

Letter Recognition

Contents

Introduction	1
Exercices	2
1. <i>Predicting B or not B</i>	2
Problem 1.1	2
Problem 1.2	4
Problem 1.3	4
2. <i>Predicting the letters A, B, P, R</i>	5
Problem 2.1	5
Problem 2.2	5
Problem 2.3	6

Introduction

One of the earliest applications of the predictive analytics methods we have studied so far in this class was to automatically recognize letters, which post office machines use to sort mail. In this problem, we will build a model that uses statistics of images of four letters in the Roman alphabet – A, B, P, and R – to predict which letter a particular image corresponds to.

Note that this is a **multiclass classification problem**. We have mostly focused on binary classification problems (e.g., predicting whether an individual voted or not, whether the Supreme Court will affirm or reverse a case, whether or not a person is at risk for a certain disease, etc.). In this problem, we have more than two classifications that are possible for each observation, like in the D2Hawkeye lecture.

The file `letters_ABPR.csv` contains 3116 observations, each of which corresponds to a certain image of one of the four letters A, B, P and R. The images came from 20 different fonts, which were then randomly distorted to produce the final images; each such distorted image is represented as a collection of pixels, each of which is “on” or “off”. For each such distorted image, we have available certain statistics of the image in terms of these pixels, as well as which of the four letters the image is. This data comes from the UCI Machine Learning Repository.

This dataset contains the following 17 variables:

letter : the letter that the image corresponds to (A, B, P or R)
xbox : the horizontal position of where the smallest box covering the letter shape begins.
ybox : the vertical position of where the smallest box covering the letter shape begins.
width : the width of this smallest box.
height : the height of this smallest box.
onpix : the total number of “on” pixels in the character image

xbar : the mean horizontal position of all of the “on” pixels
ybar : the mean vertical position of all of the “on” pixels
x2bar : the mean squared horizontal position of all of the “on” pixels in the image
y2bar : the mean squared vertical position of all of the “on” pixels in the image
xybar : the mean of the product of the horizontal and vertical position of all of the “on” pixels in the image
x2ybar : the mean of the product of the squared horizontal position and the vertical position of all of the “on” pixels
xy2bar : the mean of the product of the horizontal position and the squared vertical position of all of the “on” pixels
xedge : the mean number of edges (the number of times an “off” pixel is followed by an “on” pixel, or the image boundary is hit) as the image is scanned from left to right, along the whole vertical length of the image
xedgeycor : the mean of the product of the number of horizontal edges at each vertical position and the vertical position
yedge : the mean number of edges as the images is scanned from top to bottom, along the whole horizontal length of the image
yedgexcor : the mean of the product of the number of vertical edges at each horizontal position and the horizontal position

Exercices

1. Predicting *B* or not *B*

Problem 1.1 Let’s warm up by attempting to predict just whether a letter is B or not. To begin, load the file `letters_ABPR.csv` into R, and call it `letters`. Then, create a new variable `isB` in the dataframe, which takes the value “TRUE” if the observation corresponds to the letter B, and “FALSE” if it does not. You can do this by typing the following command into your R console:

Now split the data set into a training and testing set, putting 50% of the data in the training set. Set the seed to 1000 before making the split. The first argument to `sample.split` should be the dependent variable “`letters$isB`”. Remember that TRUE values from `sample.split` should go in the training set.

Before building models, let’s consider a baseline method that always predicts the most frequent outcome, which is “not B”.

