

D Define machine learning? Discuss with examples some applications of machine learning.

Machine learning is a branch of artificial intelligence (AI) and computer science which focuses on the use of data and algorithms to imitate the way that humans learn, gradually improving its accuracy.

Types:

1) Supervised learning:

It is when the model is getting trained on a labelled dataset. A labelled dataset is one that has both input and output parameters.

2 categories → Classification, Regression

Eg: Teaching a baby what a dog and a cat is. by telling the different characteristics of dog and cat. Here the one who is teaching is acted as supervisor and the baby acted as algorithm that had to learn

A basket full of different kinds of fruits.

2) Unsupervised learning:

It's a learning algorithm which is trained using unlabeled data. It does not need supervision, it finds the patterns from the i/p data by its own
2 types of problems: → Clustering & Association

Eg: Suppose a baby knows and identifies her family dog. Few weeks later a family friend brings along a dog. Baby has not seen this dog earlier. But it recognizes many features (Ears, eyes, 4 legs) are like her pet dog. She identifies the new animal as a dog.

3) Semi-supervised learning:

It is an approach to machine learning that combines a small amount of labeled data with a large amount of unlabeled data during training.

Eg: Web content classification is very critical & impossible to label each page on the internet because it needs more human intervention. This problem can be reduced through semi-supervised learning algorithms.

It is a feedback based ML technique in which an agent learns automatically using feedbacks without any labeled data. Agent is bound to learn by its experience only.

For each good action the agent gets positive feedback, and for each bad action, the agent gets negative feedback or penalty.

Eg: Suppose a master is training a dog to get the stick

Each time the dog gets a stick successfully, master offered him a feast. Eventually, the dog understands the pattern that whenever the master throws a stick it should get it as early as it can to gain a reward from a master.

Applications:

1) Speech Recognition → It is a process of converting voice instructions into text.

Eg: Google assistant, Siri, Cortana & Alexa

2) Traffic prediction → It predicts the traffic conditions such as whether traffic is cleared, slow moving or heavily congested

It shows us the correct path with the shortest route and predicts the traffic conditions.

3) Product recommendations → Google understands the user interest using various machine learning algorithms and suggests the product as per customer interest

4) Self driving cars → Tesla the most popular car manufacturing company is working on self driving cars. It is using unsupervised learning method to train the car models to detect people & objects while driving.

5) Medical diagnosis → In medical science, machine learning is used for diseases diagnoses.

It helps in finding brain tumors and other brain related diseases easily.

2) What do you mean by well posed learning problem.
Explain the important features that are required to define well posed learning problem.

An agent solves a problem or task T , performance P and gain some experience E .

If P is measured at T it can improve E .
(knowing by experience).
A computer program is said to learn from experience E with respect some class of tasks T and performance measure P , if its performance at tasks in T as measured by P improves with experience E ;

Three features that are required to define well posed learning problem.

- 1) The class of tasks
- 2) The measure of performance to be improved
- 3) The source of experience

Well defined learning problem

A checkers learning problem

Task T : playing checkers

Performance measure P : percent of games won against opponents

Training experience E : playing practice games against itself

A handwriting recognition learning problem

Task T : recognizing and classifying handwritten words within images

Performance measure P : percent of words correctly classified.

Training experience E : a database of handwritten words with given classification

A robot driving learning problem

Task T : driving on public four lane highways using vision sensors

Performance measure P : Average distance travelled before an error

Training experience E : a sequence of images

and steering commands recorded while observing a human driver.

3. Describe briefly the steps involved in designing a learning system.

To get a successful learning system it should be designed for a proper design several steps should be followed why? → for perfect & efficient system.

Steps for designing learning system

- 1) Choosing Training experience
- 2) Choosing target function
- 3) Choosing representation of target function
- 4) Choosing function approximation
- 5) Final design.

STEP 1

1) Choosing 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 & failure of the learner.

In choosing a training experience, 3 attributes are taken

- 1) Type of feedback (whether it is direct or indirect)
 - 2) Degree of difficulty (how hard is it to learn)
 - 3) Distribution of examples
- 1) Whether the training experience provides direct or indirect feedback regarding the choices made by performance system.

Task → P good, → good E

In order to perform good in task, we have good E

By training we can get good experience.

Direct Indirect

1) Direct Feedback:

In checkers game, when we make a move, immediately computer will tell us whether you made a correct move or wrong move.

(Immediately after 1 move)

Learning driving:

1) While learning driving trainer will tell everything accelerator, breaker.

2) Indirect feedback:

While learning driving trainer won't be physically beside you, some recorded sessions will be there. It will not spontaneously simultaneously telling you. (Feedbacks are given in intervals)

2) Degree to which learners will control the sequence of training.

Ex: learning driving

with trainer's help	with trainer's partial help	Completely on your own.
---------------------	-----------------------------	-------------------------

3) Distribution of examples:

 How well it represents the distribution of examples over which the performance of final system is measured

→ more possible combination testing, but good?

→ more situations, more examples

→ learning over distributed range of examples

Ex: learning driving  Single, 2 way, 4 way, hill, curve road, road, etc.

