

Department of Computer Science & Engineering

Final Year B. Tech. (CSE) – I : 2021-22

4CS462 : PE2 - Data Mining Lab

Assignment No. 5

By DM21G03

(2018BTECS00082 : Hritik Belani , 2019BTECS00209 : Sailee Akim)

Date: 03/10/2021

❖ Title

Data Analysis Tool

❖ Objective

Design the rule-based classifier: Extract the rules from decision tree build

❖ Specification

- Python 3.8.11

- Dataset

❖ Introduction and Theory

IF-THEN Rules

Rule-based classifier makes use of a set of IF-THEN rules for classification. We can express a rule in the following form –

IF condition THEN conclusion

Let us consider a rule R1,

R1: IF age = youth AND student = yes THEN buy_computer = yes

Points to remember –

- The IF part of the rule is called **rule antecedent** or **precondition**.
- The THEN part of the rule is called **rule consequent**.
- The antecedent part the condition consist of one or more attribute tests and these tests are logically ANDed.
- The consequent part consists of class prediction.

Note – We can also write rule R1 as follows –

R1: (age = youth) ^ (student = yes))(buys computer = yes)

If the condition holds true for a given tuple, then the antecedent is satisfied.

Rule Extraction

Here we will learn how to build a rule-based classifier by extracting IF-THEN rules from a decision tree.

Points to remember –

To extract a rule from a decision tree –

- One rule is created for each path from the root to the leaf node.
- To form a rule antecedent, each splitting criterion is logically ANDed.
- The leaf node holds the class prediction, forming the rule consequent.

Rule Induction Using Sequential Covering Algorithm

Sequential Covering Algorithm can be used to extract IF-THEN rules from the training data. We do not require to generate a decision tree first. In this algorithm, each rule for a given class covers many of the tuples of that class.

Some of the sequential Covering Algorithms are AQ, CN2, and RIPPER. As per the general strategy the rules are learned one at a time. For each time rules are learned, a tuple covered by the rule is removed and the process continues for the rest of the tuples. This is because the path to each leaf in a decision tree corresponds to a rule.

Note – The Decision tree induction can be considered as learning a set of rules simultaneously.

The Following is the sequential learning Algorithm where rules are learned for one class at a time. When learning a rule from a class C_i , we want the rule to cover all the tuples from class C only and no tuple from any other class.

Algorithm: Sequential Covering
Input: D , a data set class-labeled tuples, Att_vals , the set of all attributes and their possible values.
Output: A Set of IF-THEN rules.
Method: Rule_set = { }; // initial set of rules learned is empty
for each class c do repeat Rule = Learn_One_Rule(D , Att_vals , c);
 remove tuples covered by Rule from D ; until termination condition;
Rule_set = Rule_set + Rule; // add a new rule to rule-set
end for
return Rule_Set;

Rule Pruning

The rule is pruned is due to the following reason –

- The Assessment of quality is made on the original set of training data. The rule may perform well on training data but less well on subsequent data. That's why the rule pruning is required.
- The rule is pruned by removing conjunct. The rule R is pruned, if pruned version of R has greater quality than what was assessed on an independent set of tuples.

FOIL is one of the simple and effective method for rule pruning. For a given rule R ,

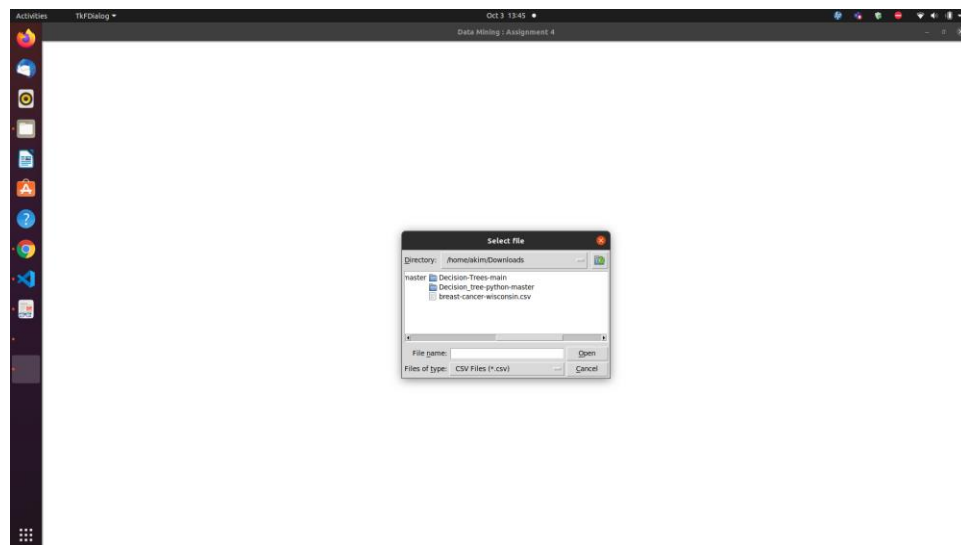
$$\text{FOIL_Prune} = \text{pos} - \text{neg} / \text{pos} + \text{neg}$$

