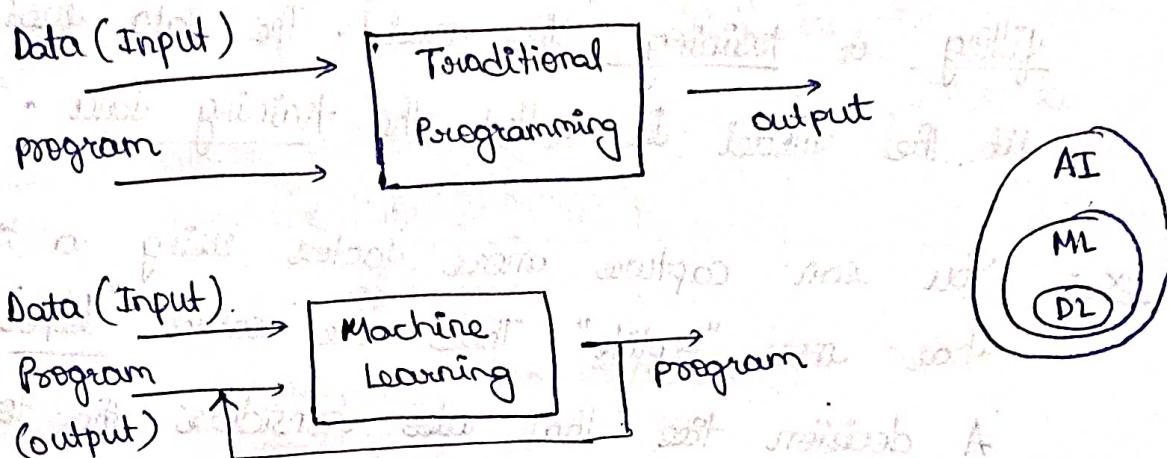


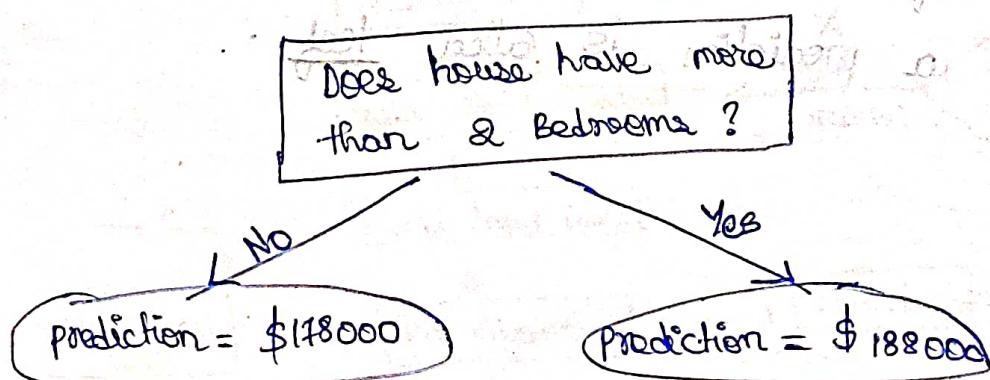
INTRODUCTION TO MACHINE LEARNING

- It is the field of study that gives computers the capability to learn without being explicitly programmed.



Ex: You ask your cousin how he's predicted real estate values in the past and he says it is just intuition. But more questioning reveals that he's identified price patterns from houses he has seen in the past, and he uses those patterns to make predictions for new houses he is considering.