→ step by step, without T

STEP 2: Choosing the target function

To determine what type of knowledge will be learned and how it will be used by the performance system

Ex: Checkers game.

→ travel only in forward direction

→ Only one move per chance & only in diagonal direction

→ Jump over opponent.

→ While moving diagonally set of all possible moves is called legal moves.

→ We need to select one move & that move is called target move.

Target function $\Rightarrow v(b)$

Board state $\Rightarrow b$ (Initial, intermediary or final state).

Legal moves set $\Rightarrow B$

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

Board state b in B is final board state, that is won, then $v(b) = 100$

► If b is a final board state, that is lost, then $v(b) = -100$

- 2) If b is a final board state that is lost, then $V(b) = -100$
 3) If b is a final board state that is drawn, then $V(b) = 0$
 4) If b is 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.

STEP 3: Choosing the representation for target function.

For any board state we calculate function c as linear combination of following board features i.e $c(b)$

Features:

- $x_1 \rightarrow$ No. of black pieces on the board
- $x_2 \rightarrow$ No. of red pieces on the board
- $x_3 \rightarrow$ No. of black kings on the board
- $x_4 \rightarrow$ No. of red kings on the board
- $x_5 \rightarrow$ No. of black pieces threatened by red
- $x_6 \rightarrow$ No. of black pieces beaten by red
- No. of red pieces threatened by black

$$V^*(b) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + w_6 x_6$$

(Representation of $c(b)$)

w_0 to $w_6 \rightarrow$ Numerical coefficient or weight of each feature
 $w_0 \rightarrow$ additive constant.

STEP 4: Choosing function approximation algorithm

To learn a target function we need set of training examples.

Describes a particular board state (b) & training value $V_{\text{train}}(b)$.

Training example representation is ordered pair = $(b, V_{\text{train}}(b))$

Example: Black won the game

i.e $x_2=0$ which means no red.

$$V_{\text{train}}(b) = +100.$$

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

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

1) Estimating training values

In every step we consider successor depending on the next step of opponent

$V_{train}(b) \leftarrow V^*(\text{successor}(b))$

\downarrow
represents the next board state

(estimating that this move will help/destroy opponent)

\uparrow represents approximation

2) Adjusting the weights:

there are some algorithm to find weights of linear functions

here we are using LMS (least mean square).

LMS is used to minimise the error.

$$\text{Error} = (V_{train}(b) - \hat{V}(b))^2$$

$$E = \sum_{(b, V_{train}(b)) \in \text{Training examples}} (V_{train}(b) - \hat{V}(b))^2$$

LMS weight update rule:

For each training example $(b, V_{train}(b))$, use the current weight s to calculate $\hat{V}(b)$. 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 that moderates the size of the weight update.

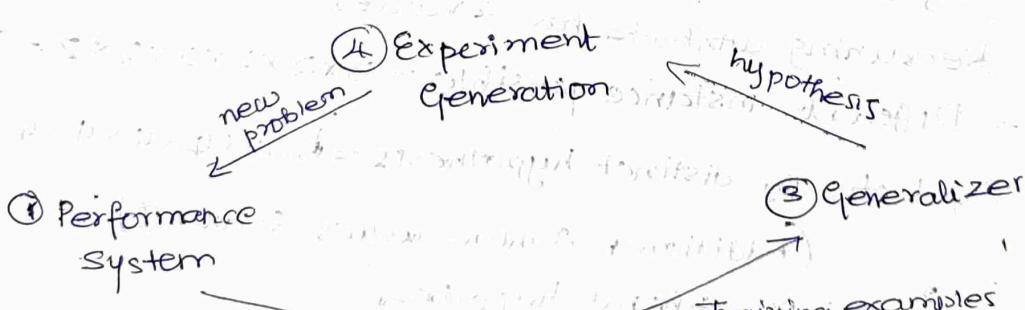
If Error = 0, no need to change weights

If Error is positive, each weight is increased in proportion

If Error is negative, each weight is decreased in proportion

STEP 5: Final design

It has 4 different modules.



Performance system:

It takes an instance of a new problem (new game)

as input and produces a trace of its solution as output.

⇒ The critic:

It takes as input the history or trace of the game and produces as output a set of strong examples of the target function.

⇒ The Generalizer:

It takes the training examples as input and produces an output hypothesis; that is its estimate of the target function.

⇒ The Experiment Generator:

It takes the current hypothesis as input and new problem as output for the performance system to explore.

Q:

Discuss the concept of learning as the task of searching with respect to generate to specific ordering of hypothesis.

Main goal of this search is to find the hypothesis that best fits the training examples.

Ex) Enjoy sport learning task.

→ 6 attributes
 rainy normal strong warm some
 Cold ↓ highest ↙ week cool ↓
 Cloudy ↓ ↓ ↓ ↓ ↓
 Stay , Air temperature, humidity, wind, water & forecast

Sky has 3 values → Rainy, cloudy & Sunny.

