

The **main help pages** for the package are at:

```
http://caret.r-forge.r-project.org/
```

Here, there are extended examples and a large amount of information that previously found in the package vignettes.

caret has several functions that attempt to streamline the model building and evaluation process, as well as feature selection and other techniques.

One of the primary tools in the package is the train function which can be used to

- evaluate, using resampling, the effect of model tuning parameters on performance
- choose the "optimal" model across these parameters
- estimate model performance from a training set

<sup>&</sup>lt;sup>1</sup>By adding formal package dependencies, the package startup time can be greatly decreased

More formally:

There are options for customizing almost every step of this process (e.g. resampling technique, choosing the optimal parameters etc). To demonstrate this function, the Sonar data from the mlbench package will be used.

The Sonar data consist of 208 data points collected on 60 predictors. The goal is to predict the two classes (M for metal cylinder or R for rock).

First, we split the data into two groups: a training set and a test set. To do this, the createDataPartition function is used:

```
> library(caret)
> library(mlbench)
> data(Sonar)
> set.seed(107)
> inTrain <- createDataPartition(y = Sonar$Class,</pre>
                                  ## the outcome data are needed
                                  p = .75,
                                  ## The percentage of data in the
                                  ## training set
                                  list = FALSE)
                                  ## The format of the results
> ## The output is a set of integers for the rows of Sonar
> ## that belong in the training set.
> str(inTrain)
 int [1:157, 1] 1 2 3 6 7 9 10 11 12 13 ...
 - attr(*, "dimnames")=List of 2
  ..$ : NULL
  ..$ : chr "Resample1"
```

By default, createDataPartition does a stratified random split of the data. To partition the data:

```
> training <- Sonar[ inTrain,]
> testing <- Sonar[-inTrain,]
> nrow(training)

[1] 157
> nrow(testing)

[1] 51
```

To tune a model using Algorithm , the train function can be used. More details on this function can be found at:

```
http://caret.r-forge.r-project.org/training.html
```

Here, a partial least squares discriminant analysis (PLSDA) model will be tuned over the number of PLS components that should be retained. The most basic syntax to do this is:

However, we would probably like to customize it in a few ways:

- expand the set of PLS models that the function evaluates. By default, the function will tune over three values of each tuning parameter.
- the type of resampling used. The simple bootstrap is used by default. We will have the function use three repeats of 10–fold cross–validation.
- the methods for measuring performance. If unspecified, overall accuracy and the Kappa statistic are computed. For regression models, root mean squared error and  $R^2$  are computed. Here, the function will be altered to estimate the area under the ROC curve, the sensitivity and specificity

To change the candidate values of the tuning parameter, either of the tuneLength or tuneGrid arguments can be used. The train function can generate a candidate set of parameter values and the tuneLength argument controls how many are evaluated. In the case of PLS, the function uses a sequence of integers from 1 to tuneLength. If we want to evaluate all integers between 1 and 15, setting tuneLength = 15 would achieve this. The tuneGrid argument is used when specific values are desired. A data frame is used where each row is a tuning parameter setting and each column is a tuning parameter. An example is used below to illustrate this.

The syntax for the model would then be:

To modify the resampling method, a trainControl function is used. The option method controls the type of resampling and defaults to "boot". Another method, "repeatedcv", is used to specify repeated K-fold cross-validation (and the argument repeats controls the number of repetitions). K is controlled by the number argument and defaults to 10. The new syntax is then:

Finally, to choose different measures of performance, additional arguments are given to trainControl. The summaryFunction argument is used to pas in a function that takes the observed and predicted values and estimate some measure of performance. Two such functions are already included in the package: defaultSummary and twoClassSummary. The latter will compute measures specific to two—class problems, such as the area under the ROC curve, the sensitivity and specificity. Since the ROC curve is based on the predicted class probabilities (which are not computed automatically), another option is required. The classProbs = TRUE option is used to include these calculations.

Lastly, the function will pick the tuning parameters associated with the best results. Since we are using custom performance measures, the criterion that should be optimized must also be specified. In the call to train, we can use metric = "ROC" to do this.

