

## ML UNIT-1 IMPORTANT QUESTIONS

Q) ~~Explain about inductive bias in decision tree learning.~~

Ans:-

### INDUCTIVE BIAS IN DECISION TREE LEARNING:

- An Inductive Bias is a set of assumptions, which along with training data, deductively justify the classification assigned by the learner to future instances.
- For a given collection of training ex. there are many decision trees consistent with it. Using Inductive bias, we select one of these consistent hypothesis over others.
- ID3 chooses, the 1<sup>st</sup> acceptable tree it encounters in simple-to-complex, Hill climbing search through the space of possible trees.
- It aims to select shorter trees over longer & trees that place the attributes which have highest info.gain, closest to the root.
- Example, in BFS-ID3 (breadth 1<sup>st</sup> search), it 1<sup>st</sup> consider all trees with depth 1 & later depth 2. It tries to search for decision tree consistent with training data & returns the smallest/consistent tree at that search depth.

## → Restriction Biases & Preference Biases:-

### IDB search

- \* searches complete hypothesis space.
- \* one capable of expressing any finite discrete valued func.
- \* searches Incompletely through this space.
- \* searches from simple to complex hypotheses until a termination condition is met (i.e. it finds a hypothesis consistent with the data).
- \* its Inductive bias is solely a consequence of hypotheses by its search strategy.
- \* Its hypothesis space introduces no additional bias.
- \* Inductive bias of IDB follows from its search strategy.
- \* Inductive bias of IDB is a Preference for certainty hypothesis over others. Thus it is called Preference Bias or Search Bias.  
Preference bias is more desirable bcoz it allows the learner to work within a complete hypothesis space i.e. assumed to contain the unknown Target func.
- \* Adv:- ~~more~~ more desirable bcoz it allows the learner to work within a complete hypothesis space i.e. assumed to contain the unknown Target func.

### Candidate elimination Algo.

- \* Searches Incomplete hypothesis space (one that can express only a subset of potentially teachable concepts).
- \* Searches completely through this space.
- \* It finds every hypothesis consistent with the data (training).
- \* Its Inductive bias is solely a consequence of expressive power of its hypothesis rep.  
Search strategy
- \* Its ~~hypothesis~~ introduces no additional bias.
- \* Its Inductive bias follows from its def<sup>n</sup> of Search space.
- \* Bias of this Algo. is in the form of a Categorical Restriction on the set of hypotheses considered. Thus it is called Restriction Bias or Language Bias.  
Disadv:- or less desirable  
It is undesirable bcoz it strictly limits the set of potential hypotheses. It introduces the possibility of excluding unknown Target func. Altogether.

→ why prefer short Hypotheses?

Accdg to Occam's Razor, prefer the simplest hypothesis that fits the data.

- Because, there are fewer short hypotheses than long ones but it can be complex as well and may fail to generalize correctly to subsequent data.
- Ex:- In Decision Tree hypotheses, Assume 2 training examples, A 5 node tree with a less likely statistical coincidence is preferred over 500 node hypothesis.
- Another problem with Occam's Razor, stmt. is the size of hypothesis. It is determined internally by the leaves. Two leaves with 2 diff. internal rep's. can occur at diff. hypotheses.
- Solution will be to have a common process of evolution that shall create an internal rep.

Q8) DESCRIBE IN DETAIL, STEPS FOR DESIGNING A LEARNING SYSTEM.

Ans:-

a. Designing a Learning system:

- Before understanding design issues & design approaches, first, we shall understand, to design a problem.
- example used: GOAL: Entering world checker tournament  
TASK: Learn to play checkers.  
PERFORMANCE MEASURE : % of games it wins in the tournaments

## → Choosing the Training Experience :-

- choose a type of training experience from which our sys will learn. It's v. imp for the success/failure of the learner.

(1) it key feature to select it will be to see if training experience will provide direct/indirect feedback.

ex:- Direct - Individual checkers board states & correct move for each.

Indirect - Move seqs. & final outcomes of various games played

- The disadv. of this is "credit assignment"

- Learning from direct feedback is easier.

## ② Degree to which learner controls the seq. of training examples.

- ex: (i) learner dependent on teacher to ~~self~~ select informative boards & to provide correct move for each.

(ii) learner can propose board states & if finds a particular move confusing then it can ask teacher.

(iii) Learner can have control on both board states & training classifications. & there is no teacher present.

## ③ Show well it. rep. distribution of examples over which the final sys. performance P must be measured:

- generally, if future test example & distribution of learning examples (train) are similar then learning is more reliable.

- In example of checkers game, If training experience is, i.e. games & played against itself & performance measure 'P' is % of games won in world tournament then there is a danger that training experience will not be similar to the situation where it will be tested in future.

- But for simplicity, we still assume a fully specified task as:

(Performance) P: % of games won in world Tournament

(Training experience) E : games played against itself

- Foll. steps are needed next to complete the design:
  - (i) exact type of knowledge to be learned.
  - (ii) "Rip" for this target knowledge
  - (iii) learning mechanism.

