

Lab2- Group A 13

Arash Haratian, Connor Turner, Yi Hung Chen

2022-11-27

Assignment 1

Assignment 2

```
## Warning: package 'tidyverse' was built under R version 4.2.1
## -- Attaching packages ----- tidyverse 1.3.2 --
## v ggplot2 3.3.6      v purrr   0.3.4
## v tibble  3.1.8      v dplyr  1.0.10
## v tidyr   1.2.1      v stringr 1.4.1
## v readr   2.1.2      v forcats 0.5.2
## Warning: package 'tibble' was built under R version 4.2.1
## Warning: package 'tidyr' was built under R version 4.2.1
## Warning: package 'dplyr' was built under R version 4.2.1
## Warning: package 'stringr' was built under R version 4.2.1
## Warning: package 'forcats' was built under R version 4.2.1
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

In this assignment we are given the data related with direct marketing campaigns of a Portuguese banking institution, and we have to use decision tree model to classify the target variable y using 15 features. The target variable has two levels: **yes** and **no**.

We begin by splitting data 40/30/30 into training, validation, and testing sets. We then use the training data to train 3 different decision trees with 3 different settings:

1. First Decision Tree:

The first model is a decision tree with default settings of the `tree::tree()`. The default values are as follows: - `mincut = 5` - `minsize = 10` - `mindev = 0.01`

The training and test errors for this model are 0.1048 and 0.1093 respectively.

2. Second Decision Tree: The second model is a decision tree with a constraint of the size of each node; In fact, the smallest allowed node size should be equal to 7000. To set this setting, we should pass `tree::tree.control(minsize = 7000)` to the `control` argument of the `tree::tree()`.

The training and test errors for the second model are 0.1048 and 0.1093 respectively.

3. Third Decision Tree:

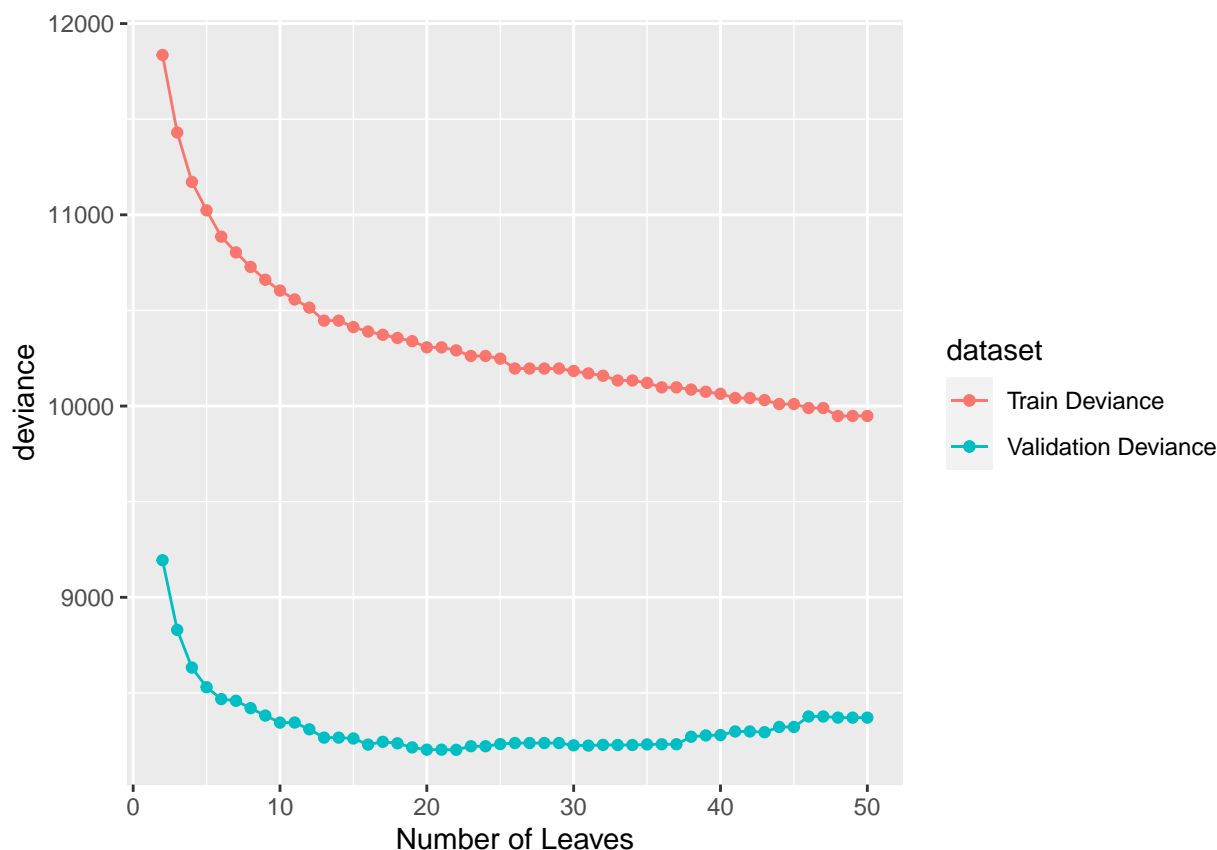
The third model is a decision we want to set the minimum deviance to 0.0005, and it can be done in the same way as the second model, by passing `tree::tree.control(mindev = 0.0005)` to the `control` argument of the `tree::tree()`.

