

# Practical No. 03

Aim: To demonstrate performing classification on data set.

☒ New Project

Project Name:  Project ID:  Project Date:  Project Status:

Project Description:

Project Details:

Attribute	Value	Weight
1. sunny	5	5
2. overcast	4	4
3. rainy	5	5

Classify Data:

Classify Data:  Classify Data:

Log

Date: 15/02/2023



### Practical No. 03

Aim: To demonstrate performing classification on data sets.

Theory:

Classification:

classification is the process for finding a model that describe the data values and concept for the purpose of prediction.

classification in data mining is a common technique that separates data points into different classes.

It allows you to organize data set of all sorts, including complex and large datasets as well as small and simple ones. It primarily involves using algorithms that you can easily modify to improve the data quality. The algorithm establishes the link between the variables for predictions.

The algorithm you use for classification in data mining is called the classifier, and observations you make through the same are called instances.

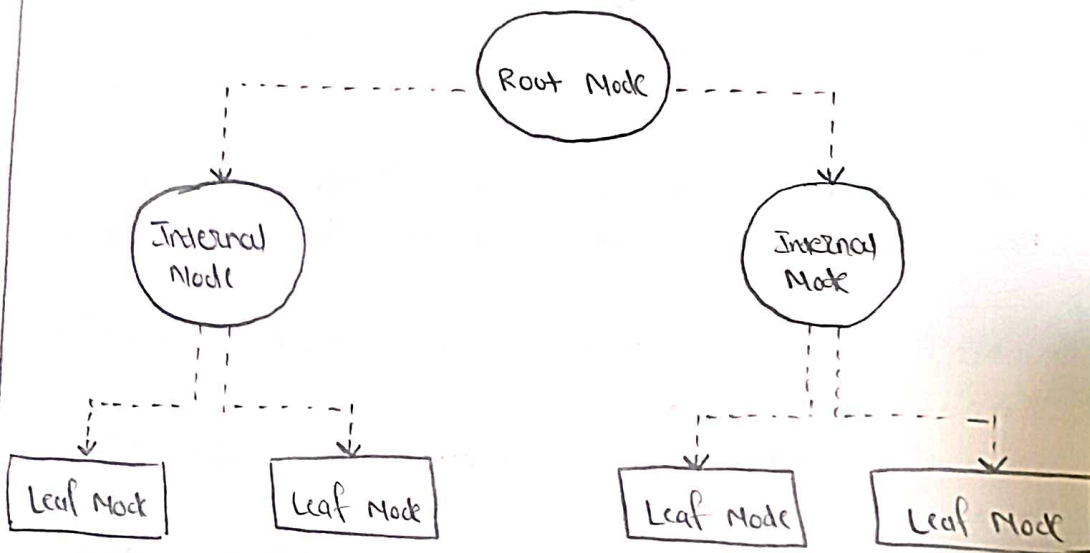
There are multiple type of classification algorithm, each with its unique functionality and application. All of those algorithm are used to extract data from a dataset.

Data Mining Algorithm for Classification:

- Decision Trees
- Logistic Regression
- Naive Bayes Classification
- k - nearest neighbors
- Support Vector Machine.

Decision Tree:

A decision tree is a structure that includes a root node, branches and leaf nodes. Each internal node depicts a test on an





Date :



attribute, each branch denotes the outcome of a test and each leaf node holds a class label. The topmost node in the tree is the root node.

A decision tree is a classification scheme to generate a tree consisting of root node, internal nodes and external nodes. Root nodes representing the attributes. Internal nodes are also the attributes. External nodes are the class and each branch represents the values of the attributes.

Decision tree also contains set of rules for a given data set; there are two subsets in decision tree, one is a training data set and second one is a testing data set. Training data set is previously classified data. Testing data set is newly generated data.

The benefits of having a decision tree are as follows :-

- It does not require any domain knowledge.
- It is easy to comprehend.
- The learning and classification steps of a decision tree are simple and fast.
- J48 classification and its decision tree
- C4.5 algorithm / J48
- The C4.5 algo. is a classification algo. which provides decision trees based on information theory. It is an extension of Ross Quinlan's earlier ID3 algorithm also known in weka as J48, J standing for Java. The decision trees generated by C4.5 are used for classification, and for this reason, C4.5 is often referred to as a statistical classifier.
- The J48 implementation of the C4.5 algo. has many additional features including accounting for missing values, decision trees pruning, continuous attribute value range, derivation of rules and etc..
- In WEKA, data mining tool, J48 is an open-source Java



Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C4.5 M2

Test options

Use training set

Supplied test set

Choose validation folds 10

Percentage split 10

More options

Prune options

Start

Result list (right-click for options)

Classifier output

```

outlook = sunny
outlook = rainy
1 windy = TRUE no (2.0)
1 windy = FALSE yes (3.0)
Number of Leaves = 5
Size of the tree = 8
Time taken to build model: 0 seconds