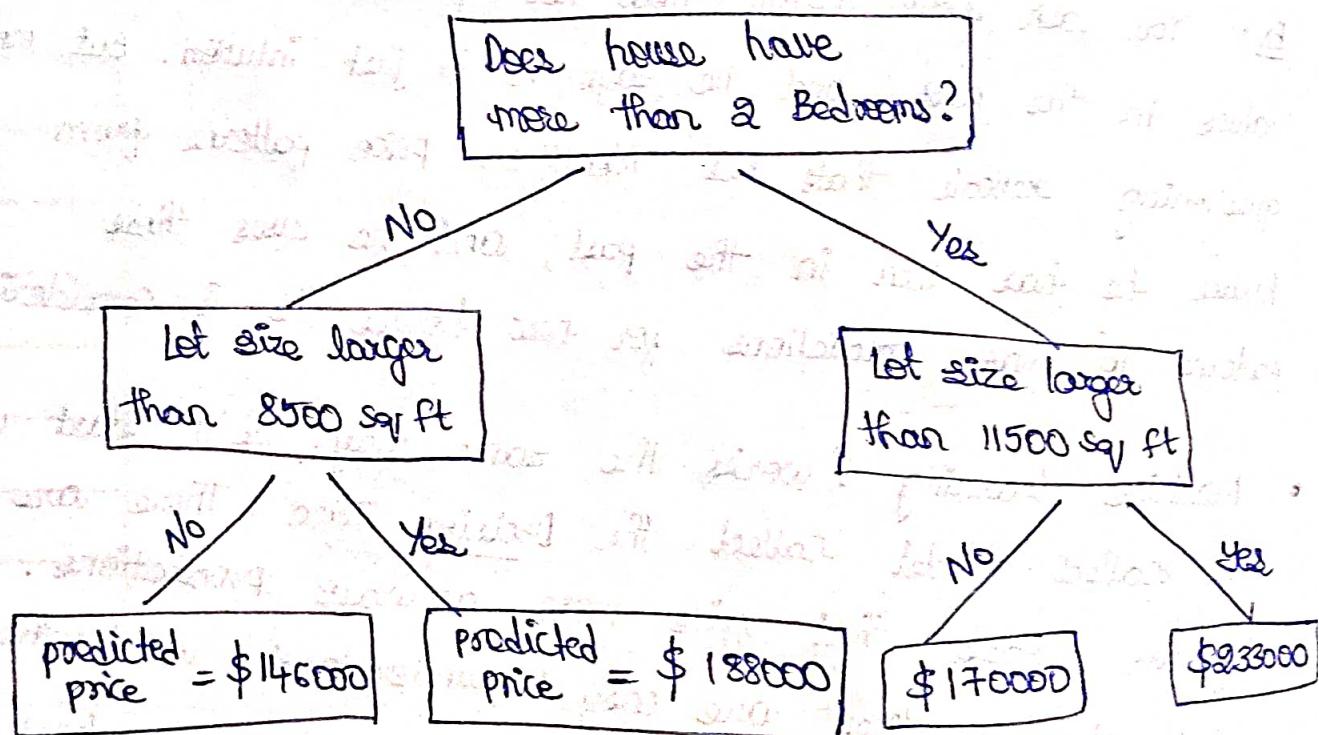
- Machine learning works the same way. We'll start with a called model called the Decision Tree. There are fancier models that give more accurate predictions. But decision trees are easy to understand, and they are the basic building block for some of the best models in data science.



GUDI VARAPRASAD

- We use data to make decision how to break the houses into two groups, and then again to determine the predicted price in each group.
- This step of capturing patterns from data is called fitting or training the model. The data used to fit the model is called the training data.

Ex: You can capture more factors using a tree that has more "splits". These are called "deeper" trees. A decision tree that also considers the total size of each house's lot might look like this.



- The predicted price for the house is at the bottom of the tree. The point at bottom where we make a prediction is called leaf.

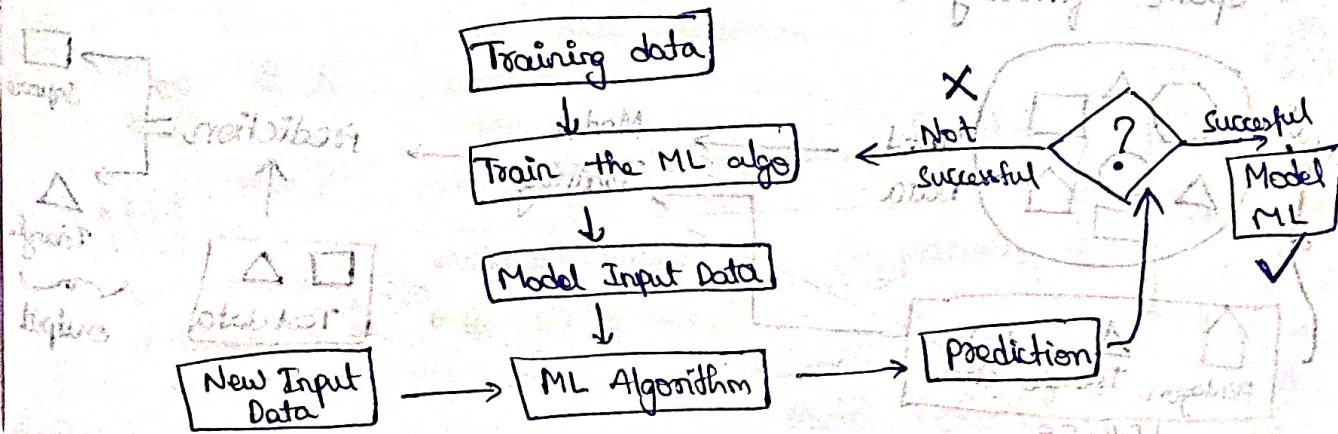
- We use scikit-learn library to create models. It is easily the most popular library for modeling the types of data typically stored in DataFrames.

