

Introduction to Machine Learning in R with caret

Keaton Wilson

8/30/2018

Introduction to Machine Learning in R with caret

Part 1 - What is machine learning? What are the tenets, what is the basic workflow?

Discussion - two questions (5-minutes with the person sitting next to you - then we'll come together and discuss as a group)

1. What is machine learning?
2. How is it different than statistics?

Some important things to know and think about:

1. Prediction is usually more important than explanation
2. Two major types of problems - regression and classification
3. Splitting the data to prevent overfitting

Classification Problem - Wine varietal identifier

Here is the scenario: we've been contacted by a famous vigner in Italy because she suspects that one of the prized varietals (a rare version of *Aglianicone* that her family has grown for 7 generations) from her vinyard has been stolen, and is being grown and sold to make competitively delicious wine in the United States. The competing winemaker claims that the varietal being grown in the US is from a closely related varietal from the same region, that he obtained legally.

Our customer has hired us to develop an algorithm to determine the likelihood that this is the wine being sold by the competitor was made from the varietal grown on her farm. Unfortunately, we don't have fancy genomic data to work with, but she has provided us with chemical profiles of a bunch of different wines made from both her grapes and two varietals that the competitor claims to be working with. The owner of the competing US vinyard has graciously provided us with the same type of data from a bunch of his wines to make comparisons on - he's looking to clear his name (and probably doesn't also believe that an algorithm can predict whether or not a given wine comes from a certain regional varietal)

Part 2 - Examining the Data

```
# Getting libraries we need loaded
library(caret)
library(tidyverse)
```

```

#Reading in the data from the github repo
wine_data = read_csv("https://raw.githubusercontent.com/keatonwilson/classification_workshop_1/master/d
#https://bit.ly/2UYsxdJ

#Overviews
glimpse(wine_data)

## Observations: 178
## Variables: 14
## $ varietal      <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
## $ alcohol       <dbl> 14.23, 13.20, 13.16, 14.37, 13.24, 14.20, 14.3...
## $ malic_acid    <dbl> 1.71, 1.78, 2.36, 1.95, 2.59, 1.76, 1.87, 2.15...
## $ ash           <dbl> 2.43, 2.14, 2.67, 2.50, 2.87, 2.45, 2.45, 2.61...
## $ alkalinity    <dbl> 15.6, 11.2, 18.6, 16.8, 21.0, 15.2, 14.6, 17.6...
## $ magnesium     <dbl> 127, 100, 101, 113, 118, 112, 96, 121, 97, 98,...
## $ total_phenol  <dbl> 2.80, 2.65, 2.80, 3.85, 2.80, 3.27, 2.50, 2.60...
## $ flavanoids    <dbl> 3.06, 2.76, 3.24, 3.49, 2.69, 3.39, 2.52, 2.51...
## $ nonflav_phenols <dbl> 0.28, 0.26, 0.30, 0.24, 0.39, 0.34, 0.30, 0.31...
## $ proantho      <dbl> 2.29, 1.28, 2.81, 2.18, 1.82, 1.97, 1.98, 1.25...
## $ color         <dbl> 5.64, 4.38, 5.68, 7.80, 4.32, 6.75, 5.25, 5.05...
## $ hue           <dbl> 1.04, 1.05, 1.03, 0.86, 1.04, 1.05, 1.02, 1.06...
## $ OD            <dbl> 3.92, 3.40, 3.17, 3.45, 2.93, 2.85, 3.58, 3.58...
## $ proline       <dbl> 1065, 1050, 1185, 1480, 735, 1450, 1290, 1295,...

summary(wine_data)

##      varietal      alcohol      malic_acid      ash
## Min.   :1.000   Min.   :11.03   Min.   :0.740   Min.   :1.360
## 1st Qu.:1.000   1st Qu.:12.36   1st Qu.:1.603   1st Qu.:2.210
## Median :2.000   Median :13.05   Median :1.865   Median :2.360
## Mean   :1.938   Mean   :13.00   Mean   :2.336   Mean   :2.367
## 3rd Qu.:3.000   3rd Qu.:13.68   3rd Qu.:3.083   3rd Qu.:2.560
## Max.   :3.000   Max.   :14.83   Max.   :5.800   Max.   :3.230
##                                     NA's   :1
##      alkalinity      magnesium      total_phenol      flavanoids
## Min.   :10.60   Min.   : 70.00   Min.   :0.980   Min.   :0.340
## 1st Qu.:17.20   1st Qu.: 88.00   1st Qu.:1.742   1st Qu.:1.220
## Median :19.50   Median : 98.00   Median :2.355   Median :2.140
## Mean   :19.49   Mean   : 99.75   Mean   :2.295   Mean   :2.036
## 3rd Qu.:21.50   3rd Qu.:107.00   3rd Qu.:2.800   3rd Qu.:2.880
## Max.   :30.00   Max.   :162.00   Max.   :3.880   Max.   :5.080
##                                     NA's   :1
##      nonflav_phenols      proantho      color      hue
## Min.   :0.1300   Min.   :0.410   Min.   : 1.280   Min.   :0.4800
## 1st Qu.:0.2700   1st Qu.:1.250   1st Qu.: 3.220   1st Qu.:0.7825
## Median :0.3400   Median :1.560   Median : 4.690   Median :0.9650
## Mean   :0.3622   Mean   :1.593   Mean   : 5.058   Mean   :0.9574
## 3rd Qu.:0.4400   3rd Qu.:1.950   3rd Qu.: 6.200   3rd Qu.:1.1200
## Max.   :0.6600   Max.   :3.580   Max.   :13.000   Max.   :1.7100
## NA's   :1       NA's   :1
##      OD      proline
## Min.   :1.270   Min.   : 278.0
## 1st Qu.:1.930   1st Qu.: 500.5

