

## VTU Solution – Jan/Feb 2023

<b>Sub:</b>	Artificial Intelligence & Machine Learning						<b>Code:</b>	18CS71	
<b>Date:</b>	16/1/2023	<b>Duration:</b>	3 Hours	<b>Max Marks:</b>	100	<b>Sem:</b>	VII	<b>Branch:</b>	ISE

<b>Q.n</b>	<b>Solutions</b>	<b>Marks</b>
1a	<p>The Water Jug Problem states: You are provided with two jugs, one with 4-gallons of capacity and the other with one with 3-gallons of capacity. Neither have any measuring markers on it. How can we get exactly two gallons of water in the 4-gallon jug</p> <ul style="list-style-type: none"> <li>• Write down the production rules for the above problem</li> <li>• Write any one solution to the above problem</li> </ul> <p>Problem is “You are given two jugs, a 4-litre one and a 3-litre one. One neither has any measuring markers on it. There is a pump that can be used to fill the jugs with water. How can you get exactly 2 litres of water into 4-litre jug?”</p> <p>Solution:</p> <p>The state space for the problem can be described as a set of states, where each state represents the number of gallons in each state. The game start with the initial state described as a set of ordered pairs of integers:</p> <ul style="list-style-type: none"> <li>• State: (x, y) <ul style="list-style-type: none"> <li>– x = number of lts in 4 lts jug</li> <li>– y = number of lts in 3 lts jug</li> </ul> </li> </ul> <p>x = 0, 1, 2, 3, or 4 y = 0, 1, 2, 3</p> <ul style="list-style-type: none"> <li>• Start state: (0, 0) i.e., 4-litre and 3-litre jugs is empty initially.</li> <li>• Goal state: (2, n) for any n that is 4-litre jug has 2 litres of water and 3-litre jug has any value from 0-3 since it is not specified.</li> <li>• Attempting to end up in a goal state.</li> </ul> <p>Production Rules: These rules are used as operators to solve the problem. They are represented as rules whose left sides are used to describe new state that result from approaching the rule.</p>	08

1	$(x, y)$ if $x < 4$	$\rightarrow (4, y)$	Fill the 4-gallon jug
2	$(x, y)$ if $y < 3$	$\rightarrow (x, 3)$	Fill the 3-gallon jug
3	$(x, y)$ if $x > 0$	$\rightarrow (x - d, y)$	Pour some water out of the 4-gallon jug
4	$(x, y)$ if $y > 0$	$\rightarrow (x, y - d)$	Pour some water out of the 3-gallon jug
5	$(x, y)$ if $x > 0$	$\rightarrow (0, y)$	Empty the 4-gallon jug on the ground
6	$(x, y)$ if $y > 0$	$\rightarrow (x, 0)$	Empty the 3-gallon jug on the ground
7	$(x, y)$ if $x + y \geq 4$ and $y > 0$	$\rightarrow (4, y - (4 - x))$	Pour water from the 3-gallon jug into the 4-gallon jug until the 4-gallon jug is full
8	$(x, y)$ if $x + y \geq 3$ and $x > 0$	$\rightarrow (x - (3 - y), 3)$	Pour water from the 4-gallon jug into the 3-gallon jug until the 3-gallon jug is full
9	$(x, y)$ if $x + y \leq 4$ and $y > 0$	$\rightarrow (x + y, 0)$	Pour all the water from the 3-gallon jug into the 4-gallon jug
10	$(x, y)$ if $x + y \leq 3$ and $x > 0$	$\rightarrow (0, x + y)$	Pour all the water from the 4-gallon jug into the 3-gallon jug
11	$(0, 2)$	$\rightarrow (2, 0)$	Pour the 2 gallons from the 3-gallon jug into the 4-gallon jug
12	$(2, y)$	$\rightarrow (0, y)$	Empty the 2 gallons in the 4-gallon jug on the ground

The solution to the water-jug problem is:

Gallons in the 4-Gallon Jug	Gallons in the 3-Gallon Jug	Rule Applied
0	0	
0	3	2
3	0	9
3	3	2
4	2	7
0	2	5 or 12
2	0	9 or 11

**1 b. Explain steepest Ascent Hill Climbing Technique with an algorithm. Comment on its drawback and how to overcome these drawbacks**

A useful variation on simple hill climbing considers all the moves from the current state and selects the best one as the next state. This method is called *steepest-ascent hill climbing* or *gradient search*

Algorithm:

1. Evaluate the initial state. If it is also a goal state then return it and quit. Otherwise continue with the initial state as the current state.
2. Loop until a solution is found or until a complete iteration produces no change to current state:

- a. Let SUCC be a state such that any possible successor of the current state will be better than SUCC.
- b. For each operator that applies to the current state do:
  - i. Apply the operator and generate a new state.
  - ii. Evaluate the new state. If it is a goal state, then return it and quit. If not compare it to SUCC. If it is better, then set SUCC to this state. If it is not better, leave SUCC alone.
- c. IF the SUCC is better than current state, then set current state to SUCC.

Both basic and steepest-ascent hill climbing may fail to find a solution. Either algorithm may terminate not by finding a goal state but by getting a state from which no better states can be generated. This will happen if the program has reached a local maximum, a plateau or a ridge.

A *local maximum* is a state that is better than all its neighbors but it not better than some other states farther away. At the local maximum, all moves appear to make things worse. Local maxima are particularly frustrating because they often occur almost within sight of a solution. In this case, they are called *foothills*.

A *plateau* is a flat area of the search space in which a whole set of neighboring states has the same value. In this, it is not possible to determine the best direction in which to move by making local comparisons.

A *ridge* is a special kind of maximum. It is an area of the search space that is higher than surrounding areas and that itself has a slope.

There are some ways of dealing with these problems, although these methods are by no means guaranteed:

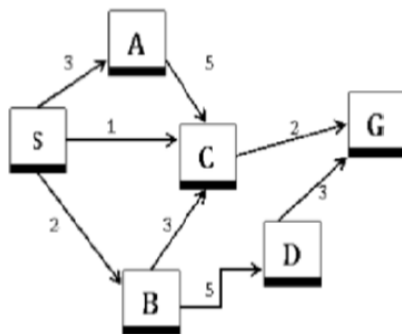
- Backtrack to some earlier node and try going in a different direction. This is particularly reasonable if at that node there was another direction that looked as promising or almost as promising as the one that was chosen earlier. This is a fairly good way to deal with local maxima.
- Make a big jump in some direction to try to get to a new section of the search space. This is a good way of dealing with plateaus.
- Apply two or more rules before doing the test. This corresponds to moving in several directions at once. This is a good strategy for dealing with ridges.



2a Explain problem reduction with respect to AND-OR graph with suitable example.

7

State	$h(n)$
S	6
A	2
B	4
C	1
D	2
G	0



#### AND-OR Graphs:

- Useful for representing the solution of problems that can be solved by **decomposing them into a set of smaller problems**, all of which must be then solved.