The training and test errors for the third model are 0.094 and 0.1119 respectively.

| different values | First Model | Second Model | Third Model |
|------------------------|-------------|--------------|-------------|
| Size of Trees | 6 | 5 | 122 |
| Number of feature used | 4 | 3 | 13 |
| Training Error | 0.1048 | 0.1048 | 0.094 |
| Testing Error | 0.1093 | 0.1093 | 0.1119 |

Based on the information in the table above, first model and second model have the same error rate. However, both the number of terminal nodes and number of feature used to split for the second model is 1 less than the first model. In our opinion the second model is more proper for interpreting the structure of the tree, but for predicting new values, the first model may perform better as it is more complex compared to the second model. The third model is overfitted on the training data and it is too complex; as the result, it may perform poorly for predicting new values.

Now we can prune the more complex model, in this case the third model (minimum deviance to 0.0005), to reach to a more optimal model. By doing so, we may end up with a model which captures the general structure better than other models while it is not overfitted on the training dataset. In other words, by pruning a fully-grown tree, we are mitigating the overfitting (or high variance) problem. The following figure shows the dependence of deviances for the training and the validation data in the number of leaves:



As it is evidence in the plot, the deviance of the training data set is decreasing by increasing the number of leaves (complexity), whereas the deviance of the validation data set is decreasing at first but then it starts to increase. This denotes that the model will overfit by increasing the number of terminal nodes. It is worth to note that the deviance of the training data set is higher than validation dataset because deviance is proportional to the number of observations in the dataset.

The best number of the leaves is 22.

