

UNIT - IAI
ML
DL

(i) what is machine learning?

Field of study that gives a computer capability to learn without being explicitly programmed.

Eg: online shopping

(show related product on my previous search)

→ Machine adapts to the user based on data.

* well posed learning problem:-

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

If P is measured at τ it can improve E .

(learning by experience)

→ Examples

1. Playing checker Problem
2. Hand written recognition problem
3. Robo driving learning problem

Problem	Task (τ)	Performance (P)	Experience (E)
1. Playing checkers learning	Playing again opponents to win game	Make perfect moves to win game	It plays itself to improve game
2. Hand writing recognition learning	classifying the images & text	Better classification	A database of homewrk test.
3. Robot driving learning problem	drive the car in a 4 line highway	Source to destination the avg distance travelled	Image, vehicles on road

* Perspective of Machine Learning :-

Perspective space of machine learning involves searching very large hypothesis space to determine one that fits the observed data and any prior knowledge held by learner.

Ex:-

→ sorting algorithm
→ Merge sort(ex.) choose one algorithm, your choice!

* Issue in Machine Learning :-

1. what algorithm should be used? (so many algorithms)
 2. which algorithm perform best for which types of problems?
 3. How much training data is sufficient? and testing data
 4. what kind of method should be used to reduce learning overhead
- what kind of method should be used?
For which type of data which method should be used

Designing a Learning System :

To get a successful learning system, it should be designed for a proper design, several steps should be followed
↳ for perfect & efficient

4 steps :-

- choosing the training experience
- choosing the target function
- choosing a representation for target function
- choosing a function for approximating target function
- final ob

Step 1

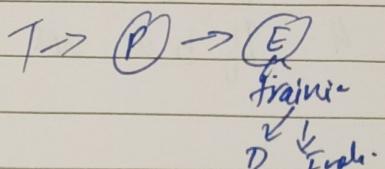
choosing the training experience :-

In choosing a training experience, 3 attributes are taken

1. type of feedback (direct, indirect)
2. degree
3. distribution example.

1) whether the training experience provides direct or indirect feedback regarding the choices made by performance system

Direct feedback - move (Physical)
 Indirect feedback (Recorded)



Eg:- checkers game, learning driving

→ You want to go with direct / Indirect. It's your choice.

2) degree to which learner will control the sequence of learning

Example:- learning driving

with trainer help with trainer partial help completely on your own

3) how well it represents the distribution of example over which the performance of final system is measured

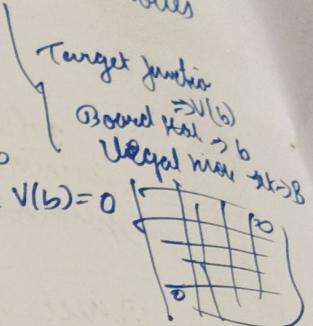
- more possible combinations
- more situations, more example
- learning over distributed range of example.

Example:- learning driving → Trainer ← ↗ ↘ ↙ ↛ without trainer

Step-2 choosing the target function

what type of knowledge is learnt and how it is used by the performance system.
Ex:- checkers game
 while moving diagonally set of all possible moves
 is called legal moves.

1. b is final board state that is won then $V(b) = 100$
2. b is final board state that is lost, then $V(b) = -100$
3. b is final board state that is draw, then $V(b) = 0$
4. If b is not final state then $V(b) = V(b')$.
 $b' \Rightarrow$ best final state.



Step-3 choosing a representation of target function

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

Features

$$x_1 - x_6$$

x_1 = No. of black pieces on board

x_2 = No. of red pieces on board

x_3 = No. of black kings on board

x_4 = No. of red kings on board

x_5 = No. of black pieces threatened by red (black beat)

x_6 = No. of red pieces threatened by black

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

Q representation of $C(b)$)

numerical coefficients (or) weight of each feature
 additional constant

Date: 25-4
choosing a learning algorithm for approximating the target function.

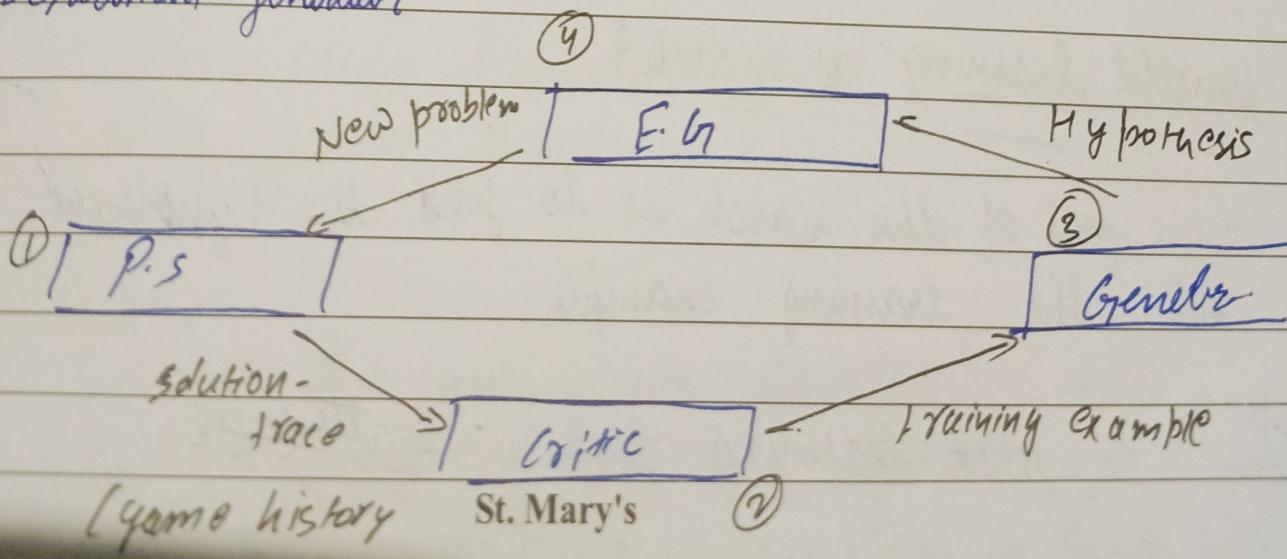
To learn a target function (1) we need a set of training examples to describe a particular board state and training value.
 (b)
 $v_{\text{train}}(b)$
 ordered pair = $(b, v_{\text{train}}(b))$
 (Training example representation)

Example:- black won the game
 i.e. $x_1=0$, which means no red
 $v_{\text{train}}(b) = +100$
 $b = (x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0)$
 $\leftarrow (x_1=3, x_2=0, x_3=1, x_4=0, x_5=0, x_6=0) +100 \rightarrow$

we need to do rule 2 steps in this phase

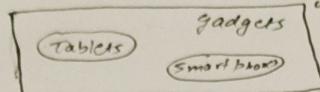
* Final Design

1. Performance system
2. Critic
3. Generalizer
4. Environment generator



* Concept learning :-

- * Features (Binary valued attribute)
 - size:- large, small $\rightarrow x_1$
 - color: Black, Blue $\rightarrow x_2$
 - screen-type: flat, folded $\rightarrow x_3$
 - shape :- square, Rectangle $\rightarrow x_4$



$$\text{concept} = \langle x_1, x_2, x_3, x_4 \rangle \rightarrow$$

$$\text{Tablet} = \langle \text{large, black, flat, square} \rangle$$

$$\text{Smartphone} = \langle \text{small, blue, folded, rectangle} \rangle$$

$$\text{No. of possible instances} = \underline{\underline{2^d}}$$

$$d = \text{no. of features}$$

$$\boxed{\text{Total possible concepts} = 2^{2^d}}$$

$$\begin{aligned} 2^4 &= 16 \\ 2^2 &= 4 \\ 2^0 &= 1 \end{aligned}$$

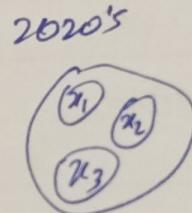
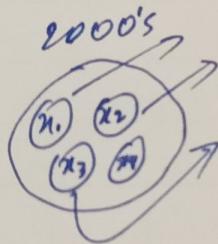
$\forall \subset \{B, F, S\}$

$\langle \phi, \phi, \phi, \phi \rangle \rightarrow \text{Reject all}$
(most hypothesis)
specific

$\langle ?, ?, ?, ?, ? \rangle \rightarrow \text{accept all}$
(most general hypothesis)

Main goal :- Finds all concept / hypothesis that are consistent.