### choosing the target function :-

- we shall now determine the Target knowledge & how will it be used.
- ex; take a checkers playing prog. Assume, it can generate "legal" moves from any board state. the "best" move among them lots say, prog needs to learn only the "best" move among them.
- Assume "chooseMove" be a func. which chooses best move for any given board state.
- "chooseMove :  $B \rightarrow M$ " is the notation used to accept any board as i/p from a set of "legal board states"  $B$  and produces " $M$ ", which is an o/p move from a set of legal move  $M$ . But this func. is v. difficult to learn since it is an "indirect" training exercise.
- Alternative Target func.: Let say, an evaluation func.  $V$  assigns a numerical score to any board state. So  $V : B \rightarrow R \Rightarrow V$  well map any legal board state from set  $B$  to some real value. The better board state the greater shall be " $R$ " value.
- This method is comparatively easier for the sys. to learn and select best move from any board state.

example:- if  $b$  is a final board state (i.e won  $\Rightarrow V(b) = 100$ )  
 if  $b$  " " " " " " " " i.e lost  $\Rightarrow V(b) = -100$   
 if  $b$  is " " " " " " " " i.e Drawn  $\Rightarrow V(b) = 0$   
 if  $b$  is not a final state in game  $\Rightarrow V(b) = V(b')$   
 where  $b'$  = best final board state that can be achieved

## → choosing a rep<sup>n</sup> for Target func:-

- We've got the ideal target func "V". Now we must choose  $\hat{V}$  via function approximation. Because learning algo. can't acquire ideal or perfect target func "V". Rather, they acquire only some approx<sup>n</sup> to the target func. so that's why we shall choose  $\hat{V}$  now.
- For the case of checkers game  $\hat{V}$  can be:
  - (i) collection of rules that match against features of board state
  - (ii) Large table with distinct entry specifying the value for each distinct board state.
- For expressing  $\hat{V}$  which is more similar to V, one needs a large amount of training data.
- $\hat{V}$  will be calculated as a linear comb<sup>n</sup> of foll. board features:
  - (i)  $x_1$ : no. of black pieces on board
  - (ii)  $x_2$ : " " red "
  - (iii)  $x_3$ : " " black kings "
  - (iv)  $x_4$ : " " Red kings "
  - (v)  $x_5$ : no. of black pieces threatened by red.
  - (vi)  $x_6$ : " " red " " " black "
- So, learning progr. is
 
$$\hat{V}(b) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + w_6 x_6$$

where,  $w_0, w_1, \dots, w_6$  = weights chosen by algo.
- weights determine the relative importance of various board features in finding out value of board  
 $w_0$  adds an additive constant to the board value.  
 After this step now, we have a partial design of checkers learning program.

(6)

- ① Task T: playing checkers
- ② performance measure P: % of games won in World tournaments
- ③ Training Experience E: games played against itself.

- ④ Target func.:  $V: \text{Board} \rightarrow \mathbb{R}$
- ⑤ Target func.:  $V(b) = w_0 + w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_5x_5 + w_6x_6$   
Ripn

### → choosing a func. Approx<sup>n</sup> Algo.:-

- For learning the Target func.  $\hat{V}$ , a set of training examples is needed, i.e.  $V_{\text{train}}(b)$ .
- Each example will describe a specific board state  $b$  & training value for it i.e.  $V_{\text{train}}(b)$ .
- Consider an example, where, black has won game.  
then  $x_5 = 0 \rightarrow$  there are no red pieces on board & value of  $V_{\text{train}}(b) = +100$ .
- So the ordered pair of training example will be in the form  $\langle b, V_{\text{train}}(b) \rangle$ .  
For above example taken it is  
 $\langle x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0, +100 \rangle$

### → Estimating Training Values:-

The Only Training info. that's available is whether game was won/lost  
So we can assign Training values v. easily.

- But, sometimes there are "Intermediate" board states that occur before game ends i.e. we can't say whether game path was good/bad or whether a subsequent poor move made the game to lose it?  
→ There is a simple approach for such

→ Intermediate board states:  
"Assign training value of  $v_{train}(b)$  for any intermediate board state 'b' to be  $\hat{v}(\text{successor}(b))$ ".  
i.e.  $v_{train}(b) \leftarrow \hat{v}(\text{successor}(b))$ .

where  $\hat{v}$  = Learner's current approx. to  $v$ .  
 $\text{successor}(b)$  = next board state following b for which  
it is again the prog's turn to move.

### → Adjusting Weights:

- we must choose weights  $w_i$  so that it fits best to the set of training examples  $\{(b, v_{train}(b))\}$
- Best fit:- First define best hypothesis or a set of weights which minimises  $E(\text{sq. error})$  b/w training values & predicted values by  $\hat{v}$ .

$$\Rightarrow E = \sum_{(b, v_{train}(b)) \in \text{Training Examples}} (v_{train}(b) - \hat{v}(b))^2.$$