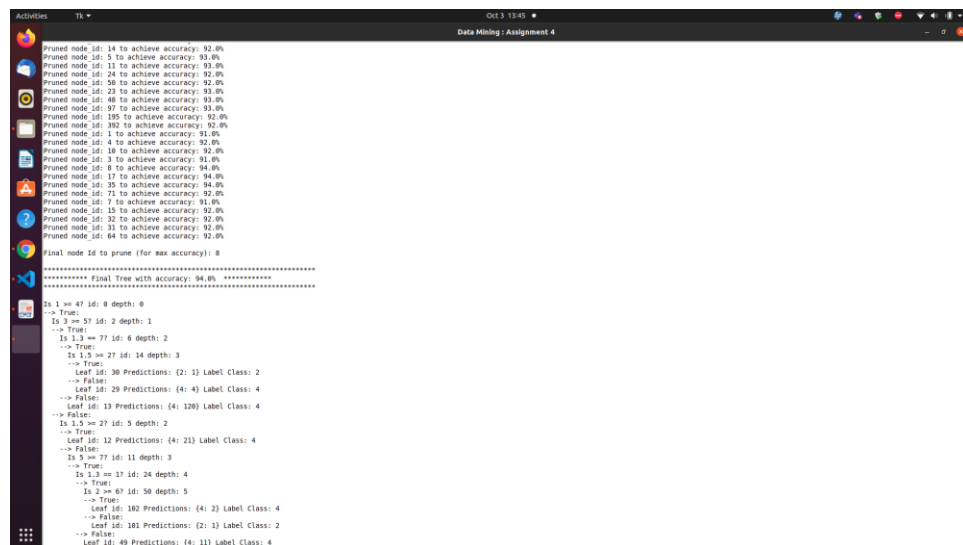
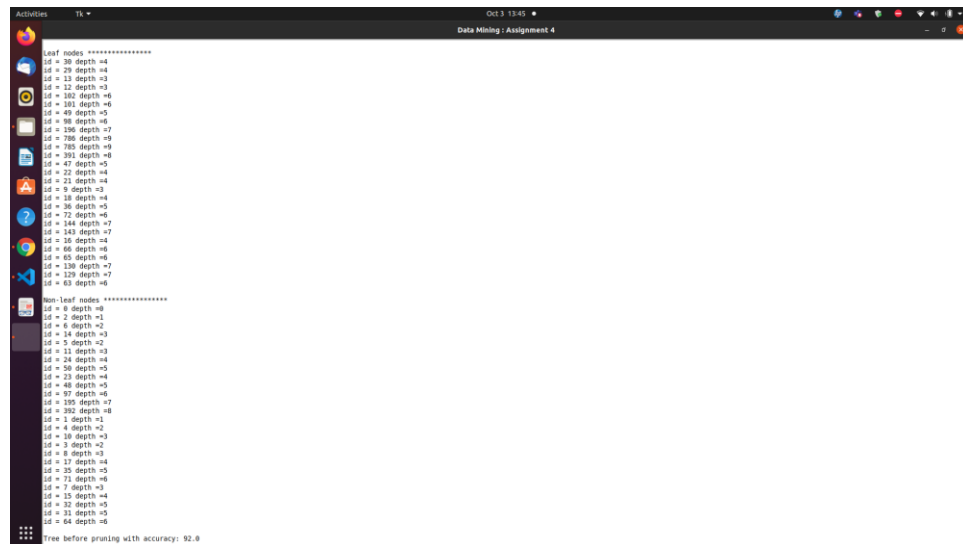
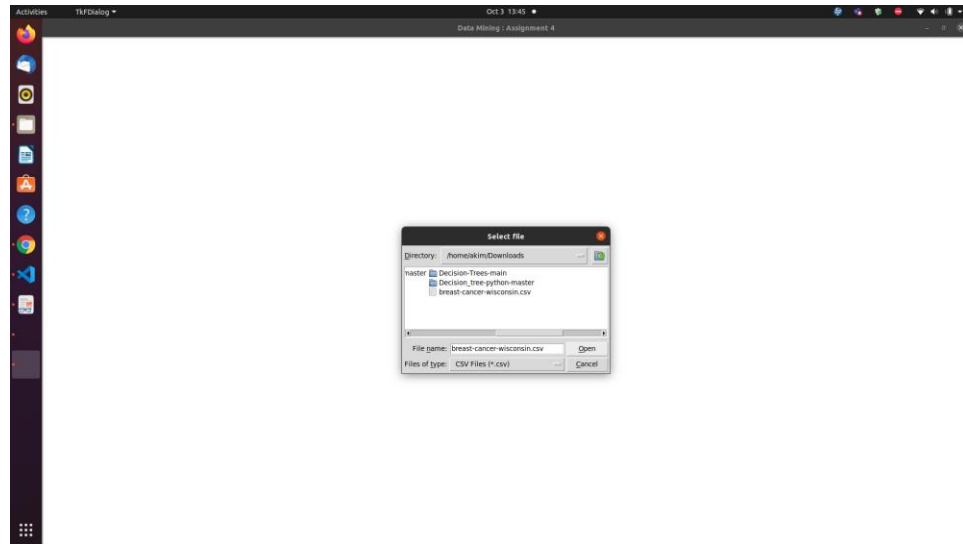
where pos and neg is the number of positive tuples covered by R, respectively.