Example



x_1, x_2, x_3
are consistent

* Concept learning as search?

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

able:

Enjoyable learning task \because 6 attributes

($\text{sky}^1, \text{Air temperature}^2, \text{humidity}^3, \text{wind}^4, \text{water}^5$ and

Date :

Sky: - 3 values - Rainy, cloudy, sunny
 Remaining attributes have only 2 values

\therefore different instances possible = $3 \times 2 \times 2 \times 2 \times 2 = 96$

Syntactically distinct hypothesis = $(5 \times 4 \times 4 \times 4 \times 4) = 5120$
 (Additionally, 2 more values - ? and Ø)

$\langle \emptyset, \text{sunny}, \text{cloudy}, \text{rainy}, ?, \rangle$

so semantically distinct hypothesis = $1 + (4 \times 3 \times 3 \times 3 \times 3) = 973$

(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
 (\because much closer to our learning problem)

Finds Algorithm :-
 $\rightarrow \emptyset$ (specific)

(Finding a maximally specific hypothesis)

- This algorithm considers only positive example

{ Yes \Rightarrow +ve
 No \Rightarrow -ve

* Representation:-

most specific hypothesis $\Rightarrow \underline{\emptyset}$

most general hypothesis $\Rightarrow \underline{?}$

\Rightarrow Algorithm:-

Step-1] Initialise the most specific hypothesis (\emptyset)

$h_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle \Rightarrow 5$ attributes

For each +ve sample | St. Mary's
 For each attribute | if (value = hypothesis value) \Rightarrow Ignore
 | else Replace with the most general hypothesis?

Example:-

	origin	manufacturer	color	year	Type	Class
1	JP	HO	Blue	1980	eco	+
2	JP	TO	green	1970	SPO	-
3	JP	TO	Blue	1990	eco	+
4	JP	AU	Red	1980	eco	-
5	JP	HO	white	1980	eco	+
6	JP	TO	green	1980	eco	-
7	JP	HO	Red	1980	eco	+

Solve

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

$$h_1 = \langle 'JP', 'HO', 'Blue', '1980', 'eco' \rangle$$

$$h_2 = h_1$$

$$h_3 = \langle JP, ?, Blue, ?, eco \rangle$$

$$h_4 = h_3$$

$$h_5 = \langle JP, ?, ?, ?, eco \rangle$$

$$h_6 = \langle JP, ?, ?, ?, eco \rangle \rightarrow \text{most general}$$

{ Only (+) taken }

* Disadvantages:-

- consider only +ve value
- h_6 may not be the sole hypothesis that fits the complete

VERSION SPACE:-

subset of hypothesis H consistent with the training examples

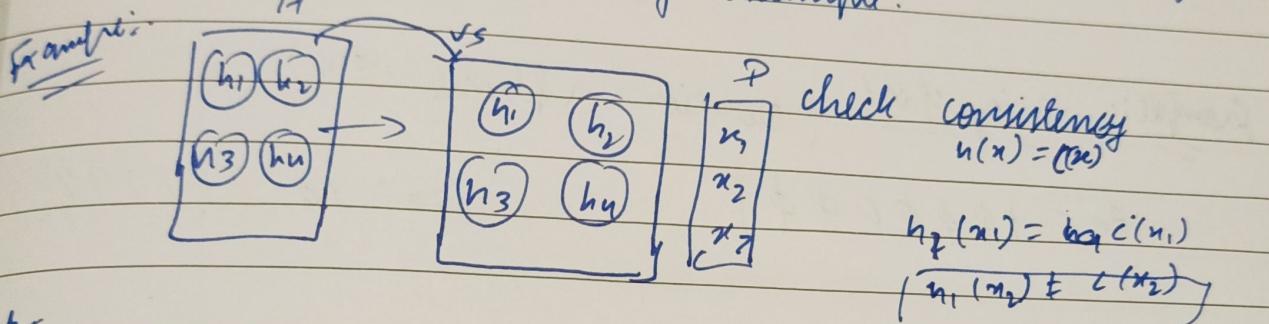
$$VS_{H,D} = \{h \in H \mid \text{consistent } (h, D)\}$$

H = hypothesis

D = training example

Algorithm to obtain version space:

1. Version space = list containing every hypothesis in H .
2. From this step, we keep on removing inconsistent hypothesis from version space.
3. for each training example, $\langle x_i, c(x_i) \rangle$ remove any hypothesis that is $h(x_i) \neq c(x_i)$
4. O/P the list of hypothesis into version space after checking for all training example.



$$\begin{aligned} h_2(x_1) &= c(x_1) \\ h_2(x_2) &= c(x_2) \\ h_3(x_3) &= c(x_3) \\ h_4(x_u) &= c(x_u) \end{aligned}$$

$$\begin{aligned} h_2(x_1) &= c(x_1) \\ h_3(x_2) &= c(x_2) \\ h_3(x_3) &= c(x_3) \\ h_4(x_u) &= c(x_u) \end{aligned}$$

are considered

CANDIDATE ELIMINATION ALGORITHM:-

- uses the concept of version space
- consider both +ve and -ve values (yes and no)
- both specific and general hypothesis
- For positive sample, move from specific to general
- For negative sample, move from general to specific

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

$$G = \{\text{?}, ?, ?, ?, ?, ?\} \uparrow$$

* Algorithm:

1. Initialise with general and specific hypothesis (S_0 and G_0)

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

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

depends on no. of attribute

2. For each example,

if example is positive

make specific to general.

else example is negative

make general to specific.

Example: enjoy short (question in phone "abc")

$$S_0 = \{\phi, \phi, \phi, \phi, \phi, \phi\}$$

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

1) +ve

$$S_1 = \{'sunny', 'warm', 'normal', 'strong', 'warm', 'same'\}$$

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

2) +ve, (specific to ~~new~~ general)

$$S_2 = \{'sunny', 'warm', '?', 'strong', 'warm', 'same'\}$$

3) -ve (general to specific)

$$S_3 = \{'sunny', 'warm', '?', 'strong', 'warm', 'same'\}$$

$$G_3 = \{<'sunny'??????> <? 'warm'?????> <????? 'same'>$$

only written

change value with specific

$$S_4 = \{'sunny', 'warm', '?', 'strong', '?', ?\}$$

$$G_4 = \{<'sunny'?????> <? 'warm'????>\}$$

Sunday
final
hypothesis

G4 doesn't
match S4
So S4 will be

INDUCTIVE BIAS :-

→ Remarks on CE and VS algorithms:-

- 1) will the CE algorithm give us correct hypothesis?
- 2) what training example should the learner request next?

→ Inductive learning :-

From example we derive rule

→ deductive learning :-

Already existing rules are applied to our examples.

* Biased hypothesis space :-

does not consider all type of training examples
solution → include all hypothesis.
Ex:- sunny & warm & normal & strong & cool & change \rightarrow Yes

* unbiased hypothesis space :-

Providing a hypothesis capable of representing all of all examples.

Possible Ensembles: $3^4 = 2 \times 2 \times 2 \times 2 = 16$

Target concept: 2^{16} (huge)
(Practically not feasible)

* Ideas of inductive bias :-

learner to infer new concepts beyond the observed training examples.