Remaining attribute has 2 values

∴ Different instance possible = $3 \times 2 \times 2 \times 2 \times 2 \times 2 = 96$

Syntactically distinct hypothesis = $5 \times 4 \times 4 \times 4 \times 4 \times 4 = 5120$

(Additional & more values ? and \emptyset)

? → most general hypothesis

\emptyset → most specific hypothesis

↳ \emptyset , Sunny, cloudy, Rainy, ?)

Semantically distinct hypothesis = $1 + (4 \times 3 \times 3 \times 3 \times 3 \times 3)$

= 973

(Null taken as common, \emptyset only 1 more value)

After finding all the syntactically & semantically distinct hypothesis.

We search the best match from all these

(ie much closer to our learning problem)

5. Illustrate find S algorithm over enjoy sport training instance given.

Find S Algorithm:

finding a maximally specific hypothesis

→ this algorithm considers only positive examples.

* Representation:

most specific hypothesis $\Rightarrow \emptyset$

most general hypothesis $\Rightarrow ?$

Algorithm:

Step 1: Initialise with most specific hypothesis (\emptyset)

$$h_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle \Rightarrow 5 \text{ attributes}$$

Step 2: For each the samples, for each attribute

if (value = hypothesis value)

 ⇒ Ignore

else

 Replace with most general hypothesis (?)

Step 3: Output the hypothesis h

Example	Sky	Airtemp	Humidity	Wind	Water	Forecast	Enjoy sport
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

1. Initialise h to the most specific hypothesis in H

$$h_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$$

$$x_1 = \langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle$$

$$h_1 = \langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle$$

$$x_2 = \langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Warm}, \text{Same} \rangle$$

$$h_2 = \langle \text{Sunny}, \text{Warm}, ?, \text{Strong}, \text{Warm}, \text{Same} \rangle$$

$$x_3 = \langle \text{Rainy}, \text{Cold}, \text{High}, \text{Strong}, \text{Warm}, \text{Change} \rangle$$

No, so ignore.

$$h_3 = \langle \text{Sunny}, \text{Warm}, ?, \text{Strong}, \text{Warm}, \text{Same} \rangle$$

$h_1 = \langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Cool}, \text{Change} \rangle$

$h_4 = \langle \text{Sunny}, \text{Warm}, ?, \text{Strong}, ?, ? \rangle \Rightarrow \text{Output}$

6. Define consistent hypothesis and version space.

With example explain version space & its representation

Consistent hypothesis:

An hypothesis h is consistent with a set of training

examples D iff $h(x) = c(x)$ for each example in D

$\text{Consistent}(h, D) \equiv \forall \langle x, c(x) \rangle \in D \quad h(x) = c(x)$

Version space VS:

Subset of hypotheses h consistent with the training examples (D)

$$VS_{h,D} = \{ h \in H \mid \text{consistent}(h, D) \}$$

h = hypothesis \exists $\forall x \in D \quad h(x) = c(x)$ \rightarrow consistent

D = training example $h(x) = c(x)$

Ex: Enjoy sport

Initialise G_0 and S_0

$$G_0 = \{ \langle ?, ?, ?, ?, ?, ?, ? \rangle \}$$

$$S_0 = \{ \langle \$, \$, \$, \$, \$, \$, \$ \rangle \}$$

For $x_1 = \langle \text{Su}, \text{Wa}, \text{H}, \text{St}, \text{Wa}, \text{Sa} \rangle$ Generalise S_0 to cover x_1

$$S_1 = \{ \langle \text{Su}, \text{Wa}, ?, ?, ?, ?, ? \rangle \}$$

$$G_1 = \{ \langle ?, ?, ?, ?, ?, ?, ? \rangle \}$$

For $x_2 = \langle ?, ?, ?, ?, ?, ?, ? \rangle$ replace attribute value by $?$ to make it consistent with x_2

$$S_2 = \{ \langle \text{Su}, \text{Wa}, ?, \text{St}, \text{Wa}, \text{Sa} \rangle \}$$

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

For $x_3 = \langle \text{Ra}, \dots \rangle$ -ve

As x_3 -ve training example, general boundary incorrectly predicts x_3 is +ve

So G_2 needs to be specialised to make it classify as -ve correctly

$$S_2 = S_3 = \{ \langle \text{Su}, \text{Wa}, - \rangle \}$$

$$G_3 = \{ \langle \text{Su}, ?, ?, ?, ?, ?, ? \rangle, \langle ?, \text{Wa}, ?, ?, ?, ?, ? \rangle \}$$

For $x_4 = \langle ?, ?, ?, ?, ?, ?, ? \rangle$ +ve, fourth training example need to be generalised to be consistent

Also G_3 needs to be updated to make consistent with the x_4

$$G_4 = \{ \langle \text{Su}, ?, ?, ?, ?, ?, ? \rangle, \langle ?, \text{Wa}, ?, ?, ?, ?, ? \rangle \}$$

Q) Describe & Explain the candidate elimination algorithm with illustration

- Uses the concept of version space
- Consider both +ve & -ve value (Yes and No)
- both specific and general hypothesis
- For specific +ve samples, move from specific to general
- For -ve samples move from general to specific

$$S = \{\phi, \phi, \phi, \phi, \phi\} + V$$

$$G_1 = \{?, P, ?, ?, ?, ?\} - V$$

Algorithm:

Step 1: Initialise both general & specific hypothesis.

$$S = \{\phi, \phi, \phi, \phi, \phi\}$$

$$G = \{?, P, ?, ?, ?, ?\}$$

Step 2: For each example,

if example is +ve → move from specific to general
make specific to general

else → move from general to specific.