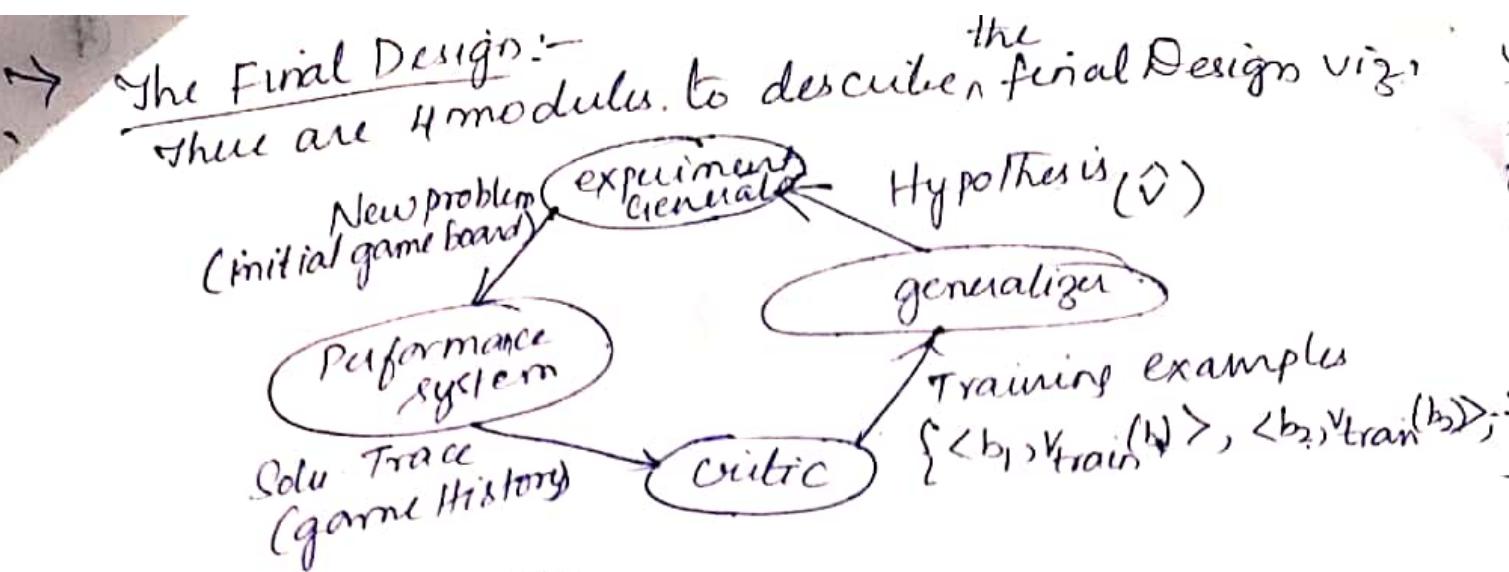
- There are many algs. for finding weights of a linear func. that minimises  $E$ . One of them is LMS (Least Mean Squares) Training Rule.

### → LMS Weight update Rule:

- For each training example  $(b, v_{train}(b))$ 
  - use current wt. to calculate  $\hat{v}(b)$ .
  - For each wt  $w_i$ , update

$$w_i \leftarrow w_i + \eta (v_{train}(b) - \hat{v}(b)) x_i$$

where  $\eta$  = small const(0.1), which modulates size of wt. update  
if  $v_{train}(b) - \hat{v}(b) = \begin{cases} +ve \Rightarrow \text{new wt. is prop. to val diff.} \\ -ve \Rightarrow \text{"No wt. change."} \\ 0 \end{cases}$



### (i) Performance sys:-

- solves given performance Task
- In this ex., it is playing the checkers game, using Learned Target function.
- Takes new game as i/p & produces a part of its solu. i.e. game history as o/p
- Strategy for performance sys:
  - select its next move at each step determined by its evaluation func.  $\hat{v}$ .
  - As  $\hat{v} \uparrow$  so  $\Rightarrow$  performance  $\uparrow$

### (ii) Critic:-

- i/p: History / part or trace of solu.
- o/p: set of training examples of the target func.
- Each training ex., is nothing but a game state  $b$ , and its estimate of target func.  $v_{train}(b)$ .