$x > y \rightarrow$ inductively inferred from
 $y \text{ is inductively inferred from } x.$

Example:

= learning algorithm $\Rightarrow L$
 training data $D_c = \{x, c(x)\}$
 New instance $= x_i$
 Represented as $L(x_i, D_c)$
 $(D_c \wedge x_i) > L(x_i, D_c)$

* Decision Tree learning:-

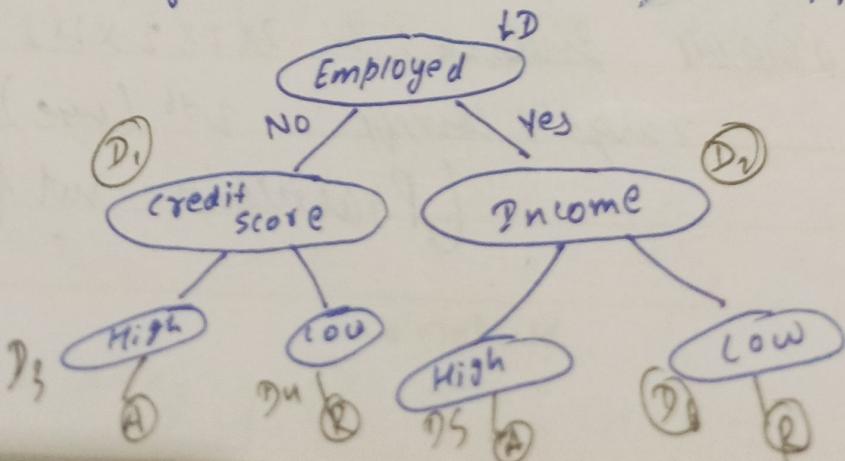
→ used in - tree structured classification & Regression

Dataset → Algorithm → classifies the data
 (decision tree algorithm)

2 types of node
 ① Decision node
 ② Leaf node

Example:- loan System:-

(decides if the loan should be approved / Rejected)



Date :

Algorithm (ID₃)

1. In the given dataset choose the target attribute
 calculate information gain of target attribute

$$IE_1 = -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

2. for remaining attributes, find entropy

$$\text{Entropy} = IE_1 \times \text{Probability.}$$

$$E(A) = \sum \frac{P_i + N_i}{P+N} I(P_i; N_i)$$

3. calculate gain = $IE_1 - E(A)$

Example:-

= step 1 Target attribute = Profit

step 2: Information gain,

$$IE_1 = -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$\Rightarrow P = \text{count (down)} = 5$$

$$N = \text{count (up)} = 5$$

$$= -\frac{5}{10} \log_2 \left(\frac{5}{10} \right) - \frac{5}{10} \log_2 \left(\frac{5}{10} \right)$$

$$= - \left(\frac{1}{2} \log_2 (2^{-1}) + \frac{1}{2} \log_2 (2^{-1}) \right)$$

$$= - \left(\frac{1}{2} \times -1 \log_2 2 + \frac{1}{2} \times -1 \log_2 2 \right) = - \left(\frac{1}{2} \times -1 + \frac{1}{2} \times -1 \right)$$

St. Mary's

$$= -(-1) = 1$$

$$\therefore IG_1 = 1$$

③ calculate entropy for remaining attributes.

$$E(A) = \sum \frac{P_i + N_i}{D+N} [IE_i \times \text{Probability}]$$

Age

- (1) Prepare a table for each attribute
 rows - values of undertaken attribute
 (old, mid, new)
 columns - values of target attribute
 (down, up)

	down	up
old	3	0
mid	2	2
new	0	3

$$\boxed{\text{Entropy} = IE_i \times \text{Probability}}$$

P = down count

N = up count

$$IE_1(\{0|10\}) = - \left(\frac{3}{5} \log \left(\frac{3}{5} \right) + \frac{2}{5} \log \left(\frac{2}{5} \right) \right) = 0$$

$$\text{Probability} = \frac{3}{10}$$

$$\text{Entropy (old)} = 0 \times \frac{3}{10} = 0$$

$$IE_1(\text{mid}) = - \left(\frac{2}{4} \log \left(\frac{2}{4} \right) + \frac{2}{4} \log \left(\frac{2}{4} \right) \right)$$

$$= - \left(\frac{1}{2} \log \left(\frac{1}{2} \right) + \frac{1}{2} \log \left(\frac{1}{2} \right) \right)$$

$$= - \left(-\frac{1}{2} + -\frac{1}{2} \right) = -(-1) = 1$$

$$\text{Probability} = \frac{4}{10}$$

$$\text{Entropy} = 1 \times \frac{4}{10} = 0.4$$

$$IE_1(\text{new}) = \left(-\frac{0}{3} \log \left(\frac{0}{3} \right) + \frac{2}{3} \log \left(\frac{2}{3} \right) \right) = 0$$

$$\text{Probability} = \frac{3}{10}$$

$$\text{Entropy} = 0 \times \frac{3}{10} = 0$$

$$\text{Entropy (Age)} = E(0) + E(m) + E(N) = 0 + 0.4 + 0 = 0.4$$

$$\text{Gain} = IE_1 - E(A) = 1 - 0.4 = 0.6 \\ \equiv$$

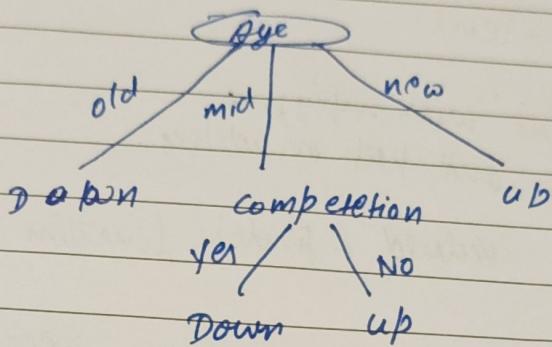
In the same way, calculate gain of other attributes

$$\text{gain (competition)} = 0.194$$

$$\text{gain (type)} = 0$$

$$\text{gain (age)} = 0.6$$

Highest gain \rightarrow root node
 (Age)



old \rightarrow all down

mid \rightarrow some down and some up

new \rightarrow all up

why competition \rightarrow next highest gain

why not other attribute (Type)

Type \Rightarrow gain = 0 and even in DT also, no requirement.
 with id

* Appropriate problem for decision tree learning:

DT algorithm is best suited to problems with the following characteristics

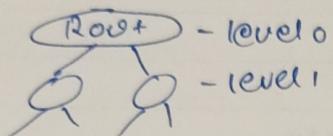
1. Instances are represented by attribute value pairs
2. The target function has discrete % values.
3. Training data can have error.
4. May contain missing attribute value also

Hypothesis space in decision tree learning:-

ID_3 can be characterized as searching a space of hypotheses for one that fits the training examples.

$ID_3 \rightarrow$ will search set of possible decision tree from available hypothesis

ID_3 perform simple to complex searching.



First start with empty tree
and keep on adding.

- every discrete valued (single) function can be described by some decision tree
- avoid major risk of searching incomplete hypothesis.
- has only single current hypothesis.
- backtracking is not possible
- can be extended easily to noisy data also.

* Inductive bias in Decision tree Algorithm:-

- set of assumptions

Inductive bias of ID_3 consists of describing the basis by which ID_3 chooses one consistent decision tree over all the possible DT's

ID_3 search strategy:

1. select in favour of shorter trees over longer ones.
2. selects element with highest IE, as root attribute over lowest IE one

types of inductive bias:

1. Restrictive bias - based on conditions
2. Red Blue Preference bias - based on priorities
 $TP_3 \Rightarrow$ Red Preference
 vs and CC \Rightarrow Restrictive

why short hypothesis?

Prefer simpler hypothesis that fits the data

Issues :-

1. overfitting the data
2. Incorporating discrete Valued algorithm
3. determining depth of tree for correct classification
4. Handling attributes with different costs.
5. Alternative measures for selecting attributes.

Introduction to artificial neural networks :- (ANN)

Example:- Neural functions are done artificially.

when we touch a hot vessel

→ A node has two parts

1. summation:

calculate the weighted sum of all the i/p

$$= x_1 w_1 + x_2 w_2 + x_3 w_3 + \dots + x_m w_m \\ = \sum x_i w_i$$

once weighted sum is calculated
sent to activation function.