```

```
## Median :2.780   Median : 673.5
## Mean   :2.608   Mean    : 746.9
## 3rd Qu.:3.170   3rd Qu.: 985.0
## Max.   :4.000   Max.    :1680.0
## NA's   :1

#making varietal a factor
wine_data = wine_data %>%
  mutate(varietal = as.factor(varietal))

#Checking for NAs
sum(is.na(wine_data))
```

```
## [1] 6
```

#Uh oh - conversation about missing data.

```
wine_data = wine_data %>%
  drop_na()
```

#or we can deal with it in pre-processing.

Ok, so this looks good. We have our item we want to classify in column 1, and all of our features in the rest. For our varietal numbers, 1 and 2 are the local varietals not owned by our customer, but varietal 3 is her special grape. So we're looking for the presence of any wines made from varietal 3 in the test set.

What do we need to do before we jump into trying to build some algorithms?

Preprocess! In particular, we need to center and scale the data. *caret* can do this for us.

Part 3 - Preprocessing

```
#Setting up the preprocessing algorithm
set.seed(42)

#Train and test split
trainindex = createDataPartition(wine_data$varietal,
                                  p = 0.8,
                                  list = FALSE,
                                  times = 1)

wine_train = wine_data[trainindex,]
wine_test = wine_data[-trainindex,]

#We could also do some sort of imputation here.
pp = preProcess(wine_data[, -1], method = c("center", "scale"), outcome = wine_data$varietal)

wine_train_pp = predict(pp, wine_train)
wine_train_pp
```

```
## # A tibble: 140 x 14
##   varietal alcohol malic_acid ash alkalinity magnesium total_phenol
##   <fct>      <dbl>      <dbl> <dbl>      <dbl>      <dbl>      <dbl>
## 1 1          1.50      -0.558 0.229      -1.15       1.91       0.778
## 2 1          0.183       0.0255 1.10       -0.252      0.0959     0.778
## 3 1          0.281       0.232 1.82        0.466      1.28       0.778
```

```
## 4 1      1.46      -0.513  0.302      -1.27      0.865      1.53
## 5 1      1.69      -0.414  0.302      -1.45     -0.254      0.299
## 6 1      1.29      -0.163  0.881      -0.551     1.49      0.459
## 7 1      1.34      -0.154 -0.241      -0.431     0.376     1.02
## 8 1      1.36      -0.764 -0.169      -0.790    -0.324    -0.180
## 9 1      0.907     -0.540  0.157      -1.03    -0.743     0.459
## 10 1     2.13      -0.540  0.0846     -2.41    -0.603     1.26
## # ... with 130 more rows, and 7 more variables: flavanoids <dbl>,
## #   nonflav_phenols <dbl>, proantho <dbl>, color <dbl>, hue <dbl>,
## #   OD <dbl>, proline <dbl>

#We also need to add this same processing algorithm to the test data.
wine_test_pp = predict(pp, wine_test)
```