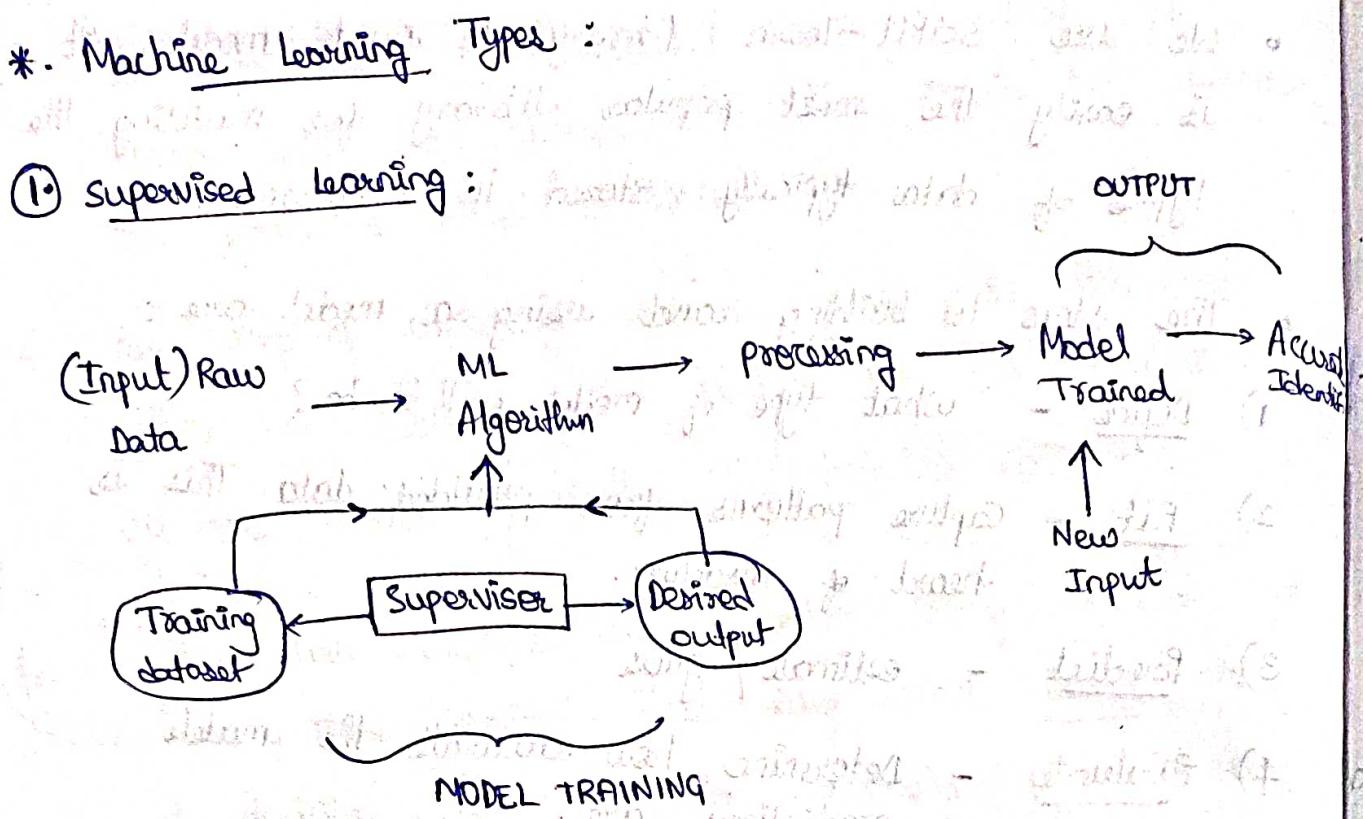
The steps to building and using a model are:

- Define - what type of model will it be?
- Fit - capture patterns from provided data. This is heart of modeling.
- Predict - estimate / guess
- Evaluate - determine how accurate the model's predictions are.

