

Internal Assessment Test –I, September 2020

Sub:	MACHINE LEARNING						Code:	18MCA53			
Date:	15-09-2020	Duration:	90 mins	Max Marks:	50	Sem:	V	Branch:	MCA		
Answer Any 5 QUESTIONS								Marks	OBE		
									CO	RBT	
1a)	Define well-posed learning problem and explain with the help of an example.						4	CO1	L1		
1b)	Discuss the following with respect to the learning task for “A checkers learning problem”. i) Choosing the training experience. ii) Choosing the target function. iii) Choosing a function approximation algorithm.						6	CO1	L2		
2)	Write FIND-S algorithm and discuss the issues with the algorithm						10	CO1	L2		
3)	Consider the given below following training example.						10	CO1	L5		
	Example	Sky	AirTemp	Humidity	Wind	Water				Forecast	EnjoySport
	1	Sunny	Warm	Normal	Strong	Warm				Same	Yes
	2	Sunny	Warm	High	Strong	Warm				Same	Yes
	3	Rainy	Cold	High	Strong	Warm				Change	No
	4	Sunny	Warm	High	Strong	Cool	Change	Yes			
	Show the general and specific boundaries of the version space after applying candidate elimination algorithm.										
4	Explain the various stages involved in designing a learning system in brief.						10	CO1	L2		
5	What is decision tree and discuss the use of decision tree for classification with an example.						10	CO2	L2		
6)	Summarize the practical issues of decision tree learning.						10	CO2	L3		
7)	Design the decision tree for the following dataset and predict whether Golf will be played on the day.						10	CO2	L5		
	Day	Outlook	Temperature	Humidity	Wind	Play Golf					
	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	Yes					
	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	Weak	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					
8)	What do you mean by gain and entropy? How it is used to build decision tree. Explain it with the help of an example.						10	CO2	L2		

--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

Internal Assessment Test 1 – Sep. 2020

Sub:	MACHINE LEARNING							Sub Code:	18MCA 53
Date:	15-09-2020	Duration:	90 min's	Max Marks:	50	Sem	5 th	Branch:	MCA

Note : Answer FIVE FULL Questions, choosing ONE full question from each Module

PART I		MAR KS	OBE	
			C O	RB T
1a)	<p>Well-Posed Learning Definition: A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.</p> <p>Examples:</p> <p>Checkers Game: A computer program that learns to play checkers might improve its performance as measured by its ability to win at the class of tasks involving playing checkers game, through experience obtained by playing games against itself:</p> <p>checkers learning problem:</p> <ul style="list-style-type: none"> <input type="checkbox"/> Task T: playing checkers <input type="checkbox"/> Performance measure P: percent of games won against opponents <input type="checkbox"/> Training experience E: playing practice games against itself <p>A handwriting recognition learning problem:</p> <ul style="list-style-type: none"> <input type="checkbox"/> Task T: recognizing and classifying handwritten words within images <input type="checkbox"/> Performance measure P: percent of words correctly classified <input type="checkbox"/> Training experience E: a database of handwritten words with given classifications <p>A robot driving learning problem:</p> <ul style="list-style-type: none"> <input type="checkbox"/> Task T: driving on public four-lane highways using vision sensors <input type="checkbox"/> Performance measure P: average distance travelled before an error (as judged by human overseer) <input type="checkbox"/> Training experience E: a sequence of images and steering commands recorded while observing a human driver 	4	CO1	L1
1b)	<p>DESIGNING A LEARNING SYSTEM</p> <p>The basic design issues and approaches to machine learning are illustrated by designing a program to learn to play checkers, with the goal of entering it in the world checkers tournament</p> <ol style="list-style-type: none"> 1. Choosing the Training Experience 2. Choosing the Target Function 3. Choosing a Function Approximation Algorithm <ol style="list-style-type: none"> 1. Estimating training values 2. Adjusting the weights <p>1. Choosing the Training Experience</p> <ul style="list-style-type: none"> <input type="checkbox"/> The first design choice is to choose the type of training experience from which the 	6	CO1	L2

system will learn.

- The type of training experience available can have a significant impact on success or failure of the learner.

There are three attributes which impact on success or failure of the learner

1. Whether the training experience provides *direct or indirect feedback* regarding the choices made by the performance system.

For example, in checkers game:

In learning to play checkers, the system might learn from *direct training examples* consisting of *individual checkers board states* and *the correct move for each*.

Indirect training examples consisting of the *move sequences* and *final outcomes* of various games played. The information about the correctness of specific moves early in the game must be inferred indirectly from the fact that the game was eventually won or lost.

Here the learner faces an additional problem of *credit assignment*, or determining the degree to which each move in the sequence deserves credit or blame for the final outcome.

2. The degree to which the *learner controls the sequence of training examples*

For example, in checkers game:

The learner might depend on the *teacher* to select informative board states and to provide the correct move for each.

Alternatively, the learner might itself propose board states that it finds particularly confusing and ask the teacher for the correct move.

The learner may have complete control over both the board states and (indirect) training classifications, as it does when it learns by playing against itself with *no teacher present*.

3. How well it represents the *distribution of examples* over which the final system performance P must be measured

For example, in checkers game:

In checkers learning scenario, the performance metric P is the percent of games the system wins in the world tournament.

If its training experience E consists only of games played against itself, there is a danger that this training experience might not be fully representative of the distribution of situations over which it will later be tested.

It is necessary to learn from a distribution of examples that is different from those on which the final system will be evaluated.

2. *Choosing the Target Function*

The next design choice is to determine exactly what type of knowledge will be learned and how this will be used by the performance program.

