

1.11219.

Machine learning:

A computer program is said to learn w.r.t task 'T' with the performance measure of 'P' and learning experience 'E' if there is improvement in learning performance measure 'P' w.r.t to task 'T' & learning experience 'E'.

* Designing a learning system:

→ machine learning:

making the machine to learn by giving data and work accordingly.

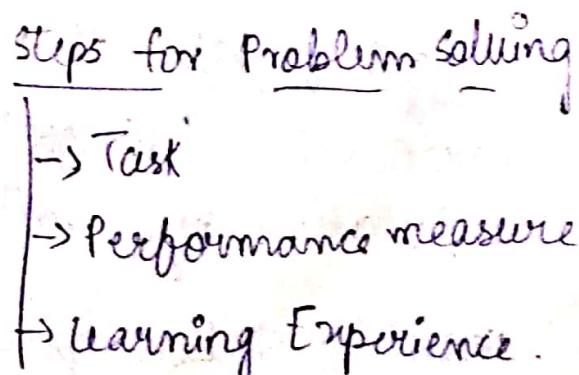
* Types of machine learnings:

- supervised learning (target point / class label is given)
- unsupervised learning (class label is not given)
- reinforcement learning (providing feedback to our model)
- semi-supervised learning (partial class labels are available)

* Well-posed problems:

problems that are solved using machine learning and they are:

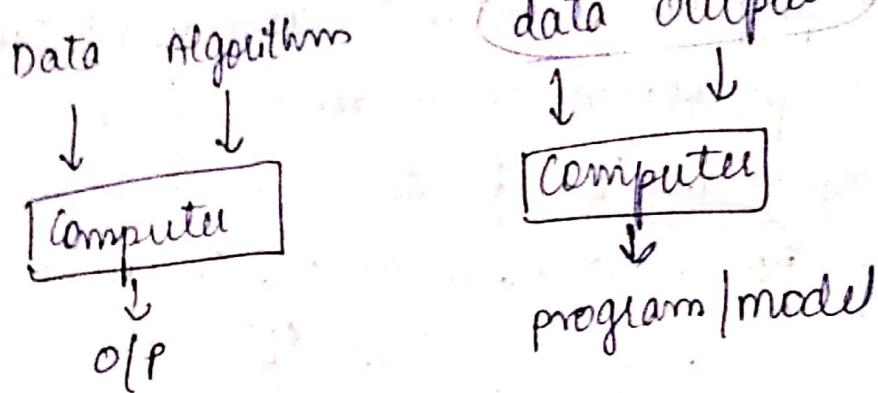
- 1) Handwritten words
- 2) Driverless cars
- 3) checkers problem
- 4) Predicting faults



*Designing a learning system:

14/12/19

together = Experience



* Regression: Predicting no. of defects.

* Classification: Predicting whether s/w has defects or not i.e. whether it is working or not.

* cluster: set of similar objects.

* Designing a learning system:

① choose the learning Experience.

Example:

- games played against themselves.
- games played against experts.
- Table of possible moves from a particular point.

② Target function: It depicts how to classify the characteristics (like for example for email system - separation of emails into inbox & spam box).

Here, in checkers problem, the target function

is to choose an efficient / appropriate move:
choose move: $B \rightarrow M$

Evaluating function: For each move, depicting its outcome, whether the move which has selected leads to either winning state or loosing state or draw state.

Generic way: $V(b) = B \rightarrow R$
 \downarrow set of all real no's.
set of all possible moves.

Example: $V(b) = 100 \Rightarrow$ winning state

$V(b) = -100 \Rightarrow$ losing state

$V(b) = 0 \Rightarrow$ draw state.

Target functions: criteria to built our systems.

③ Representation of Target Functions:

x_1 : Black cells.

x_2 : Blue cells.

x_3 : No. of black kings.

x_4 : No. of blue kings.

$$V(b) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 \quad (\text{linear model}).$$

(Similar to $y = m_1 x_1 + m_2 x_2 + \dots + c$ (i.e. $y = mx + c$))

* The previous eq is one of the various models.

* To develop a Target function, best suitable model must be chosen (in Representation of Target functions).

$$v(b) = 0.5 + 0.7x_1 \\ + 0.3x_2 \\ + 0.7x_3 \\ + 0.3x_4$$

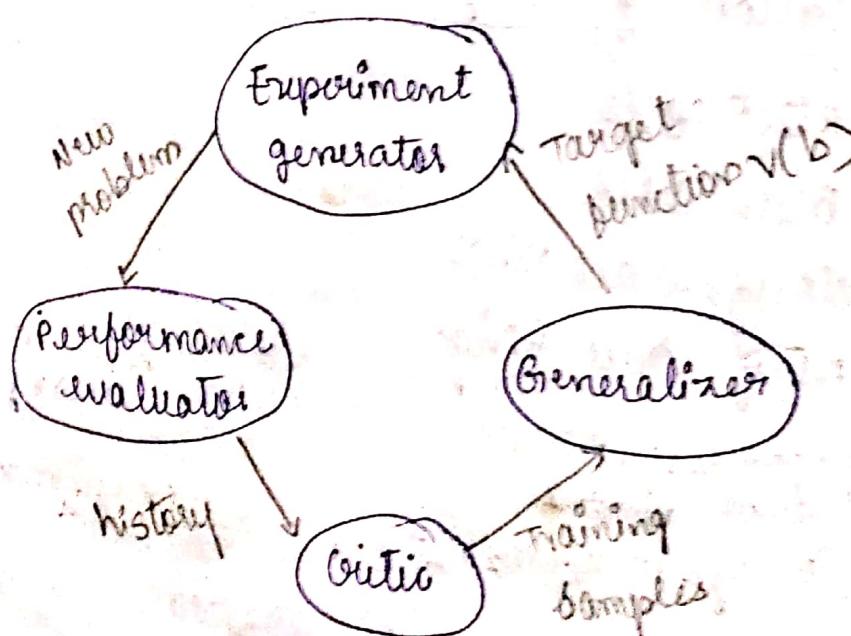
→ Approximation Algorithm: It is responsible for generating values of intercepts (w_1, w_2, w_3, w_4, w_0) or constants.

$$\text{Ex: } v(b) = 0.5 + 0.7x_1 + 0.3x_2 + 0.7x_3 + 0.3x_4$$

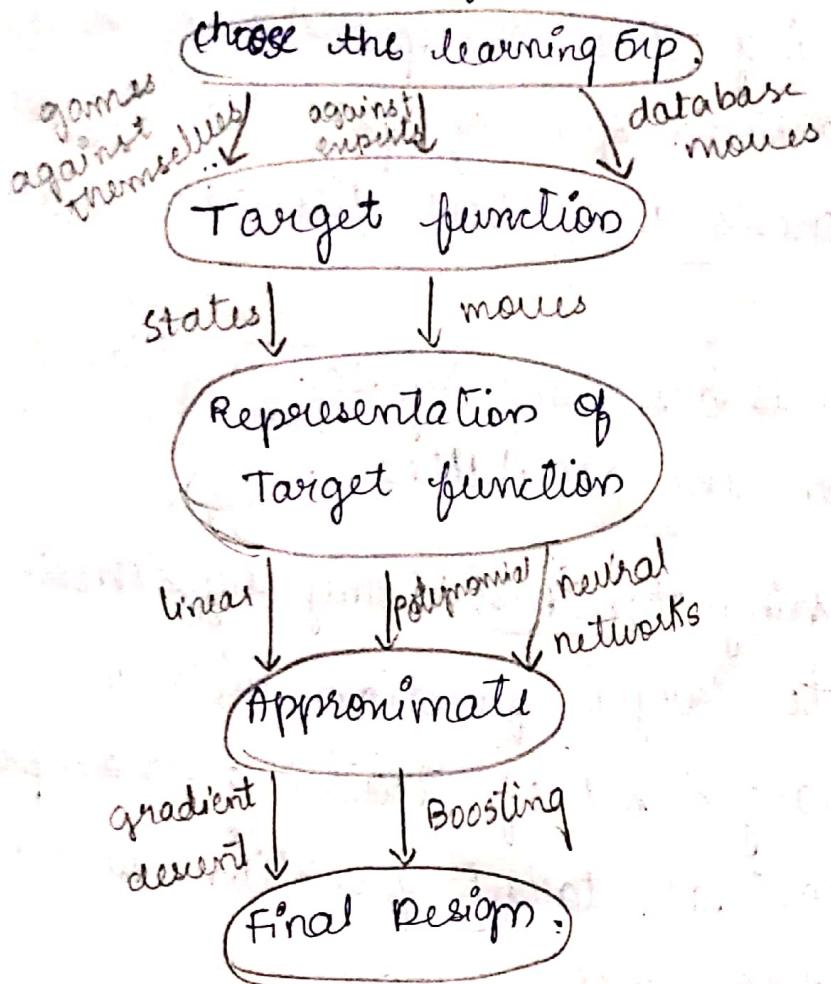
We will give x values we need to check for the $v(b)$ whether positive or negative and make a move.

Final design

Structure: Every model of mt has the same 4 modules : i) Critic module.
ii) Generalizer
iii) Experiment Generator
iv) Performance Evaluator.



* Design of learning system:



* Issues in machine learning:

1. what are the various algorithms available and how to select the best algorithm for our class..?
2. How much of training data is sufficient?
3. How to choose learning experience for the model?
4. How to define target functions to the model?
5. How the learner alters itself to better represent the target function? evaluate ML model.

Testing data: Data which is used to test our model.

Training data: Data which is used to train our model.
construct our ML

16/12/19.

* concept learning: It is the problem of searching through pre-defined space of potential hypothesis for the hypothesis that best fits the training examples.

* hypothesis: It is described by conjunction of constraints on the attributes.

* Inductive learning hypothesis: Any hypothesis formed to approximate target function well over sufficiently large set of training examples will also approximate target function well over unobserved examples.

* Positive samples: The sample for which class label has boolean value 1.

* Negative samples: The sample for which class label has boolean value 0.

→ $C(x)=1$, if training sample is +ve; $C(x)=0$ if training sample is -ve
In order to make our machine learn, both positive &

negative samples are to be considered

Ex Sky Airtemp humidity wind water forecast

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

supervised learning

Let h_j, h_k be boolean valued functions defined over \mathcal{X} , then h_j is more general than or equal to h_k which can be represented as ($h_j \geq_g h_k$)

$$\text{iff } \forall x \in \mathcal{X} [h_k(x) = 1] \rightarrow (h_j(x) = 1)$$

stands for
general

Let \mathcal{S} samples \rightarrow (these \mathcal{S} are nothing but hypothesis).

(1) <sunny, warm, ?, ?, ?, warm, same> \rightarrow specific

(2) <sunny, ?, ?, ?, ?, ?, ?> \rightarrow generic.

If sample is specific, then it also satisfies to be a generic sample.

Find S algorithm:

* Finding a maximally specific hypothesis (h)

Step-1: Initialize h to most specific hypothesis H

Step-2: For each positive training instance x
 \rightarrow For each attribute constraint a_i , if satisfied by x , then do nothing. else, replace a_i and in ' h ' by next most general constraint that is satisfied by x .

Step-3: Output the hypothesis h .

NOTE: Find S algorithm will not consider negative training samples.

