

①

Find S :-

It is a basic concept learning algorithm in machine learning.

i) Start with most specific hypothesis

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

ii) Take next example and if negative then no changes occur to hypothesis

iii) If positive and we find that our initial hypothesis is too specific then we update our current hypothesis to a general

iv) Repeat all steps until we reach last hypothesis.

Take an example

Sky	Temp	Humid	Wind	Water	Forest	Old
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

Now take most specific hypothesis as  $h_0$   
 $h_0 = \langle \phi, \phi, \phi, \phi, \phi, \phi \rangle$

I Compare with first hypothesis

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

II Compare it with second hypothesis and do accordingly

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

III Now, after comparing with third hypothesis we see that it is negative so according to algorithm we will do nothing  
 $h_3 = \langle \text{Sunny, Warm, ?, Strong, Warm, Same} \rangle$

IV Compare it with fourth hypothesis as it is positive so do accordingly  
 $h_4 = \langle \text{Sunny, Warm, ?, Strong, ?, ?} \rangle$

Now the final hypothesis is given as  
 $h_f = \langle \text{Sunny, Warm, ?, Strong, ?, ?} \rangle$

### # Limitations :-

- i) There is no way to determine if the hypothesis is consistent throughout the data
- ii) Inconsistent training set can actually mislead the find-S algorithm, since it ignores the negative examples
- iii) It does not provide a backtracking technique to determine the best possible changes that could be done to improve the result hypothesis.

### Candidate Elimination :-

It incrementally builds the version space given a hypothesis space  $H$  and a set  $E$  of examples.

- i) Load the data.
- ii) Initialize General and Specific Hypothesis

iii) For each <sup>tr</sup>e training example  
 if all value == hypo value:  
 Do nothing

else:  
 replace attribute value with ?

ii) If negative example  
 Make generalize hypothesis more specific.

Take an example

Sky	Temp	Humid	Wind	Water	Forest	off
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

I. Initialize most general and specific hypo.  
 $S_0: \langle \phi, \phi, \phi, \phi, \phi, \phi \rangle$        $G_0: \langle ?, ?, ?, ?, ?, ? \rangle$

I Compare with the first tr hypothesis.  
 $S_1: \langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle$        $G_1: \langle ?, ?, ?, ?, ?, ? \rangle$

II Now Compare it with second tr hypothesis.  
 $S_2: \langle \text{Sunny}, \text{Warm}, ?, \text{Strong}, \text{Warm}, \text{Same} \rangle$        $G_2: \langle ?, ?, ?, ?, ?, ? \rangle$

III As third hypothesis is -ve so do accordingly.

$S_3: \langle \text{Sunny}, \text{Warm}, ?, \text{Strong}, \text{Warm}, \text{Same} \rangle$        $G_3: \langle \text{Sunny}, ?, ?, ?, ?, ? \rangle$   
 $\langle ?, \text{Warm}, ?, ?, ?, ? \rangle$   
 $\langle ?, ?, \text{Normal}, ?, ?, ? \rangle$        $\langle ?, ?, ?, ?, \text{Cool}, ? \rangle$   
 $\langle ?, ?, ?, ?, ? \text{ same} \rangle$



IV Now again compare with last i.e. hyfo  
 $S_i: \langle \text{Sunny, Warm, ?, Strong, ?, ?} \rangle$   $G_i: \langle \text{Sunny, ?, ?, ?, ?} \rangle$   
 $\langle \text{?, Warm, ?, ?, ?, ?} \rangle$

Now final hypothesis is -

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

② Solve using KNN  $x(A=3, B=7)$   $k=3$

A	B	Label
7	7	False
7	4	False
3	4	True
1	4	True

Now if we want to find the label for  $x(A=3, B=7)$  then take distance formula as  $(x_1 - x_2)^2 + (y_1 - y_2)^2$

A	B	Distance	Label
7	7	16	False
7	4	25	False
3	4	9	True
1	4	18	True
3	7		?

Now, as we know that the value of  $k$  is 3 so we have to take least three distances.

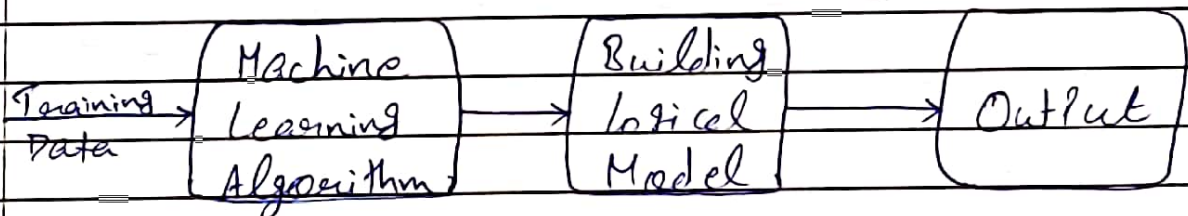
9, 13, 16 are the least three distances  
True, True, False are the respective labels

It can easily be seen that True is  
maximum time in the result, so we  
can conclude that label for  
 $x(A=3, B=7)$  is True.