Let's consider a checkers-playing program that can generate the legal moves from any board state. The program needs only to learn how to choose the best move from among these legal moves. We must learn to choose among the legal moves, the most obvious choice for the type of information to be learned is a program, or function, that chooses the best move for any given board state.

1. Let *ChooseMove* be the target function and the notation is

$$\text{ChooseMove} : B \rightarrow M$$

which indicate that this function accepts as input any board from the set of legal board states B and produces as output some move from the set of legal moves M.

ChooseMove is a choice for the target function in checkers example, but this function will turn out to be very difficult to learn given the kind of indirect training experience available to our system

2. An alternative target function is an *evaluation function* that assigns a *numerical score* to any given board state

Let the target function V and the notation

$$V : B \rightarrow R$$

which denote that V maps any legal board state from the set B to some real value. Intend for this target function V to assign higher scores to better board states. If the system can successfully learn such a target function V, then it can easily use it to select the best move from any current board position.

Let us define the target value V(b) for an arbitrary board state b in B, as follows:

- If b is a final board state that is won, then $V(b) = 100$
- If b is a final board state that is lost, then $V(b) = -100$
- If b is a final board state that is drawn, then $V(b) = 0$
- If b is a not a final state in the game, then $V(b) = V(b')$,

Where b' is the best final board state that can be achieved starting from b and playing optimally until the end of the game

3. Choosing a Function Approximation Algorithm

In order to learn the target function f we require a set of training examples, each describing a specific board state b and the training value $V_{\text{train}}(b)$ for b.

Each training example is an ordered pair of the form (b, $V_{\text{train}}(b)$).

For instance, the following training example describes a board state b in which black has won the game (note $x_2 = 0$ indicates that red has no remaining pieces) and for which the target function value $V_{\text{train}}(b)$ is therefore +100.

$$((x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0), +100)$$

Function Approximation Procedure

1. Derive training examples from the indirect training experience available to the learner
2. Adjusts the weights w_i to best fit these training examples

1. Estimating training values

A simple approach for estimating training values for intermediate board states is to assign the training value of $V_{\text{train}}(b)$ for any intermediate board state b to be

$V(\text{Successor}(b))$

Where ,

-

V is the learner's current approximation to V

- Successor(b) denotes the next board state following b for which it is again the

	<p style="text-align: center;">program's turn to move</p> <p>Rule for estimating training values</p> <p style="text-align: center;">$V_{train}(b) \leftarrow V(\text{Successor}(b))$</p> <p>2. <u>Adjusting the weights</u></p> <p>Specify the learning algorithm for choosing the weights w_i to best fit the set of training examples $\{(b, V_{train}(b))\}$</p> <p>A first step is to define what we mean by the bestfit to the training data.</p> <p>One common approach is to define the best hypothesis, or set of weights, as that which minimizes the squared error E between the training values and the values predicted by the hypothesis.</p> <p>Several algorithms are known for finding weights of a linear function that minimize E. One such algorithm is called the <i>least mean squares, or LMS training rule</i>. For each observed training example it adjusts the weights a small amount in the direction that reduces the error on this training example</p> <p>LMS weight update rule :- For each training example $(b, \hat{V}_{train}(b))$</p> <p style="padding-left: 40px;">Use the current weights to calculate $V(b)$</p> <p style="padding-left: 40px;">For each weight w_i, update it as</p> $w_i \leftarrow w_i + \eta (V_{train}(b) - \hat{V}(b)) x_i$ <p>Here η is a small constant (e.g., 0.1) that moderates the size of the weight update.</p> <p><u>Working of weight update rule</u></p> <ul style="list-style-type: none"> <input type="checkbox"/> When the error $(V_{train}(b) - \hat{V}(b))$ is zero, no weights are changed. <input type="checkbox"/> When $(V_{train}(b) - \hat{V}(b))$ is positive (i.e., when $\hat{V}(b)$ is too low), then each weight is increased in proportion to the value of its corresponding feature. This will raise the value of $\hat{V}(b)$, reducing the error. <input type="checkbox"/> If the value of some feature x_i is zero, then its weight is not altered regardless of the error, so that the only weights updated are those whose features actually occur on the training example board. 			
2)	<p>FIND-S: FINDING A MAXIMALLY SPECIFIC HYPOTHESIS</p> <p><u>FIND-S Algorithm</u></p> <ol style="list-style-type: none"> 1. Initialize h to the most specific hypothesis in H 2. For each positive training instance x <ul style="list-style-type: none"> For each attribute constraint a_i in h <ul style="list-style-type: none"> If the constraint a_i is satisfied by x <ul style="list-style-type: none"> Then do nothing Else replace a_i in h by the next more general constraint that is satisfied by x 3. Output hypothesis h <p><u>Unanswered by FIND-S</u></p> <ol style="list-style-type: none"> 1. Has the learner converged to the correct target concept? 2. Why prefer the most specific hypothesis? 3. Are the training examples consistent? 4. What if there are several maximally specific consistent hypotheses? 	10	CO1	L2

3) CANDIDATE-ELIMINATION algorithm begins by initializing the version space to the set of all hypotheses in H;

Initializing the G boundary set to contain the most general hypothesis in
 $H \quad G_0 \langle \text{?, ?, ?, ?, ?} \rangle$

Initializing the S boundary set to contain the most specific (least general) hypothesis
 $S_0 \langle \text{?, ?, ?, ?, ?} \rangle$

For training example d,
 $\langle \text{Sunny, Warm, Normal, Strong, Warm, Same} \rangle +$

