

**DEPARTMENT OF COMPUTER SCIENCE  
INSTITUTE OF MANAGEMENT AND RESEARCH, JALGAON**

Name Gadilohar Hitesh Digambar  
 Expt. Title Implement the find-s Algorithm  
 Class F.Y.MCA Batch \_\_\_\_\_ Performed on \_\_\_\_\_  
 Roll No. 38 Expt. No. 07 Submitted on \_\_\_\_\_  
 Remarks \_\_\_\_\_ Returned on \_\_\_\_\_

\* Find S Algorithm :-

The find s algorithm is a basic concept learning algorithm in machine learning. The find-s algorithm finds the most specific hypothesis that fits all the positive examples. We have to note here that the algorithm Consider only those positive training example.

Algorithm :-

Step 1 :- Initialize h to the most Specific hypothesis in H

Step 2 :- for each positive training instance x  
 for each attribute Constraint a in h If the  
 Constraint a is Satisfied by x then do  
 nothing.

else replace a in h if the next more  
 general Constraint that is Satisfied by x

Step 3 :- Output hypothesis h.

Q)

Example	Color	Toughness	Fungus	Appearance	Poisonous
1	Green	Hard	No	Wrinkled	Yes
2	Green	Hard	Yes	Smooth	No
3	Brown	Soft	No	Wrinkled	No
4	Orange	Hard	No	Wrinkled	Yes
5	Green	Soft	Yes	Smooth	Yes

Delete for :  
 algorith  
 Chart  
 ramme Listing  
 ults  
 iments

First we consider the hypothesis to be a more specific hypothesis Hence our hypothesis would be

$$h = \{ \phi, \phi, \phi, \phi, \phi, \phi, \phi \}$$

Consider example 1 :-

$$h = \{ \text{Green, Hard, No, Wrinkled} \}$$

Consider example 2 :-

$$h = \{ \text{Green, Hard, No, Wrinkled} \}$$

Consider example 3 :-

$$h = \{ \text{Green, Hard, No, Wrinkled} \}$$

Have we see that above examples has a negative outcome Hence we neglect this example & our hypothesis remains the same.

Consider Example 4 :-

$$h = \{ ?, \text{Hard, No, Wrinkled} \}$$

Consider Example 5 :-

$$h = \{ ?, ?, ?, ? \}$$

Hence for the given data find final hypothesis would be:

Final Hypothesis :  $h = \{ ?, ?, ?, ?, ? \}$

DEPARTMENT OF COMPUTER SCIENCE  
INSTITUTE OF MANAGEMENT AND RESEARCH, JALGAON

Name Gadilohar Hitesh Digambar

Expt. Title Implement the Candidate-Elimination Inductive learning algorithm

Class BT.MCA Batch \_\_\_\_\_ Performed on \_\_\_\_\_

Roll No. 38 Expt. No. 02 Submitted on \_\_\_\_\_

Remarks \_\_\_\_\_ Returned on \_\_\_\_\_

### \* Candidate - Elimination Algorithm :-

The Candidate elimination algorithm incrementally builds the version Space given a hypothesis Space  $H$  & a Set  $E$  of examples.

Algorithm :-

Step 1:- Load Data Set.

Step 2:- Initialize General Hypothesis & Specific Hypothesis.

Step 3:- for each training example.

Step 4:- if example is positive example.

    if attribute.value == hypothesis-value:

        Do nothing.

    else:

        replace attribute value with '?'

(Basically generalizing it)

Step 5:- if example is negative example make generalizing hypothesis more Specific.

Example:-

Consider the dataset given below:

Sky	Temperature	Humid	Wind	Water	forest	Output
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

Initially:-

$$G = [[?, ?, ?, ?, ?, ?], [?, ?, ?, ?, ?, ?], [?, ?, ?, ?, ?], [?, ?, ?, ?, ?, ?], [?, ?, ?, ?, ?, ?]]$$

$$S = [\text{Null}, \text{Null}, \text{Null}, \text{Null}, \text{Null}, \text{Null}]$$

for instance 1:- < 'Sunny', 'Warm', 'Normal', 'Strong', 'Same' >  
+ positive Output.

$$G_1 = G$$

$$S_1 = ['Sunny', 'Warm', 'Normal', 'Strong', 'Same']$$

for instance 2:-

$$< 'Sunny', 'Warm', '?', 'Strong', 'Warm', 'Same' >$$

for instance 3:-

$$< 'Rainy', 'Cold', 'High', 'Strong', 'Warm', 'Change' >$$

$$G_3 = [[\text{Sunny}, ?, ?, ?, ?, ?], [?, \text{Warm}, ?, ?, ?, ?], [?, ?, ?, ?, ?, ?]]$$

$$\text{Output} [?, ?, ?, ?, ?, ?] [?, ?, ?, ?, ?, ?] [?, ?, ?, ?, ?, ?]$$

$$S_3 = S_2$$

for instance 4:- < 'Sunny', 'Warm', 'High', 'Strong', 'Cool' >

$$G_4 = G_3$$

$$S_4 = ['Sunny', 'Warm', '?', 'Strong', '?', '?']$$

At last, by Synchronizing the G4 & G5 algorithm  
Output:-

$$G = [[\text{Sunny}, ?, ?, ?, ?, ?], [?, \text{Warm}, ?, ?, ?, ?]]$$

$$S = ['Sunny', 'Warm', '?', 'Strong', '?', '?']$$

**DEPARTMENT OF COMPUTER SCIENCE  
INSTITUTE OF MANAGEMENT AND RESEARCH, JALGAON**

Name Gadilohar Hitesh Digambar  
Expt. Title Write a program to implement Decision tree using python programming language of your choice.  
Class BT MCA Batch \_\_\_\_\_ Performed on \_\_\_\_\_  
Roll No. 38 Expt. No. 03 Submitted on \_\_\_\_\_  
Remarks \_\_\_\_\_ Returned on \_\_\_\_\_

\* What is decision trees?