Some problems are best represented as achieving sub goals, some of

	<p>which achieved simultaneously and independently (AND) Up to now, only dealt with OR options</p> <p><b>AO* Algorithm:</b></p> <ol style="list-style-type: none"> <li>1. Initialize the graph to start node</li> <li>2. Traverse the graph following the current path accumulating nodes that have not yet been expanded or solved</li> <li>3. Pick any of these nodes and expand it and if it has no successors call this value <i>FUTILITY</i> otherwise calculate only <math>f</math> for each of the successors.</li> <li>4. If <math>f</math> is 0 then mark the node as <i>SOLVED</i></li> <li>5. Change the value of <math>f</math> for the newly created node and let <math>f</math> reflect on its predecessors by back propagation.</li> <li>6. Wherever possible use the most promising routes and if all descendants of a node is marked as <i>SOLVED</i> then mark the parent node as <i>SOLVED</i>.</li> </ol> <p>If starting node is <i>SOLVED</i> or value greater than <i>FUTILITY</i>, stop, else repeat from 2.</p> <p>Current state is S</p> <ul style="list-style-type: none"> <li>• <math>f(A) = 3+2 = 5</math>, <math>f(B) = 2+4=6</math>, <math>f(C)=1+1=2</math></li> <li>• Since, C is having smaller distance compared to other nodes, we have chosen current state as C.</li> </ul> <p>2.Current State is C</p> <ul style="list-style-type: none"> <li>• <math>f(C) = 2+0 = 2</math></li> </ul> <p>Path = S -&gt; C -&gt; G</p>	
2b	<p>Write AO algorithm</p> <p><b>AO* Algorithm:</b></p> <ol style="list-style-type: none"> <li>1. Initialize the graph to start node</li> <li>2. Traverse the graph following the current path accumulating nodes that have not yet been expanded or solved</li> <li>3. Pick any of these nodes and expand it and if it has no successors call this value <i>FUTILITY</i> otherwise calculate only <math>f</math> for each of the successors.</li> <li>4. If <math>f</math> is 0 then mark the node as <i>SOLVED</i></li> <li>5. Change the value of <math>f</math> for the newly created node and let <math>f</math> reflect on its predecessors by back propagation.</li> <li>6. Wherever possible use the most promising routes and if all descendants of a node is marked as <i>SOLVED</i> then mark the parent node as <i>SOLVED</i>.</li> </ol> <p>If starting node is <i>SOLVED</i> or value greater than <i>FUTILITY</i>, stop, else repeat from 2.</p>	7
2c	<p>Discuss about constraint satisfaction and solve the below crypt arithmetic problem.</p> <p>CROSS+ ROADS = DANGER</p>	6

## Cryptarithmic Problem

9	6	2	3	3	
+	6	2	5	1	3
1	5	8	7	4	6
D	A	N	G	E	R

C5 C4 C3 C2 C1

C R O S S  
+ R O A D S  
-----  
D A N G E R

Character	Code
D	1
C	9
R	6
A	5
O	2
N	8
G	7
S	3
E	4

3

Consider the following sentences:

1. John likes all kinds of food
  2. Apples are food
  3. Chicken is food
  4. Anything anyone eats and isn't killed is food
  5. Bill eats peanuts and is still alive
  6. Sue eats every everything Bill eats.
- (i) Translate these sentences into formulas in predicate logic.  
(ii) Prove that John likes peanuts using backward chaining.  
(iii) Convert the formulas of (i) into clause form.  
(iv) Prove John likes peanuts using resolution.

### Step-1: Conversion of Facts into FOL

In the first step we will convert all the given statements into its first order logic.

- a.  $\forall x: \text{food}(x) \rightarrow \text{likes}(\text{John}, x)$
  - b.  $\text{food}(\text{Apple}) \wedge \text{food}(\text{vegetables})$
  - c.  $\forall x \forall y: \text{eats}(x, y) \wedge \neg \text{killed}(x) \rightarrow \text{food}(y)$
  - d.  $\text{eats}(\text{Anil}, \text{Peanuts}) \wedge \text{alive}(\text{Anil})$ .
  - e.  $\forall x : \text{eats}(\text{Anil}, x) \rightarrow \text{eats}(\text{Harry}, x)$
  - f.  $\forall x: \neg \text{killed}(x) \rightarrow \text{alive}(x)$
  - g.  $\forall x: \text{alive}(x) \rightarrow \neg \text{killed}(x)$
  - h.  $\text{likes}(\text{John}, \text{Peanuts})$
- } added predicates.

### Step-2: Conversion of FOL into CNF

In First order logic resolution, it is required to convert the FOL into CNF as CNF form makes easier for resolution proofs.

5

5

○ **Eliminate all implication ( $\rightarrow$ ) and rewrite**

- a.  $\forall x \neg \text{food}(x) \vee \text{likes}(\text{John}, x)$
- b.  $\text{food}(\text{Apple}) \wedge \text{food}(\text{vegetables})$
- c.  $\forall x \forall y \neg [\text{eats}(x, y) \wedge \neg \text{killed}(x)] \vee \text{food}(y)$
- d.  $\text{eats}(\text{Anil}, \text{Peanuts}) \wedge \text{alive}(\text{Anil})$
- e.  $\forall x \neg \text{eats}(\text{Anil}, x) \vee \text{eats}(\text{Harry}, x)$
- f.  $\forall x \neg [\neg \text{killed}(x)] \vee \text{alive}(x)$
- g.  $\forall x \neg \text{alive}(x) \vee \neg \text{killed}(x)$
- h.  $\text{likes}(\text{John}, \text{Peanuts})$ .

○ **Eliminate all implication ( $\rightarrow$ ) and rewrite**

- a.  $\forall x \neg \text{food}(x) \vee \text{likes}(\text{John}, x)$
- b.  $\text{food}(\text{Apple}) \wedge \text{food}(\text{vegetables})$
- c.  $\forall x \forall y \neg [\text{eats}(x, y) \wedge \neg \text{killed}(x)] \vee \text{food}(y)$
- d.  $\text{eats}(\text{Anil}, \text{Peanuts}) \wedge \text{alive}(\text{Anil})$
- e.  $\forall x \neg \text{eats}(\text{Anil}, x) \vee \text{eats}(\text{Harry}, x)$
- f.  $\forall x \neg [\neg \text{killed}(x)] \vee \text{alive}(x)$
- g.  $\forall x \neg \text{alive}(x) \vee \neg \text{killed}(x)$
- h.  $\text{likes}(\text{John}, \text{Peanuts})$ .