The final model fit would then be:

Partial Least Squares

```
157 samples
 60 predictor
  2 classes: 'M', 'R'
Pre-processing: centered (60), scaled (60)
Resampling: Cross-Validated (10 fold, repeated 3 times)
Summary of sample sizes: 142, 141, 141, 141, 142, 142, ...
Resampling results across tuning parameters:
        ROC
 ncomp
                    Sens
                               Spec
         0.8135582
                    0.7120370
                               0.7172619
   1
   2
         0.8713211
                    0.7782407
                               0.8130952
   3
         0.8660962 0.7699074
                               0.8339286
   4
         0.8660136 0.7740741
                               0.7714286
   5
         0.8504216
                   0.7532407
                               0.7845238
   6
         0.8352679 0.7574074
                               0.8035714
   7
         0.8093419 0.7296296
                               0.7791667
   8
         0.8152116 0.7259259
                               0.7744048
   9
         0.8194362 0.7291667
                               0.7517857
         0.8254299 0.7296296
  10
                              0.7654762
  11
         0.8303241 0.7416667
                               0.7702381
  12
         0.8275298
                   0.7458333
                               0.7738095
                               0.7744048
  13
         0.8295387
                    0.7337963
  14
         0.8225364 0.7375000
                               0.7833333
  15
         0.8150463
                    0.7337963
                               0.7744048
```

ROC was used to select the optimal model using the largest value. The final value used for the model was ncomp = 2.

In this output the grid of results are the average resampled estimates of performance. The note at the bottom tells the user that PLS components were found to be optimal. Based on this value, a final PLS model is fit to the whole data set using this specification and this is the model that is used to predict future samples.

The package has several functions for visualizing the results. One method for doing this is the plot function for train objects. The command plot(plsFit) produced the results seen in Figure 1 and shows the relationship between the resampled performance values and the number of PLS components.

To predict new samples, predict.train can be used. For classification models, the default behavior is to calculated the predicted class. Using the option type = "prob" can be used to compute class probabilities from the model. For example:

```
> plsClasses <- predict(plsFit, newdata = testing)
> str(plsClasses)
Factor w/ 2 levels "M", "R": 2 1 1 2 1 2 2 2 2 2 ...
```

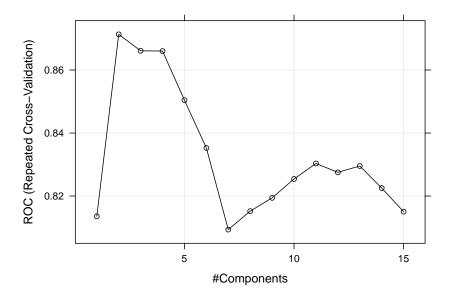


Figure 1: plot(plsFit) shows the relationship between the number of PLS components and the resampled estimate of the area under the ROC curve.

caret contains a function to compute the confusion matrix and associated statistics for the model fit:

```
> confusionMatrix(data = plsClasses, testing$Class)
```

Reference

Confusion Matrix and Statistics

Prediction M R M 20 7 R 7 17

Accuracy: 0.7255

```
95% CI : (0.5826, 0.8411)
No Information Rate : 0.5294
P-Value [Acc > NIR] : 0.003347

Kappa : 0.4491
Mcnemar's Test P-Value : 1.000000

Sensitivity : 0.7407
Specificity : 0.7083
Pos Pred Value : 0.7407
Neg Pred Value : 0.7083
Prevalence : 0.5294
Detection Rate : 0.3922
Detection Prevalence : 0.5294
Balanced Accuracy : 0.7245
```

To fit an another model to the data, train can be invoked with minimal changes. Lists of models available can be found at:

```
http://caret.r-forge.r-project.org/modelList.html
http://caret.r-forge.r-project.org/bytag.html
```

