

CSI4142 Fundamentals of Data Science Project Phase 4: Data Mining

Group 1

Professor Yazan Otoum,

TA: Lansu Dai

Jenson Wu #300166874
Victor Feng #300176400
Lootii Kiri #300189957

Date: April 9th 2024

Part A:

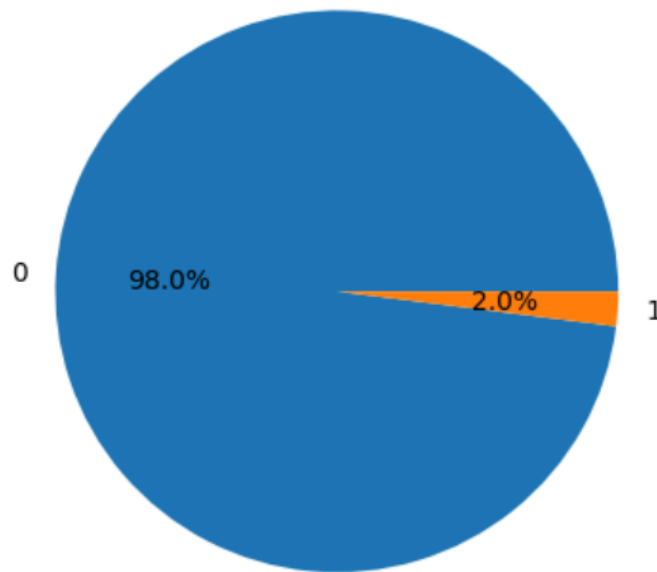


Figure 1 Class imbalance between canceled and non-cancelled flight

When looking at the dashboard that we created in phase 3 of the project, we realized that we had a huge class imbalance between canceled and non-cancelled flights, where only 2 % of all flights were canceled yearly(which meant that we had 500 000 data points for canceled flights). We decided to use undersampling for this phase because we felt that populating the minority will not produce realistic results and we also felt that 500 k for each class was sufficient to get an accurate classifier.

For the sampling technique, we decided to use stratified sampling to get the training(contain 66% of the data points) and the testing(contain 33% if the data) datasets, so that we got an even distribution of all attributes and class labels

Part B. Classification (Supervised Learning):

Decision Tree:

For the decision tree algorithm, we decided to do this in two ways to show our understanding and application of the algorithm essential for constructing decision trees. Instead of only using a decision tree library made accessible in python as seen in scikit-learn, we also implemented our own code that uses the logic of finding the Entropy of both categorical and numerical attributes of a given set. Our code also takes the entropy and calculates an information gain of all of the given attributes in relation to our class attribute.

```

1 package com.example;
2
3 import org.apache.poi.xssf.usermodel.XSSFSheet;
4 import org.apache.poi.xssf.usermodel.XSSFCell;
5 //import org.apache.poi.xssf.usermodel.XSSFWorkbook;
6 import org.apache.poi.xssf.usermodel.XSSFRow;
7
8 //import java.io.FileInputStream;
9 //import java.io.FileNotFoundException;
10 //import java.io.IOException;
11
12 public class DecisionTreeUtils {
13
14     public static int countOccurrencesOfMonth(XSSFSheet sheet, Double month) {
15
16         int numOfRows = sheet.getPhysicalNumberOfRows() - 1;
17         int count = 0;
18
19         for (int i = 1; i <= numOfRows; i++) {
20             XSSFRow row = sheet.getRow(i);
21             if (row.getCell(0).getNumericCellValue() == month) {
22                 count++;
23             }
24         }
25         return(count);
26     }
27
28     public static int countOccurrencesOfMonthCancellationType(XSSFSheet sheet, Double month, Double binary) {
29
30         int numOfRows = sheet.getPhysicalNumberOfRows() - 1;
31         int count = 0;
32
33         for (int i = 1; i <= numOfRows; i++) {
34             XSSFRow row = sheet.getRow(i);
35             XSSFCell lastCell = row.getCell(8);
36             if (row.getCell(0).getNumericCellValue() == month && lastCell.getNumericCellValue() == binary) {
37                 count++;
38             }
39         }
40         return(count);
41     }
42 }
43

```

Image 1: functions to count the occurrence of a given month. As well as one that counts the occurrence of a given month that is either canceled(binary = 1.0) or not canceled(binary = 0.0).

```

44     public static double entropy(double positive, double negative){
45         double total = positive + negative;
46         double entropy = -((positive/total) * (Math.log(positive/total) / Math.log(2)))
47         -((negative/total) * (Math.log(negative/total) / Math.log(2)));
48         return entropy;
49     }
50
51     public static double infoAttribute(XSSFSheet sheet){
52         int numOfRows = sheet.getPhysicalNumberOfRows() - 1;
53         double zeros = 0.0;
54         double ones = 0.0;
55         for(int i=1; i<= numOfRows; i++){
56             XSSFRow row = sheet.getRow(i);
57             XSSFCell lastCell = row.getCell(8);
58             if(lastCell.getNumericCellValue() == 0.0){
59                 zeros++;
60             }
61             if(lastCell.getNumericCellValue() == 1.0){
62                 ones++;
63             }
64         }
65         double infoAttribute = entropy(ones,zeros);
66
67         return infoAttribute;
68     }