○ **Move negation ( $\neg$ )inwards and rewrite**

- a.  $\forall x \neg \text{food}(x) \vee \text{likes}(\text{John}, x)$
- b.  $\text{food}(\text{Apple}) \wedge \text{food}(\text{vegetables})$
- c.  $\forall x \forall y \neg \text{eats}(x, y) \vee \text{killed}(x) \vee \text{food}(y)$
- d.  $\text{eats}(\text{Anil}, \text{Peanuts}) \wedge \text{alive}(\text{Anil})$
- e.  $\forall x \neg \text{eats}(\text{Anil}, x) \vee \text{eats}(\text{Harry}, x)$
- f.  $\forall x \neg \text{killed}(x) \vee \text{alive}(x)$
- g.  $\forall x \neg \text{alive}(x) \vee \neg \text{killed}(x)$
- h.  $\text{likes}(\text{John}, \text{Peanuts})$ .

	<ul style="list-style-type: none"> <li>◦ <b>Eliminate existential instantiation quantifier by elimination.</b> In this step, we will eliminate existential quantifier <math>\exists</math>, and this process is known as <b>Skolemization</b>. But in this example problem since there is no existential quantifier so all the statements will remain same in this step.</li> <li>◦ <b>Drop Universal quantifiers.</b> In this step we will drop all universal quantifier since all the statements are not implicitly quantified so we don't need it. <ul style="list-style-type: none"> <li>a. <math>\neg \text{food}(x) \vee \text{likes}(\text{John}, x)</math></li> <li>b. <math>\text{food}(\text{Apple})</math></li> <li>c. <math>\text{food}(\text{vegetables})</math></li> <li>d. <math>\neg \text{eats}(y, z) \vee \text{killed}(y) \vee \text{food}(z)</math></li> <li>e. <math>\text{eats}(\text{Anil}, \text{Peanuts})</math></li> <li>f. <math>\text{alive}(\text{Anil})</math></li> <li>g. <math>\neg \text{eats}(\text{Anil}, w) \vee \text{eats}(\text{Harry}, w)</math></li> <li>h. <math>\text{killed}(g) \vee \text{alive}(g)</math></li> <li>i. <math>\neg \text{alive}(k) \vee \neg \text{killed}(k)</math></li> <li>j. <math>\text{likes}(\text{John}, \text{Peanuts})</math>.</li> </ul> </li> </ul>	5
4a	<p>Distinguish forward and backward reasoning with an example</p> <p><b>Definition of Forward Reasoning</b></p> <p>The solution of a problem generally includes the initial data and facts in order to arrive at the solution. These unknown facts and information is used to deduce the result.</p> <p>For example, while diagnosing a patient the doctor first check the symptoms and medical condition of the body such as temperature, blood pressure, pulse, eye colour, blood, etcetera. After that, the patient symptoms are analysed and compared against the predetermined symptoms.</p> <p>Then the doctor is able to provide the medicines according to the symptoms of the patient. So, when a solution employs this manner of reasoning, it is known as forward reasoning.</p> <p><b>Definition of Backward Reasoning</b></p> <p>The <b>backward reasoning</b> is inverse of forward reasoning in which goal is analysed in order to deduce the rules, initial facts and data.</p> <p>We can understand the concept by the similar example given in the above definition, where the doctor is trying to diagnose the patient with the help of the inceptive data such as symptoms.</p> <p>However, in this case, the patient is experiencing a problem in his body, on the basis of which the doctor is going to prove the symptoms.</p>	04

This kind of reasoning comes under backward reasoning.

4b

Find maximally specific hypothesis for the training instances given below. Also write Find-S algorithm. The concept of this particular problem will be on what days does a person lines to go on walk

Time	Weather	Temperature	Company	Humidity	Wind	Goes
Morning	Sunny	Warm	Yes	Mild	Strong	Yes
Evening	Rainy	Cold	No	Mild	Normal	No
Morning	Sunny	Moderate	Yes	Normal	Normal	Yes
Evening	Sunny	Cold	Yes	High	Strong	Yes

Let's have a look at the algorithm of Find-S:

1. Initialize the value of the hypothesis for all attributes with the most specific one. That is,

$$h_0 = \langle \phi, \phi, \phi, \phi, \dots \rangle$$

2. Take the next example, if the taken example is negative leave them and move on to another example without changing our hypothesis for the step.

3. Now, if the taken example is a positive example, then For each attribute, check if the value of the attribute is equal to that of the value we took in our hypothesis.

So now, the general hypothesis is:

$$h_0 = \{ \text{'Morning'}, \text{'Sunny'}, \text{'Warm'}, \text{'Yes'}, \text{'Mild'}, \text{'Strong'} \}$$

This is our general hypothesis, and now we will consider each example one by one, but only the positive examples.

$$h_1 = \{ \text{'Morning'}, \text{'Sunny'}, \text{'?'}, \text{'Yes'}, \text{'?'}, \text{'?'} \}$$

$$h_2 = \{ \text{'?'}, \text{'Sunny'}, \text{'?'}, \text{'Yes'}, \text{'?'}, \text{'?'} \}$$

We replaced all the different values in the general hypothesis to get a resultant hypothesis.

8

5a

Define version space. Discuss the limitations of finds algorithm over candidate elimination algorithm

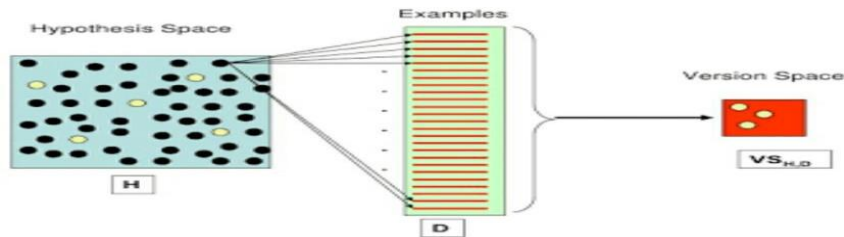
A representation of the set of all hypotheses which are consistent with D

Definition: The version space, denoted  $V_{SH,D}$  with respect to hypothesis space H and training examples D, is the subset of hypotheses from H consistent with the training examples in D

$$V_{SH,D} \equiv \{ h \in H \mid \text{Consistent}(h, D) \}$$

7





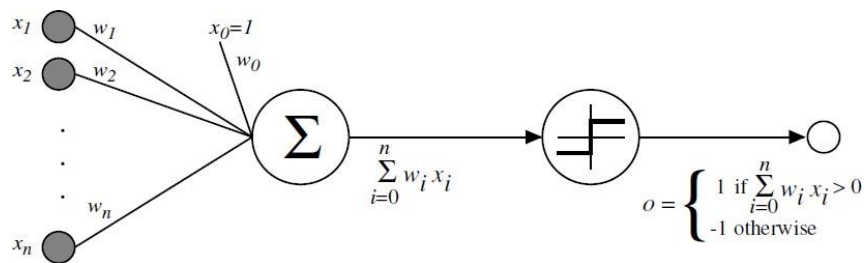
There are a few limitations of the Find-S algorithm listed down below:

1. There is no way to determine if the hypothesis is consistent throughout the data.
2. Inconsistent training sets can actually mislead the Find-S algorithm, since it ignores the negative examples.
3. Find-S algorithm does not provide a backtracking technique to determine the best possible changes that could be done to improve the resulting hypothesis.

