

assignment2

Lilach Herzog & Leslie Cohen

13 5 2022

Our whole project and code can be found on github at <https://github.com/LeslieLebon/assignment-2.git>

Our objective in this project was to analyze the “wheat-seeds” dataset (where the ‘Type’ of the seed is the column used for labeling and classifying) using two classifier models that we built. We decided on using the KNN algorithm and the decision tree (DT) algorithm.

preprocessing

Upload data and libraries

```
library(gmodels)
library(C50)
library(class)
library(tidymodels) # for the rsample package, along with the rest of tidymodels
library(modeldata) # for the cells data
library(vip) # for variable importance plots
library(pROC)

seeds <- read.csv("seeds.csv") # read data
set.seed(1234) #to initialize a pseudo random number generator.
```

initial look at the data

We know the class ‘type’ is the column (label) to classify according to. We took a quick look at the data:

```
str(seeds)

## 'data.frame': 199 obs. of 8 variables:
## $ Area : num 15.3 14.9 14.3 13.8 16.1 ...
## $ Perimeter : num 14.8 14.6 14.1 13.9 15 ...
## $ Compactness : num 0.871 0.881 0.905 0.895 0.903 ...
## $ Kernel.Length : num 5.76 5.55 5.29 5.32 5.66 ...
## $ Kernel.Width : num 3.31 3.33 3.34 3.38 3.56 ...
## $ Asymmetry.Coeff: num 2.22 1.02 2.7 2.26 1.35 ...
## $ Kernel.Groove : num 5.22 4.96 4.83 4.8 5.17 ...
## $ Type : int 1 1 1 1 1 1 1 1 1 1 ...
```

There are 3 types of seeds (in the column ‘Type’): 1, 2 and 3.

Since the whole table is sorted according to the ‘type’ column, and we want to work randomly, we will shuffle the whole table

In order to further understand the data and be able to work with it we used the “table” function. this function builds a contingency table of the counts at each combination of factor levels- what interests us is the

Type:

```
table(seeds$Type)
```

```
##
##  1  2  3
## 66 68 65
```

We can see the Type has 66 times type 1, 68 times type 2, and 65 times type 3, about a third of each type.

Before starting on teaching the algorithm, we converted the 'Type' column to a factor, as that is what is required by the C50 package.

```
seeds$Type<-as.factor(seeds$Type)
```

Now- let's get to work!!

split into training and test sets

```
length_of_training_set <-round(0.8*(nrow(seeds)),0) # 0.8% of the observations out of the total 199
seeds_train <- seeds[1:length_of_training_set, ]
seeds_test <- seeds[(length_of_training_set+1):nrow(seeds), ]
seeds_train_labels <-seeds_train[, 8]
seeds_test_labels <-seeds_test[, 8]
```

Since we know that there are about a third of each type in the data set, we want to be sure that the distribution is as we defined (about 1/3 of each type in of both training and test sets):

```
##
##      1      2      3
## 0.3081761 0.3333333 0.3584906
##
##      1      2      3
## 0.425 0.375 0.200
```

target factor (label) vector for classification

Decision tree

A Decision Tree uses a serie of questions that lead to a decision- a class. The goal is to improve purity, to try and have a higher and higher percentage of items in the node are similar, until all items in a node belong to the same class. It does so by taking the items in a node and splitting them in such a manner that they become two groups with improved purity in each (more samples in the group belong to the same class).

We will use the C5.0 method. The 8th column of the dataset is the type class variable, so we need to exclude it from the training data frame, and supply it as the target factor (label) vector for classification:

```
seeds_model <- C5.0(seeds_train[-8], seeds_train$Type)
```

Explanation of initial tree:

```
##
## Call:
## C5.0.default(x = seeds_train[-8], y = seeds_train$Type)
##
##
## C5.0 [Release 2.07 GPL Edition]      Mon Jun 13 12:45:29 2022
```

```

## -----
##
## Class specified by attribute `outcome'
##
## Read 159 cases (8 attributes) from undefined.data
##
## Decision tree:
##
## Kernel.Groove > 5.533: 2 (54/1)
## Kernel.Groove <= 5.533:
##   ...Area <= 12.7:
##     ...Kernel.Groove <= 4.703: 1 (3)
##     : Kernel.Groove > 4.703: 3 (52/1)
##     Area > 12.7:
##       ...Asymmetry.Coeff <= 4.773: 1 (43/1)
##       Asymmetry.Coeff > 4.773:
##         ...Kernel.Width <= 3.201: 3 (5)
##         Kernel.Width > 3.201: 1 (2)
##
##
## Evaluation on training data (159 cases):
##
##      Decision Tree
##      -----
##      Size      Errors
##
##          6      3( 1.9%)  <<
##
##
##      (a)   (b)   (c)   <-classified as
##      ----  ----  ----
##          47      1      1   (a): class 1
##              53             (b): class 2
##          1             56   (c): class 3
##
##
## Attribute usage:
##
## 100.00% Kernel.Groove
##  66.04% Area
##  31.45% Asymmetry.Coeff
##   4.40% Kernel.Width
##
##
## Time: 0.0 secs