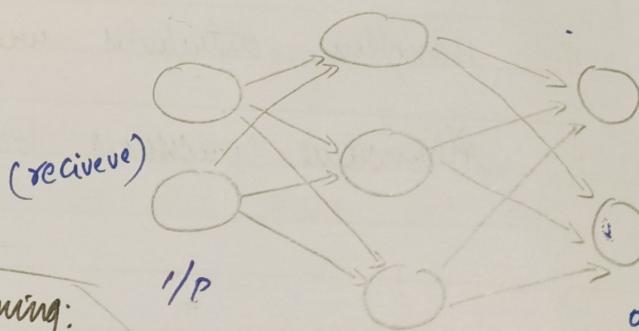
2. Activation :

generate the output based on i/p function/given
 $\rightarrow E_{new}$:

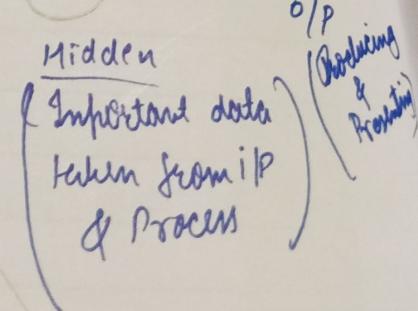
* Representation of Artificial Neural Networks:

ANN are divided into 3 parts

1. Input layer
2. hidden layers
3. output layer.



i/p

* Appropriate problem for Neural N/w learning:

Instances have many attribute value pairs
Target function has discrete values, continuous values
or combination of both.

Training example with error or missing values.

long training times are acceptable

if evaluation of the target function learnt.

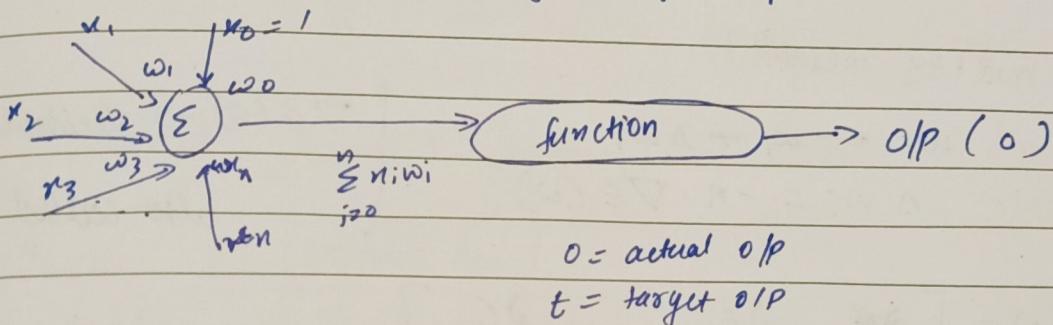
Ability for humans to learn the target function learnt by machine is not important.

PERCEPTRON:-

- * basic unit to build ANN
 - takes real valued i/p, calculate linear combination of these inputs and generate output.
- output = 1 if result > threshold
-1 Otherwise

$$o(x_1, x_2, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

$w_1, w_2, \dots, w_n \rightarrow$ weights of inputs



if actual = target \Rightarrow weights are fixed
 otherwise \Rightarrow changed

How to change weights? -

Initially, random weights

Formula to change weights:

Decide
q-14

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

$$\Delta w_i = n(t-o)x_i$$

$$n = \text{learning rate} \quad (o-1)(o-2)$$

$t = \text{target o/p}$

$o = \text{actual o/p}$

$x_i = \text{i/p associated with } w_i$

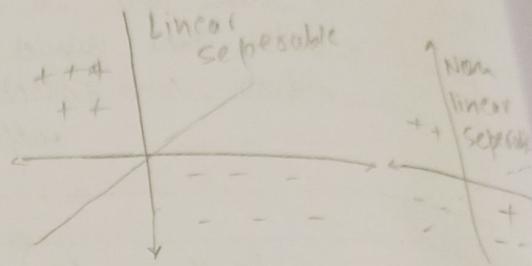
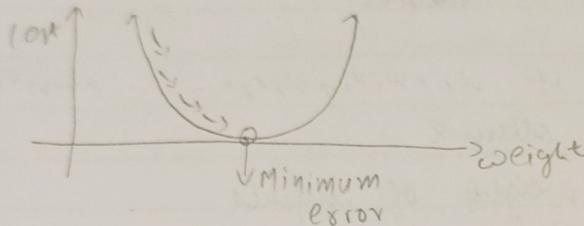
* Gradient Descent And the Delta Rule :-

Delta Rule \rightarrow used for non linearly separable data

Perception training rule \rightarrow works good only for linearly separable.

Main Idea of Delta Rule :-

\rightarrow uses gradient descent rule and finds out the best weight. \rightarrow (Reduce) \downarrow



How to modify weights?

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

$$\Delta w_i = -n \nabla E(w)$$

$\nabla E(w) \rightarrow$ derivative of error w.r.t weight
also called as gradient.

$$\nabla E(w) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \left(\frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \right) \Rightarrow \text{formula to find error}$$

$$\frac{d}{dn}(a^n) = n \cdot a^{n-1}$$

$$E = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{1}{2} \frac{\partial}{\partial w_i} \sum (t_d - o_d)^2 \\ &= \sum \frac{1}{2} \times 2x (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - w_i x_i) \end{aligned}$$

$$o_d = \sum w_i x_i$$

$$w_0 x_0 + w_1 x_1 + \dots$$

$$o_i = x_i w_i$$

$$o_d = \sum w_i x_i$$

$$\therefore \frac{\partial}{\partial w} (w x) = x \frac{\partial}{\partial w} (w)$$

$$\nabla E(w) = \sum (t_d - o_d) (-x_i)$$

$$\therefore \Delta w_i = -n \sum (t_d - o_d) x_i$$

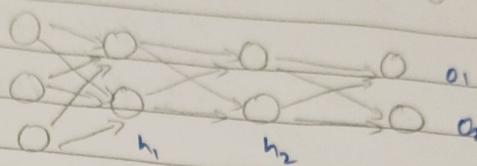
$$\Delta w_i = n \sum (t_d - o_d) x_i$$

$$\therefore \Delta w_i$$

Multi layer Neural Networks:-

→ non linearly separable data

has 2 hidden layers



BACK PROPAGATION ALGORITHM

(backward propagation of error)

[go back and correct the errors]

when error occurs, we go in backward direction

i.e. output \rightarrow hidden layer \rightarrow input layer.

Part 1: calculate forward propagation error1) calculate h_1 (in and out)

Question in file no. 3

$$h_1(\text{in}) = x_1 w_1 + x_2 w_2 + b_1$$

b_1 = bias b/w I/P \Rightarrow h

$$= 0.15 \times 0.05 + 0.2 \times 0.1 + 0.35$$

b_2 = bias b/w h \Rightarrow o/p

$$= 0.377$$

$$h_1(\text{out}) = \frac{1}{1 + e^{-h_1(\text{in})}}$$

$$= \frac{1}{1 + e^{-(0.377)}} = 0.5932$$

2) calculate h_2 (in and out)

$$h_2(\text{in}) = x_1 w_3 + x_2 w_4 + b_2$$

$$h_2(\text{out}) = 0.5968$$

$$= 0.3925$$

3) calculate o_1 (in and out)

$$o_1(\text{in}) = h_1(\text{out}) w_5 + h_2(\text{out}) w_6 + b_3$$

$$o_1(\text{out}) = \frac{1}{1 + e^{-o_1(\text{in})}}$$

$$= 1.105$$

$$= 0.7513$$

4) calculate o_2 (in and out)

$$o_2(\text{in}) = h_1(\text{out}) w_7 + h_2(\text{out}) w_8 + b_4$$

$$o_2(\text{out}) = \frac{1}{1 + e^{-o_2(\text{in})}}$$

$$= 1.12984$$

$$= 0.7729$$

5) calculate E_{Total} :

$$\begin{aligned}
 E_{\text{Total}} &= \frac{1}{2} (\text{target} - o/p)^2 \\
 &= E_{o1} + E_{o2} \quad (\text{error at } o_1 \text{ & } o_2) \\
 &= \frac{1}{2} (0.01 - 0.7513)^2 + \frac{1}{2} (0.99 - 0.7729)^2 \\
 &= \frac{1}{2} (-0.7413)^2 + \frac{1}{2} (0.2171)^2 \\
 &= 0.274 + 0.0235 = 0.29832 \quad (\text{approx})
 \end{aligned}$$