5b

How a single perceptron can be used to represent the Boolean functions such as AND, OR.

- One type of ANN system is based on a unit called a perceptron. Perceptron is a single layer neural network.



- A perceptron takes a vector of real-valued inputs, calculates a linear combination of these inputs, then outputs a 1 if the result is greater than some threshold and -1 otherwise.
- Given inputs  $x_1$  through  $x_n$ , the output  $O(x_1, \dots, x_n)$  computed by the perceptron is

$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1x_1 + \dots + w_nx_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

- Where, each  $w_i$  is a real-valued constant, or weight, that determines the contribution of input  $x_i$  to the perceptron output.
- $-w_0$  is a threshold that the weighted combination of inputs  $w_1x_1 + \dots + w_nx_n$  must surpass in order for the perceptron to output a 1.

Sometimes, the perceptron function is written as,

$$O(\vec{x}) = \text{sgn}(\vec{w} \cdot \vec{x})$$

Where,

$$\text{sgn}(y) = \begin{cases} 1 & \text{if } y > 0 \\ -1 & \text{otherwise.} \end{cases}$$

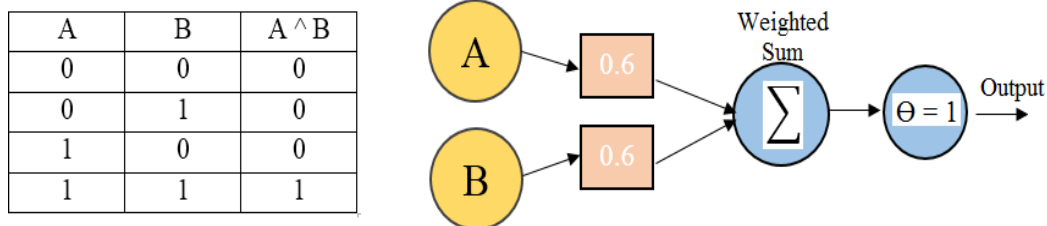
Learning a perceptron involves choosing values for the weights  $w_0, \dots, w_n$ . Therefore, the space  $H$  of candidate hypotheses considered in perceptron learning is the set of all possible real-valued weight vectors

$$H = \{\vec{w} \mid \vec{w} \in \mathcal{R}^{(n+1)}\}$$

Perceptrons can represent all of the primitive Boolean functions AND, OR, NAND ( $\sim$  AND), and NOR ( $\sim$ OR)

Some Boolean functions cannot be represented by a single perceptron, such as the XOR function whose value is 1 if and only if  $x_1 \neq x_2$

Example: Representation of AND functions



If  $A=0$  &  $B=0 \rightarrow 0*0.6 + 0*0.6 = 0$ .

This is not greater than the threshold of 1, so the output = 0.

If  $A=0$  &  $B=1 \rightarrow 0*0.6 + 1*0.6 = 0.6$ .

This is not greater than the threshold, so the output = 0.

If  $A=1$  &  $B=0 \rightarrow 1*0.6 + 0*0.6 = 0.6$ .

This is not greater than the threshold, so the output = 0.

If  $A=1$  &  $B=1 \rightarrow 1*0.6 + 1*0.6 = 1.2$ .

This exceeds the threshold, so the output = 1.

5c

Write Gradient Descent algorithm to train a linear unit along with the derivation

Consider locally weighted regression in which the target function  $f$  is approximated near  $x_q$  using a linear function of the form

$$\hat{f}(x) = w_0 + w_1 a_1(x) + \dots + w_n a_n(x)$$

Where,  $a_i(x)$  denotes the value of the  $i$ th attribute of the instance  $x$

Derived methods are used to choose weights that minimize the squared error summed over the set  $D$  of training examples using gradient descent

$$E \equiv \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x))^2$$

Which led us to the gradient descent training rule

Where,  $\eta$  is a constant learning rate

$$\Delta w_j = \eta \sum_{x \in D} (f(x) - \hat{f}(x)) a_j(x)$$

- Need to modify this procedure to derive a local approximation rather than a global one. The simple way is to redefine the error criterion E to emphasize fitting the local training examples. Three possible criteria are given below.

1. Minimize the squared error over just the k nearest neighbors:

$$E_1(x_q) \equiv \frac{1}{2} \sum_{x \in k \text{ nearest nbrs of } x_q} (f(x) - \hat{f}(x))^2 \quad \text{equ(1)}$$

2. Minimize the squared error over the entire set D of training examples, while weighting the error of each training example by some decreasing function K of its distance from  $x_q$  :

$$E_2(x_q) \equiv \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x))^2 K(d(x_q, x)) \quad \text{equ(2)}$$

3. Combine 1 and 2:

$$E_3(x_q) \equiv \frac{1}{2} \sum_{x \in k \text{ nearest nbrs of } x_q} (f(x) - \hat{f}(x))^2 K(d(x_q, x)) \quad \text{equ(3)}$$

1. Combine 1 and 2:

$$E_3(x_q) \equiv \frac{1}{2} \sum_{x \in k \text{ nearest nbrs of } x_q} (f(x) - \hat{f}(x))^2 K(d(x_q, x)) \quad \text{equ(3)}$$

If we choose criterion three and re-derive the gradient descent rule, we obtain the following training rule

$$\Delta w_j = \eta \sum_{x \in k \text{ nearest nbrs of } x_q} K(d(x_q, x)) (f(x) - \hat{f}(x)) a_j(x)$$

The differences between this new rule and the rule given by Equation (3) are that the contribution of instance x to the weight update is now multiplied by the distance penalty  $K(d(x_q, x))$ , and that the error is summed over only the k nearest training examples

6a

What do you mean by Gain and entropy? How it is used to build the decision tree  
Information Gain:

8

- Information gain is the expected reduction in entropy caused by partitioning the examples on an attribute.
- The higher the information gain the more effective the attribute in classifying training data.
- Expected reduction in entropy knowing A

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

Values(A) possible values for A

$S_v$  subset of S for which A has value v

Entropy:

- Entropy measures the impurity of a collection of examples. It depends from the distribution of the random variable p.
  - S is a collection of training examples
  - $p_+$  the proportion of positive examples in S
  - $p_-$  the proportion of negative examples in S

$$\text{Entropy}(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_- \quad [0 \log_2 0 = 0]$$

$$\text{Entropy}([14+, 0-]) = -14/14 \log_2 (14/14) - 0 \log_2 (0) = 0$$

$$\text{Entropy}([9+, 5-]) = -9/14 \log_2 (9/14) - 5/14 \log_2 (5/14) = 0,94$$

$$\begin{aligned} \text{Entropy}([7+, 7-]) &= -7/14 \log_2 (7/14) - 7/14 \log_2 (7/14) = \\ &= 1/2 + 1/2 = 1 \quad [\log_2 1/2 = -1] \end{aligned}$$

Note: the log of a number < 1 is negative,  $0 \leq p \leq 1$ ,  $0 \leq \text{entropy} \leq 1$