Part 4 - Model Testing and Tuning

There are a ton of classification models to choose from - when starting ML stuff, this can be a really daunting part of the thing. Today, we're going to explore a couple of bread-and-butter models:

1. k-nn - nearest neighbor classifier
2. Naive Bayes
3. Decision Trees
4. Support Vector Machines

I'm not going to go into the math of how all of these operate at all. It's beyond the scope of this workshop, but here is a good overview: <https://medium.com/@sifium/machine-learning-types-of-classification-9497bd4f2e14>

One thing that we need to talk about briefly is resampling - this is the method we're going to use to assess how 'good' a model is, without applying it to the test data. There are a couple of main ways to do this:

1. bootstrapping - random sampling within the dataset with replacement. Pulling a bunch of subsets of the data and looking at how the model performs across these subsets.
2. Repeated n-fold cross-validation - does a bunch of splitting into training and test data **within** the training set, and then averages accuracy or RMSE across all these little mini-sets.

We're going to use the second type.

```
# setting up the control object to feed to all of the subsequent models
fit_control = trainControl(method = "repeatedcv", number = 5, repeats = 5)

#models
#knn
knn_model = train(varietal ~ ., data = wine_train_pp,
                  method = "knn", trControl = fit_control)

#naive bays
bayes_model = train(varietal ~ ., data = wine_train_pp,
                   method = "nb", trControl = fit_control)

#CART
cart_model = train(varietal ~ ., data = wine_train_pp,
                  method = "rpart", trControl = fit_control)

#svm
svm_model = train(varietal ~ ., data = wine_train_pp,
                 method = "svmRadial", trControl = fit_control)
```

The models default to using accuracy as the score to determine how good they are. Accuracy is the

percentage of the predictions made by the model that are correct.