$S_0 \quad \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$

$S_1 \quad \langle \text{Sunny, Warm, Normal, Strong, Warm, Same} \rangle$

$G_0, G_1 \quad \langle \text{?, ?, ?, ?, ?, ?} \rangle$

□ When the second training example is observed, it has a similar effect of generalizing S further to S2, leaving G again unchanged i.e., $G_2 = G_1 = G_0$

For training example d,
 $\langle \text{Sunny, Warm, High, Strong, Warm, Same} \rangle +$

$S_1 \quad \langle \text{Sunny, Warm, Normal, Strong, Warm, Same} \rangle$

$S_2 \quad \langle \text{Sunny, Warm, ?, Strong, Warm, Same} \rangle$

$G_1, G_2 \quad \langle \text{?, ?, ?, ?, ?, ?} \rangle$

Consider the third training example.

For training example d,
 $\langle \text{Rainy, Cold, High, Strong, Warm, Change} \rangle -$

$S_2, S_3 \quad \langle \text{Sunny, Warm, ?, Strong, Warm, Same} \rangle$

$G_3 \quad \langle \text{Sunny, ?, ?, ?, ?, ?} \rangle \langle \text{?, Warm, ?, ?, ?, ?} \rangle \langle \text{?, ?, ?, ?, ?, Same} \rangle$

$G_2 \quad \langle \text{?, ?, ?, ?, ?, ?} \rangle$

Consider the fourth training example

For training example d,

$\langle \text{Sunny, Warm, High, Strong, Cool Change} \rangle +$

S_3

$\langle \text{Sunny, Warm, ?, Strong, Warm, Same} \rangle$

S_4

$\langle \text{Sunny, Warm, ?, Strong, ?, ?} \rangle$

G_4

$\langle \text{Sunny, ?, ?, ?, ?, ?} \rangle \langle \text{?, Warm, ?, ?, ?, ?} \rangle$

G_3

$\langle \text{Sunny, ?, ?, ?, ?, ?} \rangle \langle \text{?, Warm, ?, ?, ?, ?} \rangle \langle \text{?, ?, ?, ?, ?, Sam} \rangle$

After processing these four examples, the boundary sets S_4 and G_4 delimit the version space of all hypotheses consistent with the set of incrementally observed training examples.

S_4

$\langle \text{Sunny, Warm, ?, Strong, ?, ?} \rangle$

$\langle \text{Sunny, ?, ?, strong, ?, ?} \rangle \langle \text{Sunny, Warm, ?, ?, ?, ?} \rangle \langle \text{?, Warm, ?, Strong, ?} \rangle$

G_4

$\langle \text{Sunny, ?, ?, ?, ?, ?} \rangle \langle \text{?, Warm, ?, ?, ?, ?} \rangle$

4) DESIGNING A LEARNING SYSTEM

The basic design issues and approaches to machine learning are illustrated by designing a program to learn to play checkers, with the goal of entering it in the world checkers tournament

4. Choosing the Training Experience
5. Choosing the Target Function
6. Choosing a Representation for the Target Function
7. Choosing a Function Approximation Algorithm
 1. Estimating training values
 2. Adjusting the weights
8. The Final Design

2. Choosing the Training Experience

- The first design choice is to choose the type of training experience from which the system will learn.
- The type of training experience available can have a significant impact on success or failure of the learner.

There are three attributes which impact on success or failure of the learner

2. Whether the training experience provides *direct or indirect feedback* regarding the choices made by the performance system.

10

CO1

L2

For example, in checkers game:

In learning to play checkers, the system might learn from *direct training examples* consisting of *individual checkers board states* and *the correct move for each*.

Indirect training examples consisting of the *move sequences* and *final outcomes* of various games played. The information about the correctness of specific moves early in the game must be inferred indirectly from the fact that the game was eventually won or lost.

Here the learner faces an additional problem of *credit assignment*, or determining the degree to which each move in the sequence deserves credit or blame for the final outcome.

4. The degree to which the *learner controls the sequence of training examples*

For example, in checkers game:

The learner might depend on the *teacher* to select informative board states and to provide the correct move for each.

Alternatively, the learner might itself propose board states that it finds particularly confusing and ask the teacher for the correct move.

The learner may have complete control over both the board states and (indirect) training classifications, as it does when it learns by playing against itself with *no teacher present*.

5. How well it represents the *distribution of examples* over which the final system performance P must be measured

For example, in checkers game:

In checkers learning scenario, the performance metric P is the percent of games the system wins in the world tournament.

If its training experience E consists only of games played against itself, there is a danger that this training experience might not be fully representative of the distribution of situations over which it will later be tested.

It is necessary to learn from a distribution of examples that is different from those on which the final system will be evaluated.

3. *Choosing the Target Function*

The next design choice is to determine exactly what type of knowledge will be learned and how this will be used by the performance program.

Let's consider a checkers-playing program that can generate the legal moves from any board state. The program needs only to learn how to choose the best move from among these legal moves. We must learn to choose among the legal moves, the most obvious choice for the type of information to be learned is a program, or function, that chooses the best move for any given board state.

1. Let *ChooseMove* be the target function and the notation is

***ChooseMove* : $B \rightarrow M$**

which indicate that this function accepts as input any board from the set of legal board states B and produces as output some move from the set of legal moves M.