Example:

- Let
  - Values(Wind) = {Weak, Strong}
  - $S = [9+, 5-]$
  - $S_{\text{Weak}} = [6+, 2-]$
  - $S_{\text{Strong}} = [3+, 3-]$

- Information gain due to knowing Wind:

$$\begin{aligned} \text{Gain}(S, \text{Wind}) &= \text{Entropy}(S) - 8/14 \text{Entropy}(S_{\text{Weak}}) - 6/14 \text{Entropy}(S_{\text{Strong}}) \\ &= 0,94 - 8/14 \times 0,811 - 6/14 \times 1,00 \\ &= 0,048 \end{aligned}$$

	<p style="text-align: center;"><b>Which attribute is the best classifier?</b></p> <div style="display: flex; justify-content: space-around;"> <div style="text-align: center;"> <p>S: [9+,5-] E = 0.940</p> <div style="border: 1px solid black; padding: 5px; width: 100px; margin: 0 auto;">Humidity</div> <p>High                  Normal</p> <div style="display: flex; justify-content: space-around; margin-top: 20px;"> <div style="text-align: center;"> <p>[3+,4-] E = 0.985</p> </div> <div style="text-align: center;"> <p>[6+,1-] E = 0.592</p> </div> </div> <p>Gain (S, Humidity ) = .940 - (7/14).985 - (7/14).592 = .151</p> </div> <div style="text-align: center;"> <p>S: [9+,5-] E = 0.940</p> <div style="border: 1px solid black; padding: 5px; width: 100px; margin: 0 auto;">Wind</div> <p>Weak                  Strong</p> <div style="display: flex; justify-content: space-around; margin-top: 20px;"> <div style="text-align: center;"> <p>[6+,2-] E = 0.811</p> </div> <div style="text-align: center;"> <p>[3+,3-] E = 1.00</p> </div> </div> <p>Gain (S, Wind) = .940 - (8/14).811 - (6/14)1.0 = .048</p> </div> </div>	
6b	<p><b>Explain Back propagation algorithm? Why is it not likely to be trapped in local minima?</b></p> <p><i>Algorithm:</i></p> <p>Initialize all weights to small random numbers. Until satisfied, Do</p> <ul style="list-style-type: none"> <li>• For each training example, Do <ol style="list-style-type: none"> <li>1. Input the training example to the network and compute the network outputs</li> <li>2. For each output unit <math>k</math> <math display="block">\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)</math> </li> <li>3. For each hidden unit <math>h</math> <math display="block">\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k</math> </li> <li>4. Update each network weight <math>w_{i,j}</math> <math display="block">w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}</math> <p style="margin-left: 40px;">where</p> <math display="block">\Delta w_{i,j} = \eta \delta_j x_{i,j}</math> </li> </ol> </li> </ul> <p>It is not likely to be trapped in a local minima because of the following reasoning:</p>	8

	<p>Gradient descent to some local minimum</p> <ul style="list-style-type: none"> <li>• Perhaps not global minimum...</li> <li>• Add momentum</li> <li>• Stochastic gradient descent</li> <li>• Train multiple nets with different initial weights</li> </ul> <p>Nature of convergence</p> <ul style="list-style-type: none"> <li>• Initialize weights near zero</li> <li>• Therefore, initial networks near-linear</li> <li>• Increasingly non-linear functions possible as training progresses</li> </ul>	
6c	<p><b>Discuss the perceptron training rule and delta rule that solves the learning problem of perceptron</b></p> <p><i>Perceptron training rule:</i></p> $w_i \leftarrow w_i + \Delta w_i$ <p>Where</p> $\Delta w_i = \eta(t - o)x_i$ <p>Where</p> <p><math>\eta</math> – is the learning rate  <math>t</math> – is the target O/P  <math>o</math> – is the output generated by the perceptron  <math>x_i</math> – is the input</p> <p><i>Delta rule:</i></p> $E(\vec{w}) \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$ <p>Where</p> <p><math>D</math> – is the training samples  <math>td</math> – is the target O/P for the training examples  <math>o_d</math> – is the O/P of the linear unit for training examples  <math>E(w)</math> – is half the squared difference between the target O/P, <math>td</math> and the linear unit output, <math>o_d</math></p>	4
7a	<p><b>Explain Naïve Bayes Classifier</b></p> <p>The naive Bayes classifier applies to learning tasks where each instance <math>x</math> is described by a conjunction of attribute values and where the target function <math>f(x)</math> can take on any value from some finite set <math>V</math>.</p> <p>A set of training examples of the target function is provided, and a new instance is presented, described by the tuple of attribute values <math>(a_1, a_2, \dots, a_m)</math>.</p> <p>The learner is asked to predict the target value, or classification, for this new instance.</p>	8

The Bayesian approach to classifying the new instance is to assign the most probable target value, VMAP, given the attribute values (a1, a2.. .am) that describe the instance

The naive Bayes classifier is based on the assumption that the attribute values are conditionally independent given the target value. Means, the assumption is that given the target value of the instance, the probability of observing the conjunction (a1, a2.. .am), is just the product of the probabilities for the individual attributes:

Naive Bayes classifier:

Where, VNB denotes the target value output by the naive Bayes classifier

An Illustrative Example

Let us apply the naive Bayes classifier to a concept learning problem i.e., classifying days according to whether someone will play tennis.

The below table provides a set of 14 training examples of the target concept PlayTennis, where each day is described by the attributes Outlook, Temperature, Humidity, and Wind

Use the naive Bayes classifier and the training data from this table to classify the following novel instance:

< Outlook = sunny, Temperature = cool, Humidity = high, Wind = strong > · Our task is to predict the target value (yes or no) of the target concept PlayTennis for this new instance

The probabilities of the different target values can easily be estimated based on their frequencies over the 14 training examples

·  $P(\text{PlayTennis} = \text{yes}) = 9/14 = 0.64$

·  $P(\text{PlayTennis} = \text{no}) = 5/14 = 0.36$

Similarly, estimate the conditional probabilities. For example, those for Wind = strong

·  $P(\text{Wind} = \text{strong} \mid \text{PlayTennis} = \text{yes}) = 3/9 = 0.33$

·  $P(\text{Wind} = \text{strong} \mid \text{PlayTennis} = \text{no}) = 3/5 = 0.60$

Thus, the naive Bayes classifier assigns the target value PlayTennis = no to this new instance, based on the probability estimates learned from the training data.

By normalizing the above quantities to sum to one, calculate the conditional probability that the target value is no, given the observed attribute values

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
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

7b

**Explain Bayesian Belief Networks and conditional independence with an example.**

**BAYESIAN BELIEF NETWORKS**

The naive Bayes classifier makes significant use of the assumption that the values of the attributes  $a_1 \dots a_n$  are conditionally independent given the target value  $v$ .