Part 2:- calculate backward propagation error
 w_5, w_6, w_7 and w_8 (output layer \rightarrow hidden layer)

First let us adjust w_5

$$w_5^* = w_5 - n \frac{\partial E_{\text{total}}}{\partial w_5}$$

$n = 0.6$
learning rate

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial \text{out } o_1} \times \frac{\partial \text{out } o_1}{\partial \text{net } o_1} \times \frac{\partial \text{net } o_1}{\partial w_5}$$

$$\begin{aligned}
 \frac{\partial E_{\text{total}}}{\partial \text{out } o_1} &= \frac{\text{out } o_1 - \text{target } o_1}{(0.1, \text{out})} = 0.751365 - 0.01 = 0.7413565 \\
 &\quad (\text{original})
 \end{aligned}$$

$$\frac{\partial \text{out } o_1}{\partial \text{net } o_1} = \text{out } o_1 \times (1 - \text{out } o_1) = 0.751365 \times (1 - 0.751365) = 0.186$$

$$\frac{\partial \text{net } o_1}{\partial w_5} = \text{out } h_1 = 0.59326992$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \dots \times \dots \times \dots = 0.8216704$$

$$w_5^* = w_5 - n \frac{\partial E_{\text{total}}}{\partial w_5}$$

$$= 0.4 - 0.6 \times 0.8216704 = 0.35069972$$

Part 3: calculating backward propagation of error.

(hidden - input layer)

(w_1, w_2, w_3, w_4)

just lets adjust w_1 ,

$n = 0.6$

$$w_1^* = w_1 - n \frac{\partial E_{\text{total}}}{\partial w_1}$$

$$\frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_{\text{total}}}{\partial \text{out } h_1} \times \frac{\partial \text{out } h_1}{\partial \text{net } h_1} \times \frac{\partial \text{net } h_1}{\partial w_1}$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out } h_1} = \frac{\partial E_{01}}{\partial \text{out } h_1} + \frac{\partial E_{02}}{\partial \text{out } h_1} \Rightarrow \frac{\partial E_{02}}{\partial \text{net } h_1} \times \frac{\partial \text{net } h_1}{\partial \text{out } h_1}, \frac{\partial E_{02}}{\partial w_1}$$

$$\begin{aligned} \frac{\partial \text{out } h_1}{\partial \text{net } h_1} &= \text{out } h_1 (1 - \text{out } h_1) \\ \frac{\partial \text{net } h_1}{\partial w_1} &= 0.175510052 \\ \frac{\partial E_{02}}{\partial \text{out } h_1} &= 0.1729(1 - 0.1729) \\ &= 0.175510052 \\ \frac{\partial E_{02}}{\partial \text{out } h_1} &= 0.172928465 - 0.99 \\ &= 0.217071535 \end{aligned}$$

$$\frac{\partial E_{02}}{\partial \text{out } h_1} = 0.217071535 \times 0.175510052 = \dots$$

same process for $\frac{\partial E_{01}}{\partial \text{out } h_1}$

other rule:

$$\frac{\partial \text{out } h_1}{\partial \text{net } h_1} = \text{out } h_1 \text{ from } h_1 \rightarrow w_1$$

also call this

$$\begin{aligned} \frac{\partial \text{out } h_1}{\partial \text{net } h_1} &= \text{out } h_1 (1 - \text{out } h_1) \\ &= 0.217071535 \end{aligned}$$

$$\frac{\partial \text{net } h_1}{\partial w_1} = \frac{\partial}{\partial w_1} (w_1 x_1 + w_2 x_2 + b_1) = x_1 = 0.05$$

$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial w_1} &= 0.036350306 \times 0.217071535 \times 0.05 \\ &= 0.00438568 \end{aligned}$$

$$\begin{aligned} w_1^* &= w_1 - n \frac{\partial E_{\text{total}}}{\partial w_1} = 6.15 - 0.6 \times 0.00438568 \\ &= 0.1499368592 \end{aligned}$$

* Remarks on Back Propagation Algorithm:

5 Remarks

1. Convergence of local minima (error backpropagation level)
↳ interconnected to each other
2. Representation power of feed forward networks. (depth nodes)
3. Hypothesis space search and inductive bias
4. hidden layer representation.
5. generalisation, overfitting and stopping criterion.
(when to stop algorithm) generalise the curve and overfit

function -
Boolean
constant
subtleties

* AN ILLUSTRATIVE EXAMPLE: FACE RECOGNITION:

- Face Recognition is a category of biomorphic software that maps individuals facial features and stores the data as a face print.
- Machine learning and image processing are backbones of FR.
- * Image Processing:-
Involves the process of computer vision
- Image Reading:- Read any image between the range of 0-255
for any image, there are 3 colors - R G B
- open cv :- Python library to solve computer vision

* Machine Learning:-

takes data as input and learn from that data

- Identified the pattern

- ML does 3 major function in Face Recognition

1. Deriving the feature vector

2. Matching algorithm

3. Face Recognition operation

→ 1) Face detection
2) Face Analysis
3) Image to data converter
4) Match finding.

* Face Recognition S/w:-

Deepvision AI, Small Time, Amazon Recognition, Facefirst, True face,

Application:- Hospital Mentity, Airline industry etc.

cognitve --

Advanced topic in Neural networks.
3 advanced topics

1) Alternative error functions :-

Adding an extra penalty term to $F(\omega)$ can reduce the problem of overfitting.

- each weight is multiplied with a constant in each iteration

: we can redefine G as.

$$E(\bar{\omega}) = \frac{1}{2} \sum_{t \in D} \sum_{R \in Q^P} (t_{R,t} - o_{R,t})^2 + \gamma \sum \omega_j^2$$

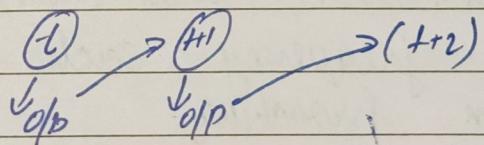
2) Alternative error minimization procedures:-
 2 procedures

1. line search

2. conjugate gradient descent method

3) Recurrent Neural N/w (RNN):-

ANN - that can be applied to time series data and uses the o/p of N/w units at time t as input to another unit at time $t+1$.



* Evaluating the hypothesis:-

Evaluating the accuracy of hypothesis as of basis of ML

For evaluating the accuracy of hypothesis, we mainly focus on 3 questions.

- 1) Hypothesis is accurate over limited sample of data, then what about additional data.
- 2) we have 2 hypothesis - one is better than other when used on sample of data.
Then will better hypothesis work good for general data also!
- 3) when we have limited data, what is the best way to use this data for both learning and estimating hypothesis.

* Motivation:-

sometime, very precise hypothesis is required. In that case estimating the accuracy is very important.

Example:- Medical treatment.

* Estimating hypothesis accuracy:-

let us make some assumptions,

- There is some space of possible instances X over which various target functions may be defined.
- different instances will have different frequencies.
- Based on frequency, each instances in X will have some unknown probability.
- Trainer will teach the machine about the training examples ↓ of target function $f(x) = y$.
Send a particular instance (x_i) along with target value y_i
Target function: $f: X \rightarrow \{0, 1\}$

Classify each instance into diff. categories based on req.

Now, after classification, what is the probable error in this estimation we made.