ChooseMove is a choice for the target function in checkers example, but this function will turn out to be very difficult to learn given the kind of indirect training experience available to our system

3. An alternative target function is an *evaluation function* that assigns a *numerical score* to any given board state

Let the target function V and the notation

$$V: B \rightarrow R$$

which denote that V maps any legal board state from the set B to some real value. Intend for this target function V to assign higher scores to better board states. If the system can successfully learn such a target function V , then it can easily use it to select the best move from any current board position.

Let us define the target value $V(b)$ for an arbitrary board state b in B , as follows:

- If b is a final board state that is won, then $V(b) = 100$
- If b is a final board state that is lost, then $V(b) = -100$
- If b is a final board state that is drawn, then $V(b) = 0$
- If b is a not a final state in the game, then $V(b) = V(b')$,

Where b' is the best final board state that can be achieved starting from b and playing optimally until the end of the game

3. Choosing a Representation for the Target Function

Let's choose a simple representation - for any given board state, the function c will be calculated as a linear combination of the following board features:

- x_1 : the number of black pieces on the board
- x_2 : the number of red pieces on the board
- x_3 : the number of black kings on the board
- x_4 : the number of red kings on the board
- x_5 : the number of black pieces threatened by red (i.e., which can be captured on red's next turn)
- x_6 : the number of red pieces threatened by black

Thus, learning program will represent as a linear function of the form

$$\hat{V}(b) = w_0 + w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_5x_5 + w_6x_6$$

Where,

- w_0 through w_6 are numerical coefficients, or weights, to be chosen by the learning algorithm.
- Learned values for the weights w_1 through w_6 will determine the relative importance of the various board features in determining the value of the board
- The weight w_0 will provide an additive constant to the board value

4. Choosing a Function Approximation Algorithm

In order to learn the target function f we require a set of training examples, each describing a specific board state b and the training value $V_{train}(b)$ for b .

Each training example is an ordered pair of the form $(b, V_{train}(b))$.

For instance, the following training example describes a board state b in which black has won the game (note $x_2 = 0$ indicates that red has no remaining pieces) and for which the target function

value $V_{train}(b)$ is therefore +100.

$$((x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0), +100)$$

Function Approximation Procedure

3. Derive training examples from the indirect training experience available to the learner
4. Adjusts the weights w_i to best fit these training examples

2. Estimating training values

A simple approach for estimating training values for intermediate board states is to assign the training value of $V_{train}(b)$ for any intermediate board state b to be

$V(\text{Successor}(b))$

Where ,

□

V is the learner's current approximation to V

□ $\text{Successor}(b)$ denotes the next board state following b for which it is again the program's turn to move

Rule for estimating training values

$$V_{train}(b) \leftarrow V(\text{Successor}(b))$$

3. Adjusting the weights

Specify the learning algorithm for choosing the weights w_i to best fit the set of training examples $\{(b, V_{train}(b))\}$

A first step is to define what we mean by the bestfit to the training data.

One common approach is to define the best hypothesis, or set of weights, as that which minimizes the squared error E between the training values and the values predicted by the hypothesis.

Several algorithms are known for finding weights of a linear function that minimize E . One such algorithm is called the *least mean squares, or LMS training rule*. For each observed training example it adjusts the weights a small amount in the direction that reduces the error on this training example

LMS weight update rule :- For each training example $(b, \hat{V}_{train}(b))$

Use the current weights to calculate $V(b)$

For each weight w_i , update it as

$$w_i \leftarrow w_i + \eta (V_{train}(b) - \hat{V}(b)) x_i$$

Here η is a small constant (e.g., 0.1) that moderates the size of the weight update.

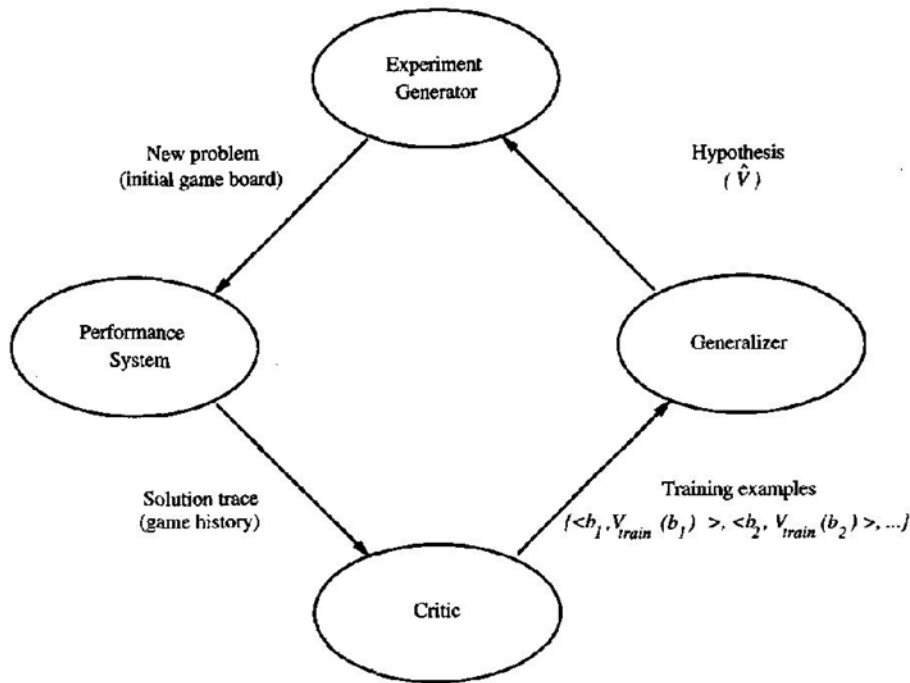
Working of weight update rule

- When the error $(V_{train}(b) - \hat{V}(b))$ is zero, no weights are changed.
- When $(V_{train}(b) - \hat{V}(b))$ is positive (i.e., when $\hat{V}(b)$ is too low), then each weight is increased in proportion to the value of its corresponding feature. This will raise the value of $\hat{V}(b)$, reducing the error.
- If the value of some feature x_i is zero, then its weight is not altered regardless of the error, so that the only weights updated are those whose features x_i actually occur on the training example board.