This assumption dramatically reduces the complexity of learning the target function  $A$ . A Bayesian belief network describes the probability distribution governing a set of variables by specifying a set of conditional independence assumptions along with a set of conditional probabilities

Bayesian belief networks allow stating conditional independence assumptions that apply to subsets of the variables

**Representation**

A Bayesian belief network represents the joint probability distribution for a set of variables. Bayesian networks (BN) are represented by directed acyclic graphs.

The Bayesian network in above figure represents the joint probability distribution over the boolean variables *Storm*, *Lightning*, *Thunder*, *ForestFire*, *Campfire*, and *BusTourGroup*

A Bayesian network (BN) represents the joint probability distribution by specifying a

8



	<p>set of <i>conditional independence assumptions</i></p> <p>BN represented by a directed acyclic graph, together with sets of local conditional probabilities</p> <p>Each variable in the joint space is represented by a node in the Bayesian network</p> <p>The network arcs represent the assertion that the variable is conditionally independent of its non-descendants in the network given its immediate predecessors in the network.</p> <p>A <b>conditional probability table (CPT)</b> is given for each variable, describing the probability distribution for that variable given the values of its immediate predecessors</p> <p>The joint probability for any desired assignment of values (<math>y_1, \dots, y_n</math>) to the tuple of network variables (<math>Y_1 \dots Y_m</math>) can be computed by the formula  Where, <math>Parents(Y_i)</math> denotes the set of immediate predecessors of <math>Y_i</math> in the network. <b>Example:</b></p> <p>Consider the node <i>Campfire</i>. The network nodes and arcs represent the assertion that <i>Campfire</i> is conditionally independent of its non-descendants <i>Lightning</i> and <i>Thunder</i>, given its immediate parents <i>Storm</i> and <i>BusTourGroup</i>.</p> <p>This means that once we know the value of the variables <i>Storm</i> and <i>BusTourGroup</i>, the variables <i>Lightning</i> and <i>Thunder</i> provide no additional information about <i>Campfire</i>. The conditional probability table associated with the variable <i>Campfire</i>. The assertion is</p> <p><math>P(\text{Campfire} = \text{True} \mid \text{Storm} = \text{True}, \text{BusTourGroup} = \text{True}) = 0.4</math></p>	
7c	<p><b>Let us assume dangerous fires are rare (1%) but smoke is fairly common (10%) due to the barbecues, and 90% of dangerous fires make smokes. Find the probability of dangerous fire when there is smoke.</b></p> <p>Let us say <math>P(\text{Fire})</math> means how often there is fire, and <math>P(\text{Smoke})</math> means how often we see smoke, then:</p> <p><math>P(\text{Fire} \mid \text{Smoke})</math> means how often there is fire when we can see smoke</p> <p><math>P(\text{Smoke} \mid \text{Fire})</math> means how often we can see smoke when there fire</p> <p>Given,</p> <p>dangerous fires are rare (1%) i.e. <math>P(\text{Fire})</math></p> <p>but smoke is fairly common (10%) due to barbecues i.e. <math>P(\text{Smoke})</math></p> <p>and 90% of dangerous fires make smoke i.e. <math>P(\text{Smoke} \mid \text{Fire})</math></p> <p>Probability of dangerous fire when there is smoke:</p> <p><math>P(\text{Fire} \mid \text{Smoke}) = P(\text{Fire}) \times P(\text{Smoke} \mid \text{Fire}) / P(\text{Smoke})</math></p> <p><math>= 1\% \times 90\% / 10\%</math></p> <p><math>= 9\%</math></p>	4

### Discuss minimum description length principle in brief.

The Minimum Description Length principle is motivated by interpreting the definition of  $h_{MAP}$  in the light of basic concepts from information theory. Consider again the now familiar definition of  $h_{MAP}$ .

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(D|h)P(h)$$

which can be equivalently expressed in terms of maximizing the  $\log_2$

$$h_{MAP} = \operatorname{argmax}_{h \in H} \log_2 P(D|h) + \log_2 P(h)$$

or alternatively, minimizing the negative of this quantity

$$h_{MAP} = \operatorname{argmin}_{h \in H} -\log_2 P(D|h) - \log_2 P(h) \quad (6.16)$$

Somewhat surprisingly, Equation (6.16) can be interpreted as a statement that short hypotheses are preferred, assuming a particular representation scheme for encoding hypotheses and data. To explain this, let us introduce a basic result from information theory: Consider the problem of designing a code to transmit messages drawn at random, where the probability of encountering message  $i$  is  $p_i$ . We are interested here in the most compact code; that is, we are interested in the code that minimizes the expected number of bits we must transmit in order to encode a message drawn at random. Clearly, to minimize the expected code length we should assign shorter codes to messages that are more probable. Shannon and Weaver (1949) showed that the optimal code (i.e., the code that minimizes the expected message length) assigns  $-\log_2 p_i$  bits<sup>†</sup> to encode message  $i$ . We will refer to the number of bits required to encode message  $i$  using code  $C$  as the *description length of message  $i$  with respect to  $C$* , which we denote by  $L_C(i)$ .

Let us interpret Equation (6.16) in light of the above result from coding theory.

- $-\log_2 P(h)$  is the description length of  $h$  under the optimal encoding for the hypothesis space  $H$ . In other words, this is the size of the description of hypothesis  $h$  using this optimal representation. In our notation,  $L_{C_H}(h) = -\log_2 P(h)$ , where  $C_H$  is the optimal code for hypothesis space  $H$ .
- $-\log_2 P(D|h)$  is the description length of the training data  $D$  given hypothesis  $h$ , under its optimal encoding. In our notation,  $L_{C_{D|h}}(D|h) = -\log_2 P(D|h)$ , where  $C_{D|h}$  is the optimal code for describing data  $D$  assuming that both the sender and receiver know the hypothesis  $h$ .
- Therefore we can rewrite Equation (6.16) to show that  $h_{MAP}$  is the hypothesis  $h$  that minimizes the sum given by the description length of the hypothesis plus the description length of the data given the hypothesis.

$$h_{MAP} = \operatorname{argmin}_h L_{C_H}(h) + L_{C_{D|h}}(D|h)$$

where  $C_H$  and  $C_{D|h}$  are the optimal encodings for  $H$  and for  $D$  given  $h$ , respectively.