What is the accuracy of this baseline method on the test set?

```
## 'data.frame': 3116 obs. of 17 variables:
## $ letter : chr "B" "A" "R" "B" ...
## $ xbox : int 4 1 5 5 3 8 2 3 8 6 ...
## $ ybox : int 2 1 9 9 6 10 6 7 14 10 ...
## $ width : int 5 3 5 7 4 8 4 5 7 8 ...
## $ height : int 4 2 7 7 4 6 4 5 8 8 ...
## $ onpix : int 4 1 6 10 2 6 3 3 4 7 ...
## $ xbar : int 8 8 6 9 4 7 6 12 5 8 ...
## $ ybar : int 7 2 11 8 14 7 7 2 10 5 ...
## $ x2bar : int 6 2 7 4 8 3 5 3 6 7 ...
## $ y2bar : int 6 2 3 4 1 5 5 2 3 5 ...
## $ xybar : int 7 8 7 6 11 8 6 10 12 7 ...
## $ x2ybar : int 6 2 3 8 6 4 5 2 5 6 ...
## $ xy2bar : int 6 8 9 6 3 8 7 9 4 6 ...
## $ xedge : int 2 1 2 6 0 6 3 2 4 3 ...
## $ xedgeycor: int 8 6 7 11 10 6 7 6 10 9 ...
## $ yedge : int 7 2 5 8 4 7 5 3 4 8 ...
## $ yedgexcor: int 10 7 11 7 8 7 8 8 8 9 ...
```

```
##      letter      xbox      ybox      width
## Length:3116      Min.   : 0.000      Min.   : 0.000      Min.   : 1.000
## Class :character 1st Qu.: 3.000      1st Qu.: 5.000      1st Qu.: 4.000
## Mode  :character Median : 4.000      Median : 7.000      Median : 5.000
##                Mean  : 3.915      Mean  : 7.051      Mean  : 5.186
##                3rd Qu.: 5.000      3rd Qu.: 9.000      3rd Qu.: 6.000
##                Max.   :13.000      Max.   :15.000      Max.   :11.000
##      height      onpix      xbar      ybar
## Min.   : 0.000      Min.   : 0.000      Min.   : 3.000      Min.   : 0.000
## 1st Qu.: 4.000      1st Qu.: 2.000      1st Qu.: 6.000      1st Qu.: 6.000
## Median : 6.000      Median : 4.000      Median : 7.000      Median : 7.000
## Mean   : 5.276      Mean   : 3.869      Mean   : 7.469      Mean   : 7.197
## 3rd Qu.: 7.000      3rd Qu.: 5.000      3rd Qu.: 8.000      3rd Qu.: 9.000
## Max.   :12.000      Max.   :12.000      Max.   :14.000      Max.   :15.000
##      x2bar      y2bar      xybar      x2ybar
## Min.   : 0.000      Min.   :0.000      Min.   : 3.000      Min.   : 0.00
## 1st Qu.: 3.000      1st Qu.:2.000      1st Qu.: 7.000      1st Qu.: 3.00
## Median : 4.000      Median :4.000      Median : 8.000      Median : 5.00
## Mean   : 4.706      Mean   :3.903      Mean   : 8.491      Mean   : 4.52
## 3rd Qu.: 6.000      3rd Qu.:5.000      3rd Qu.:10.000      3rd Qu.: 6.00
## Max.   :11.000      Max.   :8.000      Max.   :14.000      Max.   :10.00
##      xy2bar      xedge      xedgeycor      yedge
## Min.   : 0.000      Min.   : 0.000      Min.   : 1.000      Min.   : 0.0
## 1st Qu.: 6.000      1st Qu.: 2.000      1st Qu.: 7.000      1st Qu.: 3.0
## Median : 7.000      Median : 2.000      Median : 8.000      Median : 4.0
## Mean   : 6.711      Mean   : 2.913      Mean   : 7.763      Mean   : 4.6
## 3rd Qu.: 8.000      3rd Qu.: 4.000      3rd Qu.: 9.000      3rd Qu.: 6.0
## Max.   :14.000      Max.   :10.000      Max.   :13.000      Max.   :12.0
##      yedgexcor
## Min.   : 1.000
## 1st Qu.: 7.000
## Median : 8.000
## Mean   : 8.418
## 3rd Qu.:10.000
## Max.   :13.000

##
## FALSE TRUE
## 1175 383

##
## FALSE TRUE
## 1175 383

## [1] 0.754172
```

Answer : 0.754172

Explanation :

Load the csv file:

To compute the accuracy of the baseline method on the test set, we first need to see which outcome value is more frequent in the training set, by using the table function. The output of `table(train$isB)` tells us that “not B” is more common. So our baseline method is to predict “not B” for everything. How well would this do on the test set? We need to run the table command again, this time on the test set:

There are 1175 observations that are not B, and 383 observations that are B. So the baseline method accuracy on the test set would be $1175/(1175+383) = 0.754172$

Problem 1.2 Now build a classification tree to predict whether a letter is a B or not, using the training set to build your model. Remember to remove the variable “letter” out of the model, as this is related to what we are trying to predict! To just remove one variable, you can either write out the other variables, or remember what we did in the Billboards problem in Week 3, and use the following notation:

We are just using the default parameters in our CART model, so we don’t need to add the minbucket or cp arguments at all. We also added the argument method=“class” since this is a classification problem.

What is the accuracy of the CART model on the test set?

(Use type=“class” when making predictions on the test set.)

```
##          CARTletterB.pred
##          FALSE TRUE
##  FALSE  1118   57
##   TRUE    43  340
```

```
## [1] 0.9358151
```

Answer : 0.9358151

Explanation :

You can build the CART tree with the following command:

To compute the accuracy on the test set, we need to divide the sum of the true positives and true negatives by the total number of observations: $(1118+340)/\text{nrow}(\text{test}) = 0.9358151$

Problem 1.3 Now, build a random forest model to predict whether the letter is a B or not (the isB variable) using the training set. You should use all of the other variables as independent variables, except letter (since it helped us define what we are trying to predict!). Use the default settings for ntree and nodesize (don’t include these arguments at all). Right before building the model, set the seed to 1000. (NOTE: You might get a slightly different answer on this problem, even if you set the random seed. This has to do with your operating system and the implementation of the random forest algorithm.)