The Final Design

The final design of checkers learning system can be described by four distinct program modules that represent the central components in many learning systems

1. **The Performance System** is the module that must solve the given performance task by using the learned target function(s). It takes an instance of a new problem (new game) as input and produces a trace of its solution (game history) as output.
2. **The Critic** takes as input the history or trace of the game and produces as output a set of training examples of the target function
3. **The Generalizer** takes as input the training examples and produces an output hypothesis that is its estimate of the target function. It generalizes from the specific training examples, hypothesizing a general function that covers these examples and other cases beyond the training examples.
4. **The Experiment Generator** takes as input the current hypothesis and outputs a new problem (i.e., initial board state) for the Performance System to explore. Its role is to pick new practice problems that will maximize the learning rate of the overall system.



5)

Decision tree learning is a method for approximating discrete-valued target functions, in which the learned function is represented by a decision tree.

DECISION TREE REPRESENTATION

- Decision trees classify instances by sorting them down the tree from the root to some leaf node, which provides the classification of the instance.
- Each node in the tree specifies a test of some attribute of the instance, and each branch descending from that node corresponds to one of the possible values for this attribute.
- An instance is classified by starting at the root node of the tree, testing the attribute specified by this node, then moving down the tree branch corresponding to the value of the attribute in the given example. This process is then repeated for the subtree rooted at the new node.

10

CO2

L2

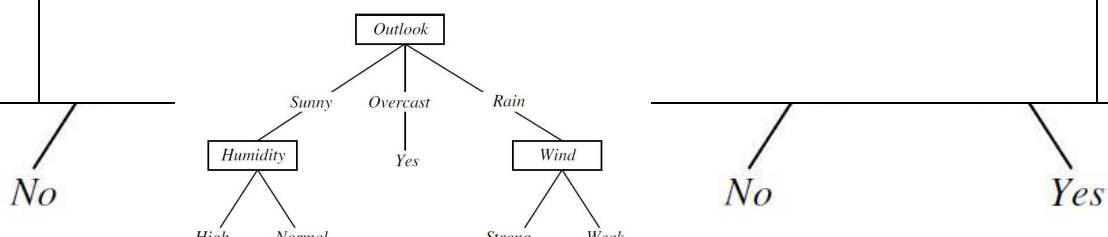


FIGURE: A decision tree for the concept *PlayTennis*. An example is classified by sorting it through the tree to the appropriate leaf node, then returning the classification associated with this leaf.

Example:-

Let's say we have a sample of 50 students with three variables Gender (Boy/ Girl), Class(X/ XI) and Height (5 to 6 ft). 20 out of these 50 play cricket in rest time. Suppose you want to find on unknown dataset which contains all the features(Gender, class, height) that he/she will play or not in rest time.

This is where decision tree supports, it will separate the students based on all values of three variable and identify the variable, which creates the best uniform sets of students

6)

Issues in learning decision trees include

- 1 Avoiding Overfitting the Data
 - 2 Reduced error pruning
- 3 Rule post-pruning
- 4 Incorporating Continuous-Valued Attributes
- 5 Alternative Measures for Selecting Attributes
- 6 Handling Training Examples with Missing Attribute Values
- 7 Handling Attributes with Differing Costs

Avoiding Overfitting the Data

- 8 The ID3 algorithm grows each branch of the tree just deeply enough to perfectly classify the training examples but it can lead to difficulties when there is noise in the data, or when the number of training examples is too small to produce a representative sample of the true target function. This algorithm can produce trees that overfit the training examples.

How can it be possible for tree h to fit the training examples better than h' , but for it to perform more poorly over subsequent examples?

- 9 Overfitting can occur when the training examples contain random errors or noise
- 10 When small numbers of examples are associated with leaf nodes.

Approaches to avoiding overfitting in decision tree learning

- 11 Pre-pruning (avoidance): Stop growing the tree earlier, before it reaches the point where it perfectly classifies the training data
- 12 Post-pruning (recovery): Allow the tree to overfit the data, and then post-prune the tree

Criterion used to determine the correct final tree size

- 13 Use a separate set of examples, distinct from the training examples, to evaluate the utility of post-pruning nodes from the tree
- 14 Use all the available data for training, but apply a statistical test to estimate whether expanding (or pruning) a particular node is likely to produce an improvement beyond the training set
- 15 Use measure of the complexity for encoding the training examples and the decision tree, halting growth of the tree when this encoding size is minimized. This approach is called the Minimum Description Length

$$MDL - \text{Minimize} : \text{size}(\text{tree}) + \text{size}(\text{misclassifications}(\text{tree}))$$

Reduced-Error Pruning

- 16 Reduced-error pruning, is to consider each of the decision nodes in the tree to be candidates for pruning
- 17 **Pruning** a decision node consists of removing the subtree rooted at that node, making it a

10 CO2 L3

leaf node, and assigning it the most common classification of the training examples affiliated with that node

- 18 Nodes are removed only if the resulting pruned tree performs no worse than the original over the validation set.
- 19 Reduced error pruning has the effect that any leaf node added due to coincidental regularities in the training set is likely to be pruned because these same coincidences are unlikely to occur in the validation set

2. Incorporating Continuous-Valued Attributes

Continuous-valued decision attributes can be incorporated into the learned tree.