Ex: Enjoy Sport.

$$S_0 = \{\phi, \phi, \phi, \phi, \phi, \phi\} \quad G_0 = \{?, ?, ?, ?, ?, ?\}$$

1) +ve

$$S_1 = \{\text{Sunny, Warm, Normal, Strong, Warm, Same}\}$$

$$G_1 = \{?, ?, P, ?, ?, ?\}$$

2) +ve (specific to general) Only change is not G_1 .

$$S_2 = \{\text{Sunny, warm, ?, strong, Warm, Same}\}$$

$$\begin{aligned} S_1 &= \text{Normal} \\ x_2 &= \text{High} \end{aligned}$$

$$G_2 = \{?, ?, ?, ?, ?, ?\}$$

keeps as its. Only change G_1

3) +ve (general to specific)

$$S_3 = \{\text{Sunny, Warm, ?, Strong, Warm, Same}\}$$

$$\begin{cases} x_1 = \text{Sunny} & x_2 = \text{Wom} & x_2 = \text{High} & \text{no need to write} \\ x_3 = \text{Rainy} & x_3 = \text{Old} & x_3 = \text{High} & \end{cases}$$

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

$$G_3 = \{\text{Sunny, ?, ?, ?, ?}\}$$

$$\langle ?, ?, ?, ?, \text{Same} \rangle$$

4) +ve $S_3 = \{ ? , ? , ? , ? , ? \}$ no companion

$S_4 = \{ \langle \text{Sunny}, \text{?}, \text{?}, \text{?}, \text{?}, \text{?} \rangle \}$

$G_4 = \{ \langle \text{Sunny}, ? ? ? ? ? \rangle, \langle ? \text{ warm} ? ? ? ? \rangle \}$

Don't include $\langle ? ? ? ? ? \rangle$
Because in S_4 there's
no some

S_4 and $G_4 \Rightarrow$ final hypothesis.

- Q. Consider the given below training example, which finds Malignant tumours from the MRI scans.

Example Shape Size Color Surface Thickness Target concept

	Shape	Size	Color	Surface	Thickness	Target concept
1	Circular	Large	Light	Smooth	Thick	Malignant
2	Circular	Large	Light	Irregular	Thick	Malignant
3	Oval	Large	Dark	Smooth	Thin	Benign
4	Oval	Large	Light	Irregular	Thick	Malignant
5	Circular	Small	Light	Smooth	Thick	Benign

Show the specific and general boundaries of the version space after applying candidate elimination algorithm (Note: Malignant is +ve & Benign -ve)

$$\Rightarrow S_0 = \{ \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \} \quad G_0 = \{ ?, ?, ?, ?, ? \}$$

+ve

$$S_1 = \{ \langle \text{Circular}, \text{large}, \text{light}, \text{smooth}, \text{thick} \rangle \}$$

$$G_1 = \{ \langle ?, ?, ?, ?, ? \rangle \}$$

2) +ve ($S_1 \rightarrow G_1$)

$$S_2 = \{ \langle \text{Circular}, \text{large}, \text{light}, ?, \text{thick} \rangle \}$$

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

3) -ve ($G_1 \rightarrow S_1$)

$$S_3 = \{ \langle \text{Circular}, \text{large}, \text{light}, ?, \text{thick} \rangle \}$$

$$G_3 = \{ \langle \text{Circular}, ?, ?, ?, ?, ? \rangle \}$$

(Check consistency till S_3)

4) +ve ($S_1 \rightarrow G_1$)

$$S_4 = \{ \langle ?, \text{large}, \text{light}, ?, \text{thick} \rangle \}$$

$$G_4 = \{ \langle ?, ?, \text{light}, ?, ?, ? \rangle \}$$

? -ve

$S_5 = \{ \leftarrow ? \text{ Large } | \text{ light } , ? \text{ thick } \rightarrow \}$

$G_{15} = \{ \leftarrow ? \text{ Large } | \text{ light } , ? \text{ light } \rightarrow ? \text{ Large } | ? \text{ Thick } \}$

NOTE: Final version space for the enjoy Sport concept learning problem and training examples.