The important variables in the optimal tree are as follows: - poutcome - month - contact - pdays - age - day - balance - housing - job Also we can see the the structure of the tree:

```
## node), split, n, deviance, yval, (yprob)
##      * denotes terminal node
##
##      1) root 18084 12850.00 no ( 0.88576 0.11424 )
##      2) poutcome: failure,other,unknown 17468 11030.00 no ( 0.90422 0.09578 )
##      4) month: apr,aug,feb,jan,jul,jun,may,nov 16828 9772.00 no ( 0.91520 0.08480 )
##      8) contact: unknown 5130 1599.00 no ( 0.96374 0.03626 )
##      16) month: jul,jun,may 5074 1502.00 no ( 0.96610 0.03390 ) *
##      17) month: apr,aug,feb,jan,nov 56 62.98 no ( 0.75000 0.25000 ) *
##      9) contact: cellular,telephone 11698 7914.00 no ( 0.89391 0.10609 )
##      18) month: aug,jan,jul,may,nov 9284 5503.00 no ( 0.91265 0.08735 )
##      36) pdays < 383.5 9246 5373.00 no ( 0.91510 0.08490 )
##      72) age < 60.5 9097 5107.00 no ( 0.91920 0.08080 )
##      144) day < 27.5 7670 4588.00 no ( 0.91147 0.08853 )
##      288) age < 29.5 754 637.10 no ( 0.85013 0.14987 )
##      576) month: jan,jul,may 681 528.00 no ( 0.86931 0.13069 ) *
##      577) month: aug,nov 73 92.46 no ( 0.67123 0.32877 ) *
##      289) age > 29.5 6916 3918.00 no ( 0.91816 0.08184 )
##      578) balance < 1493 5180 2635.00 no ( 0.92973 0.07027 )
##      1156) month: aug,jul,may,nov 5164 2596.00 no ( 0.93087 0.06913 ) *
##      1157) month: jan 16 21.93 no ( 0.56250 0.43750 ) *
##      579) balance > 1493 1736 1249.00 no ( 0.88364 0.11636 ) *
##      145) day > 27.5 1427 472.40 no ( 0.96076 0.03924 )
##      290) pdays < 184.5 1308 462.50 no ( 0.95719 0.04281 )
##      580) pdays < 80.5 1277 402.60 no ( 0.96319 0.03681 ) *
##      581) pdays > 80.5 31 37.35 no ( 0.70968 0.29032 ) *
##      291) pdays > 184.5 119 0.00 no ( 1.00000 0.00000 ) *
##      73) age > 60.5 149 190.10 no ( 0.66443 0.33557 ) *
##      37) pdays > 383.5 38 47.40 yes ( 0.31579 0.68421 ) *
##      19) month: apr,feb,jun 2414 2262.00 no ( 0.82187 0.17813 )
##      38) housing: no 1123 1321.00 no ( 0.72484 0.27516 )
##      76) day < 9.5 691 668.20 no ( 0.81187 0.18813 )
##      152) month: feb 505 364.20 no ( 0.88317 0.11683 ) *
##      153) month: apr,jun 186 247.30 no ( 0.61828 0.38172 ) *
##      77) day > 9.5 432 586.10 no ( 0.58565 0.41435 ) *
##      39) housing: yes 1291 803.20 no ( 0.90627 0.09373 )
##      78) day < 20.5 1210 655.90 no ( 0.92314 0.07686 )
##      156) month: apr,feb 1154 565.70 no ( 0.93328 0.06672 ) *
##      157) month: jun 56 67.01 no ( 0.71429 0.28571 ) *
##      79) day > 20.5 81 104.40 no ( 0.65432 0.34568 ) *
##      5) month: dec,mar,oct,sep 640 852.70 no ( 0.61562 0.38438 ) *
##      3) poutcome: success 616 806.40 yes ( 0.36201 0.63799 )
##      6) pdays < 94.5 170 185.50 yes ( 0.23529 0.76471 ) *
##      7) pdays > 94.5 446 603.90 yes ( 0.41031 0.58969 )
##      14) job: admin.,blue-collar,entrepreneur,services,technician 213 295.20 no ( 0.50704 0.49296 )
##      15) job: housemaid,management,retired,self-employed,student,unemployed,unknown 233 292.80 yes ( 0.50704 0.49296 )
```

Based on the structure of the tree, the first variable that has been used to split the data, is poutcome which is the outcome of the previous marketing campaign. Also it is worth to note that the variable month (last contact month of year) has been used frequently for splitting the data.

For evaluating the performance of the optimal tree we can check the confusion matrix, accuracy and F1-score for the testing data. F1-score is better metric to measure the performance of the model as the target variable **y** is unbalanced.