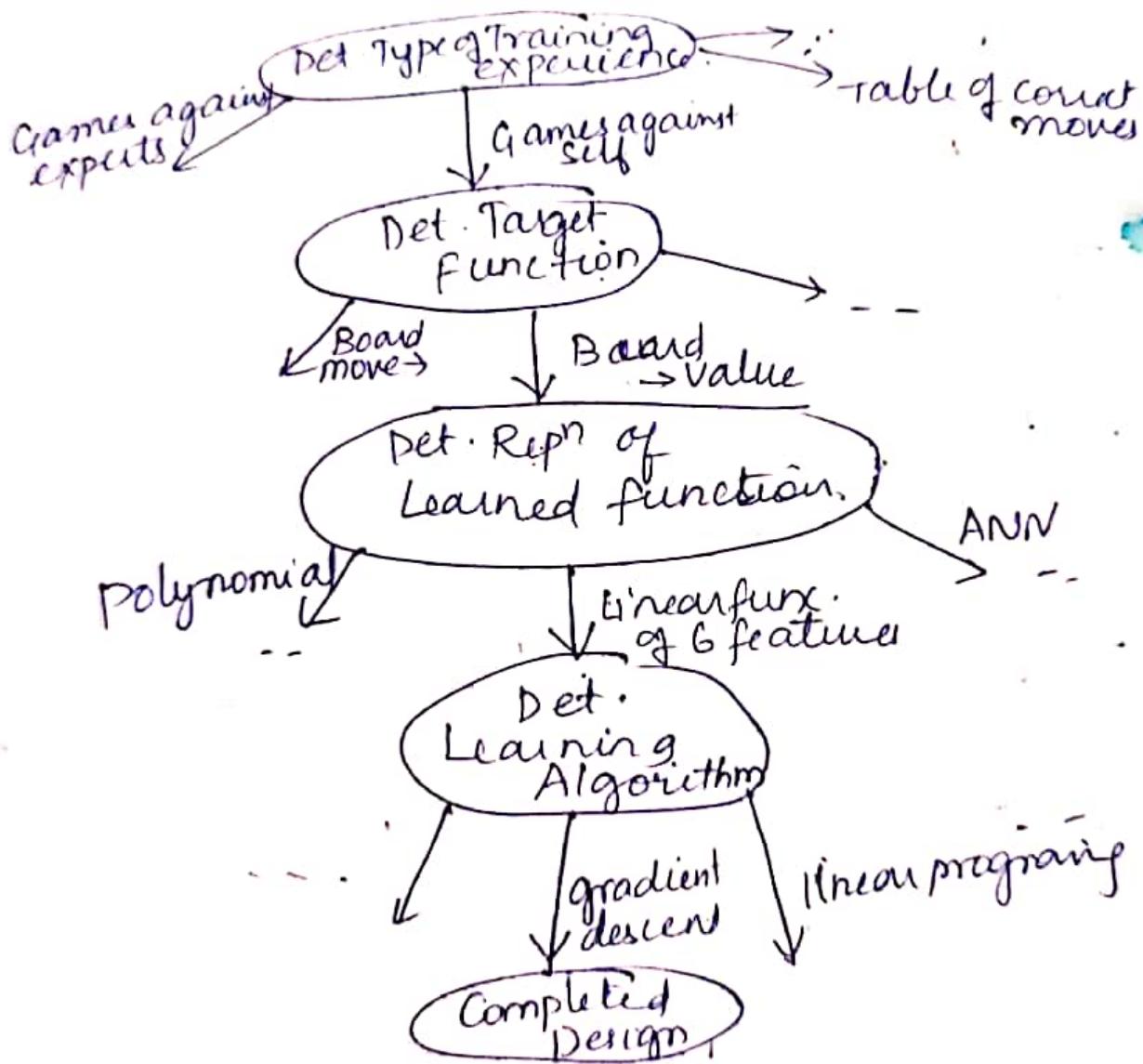
### (iii) Generalizer:-

- i/p: Training examples
- o/p: Hypothesis  $\hat{v}$  = estimate of Target func.
- It takes specific training examples & generalises a general func. (Hypothesis) that covers these examples.
- In this ex., it has followed LMS algo. It uses learned weights  $w_0, w_1, \dots, w_6$  to describe the o/p hypothesis  $\hat{v}$ .

#### (iv) Experiment evaluation

i/p: current Hypothesis (*i.e.* initial board state) for the O/P: new problem (*i.e.* initial board state) for the performance *etc.* to explore.  
Its role is to pick new practice problems in order to maximize learning rate of overall system.

- The seq. of design choices made for the checkers program summarised below:
- Assumptions: -  
- Single linear evaluation func.  
- 6 specific board features



(Q3) DISCUSS ABOUT FIND-S ALGORITHM IN FINDING A MAXIMUM SPECIFIC HYPOTHESIS.

### FIND-S: Finding a Maximally Specific Hypothesis:

- It is an app'n of more general than partial ordering
- Here, we begin with most specific possible hypothesis in  $H_0$  & generalize it, when it fails to cover an observed +ve training example (classifies as +ve)

#### FIND-S Algorithm:

1. Initialise  $h$  to most specific hypothesis in  $H_0$ .

2. For each +ve Training instance  $x$ .  
    - For each attribute constraint  $a_i$  in  $h$ .  
        if constraint  $a_i$  is satisfied by  $x$   
            then do nothing  
        else  
            Replace  $a_i$  in  $h$  with next more general constraint satisfied by  $x$ .

3. O/P hypothesis  $h$ .

#### Example: -

- consider "Enjoy Sport task".

(i) Initialise  $h$  to most specific hypothesis in  $H_0$ .

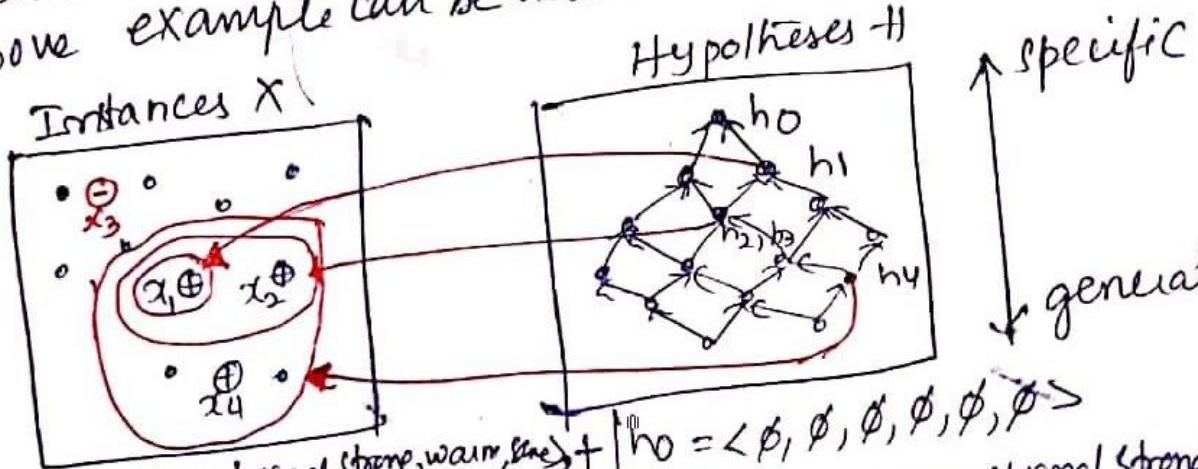
$$h \leftarrow \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$$

(ii) From table 2.1, none of constraints are satisfied by this so it is replaced by next more general constraint that fits this example.

$$h \leftarrow \langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle$$

[It's too specific b'coz all instances are -ve except one +ve Training example, (See Table 2.1 to relate)]