```

Image 2: function that calculates the entropy and method that calculates Info(Class). As in the information or entropy of the entire class which will be later used to subtract from attribute entropies to receive information gain.

```

1 package com.example;
2
3 import org.apache.poi.xssf.usermodel.XSSFSheet;
4
5 public class MonthInfoGain {
6
7     public static double calculateEntropy(XSSFSheet sheet, Double month) {
8
9         int occurrencesOfCancelled = DecisionTreeUtils.countOccurrencesOfMonthCancellationType(sheet, month, 1.0);
10        int occurrencesOfNotCancelled = DecisionTreeUtils.countOccurrencesOfMonthCancellationType(sheet, month, 0.0);
11        int totalOccurrences = DecisionTreeUtils.countOccurrencesOfMonth(sheet, month);
12
13        double probabilityOfCancelled = (double) occurrencesOfCancelled / totalOccurrences;
14        double probabilityOfNotCancelled = (double) occurrencesOfNotCancelled / totalOccurrences;
15
16        double entropy = -(probabilityOfCancelled * (Math.log(probabilityOfCancelled) / Math.log(2)))
17        -(probabilityOfNotCancelled * (Math.log(probabilityOfNotCancelled) / Math.log(2)));
18
19        return entropy;
20    }
21
22
23
24
25    public static double infoGainMonth(XSSFSheet sheet){
26        double IGM = 0;
27        for(double i =1.0; i<=12.0; i++){
28            int totalOccurrences = DecisionTreeUtils.countOccurrencesOfMonth(sheet, i);
29            int numOfRows = sheet.getPhysicalNumberOfRows() - 1;
30            double entropy = calculateEntropy(sheet, i);
31            IGM += (((double) totalOccurrences/numOfRows)*entropy);
32        }
33
34        double classEntropy = DecisionTreeUtils.infoAttribute(sheet);
35
36        return classEntropy - IGM;
37    }
38
39
40 }

```

image 3: since our “Month” attribute is categorical, it is handled differently then the numerical attributes. Here are functions that calculate the entropy and information gain for “Month”.

```

62 public static double[] countOfZeroAndOne(double[][] sortedArray, int number){
63     int numOfRows = sortedArray.length -1;
64     double zeroCountLessThan = 0;
65     double oneCountLessThan = 0;
66     double zeroCountGreaterThan = 0;
67     double oneCountGreaterThan = 0;
68
69     for(int j=0; j<=numOfRows; j++){
70         if(j<=number && sortedArray[j][1] == 0.0){
71             zeroCountLessThan++;
72         }
73         if(j<=number && sortedArray[j][1] == 1.0){
74             oneCountLessThan++;
75         }
76         if(j>number && sortedArray[j][1] == 0.0){
77             zeroCountGreaterThan++;
78         }
79         if(j>number && sortedArray[j][1] == 1.0){
80             oneCountGreaterThan++;
81         }
82     }
83     double[] values = {oneCountLessThan, zeroCountLessThan-1, oneCountGreaterThan, zeroCountGreaterThan};
84
85     return values;
86 }
87
88 public static double[][] sortedAttributesClass(XSSFSheet sheet, int coloumn){
89     int numOfRows = sheet.getPhysicalNumberOfRows() - 1;
90     double[][] sorted = new double[numOfRows+1][2];
91     for(int i=1; i<=numOfRows; i++){
92         XSSFRow row = sheet.getRow(i);
93         XSSFCell cellOfAttribute = row.getCell(coloumn);
94         XSSFCell cellOfClass = row.getCell(8);
95         sorted[i][0] = cellOfAttribute.getNumericCellValue();
96         sorted[i][1] = cellOfClass.getNumericCellValue();
97     }
98     // Sort the 2D array based on the first element of each row
99     Arrays.sort(sorted, Comparator.comparingDouble(a -> a[0]));
100
101     return sorted;
102 }
103 }