Example:

→ 1) $h < \theta, \theta, \theta, \theta, \theta, \theta >$

~~h(x)~~

1st sample $\leftarrow h < \text{sunny}, \text{warm}, \text{normal}, \text{strong}, \text{warm}, \text{same} \right>$

2nd sample $\leftarrow h < \text{sun}, \text{warm}, ?, \text{strong}, \text{warm}, \text{same} \right>$

3rd sample $\leftarrow h < ?, \text{warm} \right>$

is -ve

so ignored.

4th sample $\leftarrow h < \text{sunny}, \text{warm}, ?, \text{strong}, ?, ? \right>$

align:

* A hypothesis 'h' is consistent with set of Training examples D, iff $h(x) = c(x)$.

→ for each example $\langle x, c(x) \rangle$ in D.

consistent (h, D) = $(\forall \langle x, c(x) \rangle \in D)$

1st sample circle big blue Yes.

2nd sample triangle small red NO.

Hypothesis - 1 $\rightarrow \boxed{\langle ?, \text{big}, ? \rangle \rightarrow \text{Yes}}$

Hypothesis - 2 $\rightarrow \boxed{\langle ?, \text{small}, ? \rangle \rightarrow \text{Yes}}$

Now, 1st hypothesis is consistent because, it is correctly classifying the above samples (ie if 2nd attribute is big, it has class label as Yes, otherwise → NO).
∴ 1st sample & 2nd are classified correctly w.r.t Hyp-1).

* Version Space: Denoted by $V_{H,D}$, whereas 2nd hyp is not consistent because it is not classifying given samples correctly. (i.e. hyp-2 says, if 2nd attribute is small, then class label is yes, but the samples are not satisfied).

∴ Hyp-1 falls under version space.

* Version Space: Denoted by $V_{H,D}$, w.r.t hypothesis space H , if training examples D is subset of hypothesis from H consistent in Training

example in D : $\{V_{H,D} = \{h \in H \mid \text{Consistent}(h, D)\}\}$.

* List then Eliminate algorithm:

Step 1: Initialize Version space with a list containing every hypothesis in H .

Step 2: For every training example $x_i, c(x_i)$ remove from version space, any hypothesis h , for which

$$h(x_i) \neq c(x_i)$$

Step 3: Output the list of hypothesis in Version Space.

* General boundary:

→ General boundary 'G' w.r.t hypothesis space H , and training data D is set of maximally general members of H consistent with D .

* Specific boundary:

→ Specific boundary 'S' w.r.t hypothesis space H , and training data D is set of minimally general.

(maximally specific) members of H which is consistent with D .

* Candidate elimination algorithm:

- 1) It is applicable for both positive and negative samples.
- 2) If it is +ve sample, we would work with specific boundary $\rightarrow S_0 = \{0, 0, 0, 0, 0, 0\}$.
- 3) If it is -ve sample, we would work with general boundary $\rightarrow G_0 = \{?, ?, ?, ?, ?, ?\}$.

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

Initialize specific boundary.

Step 1: $S_0 = \{0, 0, 0, 0, 0, 0\}$

$G_0 = \{?, ?, ?, ?, ?, ?\} \rightarrow$ Initialize generic boundary.

Step 2: For first training sample,
it is +ve sample, so work with specific boundary

$$S_1 = \{\text{Sunny, warm, high, normal, strong, warm, same}\}$$

$$G_1 = \{?, ?, ?, ?, ?, ?\} \rightarrow \text{Same as } G_0 \text{ (no change)}$$

Now, check whether S_1 & G_1 is consistent with 1st sample or not.

Step 3: For 2nd sample:

→ it is true, work with specific boundary.

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

$G_{12} = \{ ?, ?, ?, ?, ?, ? \} \rightarrow \text{same as } G_1 \text{ (no change)}$

→ check whether S_2, G_{12} is consistent with 1st & 2nd samples (or) not

Step 4: For 3rd training sample:

→ it is 3rd sample, work with generic boundary.

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

(no change, write same as S_2)

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

$\langle ?, ?, \text{normal, ?, ?, ?} \rangle, \langle ?, ?, ?, \text{cool, ?, ?} \rangle \}$

(it is not satisfying 3rd sample)

(e.g. high → yes), so inconsistent,

hence it is removed

(it is not satisfying 1st sample)

(e.g. warm → yes) So inconsis-

tent, hence it is removed

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

[we are checking only for 1st, 2nd & 3rd samples so it is consistent]

[∴ for 4th sample (cool → yes) it becomes inconsistent]

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

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

{

Step 5: For fourth sample
it is the sample, so work with specific boundary

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

G_{14} = 4th sample is the, we don't apply general boundary.

We consider same G_3 & check whether consistent for all 4 samples (or) not

$\therefore G_{14} = \{ \langle \text{sunny, ?, ?, ?, ?, ?} \rangle, \langle ?, \text{warm, ?, ?, ?, ?} \rangle,$
 $\langle ?, ?, ?, ?, ?, \text{same} \rangle \}$

(it is satisfying 1, 2, 3 samples, but not 4th sample.
Hence eliminated)

~~G_{14}~~ : $G_{14} = \{ \langle \text{sunny, ?, ?, ?, ?, ?} \rangle, \langle ?, \text{warm, ?, ?, ?, ?} \rangle \}$

→ Finally :

$S = \{ \text{sunny, warm, ?, strong, ?, ?} \}$

$G_1 = \{ \langle \text{sunny, ?, ?, ?, ?, ?} \rangle, \langle ?, \text{warm, ?, ?, ?, ?} \rangle \}$

Version Space:

$S \quad \boxed{\langle \text{sunny, warm, ?, strong, ?, ?} \rangle}$

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

$\langle ?, \text{warm, ?, strong, ?, ?} \rangle$ → intermediate hypothesis

$G_1 \quad \boxed{\langle \text{sunny, ?, ?, ?, ?, ?} \rangle \quad \langle ?, \text{warm, ?, ?, ?, ?} \rangle}$

Q) Size color shape label

- big red circle NO
small red triangle NO
small red circle YES (wrong)
big blue circle NO
small blue circle YES

for step 1: $s_0 = \{0, 0, 0\}$

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

Step 2: Take 1st sample

→ it is -ve sample so, we work on general boundary

$$s_1 = \{0, 0, 0\} \rightarrow \text{same as } s_0$$

$$G_1 = \{<\text{small}, ?, ?>, <?, \text{blue}, ?>, <?, ?, \text{triangle}>\}$$

→ we need to check consistency of all these attributes

for only 1st sample

Step 3: Take 3rd sample

→ it is -ve sample, so we work on general boundary.

$$s_2 = \{0, 0, 0\} \rightarrow \text{same as } s_1$$

$$G_2 = \{<\text{big}, ?, ?>, <?, \text{blue}, ?>, <?, ?, \text{circle}>\}$$

(it is not consistent in
(1st sample since eliminated))

(it is not consistent in
(1st sample i.e. (circle → NO))
hence eliminated)

∴ By checking consistency on 1st two samples we get

$$G_2 = \{<?, \text{blue}, ?>\}$$

Step 4: Take 3rd sample.

→ it is the sample, work on specific boundary.

$S_3 = \{ \text{small, red, circle} \}$.

$G_3 = \{ ?, \text{blue}, ? \} \rightarrow \text{same as } G_2$.

Check consistency for 1st three samples.

Here $G_3 = \{ ?, \text{blue}, ? \}$

(it is not satisfying Sample 3 ie red → Yes)
hence eliminated

$\therefore G_3 = \{ ?, ?, ? \}$

Step 5: Take 4th sample.

→ it is -ve sample, work on general boundary.

$S_4 = \{ \text{small, red, circle} \} \rightarrow \text{same as } S_3$

$G_4 = \{ \text{small, ?, ?} \}$

it is not

solving:

28) size color shape label.

big red circle NO.

Small red Δ NO.

small red circle Yes

big blue circle NO.

small blue circle Yes

def: step 1: Initially;

$$S_0 = \{\emptyset, \emptyset, \emptyset\}$$

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

step 2: Take 1st Sample

→ it is -ve sample, so work with Generic boundary.

$$S_0 = \{\emptyset, \emptyset, \emptyset\} \Rightarrow \text{same as } S_0$$

$$G_1 = \{<\text{small}, ?, ?> <?, \text{blue}, ?> <?, ?, \text{triangle}>\}$$

satisfies 2nd sample also

step 3: Take 2nd sample.

→ it is -ve sample, so work with Generic boundary.

$$S_1 = \{\emptyset, \emptyset, \emptyset\} \Rightarrow \text{same as } S_1$$

$$G_2 = \{<\text{small}, \text{blue}, ?> <\text{small}, ?, \text{circle}> <?, \text{blue}, ?> \\ <\text{big}, ?, \Delta> <?, \text{blue}, \Delta>\}$$

step 4: 3rd sample is +ve sample, so Specific boundary.

$$S_2 = \{\text{small, red, circle}\}$$

$$G_3 = \{<\text{small}, \text{blue}, ?> <\text{small}, ?, \text{circle}> <?, \text{blue}, ?>$$

$\langle \text{big, ?, } \Delta^{\text{le}} \rangle \langle ?, \text{blue, } \Delta^{\text{le}} \rangle \}$

$\therefore G_3 = \{ \text{small, ?, circle} \}$

Step 5: It is -ve sample, so Generic boundary.

$G_4 = \text{Same} = \{ \text{small, red, circle} \}$

$G_4 = \{ \langle \text{small, ?, circle} \rangle \}$

Step 6: +ve sample $\therefore S_5 = \{ \text{small, ?, circle} \}$

satisfies all samples.

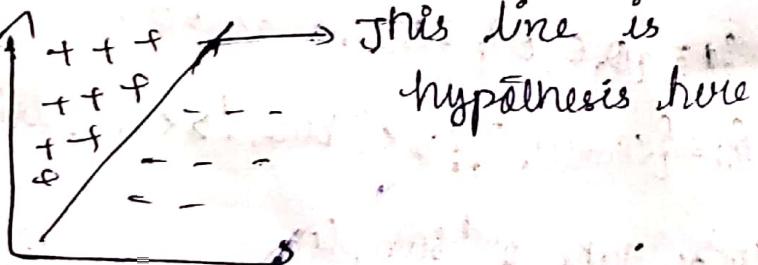
$G_5 = \{ \text{small, ?, circle} \} \rightarrow$ satisfies all samples.

finally intermediate sample = $\{ \langle \text{small, ?, circle} \rangle \}$

it is the best hypothesis.

* The final o/p of machine learning model is

Hypothesis: Eq:



* Remarks on Version Space & Candidate Elimination.

Will the candidate elimination algorithm converges to the correct hypothesis (Yes it converges, but there's a problem when there is an error in the data set).

- 1) Which training example should the learner request next.
- 2) How can partially learned concepts be used.

21/12/19:

* Inductive Bias:

→ Using find's find the hypothesis.

Sunny Warm Normal Strong Cool Change Yes

Cloudy Warm Normal Strong Cool Change Yes.

Rainy Warm Normal Strong Cool Change No

∴ Find S can be applied only for the samples.

Sample 1: $s_1 = \langle \text{sunny}, \text{warm}, \text{normal}, \text{strong}, \text{cool}, \text{change} \rangle$

Sample 2: $s_2 = \langle ?, \text{warm}, \text{normal}, \text{strong}, \text{cool}, \text{change} \rangle$

Sample 3: ∵ It is -ve sample, so find S can't be applied.