The Minimum Description Length (MDL) principle recommends choosing the hypothesis that minimizes the sum of these two description lengths. Of course to apply this principle in practice we must choose specific encodings or representations appropriate for the given learning task. Assuming we use the codes  $C_1$  and  $C_2$  to represent the hypothesis and the data given the hypothesis, we can state the MDL principle as

**Minimum Description Length principle:** Choose  $h_{MDL}$  where

$$h_{MDL} = \operatorname{argmin}_{h \in H} L_{C_1}(h) + L_{C_2}(D|h) \quad (6.17)$$

	<p>The above analysis shows that if we choose <math>C_1</math> to be the optimal encoding of hypotheses <math>C_H</math>, and if we choose <math>C_2</math> to be the optimal encoding <math>C_{D h}</math>, then <math>h_{MDL} = h_{MAP}</math>.</p>	
8b	<p><b>Explain brute force MAP learning algorithm.</b></p> <p>In many learning scenarios, the learner considers some set <math>\mathcal{H}</math> of candidate hypotheses and is interested in finding the most probable hypothesis <math>h \in \mathcal{H}</math> given the observed data <math>D</math>. Any such maximally probable hypothesis is called a maximum a posteriori (MAP) hypothesis.</p> <p>Bayes theorem to calculate the posterior probability of each candidate hypothesis is <math>h_{MAP}</math> is a MAP hypothesis provided</p> <ul style="list-style-type: none"> <li>· <math>P(D)</math> can be dropped, because it is a constant independent of <math>h</math></li> </ul> <p>Maximum Likelihood (ML) Hypothesis</p> <ul style="list-style-type: none"> <li>· In some cases, it is assumed that every hypothesis in <math>\mathcal{H}</math> is equally probable a priori (<math>P(h_i) = P(h_j)</math> for all <math>h_i</math> and <math>h_j</math> in <math>\mathcal{H}</math>).</li> <li>· In this case the below equation can be simplified and need only consider the term <math>P(D h)</math> to find the most probable hypothesis.</li> </ul> <p><math>P(D h)</math> is often called the likelihood of the data <math>D</math> given <math>h</math>, and any hypothesis that maximizes <math>P(D h)</math> is called a maximum likelihood (ML) hypothesis</p> <p>Example:</p> <p>Consider a medical diagnosis problem in which there are two alternative hypotheses: (1) that the patient has particular form of cancer, and (2) that the patient does not. The available data is from a particular laboratory test with two possible outcomes: + (positive) and - (negative).</p> <p>We have prior knowledge that over the entire population of people only .008 have this disease. Furthermore, the lab test is only an imperfect indicator of the disease.</p> <p>The test returns a correct positive result in only 98% of the cases in which the disease is actually present and a correct negative result in only 97% of the cases in which the disease is not present. In other cases, the test returns the opposite result.</p> <p>The above situation can be summarized by the following probabilities:</p> <p>Suppose a new patient is observed for whom the lab test returns a positive (+) result. Should we diagnose the patient as having cancer or not?</p> <p>The exact posterior probabilities can also be determined by normalizing the above quantities so that they sum to 1</p>	8
8c	<p><b>Explain EM algorithm.</b></p>	4

**Step 1:** Calculate the expected value  $E[z_{ij}]$  of each hidden variable  $z_{ij}$ , assuming the current hypothesis  $h = \langle \mu_1, \mu_2 \rangle$  holds.

**Step 2:** Calculate a new maximum likelihood hypothesis  $h' = \langle \mu'_1, \mu'_2 \rangle$ , assuming the value taken on by each hidden variable  $z_{ij}$  is its expected value  $E[z_{ij}]$  calculated in Step 1. Then replace the hypothesis  $h = \langle \mu_1, \mu_2 \rangle$  by the new hypothesis  $h' = \langle \mu'_1, \mu'_2 \rangle$  and iterate.

Let us examine how both of these steps can be implemented in practice. Step 1 must calculate the expected value of each  $z_{ij}$ . This  $E[z_{ij}]$  is just the probability that instance  $x_i$  was generated by the  $j$ th Normal distribution

$$\begin{aligned} E[z_{ij}] &= \frac{p(x = x_i | \mu = \mu_j)}{\sum_{n=1}^2 p(x = x_i | \mu = \mu_n)} \\ &= \frac{e^{-\frac{1}{2\sigma^2}(x_i - \mu_j)^2}}{\sum_{n=1}^2 e^{-\frac{1}{2\sigma^2}(x_i - \mu_n)^2}} \end{aligned}$$

Thus the first step is implemented by substituting the current values  $\langle \mu_1, \mu_2 \rangle$  and the observed  $x_i$  into the above expression.

In the second step we use the  $E[z_{ij}]$  calculated during Step 1 to derive a new maximum likelihood hypothesis  $h' = \langle \mu'_1, \mu'_2 \rangle$ . maximum likelihood hypothesis in this case is given by

$$\mu_j \leftarrow \frac{\sum_{i=1}^m E[z_{ij}] x_i}{\sum_{i=1}^m E[z_{ij}]}$$

9a

Explain k- Nearest Neighbor learning algorithm  
The k- Nearest Neighbor algorithm for approximation a discrete-valued target function is given below:

Training algorithm:

- For each training example  $\langle x, f(x) \rangle$ , add the example to the list *training\_examples*

Classification algorithm:

- Given a query instance  $x_q$  to be classified,
  - Let  $x_1 \dots x_k$  denote the  $k$  instances from *training\_examples* that are nearest to  $x_q$
  - Return

$$\hat{f}(x_q) \leftarrow \operatorname{argmax}_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$$

where  $\delta(a, b) = 1$  if  $a = b$  and where  $\delta(a, b) = 0$  otherwise.

The K- Nearest Neighbor algorithm for approximation a real-valued target function is given below  $f : \mathfrak{R}^n \rightarrow \mathfrak{R}$

8

Training algorithm:

- For each training example  $(x, f(x))$ , add the example to the list *training\_examples*

Classification algorithm:

- Given a query instance  $x_q$  to be classified,
  - Let  $x_1 \dots x_k$  denote the  $k$  instances from *training\_examples* that are nearest to  $x_q$
  - Return

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$$

9b

- Consider locally weighted regression in which the target function  $f$  is approximated near  $x_q$  using a linear function of the form

$$\hat{f}(x) = w_0 + w_1 a_1(x) + \dots + w_n a_n(x)$$

Where,  $a_i(x)$  denotes the value of the  $i^{\text{th}}$  attribute of the instance  $x$

- Derived methods are used to choose weights that minimize the squared error summed over the set  $D$  of training examples using gradient descent

$$E \equiv \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x))^2$$

Which led us to the gradient descent training rule

$$\Delta w_j = \eta \sum_{x \in D} (f(x) - \hat{f}(x)) a_j(x)$$

Where,  $\eta$  is a constant learning rate

- Need to modify this procedure to derive a local approximation rather than a global one. The simple way is to redefine the error criterion  $E$  to emphasize fitting the local training examples. Three possible criteria are given below.
- Minimize the squared error over just the  $k$  nearest neighbors:

$$E_1(x_q) \equiv \frac{1}{2} \sum_{x \in k \text{ nearest nbrs of } x_q} (f(x) - \hat{f}(x))^2 \quad \text{equ(1)}$$

- Minimize the squared error over the entire set  $D$  of training examples, while weighting the error of each training example by some decreasing function  $K$  of its distance from  $x_q$  :

$$E_2(x_q) \equiv \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x))^2 K(d(x_q, x)) \quad \text{equ(2)}$$

- Combine 1 and 2:

$$E_3(x_q) \equiv \frac{1}{2} \sum_{x \in k \text{ nearest nbrs of } x_q} (f(x) - \hat{f}(x))^2 K(d(x_q, x)) \quad \text{equ(3)}$$