Error types ① Sample error ② True error

Sample Error

The error rate of hypothesis over a sample of data

$$\text{sample error} = \text{error}_s(n) = \frac{1}{n} \sum_{i=1}^n S(f(n), h(n))$$

$$(f(n), h(n)) = 1 \text{ if } f(n) \neq h(n)$$

$$= 0 \text{ otherwise}$$

True Error

The error rate over the entire distribution D

$$\text{True error} = \text{error}_T(n)$$

$$= P[f(n) \neq h(n)]$$

Basics of Sampling Theory :-

(1) Definitions

1) Random Variable :- Name of an experiment with probabilistic outcomes.
 ex:- $P(X=H) = 0.5$ $P(X=T) = 0.5$ $y_i \rightarrow$ random variable

2) Probability Distribution :- Probability distribution for a random variable Y $P(Y=y_i)$ that Y will take value y_i for each possible y_i .

[ex] when two coins are tossed, let Y represent no. of heads

$$P(Y \leq 1) = P(Y=0) + P(Y=1) = 0.25 + 0.5 = 0.75$$

3) Expected value (or) mean

$$E(Y) = \sum y_i P(Y=y_i)$$

4) Standard deviation :-

$$= \sqrt{\text{Variance}}$$

$$= \sqrt{\text{Var}(Y)}$$

5) Variance :- shows dispersion of data from its mean.

$$\text{Var}(Y) = E[(Y - \mu_Y)^2]$$

Random Variable \rightarrow mean

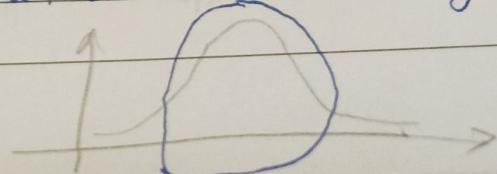
6) Binomial distribution :-

Gives probability of observing r heads in a series of n independent coin tosses

Ex:- finding prob. of 2 heads in 3 tosses

7) Normal distribution :-

bell shaped probability distribution that covers many real life values.

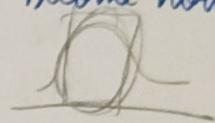


8) Central Limit Theorem :-

If you are having population with mean μ and SD as σ , and you take large random samples from that population with replacement, then the distribution will approximately become normal.

Example: $80, 40, 60 \Rightarrow \underline{\underline{60}}$

(50)



9) Estimation :-

If y is a random variable Y used to estimate some parameter P of an underlying population. ex: $P(Y = 50)$

10) Estimation Bias :-

Sample mean = \bar{y}_s
Population mean = μ_p

$$EB = \bar{y}_s - \mu_p$$

11) Confidence Interval

Probability (P) at $N\%$.

(Area under normal distribution curve in which data will fall 95% of the times)

* Difference in error of two hypothesis :-

* Hypothesis testing :

we calculate 2 mutually exclusive statements

true or true

↓ on population

↓ using sample data

and decide upon population belongs to which statement.

* Steps in hypothesis testing :-

1) make initial assumption (H_0)

2) collect data

3) Gather evidence to accept or reject null hypothesis

(we should also assume an alternate hypothesis - H_1)

→ 2 errors are possible

(Type 1 and Type 2 error)

Date: _____
→ Type I error

we implicitly know that h_0 is correct but we do not have enough evidence to prove that item, automatically, h_1 is taken as correct

→ Type II error

when we do not have enough evidence to prove the truth of h_0 , without any chance we approve h_1 .

(we do not know about truth of h_0 & h_1)

	h_0	h_1
do not reject	OK	Type I error
Reject	Type II error	OK

Comparing Learning algorithm :-

- we have so many algorithm in ML
 There is always a confusion about which one to choose.
 : we consider some factors:
- 1) time complexity - training & testing
 - 2) space complexity - less space
 - 3) sample complexity - no. of examples to train the n
 - 4) unbiased data - not learn about everything $60-30\%$
 - 5) online and offline algorithm
 Online → new data automatically updated
 Offline → obtain the entire machine
 - 6) Parallelizability → parallel algorithm → complete multiple algorithms at once
 - 7) Parametricity
 ↳ {Parametric - fixed data
 Nonparametric - change size data.}

* BAYES THEOREM:-

Bayes theorem gives the probability of an event based on prior knowledge of conditions.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

↓
hypothesis

Likelihood
Marginal
 $P(A)$ - Prior

Proof:- Let us take events : A and B

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$P(A|B)$ - Conditional Prob.

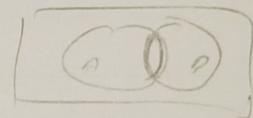
$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$P(A|B) \cdot P(B) = P(A \cap B) \quad \text{--- (1)}$$

$$P(B|A) \cdot P(A) = P(B \cap A) \quad \text{--- (2)}$$

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \rightarrow \text{Bayes thm}$$



from (1) and (2)

Terms explanation :-

A = hypothesis $P(A)$ = Prob. of hypothesis before observing given data

B = given data $P(B)$ = Prob. of given data

$P(A|B)$ = finding prob. of hypothesis when probability of training examples are given.

$P(B|A)$ = finding prob. of given data provided with prob. of hypothesis that is true.

Ex:- $P(\text{King}|\text{Face})$

\Rightarrow Prob. of face card when it king is true.

$$= \frac{P(\text{Face}|\text{King}) \cdot P(\text{King})}{P(\text{Face})}$$

Face - J, K, Q
 $3 \times 3 = 12$

$$= \frac{\frac{9}{4} \times \frac{4}{52}}{\frac{12}{52}} = \frac{1 \times \frac{1}{3}}{\frac{3}{13}} = \frac{1}{3}$$

BAYES THEOREM AND CONCEPT LEARNING :-

* what is the relation between bayes theorem and concept learning?

→ Bayes theorem calculates the probability of each possible hypothesis and outputs the most probable one.

* Before Bayes concept learning.

1) for each hypothesis in H calculate posterior probability

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)} \quad \text{--- (1)}$$

② output the hypothesis to $h_{MAP}(h_m)$ - maximum likelihood with highest probability,

$$h_{MAP} = \arg \max_{h \in H} P(h|D)$$

some assumptions:

↳ we need values $P(h)$ and $P(D|h)$

1) training data D is noise free

2) target concept c is present in hypothesis space H

3) we have no prior reason to believe that any hypothesis is more probable than any other.

$$P(h) = \frac{1}{|H|} \text{ for all } h \in H$$

$$P(D|h) = \begin{cases} 1 & \text{if } d_i = h(x_i) \text{ for all } d_i \in D \\ 0 & \text{otherwise} \end{cases}$$

From (1)

case 1: Hypothesis h is inconsistent i.e $P(D|h) = 0$

$$P(h|D) = \frac{0 \times P(h)}{P(D)} = 0 \quad \text{if } h \text{ is inconsistent}$$

case 2: Hypothesis h is consistent i.e $P(D|h) = 1$

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)} = \frac{1 \times P(h)}{P(D)} = \frac{\frac{1}{|H|}}{\frac{|S|}{|S \cap H|}} = \frac{1}{\frac{|S|}{|S \cap H|}}$$

$P(h|D) = \begin{cases} 1 & \text{if } h \text{ is consistent with } D \\ 0 & \text{otherwise} \end{cases}$

* MAXIMUM LIKELIHOOD AND LEAST SQUARED ERROR HYPOTHESIS

$D \rightarrow$ set of target values.

To find the maximum likelihood hypothesis in bayesian learning.

$$h_{MAP}(h|D) = \underset{h \in H}{\operatorname{argmax}} P(h|D) \quad \begin{cases} \text{is the one which has the Max} \\ \text{at the min squared error b/w the} \\ \text{hypothesis} \end{cases}$$

Let us take training instances (x_1, x_2, \dots, x_n) and consider target values $D = (d_1, d_2, \dots, d_m)$

\therefore we write $P(D|h)$ as product of $P(d_i|h)$

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m P(d_i|h)$$

\prod - Product
 \sum - Summation
 $\sum_{i=1}^m$

By assuming ~~great~~ ~~formal~~ normal distribution

$$f(x|u) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-u)^2}{2\sigma^2}}$$

$$\rightarrow f(d_i|h) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(d_i-h(x_i))^2}{2\sigma^2}}$$

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} (d_i - h(x_i))^2}$$

\Rightarrow use \ln ^{logarithm} function

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

\hookrightarrow constant term (so delete this)

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m -\frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

\hookrightarrow (add - since the term is we to min)

$$= \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

\hookrightarrow (constant we delete)

$$= \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m (d_i - h(x_i))^2$$

$$\therefore h_{MAP} \text{ or } h_{ML} = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m (d_i - h(x_i))^2$$

\hookrightarrow least square error.