But s_2 hyp is wrongly classifying the 3rd sample.

∴ Hence the hypothesis s_2 is not consistent.

∴ we use disjunction hypothesis.

i.e. $\{ \langle \text{sunny}, ?, ?, ?, ?, ? \rangle \vee \langle \text{cloudy}, ?, ?, ?, ?, ? \rangle \}$

Now this is a consistent hyp. ∵ it satisfies the given OR operation whole data set.

Inductive bias: Inductive bias of L ^{learning alg.} is any minimal set of assertions 'B' such that for any Target concept 'c' & corresponding Training Samples D_c ,

$$(\forall x_i \in x) [(B \cap D_c \wedge x_i) \vdash L(x_i, D_c)]$$

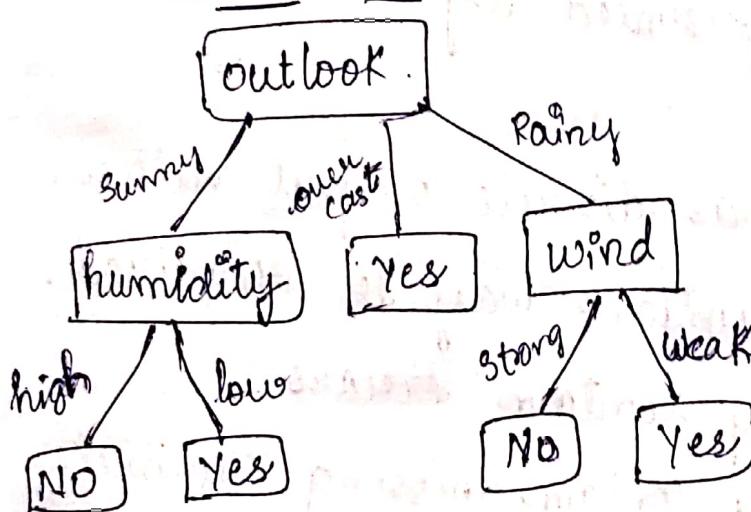
set of all training samples set of assertions

NOTE: $y \vdash z$ means $\Rightarrow z$ is provable by using y .

MODULE - 2 DECISION TREES.

- Decision Tree is also a supervised learner.
- Definition: Decision tree is a classifier which consists of nodes, where each node is representing an attribute and branch represents test condition and leaf node represents the outcome.

Decision tree:



depicting whether
a person can
play game or
not.

Hypotheses for above decision tree:

$$\langle \text{outlook} = \text{sunny} \wedge \text{humidity} = \text{low} \rangle \vee \langle \text{outlook} = \text{overcast} \rangle \\ \vee \langle \text{outlook} = \text{Rainy} \wedge \text{wind} = \text{weak} \rangle$$

formal definition: Decision tree learning is a method of approximating ~~discrete~~^(available kind) value target function in which learned functions is represented by decision tree.

*Appropriate problems for decision tree learning

- 1) Instances are represented by attribute value pairs.
- 2) Target function has discrete output values.
- 3) Disjunctive descriptions may be required.
- 4) Training data may contain errors.
- 5) Training data may contain missing attribute values.

*Data preprocessing: process of data cleaning, finding values of missed attributes and Error handling.

Feature extraction: Extracting the required feature among various features.

8/20/2020

\rightarrow ID₃ algorithm is used.

Entropy: It characterizes the impurity of (randomness in data) arbitrary collection of examples (samples).

high entropy \Rightarrow high impurity \rightarrow (means not a best classifier).

Formula:

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

P_+ \Rightarrow probability of +ve sample.

P_- \Rightarrow probability of -ve sample.

* Entropy: is the measure of randomness.

* Information Gain: Effectiveness of attribute in classifying the data (mainly used for splitting the nodes).

It is an estimation of reduction in entropy w.r.t that attribute.

(For every attribute, information gain is calculated. The attribute which has highest information gain will be selected as the root node).

Formula:

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

A \rightarrow attribute S \rightarrow entire sample.

(i.e either +ve samples or -ve samples)

NOTE:

i) If all samples belong to the same class, then entropy is 0, (stopping condition in decision tree).

2) Equal no. of the +ve & -ve samples, then
Entropy = 1

3) Unequal no. of the +ve & -ve samples, then
entropy lies b/w 0 & 1.

Q) Construction of decision tree: (using ID3 algorithm)

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D ₁	Sunny	Hot	High	Weak	No
D ₂	Sunny	Hot	High	Strong	No
D ₃	Overcast	Hot	High	Weak	Yes
D ₄	Rain	Mild	High	Weak	Yes
D ₅	Rain	Cool	Normal	Weak	Yes
D ₆	Rain	Cool	Normal	Strong	No
D ₇	Overcast	Cool	Normal	Strong	Yes
D ₈	Sunny	Mild	High	Weak	No
D ₉	Sunny	Cool	Normal	Weak	Yes
D ₁₀	Rain	Mild	Normal	Weak	Yes
D ₁₁	Sunny	Mild	Normal	Strong	Yes
D ₁₂	Overcast	Mild	High	Strong	Yes
D ₁₃	Overcast	Hot	Normal	Weak	Yes
D ₁₄	Rain	Mild	High	Strong	No

Step-1: Entropy of entire sample i.e. Entropy (S) is to be calculated.
Step-2: Find out the root node

→ To find the root node, first we calculate Information gains of each attribute, and the attribute which is having high Information gains is taken as root

node:

Information Gain:

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

$$\rightarrow Entropy(S) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

$$= -\frac{9}{14} \times \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) = 0.9403$$

$$Entropy(S) = 0.9403 \quad (\text{remains same for all attributes})$$

(attribute) (values)

i) For $A = \text{outlook}$, $v = \text{sunny, overcast, rain}$.
probability of sunny in total sample

$$Gain(S, \text{outlook}) = Entropy(S) - \left[\frac{5}{14} \underset{\substack{\text{prob. of overcast}}}{\text{Entropy(outlook, sunny)}} + \frac{4}{14} \underset{\substack{\text{prob. of overcast}}}{\text{Entropy(outlook, overcast)}} + \frac{5}{14} \underset{\substack{\text{prob. of rain}}}{\text{Entropy(outlook, rain)}} \right]$$

$$\rightarrow Gain(S, \text{outlook}) = 0.9403 - \left[\frac{5}{14} \underset{\substack{\text{prob. of overcast}}}{\text{Entropy(outlook, sunny)}} + \frac{4}{14} \underset{\substack{\text{prob. of overcast}}}{\text{Entropy(outlook, overcast)}} + \frac{5}{14} \underset{\substack{\text{prob. of rain}}}{\text{Entropy(outlook, rain)}} \right]$$

$$i) Entropy(\text{outlook, sunny}) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

$$= -\frac{2}{5} \times \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \times \log_2 \left(\frac{3}{5}\right) = 0.9710$$

$$ii) Entropy(\text{outlook, overcast}) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

\Rightarrow Here since, all the samples (overcast) belong to same class (i.e. +ve samples), thus entropy is '0'.

$$iii) Entropy(\text{outlook, rain}) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

$$= -\frac{3}{5} \times \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \times \log_2 \left(\frac{2}{5}\right) = 0.9710$$

$$\therefore Gain(S, \text{outlook}) = 0.9403 - \left[\frac{5}{14}(0.9710) + \frac{4}{14}(0) + \frac{5}{14}(0.9710) \right]$$

$$\therefore \text{Gain}(S, \text{outlook}) = 0.2467$$

③ A = Temperature, v = hot, mild, cold

$$\begin{aligned} \rightarrow \text{Gain}(S, \text{Temperature}) &= \text{Entropy}(S) - \left[\frac{4}{14} \text{Entropy}(\text{Temperature}_{\text{hot}}) + \right. \\ &\quad \left. + \frac{6}{14} \text{Entropy}(\text{Temperature}_{\text{mild}}) + \frac{4}{14} \text{Entropy}(\text{Temperature}_{\text{cold}}) \right] \end{aligned}$$

$$\text{i) } \text{Entropy}(\text{Temperature}_{\text{hot}}) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

$$= -\frac{2}{4} \log_2 \left(\frac{2}{4} \right) = \frac{2}{4} \log_2 \left(\frac{1}{2} \right) \Rightarrow \text{they are equal}$$

no. of the 6s - 1s samples, so Entropy = 1

$$\text{ii) } \text{Entropy}(\text{Temperature}_{\text{mild}}) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

$$= -\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) = \boxed{0.9183}$$

$$\text{iii) } \text{Entropy}(\text{Temperature}_{\text{cold}}) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

$$= -\frac{8}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) = \boxed{0.5113}$$

$$\therefore \text{Gain}(S, \text{Temperature}) = 0.9403 - \left[\frac{4}{14}(1) + \frac{6}{14}(0.9183) + \frac{4}{14}(0.5113) \right]$$

$$\boxed{\therefore \text{Gain}(S, \text{Temperature}) = 0.0292}$$

④ A = Humidity, v = High, Normal.

$$\text{Gain}(S, \text{Humidity}) = \text{Entropy}(S) - \left[\frac{7}{14} \text{Entropy}(\text{Humidity}_{\text{High}}) + \right. \\ \left. + \frac{7}{14} \text{Entropy}(\text{Humidity}_{\text{Normal}}) \right]$$

$$\text{i) } \text{Entropy}(\text{Humidity}_{\text{High}}) = -P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$$

$$= -\frac{3}{7} \log_2 \left(\frac{3}{7} \right) - \frac{4}{7} \log_2 \left(\frac{4}{7} \right) = \boxed{0.4852}$$

i) Entropy(Humidity Normal) = $-P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$

$$= -\frac{6}{7} \log_2 \left(\frac{6}{7}\right) - \frac{1}{7} \log_2 \left(\frac{1}{7}\right) = [0.5917]$$

$$\therefore \text{Gain}(S, \text{Humidity}) = 0.9403 - \left[\frac{8}{14} (0.985^2) + \frac{7}{14} (0.5917) \right]$$

$$\therefore \text{Gain}(S, \text{Humidity}) = 0.1519$$

ii) For A = Wind, v = weak, strong.

$$\text{Gain}(S, \text{Wind}) = \text{Entropy}(S) - \left[\frac{8}{14} \text{Entropy}(\text{Wind weak}) + \frac{6}{14} \text{Entropy}(\text{Wind strong}) \right]$$

i) Entropy(Wind weak) = $-P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$

$$= -\frac{6}{8} \times \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \times \log_2 \left(\frac{2}{8}\right) = [0.8113]$$

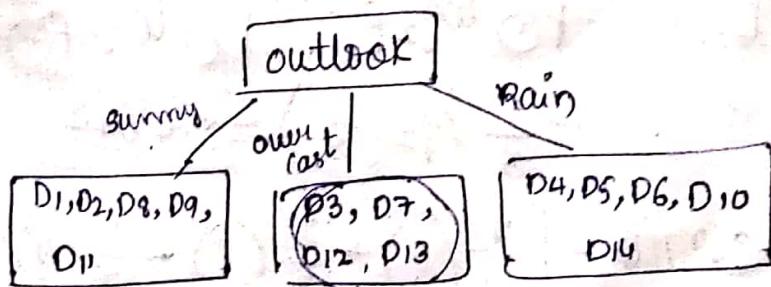
ii) Entropy(Wind strong) = $-P_{\oplus} \log_2 P_{\oplus} - P_{\ominus} \log_2 P_{\ominus}$