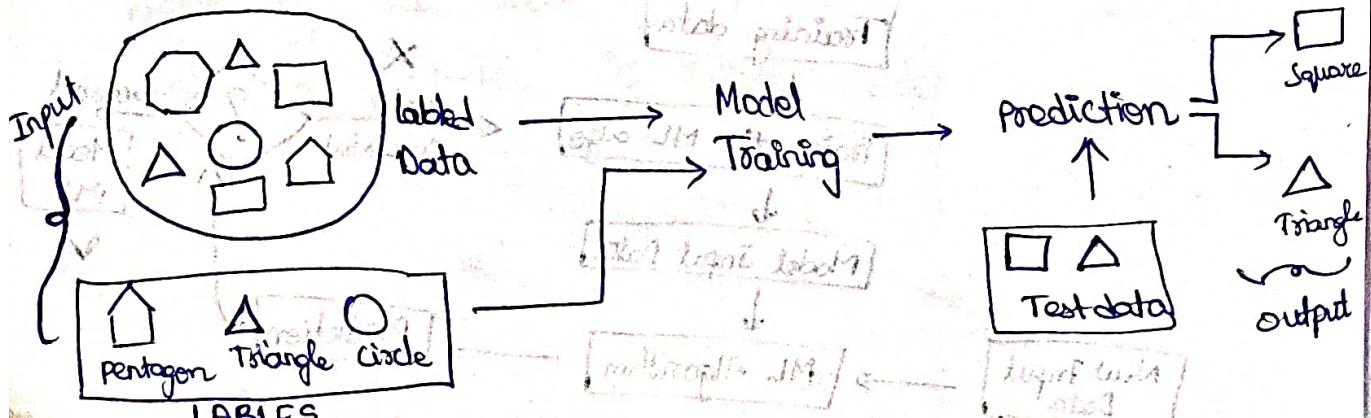
* INTRODUCTION

- Machine learning is a subset of Artificial Intelligence.
- It focuses mainly on the designing of systems, thereby allowing them to learn and make predictions based on some experience which is data in case of machines.
- AI - A technique which enables machines to mimic human behaviour.
- ML - subset of AI techniques which use statistical methods to enable machines to improve with experience.
- DL - subset of ML which make the computation of multi-layer neural network feasible.





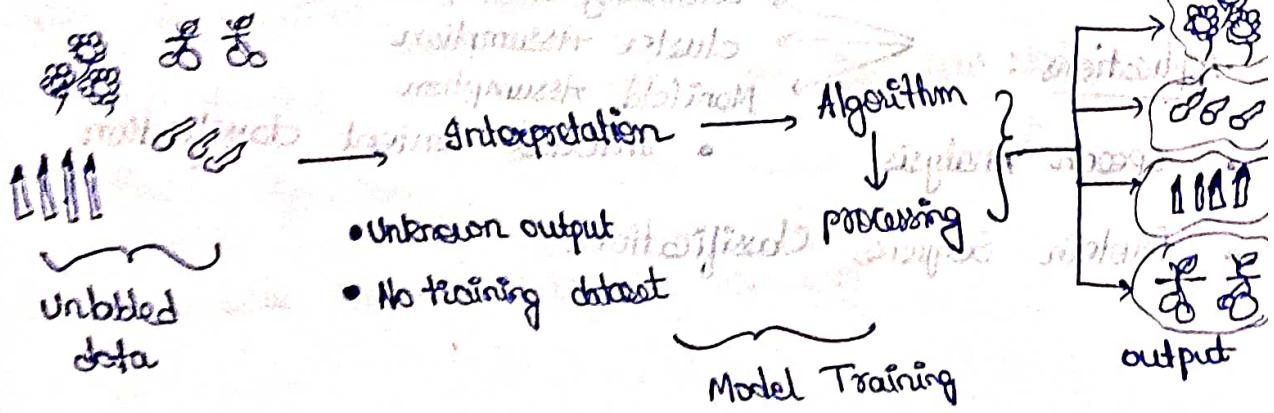
- It is a type of ML in which machines are trained using well "labelled" training data, and on basis of the data, machines predict the output.
- The "labelled" data means some input data is already tagged with the correct output.
- The aim of supervised learning algorithm is to find a mapping function to map the input variable (x) with the output variable (y).
- In real world, supervised learning can be used for Risk Assessment, Image Classification, Fraud Detection, spam filtering (etc).



- Supervised learning is a type of ~~but independent~~ learning.
- 1) Regression - relationship between two variables.
- 2) Classification - there are two classes (output variable is categorical).
- Advantages:
 - + The model can predict the output on basis of prior experiences.
 - + We have exact idea about classes of objects.
- Disadvantages:
 - Not suitable for handling complex tasks.
 - Cannot predict correct output if the test data is different from the training dataset.
 - Training required lots of computation times.
 - We need enough knowledge about classes of object.