$S_4 : \{ \leftarrow \text{sunny, warm, ?}, \text{Strong, ?, ?, ?} \}$

$S_4 \cap G_4 = \{ \leftarrow \text{sunny, warm, ?}, \text{Strong, ?, ?} \} \cap \{ \leftarrow \text{sunny, warm, ?}, \text{Strong, ?, ?} \} = \{ \leftarrow \text{sunny, warm, ?}, \text{Strong, ?, ?} \}$

$G_4 : \{ \leftarrow \text{sunny, ?, ?, ?, ?} \} \cap \{ \leftarrow \text{?, warm, ?, ?, ?, ?} \} = \{ \leftarrow \text{?, warm, ?, ?, ?, ?} \}$

If mismatch → replace with ? and copy remaining of S_4 instances with specific hypothesis of copy remaining of G_4 .

g. Describe ID3 algorithm. Calculate entropy and information gain for the following dataset.

Instance	Classification	A1	A2
1	+	T	T
2	-	T	F
3	-	T	F
4	+	F	F
5	-	F	T
6	-	-	-

ID3 standards for iterative Dichotomised 3. Used to generate a decision tree from a given dataset by employing a top-down greedy search to test each attribute at every node of the tree. The resulting tree is used to classify future samples.

Algorithm:

ID3 (Examples, Target-attribute, Attributes)
Examples are the training examples. Target-attribute is the attribute whose value is to be predicted by the tree. Attributes is a list of other attributes. Returns a decision tree that correctly classifies the given Examples.

④

- ▷ Create a root node for the tree
- 2) If all Examples are positive, return the single-node tree root, with label = +
- 3) If all Examples are -ve, return the single-node tree root, with label = -
- 4) If attributes is empty, return the single-node tree root = most common value of target attribute in Examples.
- 5) Otherwise Begin
 - i) A \leftarrow the attribute from attributes that best-classifies examples.
 - ii) The decision attribute for root \leftarrow A
 - iii) For each possible value v_i of A
 - a) Add a new tree branch below root corresponding to the test $A = v_i$
 - b) Let examples v_i , be the subset of examples that have value v_i for A
 - c) If Examples v_i is empty
 - Then below this add new branch add a leaf node with label = most common value of target attribute in Examples
 - Else below this new branch add the subtree ID3 (Examples v_i , Target-attribute, Attributes - {A})
- 6) End
- 7) Return root

Entropy \rightarrow measures the impurity of a collection of examples.
Information gain \rightarrow that measures how well a given attribute separates the training examples according to their target classification.

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

Where P_+ is the proportion of the examples in S
 P_- is the proportion of -ve examples in S

$$\text{Information gain} = \text{Gain}(S, A) = \text{Entropy}(S) - \sum_{S'} [P(S') \text{Entropy}]$$

$A \rightarrow \text{attribute}$

$\text{Gain}(S, A)$ of an attribute A , relative to a collection of examples S .

Problem solution:

$$\log_2 x = \frac{\log_{10} x}{\log_2}$$

Attribute : a₁

$$\text{Values } (a_1) = T, F \quad \begin{array}{c} T \rightarrow 3 \\ \diagdown \\ F \rightarrow 3 \end{array} \quad \begin{array}{c} S \rightarrow +1/3 \\ -1/3 \end{array}$$

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

$$\begin{aligned} &= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ &= -\frac{2}{3} \times (-0.5849) - \frac{1}{3} \times (-1.5849) \\ &= 0.9183 = \frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1 \end{aligned}$$

$$S_T \begin{cases} + \rightarrow 2 \\ - \rightarrow 1 \end{cases}$$

$$\begin{aligned} \text{Entropy } (S_T) &= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ &= -\frac{2}{3} \times (-0.58496) - \frac{1}{3} \times (-1.58496) \\ &= 0.9183 \end{aligned}$$

$$S_F \begin{cases} + \rightarrow 1 \\ - \rightarrow 2 \end{cases}$$

$$\begin{aligned} \text{Entropy } (S_F) &= -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \\ &= 0.9183 \end{aligned}$$

Gain of a₁:

$$\text{Gain } (S, a_1) = \text{Entropy } (S) - \sum_{v \in \{T, F\}} \frac{|S_v|}{|S|} \text{Entropy } (S_v)$$

$$\begin{aligned} \text{Gain } (S, a_1) &= \cancel{\frac{1}{3}} \log_2 \frac{1}{3} - \frac{3}{6} \text{Entropy } (S_T) - \frac{3}{6} \text{Entropy } (S_F) \\ &\quad \cancel{+\frac{2}{3} \log_2 \frac{2}{3}} \\ &= -\frac{1}{3} \times 0.9183 - \frac{3}{6} \times 0.9183 \\ &= 0.0817 \end{aligned}$$