What is the accuracy of the model on the test set?

```
##          LetterForest.pred
##          FALSE TRUE
##  FALSE  1163   12
##   TRUE    9  374
```

```
## [1] 0.9865212
```

Answer : 0.9865212

Explanation :

To build the random forest model, first set the seed to 1000:

The accuracy of the model on the test set is the sum of the true positives and true negatives, divided by the total number of observations in the test set:

In lecture, we noted that random forests tends to improve on CART in terms of predictive accuracy. Sometimes, this improvement can be quite significant, as it is here.

2. Predicting the letters A, B, P, R

Problem 2.1 Let us now move on to the problem that we were originally interested in, which is to predict whether or not a letter is one of the four letters A, B, P or R.

As we saw in the D2Hawkeye lecture, building a multiclass classification CART model in R is no harder than building the models for binary classification problems. Fortunately, building a random forest model is just as easy.

The variable in our data frame which we will be trying to predict is “letter”. Start by converting letter in the original data set (letters) to a factor by running the following command in R:

Now, generate new training and testing sets of the letters data frame using letters\$letter as the first input to the sample.split function. Before splitting, **set your seed to 2000**. Again put **50%** of the data in the training set. (Why do we need to split the data again? Remember that sample.split balances the outcome variable in the training and testing sets. With a new outcome variable, we want to re-generate our split.)

In a multiclass classification problem, a simple baseline model is to predict the most frequent class of all of the options.

What is the baseline accuracy on the testing set?

```
##
##   A   B   P   R
## 385 383 396 394

## [1] 0.2580231
```

Answer : 0.2580231

Explanation :

After converting the variable “letter” to a factor, set the seed to 2000 and generate the new split:

Then to compute the accuracy of the baseline method on the test set, we need to first figure out the most common outcome in the training set. The output of

tells us that “P” has the most observations. So we will predict P for all letters. On the test set, we can run the table command table(test2\$letter) to see that it has 401 observations that are actually P. So the test set accuracy of the baseline method is $401/\text{nrow}(\text{test}) = 0.2573813$.

Problem 2.2 Now build a classification tree to predict “letter”, using the training set to build your model. You should use all of the other variables as independent variables, except “isB”, since it is related to what we are trying to predict! Just use the default parameters in your CART model. Add the argument method=“class” since this is a classification problem. Even though we have multiple classes here, nothing changes in how we build the model from the binary case.

What is the test set accuracy of your CART model?

Use the argument type=“class” when making predictions.

(HINT: When you are computing the test set accuracy using the confusion matrix, you want to add everything on the main diagonal and divide by the total number of observations in the test set, which can be computed with nrow(test), where test is the name of your test set).

```
##   CARTLetters.pred
##      A   B   P   R
## A 348   4   0  43
## B   8 318  12  45
## P   2  21 363  15
## R  10  24   5 340
```

```
## [1] 0.8786906
```

Answer : 0.8786906

Explanation :

You can build the CART tree with the following command:

Looking at the confusion matrix, `table(test2$letter, predictLetter)`, we want to sum the main diagonal (the correct predictions) and divide by the total number of observations in the test set:

Problem 2.3 Now build a random forest model on the training data, using the same independent variables as in the previous problem – again, don't forget to remove the `isB` variable. Just use the default parameter values for `ntree` and `nodesize` (you don't need to include these arguments at all). Set the seed to 1000 right before building your model. (Remember that you might get a slightly different result even if you set the random seed.)

What is the test set accuracy of your random forest model?

```
##      LettersForest.pred
##      A  B  P  R
## A 391  0  3  1
## B  0 380  1  2
## P  0  6 394  1
## R  3 14  0 362
```

```
## [1] 0.9801027
```

Answer : 0.9801027

Explanation :

First set the seed, and then build the random forest model:

And then we can compute the test set accuracy by looking at the confusion matrix `table(test2$letter, predictLetter)`. The test set accuracy is the sum of the numbers on the main diagonal, divided by the total number of observations in the test set:

You should find this value rather striking, for several reasons. The first is that it is significantly higher than the value for CART, highlighting the gain in accuracy that is possible from using random forest models. The second is that while the accuracy of CART decreased significantly as we transitioned from the problem of predicting B/not B (a relatively simple problem) to the problem of predicting the four letters (certainly a harder problem), the accuracy of the random forest model decreased by a tiny amount.