② UNSUPERVISED LEARNING

- It is a type of ML in which models are trained using unlabeled datasets and are allowed to act on the data without any supervision.
- The goal of unsupervised learning is to find the underlying structure of dataset, group that data according to similarities, and represent that dataset in a compressed format.



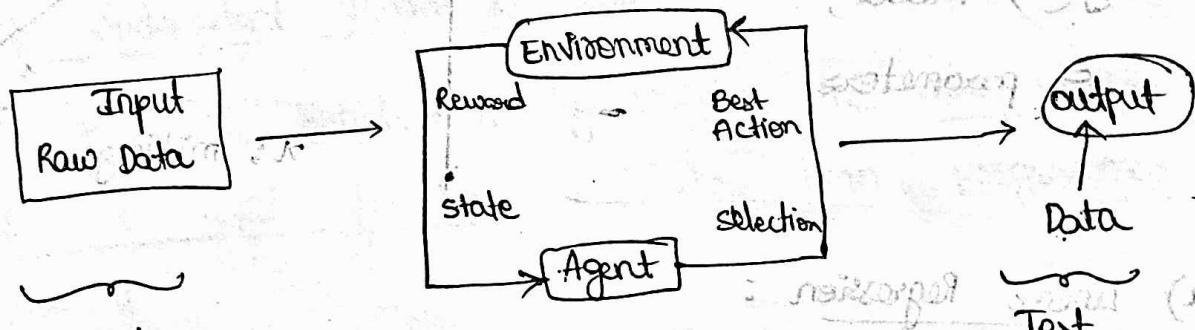
- Types of Unsupervised Learning :
 - ① clustering : similarities & non-similarities are made into different groups.
 - ② Association : determine set of items that occurs together in dataset & finding the relationships between variables.
- Advantages :
 - + More complex tasks.
 - + Market Basket Analysis.
 - + easy to get unlabeled data in comparison to labeled data.
- Disadvantages :
 - More difficult.
 - Result may be less accurate.
 - Algorithms don't know exact output in advance.

- Note : Semi Supervised Learning :
- The algorithm is trained upon a combination of labeled and unlabeled data.
 - Typically, this combination will contain a very small amount of labeled data and a very large amount of unlabeled data.

- Applications :
- Speech Analysis
 - Protein Sequence Classification
 - Internet Content classification

→ continuity Assumption
 → cluster Assumption
 → Manifold Assumption

- ③ Reinforcement Learning: It is a sub-domain of ML.
- It is a feedback based ML in which a machine learnt to behave in an environment by performing the actions and seeing the results of action.
 - For each good action, the machine gets positive feedback and for each bad action, the machine gets negative feedback.
 - Here machine learns automatically using feedbacks without any labeled data. It is bound to learn by its experience only.
 - Applications: Robotics, Control, Game Playing, Chemistry, Business, Manufacturing, Finance Sector.



* REGRESSION:

- Supervised learning has been broadly classified into 2:

 - Regression
 - Classification

- Regression is a kind of supervised learning that learns from the labeled datasets and is then able to predict a continuous valued output for the new data given to the algorithm.
- Ex: Linear Regression, Logistic Regression

GUDI VARAPRASAD

- Ex: We want to have a system that can predict the price of a used car.
- Inputs are the car attributes (brand, year, engine capacity, mileage, etc...)
- The output is the price of the car. Such problems where the output is a number are regression problems.