* MINIMUM DESCRIPTION LENGTH PRINCIPLE! -

Axiom:

Representing a concept in minimum possible way
 - then it is said to be good one

Mathematical /

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

Applying algorithm.

$$\arg \max_{h \in H} \log P(D|h) + \log p(h) \quad \log(p_h) - \log a + \log b$$

$$\arg \min_{h \in H} -\log P(D|h) - \log p(h)$$

(minimum length / short hypothesis is preferred)

Example:-

- 1) let us consider a probability of designing a code to transmit message drawn at random from a set D , where probability of drawing an i^{th} msg = P_i
- 2) while transmitting, we want a code that minimizes the expected no. of bits
To do this, we should assign shorter codes to the most probable we represent the length of message i.w.r.t c as $L_C(i)$
 \hookrightarrow increase.

$$h_{MAP} = \arg \min L_C(h) + L_C D/h \quad (D/h)$$

$L_H \rightarrow$ optimal encoding for H

$L_D/h \rightarrow$ " " " " \Rightarrow given h

$$\therefore h_{MDL} = h_{MAP}$$

Bayes OPTIMAL CLASSIFIER :-

Bayes optimal classifier is a probabilistic model that makes the most probable prediction for a new example.

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

for a dataset,

$$x = \{x_1, x_2, x_3, \dots, x_n\} \{y\}$$

$$P(Y|x_1, x_2, \dots, x_n) = [P(x_1|Y) \cdot P(x_2|Y) \cdots P(x_n|Y)] \cdot P(Y)$$

$$\text{St. Mary's } P(x_1), P(x_2), \dots, P(x_n)$$

$$= P(y) \prod_{i=1}^n P(x_i|y)$$

$$\Rightarrow \frac{P(x_1)P(x_2)\dots P(x_n)}{P(y)} \xrightarrow{\text{Probability of this value will remain constant, so it will cancel out}} P(x_1)P(x_2)\dots P(x_n)$$

Ques
Play

yes	9	9/14
no	5	5/14
Total	14	100%
Total (sunny, hot)		

Pic - 45

$$P(\text{Yes} / \text{sunny, hot}) = P(\text{sunny}) \times P(\text{hot}) \times P(\text{Yes})$$

$$= \frac{2}{9} \times \frac{2}{9} \times \frac{9}{14} = 0.031$$

* $P(\text{No} / \text{sunny, hot}) = P(\text{sunny}) \times P(\text{hot}) \times P(\text{No})$

$$= \frac{3}{5} \times \frac{2}{5} \times \frac{5}{14} = 0.08571$$

$$\text{Total} = 0.031 + 0.08571 = 0.27$$

$$P(\text{Yes}) = \frac{0.031}{0.27} = 0.114$$

$$P(\text{No}) = \frac{0.08571}{0.27} = 0.317$$

} Probability of No is more.
; Player will not enjoy game.

Gibbs algorithm:
(similar to Bayes rule.)

▷ choose one hypothesis at random, according to $P(h/D)$
⇒ use this to classify new instance.

$$2[\text{error giv}] \leq 2E[\text{error Bayes optimal}]$$

NAIVE BAYES CLASSIFIER:

→ classification technique based on Bayes theorem with an assumption of independence among features

$$P(A/B) = \frac{P(B/A)P(A)}{P(B)}$$

Pic - 6

fruit: yellow, sweet, bad

Ques Pic - b
 $P(\text{yellow/orange}) = \frac{P(\text{orange/yellow})P(\text{yellow})}{P(\text{orange})}$

$$= \frac{\frac{350}{800} \times \frac{900}{1200}}{\frac{650}{1200}} = 0.53$$

$$P(\text{sweet} / \text{orange}) = \frac{P(O/S) P(S)}{P(O)} = 0.69$$

$$P(\text{Fruit} / \text{orange}) = P(Y/O) \times P(S/Y) \times P(L/Y) \\ = 0.53 \times 0.69 \times 0 = 0$$

$$P(\text{Fruit} / \text{Banana}) = P(Y/B) \times P(S/B) \times P(L/B) \\ = 1 \times 0.75 \times 0.89 = 0.65 > 0.69$$

$$P(\text{Fruit} / \text{Others}) = P(Y/O) \times P(S/O) \times P(L/O) \\ = 0.33 \times 0.66 \times 0.33 = 0.072 \quad \text{more probability}$$

Banana is the yellow, long, sweet, and related to the most probable one.

BAYESIAN BELIEF NETWORKS:-

- 2 important concepts

- (1) directed acyclic graph (DAG)
- (2) conditional probability Table (CPT)

* Directed acyclic graph.

RAIN \rightarrow Node, R \downarrow hydr.

\downarrow dry bark/reid

dry bark \downarrow cathode/bogbam

catmilk

* C.P.T

calculate
(always respect to parent node)

R \bullet $\sim R$ \bullet

B $9/48$ $18/48$

$\sim B$ $3/48$ $18/48$

$(B = T \& R = T) = 9/48 = 0.19$

$(B = T \& R = F) = 18/48 = 0.375$

$(B = F \& R = T) = 3/48 = 0.06$

$(B = F \& R = F) = 18/48 = 0.375$

Bayesian belief nw is a probabilistic graphical model (PBM)
that represents conditioned dependence between random variable
through DAG.

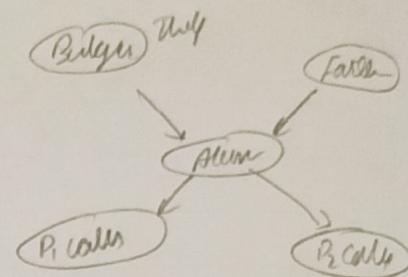
also suitable for representing probabilistic relation between multiple
events (more than 2 events)

Given Probabilities are:

$$\begin{array}{ll} P(B=T) = 0.001 & P(E=T) = 0.002 \\ P(B=F) = 0.999 & P(E=F) = 0.998 \end{array}$$

Probability of alarm:- (2 parents)

B	E	$P(A=T)$	$P(A=F)$
T	T	0.95	0.05
T	F	0.99	0.06
F	T	0.29	0.71
F	F	0.001	0.999



Probability of P_1 (1 parent)

A	$P(P_1=T)$	$P(P_1=F)$
T	0.90	0.10
F	0.05	0.95

Probability of P_2

A	$P(P_2=T)$	$P(P_2=F)$
T	0.70	0.30
F	0.01	0.99

Example

Find the probability of P_1 is T, P_2 is T, A is T, B is F and E is F
i.e $P(P_1, P_2, A, \sim B, \sim E)$

$$= P(P_1/A) P(P_2/A) P(A/\sim B, \sim E) P(\sim B) \cdot P(\sim E)$$

$$= 0.90 \times 0.70 \times 0.001 \times 0.999 \times 0.998 = 0.00062$$

EM ALGORITHM: (Expectation - Maximisation)

- used to find latent variable → not directly observed.
- basic for many unsupervised clustering algorithm

* Steps involved in EM algorithm:

1. Initially, a set of initial values are considered
2. Next step - expectation step - e-step
Here, we use observed data to estimate or guess the values of missing/unobserved data.
3. Maximisation step or m-step
Here, we use the complete data generated in preceding e-step to update the values.
4. We check if values are converging / not
If converging - stop. Otherwise, Repeat step 2 and still the convergence check

Date :

Page No.