Attribute: a₂

$$\text{Values } (a_2) = T, F \quad \begin{array}{c} T \rightarrow 4 \\ \diagdown \\ F \rightarrow 2 \end{array}$$

$$S \begin{cases} + \rightarrow 3 \\ - \rightarrow 3 \end{cases}$$

$$\begin{aligned} \text{Entropy } (S) &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\ &= 1 \end{aligned}$$

$$S_T \begin{cases} + \rightarrow 2 \\ - \rightarrow 2 \end{cases}$$

$$\text{Entropy } (S_T) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1.0$$

$$S_F \leftarrow \begin{array}{c} + \\ - \end{array} \rightarrow 1$$

$$\text{Entropy}(S_F) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1.0$$

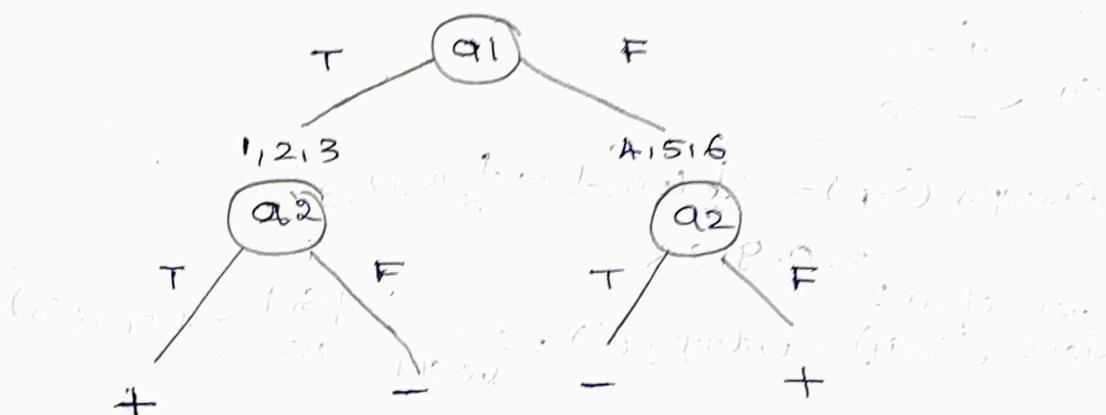
Gain of a_2 :

$$\text{Gain}(S, a_2) = \text{Entropy}(S) - \sum_{V \in T_F} \frac{|S_V|}{|S|} \text{Entropy}(S_V)$$

$$\begin{aligned} \text{Gain}(S, a_2) &= 1 - \frac{4}{6} \text{Entropy}(ST) - \frac{2}{6} \text{Entropy}(S_F) \\ &= 1 - \frac{4}{6} \times 1.0 - \frac{2}{6} \times 1.0 \\ &= \underline{\underline{0.0}} \end{aligned}$$

$\text{Gain}(S, a_1) = 0.0817 \rightarrow \text{Maximum Gain}$

$\text{Gain}(S, a_2) = 0.0$



1) Describe inductive bias in decision tree learning.

Inductive bias is the set of assumptions that, together with the training data, deductively justify the classifications assigned by the learner to future instances.

Inductive bias of ID3 consists of describing the basis by which ID3 chooses one consistent decision tree over all the possible ~~one~~ decision tree's

ID3 search strategy:

- > Selects in favour of shorter trees over longer ones
- > Selects trees that place the attribute with highest information gain closest to the root

Types of inductive bias

> Restrictive bias \rightarrow based on conditions

> Preference bias \rightarrow based on priorities

ID₃ \Rightarrow Preference

Version space and candidate-elimination \Rightarrow Restrictive

Why short hypothesis? According to Ockham's Razor, prefer simplest hypothesis that fits the data.

Difference b/w the types of inductive bias exhibited by ID₃ and by the CANDIDATE-ELIMINATION algorithm.

ID₃:

- The inductive bias of ID₃ is a preference bias for certain hypotheses over others, with no hard restriction on the hypotheses that can be eventually enumerated. This form of bias is called a preference bias or a search bias.
- ID₃ searches a complete hypothesis space. set of possible DT's
- It searches incompletely through this space from simple to complex hypotheses until its termination condition is met.

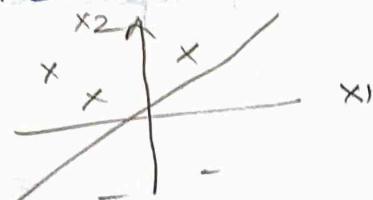
CANDIDATE-ELIMINATION

- The bias of the CE algorithm is in the form of a categorical restrictions on the set of hypotheses considered.
- The version space CE searches an incomplete hypothesis space.
- It searches this space completely, finding every hypothesis consistent with the training data.

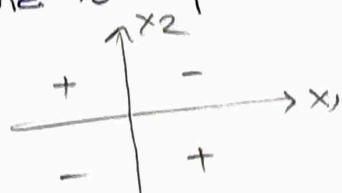
(ii) What is linearly separable problem. Design a network of perceptron to implement X AND Y.