There are two methods for Handling Continuous Attributes

20 Define new discrete valued attributes that partition the continuous attribute value into a discrete set of intervals.

E.g., {high \equiv Temp > 35° C, med \equiv 10° C < Temp \leq 35° C, low \equiv Temp \leq 10° C}

21 Using thresholds for splitting nodes

e.g., $A \leq a$ produces subsets $A \leq a$ and $A > a$

What threshold-based Boolean attribute should be defined based on Temperature?

Temperature:	40	48	60	72	80	90
PlayTennis:	No	No	Yes	Yes	Yes	No

- 22 Pick a threshold, c, that produces the greatest information gain
- 23 In the current example, there are two candidate thresholds, corresponding to the values of Temperature at which the value of PlayTennis changes: $(48 + 60)/2$, and $(80 + 90)/2$.
- 24 The information gain can then be computed for each of the candidate attributes, Temperature >54, and Temperature >85 and the best can be selected (Temperature >54)

3. Alternative Measures for Selecting Attributes

- The problem is if attributes with many values, Gain will select it ?
- Example: consider the attribute Date, which has a very large number of possible values. (e.g., March 4, 1979).
- If this attribute is added to the PlayTennis data, it would have the highest information gain of any of the attributes. This is because Date alone perfectly predicts the target attribute over the training data. Thus, it would be selected as the decision attribute for the root node of the tree and lead to a tree of depth one, which perfectly classifies the training data.
- This decision tree with root node Date is not a useful predictor because it perfectly separates the training data, but poorly predict on subsequent examples.

4. Handling Training Examples with Missing Attribute Values

The data which is available may contain missing values for some attributes Example: Medical diagnosis

- <Fever = true, Blood-Pressure = normal, ..., Blood-Test = ?, ...>
- Sometimes values truly unknown, sometimes low priority (or cost too high)

Strategies for dealing with the missing attribute value

- If node n test A, assign most common value of A among other training examples sorted to node n
- Assign most common value of A among other training examples with same target value
- Assign a probability p_i to each of the possible values v_i of A rather than simply assigning the most common value to $A(x)$

5. Handling Attributes with Differing Costs

- In some learning tasks the instance attributes may have associated costs.
- For example: In learning to classify medical diseases, the patients described in terms of attributes such as Temperature, BiopsyResult, Pulse, BloodTestResults, etc.
- These attributes vary significantly in their costs, both in terms of monetary cost and cost to patient comfort
- Decision trees use low-cost attributes where possible, depends only on high-cost attributes only when needed to produce reliable classifications.

7)

Step 1:

Total = 14 Yes(p) = 9 No(n) = 5
 Attributes: Outlook = {Sunny, Overcast, Rain}
 Mild, Cool
 Humidity = {High, Normal}

Temperature = {Hot,

Wind = {Weak, Strong}

Step 2: Calculate the entropy of the dataset

$$\begin{aligned}
 \text{Entropy}(S) &= -p_+ \log_2 p_+ - p_- \log_2 p_- \\
 &= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right) \\
 &= -(9/(9+5)) \log(9/(9+5)) - (5/(9+5)) \log(5/(9+5)) \\
 &= -(9/14) \log(9/14) - (5/14) \log(5/14) \\
 &= (-0.643)(-0.637) - (0.357)(-1.486) \\
 &= 0.940
 \end{aligned}$$

Step 3:

i)

Select Outlook attribute

Outlook = {Sunny, Overcast, Rain}

Sunny : Yes(p) = 2

No(n) = 3

Overcast: Yes(p) = 4

No(n) = 0

Rain: Yes(p) = 3

No(n) = 2

a) Entropy of Outlook attribute

$$\begin{aligned}
 &= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right) \\
 \text{Entropy}(\text{Outlook} = \text{Sunny}) &= -(2/5) \log(2/5) - (3/5) \log(3/5) \\
 &= -(0.4)(-1.322) - (0.6)(-0.737) = 0.971 \\
 \text{Entropy}(\text{Outlook} = \text{Overcast}) &= -(4/4) \log(4/4) - 0 = 0 \\
 \text{Entropy}(\text{Outlook} = \text{Rain}) &= -(3/5) \log(3/5) - (2/5) \log(2/5) \\
 &= -(0.6)(-0.737) - (0.4)(-1.322) = 0.971
 \end{aligned}$$

b) Average Information Entropy(I)

$$\begin{aligned}
 I(\text{Outlook}) &= ((2+3)/(9+5)) * 0.971 + ((3+2)/(9+5)) * 0.971 + 0 \\
 &= (5/14) * 0.971 + 0.3571 * 0.971 \\
 &= 0.693
 \end{aligned}$$

c) Information Gain(Outlook) = Entropy(S) - I(Outlook)

$$= 0.940 - 0.693 = 0.247$$

ii)

Select Temperature attribute

Temperature = {Hot, Mild, Cool}

Hot : p=2

n=2

Mild: p=4

n=2

Cool: p=3

n=1

a. Calculate the entropy for Temperature

$$\begin{aligned}
 &= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right) \\
 \text{Entropy}(\text{Temperature} = \text{Hot}) &= -(2/4) \log(2/4) - (2/4) \log(2/4) \\
 &= -(0.5)(-1) - (0.5)(-1) \\
 &= 1
 \end{aligned}$$