8

If we choose criterion three and re-derive the gradient descent rule, we obtain the following training rule

$$\Delta w_j = \eta \sum_{x \in k \text{ nearest nbrs of } x_q} K(d(x_q, x)) (f(x) - \hat{f}(x)) a_j(x)$$

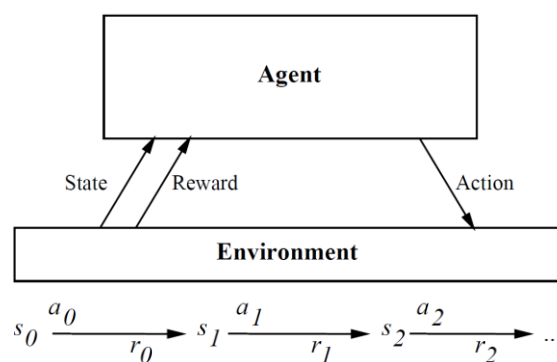
The differences between this new rule and the rule given by Equation (3) are that the contribution of instance  $x$  to the weight update is now multiplied by the distance penalty  $K(d(x_q, x))$ , and that the error is summed over only the  $k$  nearest training examples.

9c

What is Reinforcement learning?

Reinforcement learning is an area of Machine Learning. It is about taking suitable action to maximize reward in a particular situation. It is employed by various software and machines to find the best possible behaviour or path it should take in a specific situation.

Reinforcement learning differs from the supervised learning in a way that in supervised learning the training data has the answer key with it so the model is trained with the correct answer itself whereas in reinforcement learning, there is no answer but the reinforcement agent decides what to do to perform the given task. In the absence of training dataset, it is bound to learn from its experience.



Goal: Learn to choose actions that maximize

$$r_0 + \gamma r_1 + \gamma^2 r_2 + \dots, \text{ where } 0 \leq \gamma < 1$$

4

10a

Distinguish Eager learning vs Lazy learning algorithms.

Eager learning method constructs general explicit description of the target function based on the provided training examples. Eg. Decision Tree based Learning  
 Lazy learning method simply store the data and generalizes beyond these post point until an explicit request is made. Eg. KNN algorithm

4

10  
b

Write short notes on Q-learning

**The Q Function**

8

The value of Evaluation function  $Q(s, a)$  is the reward received immediately upon executing action  $a$  from state  $s$ , plus the value (discounted by  $\gamma$ ) of following the optimal policy thereafter

$$Q(s, a) \equiv r(s, a) + \gamma V^*(\delta(s, a)) \quad \text{equ (4)}$$

Rewrite Equation (3) in terms of  $Q(s, a)$  as

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q(s, a) \quad \text{equ (5)}$$

Equation (5) makes clear, it need only consider each available action  $a$  in its current state  $s$  and choose the action that maximizes  $Q(s, a)$ .

### An Algorithm for Learning $Q$

- Learning the  $Q$  function corresponds to learning the **optimal policy**.
- The key problem is finding a reliable way to estimate training values for  $Q$ , given only a sequence of immediate rewards  $r$  spread out over time. This can be accomplished through *iterative approximation*

$$V^*(s) = \max_{a'} Q(s, a')$$

Rewriting Equation

$$Q(s, a) = r(s, a) + \gamma \max_{a'} Q(\delta(s, a), a')$$

---

#### $Q$ learning algorithm

For each  $s, a$  initialize the table entry  $\hat{Q}(s, a)$  to zero.

Observe the current state  $s$

Do forever:

- Select an action  $a$  and execute it
- Receive immediate reward  $r$
- Observe the new state  $s'$
- Update the table entry for  $\hat{Q}(s, a)$  as follows:

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a'} \hat{Q}(s', a')$$

- $s \leftarrow s'$
- 

10c

Discuss about Radial basis function in detail

- One approach to function approximation that is closely related to distance-weighted regression and also to artificial neural networks is learning with radial basis functions

8

- In this approach, the learned hypothesis is a function of the form

$$\hat{f}(x) = w_0 + \sum_{u=1}^k w_u K_u(d(x_u, x)) \quad \text{equ (1)}$$

- Where, each  $x_u$  is an instance from  $X$  and where the kernel function  $K_u(d(x_u, x))$  is defined so that it decreases as the distance  $d(x_u, x)$  increases.
- Here  $k$  is a user provided constant that specifies the number of kernel functions to be included.
- $\hat{f}$  is a global approximation to  $f(x)$ , the contribution from each of the  $K_u(d(x_u, x))$  terms is localized to a region nearby the point  $x_u$ .

Choose each function  $K_u(d(x_u, x))$  to be a Gaussian function centred at the point  $x_u$  with some variance  $\sigma_u^2$

$$K_u(d(x_u, x)) = e^{-\frac{1}{2\sigma_u^2} d^2(x_u, x)}$$

- The functional form of equ(1) can approximate any function with arbitrarily small error, provided a sufficiently large number  $k$  of such Gaussian kernels and provided the width  $\sigma^2$  of each kernel can be separately specified
- The function given by equ(1) can be viewed as describing a two layer network where the first layer of units computes the values of the various  $K_u(d(x_u, x))$  and where the second layer computes a linear combination of these first-layer unit values

**Example: Radial basis function (RBF) network**

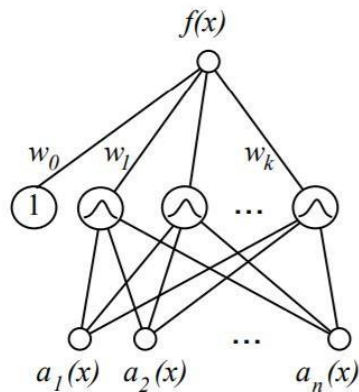
Given a set of training examples of the target function, RBF networks are typically trained in a two-stage process.

1. First, the number  $k$  of hidden units is determined and each hidden unit  $u$  is defined by choosing the values of  $x_u$  and  $\sigma_u^2$  that define its kernel function  $K_u(d(x_u, x))$
2. Second, the weights  $w$ , are trained to maximize the fit of the network to the training data, using the global error criterion given by

$$E \equiv \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x))^2$$

Because the kernel functions are held fixed during this second stage, the linear weight values  $w$ , can be trained very efficiently





Several alternative methods have been proposed for choosing an appropriate number of hidden units or, equivalently, kernel functions.

- One approach is to allocate a Gaussian kernel function for each training example  $(x_i, f(x_i))$ , centring this Gaussian at the point  $x_i$ . Each of these kernels may be assigned the same width  $\sigma^2$ . Given this approach, the RBF network learns a global approximation to the target function in which each training example  $(x_i, f(x_i))$  can influence the value of  $f$  only in the neighbourhood of  $x_i$ .
- A second approach is to choose a set of kernel functions that is smaller than the number of training examples. This approach can be much more efficient than the first approach, especially when the number of training examples is large.