~~$\frac{3}{6} \log$~~ [Here equal no. of the ~~class~~ samples, so Entropy is 1] //

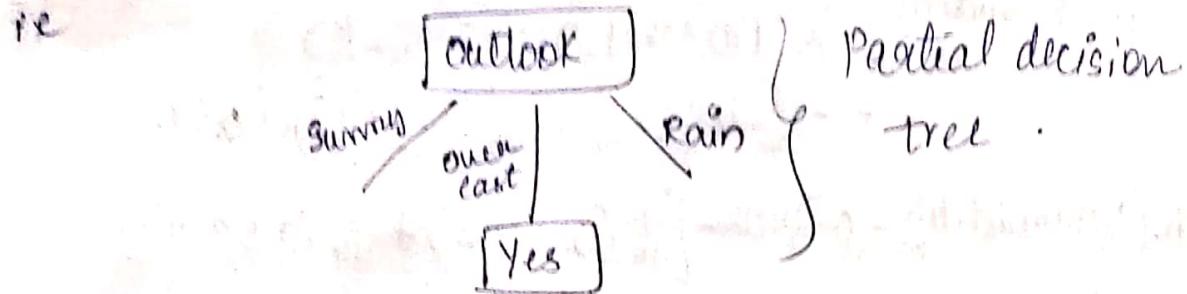
$$\therefore \text{Gain}(S, \text{Wind}) = 0.9403 - \left[\frac{8}{14} (0.8113) + \frac{6}{14} (1) \right] = [0.0481]$$

$$\therefore \text{Gain}(S, \text{Wind}) = 0.0481$$

∴ Among four attributes, outlook attribute is having the highest Information Gain value. Hence outlook is taken as Root Node.



All four are (ie the class)
of same class.
Hence stopping condition (acc to NOTE(i))
∴ Hence it becomes leaf node.



→ STEP 3

NOW for left side i.e. for $D_1, D_2, D_8, D_9, D_{11}$.

Find the entropy for total sample $(D_1, D_2, D_8, D_9, D_{11})$.

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

$$= -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) = [0.9710]$$

→ Now find information Gain for attributes Temperature, Humidity, Wind. (because Outlook is already a root node).

i) $A = \text{Temperature}, V = \text{hot, cold, mild}$.

$$\text{Gain}(S, \text{Temperature}) = \text{Entropy}(S) - \left[\frac{2}{5} \text{Entropy}(\text{Temp hot}) + \frac{2}{6} \text{Entropy}(\text{Temp mild}) + \frac{1}{5} \text{Entropy}(\text{Temp cold}) \right]$$

i) $\text{Entropy}(\text{Temp hot}) = 0$ [because all samples (hot) belong to same class (-ve)]

ii) $\text{Entropy}(\text{Temp mild}) = 1$ [because eq no. of true & -ve samples]

iii) $\text{Entropy}(\text{Temp cool}) = -P_+ \log_2 P_+ - P_- \log_2 P_-$

$$= -\frac{1}{5} \log_2 \left(\frac{1}{5}\right) = [0.4644] = 0$$

[because only one class is true]

$$\text{Gain}(S, \text{Temp}) = 0.9710 - \left[\frac{2}{5}(0) + \frac{2}{5}(1) + \frac{1}{5}(0) \right]$$

$$\Rightarrow \boxed{\text{Gain}(S, \text{Temp}) = 0.571}$$

② A = Humidity; v = high, normal

$$\text{Gain}(S, \text{Humidity}) = \text{Entropy}(S) - \left[\frac{3}{5} \text{Entropy}(\text{Humidity high}) + \frac{2}{5} \text{Entropy}(\text{Humidity normal}) \right]$$
$$= 0.971 - \left[\frac{3}{5} (0) + \frac{2}{5} (0) \right]$$

$$\boxed{\text{Gain}(S, \text{Humidity}) = 0.9710}$$

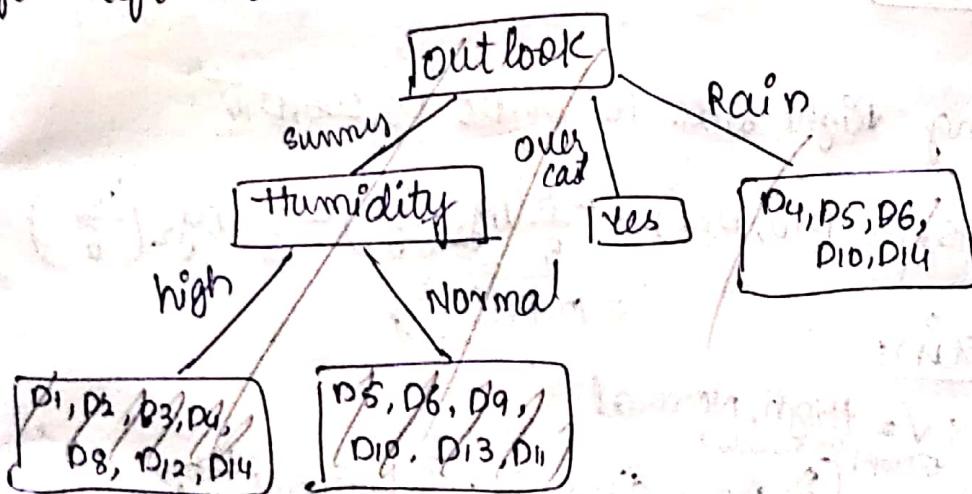
③ A = wind, v = weak, strong

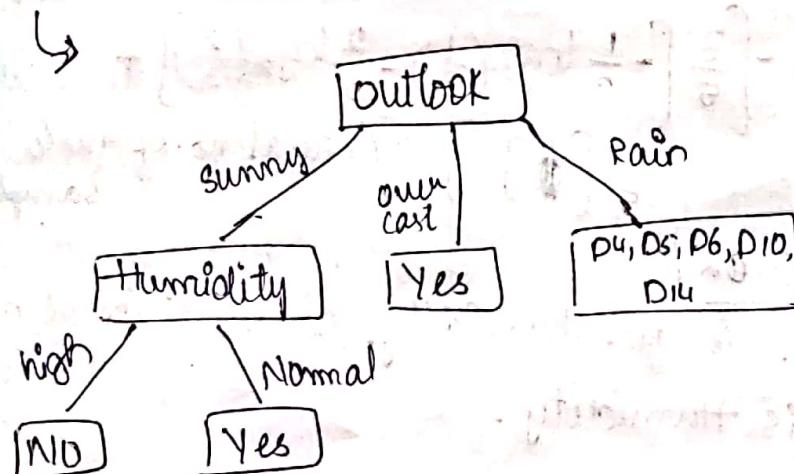
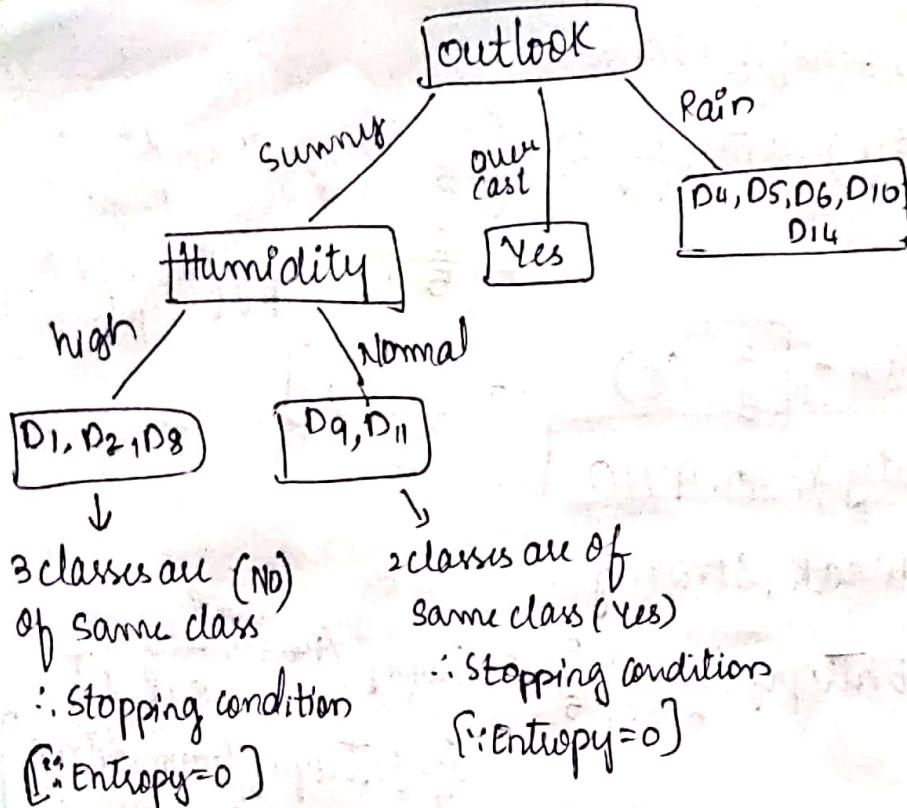
$$\text{Gain}(S, \text{Wind}) = \text{Entropy}(S) - \left[\frac{3}{5} \text{Entropy}(\text{Humidity weak}) + \frac{2}{5} \text{Entropy}(\text{Humidity strong}) \right]$$
$$= 0.971 - \left[\frac{3}{5} \left[-\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] + \frac{2}{5} [0] \right]$$

\rightarrow equal no. of true & false samples

$$\boxed{\text{Gain}(S, \text{Wind}) = 0.02}$$

∴ Highest Gain is for humidity. so, next (root node)
for left side is humidity.