| | no | yes |
|-----|-------|-----|
| no | 11872 | 107 |
| yes | 1371 | 214 |

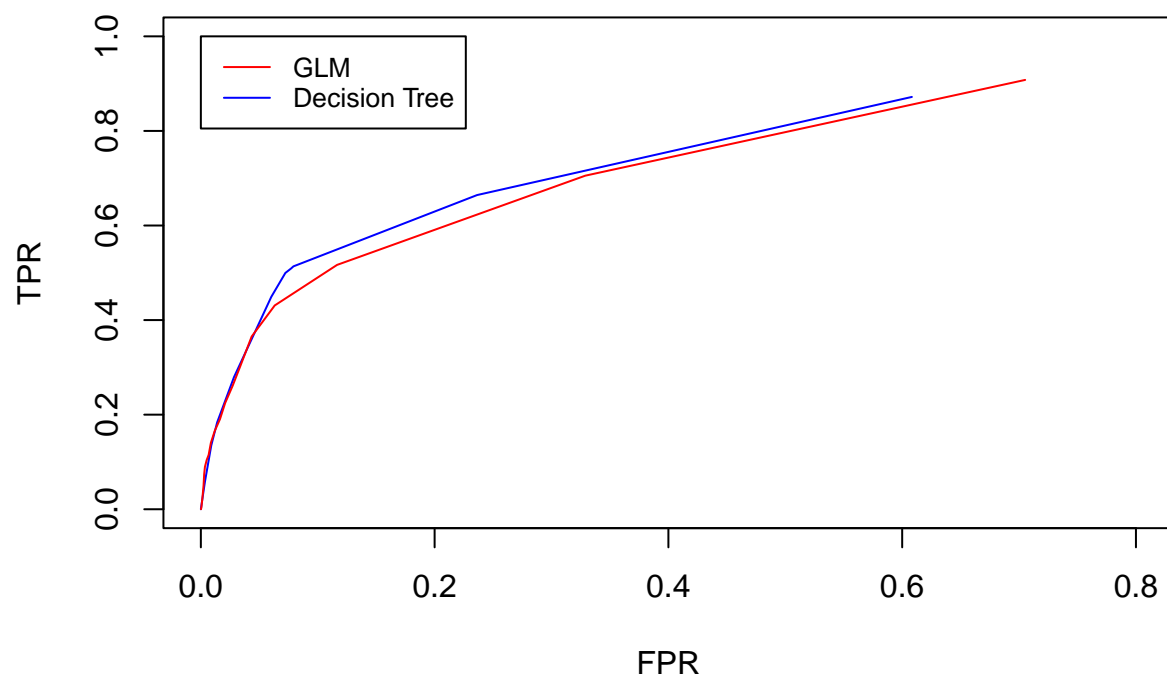
The accuracy is 0.891 which indicates that model is very accurate in predicting new values. However, as mentioned before, since the target variable has mostly **no** values, a model which predicts all the values as **no** will yield a high accuracy as well. A good model is a model which can predict **yes** more accurately in this dataset. F1-score for this model is 0.224554 which is very low and it indicates that model is not performing good in predicting true positives.

One way to increase the performance of the model is to use a loss matrix to predict the values based on their probabilities (instead of using 0.5 as threshold, we are using $1/5$). In other words, we are penalizing the model if it predicts positives wrongly. The confusion matrix for the testing data and using the loss matrix is as follows:

| | 0 | 1 |
|-----|-------|-----|
| no | 11030 | 949 |
| yes | 771 | 814 |

The new F1-score is 0.4862605 which means now the model predict true positives more accurately although the accuracy has decreased to 0.8732.

One way to compare two different models is to use a ROC curve. Here we want to compare the optimal decision tree with a logistic regression.



Based on the ROC curve, both models are almost perform in the same way, but decision tree is slightly better. For the same value of π and same value of False Positive Rate (FPR), the decision tree has higher value of True Positive Rate.

The precision-recall curve is a better choice to compare two models as the classes of target data is imbalanced. The model which has the higher recall for the same value of precision, is more accurate to predict positive values.