- The second training example forces algo. to further generalise  $h$ , which can be done by substituting any attribute value in ' $h$ ' by a '?'.
- $$h \leftarrow \langle \text{sunny}, \text{warm}, ?, \text{strong}, \text{warm}, \text{same} \rangle$$
- Upon encountering a -ve Training example (3rd one), Find-S algo simply ignore it. So no revision of ' $h$ ' is needed.
- the 4th example (+ve) leads to further generalisation of ' $h$ '.
- $$h \leftarrow \langle \text{sunny}, \text{warm}, ?, \text{strong}, ?, ? \rangle$$
- In this Algo., "More-general-than" partial ordering is thus used to organise the search for an acceptable hypothesis.
- At each step, the hypothesis is generalised as far as necessary to cover the new +ve examples thus making the hypothesis more consistent with Training examples. Hence it is named as "Find-S" Algorithm.
- Above example can be illustrated further with foll. fig.



$$\begin{aligned}
 x_1 &= \langle \text{sunny}, \text{warm}, \text{normal}, \text{strong}, \text{warm}, \text{same} \rangle + h_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle \\
 x_2 &= \langle \text{sunny}, \text{cootn high}, \text{strong}, \text{warm}, \text{same} \rangle + h_1 = \langle \text{sunny}, \text{warm}, \text{Normal}, \text{strong}, \text{warm}, \text{same} \rangle \\
 x_3 &= \langle \text{rainy}, \text{cold}, \text{high}, \text{strong}, \text{warm}, \text{cha} \rangle - h_2 = \langle \text{sunny}, \text{warm}, ?, \text{strong}, \text{warm}, \text{same} \rangle \\
 x_4 &= \langle \text{sunny}, \text{warm}, \text{High}, \text{strong}, \text{cootday} \rangle + h_3 = \langle \text{sunny}, \text{warm}, ?, \text{strong}, \text{warm}, \text{same} \rangle \\
 &\quad + h_4 = \langle \text{sunny}, \text{warm}, ?, \text{strong}, ?, ? \rangle
 \end{aligned}$$

- "Find -S" is guaranteed to o/p most specific hypothesis within  $H$  that is consistent with +ve training example. Finally, it is also consistent with -ve examples provided the correct target concept is contained in  $H$  & also provided the training examples are correct.
- few voids in this Algo. in the form of Questions are:
  - \* Has learner converged to correct target concept?
  - \* Why prefer most specific hypothesis?
  - \* Are Training examples consistent?
  - \* What if there are several maximally specific consistent hypotheses?

Q4. EXPLAIN IN DETAIL "CANDIDATE ELIMINATION" LEARNING ALGORITHM.

A) Candidate Elimination Learning Algorithm

- Initialise  $G$  to set of maximally general Hypotheses int "specific"
  - " S " " "
  - For each training example  $d$ , do
    - If  $d$  is positive example
      - Remove from  $G$  any hypothesis inconsistent with  $d$ .
      - For each hypothesis  $h$  in  $S$  that's inconsistent with  $d$ 
        - Remove  $h$  from  $S$ .
        - Add to  $S$  all minimal generalization of  $h$  such that
          - $h$  is consistent with  $d$
          - some member of  $G$  is more general than  $h$ .
        - Remove from  $S$  any hypothesis i.e. more general than another hypothesis in  $S$ .
    - If  $d$  is Negative example
      - Remove from  $S$  any hypothesis inconsistent with  $d$ .
      - For each hypothesis  $g$  in  $G$  i.e. not consistent with  $d$ 
        - Remove  $g$  from  $G$ .
        - Add to  $G$  all minimal generalization of  $g$  such that
          - $h$  is consistent with  $d$
          - some member of  $S$  is more specific than  $h$ .
        - Remove from  $G$  any hypothesis i.e. less general than another hypothesis in  $G$ .

Examp:-

(24)

Initialise G & S

$G_0 \leftarrow \{<?, ?, ?, ?, ?, ?, ?>\}$  most general hypothesis  
 $S_0 \leftarrow \{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset\}$  " specific " "

Candidate Elimination Trace :-

Training Examples :

1.  $\langle \text{sunny}, \text{warm}, \text{Normal}, \text{strong}, \text{warm}, \text{same}, \text{high}, \text{Enjoy sport} = \text{yes} \rangle$

2.  $\langle \text{sunny}, \text{warm}, \cancel{\text{Normal}}, \text{strong}, \text{warm}, \text{same}, \text{Enjoy sport} = \text{yes} \rangle$

-  $S_0$  &  $G_0$  are initial boundary sets.

- Training Examples 1 & 2 force  $S$  boundary to become more general, as in "Find-S" Algo.

- They don't have any effect on  $G$  boundary.

$S_0: \boxed{\{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset\}}$

$S_1: \boxed{\{\text{sunny}, \text{warm}, \text{Normal}, \text{strong}, \text{warm}, \text{same}\}}$

$S_2: \boxed{\{\text{sunny}, \text{warm}, ?, \text{strong}, \text{warm}, \text{same}\}}$

$G_0, G_1, G_2: \boxed{\{<?, ?, ?, ?, ?, ?, ?>\}}$

## → Candidate Elimination Trace 2 :-

(25)

Training Example:

3.  $\langle \text{Rainy}, \text{cold}, \text{High}, \text{strong}, \text{Warm}, \text{change} \rangle$   
    Enjoy sport

∴ it's a -ve example, it will now force  $G_2$  boundary to be specialised to  $G_3$ . & there can be several alternative maximally general hypotheses that are included.