Step 4: Considering right side ie D₄, D₅, D₆, D₁₀, D₁₄

$$\text{Entropy}(D_4, D_5, D_6, D_{10}, D_{14}) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$

$$= 0.9710 \dots$$

* Information Gain:

(D_A) = Humidity ; V = High, Normal.

$$IG_1(S, \text{Humidity}) = 1 - \left[\frac{2}{5} \text{Entropy}(\text{Hum. High}) + \frac{3}{5} \text{Entropy}(\text{Hum. Normal}) \right]$$

$$IG_1(S, \text{Humidity}) = 1 - \left[\frac{2}{5} (1) + \frac{3}{5} \text{Entropy}(\text{Hum. Normal}) \right]$$

$$\text{Entropy}(\text{Hum. Normal}) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.9183$$

$$IG_1(S, \text{Humidity}) = 0.0200$$

④ $N = \text{Temp}$; $V = \text{Hot, mild, cool}$

$$IG(S, \text{Temp}) = H(6, 10, 14) - \left[\frac{3}{5} E(\text{Temp, mild}) + \frac{2}{5} E(\text{Temp, cool}) \right]$$

$$E(\text{Temp, mild}) = -\frac{2}{3} \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) = 0.9183$$

$$E(\text{Temp, cool}) = 1 \quad (\text{equal no. of samples +ve & -ve})$$

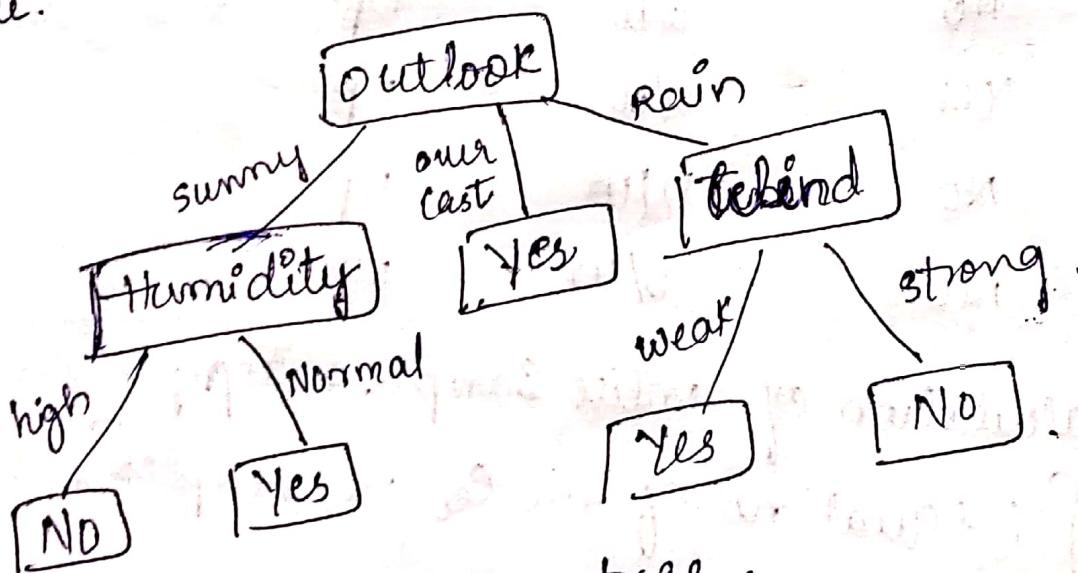
$$\boxed{IG(S, \text{Temp}) = 0.02}$$

$$⑤ IG(S, \text{wind}) = 0.9710 - \left[\frac{3}{5} E(\text{wind, weak}) + \frac{2}{5} E(\text{wind, strong}) \right]$$

$$IG(S, \text{wind}) = 0.9710 - \left[\frac{3}{5}(0) + \frac{2}{5}(0) \right]$$

$$\boxed{IG(S, \text{wind}) = 0.9710}$$

Here IG_1 is more for wind; so it is taken as root node.



final decision tree

21/20

2) Construction of Decision Tree:

	Age	Competition	Type	Profit
1.	old	Yes	s/w	down
2.	old	No	s/w	down
3.	old	No	h/w	down
4.	mild	Yes	s/w	down
5.	mild	Yes	h/w	down
6.	mild	No	h/w	up
7.	mild	No	s/w	up
8.	new	Yes	s/w	up
9.	new	No	h/w	up
10.	new	No	s/w	up

Sol: Step-1: Calculation of entire Sample entropy

$$E(S) = 1 \left\{ \because \text{equal no. of true Neg. examples} \right\}$$

Step-2: Information Gain: Calculation

$$\begin{aligned} \text{IG}(S, \text{Age}) &= 1 - \left[\frac{3}{10} \text{Entropy}(\text{Age}_{\text{old}}) + \frac{4}{10} E(\text{Age}_{\text{mild}}) + \frac{3}{10} E(\text{Age}_{\text{new}}) \right] \\ &= 1 - \left[\frac{3}{10} (0) + \frac{4}{10} E(\text{Age}_{\text{mild}}) + \frac{3}{10} (0) \right] \end{aligned}$$

$$E(\text{Age}_{\text{mild}}) = -\frac{2}{3} \left[\log_2 \frac{2}{3} \right] - \frac{1}{3} \log_2 \frac{1}{3} = 0.9183.$$

$$\text{IG}_1(S, \text{Age}) = 1 - \left[\frac{3}{10} (0) + \frac{4}{10} (0.9183) + \frac{3}{10} (0) \right]$$

$$\boxed{\text{IG}_1(S, \text{Age}) = 0.6.}$$

for A = competition ; v = Yes, No

$$IG(S, \text{competition}) = 1 - \left[\frac{4}{10} \text{Entropy}(\text{comp}_{\text{Yes}}) + \frac{6}{10} E(\text{comp}_{\text{No}}) \right].$$

$$\text{Entropy}(\text{comp}_{\text{Yes}}) = -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{3}{4} \log_2 \left(\frac{3}{4} \right) = 0.8113.$$

$$\text{Entropy}(\text{comp}_{\text{No}}) = -\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) = 0.9183$$

$$IG(S, \text{comp}) = 0.1245.$$

for A = Type ; v = S1W, H1W

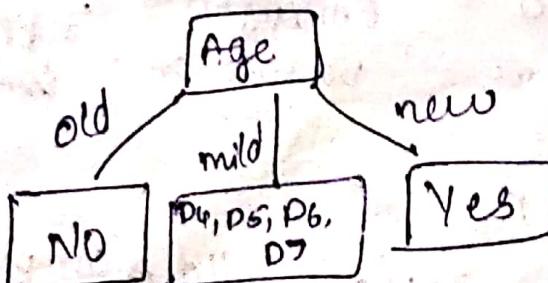
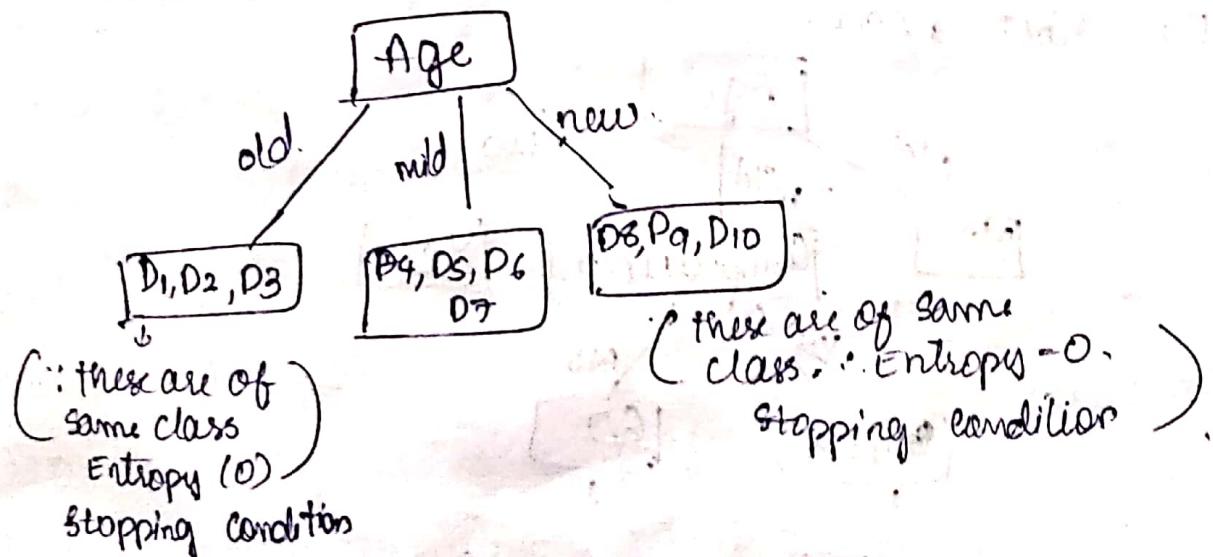
$$IG(S, \text{Type}) = 1 - \left[\frac{6}{10} E(\text{Type}_{S1W}) + \frac{4}{10} E(\text{Type}_{H1W}) \right].$$

$$E(\text{Type}_{S1W}) = 0 \quad (\because \text{equal no. of samples (one G - one H)})$$

$$E(\text{Type}_{H1W}) = 0 \quad (\because \text{equal no. of samples (one G - one H)})$$

$$IG(S, \text{Type}) = 0$$

∴ highest Information Gain is for Age attribute. Hence it becomes the root node.



Step 3: for the instances D_4, D_5, D_6, D_7 .

Calculate Entropy $\Rightarrow E(4, 5, 6, 7) = 1$ [P: equal no. of +ve & -ve samples]

Now calculate IG_1 for Competition, C₁ Type.

① for A₁ = competition : V = Yes, No [for (D_4, D_5, D_6, D_7)]

$$IG_1(S, \text{competition}) = 1 - \left[\frac{2}{4} E(\text{competition})_{\text{Yes}} + \frac{2}{4} E(\text{competition})_{\text{No}} \right]$$

$$IE_1(S, \text{comp}) = 1 - \left[\frac{2}{4}(1) + \frac{2}{4}(1) \right]$$

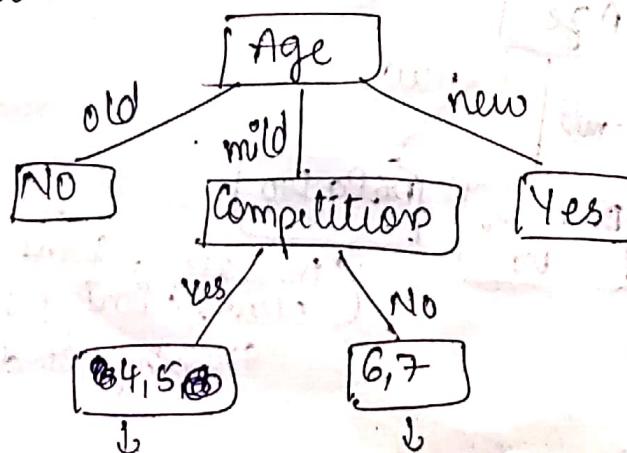
$$\boxed{IG_1(S, \text{comp}) = 1}$$

② for A₂ Type : V = H/W, S/W.

$$IG_1(S, \text{Type}) = 1 - \left[\frac{2}{4} E(\text{Type}_{H/W}) + \frac{2}{4} E(\text{Type}_{S/W}) \right]$$
$$= 1 - \left[\frac{2}{4}(1) + \frac{2}{4}(1) \right]$$

$$\boxed{IG_1(S, \text{Type}) = 0}$$

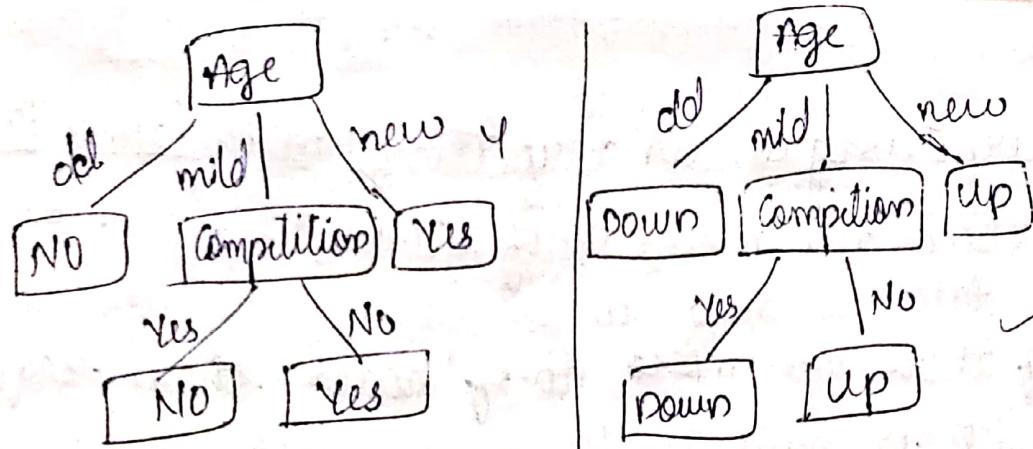
\therefore Highest Information Gain is for Competition. Hence it is the root node.