- * usage
 - 1) used to fill missing data
 - 2) used for unsupervised clustering
 - 3) used to calculate values of related variables

+ Advantages	- Disadvantages
With each iteration, likelihood increases	1) slow convergence
2) step wise methods	2) makes convergence to local optimum only
are easy to implement	

* Disadvantages

- 1) slow convergence
- 2) makes convergence to local optimum only

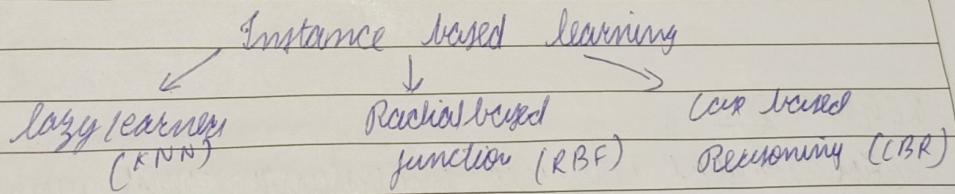
Instance Based Learning :-

memorise and then apply (not use brain)

- > Instead of performing explicit generalisation, it compares new problems with instances in training, which are stored in memory.

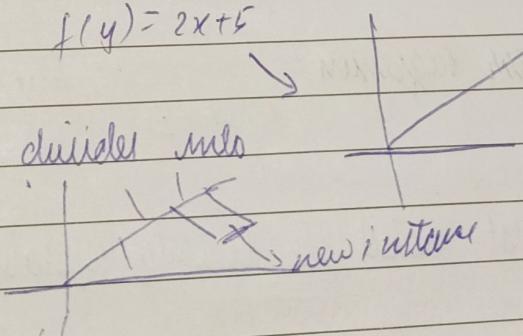
Example: spam mail - (spam filter all spam emails goes to spam folder after checking)

- > also called as memory based learning / lazy learning
- > done with 3 different approaches.



Example:- Initially, we have $f(y) = 2x + 5$

If we know $f(y) = 2x^2 + 3x^2$ and divide into pocket



↳ This graph will guide by previous memory.

K-Nearest Neighbour Algorithm (KNN) :- (example for lazy learning)

Example:-

Given data query $\Rightarrow x = (\text{math} = 6, \text{CS} = 8)$

and $k=3$ \rightarrow nearest neighbor
classification - Pass/Fail

	Math	CS	Result
1)	4	3	F
2)	6	7	P
3)	7	8	P
4)	5	5	F
5)	8	8	P

Euclidean distance (d)

$$d = \sqrt{(x_0 - x_1)^2 + (x_0 - x_2)^2}$$

o = observed value

a = actual value

\therefore calculated $\Rightarrow d_1 = \sqrt{(6-4)^2 + (8-3)^2} = \sqrt{2^2 + 5^2} = \sqrt{29} = 5.38$

$\therefore d_2 = \sqrt{(6-6)^2 + (8-7)^2} = \sqrt{0+1} = \text{D}_2$ \rightarrow 1st nearest value

$\therefore d_3 = \sqrt{(6-7)^2 + (8-8)^2} = \sqrt{1+0} = \text{D}_3$ \rightarrow 2nd nearest value

$\therefore d_4 = \sqrt{(6-5)^2 + (8-5)^2} = \sqrt{1+9} = \sqrt{10} = 3.16$

$\therefore d_5 = \sqrt{(6-8)^2 + (8-8)^2} = \sqrt{4+0} = \sqrt{4} = 2$ \rightarrow 3rd nearest value

3 P
of
3 P \rightarrow Pass
 \rightarrow Majority go for it

REGRESSION:-

Statistical tool used to understand and quantify the relation between 2 more variables.

* linear Regression:

$$y = \beta_0 + \beta_1 x + \epsilon$$

\rightarrow best suited for linearly separable data

y = dependent variable

x = independent variable

β_0 = constant / Intercept

β_1 = slope / coefficient

ϵ = error.

lazy learning

Date :

Page No.

* Locally weighted regression:-

- To overcome the problem of non linearly separable data
- LWR algorithm assign weight to data to overcome the problem
- Computationally more expensive.

↳ finding weight for every data

Finding weights? By kernel smoothing

$$D = \alpha c \frac{e^{-\frac{\|x - x_0\|^2}{2c^2}}}{2c^2} \quad x \rightarrow \text{each training i/p}$$

$x_0 \rightarrow \text{value we are predicting}$

$c \rightarrow \text{constant}$

- we construct a weight matrix (w)

→ for each training P/P(x) and for the value we are trying to predict (x_0)

weight matrix - diagonal matrix

$$\beta = (X' W X)^{-1} X' W Y$$

β - model parameter

Then, prediction can be defined as

$$Y = \beta X_0, \quad Y \Rightarrow \text{Prediction.}$$

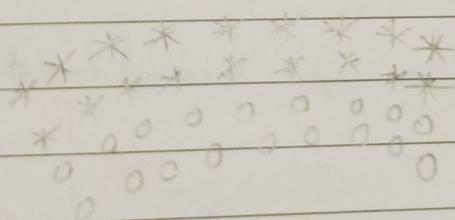
* Drawback:-

- 1) Need to calculate whole dataset everytime
- 2) Computation cost is more.
- 3) memory requirement is more.

RADIAL BASIS FUNCTIONS:- \rightarrow used in ANN

\rightarrow has only one hidden nodes.

Example:



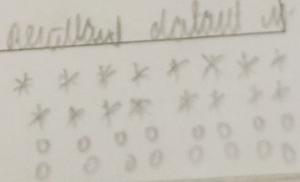
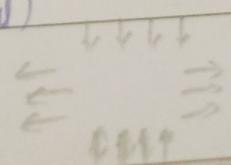
→ 2 classes
(Standard circle)

→ data is not linearly separable.

\rightarrow 2 ways.

1. Increase the dimensionality (2D-3D) (not mentioned, only based on diagram)
2. expand the direction (Horizontal)
compress the direction (Vertical)

St. Mary's



now RBF does this

→ consider one center randomly
→ draw concentric circle → one centre so many units

to expand / compress, we use ③ functions.

▷ multiquadrature

$$\phi(r) = (r^2 + c^2)^{1/2}$$
$$c > 0 \Rightarrow \text{constant}$$

▷ mean multiquadrature

$$\phi(r) = \frac{1}{(r^2 + c^2)^{1/2}}$$



▷ Exponential function

$$\phi(r) = e^{-\frac{r^2}{2c^2}}$$
$$c > 0 \Rightarrow \text{exp}$$

→ use any one of these
for combination

CASE BASED REASONING:- (CBR)

All instances based learner have 3 properties

1) fuzzy learners

2) classification is different for each instance

3) instance are represented with n dimensions euclidean space

For CBR,

everything is considered as case and based on previous case

we propose a solution

- instances are represented as symbols (not values)

↳ even

⇒ CBR has 3 components

1. similarity function or distance measure

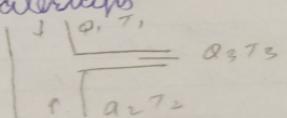
2. Approximation / Adjustment of instances

3. symbolic representations of instances

for modelling CBR, we use CADET system

↳ core visual design tool { has 75 predefined libraries }

Ex- Modern waterfalls



junction

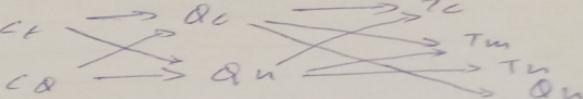
$$Q_1 \rightarrow Q_3$$

$$Q_2 \rightarrow Q_3$$

$$T_1 \rightarrow T_3$$

$$T_2 \rightarrow T_3$$

use this
make sure



REMARKS ON LAZY AND EAGER ALGORITHMS:-

* Lazy learning:-

1. Simplicity stores training data and waits until it gets a test sample.
2. less training time, more prediction time.
3. ex: All instance based learning algorithms

* Eager learning:-

1. When we give a training set, it computes a model for classification before getting new example.
2. more training time, less prediction time.
3. ex:- Decision tree, Naive Bayes, ANN etc.