For example, to fit a regularized discriminant model to these data, the following syntax can be used:

```
> ## To illustrate, a custom grid is used
> rdaGrid = data.frame(gamma = (0:4)/4, lambda = 3/4)
> set.seed(123)
> rdaFit <- train(Class ~ .,</pre>
                  data = training,
                  method = "rda",
                  tuneGrid = rdaGrid,
                  trControl = ctrl,
                  metric = "ROC")
> rdaFit
Regularized Discriminant Analysis
157 samples
60 predictor
 2 classes: 'M', 'R'
No pre-processing
Resampling: Cross-Validated (10 fold, repeated 3 times)
Summary of sample sizes: 142, 141, 141, 141, 142, 142, ...
Resampling results across tuning parameters:
```

```
gamma ROC
                   Sens
                             Spec
 0.00
        0.8448826 0.7884259
                             0.7625000
 0.25
        0.8860119 0.8060185
                             0.8035714
 0.50
        0.8851190 0.8097222
                             0.7666667
 0.75
        0.8685847 0.7745370 0.7529762
 1.00
        Tuning parameter 'lambda' was held constant at a value of 0.75
ROC was used to select the optimal model using the largest value.
The final values used for the model were gamma = 0.25 and lambda = 0.75.
> rdaClasses <- predict(rdaFit, newdata = testing)</pre>
> confusionMatrix(rdaClasses, testing$Class)
Confusion Matrix and Statistics
         Reference
Prediction M R
        M 22 5
        R 5 19
              Accuracy: 0.8039
                95% CI: (0.6688, 0.9018)
   No Information Rate: 0.5294
   P-Value [Acc > NIR] : 4.341e-05
                 Kappa: 0.6065
Mcnemar's Test P-Value : 1
           Sensitivity: 0.8148
           Specificity: 0.7917
        Pos Pred Value: 0.8148
        Neg Pred Value: 0.7917
            Prevalence: 0.5294
        Detection Rate: 0.4314
  Detection Prevalence: 0.5294
     Balanced Accuracy: 0.8032
       'Positive' Class : M
```

How do these models compare in terms of their resampling results? The resamples function can be used to collect, summarize and contrast the resampling results. Since the random number seeds were initialized to the same value prior to calling train, the same folds were used for each model. To assemble them:

```
> resamps <- resamples(list(pls = plsFit, rda = rdaFit))
> summary(resamps)
```

```
Call:
summary.resamples(object = resamps)
Models: pls, rda
Number of resamples: 30
ROC
         Min.
                1st Qu.
                           Median
                                        Mean
                                               3rd Qu. Max. NA's
pls 0.5396825 0.8333333 0.8671875 0.8713211 0.9508929
rda 0.6984127 0.8397817 0.9027778 0.8860119 0.9786706
Sens
         Min. 1st Qu.
                                             3rd Qu. Max. NA's
                         Median
                                      Mean
pls 0.3333333  0.7500 0.7777778 0.7782407 0.8750000
rda 0.4444444 0.6875 0.8750000 0.8060185 0.8888889
                                                             0
Spec
                           Median
                                               3rd Qu. Max. NA's
         Min.
                1st Qu.
                                        Mean
pls 0.5000000 0.7142857 0.8571429 0.8130952 0.9687500
                                                          1
rda 0.1428571 0.7232143 0.8571429 0.8035714 0.8571429
```

There are several functions to visualize these results. For example, a Bland-Altman type plot can be created using xyplot(resamps, what = "BlandAltman") (see Figure 2). The results look similar. Since, for each resample, there are paired results a paired t-test can be used to assess whether there is a difference in the average resampled area under the ROC curve. The diff.resamples function can be used to compute this:

```
> diffs <- diff(resamps)</pre>
> summary(diffs)
Call:
summary.diff.resamples(object = diffs)
p-value adjustment: bonferroni
Upper diagonal: estimates of the difference
Lower diagonal: p-value for HO: difference = 0
ROC
    pls
           rda
pls
           -0.01469
rda 0.2975
Sens
    pls
          rda
          -0.02778
pls
rda 0.125
Spec
    pls
           rda
```

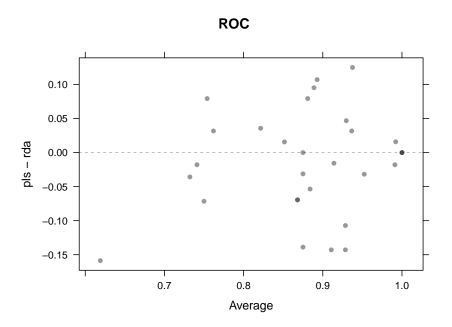


Figure 2: A Bland-Altman plot of the resampled ROC values produced using xyplot(resamps, what = "BlandAltman").

pls 0.009524 rda 0.7348

Based on this analysis, the difference between the models is -0.015 ROC units (the RDA model is slightly higher) and the two–sided p–value for this difference is 0.29749.