Entropy = 0
because same class
(stopping cond)

Entropy = 0 (stopping)
because same class
(cond)

\therefore Hypothesis = $(\text{age} = \text{new}) \vee (\text{age} = \text{mild} \wedge \text{competition} = \text{No})$



final decision tree

* Hypothesis Space Search in Decision Tree Learning

- * In decision tree: (bias means preference / constraint)
- 1) One of the bias is selecting simple tree rather than all the other complex tree (ie a preference).
- 2) Another bias is using information gain - to select best possible node (ie here the bias is a constraint).

* Hypothesis Space Search in Decision Tree Learning

- 1) Stereotype Dicotomy: ID3 is complete space for finite discrete valued functions relative to available attributes.

- 2) ID3 maintains only single current hypothesis as it searches through space of decision trees.
- 3) ID3 in its pure form performs no back tracking.
- 4) ID3 uses all training examples at each step (information gain) in search to make statistically based decisions regarding how to refine current hypothesis

Q1/20

Overfitting & Underfitting: in Decision Tree:

If there are more no. of attributes in a sample, then the data is said to be overfitting the model.

→ If there are more no. of nodes than required (which is said to be more specific), then it is called as overfitting.

If there are less no. of attributes in a sample, than required, then the data is said to be underfitting the model. ∵ If there are less no. of nodes/attributes (which is said to be more generic), then it is called underfitting.

Q1/20

* Inductive bias in Decision tree:

1) Restriction Bias: Categorical restriction on set

* of hypothesis considered.

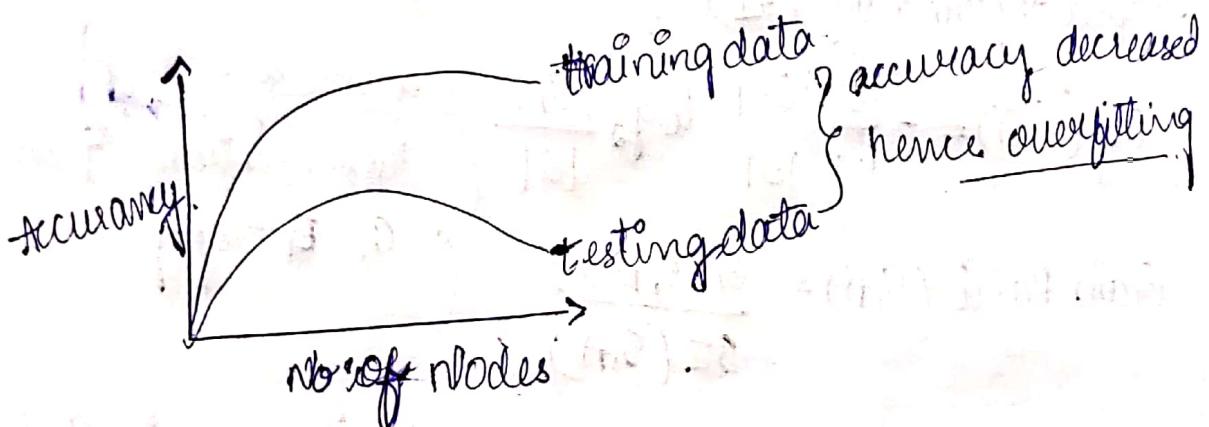
2) Preference Bias: Preference for certain hypotheses over others. (In decision tree → preferring ^{shorter} trees i.e. less no. of nodes).

→ most of the ML models use only any one of the biases.

* Issues in Decision tree:

① Avoid overfitting in the data.

Overfitting: Given a hypothesis space (H). A hypothesis (h) belongs to H ($h \in H$) is said to overfit the training data if there exists some alternative hypothesis $h' \notin H$ such that h has smaller error than h' over training examples, but h' has smaller error than h over entire distribution of instances.



* How to Avoid Overfitting: The method is - Reducing error pruning

i) Reduce error pruning (Reducing the no. of nodes in decision tree)

 | Post pruning (Simpler but memory wastage)

 | Pre pruning (Cost-efficient)

ii) Rule Post pruning:

1) Infer decision tree from training set until the tree fits the data.

2) Convert the learned tree into equivalent set of

rules.

B) Prune the rules by removing any pre-condition that improves accuracy.

④ Sort and value the pruned rules by their estimated accuracy.

⑤ Incorporating continuous valid valued attributes.
(like converting continuous values to discrete).

⑥ Handling training examples with missing attribute values.

⑦ Alternative measures for selecting attributes:

→ split information (SI)

$$SI(S, A) = - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

alternatives for
Information Gain
& Entropy.

$$\text{Gain Ratio}(S, A) = \frac{G(S, A)}{SI(S, A)}$$

⑧ Handling the attributes with different cost:

$$\frac{\text{Gain}(S, A)}{\text{Cost}(A)}$$

Time required for finding the value of
an attribute.

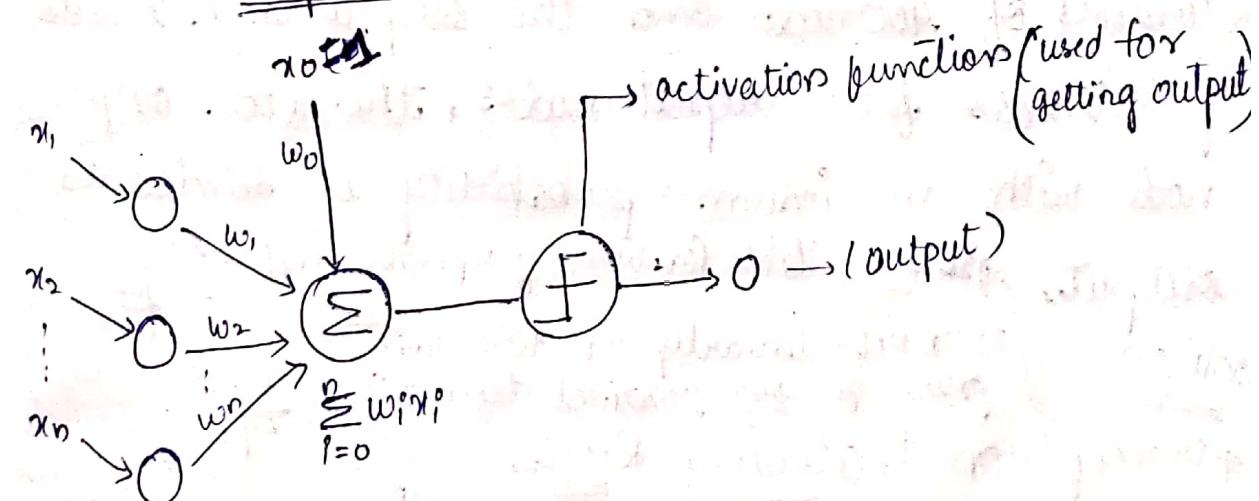
* Use of Bias: It is used to fit the training data in constructing the model.

Artificial Neural Networks (ANN)

* Problems suited for Neural Networks:

- 1) Instances are represented by many attribute value pairs.
- 2) Target functions output may be discrete, real valued, vector of several real valued (or) discrete valued attributes.
- 3) Training examples may contain errors.
- 4) Long training times are acceptable.
- 5) Fast evaluation of learned target function may be required.
- 6) Ability of humans to understand learned target function is not important.

* Single Perceptron (one type of Neural Network)



$$\text{output}(O) = \begin{cases} 1 & \text{if } \sum_{i=0}^n w_i x_i^0 > 0 \\ 0 & \text{otherwise} \end{cases}$$

→ Single Perceptron is used for ~~also~~ representing single output. Output depends on activation function.

→ Perception is used for performing logical operations
(Best suitable for logical op. except for XOR)

Example: (i) Perform AND operation on

$T=1$ if weights w_1, w_2 are 0.5 and $w_0=-0.8$
 $F=-1$

for performing we use $\sum_{i=1}^2 x_i w_i$

fix $x_0=1$ (always) inputs

$$= x_0 w_0 + x_1 w_1 + x_2 w_2 \\ = 1(-0.8) + (1)(0.5) + (-1)(0.5) \Rightarrow -0.8 + 0.5 - 0.5 \\ = -0.8$$

i.e if $\sum_{i=1}^2 x_i w_i = -1$, if $\sum_{i=1}^2 x_i w_i < 0$. i.e F

Representation of Neural Network:

Hidden layer: performs all the computations on the input data to be processed.

→ If the image size is 30×30 , the input layer consists of 900 nodes and the output layer consists of few output nodes, the node output node with maximum probability is selected as output. ~~It is suitable for linearly separable unit~~ → $\begin{array}{c} + \\ + \\ - \\ - \end{array}$

Activation: after not-linearly separable unit

we go for gradient descent → $\begin{array}{c} + \\ - \\ + \\ - \end{array}$

Perception Learning Rules:

$$w_i \leftarrow w_i + \Delta w_i \quad \text{target off} \quad \text{output given by model}$$

where as $\Delta w_i \leftarrow \eta(t-O)x_i$ inputs.

Perception

constant learning rate.

→ This learning rule is best suitable for linearly separable unit.

Gradient descent or delta Rule:

In gradient descent approach, we repeatedly calculate error which is given by -

$$\text{error function } E(\bar{w}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$t_d \rightarrow$ target O/P of a sample $D \rightarrow$ Data set
 $o_d \rightarrow$ O/P of a sample given by $d \rightarrow$ data sample
 our model (actual O/P)

$$\Delta w_i^o = -\eta \frac{\partial E}{\partial w_i}, \quad w_i^o \leftarrow w_i^o + \Delta w_i^o$$

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

$$\Rightarrow \text{then } \frac{\partial E}{\partial w} = \frac{1}{2} \frac{\partial}{\partial w} \left(\sum_{d \in D} (t_d - o_d)^2 \right) \quad (\because o_d = w_i^o x_i)$$

$$= \frac{1}{2} \times \sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \quad \frac{\partial}{\partial w} (t_d - o_d)$$