```

The algorithm did six iterations. In the first iteration it divided all the seeds into two subgroups (nodes) according to whether their Groove was higher or lower than 5.533. One node was pure (all seeds in the node were Class 2), and the other node went through another iteration according to the Area, where both nodes were not pure and needed to go through another iteration. One node was split into two pure nodes-again according to Groove-which yielded 2 pure nodes (Class 3 or class 1). The other node was split according to the asymmetry coefficient which resulted in one pure node (of Class 1) whereas the other node went into another iteration according to the width which yielded 2 pure nodes (classes 3 and 1).

The percentage of error is relatively low: 1.3%. The algorithm uses the feature that provides the most

Information Gain in every split. Since it used 100.00% of the information from the Kernel.Groove column, we can assume that that column provided the most information gain (at least in the first iteration). The algorithm also used 66.04% of the Area column; 31.45% of the Asymmetry.Coeff; and 4.40% of the Kernel.Width column, so that also gives us some information about how much these features influenced the Information Gain function.

After teaching the model on the training set, we want to use it to predict the class is on the test set, and then compare the test set to what we predicted for the test said using our predictive model.

```
# apply model on test data
seeds_pred <- predict(seeds_model, seeds_test)
CrossTable(seeds_test$Type, seeds_pred, prop.chisq = FALSE, prop.c = FALSE, prop.r = FALSE, dnn = c('actual type', 'predicted type'))
```

```
##
##
##      Cell Contents
## |-----|
## |                      N |
## |      N / Table Total |
## |-----|
##
##
## Total Observations in Table:  40
##
##
##      | predicted type
## actual type |          1 |          2 |          3 | Row Total |
## -----|-----|-----|-----|-----|
##          1 |          17 |          0 |          0 |          17 |
##          |          0.425 |          0.000 |          0.000 |
## -----|-----|-----|-----|
##          2 |          1 |          14 |          0 |          15 |
##          |          0.025 |          0.350 |          0.000 |
## -----|-----|-----|-----|
##          3 |          2 |          0 |          6 |          8 |
##          |          0.050 |          0.000 |          0.150 |
## -----|-----|-----|-----|
## Column Total |          20 |          14 |          6 |          40 |
## -----|-----|-----|-----|
##
##
```

Type seeds 2 and 3 both have no false positives (all predictions of type two or three were actually typed two or 3) with a few false negatives (one type 2 and 2 type 3s that were predicted as type one). Type 1 on the other hand has no false negative (all type 1 seeds were predicted as type 1) but it had a few false positives (as mentioned before- one that was actually type 2 and two were actually type 3).

We will try to improve our results, using adaptive boosting to have a lower percentage of error. This is a process in which many decision trees are built and the trees vote on the best class for each example.

Adaptive Boosting

We'll start with 10 trials, a number that has become the de facto standard, as research suggests that this reduces error rates on test data by about 25 percent:

The classifier made 1 mistake for an error rate of 0.6% percent. This is an improvement over the previous training error rate before adding boosting! However, it remains to be seen whether we see a similar

improvement on the test data. Let's take a look:

```
# boosting on test data
seeds_boost_pred10 <- predict(seeds_boost10, seeds_test)
```

```
CrossTable(seeds_test$Type, seeds_boost_pred10,
prop.chisq = FALSE, prop.c = FALSE, prop.r = FALSE,
dnn = c('actual Type', 'predicted Type'))
```

```
##
##
##      Cell Contents
## |-----|
## |                      N |
## |      N / Table Total |
## |-----|
##
##
## Total Observations in Table:  40
##
##
##      | predicted Type
## actual Type |          1 |          2 |          3 | Row Total |
## -----|-----|-----|-----|-----|
##          1 |          17 |          0 |          0 |          17 |
##          |          0.425 |          0.000 |          0.000 |          |
## -----|-----|-----|-----|-----|
##          2 |          1 |          14 |          0 |          15 |
##          |          0.025 |          0.350 |          0.000 |          |
## -----|-----|-----|-----|-----|
##          3 |          1 |          0 |          7 |          8 |
##          |          0.025 |          0.000 |          0.175 |          |
## -----|-----|-----|-----|-----|
## Column Total |          19 |          14 |          7 |          40 |
## -----|-----|-----|-----|-----|
##
##
```