Decision tree is a Supervised learning technique that can be used for both classification & regression problems, but mostly it is preferred for solving classification.

program. It is a tree structure classifier where internal nodes represent the features of a dataset. branches represent the decision rules & each leaf node represent the outcome.

It is a graphical represent of for getting all the possible solution to problem / decision based on given conditions. This works surprisingly well in most of class.

\* ID3 Algorithm:-

- If all examples have same label:-
  - return a leaf with that label
- Else if there are no feature leaf to test:-
  - return a leaf with the most common label
- Else:-
  - chose the feature  $F$  that maximise the information gain of  $S$  to be the next node using function  $IG(F, S)$ .
    - add branch from the node for each possible value  $f$  in  $F$ .

- For each Branch:-

- \* Calculate SF by removing  $\vec{f}$  from the set of
- \* recursively call the algorithm with SF to compute the gain relative to the current set of each

\* formula of Information Gain :-

$$1) \text{ Entropy} := -p \log_2 p - q \log_2 q.$$

To build a decision tree, we can need to calculate two types of entropy using frequency tables:

i) Entropy Using the frequency table of one

$$E(s) = \sum_{i=1}^c -p_i \log_2 p_i$$

ii) Entropy Using the frequency table of two attr.

$$E(T, x) = \sum_{C \in X} p(c) \cdot E(c)$$

$$2) \text{ Gini} = 1 - \sum_{i=1}^n (p_i)^2$$

$$3) \text{ Information Gain} = 1 - \text{Entropy}.$$

**DEPARTMENT OF COMPUTER SCIENCE  
INSTITUTE OF MANAGEMENT AND RESEARCH, JALGAON**

Name Gadilohar Hitesh Digambar

Expt. Title White Program to calculate popular Attribute.

Class MCA I Batch B1 Performed on \_\_\_\_\_

Roll No. 38 Expt. No. 04 Submitted on \_\_\_\_\_

Remarks \_\_\_\_\_ Returned on \_\_\_\_\_

### Attribute Selection Measure: (ASM)

The best attribute or feature is Selected using the attribute Selection measure (asm) the attribute selection is the root node lecture.

Asm is a technique used for the selecting best attribute for discrimination among tuples, it give rank to each attribute & the best attribute is Selected as splitting criterion.

The most popular methods of Selection one:-

#### 1] Information Gain:-

Information gain is a decision in decision in entropy Decision tree none use of information gain will entropy to determine which feature to split into nodes to get closer to predicting the target & also to determine when to stop splitting.

$$\text{Information Gain} = \text{Entropy}(S) - [\text{Weighted Avg}]$$

\* Entropy (each features).

#### 2) Gini Index:-

Gini index is a measure of impurity or purity used while creating a decision tree in the CART model can be calculated using the below formula

$$\text{Gini index} = 1 - \sum p_i^2$$

- complete for:  
Algorithm  
Flow Chart  
Programme Listing  
Results  
Comments

**DEPARTMENT OF COMPUTER SCIENCE  
INSTITUTE OF MANAGEMENT AND RESEARCH, JALGAON**

Name	<u>Gadilohar Hitesh Digambar</u>		
Expt. Title			
Class	<u>F.Y. MCA</u>	Batch	
Roll No.	<u>38</u>	Expt. No.	<u>05</u>
Remarks			
	Performed on _____		
	Submitted on _____		
	Returned on _____		

### \* K - Nearest Neighbour (KNN) Algorithm :-

K - Nearest Neighbour is one of the simple machine learning algorithm based on Supervised learning technique. In NN algorithm store all the available data & classification a new data point based on the Similarity. This mean when new data appears, then it can be easily classified into a well suit Category by using K-NN algorithm.

K-NN is a non-parametric algorithm which means it does not make any assumption on underlying data. It is a lazy learning algorithm where all computation is deferred until classification.

### \* Algorithm:-

- The KNN Classification is performed using the following four steps

Compute the distance metric bet<sup>n</sup> the test data point & all the labeled data point.

- Order the labeled data point in the increasing order of the distance metric.

for:  
Q Select the top k labeled data point & look at the class tables.

- find the class label that the majority of these labeled data points have & assign it to the test data point.

\* Distance Calculation formula :-

1) Euclidean Distance:-

It is generally used to find the distance b/w real-valued vectors.

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

2) Manhattan Distance:-

This is simple way or technique to calculate distance b/w two point of ten called Taxicab dis. or city block distance.

Manhattan Distance = Sum for i to N

$$\sum |x_{1i} - x_{2i}|$$

3) Hamming distance :-

The hamming distance is mostly used in processing or having the boolean vector Boolean.

Mean the data is in the form of binary.

Hamming distance ( $x_1, x_2$ ) =

4) Minkowski Distance:-

Minkowski distance is a generalization of euclidean & manhattan distance.

$$||x_1 - x_2|| = \left( \sum_{i=1}^n |x_{1i} - x_{2i}|^p \right)^{1/p}$$