```

Image 4: Method to count the numbers of flights canceled (oneCount..) and not canceled (zeroCount..) for values less than and equal to a particular row and greater than. Another method that sorts the Attributes in order for numeric attributes.

```

106 public static double numericalSplitEntropy(double[][] sortedList){
107
108     double[] values = countOfZeroAndOne(sortedList, 1);
109     values[0]--;
110     values[2]++;
111     double maxGain = 0.0;
112     for(int i=1; i<sortedList.length; i++){
113         double currentGain;
114
115         if(sortedList[i][1] == 1.0){
116             values[0]++;
117             values[2]--;
118         }
119         if(sortedList[i][1] == 0.0){
120             values[1]++;
121             values[3]--;
122         }
123         double positiveValuesLessThan = values[0];
124         double negativeValuesLessThan = values[1];
125         double positiveValuesGreaterThan = values[2];
126         double negativeValuesGreaterThan = values[3];
127
128         double totalPos = positiveValuesGreaterThan + positiveValuesLessThan;
129         double totalNeg = negativeValuesGreaterThan + negativeValuesLessThan;
130         double attributeEntropy = DecisionTreeUtils.entropy(totalPos, totalNeg);
131
132         double lessThanEntropy = DecisionTreeUtils.entropy(positiveValuesLessThan, negativeValuesLessThan);
133         double greaterThanEntropy = DecisionTreeUtils.entropy(positiveValuesGreaterThan, negativeValuesGreaterThan);
134
135         double total = totalNeg + totalPos;
136         currentGain = attributeEntropy - (((positiveValuesLessThan+negativeValuesLessThan)/total) * lessThanEntropy)
137         - (((positiveValuesGreaterThan + negativeValuesGreaterThan)/ total) *greaterThanEntropy);
138         if(currentGain > maxGain){
139             maxGain = currentGain;
140         }
141     }
142
143     return maxGain;
144 }
145
146 }

```

Image 5: Finally we have the method that calculates and returns the information gain of a numerical attribute. It goes through the whole dataset, calculates each split points entropy, subtracts it from the class entropy for information gain and keeps track of the maximum one as it traverses the dataset. In the end it will return the largest info gain.

```

17 public class App
18 {
19
20     private static XSSFSheet readXLSXFile(String file) {
21         try {
22             XSSFWorkbook work = new XSSFWorkbook(new FileInputStream(file));
23             XSSFSheet sheet = work.getSheet("in");
24             return sheet;
25
26         } catch (FileNotFoundException e) {
27             // TODO Auto-generated catch block
28             e.printStackTrace();
29         } catch (IOException e) {
30             // TODO Auto-generated catch block
31             e.printStackTrace();
32         }
33         return null;
34     }
35
36
37     Run | Debug
38     public static void main( String[] args ) {
39         System.out.println( "Hello World!" );
40
41         XSSFSheet sheet = readXLSXFile("/Users/lootiikiri/Desktop/CSI4142/Training_DataSet.xlsx");
42
43         double[][] arr = ContinuousAttributesUtil.sortedAttributesClass(sheet, 1);
44         double[][] a = ContinuousAttributesUtil.sortedAttributesClass(sheet, 2);
45
46         System.out.println(ContinuousAttributesUtil.numericalSplitEntropy(arr));
47         System.out.println(ContinuousAttributesUtil.numericalSplitEntropy(a));
48         System.out.println(MonthInfoGain.infoGainMonth(sheet));
49
50
51     }
52
53
54 }

```

```

0.0023154833669960656
0.0017042592922772715
0.047464345639343986

```

Performance Measurement for all the algorithm

Classifier	Construction Time	Accuracy	Precision	Recall	F_Score
Decision Tree	2sec	63.6%	64.23%	61.37%	62.77%

Random Forest	1m 24sec	64.21%	63.61%	66.4%	64.98%
Gradient Boosting	3m 28sec	62.57%	61.89%	65.45%	63.62%

Analysis

1. Construction Time

The construction time is very coherent to our understanding of how each of the classifiers are ensembled. First the decision was expected to have by far the lowest construction time because only 1 tree is being constructed for the classification, where the other classifier algorithms are using multiple subtrees and trees to perform the classification. The Random Forest classifier was also expected to have a lower construction time than the Gradient boosting because unlike Gradient boosting, it only utilize one decision tree to create multiple subtree that will perform different classification algorithm creating this forest type structure, where Gradient boosting create multiple decision tree that would gradually correct the mistakes made in the prior trees.

2. Accuracy

The accuracy is not coherent with our understanding of the behavior of the classifier, but It does make sense given the context of our dataset. Hypothetically, the gradient boosting should have a superior accuracy to the decision tree because it uses multiple trees that will correct itself after each iteration of the decision tree. If you look at the recall, it has a superior recall to the decision tree, which implies that the gradient boosting is very good at predicting all of the relevant cases. The problem that may have caused this situation is the class imbalance that we had, where we utilize under sampling (removing non-cancelled flight from our dataset) in order to remove the class imbalance, which the data points that were removed may have been useful to accurately predicting the negative cases. Another potential issues with the gradient boosting is that since our datasets is being reduced, the value of each attribute may be very sparse which could impact the accuracy of the classifier

3. Precision, Recall, F-Score

The F-score and the recall was coherent with our understanding of the behavior of the classifier, but the precision of the gradient is substantially lower than we anticipated. Similar to the accuracy, this may have been caused by the class imbalance that we had and how we removed many negative cases (non-cancelled flights) in order to resolve the class imbalance.