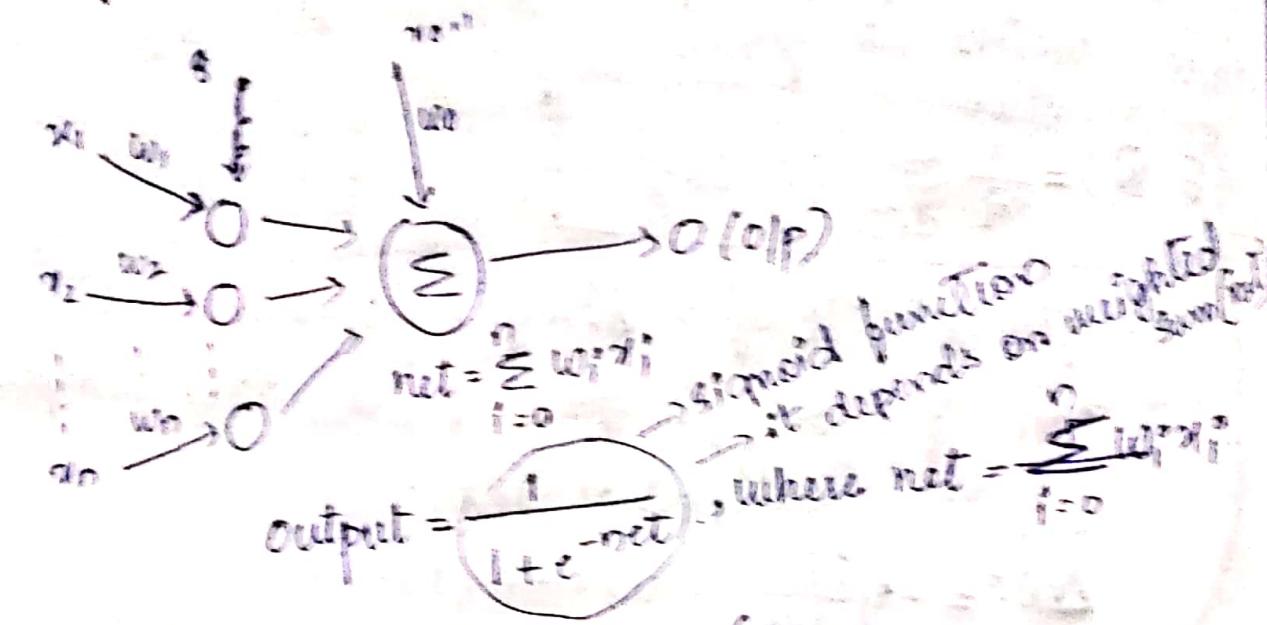
$$= \sum_{d \in D} (t_d - o_d) (-x_i) \quad = 0 - \frac{\partial}{\partial w} (w_i^o x_i)$$

$$= -x_i$$

$$\Delta w_i^o = \eta \sum_{d \in D} (t_d - o_d) x_i^o$$

whl20

* Sigmoid: It is an activation function.



→ here output lies bw 0 & 1

NOTE: In neural networks, output at a particular neuron, is considered as the input to the next layer neuron.

* Back Propagation Algorithm:

- 1) Create a feed forward network with 'n' inputs, ' n_h ' hidden units and 'o' output units.
- 2) Initialize all network weights to small random numbers.
Example: Between -0.5 and +0.5.
- 3) Until the termination condition is met, perform the following steps.
 - For each (\vec{x}, \vec{t}) in training examples
 - i) Input the instance \vec{x} to the network and compute the output of every unit in the network.

ii) For each network output unit 'K' calculate its error δ_K .

$$\boxed{\delta_K = o_K (1 - o_K) (t_K - o_K)}$$

iii) For each hidden unit calculate its error.

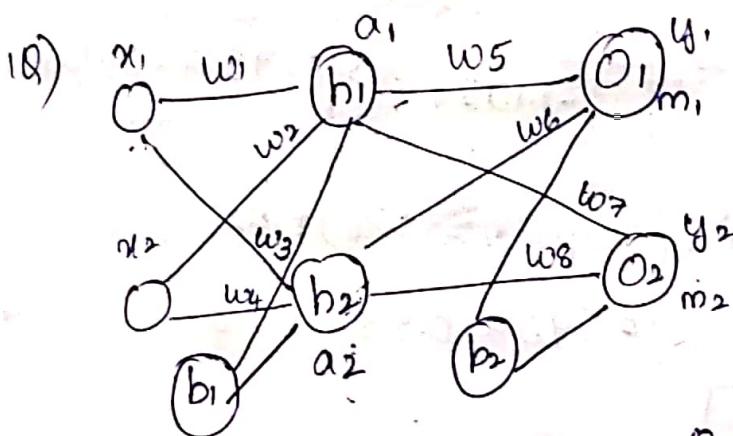
$$\boxed{\delta_h = o_h (1 - o_h) \sum_{K \in \text{output}} w_{Kh} \delta_K}$$

o_h = output at 'h'

Update each network weight

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

$$\Delta w_{ji} \leftarrow \eta \delta_j x_{ji}$$



$$x_1 = 0.05$$

$$b_1 = 0.35$$

$$o_1 = 0.01$$

$$x_2 = 0.10$$

$$b_2 = 0.60$$

$$o_2 = 0.99$$

$$w_1 = 0.15$$

$$w_5 = 0.40$$

$$\eta = 0.6$$

$$w_2 = 0.20$$

$$w_6 = 0.45$$

$$w_3 = 0.25$$

$$w_7 = 0.50$$

$$w_4 = 0.30$$

$$w_8 = 0.55$$

Sol: i) $h_{\text{net}} = x_1 w_1 + x_2 w_2 + b_1$

$$= (0.05)(0.15) + (0.10)(0.20) + 0.35$$

$$\boxed{h_{\text{net}} = 0.3775}$$

$$h_1 \text{Output} = \frac{1}{1 + e^{-\text{hnet}}} = \frac{1}{1 + e^{-0.3775}} \Rightarrow 0.59337 \Rightarrow h_1 \text{out}$$

$$\rightarrow \text{net } h_2 = x_1 w_3 + x_2 w_4 + b_1 = (0.05 \times 0.25) + (0.10 \times 0.30) + 0.35$$

| net $h_2 = 0.3925$

$$\text{Output } h_2 = \frac{1}{1 + e^{-\text{net } h_2}} = \frac{1}{1 + e^{-0.3925}} \Rightarrow 0.5969 = \text{output } h_2$$

$$\begin{aligned} \rightarrow \text{net } O_1 &= (\text{output } h_1 \times w_5) + (\text{output } h_2 \times w_6) + b_2 \\ &= (0.5933 \times 0.4) + (0.5969 \times 0.45) + 0.60 \end{aligned}$$

| net $O_1 = 1.1059$

$$\text{output } O_1 = \frac{1}{1 + e^{-\text{net } O_1}} = \frac{1}{1 + e^{-1.1059}} \Rightarrow 0.7514 = \text{output } O_1$$

$$\begin{aligned} \rightarrow \text{net } O_2 &= (\text{output } h_2 \times w_7) + (\text{output } h_2 \times w_8) + b_2 \\ &= (0.5933 \times 0.5) + (0.5969 \times 0.55) + 0.6 \end{aligned}$$

| net $O_2 = 1.2249$

$$\text{output } O_2 = \frac{1}{1 + e^{-\text{net } O_2}} = \frac{1}{1 + e^{-1.2249}} \Rightarrow 0.7729 = \text{output } O_2$$

$$\Rightarrow \text{actual outputs: } O_1 = 0.7514$$

$$O_2 = 0.7729$$

But given target O/p's = $O_1 = 0.01$
 $O_2 = 0.99$

∴ checking the error values.

E_{O1} = Error at output 1 (O_1)

$$= \frac{1}{2} (t_d - o_d)^2 = \frac{1}{2} (0.01 - 0.7514)^2 = 0.2748$$

$$\text{similarly } E_{O_2} = \frac{1}{2}(0.99 - 0.7729)^2 = 0.0235$$

since the error values are not approx to '0', we try to adjust the weights (back propagation).

$$\frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial \text{out } O_1} \times \frac{\partial \text{out } O_1}{\partial \text{net } O_1} \times \frac{\partial \text{net } O_1}{\partial w_5}$$

By some derivations $\Rightarrow \frac{\partial E_{\text{total}}}{\partial \text{out } O_1} = \text{out } O_1 - \text{Target } O_1 = 0.7414$

$$\frac{\partial \text{out } O_1}{\partial \text{net } O_1} = \text{out } O_1(1 - \text{out } O_1) = 0.1868$$

$$\frac{\partial \text{net } O_1}{\partial w_5} = \text{out } h_1 = 0.5969$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = 0.7414 \times 0.1868 \times 0.5969 = 0.0826$$

23/1/20

* Remarks on Back Propagation Algorithm:

1) Convergence and local minima -

In order to converge, add momentum in every iteration.

$$\Delta w_{ji} = \eta \delta_j x_{ji} + \alpha \Delta w_{ji}(n-1); \text{ where } 0 < \alpha < 1$$

$n = \text{no. of iterations}$
 $(n^{\text{th}} \text{ iteration})$

2) Representational power of feed forward network.

3) Hypothesis Space

4) Generalization and overfitting.

* Evaluating Hypotheses:

Sample Error: Denoted by $\text{error}(h)$. of hypothesis

h w.r.t target function f and data sample 'S' is

$$\text{error}(h) = \frac{1}{n} \sum_{x \in S} \delta(f(x), h(x))$$

where,

$n = \text{no. of samples in } S.$

$f(x) = \text{target function.}$

$h(x) = \text{actual O/P given by model.}$

$\delta(f(x), h(x)) = 1, \text{ if } f(x) \neq h(x); \text{ else } 0.$

a) True Error: Denoted by $\text{error}(h)$ of hypothesis h w.r.t target function ' f ' and distribution ' D ', probability that ' h ' will misclassify an instance drawn at random according ' D '.

$$\Rightarrow \text{error}(h) = \text{probability}_{D, x \in D} [f(x) \neq h(x)]$$

~~aaalilao~~

$$\text{error}(h) = \frac{\text{error}(h)}{s} \pm z_n$$

$$\sqrt{\frac{\text{error}(h) * (1 - \text{error}(h))}{n}}$$

*Thumb rule for Estimation: Error & here Estimator is Sample Error.

$$n \frac{\text{error}(h)}{s} (1 - \frac{\text{error}(h)}{s}) \geq 5.$$

*Definitions:

① A random variable can be viewed as name of an experiment with a probabilistic outcome. It's values is the outcome of the experiment.

② Probability distributions for a random variable by specifies the probability $\Pr(Y=y_i)$ that Y will take that on the value y_i for each possible value of y_i .

③ Expected value (or) mean of a random variable

$$y \text{ is } E(Y) = \sum y_i p_r(y=y_i)$$

④ Variance of random variable is the width or dispersion of the distribution about its mean.

⑤ Central limit theorem:

It is the theorem stating that, sum of large no. of distributed random variables approximately follow normal distribution.

⑥ Estimator: A random variable that estimates a random variable 'P'.

⑦ n% confidence interval, estimate for parameter 'p' is an interval that includes 'p' with n% of probability.

* General approach for determining confidence intervals

① Identify the underline population 'P' to be estimated. Eg: True Error

② Define the estimator 'y': \rightarrow E_y: Sample Error.

③ Determine the probability distribution by that the once estimator 'y' including mean μ_y variance.

④ Determine the probability distribution by that the once estimator 'y' including mean μ_y variance.

④ Determine n% confidence interval by finding the threshold values l and u such that D_p falls between l and u

in this case (Normal distribution) ^{continuous data}

i.e if $n \rightarrow \infty$

then Interval is

$$y \pm z_n \sigma$$

y = estimator

z_n = constant

σ is variance.

Confidence interval

ask

* Bayesian learning and Instance based learning,

Bayesian learning:

$$P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

where $P(h|D)$ = post probability

$P(h)$ = prior probability.

$P(D|h)$ = likelihood

$P(D)$ = marginal value.