$S_2, S_3$ :  $\{\langle \text{sunny}, \text{warm}, ?, \text{strong}, \text{Warm}, \text{same} \rangle\}$

$G_3$ :  $\{\langle \text{sunny}, ?, ?, ?, ?, ?, ? \rangle, \langle ?, \text{warm}, ?, ?, ?, ?, ? \rangle, \langle ?, ?, ?, ?, ?, ?, ? \rangle\}$

$G_2$ :  $\{\langle ?, ?, ?, ?, ?, ?, ? \rangle\}$

## → Candidate Elimination Trace 3:-

Training Example:

4.  $\langle \text{sunny}, \text{warm}, \text{high}, \text{strong}, \text{cool}, \text{change} \rangle$ ,  
    Enjoy sport Yes.

∴ the Training example, it generalizes 's' boundary from  $S_3$  to  $S_4$ .

Also, one member of  $G_3$  should be deleted since it is no longer more general than  $S_4$  boundary.

S3:  $\{ \langle \text{sunny}, \text{warm}, ?, \text{strong}, \text{warm}, \text{same} \rangle \}$

S4:  $\{ \langle \text{sunny}, \text{warm}, ?, \text{strong}, ?, ?, ? \rangle \}$

G4:  $\{ \langle \text{sunny}, ?, ?, ?, ?, ? \rangle \} \subset \{ \langle ?, \text{warm}, ?, ?, ?, ? \rangle \}$

G3:  $\{ \langle \text{sunny}, ?, ?, ?, ?, ? \rangle \} \subset \{ \langle ?, \text{warm}, ?, ?, ?, ? \rangle \} \subset \{ \langle ?, ?, ?, ?, ?, ? \rangle \}$

Final version :-

- Set S4 & G4 delimit version space of all hypotheses consistent with set of incrementally observed training examples.
- The foll. learned Version space is independent of the seq. in which training examples present.
- S & G Boundaries move monotonically closer, delimiting smaller & smaller version space of candidate hypotheses.

S4:  $\{ \langle \text{sunny}, \text{warm}, ?, \text{strong}, ?, ?, ? \rangle \}$

$\langle \text{sunny}, ?, ?, \text{strong}, ?, ?, ? \rangle \subset \langle \text{sunny}, \text{warm}, ?, ?, ?, ? \rangle \subset \langle ?, \text{warm}, ?, \text{strong}, ?, ? \rangle$

G4:  $\{ \langle \text{sunny}, ?, ?, ?, ?, ? \rangle \} \subset \{ \langle ?, \text{warm}, ?, ?, ?, ? \rangle \}$

## Remarks On Version space & "Candidate-Elimination".

- "Version space" learned by the Candidate-Elimination algorithm will converge towards the hypothesis that correctly describes the target concept, provided
  - (i) There are no errors in the training examples.
  - (ii) There is some hypothesis in  $H$  that correctly describes the target concept.
- Also if training data contains errors, Algo. will remove correct target concept from Version space & may create an empty version space.
- If the learner is allowed to conduct experiments in which it chooses next instance (a.k.a. Query) then it obtains correct classification for this instance from an external oracle (ex: Nature or Teacher).  
Example:-  
For "enjoy sport" task, the "Query strategy" would be to choose an instance that would be classified negative by some of these hypotheses, but positive by others such as:  
 $\langle \text{sunny}, \text{warm}, \text{Normal}, \text{light}, \text{warm}, \text{same} \rangle$ .
- With optimal query strategy, size of version space is reduced by half with each new example. The correct target concept can therefore be found with only  $\lceil \log_2 |V_S| \rceil$  experiments.

- (26)
- It is possible to classify certain examples, with same degree of confidence as if the target concept had been uniquely identified, even if it were partially learned concept.
  - Consider an instance A, which was not amongst the training examples, the learner will classify it as +ve since, condition will be met iff instance satisfies every member of S. (with definition of more general than).
  - <sup>ifly</sup> For instance B, it will be classified Negative by every hypothesis in version space; gives a partially learned concept. (Just check if instance satisfies none of members of G).
  - For instance C, Half of version space hypotheses classify it +ve &  $\frac{1}{2}$  as -ve. Hence Learner can't classify this instance confidently.
  - In case of instance D, it is classified +ve by 2 of Version space hypotheses & -ve by other 4  $\Rightarrow$  Negative classification it is.

<u>Instance</u>	<u>Sky</u>	<u>AirTemp</u>	<u>Humidity</u>	<u>Wind</u>	<u>Water</u>	<u>Forecast</u>	<u>Enjoy Sport</u>
A	sunny	warm	normal	strong	cool	same	?
B	rainy	cold	"	light	warm	"	?
C	sunny	warm	"	"	"	"	-
D	sunny	cold	"	strong	"	"	?

Q5. ILLUSTRATE ABOUT BASIC DECISION TREE ALGORITHM WITH A SUITABLE EXAMPLE.