### ③ a) Design a Learning System

When we feed the  
training data to machine learning algorithm  
this algorithm will produce a mathematical  
model and with help of model, the  
machine will make a prediction and  
take a decision without being explicitly  
programmed.

Also during training data, the  
more machine will work with it more  
it will get experience and the more  
it will get experience the more efficient  
result is produced.



Learn From Data

### Designing a System

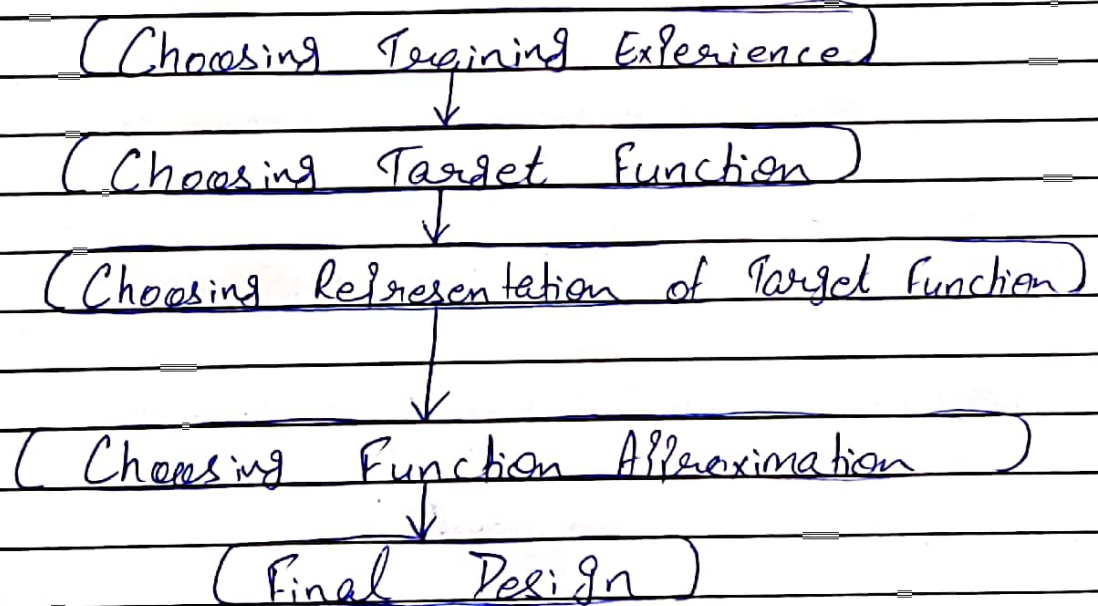
According to Tom Mitchell  
A computer program is said to be learning

from experience ( $E$ ), with respect to some task ( $T$ ). Thus, the performance measure ( $P$ ) is the performance at task  $T$ , which is measured by  $P$  and it improves with experience  $E$ .

Take an example:

Suppose we are designing a system for Spam E-mail Detection.

- Task,  $T$ : To classify mails into spam or Not spam.
- Performance measure,  $P$ : Total percent of mails being correctly classified as being 'spam' or 'Not spam'.
- Experience,  $E$ : Set of mails with label 'spam'.





	Supervised	Unsupervised
i)	They are trained using labeled data	i) They are trained using unlabeled data
ii)	It takes direct feedback to check if it is predicting correct output or not	ii) It does not take any feedback.
iii)	These model predicts the output.	iii) These model finds the hidden pattern in data.
iv)	In this, input data is provided to the model along with the output.	iv) In this, only input data is provided.
v)	The goal of this learning is to train the model so that it can predict the output when it is given new data.	v) The goal of this model is to find the hidden pattern and useful insights from unknown datasets.
vi)	It needs supervision to train the model.	vi) It does not need any supervision to train the model.
vii)	It can be categorized in Classification and Regression problems.	vii) It can be classified in clustering and Association problem.

④ Information gain calculates the reduction in entropy of an attribute. It is used for the construction of the decision tree from the training dataset by evaluating the information gain of each variable and collecting the variable that maximises this information gain which in turn minimises the entropy and splits the dataset into groups for effective classification.

A skewed distribution has low probability whereas a distribution where centres have equal to have larger entropy.

Entropy is used to calculate the purity of dataset and information gain precedes a way to use entropy to calculate how a change to the dataset impacts the purity of the dataset.