Q) Find the probability of King in face card;

sol: $P(\text{King}|\text{face card}) = \frac{P(\text{face}|\text{King}) * P(\text{King})}{P(\text{face})}$

posterior probability = $\frac{1 * \frac{4}{52}}{\frac{12}{52}} = \frac{4}{12} = \frac{1}{3}$.

* Features of Bayesian learning:

- 1) Each observed training example can incrementally decrease or increase the estimated probability that a hypothesis is correct.
- 2) Prior knowledge can be combined with observed data to determine final probability of the hypothesis.
- 3) Machine Bayesian methods can accommodate hypothesis that make probabilistic predictions.

4) New instances can be classified by combining the predictions of multiple hypothesis weighted by their probabilities.

Q)

Fruit	yellow	sweet	long	Total
Orange	350	450	0	650
Banana	400	300	350	400
Others	50	100	50	150
total	800	850	400	1200

Ans: The new instance (is) ~~is~~ is

$$t = \{ \text{yellow, sweet, long} \}$$

First we have to find out probability with every fruit attribute and then with fruit.

$$\begin{aligned} ① P(\text{Yellow/orange}) &= \frac{P(\text{Orange}/\text{yellow}) \times P(\text{yellow})}{P(\text{orange})} \\ &= \frac{\frac{350}{800} \times \frac{800}{1200}}{\frac{650}{1200}} = 0.66 // \end{aligned}$$

$$\begin{aligned} P(\text{Sweet/orange}) &= \frac{P(\text{Orange}/\text{sweet}) \times P(\text{Sweet})}{P(\text{orange})} \\ &= \frac{\frac{450}{850} \times \frac{850}{1200}}{\frac{650}{1200}} = 0.90 // \end{aligned}$$

$$P(\text{long/orange}) = 0,$$

$$\text{Now } p(\text{-fruit/orange}) = 0.66 \times 0.90 \times 0 \\ = 0.594$$

$$② p(\text{yellow/banana}) = \frac{p(\text{banana/Yellow}) \times p(\text{Yellow})}{p(\text{banana})} \\ = \frac{\frac{400}{800} \times \frac{800}{1200}}{\frac{400}{1200}} = 1.$$

$$p(\text{Sweet/banana}) = \frac{p(\text{banana/Sweet}) \times p(\text{Sweet})}{p(\text{banana})} \\ = \frac{\frac{300}{800} \times \frac{800}{1200}}{\frac{400}{1200}} = \frac{3}{4}$$

$$p(\text{long/banana}) = \frac{p(\text{banana/long}) \times p(\text{long})}{p(\text{banana})} \\ = \frac{\frac{350}{400} \times \frac{400}{1200}}{\frac{400}{1200}} = \frac{7}{8}$$

$$\text{Now } p(\text{fruit/banana}) = 1 \times \frac{3}{4} \times \frac{7}{8} = 0.65,$$

$$③ p(\text{yellow/others}) = \frac{p(\text{others/Yellow}) \times p(\text{Yellow})}{p(\text{others})} \\ = \frac{\frac{50}{800} \times \frac{800}{1200}}{\frac{100}{1200}} = \frac{1}{3}$$

$$P(\text{Sweet/others}) = \frac{P(\text{others}|\text{sweet}) \times P(\text{sweet})}{P(\text{others})}$$
$$= \frac{\frac{100}{850} \times \frac{350}{1200}}{\frac{150}{1200}} = \frac{1}{3}$$

$$P(\text{long/others}) = \frac{P(\text{others}|\text{long}) \times P(\text{long})}{P(\text{others})}$$
$$= \frac{\frac{50}{400} \times \frac{400}{1200}}{\frac{150}{1200}} = \frac{1}{3}$$

Now ~~p(fruit)~~ $P(\text{fruit/others}) = \frac{1}{3} \times \frac{2}{3} \times \frac{1}{3} = 0.07$
∴ probability of banana is higher.

- Q) Given we have Prior knowledge that over the entire population only 0.008 have cancer.
- * The lab test returns the correct positive results in only 98% of the cases, and the correct negative results in only 97% of the cases.
 - Suppose we now observe the new patient for whom lab test returns the positive result. Should we diagnose the patient or not?

$$\text{Sol: } p(\text{cancer}) = 0.008 \quad p(\text{+}/\sim \text{cancer}) = 0.97$$

$$p(\sim \text{cancer}) = 0.992 \quad p(\text{-}/\sim \text{cancer}) = 0.03$$

$$p(\text{+}/\text{cancer}) = 0.98$$

$$p(\text{-}/\text{cancer}) = 0.02$$

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

$$h_{MAP} = \arg \max_{h \in H} p(\text{+}/\text{cancer}) p(\text{cancer})$$

$$h_{MAP} = \arg \max_{h \in H} p(\text{+}/\sim \text{cancer}) p(\sim \text{cancer})$$

$$\textcircled{1} \quad p(\text{+}/\text{cancer}) p(\text{cancer}) = 0.98 \times 0.008$$

$$= 0.0078$$

$$\textcircled{2} p(+|\sim \text{cancer}) \times p(\sim \text{cancer}) = 0.03 \times 0.992 \\ = 0.0298$$

$h_{MAP} = \sim \text{cancer}$.

we have the 2nd case as high probability
so, we don't need to diagnose the patient.

$$* p(h|D) = \frac{p(D|h) * p(h)}{p(D)}$$

$p(h)$ = prior probability that hypothesis h hold

$p(D|h)$ = probability of data D given some world
in which hypothesis ' h ' holds.

$p(D)$ = probability of D given no knowledge
about which hypothesis holds.

* Example of Text classification:

<u>Text</u>	<u>Category</u>
1) A great game	Sports
2) The election was over	Non Sport
3) Very clean match.	Sports
4) A clean but forgettable game	Sports
5) It was a close election	Non Sports
<u>Qn:</u> $\text{Sen} \mid \text{Sports} = ?$	

$$P(Sem | Non-Sports) = ?$$

$$P(A | Sports) = ?$$

$$P(veery | Sports) = ?$$

$$P(close | Sports) = ?$$

$$P(game | Sports) = ?$$

Total no. of words in sports = 11

Total no. of words in non-sports = 9

Total no. of unique words = 14

* Laplace Smoothing

$$\text{Probability of word} = \frac{\text{Word Count} + 1}{\text{Total no. of words} + \text{no. of unique words}}$$

$$P(A | Sports) = \frac{\text{No. of times word 'A' occurred in Sports} + 1}{\text{Total no. of words in Sports} + \text{unique word}}$$

$$= \frac{2+1}{11+14} = \frac{3}{25} = 0.12$$

$$P(veery | Sports) = \frac{1+1}{11+14} = \frac{2}{25} = 0.08$$

$$P(close | Sports) = \frac{0+1}{11+14} = \frac{1}{25} = 0.04$$

$$P(game | Sports) = \frac{2+1}{11+14} = \frac{3}{25} = 0.12$$

$$= P(A | Sports) * P(veery | Sports) * P(close | Sports) * P(game | Sports)$$

$$= 0.12 \times 0.08 \times 0.04 \times 0.12 \\ = 0.00004608$$

$$P(A/\text{non-sports}) = \frac{1+1}{9+14} = \frac{2}{23} = 0.086.$$

$$P(\text{Very}/\text{non-sports}) = \frac{0+1}{9+14} = \frac{1}{23} = 0.043$$

$$P(\text{Close}/\text{non-Sports}) = \frac{1+1}{9+14} = \frac{2}{23} = 0.086.$$

$$P(\text{game}/\text{non-Sports}) = \frac{0+1}{9+14} = \frac{1}{23} = 0.043$$

$$= P(A/\text{Sports}) \times P(\text{Very}/\text{Sports}) + P(\text{Close}/\text{Sports}) \times P(\text{game})$$

$$= 0.12 \times 0.08 \times 0.04 \times 0.12$$

$$= 0.00004608$$

$$P(A/\text{non-Sports}) = \frac{1+1}{9+14} = \frac{2}{23} = 0.086.$$

$$P(\text{Very}/\text{non-sports}) = \frac{0+1}{9+14} = \frac{1}{23} = 0.043$$

$$P(\text{close}/\text{non-Sports}) = \frac{1+1}{9+14} = \frac{2}{23} = 0.086.$$

$$P(\text{game}/\text{non-Sports}) = \frac{0+1}{9+14} = \frac{1}{23} = 0.043$$

$$= 0.086 \times 0.043 \times 0.086 \times 0.043$$

$$= 0.0000136752$$

This sentence belongs to sports class.

* maximum Posterioric hypothesis: (h_{MAP})

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

$$= \underset{h \in H}{\operatorname{argmax}} \frac{p(D|h) * p(h)}{p(D)}$$

$$p(D) = \text{constant}$$

$$= \underset{h \in H}{\operatorname{argmax}} p(D|h) * p(h)$$

when prior probability $p(h)$ is same for every hypothesis

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

30/1/2020

* Route - Force Bayes Concept learning:

1) For each hypothesis h in H , calculate posterior probability $p(h|D) = \frac{p(D|h) * p(h)}{p(D)}$

2) output the hypothesis h_{MAP} with highest posterior probability $h_{MAP} = \underset{h \in H}{\operatorname{argmax}} p(D|h)$

3) Assumptions in this approach:

i) Training data T is noise free (no error ~~not~~)

-) Target concept c is contained in hypothesis space
) we have no prior reason to believe any hypothesis, is more probable than any other.

$$\text{Explanation: } P(h|D) = \frac{P(D|h) * p(h)}{P(D)}$$

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

if $P(D|h) = 1$ (i.e. when hyp is consistent)
 then $P(h|D) = \frac{1 * p(h)}{P(D)}$

$$p(h) = \text{prob. of a hyp that it classifies the dataset correctly} = \frac{1}{|H|}$$

$$P(D) = \text{prob. of a data set that it holds a hyp} = \frac{|\text{Consistent hyp}|}{|H|}$$

$$(\because \text{set of consistent hyp} = \text{Version Space}) \therefore P(D) = \frac{|\text{VSH}, D|}{|H|}$$

$$\therefore P(D|h) = \frac{1 \times \frac{1}{|H|}}{\frac{|\text{VSH}, D|}{|H|}} = \frac{1}{|\text{VSH}, D|}$$

else if $p(h|D) = 0$ (i.e. for inconsistent hyp)

then $P(D|h) = P(h|D) * p(h) = \frac{0 * p(h)}{P(D)} = 0$

$\boxed{P(D|h) = 0 \text{ (for inconsistent hyp)}}$ & $\boxed{P(D|h) = \frac{1}{|\text{VSH}, D|} \text{ for consistent hyp}}$

~~* maximum likelihood and least square error hypotheses.~~

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

$$D = \{d_1, d_2, \dots, d_m\}$$

\rightarrow If instances are mutually independent it follows normal distribution, hence we consider probability density function.

$$= \arg \max_{h \in H} \prod_{i=1}^m P(d_i|h)$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2}$$

Apply log

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

maximum likelihood hypothesis for predicting hypothesis:

$$p(D|h) = P(d_i|h(x_i)) = \begin{cases} h(x_i) & \text{if } d_i = 1 \\ 1 - h(x_i) & \text{if } d_i = 0 \end{cases}$$

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i}$$

minimum description length principle:

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

$$h_{MDL} = \underset{h \in H}{\operatorname{argmin}} LC_1(\#) + LC_2(D|\#).$$