The model made 2 errors instead of 3 before the boost, on seeds of type 1. WE will try to improve the result, with more trials than 10, for example 15.

```
# boosting with 20 trials (on training)
seeds_boost20 <- C5.0(seeds_train[-8], seeds_train$Type, trials = 20)
summary(seeds_boost20 )
```

```
# boosting on test data
seeds_boost_pred20 <- predict(seeds_boost20, seeds_test)
CrossTable(seeds_test$Type, seeds_boost_pred20,
prop.chisq = FALSE, prop.c = FALSE, prop.r = FALSE,
dnn = c('actual Type', 'predicted Type'))
```

With 10 trials, we obtained the best result. Indeed, choosing 20 trials in the boost does not improve the result. The error rate stay similar than with 10 trials.

DT conclusion

Then we should write our conclusion about the different predictions obtained. Compare these results to the boosted model; this version makes more mistakes overall, but the types of mistakes are very different. Where the previous models incorrectly classified a small number of defaults correctly, our weighted model has does much better in this regard. This trade resulting in a reduction of false negatives at the expense of increasing false positives may be acceptable if our cost estimates were accurate. To create our decision trees in this practice we used the C5.0 package.

KNN Algorithm

In the KNN Algorithm, k is a user-defined constant, and an unlabeled vector (a query or test point) is classified by assigning the label which is most frequent among the k training samples nearest to that query point. *## Normalization* In order to compare the distances of different features with different scales, we made sure the range of values of each parameter is similar by using a simple normalization function (min/max normalization):

and to test our function:

We used the normalization function on our training and test sets, and made sure it worked:

```
normalized_seeds_train <- as.data.frame(lapply(seeds_train[1:7], normalize))
normalized_seeds_test <- as.data.frame(lapply(seeds_test[1:7], normalize))
summary(normalized_seeds_train$Area)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.0000  0.1331  0.3228  0.3889  0.6470  1.0000
```

```
summary(normalized_seeds_test$Area)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.0000  0.3613  0.5135  0.5411  0.7272  1.0000
```

run the KNN algorithm

Now, we will run the KNN algorithm and We will start with K of 3, because we have 3 types of seeds

And finally, we are ready to run our algorithm. We'll start with K of 21, the square root of 455 and also an odd number, reducing the change of a tie vote.

```
seeds_test_pred <- knn(train = normalized_seeds_train, test = normalized_seeds_test, cl = seeds_train_1
```

Since the knn function returns a factor vector of the predicted values, we'll compare that vector with the true labels we saved in advance. We'll do the comparison with the CrossTable function, from the gmodels package loaded:

```
# it is not running
CrossTable(x = seeds_test_labels, y = seeds_test_pred, prop.chisq=FALSE)
```

```
##
##
##      Cell Contents
## |-----|
## |                      N |
## |      N / Row Total |
## |      N / Col Total |
## |      N / Table Total |
## |-----|
```

```
##
##
## Total Observations in Table:  40
##
##
## seeds_test_labels | seeds_test_pred
## seeds_test_labels |          1 |          2 |          3 | Row Total |
## -----|-----|-----|-----|-----|
##          1 |          17 |          0 |          0 |          17 |
##          |          1.000 |          0.000 |          0.000 |          0.425 |
##          |          0.850 |          0.000 |          0.000 |          |
##          |          0.425 |          0.000 |          0.000 |          |
## -----|-----|-----|-----|-----|
##          2 |          1 |          14 |          0 |          15 |
##          |          0.067 |          0.933 |          0.000 |          0.375 |
##          |          0.050 |          1.000 |          0.000 |          |
##          |          0.025 |          0.350 |          0.000 |          |
## -----|-----|-----|-----|-----|
##          3 |          2 |          0 |          6 |          8 |
##          |          0.250 |          0.000 |          0.750 |          0.200 |
##          |          0.100 |          0.000 |          1.000 |          |
##          |          0.050 |          0.000 |          0.150 |          |
## -----|-----|-----|-----|-----|
##      Column Total |          20 |          14 |          6 |          40 |
##          |          0.500 |          0.350 |          0.150 |          |
## -----|-----|-----|-----|-----|
##
##
```