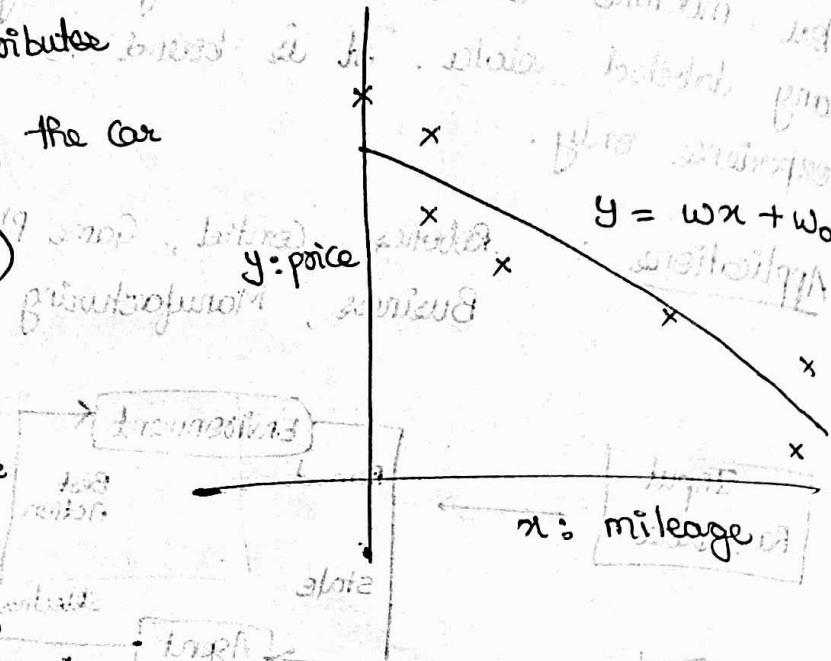
x : car attributes

y : price of the car

$$y = g(x|\theta)$$

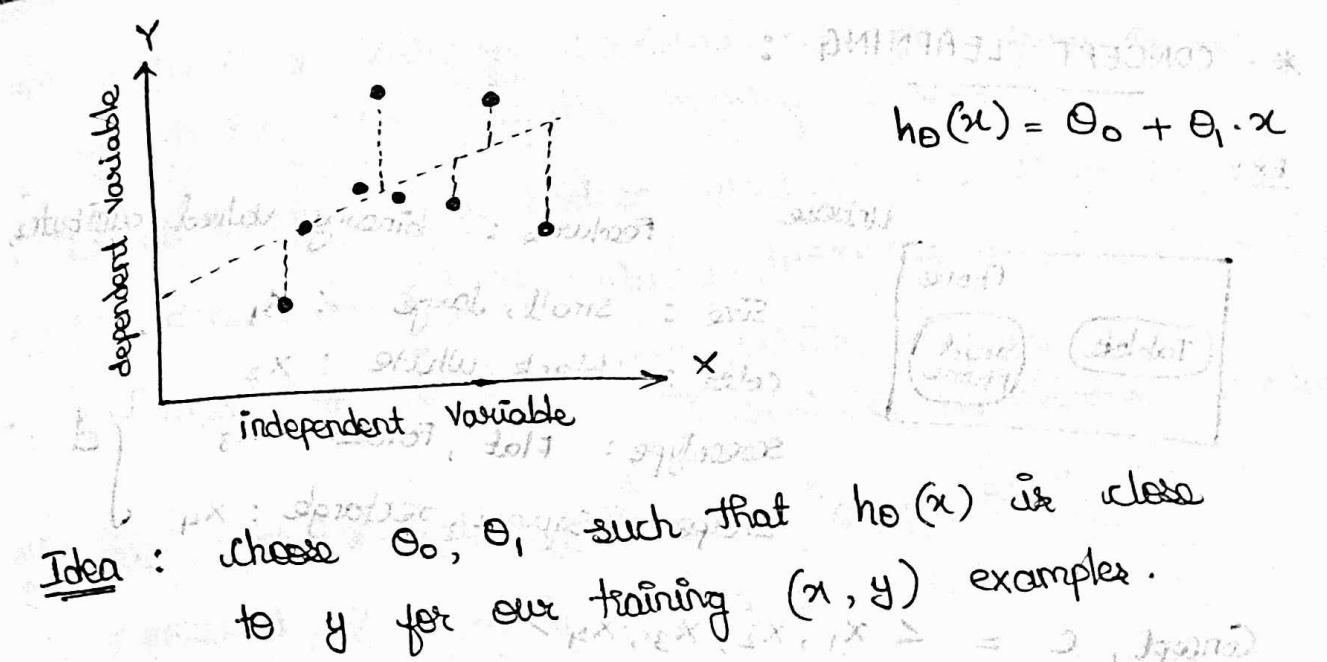
$g(\cdot)$ model,

θ parameters



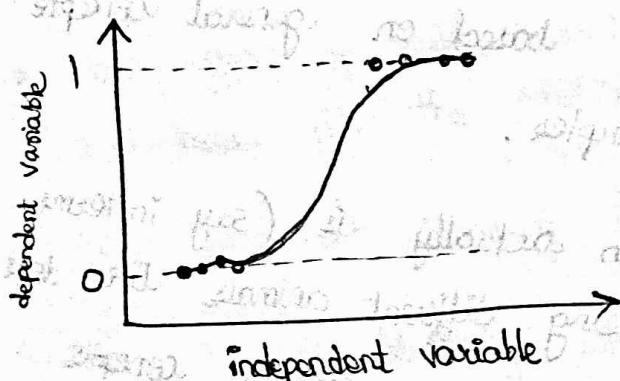
a) Linear Regression:

- This algorithm assumes that there is a linear relationship between the 2 variables, input (x) and output (y), of the data it has learnt from.
- The input variable is called the "Independent Variable" and the output variable is called Output/Dependent Variable". When unseen data is passed to the algorithm, it uses the function, calculates and maps the input to a continuous value for the output.



b) Logistic Regression :

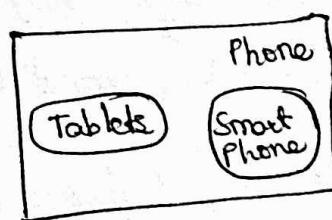
- This algorithm predicts discrete values for the set of Independent variables that have been passed to it.
- It does the prediction by mapping the unseen data that has been programmed into it.
- The algorithm predicts the probability of new data and as its' output lies between the range of 0 and 1.