Formula is given as.

$$\text{Entropy}(S) = -P_0 \log_2(P_0) - P_1 \log_2(P_1)$$

$$\text{Information Gain}(S, x) = E(S) - \sum_{i=1}^n E(S_i)$$

where  $S$  is target population and  $x$  is the attributes

Now, take the example



A	B	Label
a1	b1	No
a1	b2	Yes
a2	b3	Yes
a2	b2	No
a2	b1	Yes

We can see that it contains 3 true values and 2 negative values.

$$\text{Entropy} \rightarrow -p_0 \log_2 p_0 - p_1 \log_2 p_1$$

$$E(S, -) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$$

$$\text{Gain}(S, A) = E(S) - \sum_{v \in \{a_1, a_2\}} \frac{|S_v|}{|S|} E(S_v)$$

$$= E(S) - \frac{2}{5} E(S_{a_1}) - \frac{3}{5} E(S_{a_2})$$

$$E(S_{a_1}) = 1 \quad (\because \text{equal true \& false examples})$$

$$E(S_{a_2}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.917$$

$$\text{Gain}(S, A) = 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.917$$

$$\boxed{\text{Gain}(S, A) = 0.0198}$$

$$\text{Now, Gain}(S, B) = E(S) - \sum_{v \in \{b_1, b_2, b_3\}} \frac{|S_v|}{|S|} E(S_v)$$

$$E(S) - \frac{2}{5} E(S_{b_1}) - \frac{2}{5} E(S_{b_2}) - \frac{1}{5} E(S_{b_3})$$

$$E(S_{b_1}) = 0$$

$$\text{Gain}(S, B) = 0.17$$

5 a)

10/13

- |      |   |      |   |
|------|---|------|---|
| vi)  | It is driven by binary differences b/w correct and predicted output | vi)  | It is driven by continuous differences between correct and predicted output |
| vii) | In this rule, the weight is modified by incremental updates.        | vii) | In this rule, the weight is modified by gradient descent                    |

## b) Back propagation Algorithm :-

It is used to effectively train a neural network through a method called chain rule. This algorithm performs learning on a multilayer feed forward neural network. It iteratively learns a set of weights for predicting a class label of tuples.

The multilayer feed forward neural network has three layers  $\rightarrow$  input, hidden and output layer.

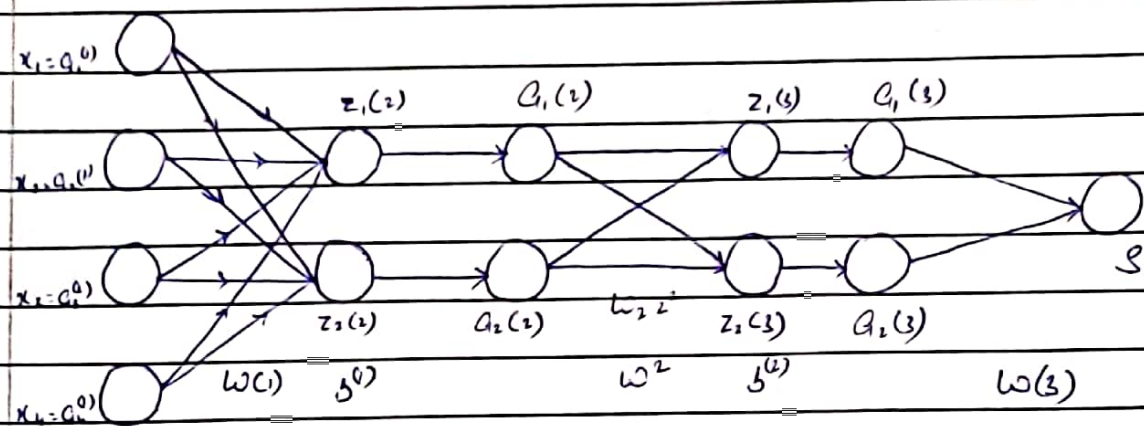
### Input Layer

A layer from where an input is fed to units.

$$x_i = a_i^{(1)}, i \in 1, 2, 3, 4$$

$a$  is a set of activations equal to the input values.





### Hidden Layer :-

These inputs then passed through the input layers and the weighted and feed to hidden layers, a layer of neuron like units. The final values are computed using  $z^{(l)}$  in layer and  $a^{(l)}$  - activations in layer.

$$l=2, \quad z^{(2)} = w^{(1)} x + b^{(1)}, \quad a^{(2)} = f(z^{(2)})$$

$$l=3$$

$$z^{(3)} = w^{(2)} a^{(2)} + b^{(2)}$$

$$a^{(3)} = f(z^{(3)})$$

where  $w^{(1)}$  &  $w^{(2)}$  are weights of layers 1 & layer 2 and  $b^{(1)}$  and  $b^{(2)}$  are the biases

### Output Layer :-

The final part of the

neural network is output layer which produces the predicted value

$$s = w^{(s)} a^{(s)}$$

### 6) Brute Force Bayes Learning

- For each hypothesis  $h$  in  $H$ , calculate the posterior probability.

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

- Output the hypothesis  $h_{max}$  with the highest posterior probability

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

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

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

To summarize, Bayes theorem implies that the posterior probability  $P(h|D)$  under our assumed  $P(h)$  and  $P(D|h)$  is

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