```

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances 7 50 %

Incorrectly Classified Instances 2 50 %

Kappa statistic -0.0426

Mean absolute error 0.6147

Root mean squared error 0.5584

Relative absolute error 87.5 %

Root relative squared error 121.2987 %

Total Number of Instances 14

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.500	0.500	0.500	0.500	0.500	-0.043	0.433	0.750	yes
	0.500	0.500	0.500	0.500	0.500	-0.043	0.433	0.437	no
Weighted Avg.	0.500	0.500	0.500	0.500	0.500	-0.043	0.433	0.437	

=== Confusion Matrix ===

a b <-- classified as

5 4 | a = yes

3 2 | b = no

Status OK

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 - C4.5 M2

Test options

Use training set

Supplied test set

Choose validation folds 10

Percentage split 10

More options

Prune options

Start

Result list (right-click for options)

19/12/03: Vecs343

Classifier output

```

outlook = sunny
outlook = rainy
1 windy = TRUE no (2.0)
1 windy = FALSE yes (3.0)
Number of Leaves = 5
Size of the tree = 8
Time taken to build model: 0 seconds

```

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances 7 50 %

Incorrectly Classified Instances 2 50 %

Kappa statistic -0.0426

Mean absolute error 0.6147

Root mean squared error 0.5584

Relative absolute error 87.5 %

Root relative squared error 121.2987 %

Total Number of Instances 14

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.500	0.500	0.500	0.500	0.500	-0.043	0.433	0.750	yes
	0.500	0.500	0.500	0.500	0.500	-0.043	0.433	0.437	no
Weighted Avg.	0.500	0.500	0.500	0.500	0.500	-0.043	0.433	0.437	

=== Confusion Matrix ===

a b <-- classified as

5 4 | a = yes

3 2 | b = no

Weka Classifier Tree Visualizer: 19/12/03 - trees.J48 (recursive symbolic)

Tree View

```

graph TD
    outlook((outlook)) -- sunny --> humidity((humidity))
    outlook -- overcast --> yes40[yes 4/0]
    outlook -- rainy --> wind((wind))
    humidity -- high --> no31[no 3/1]
    humidity -- normal --> yes20[yes 2/0]
    wind -- TRUE --> yes20
    wind -- FALSE --> yes10[yes 1/0]

```

Status OK



Date :



implementation of the C4.5 algorithm. It allows classification via either decision trees or rules generated from them.

Generating a decision tree from training tuples of data partition  $D$   
Algorithm : generate-decision-tree

Input :

Data partition  $D$ , which is a set of training tuples and their associated class labels.

attribute-list, the set of candidate attributes.

Attribute selection method, a procedure to determine the splitting criterion that best partitions the data tuples into individual classes. This criterion includes a splitting-attribute and either a splitting point or splitting subset.

Output :

A Decision Tree

Method :

create a node  $N$ ;

if tuples in  $D$  are all of the same class,  $c$  then  
return  $N$  as leaf node labeled with class  $c$ ;

if attribute-list is empty then  
return  $N$  as leaf node with labeled  
with majority class in  $D$ ; // majority voting

apply attribute-selection method ( $D$ , attribute-list)  
to find the best splitting-criterion;  
label node  $N$  with splitting-criterion;

if splitting-attribute is discrete-valued and  
multiway split allowed then // no restriction to binary trees.

Result :

Then, the classification on data set is performed by decision tree (J48) method.



Date :



attribute - list = splitting attribute // remove splitting attribute  
for each outcome  $j$  of splitting criterion

// partition the tuples and grow subtrees for each partitions  
let  $D_j$  be the set of data tuples in a satisfying outcome  $j$ ;  
// a partition

if  $D_j$  is empty then

attach a leaf labeled with the majority  
class in  $D$  to node  $N_j$ ;

else

attach the node returned by generate  
decision tree ( $D_j$ , attribute list) to node  $N_j$ ;

end for

return  $N_j$ ;

Result:

Thus the classification on data set is performed by decision  
tree (J48) method.

Viva Question :

① what is classification?

→ Classification is the process for finding a model that describe  
the data values and concept for the purpose of prediction.

② what is the need of classification?

→ The need of classification is to accurately predict the target  
class for each case in the data, it allows you to organize  
data set of all sorts, including complex and large datasets as  
well as small and simple ones.





Q.3 What are the different methods of classification?

→ Logistic regression, Naive Bayes, Decision tree, k-Nearest Neighbours, Support vector machines, Bayes classifiers, Function classifier, Lazy classifier, meta classifier and so on.

Q.4 What are the advantages of a decision tree classifier?

- ① It is easy to comprehend
- ② It does not require any domain knowledge.
- ③ The learning and classification steps of decision tree are simple and fast
- ④ Less data preparation, Non-parametric versatility, Non-Linearity.