10

CO2

L5

$$\begin{aligned} \text{Entropy(Temperature = Mild)} &= -(4/6)\log(4/6) - (2/6)\log(2/6) \\ &= -(0.66)(-0.599) - (0.33)(-1.599) \\ &= 0.923 \end{aligned}$$

$$\begin{aligned} \text{Entropy(Temperature = Cool)} &= -(3/4)\log(3/4) - (1/4)\log(1/4) \\ &= -(0.75)(-0.415) - (0.25)(-2) \\ &= 0.811 \end{aligned}$$

b. Average Information Entropy(I)

$$I(\text{Temperature}) = (4/14)*1 + (6/14)*0.923 + (4/14)*0.811 = 0.913$$

c. Information Gain(Temperature)

$$\begin{aligned} \text{IG(Temperature)} &= \text{Entropy(S)} - I(\text{Temperature}) \\ &= 0.940 - 0.913 \\ &= 0.027 \end{aligned}$$

iii) Select Humidity attribute

Humidity = {High, Normal}

High: p: 3 n:4

Normal: p: 6 n: 1

a. Calculate the entropy for Temperature

$$\begin{aligned} &= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right) \\ \text{Entropy(Humidity = High)} &= -(3/7)\log(3/7) - (4/7)\log(4/7) \\ &= -(0.4286)(-1.2223) - (0.5714)(-0.8074) \\ &= 0.985 \end{aligned}$$

$$\begin{aligned} \text{Entropy(Humidity = Normal)} &= -(6/7)\log(6/7) - (1/7)\log(1/7) \\ &= -(0.8571)(-0.2225) - (0.1429)(-2.8069) = \\ &0.591 \end{aligned}$$

b. Average Information Entropy

c. Average Information Entropy(I)

$$\begin{aligned} I(\text{Humidity}) &= (7/14)*0.985 + (7/14)*0.591 \\ &= 0.788 \end{aligned}$$

d. Information Gain(Humidity)

$$\begin{aligned} \text{IG(Humidity)} &= \text{Entropy(S)} - I(\text{Humidity}) \\ &= 0.940 - 0.788 \\ &= 0.152 \end{aligned}$$

iv) Select Windy attribute

Wind = {Weak, Strong}

Weak: p:6 n:2

Strong: p:3 n:3

a. Calculate the entropy for Windy

$$\begin{aligned} &= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right) \\ \text{Entropy(Windy = Weak)} &= -(6/8)\log(6/8) - (2/8)\log(2/8) \\ &= 0.811 \end{aligned}$$

$$\begin{aligned} \text{Entropy(Windy = Strong)} &= -(3/6)\log(3/6) - (3/6)\log(3/6) \\ &= 1 \end{aligned}$$

b. Average Information Entropy(I)

$$\begin{aligned} I(\text{Windy}) &= (8/14)*0.811 + (6/14)*1 \\ &= 0.892 \end{aligned}$$

c. Information Gain(Windy)

$$\begin{aligned} \text{IG(Windy)} &= \text{Entropy(S)} - I(\text{Windy}) \\ &= 0.940 - 0.892 = 0.048 \end{aligned}$$

IG(Outlook) = 0.247
IG(Temperature) = 0.27
IG(Humidity) = 0.152
IG(Windy) = 0.048

Highest Information Gain is 0.247 -> Outlook

P:2	N:3	Total:5
Temperature= {hot, cool, mild}		
Hot:p:0		n:2
Cool: p:1		n:0
Mild: p:1		n:1
Humidity={High, Normal}		
High: p:0		n:3
Normal:p:2		n:0
Windy:{Weak, Strong}		
Weak: p:1		n:2
Strong: p:1		n:1

1. Calculate the entropy of Dataset(S)

$$= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right)$$

$$\text{Entropy} = -(2/5)\log(2/5) - (3/5)\log(3/5) = 0.971$$

2. Calculate the Information Gain

a. Calculate entropy of humidity

$$\text{Entropy(Humidity = High)} = 0 - (3/3)\log(3/3) = 0$$

$$\text{Entropy(Humidity = Normal)} = 0$$

b. Calculate Average information entropy(I) of humidity

$$I(\text{Humidity}) = 0$$

c. Information gain of humidity

$$\text{IG(Humidity)} = \text{Entropy(S)} - I(\text{Humidity})$$

$$= 0.971 - 0 = 0.971$$

d. Calculate entropy of Windy

Windy:{Weak, Strong}

Weak: p:1 **n:2**

Strong: p:1 **n:1**

$$= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right)$$

$$\text{Entropy(Windy = Weak)} = -(1/3)\log(1/3) - (2/3)\log(2/3) = 0.918$$

$$\text{Entropy(wind = Strong)} = -(1/2)\log(1/2) - (1/2)\log(1/2) = 1$$

e. Calculate Average information entropy of windy

$$I(\text{Windy}) = (3/5) * 0.918 + (2/5) * 1 = 0.951$$

f. Information gain of windy

$$\text{IG(Windy)} = \text{Entropy(S)} - I(\text{Windy})$$

$$= 0.971 - 0.951 = 0.020$$

g. Calculate entropy of temperature

Temperature = {hot, cool, mild}

Hot: p:0 n:2

Cool: p:1 n:0

Mild: p:1 n:1

$$= -\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right)$$

Entropy (Temperature = hot) = 0

Entropy (Temperature = Cool) = 0

Entropy (Temperature = mild) = $-(1/2)\log(1/2) - (1/2)\log(1/2) = 1$

h. Calculate Average information entropy of temperature

$I(\text{temperature}) = (2/5) * 0 + (1/5)*0 + (2/5)*1 = 0.4$

i. Information gain of temperature

$IG(\text{Temperature}) = 0.971 - 0.4 = 0.571$

3. Select the attribute with highest information gain

$IG(\text{Temperature}) = 0.971 - 0.4 = 0.571$

$IG(\text{Windy}) = 0.020$

$IG(\text{Humidity}) = 0.971$

Total = 5

P=3

N= 2

1. Calculate the entropy of the dataset(S)

$$\frac{p}{p+n} \log\left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log\left(\frac{n}{p+n}\right)$$