Linearly inseparable means if there's a function that is not linearly separable it by definition means that there's no line such that points of one type are on one side and the other type on the other side.

i.e. there exists no line to separate the points.



linearly separable

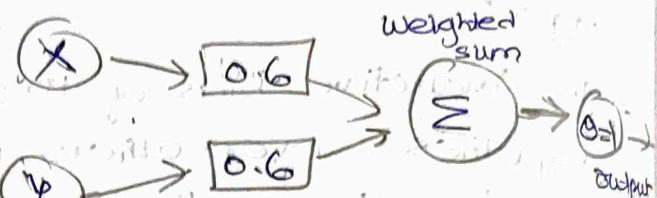


linearly non separable

Decision boundaries could have linear or non-linear form. If it is not possible to have straight line separating the regions then, the patterns are said to be linearly Non-Separable. Such patterns then may be separated by decision boundaries that are non-linear.

logical AND operation

$X \oplus Y$	
0	0
0	1
1	0
1	1



$$\text{If } X=0 \text{ & } Y=0 \rightarrow 0 \times 0.6 + 0 \times 0.6 = 0$$

This is not greater than the threshold of 1 so output

$$\text{If } X=0 \text{ & } Y=1 \rightarrow 0 \times 0.6 + 1 \times 0.6 = 0.6$$

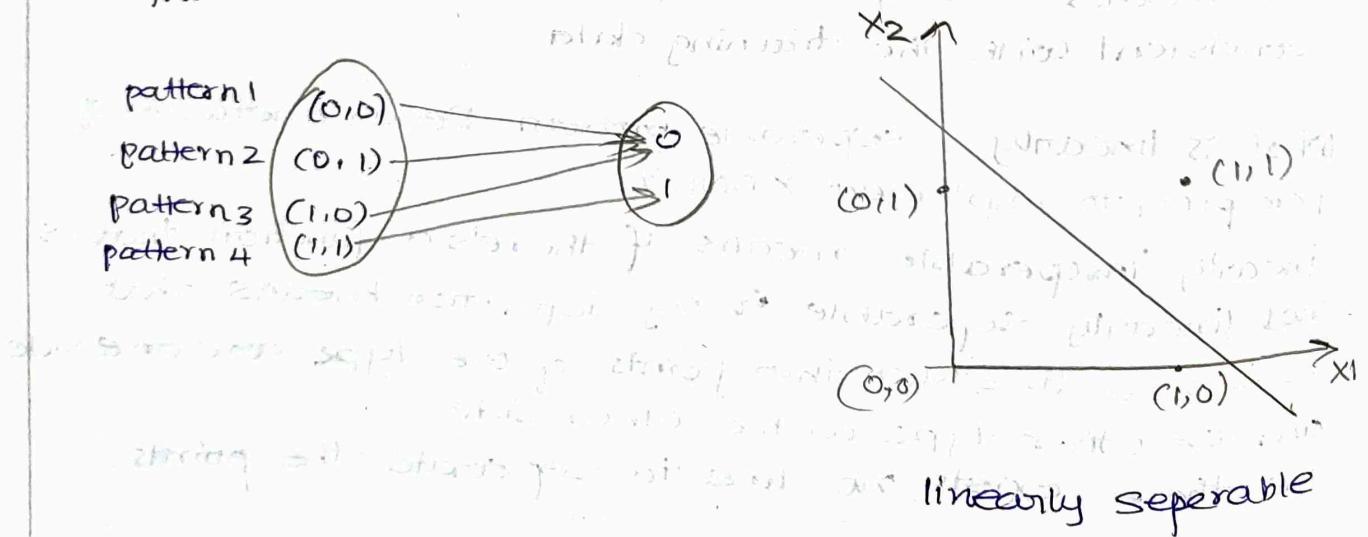
This is not greater than the threshold of 1 so output

$$\text{If } X=1 \text{ & } Y=0 \rightarrow 1 \times 0.6 + 0 \times 0.6 = 0.6$$

This is not greater than the threshold of 1 so output

$$\text{If } X=1 \text{ & } Y=1 \rightarrow 1 \times 0.6 + 1 \times 0.6 = 1.2$$

This exceeds the threshold so output = 1



12) What is perceptron? Discuss perceptron training rule.

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.

Perceptron training rule:

The learning problem is to determine a weight vector that causes the perceptron to produce the correct +1 or -1 output for each of the ~~following~~ training examples.

To learn an acceptable weight vector:

1) Begin with random weights then iteratively apply the perceptron to each training example modifying the perceptron weights whenever it misclassifies an example.

2) This process is repeated iterating through the training examples as many times as needed until the perceptron classifies all training examples correctly.

3) Weights are modified at each step according to the perceptron training rule which revises the weight w_i associated with input x_i according to the rule.

$$w_i \leftarrow w_i + \Delta w_i$$

$$\text{where } \Delta w_i = \eta(t - \theta)x_i$$

here t is the target output for the current training example.

θ is the output generated by the perceptron.

η is the learning rate.

η is a positive constant called the learning rate.

4) The role of the learning rate is to moderate the degree to which weights are changed at each step.

It is usually set to some small value and is sometimes made to decay as the number of weight-stunning iterations increases.

DRAWBACK:

The perceptron rule finds a successful weight vector when the training examples are linearly separable, it can fail to converge if the examples are not linearly separable.