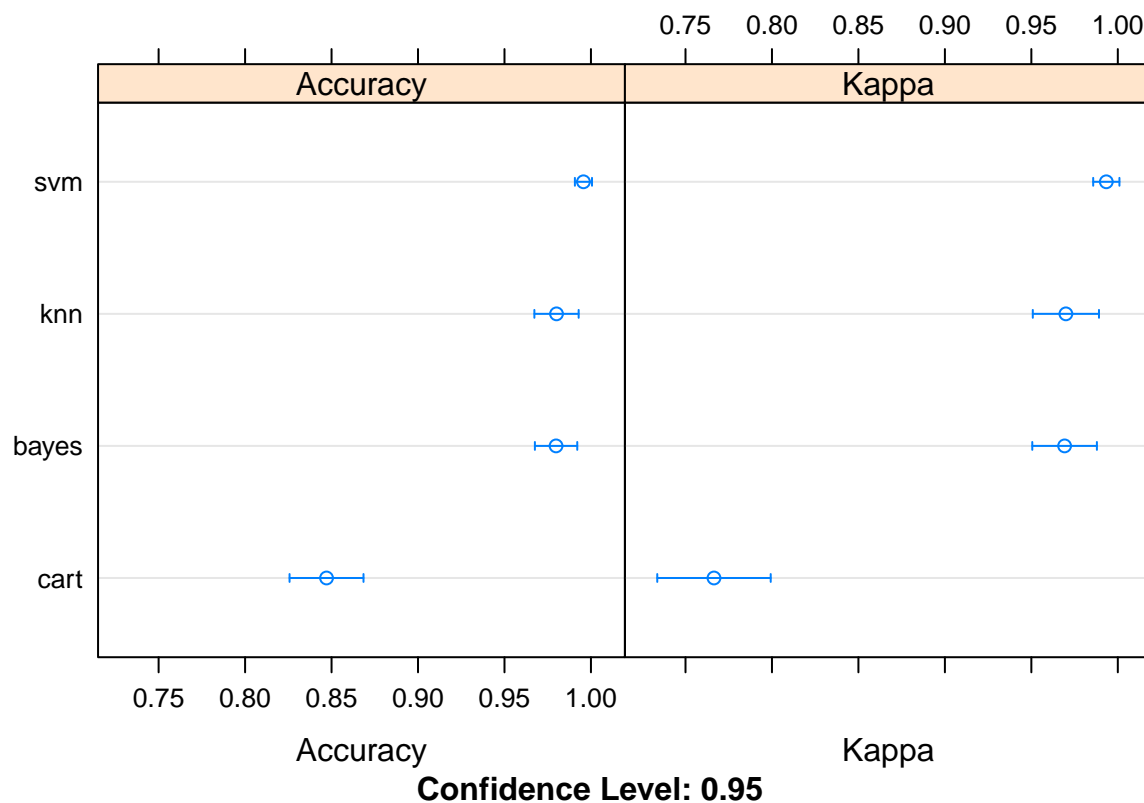
Kappa compares observed accuracy with expected accuracy (chance). Less misleading than simply accuracy.

Let's compare models

```
results = resamples(list(knn = knn_model, bayes = bayes_model, cart = cart_model, svm = svm_model))
summary(results)
```

```
##
## Call:
## summary.resamples(object = results)
##
## Models: knn, bayes, cart, svm
## Number of resamples: 25
##
## Accuracy
##      Min.    1st Qu.    Median      Mean   3rd Qu.      Max. NA's
## knn  0.8571429 0.9642857 1.0000000 0.9800912 1.0000000 1.0000000    0
## bayes 0.9259259 0.9642857 1.0000000 0.9798340 1.0000000 1.0000000    0
## cart  0.7407407 0.8214286 0.8571429 0.8470929 0.8888889 0.9642857    0
## svm   0.9629630 1.0000000 1.0000000 0.9956614 1.0000000 1.0000000    0
##
## Kappa
##      Min.    1st Qu.    Median      Mean   3rd Qu.      Max. NA's
## knn  0.7874763 0.9459459 1.0000000 0.9699715 1.0000000 1.0000000    0
## bayes 0.8860759 0.9447732 1.0000000 0.9692260 1.0000000 1.0000000    0
## cart  0.6078838 0.7233202 0.7821012 0.7664807 0.8298319 0.9456311    0
## svm   0.9432773 1.0000000 1.0000000 0.9933518 1.0000000 1.0000000    0
```

```
dotplot(results)
```



```
svm_model
```

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 140 samples
## 13 predictor
## 3 classes: '1', '2', '3'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold, repeated 5 times)
## Summary of sample sizes: 111, 112, 113, 112, 112, ...
## Resampling results across tuning parameters:
##
##  C      Accuracy  Kappa
##  0.25  0.9799927  0.9693369
##  0.50  0.9886170  0.9826210
##  1.00  0.9956614  0.9933518
##
## Tuning parameter 'sigma' was held constant at a value of 0.08359607
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.08359607 and C = 1.
```

We can also look at the relative amount of false negatives and false positives with a confusion matrix.