Entropy = $-(3/5)\log(3/5) - (2/5)\log(2/5) = 0.971$

2. Calculate the information gain

a. Calculate entropy of temperature

Temperature = {mild, cool}

Mild: p:2 n:1

Cool: p:1 n:1

Entropy (temperature = mild) = $-(2/3)\log(2/3) - (1/3)\log(1/3) = 0.918$

Entropy (temperature = cool) = $-(1/2)\log(1/2) - (1/2)\log(1/2) = 1$

b. Calculate average information entropy of temperature

$I(\text{Temperature}) = 0.951$

c. Information gain of temperature

$0.971 - 0.951 = 0.020$

d. Calculate entropy of Humidity

Entropy (Humidity = High) = 1

Entropy (Humidity = Normal) = 0.918

e. Calculate average information entropy of humidity

$I(\text{Humidity}) = 0.951$

f. Information gain of humidity

Gain = $0.971 - 0.951 = 0.020$

g. Calculate entropy of Windy

Entropy (Windy = Strong) = 0

Entropy (Windy = Weak) = 0

h. Calculate average information entropy of Windy

$I(\text{Windy}) = 0$

i. Information gain of Windy

$$\text{Gain} = 0.971 - 0 = 0.971$$

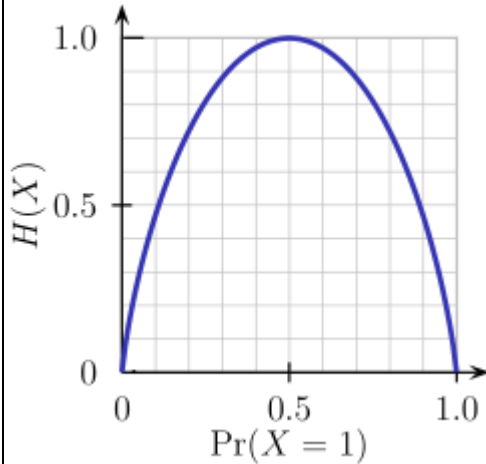
3. Select the attribute with highest information gain

Select Windy.

8)

Entropy

Entropy is a measure of the randomness in the information being processed. The higher the entropy, the harder it is to draw any conclusions from that information. Flipping a coin is an example of an action that provides information that is random.



From the above graph, it is quite evident that the entropy $H(X)$ is zero when the probability is either 0 or 1. The Entropy is maximum when the probability is 0.5 because it projects perfect randomness in the data and there is no chance if perfectly determining the outcome.

ID3 follows the rule — A branch with an entropy of zero is a leaf node and A branch with entropy more than zero needs further splitting.

Mathematically Entropy for 1 attribute is represented as:

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

Play Golf	
Yes	No
9	5



$$\begin{aligned} \text{Entropy(PlayGolf)} &= \text{Entropy}(5,9) \\ &= \text{Entropy}(0.36, 0.64) \\ &= - (0.36 \log_2 0.36) - (0.64 \log_2 0.64) \\ &= 0.94 \end{aligned}$$

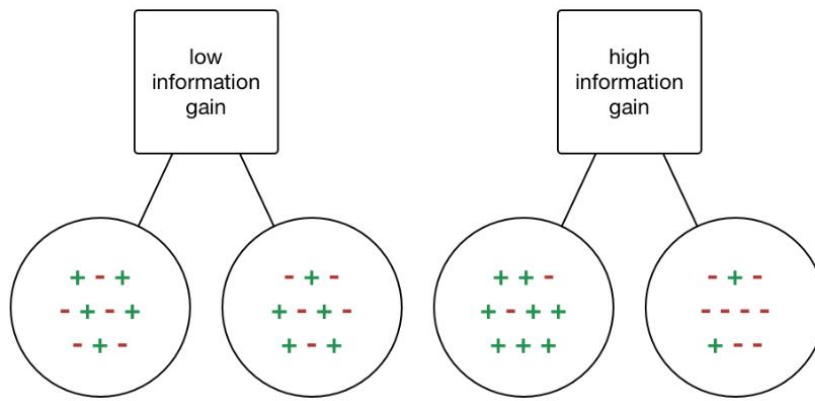
Information Gain

Information gain or **IG** is a statistical property that measures how well a given attribute separates the training examples according to their target classification. Constructing a decision tree is all about finding an attribute that returns the highest information gain and the smallest entropy.

10

CO2

L2



Information gain is a decrease in entropy. It computes the difference between entropy before split and average entropy after split of the dataset based on given attribute values. ID3 (Iterative Dichotomiser) decision tree algorithm uses information gain.

Mathematically, IG is represented as:

$$\text{Information Gain}(T, X) = \text{Entropy}(T) - \text{Entropy}(T, X)$$

$$\begin{aligned} \text{IG}(\text{PlayGolf}, \text{Outlook}) &= E(\text{PlayGolf}) - E(\text{PlayGolf}, \text{Outlook}) \\ &= 0.940 - 0.693 \\ &= 0.247 \end{aligned}$$