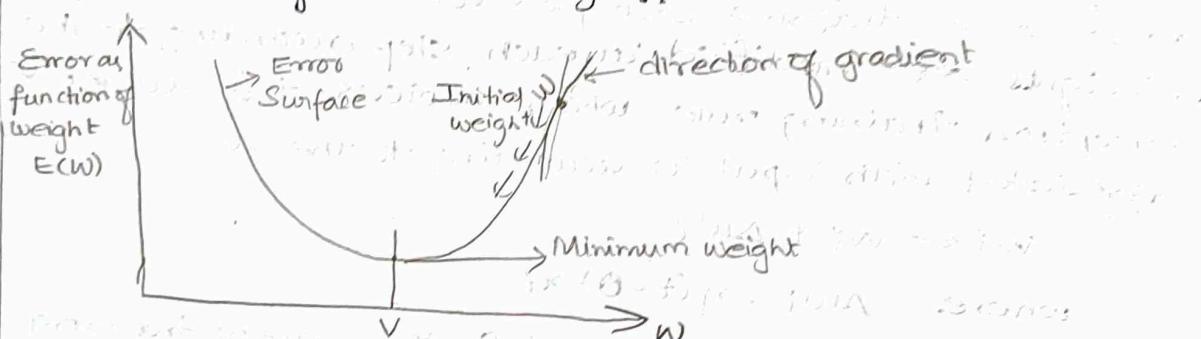
13. What do you mean by gradient descent? What are the conditions in which Gradient descent is applied?

Gradient Descent is an optimization algorithm for finding a local minimum of a differentiable function.

Gradient descent is simply used in machine learning to find the values of a function's parameters that minimize a cost function as far as possible.

It is used to determine how much a weight should be changed.

At every stage of the computation, the error E is a function of the weight. If we plot the error against the weight, we get a higher dimensional analog of something like a curve. At any point on this surface the gradient suggests how steeply the error will be reduced/increased for a change in the weights.



The gradient specifies the direction of steepest rise in E . The training rule for gradient descent is given as

$$\vec{w} \leftarrow \vec{w} + \Delta \vec{w} \text{ where } \Delta \vec{w} = -\eta \nabla E(\vec{w})$$

$\eta \rightarrow$ Learning rate

\rightarrow The constant η

\rightarrow Is a hyperparameter in the optimization algorithm that determines the step size for iterative solution to find minimum of a loss function



It is the strategy for searching through a large hypothesis space that can be applied whenever

(1) The hypothesis Space contains continuously parameterized hypotheses.

(2) The error can be differentiated wrt these hypothesis parameters.

Q) Write algorithm for Back propagation.

The BACKPROPAGATION Algorithm learns the weights for a multilayer network, given a network with a fixed set of units and interconnections.

It employs gradient descent to attempt to minimize the squared error between the network output values and the target values for these outputs.

In BACKPROPAGATION algorithm we consider networks with multiple output units rather than single units as before so we redefine E to sum the errors over all of the n/w output units.

$$E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - O_{kd})^2$$

where

Outputs — is the set of output units in the n/w

t_{kd} and O_{kd} — the target and output values associated with the k th output unit

d — training example.

BACKPROPAGATION (Training example, η , m in, n_{out} , n_{hidden})

Each training example is a pair of the form (\vec{x}, \vec{t}) where (\vec{x}) is the vector of network input value (\vec{t}) and is the vector of target network output values.

η is learning rate, n_b is the number of network inputs

n_{hidden} the number of units in the hidden layer

and n_{out} the number of output units.

The input from unit i into unit j is denoted x_{jb} and the weight from unit i to unit j is denoted w_{ji} .

1) Create a feed-forward network with n_b inputs, n_{hidden} hidden units and n_{out} output units.

2) Initialise all network weights to small random numbers.

3) Until the termination condition is met, Do

 a) For each (\vec{x}, \vec{t}) in training examples, Do

 i) Propagate the input forward through the network

 ii) Input the instance \vec{x} to the network &

 Compute the output O_u of every unit u in the network

propagate the ~~input~~ errors backward through the network.

2) For each network output unit k , calculate its error term δ_k

$$\delta_k \leftarrow O_k (1 - O_k) (t_k - O_k)$$

3) For each hidden unit h , calculate its error term δ_h

$$\delta_h \leftarrow O_h (1 - O_h) \sum_{k \text{ outputs}} w_{h,k} \delta_k$$

4) Update each network weight w_{ji}

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

where

$$\Delta w_{ji} = \eta \delta_j x_{ij}$$

15) Explain the importance of stochastic gradient descent.

Stochastic gradient descent is an optimization algorithm often used in machine learning applications to find the model parameters that correspond to the best fit between predicted and actual outputs.

- Suppose you have one million samples in your dataset (so if you use a typical GD optimization technique you will have to use all of the one million samples for 1 iteration). This problem is solved by SGD which uses only a single sample to perform each iteration. The sample is randomly shuffled & selected for performing the iteration.
- In GD the error is summed over all examples before updating weight whereas in SGD weight are updated upon examining each training example.
- Storage: It is easier to fit in the memory due to a single training example being processed by the rule.
- It is computationally fast as only 1 sample is processed at a time.
- For large datasets it can converge faster as it causes updates to the parameters more frequently.
- Due to frequent updates, the steps taken towards the minima of the ~~loss~~ loss function have oscillations that can help to get out of the local minimum of the loss function.