- THE BASIC DECISION TREE LEARNING ALGORITHM:
- Most Algo's for learning decision trees are built by varying the core Algo.
  - The core Algo. uses Top-down, greedy search through the space of possible decision trees.
  - Few ex. are: ID3 Algo., C4.5 Algo.
  - The Basic Algo. is ID3 Algo. in which it starts by questioning which Attribute to be tested at Root?
  - In order to answer it, first, all the instance attributes are evaluated using a statistical test which determines how well it classifies the training examples.
  - Later, the best attribute is selected & used as test at "Root" node of the tree.
  - Further, a descendant to the root node will be created & training examples are sorted w.r.t. it. Entire process is repeated using training examples assoc. with each descendant node to select the best attribute to test at that particular point of a tree.
  - Greedy Search:- The above process forms a greedy search since, algo. never back tracks to reconsider earlier choices made.
  - Summary of ID3 Algo: (Training Examples, Target attribute, Attributes)  
ID3 (Examples, Target attribute, Attributes)
    1. Create Root Node for Tree
    2. If all Examples are +ve, Return single node tree Root, with label = +
    3. " " " are -ve, single node tree Root, with label = -
    4. If Attribute = empty  $\Rightarrow$  Return single node tree Root  
label = most Common value of Target attribute

else Begin  
(i)  $A \leftarrow$  attribute from "Attribute" that best classifies "Example"  
(ii) Decision Attr rule for Root  $\leftarrow A$   
(iii) for all,  $v_i$  of  $A$   
    + Add new tree branch below Root corresponding to test  $A = v_i$   
    + let  $Example_{v_i}$  be subset of Example that have the value  $v_i$  for  $A$ .

- \* If  $\text{Example}_{i,j}$  is empty  
 $\Rightarrow$  Below this new branch add a leaf node  
 with label = most common value of  
 "Target\_attribute" in "Example"
- Else  
 $\Rightarrow$  Below this new branch add the SubTree  
 $\text{ID3}(\text{Example}_{i,j}, \text{Target\_attribute}, \text{Attribute} - \{A\})$

End

Return Root.

- $\rightarrow$  which Attribute is the Best classifier?
- An Attribute that is most useful for classifying examples
  - An Attribute that has "highest" information gain
  - Information gain: It is a measure of how well a given attribute separates the training examples accdg to their Target classification.

$\rightarrow$  Entropy:

- Measures Homogeneity of Examples.
- Entropy is needed to define Info.gain.
- It characterizes purity of an arbitrary collection of examples.
- consider 'S' = collection of +ve & -ve Examples of a target concept.

$$\text{Entropy}(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

Where,  $P_+$  = proportion of +ve examples in S  
 $\& P_- = \frac{\text{No. of -ve}}{\text{Total No. of examples}}$

(Note: -  $\log_0 = 0$ )

Ex: - If S has 14 ex. & 9 are +ve & 5 are -ve, then,  $\text{Entropy}(S) = -(9/14) \log_2(9/14) - (5/14) \log_2(5/14) = 0.940$ .

- If the members of  $S'$  belong to same class then Entropy = 0.  
 [soln:- If 'S' has 14 exs,  $P_{(+)} = 14/14 = 1$  &  $P_{(-)} = 0$   
 $\rightarrow H(14, 0) = -\frac{14}{14} \log_2 \frac{14}{14} - \frac{0}{14} \log_2 0$   
 $= -1 \log_2 1 - 0 = 0$  .]
- If there are equal no. of +ve & -ve ex then Entropy = 1.
- If there are Unequal no. of +ve & -ve ex then Entropy lies b/w 0 & 1.

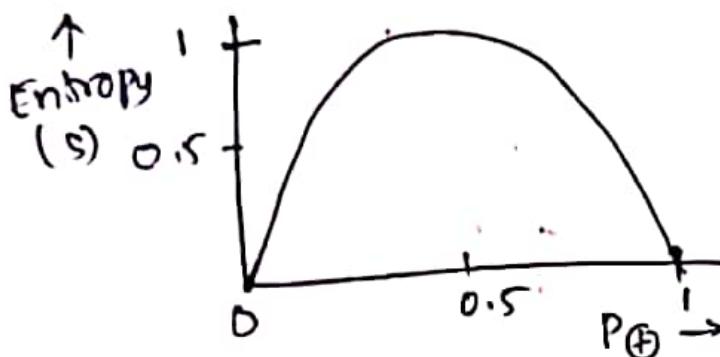


Fig:- Entropy for a classification for a boolean classification as the proportion,  $P_{(+)}$ , of +ve examples varies b/w 0 & 1.

- "Entropy", accdg to info. theory, specifies min. no. of bits of info. needed in order to encode the classification of an arbitrary member of  $S$ .
  - \* If  $P_{(+)} = 1 \Rightarrow$  Receiver knows that drawn ex. will be +ve  
~~ex~~  $\Rightarrow$  no msg need to be sent  
 $\Rightarrow$  Entropy = 0.
  - \* If  $P_{(+)} = 0.5 \Rightarrow$  1 bit is needed to indicate if drawn ex. is +ve or -ve.
  - \* If  $P_{(+)} = 0.8 \Rightarrow$  A collection of msgs. can be encoded using almost  $< 1$  bit per msg by assigning shorter codes for +ve exs. & longer codes for less likely, -ve exs.
- Generally, (in classification which is not boolean), the target attribute can take on 'c' diff. values, then the entropy of S relative to this c-wise classification is defined as

$$\text{Entropy}(S) = \sum_{i=1}^c -P_i \log_2 P_i$$

where,  
 $P_i$  = proportion of S belonging to class i

i.e.  $\log_2$  is used as encoded length in bits. If target attribute can take c values, then entropy =  $\log_2 c$ .

→ Info. gain Measures the expected reduction in Entropy: (31)

- WKT Entropy is a measure of purity/impurity of the collection of Training exs. (examples)
- Info. gain is a measure of effectiveness of an attribute in classifying the Training data.
- Info. gain, measures the expected reduction in entropy which is caused due to the partitioning of exs. accdg to this attribute.
- $\text{Gain}(S, A) = \text{Info. gain of an attribute } A \text{ relative to a collection of examples } S$ .