In the above table we see that all 17 seeds(100%) of type 1 were predicted correctly, but it has a few false positives (3 seeds that were wrongly predicted as Type one: one seed was actually type 2 and 2 were type 3). on the other hand- types two and three which had no false positives (all seeds that were predicted as two or three were actually two or 3 respectively) and only a few false negatives (14/15 type 2 seeds (93.3%) and 6/8 type 3 seeds (75%) were identified correctly). these are really good results!! Lets see if we can improve our results. First, lets try z-score standardization instead of normalization:

```
seeds_train_z <- as.data.frame(scale(seeds_train[1:7]))
seeds_test_z <- as.data.frame(scale(seeds_test[1:7]))
seeds_test_pred_z <- knn(train = seeds_train_z, test = seeds_test_z, cl = seeds_train_labels, k=3)
CrossTable(x = seeds_test_labels, y = seeds_test_pred_z, prop.chisq=FALSE)
```

Not much of an improvement, maybe even worse. We tried several more times with different K's:

```
seeds_test_pred_k3 <- knn(train = normalized_seeds_train, test = normalized_seeds_test, cl = seeds_train_labels, k=3)
CrossTable(x = seeds_test_labels, y = seeds_test_pred_k3, prop.chisq=FALSE)

seeds_test_pred_k9 <- knn(train = normalized_seeds_train, test = normalized_seeds_test, cl = seeds_train_labels, k=9)
CrossTable(x = seeds_test_labels, y = seeds_test_pred_k9, prop.chisq=FALSE)

seeds_test_pred_k15 <- knn(train = normalized_seeds_train, test = normalized_seeds_test, cl = seeds_train_labels, k=15)
CrossTable(x = seeds_test_labels, y = seeds_test_pred_k15, prop.chisq=FALSE)

seeds_test_pred_k21 <- knn(train = normalized_seeds_train, test = normalized_seeds_test, cl = seeds_train_labels, k=21)
CrossTable(x = seeds_test_labels, y = seeds_test_pred_k21, prop.chisq=FALSE)
```

It looks like the results are exactly the same, except for K = 9 we got 7/8 type 3 seeds (87.5%) that were

predicted correctly and only once predicted falsely as one (which improved both the true positive on type 3 and the false negative on type one).

feature selection

The DT algorithm uses the feature that provides the most Information Gain in every split. Using that logic, the feature that is used the least provides the least information gain, and if it doesn't add much information it might disturb as noise. Thus, since we know that the accuracy of the k-NN algorithm can be severely degraded by the presence of noisy or irrelevant features, we tried the KNN algorithm once more, with K=9 and without the "Kernel.Length" column (which attributed the least to the building of the DT).

```
# z-score normalization
feature_selected_train <- seeds_train[-4]
feature_selected_test <- seeds_test[-4]
feature_selected_test_pred <- knn(train = seeds_train_z, test = seeds_test_z, cl = seeds_train_labels,
CrossTable(x = seeds_test_labels, y = feature_selected_test_pred, prop.chisq=FALSE)
```

```
##
##
##      Cell Contents
## |-----|
## |                N |
## |      N / Row Total |
## |      N / Col Total |
## |      N / Table Total |
## |-----|
##
##
## Total Observations in Table:  40
##
##
##      | feature_selected_test_pred
## seeds_test_labels |      1 |      2 |      3 | Row Total |
## -----|-----|-----|-----|-----|
##           1 |      17 |       0 |       0 |      17 |
##           |      1.000 |      0.000 |      0.000 |      0.425 |
##           |      0.944 |      0.000 |      0.000 |      |
##           |      0.425 |      0.000 |      0.000 |      |
## -----|-----|-----|-----|-----|
##           2 |       1 |      14 |       0 |      15 |
##           |      0.067 |      0.933 |      0.000 |      0.375 |
##           |      0.056 |      1.000 |      0.000 |      |
##           |      0.025 |      0.350 |      0.000 |      |
## -----|-----|-----|-----|-----|
##           3 |       0 |       0 |       8 |       8 |
##           |      0.000 |      0.000 |      1.000 |      0.200 |
##           |      0.000 |      0.000 |      1.000 |      |
##           |      0.000 |      0.000 |      0.200 |      |
## -----|-----|-----|-----|-----|
##      Column Total |      18 |      14 |       8 |      40 |
##           |      0.450 |      0.350 |      0.200 |      |
## -----|-----|-----|-----|-----|
##
##
```

We managed to improve the true positive of type three to 100%! We have only one mistake left (a type two

seed that was predicted as type 1). `## KNN conclusions` We managed to classify our data using the KNN algorithm. We saw that the best results were with $K=9$ and without the “Kernel.Length” column, and we managed to classify correctly all seeds except for one.

total conclusions