GUDI VARAPRASAD

* CONCEPT LEARNING :

Ex:



Features : binary valued attributes

size : small, large	: x_1
color : black, white	: x_2
screen type : Flat, folded	: x_3
shape : square, rectangle	: x_4

$\} d$

$$\text{Concept, } C = \langle x_1, x_2, x_3, x_4 \rangle$$

$$\text{Tablet} = \langle \text{large, black, Flat, Square} \rangle$$

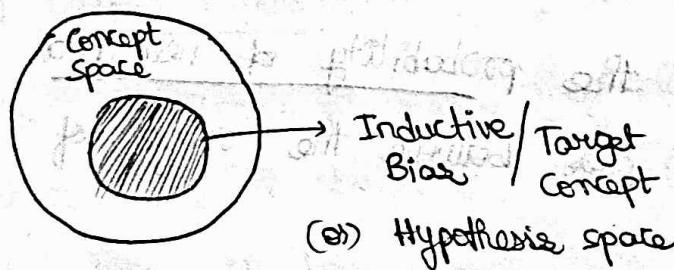
$$\text{SmartPhone} = \langle \text{small, grey, folded, rect} \rangle$$

$$\begin{aligned} 2^d &= \text{possible instances} \\ 2^d &= \text{possible concepts} \end{aligned}$$

$$\text{Here, } d = 4$$

$$\text{possible instances} = 2^d = 16$$

$$\text{possible concepts} = 2^{16}$$



- It is a type of learning based on general concept from specific training examples.
- Ex: People learn what an actually is (say in terms of its features) by seeing different animals like dogs, cats, etc.. Here dogs, cats are specific concepts whereas Animal is general concept.

$$C(x) = x \text{ is bird}$$

$$\text{If } x \text{ is cat} \Rightarrow C(x) = 0$$

$$\text{If } x \text{ is parrot} \Rightarrow C(x) = 1$$

- Learning a category description (concept) from a set of positive and negative training examples.
- Target Function : a boolean function $c: X \rightarrow \{0, 1\}$.
- Experience : a (set) of training instances $D: \{<x, c(x)>\}$
- Instances of $c(x) = 1 \rightarrow$ positive example
 $c(x) = 0 \rightarrow$ negative example
- IMP: The most general hypothesis - "accept all" is represented by $\langle ?, ?, ?, ?, ? \rangle$
- IMP: The most specific possible hypothesis - "reject all" is represented by $\langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$
- $\langle x_i, \Theta_i \rangle; \Theta_i = \{+, -\}$ so, $\langle x_i, c(x_i) \rangle$ is the training sample with Θ_i as attributes like example $\langle \text{Redmi}, 1 \rangle$, $\langle \text{Apple}, 0 \rangle$.

- Inductive learning hypothesis : An hypothesis that approximates well over a sufficiently large set of training examples will also approximate the target function (hypothesis) over other unobserved examples.

* CONCEPT LEARNING AS SEARCH :

- Main goal of this search is to find the hypothesis that best fits the training example.
- Ex : Enjoy sport learning Task (6 attributes)
 - $\langle \text{Sky}, \text{Air Temp}, \text{Humidity}, \text{Wind}, \text{Water}, \text{Forecast} \rangle$
 - Sky has values (3) : (Foggy, cloudy, sunny)
 - remaining has only 2 values.

\therefore Different instances possible = $3 \times 2 \times 2 \times 2 \times 2 \times 2 = 96$,
 \Rightarrow Syntactically distinct possible hypothesis (Additional 2 & 4 are added)

$$\text{So, } = (3+2) \times (2+2) \times (2+2) \times (2+2) \times (2+2) \times (2+2) \\ = 5 \times 4 \times 4 \times 4 \times 4 \times 4 = 5120$$

\Rightarrow Semantically distinct possible hypothesis (Null is taken care of)

$$= 1 + ((5-1) \times (4-1) \times (4-1) \times (4-1) \times (4-1) \times (4-1)) \\ = 1 + (4 \times 3 \times 3 \times 3 \times 3 \times 3) = 973$$

After finding all the syntactically & semantically distinct hypothesis, we search the best match from all those (i.e. much closer to our learning process).