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

where  $\text{Values}(A)$  = set of all possible values for the attribute  $A$ .

$S_v$  = subset of 'S' for which attribute 'A' has the value  $v$   
 $= \{s \in S \mid A(s) = v\}$ .

The  $\text{Entropy}(S)$  = Entropy of original collection 'S'.

The second term in above eqn. is Expected Entropy.  
 $\text{Expected Entropy} = \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$ .

It is the sum of entropies of each subset  $S_v$ , weighted by the fraction of examples  $\frac{|S_v|}{|S|}$  that belong to  $S_v$ .

$\text{Gain}(S, A)$  can also be defined as info. provided about "target func. value" given value of some other attribute.

A - It is also the no. of bits saved when encoding target value of an arbitrary member of  $S$ , by knowing the value of attribute 'A'.

### Example:-

- let  $S$  = collection of training example described by attributes (including like wind that can take 2 values, weak or strong).
- Assuming, as previous,  $S[9+, 5-]$  & with that now add that, 6 tve  $\overset{k=2}{\text{ex-ve}}$  examples that have  $\text{wind} = \text{weak}$  & remaining all have  $\text{wind} = \text{strong}$   
 $\rightarrow \text{values}(\text{wind}) = \text{Weak, strong}$ .
- $S = [9+, 5-]$   
 $S_{\text{Weak}} \leftarrow [6+, 2-]$   
 $S_{\text{Strong}} \leftarrow [3+, 3-]$

For such situation Infogain, will be,

$$\begin{aligned}\text{Gain}(S, \text{wind}) &= \text{Entropy}(S) - \sum_{V \in \{\text{weak, strong}\}} \frac{|S_V|}{|S|} \text{Entropy}(S_V) \\ &= \text{Entropy}(S) - \frac{8}{14} \text{Entropy}(S_{\text{weak}}) - \frac{6}{14} \text{Entropy}(S_{\text{strong}}) \\ &= 0.940 - \frac{8}{14} \left[ -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right] \\ &\quad - \frac{6}{14} \left[ \frac{3}{3} \right] \left[ \because \text{Entropy}(S_{\text{strong}}) \text{ coz } 3+3 \right] \\ &= 0.940 - \frac{8}{14} (0.811) - \frac{6}{14} (1) \\ &= 0.048.\end{aligned}$$

||ly we calculate  $\text{Gain}(S, \text{Humidity})$  considering  $[3+, 4-]$  ex-  
with Humidity = High &  $[6+, 1-]$  as Normal  $\Rightarrow \text{Gain}(S, \text{Humidity}) = 0.151$

$$\begin{array}{c} S: [9+, 5-] \\ E = 0.940 \end{array}$$

<b>Humidity</b>	
High	Normal
$[3+, 4-]$	$[6+, 1-]$
$E = 0.985$	$E = 0.592$

 $\text{Gain}(S, \text{Humidity}) = 0.15$

$$\begin{array}{c} S: [9+, 5-] \\ E = 0.940 \end{array}$$

<b>wind</b>	
weak	strong
$[6+, 2-]$	$[3+, 3-]$
$E = 0.811$	$E = 1.0$

 $\text{Gain}(S, \text{wind}) = 0.048.$

$$\begin{array}{l} \text{Gain}(S, \text{outlook}) = 0.246 \\ \text{Gain}(S, \text{wind}) = 0.048 \\ \text{Gain}(S, \text{Temperature}) = 0.029 \end{array}$$

(4)

→ To understand working of ID3, let's consider following training exs.

Day	outlook	Temperature	Humidity	wind	Play Tennis
D1	sunny	Hot	High	weak	NO
D2	sunny	Hot	High	strong	NO
D3	overcast	Hot	High	weak	YES
D4	Rain	mild	High	weak	YES
D5	Rain	cool	Normal	weak	NO
D6	Rain	cool	Normal	strong	NO
D7	overcast	cool	Normal	strong	YES
D8	sunny	mild	High	weak	NO
D9	sunny	cool	Normal	weak	YES
D10	Rain	mild	Normal	strong	YES
D11	sunny	mild	Normal	strong	YES
D12	overcast	mild	High	strong	YES
D13	overcast	Hot	Normal	weak	YES
D14	Rain	mild	High	strong	NO

→ Topmost node of decision Tree is created -

→ which attribute should be tested 1st .

→ ID3 determines Info. gain for each candidate attribute i.e outlook, Temp., Humidity & Wind & selects one with highest Info. gain i.e Here, "Outlook".

→ Thus "outlook" is selected as decision attribute for root node and branches are created for each possible value

Here 3 vif. Sunny, Rain, & overcast .

$\{D1, D2, \dots, D14\}$

$[9+, 5-]$

outlook

Sunny

overcast

Rain

$\{D1, D2, D8, D11\}$

$[2+, 3-]$

$\{D3, D7, D12, D9\}$

$[4+, 0-]$

$\{D4, D5, D6, D10, D14\}$

$[3+, 2-]$

YES

?

- For the Outlook = overcast, "playTennis" is always a tve example. (42)  
Therefore it became leaf node ..
- But the process must be further repeated for new attributes at the descendant nodes of outlook = sunny & rain.
- Attributes that are included in higher part of tree are excluded. process continues for each leaf node until
  - ① every attribute has already been included through the tree (or)
  - ② Training exs. assoc. with this leaf node have all same target attribute value (Entropy = 0)