```
predictions = predict(svm_model, wine_train_pp)
confusionMatrix(predictions, wine_train_pp$varietal)
```

```
## Confusion Matrix and Statistics
##
##           Reference
```

```
## Prediction  1  2  3
##           1 48  0  0
##           2  0 56  0
##           3  0  0 36
##
## Overall Statistics
##
##           Accuracy : 1
##           95% CI : (0.974, 1)
##           No Information Rate : 0.4
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 1
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: 1 Class: 2 Class: 3
## Sensitivity      1.0000      1.0   1.0000
## Specificity      1.0000      1.0   1.0000
## Pos Pred Value   1.0000      1.0   1.0000
## Neg Pred Value   1.0000      1.0   1.0000
## Prevalence       0.3429      0.4   0.2571
## Detection Rate   0.3429      0.4   0.2571
## Detection Prevalence 0.3429      0.4   0.2571
## Balanced Accuracy 1.0000      1.0   1.0000
```

Not particularly informative, given the model did a 100% accurate job of predicting on the training-data, but you get the gist.

Part 5 - Using the model on out of sample test data.

```
wine_test_pp$pred = predict(svm_model, wine_test_pp)

wine_test_pp %>%
  select(varietal, pred) %>%
  print(n = 32)
```

```
## # A tibble: 32 x 2
##   varietal pred
##   <fct>    <fct>
## 1 1      1
## 2 1      1
## 3 1      1
## 4 1      1
## 5 1      1
## 6 1      1
## 7 1      1
## 8 1      1
## 9 1      1
## 10 1     1
## 11 1     1
```

```
## 12 2      2
## 13 2      2
## 14 2      2
## 15 2      3
## 16 2      2
## 17 2      2
## 18 2      2
## 19 2      2
## 20 2      2
## 21 2      2
## 22 2      2
## 23 2      2
## 24 2      2
## 25 3      3
## 26 3      3
## 27 3      3
## 28 3      3
## 29 3      3
## 30 3      3
## 31 3      3
## 32 3      3
```

```
confusionMatrix(wine_test_pp$pred, wine_test_pp$varietal)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction  1  2  3
```

```
##           1 11  0  0
```

```
##           2  0 12  0
```

```
##           3  0  1  8
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##           Accuracy : 0.9688
```

```
##           95% CI : (0.8378, 0.9992)
```

```
## No Information Rate : 0.4062
```

```
## P-Value [Acc > NIR] : 1.447e-11
```

```
##
```

```
##           Kappa : 0.9526
```

```
##
```

```
## McNemar's Test P-Value : NA
```

```
##
```

```
## Statistics by Class:
```

```
##
```

```
##           Class: 1 Class: 2 Class: 3
```

```
## Sensitivity      1.0000  0.9231  1.0000
```

```
## Specificity      1.0000  1.0000  0.9583
```

```
## Pos Pred Value   1.0000  1.0000  0.8889
```

```
## Neg Pred Value   1.0000  0.9500  1.0000
```

```
## Prevalence       0.3438  0.4062  0.2500
```

```
## Detection Rate   0.3438  0.3750  0.2500
```

```
## Detection Prevalence 0.3438  0.3750  0.2812
```

```
## Balanced Accuracy 1.0000  0.9615  0.9792
```


#we can even look at that wine

```
wine_test_pp %>%  
  filter(pred != varietal)
```

```
## # A tibble: 1 x 15  
##   varietal alcohol malic_acid   ash alkalinity magnesium total_phenol  
##   <fct>      <dbl>      <dbl> <dbl>      <dbl>      <dbl>      <dbl>  
## 1 2          0.0482        1.37 -0.169      0.915      -1.02      -1.06  
## # ... with 8 more variables: flavanoids <dbl>, nonflav_phenols <dbl>,  
## #   proantho <dbl>, color <dbl>, hue <dbl>, OD <dbl>, proline <dbl>,  
## #   pred <fct>
```

Part 6 - Continuing Practice

Some resources if you want to get better at this:

1. Kaggle - an online community of data scientists - lots of cool datasets to play with, and competitions!
2. https://kbroman.org/pkg_primer/pages/resources.html - great list of resources!
3. Machine Learning with R - Brett Lantz. Great book!