Note – This value will increase with the accuracy of R on the pruning set. Hence, if the FOIL_Prune value is higher for the pruned version of R, then we prune R.

❖ Procedure

1. Design the rule based classifier : Extract the rules from decision tree build in assignment no. 4.
2. Tabulate the results and evaluate the performance of rules generated using following metrics :
 - a. Coverage
 - b. Accuracy
 - c. Toughness (size)
3. Use the following categorical data sets from UCI machine learning repository :
 - a. Balance Scale data set
 - b. Car evaluation data set
 - c. Breast-cancer data set





```
Activities TK - Oct 3 13:45 Data Mining : Assignment 4

--> False:
Is 1.2 == 27 id: 23 depth: 4
--> True:
Is 1.2 == 187 id: 48 depth: 5
--> True:
Leaf id: 98 Predictions: (4: 4) Label Class: 4
--> False:
Is 1.4 == 67 id: 97 depth: 6
--> True:
Leaf id: 196 Predictions: (2: 2) Label Class: 2
--> False:
Is 3 == 37 id: 195 depth: 7
--> True:
Is 5 == 67 id: 392 depth: 8
--> True:
Leaf id: 786 Predictions: (2: 1) Label Class: 2
--> False:
Leaf id: 785 Predictions: (4: 5) Label Class: 4
--> False:
Leaf id: 393 Predictions: (2: 1) Label Class: 2
--> False:
Leaf id: 47 Predictions: (2: 4) Label Class: 2
Is 1.3 == 17 id: 3 depth: 1
--> True:
Is 5 == 87 id: 4 depth: 2
--> True:
Is 1.5 == 47 id: 18 depth: 3
--> True:
Leaf id: 22 Predictions: (4: 1) Label Class: 4
--> False:
Leaf id: 21 Predictions: (2: 1) Label Class: 2
--> False:
Leaf id: 9 Predictions: (2: 383) Label Class: 2
--> False:
Is 1.4 == 37 id: 3 depth: 2
--> True:
Leaf id: 8 Predictions: (4: 21, 2: 4) Label Class: 4
--> False:
Is 1.3 == 187 id: 7 depth: 3
--> True:
Leaf id: 16 Predictions: (4: 4) Label Class: 4
--> False:
Is 5 == 77 id: 15 depth: 4
--> True:
Is 1.3 == 27 id: 32 depth: 5
--> True:
Leaf id: 66 Predictions: (2: 1) Label Class: 2
--> False:
Leaf id: 65 Predictions: (4: 3) Label Class: 4
--> False:
Is 1.3 == 57 id: 31 depth: 5
--> True:
Is 2 == 27 id: 64 depth: 6
--> True:
Leaf id: 130 Predictions: (2: 5) Label Class: 2
--> False:
Leaf id: 129 Predictions: (4: 1) Label Class: 4
--> False:

Confusion Matrix:
[[ [ 5 3]
  [ 8 38]]]

[[10 9]
 [ 3 31]]

precision recall f1-score support
2 0.77 0.53 0.62 19
4 0.36 0.62 0.45 8
accuracy 0.56 0.58 0.54 27
macro avg 0.65 0.56 0.57 27
```

❖ Conclusion

Thus, implemented and discovered several ways to build decision trees along with rule-based Classification and their metrics such as recall, precision, support etc as well as coverage, accuracy and so on.

❖ References

https://www.tutorialspoint.com/data_mining/dm_rbc.htm

https://link.springer.com/referenceworkentry/10.1007%2F978-0-387-39940-9_559