- * FIND-S ALGORITHM: (finding a maximally specific hypothesis)
- S means most specific, this algorithm considers only positive examples.
- Algorithm :

Step 1 : Initialize with most specific hypothesis $h_0 (\phi)$

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

Step 2 : For each +ve sample,

for each attribute

if (Value = hypothesis value) \Rightarrow Ignore

else Replace with most general hypothesis (?)

GUDI VARAPRASAD

Example:

origin	Manufacturer	color	Year	Type	Class
Japan	Honda	Blue	1980	economy	+ (Yes)
Japan	Toyota	Green	1970	Sport	- (No)
Japan	Toyota	Blue	1990	economy	+ (Yes)
Japan	Audi	Red	1980	economy	- (No)
Japan	Honda	White	1980	economy	+ (Yes)
Japan	Toyota	Green	1980	economy	+ (Yes)
Japan	Honda	Red	1980	economy	- (No)

solution:

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

economy

$$h_1 = \langle \text{Japan}, \text{Honda}, \text{Blue}, 1980, \text{economy} \rangle$$

$$h_2 = h_1 \quad (\text{since } -\text{ve class})$$

$$h_3 = \langle \text{Japan}, ?, \text{Blue}, ?, \text{economy} \rangle$$

$$h_4 = h_3 \quad (\text{since } -\text{ve class})$$

$$h_5 = \langle \text{Japan}, ?, ?, ?, \text{economy} \rangle$$

$$h_6 = \langle \text{Japan}, ?, ?, ?, \text{economy} \rangle$$

$$h_7 = \langle \text{Japan}, ?, ?, ?, \text{economy} \rangle$$

This is most specific hypothesis

- Disadvantages:

- Consider only +ve values / ignores -ve values. Complete
- Here, h_7 may not be sole hypothesis that fits the data.
- There can be more than one "most specific hypothesis"
- If training data is inconsistent \Rightarrow algorithm will mislead.

* VERSION SPACES :

- subset of hypotheses H consistent with the training examples (d)

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

$H \rightarrow$ hypothesis (entire) $d =$ -training example

if $h(x) = c(x)$ \Rightarrow the hypothesis is consistent with the training example

Algorithm: (List then Eliminate)

- Version Space \leftarrow list containing every hypothesis in H
- From this step, we keep on removing inconsistent hypotheses from version space
for each training example, $\langle x, c(x) \rangle$ remove any hypothesis that is $h(x) \neq c(x)$
- output the list of hypotheses into version space after checking for all training examples.

Advantages:

- + It guarantees to output all hypotheses consistent with training data.

Disadvantages:

- The hypothesis space must be finite.
- Enumeration of all the hypotheses, rather inefficient.

* CANDIDATE ELIMINATION ALGORITHM

- It uses the concept of version space.
- Considers both +ve and -ve values (Yes, No).
- Both specific and general hypotheses.
 - For positive samples, move from specific to general.
 - For negative samples, move from general to specific.

$$S = \{\emptyset, \emptyset, \emptyset, \emptyset, \emptyset\} \downarrow \text{ (Step 1)}$$

$$G = \{?, ?, ?, ?, ?\} \uparrow \text{ (Step 2)}$$

Algorithm :

1. Initialize both general and specific hypothesis (S, G)

$$S = \{\emptyset, \emptyset, \emptyset, \dots, \emptyset\} \quad G = \{?, ?, ?, \dots, ?\}$$
2. For each example,
 - if example is positive \rightarrow make specific to general
 - else example is negative \rightarrow make general to specific

Example : Enjoy Sport Learning Task.

No	Sky	Temperature	Humidity	Wind	Water	Forecast	Enjoy
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

GUDI VARAPRASAD

No. of Attributes = 6 VARIOBLES ELIMINATION GENERALISATION

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

(1) +ve (positive) output any basic out. After elimination,

$$S_1 = \{ \text{sunny, warm, Normal, strong, warm, same?} \}$$

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

(2) +ve, (specific to general)

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

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

(3) -ve, (general to specific)

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

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

usage of $\langle ?, ?, ?, ?, ?, ? \rangle$, same $\rangle \{ \}$

(4) +ve (specific to general)

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

Some is removed because is changed

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

Finally, S_4 & G_4 are Final Hypothesis

* LEARNED VERSION SPACE :

- The learned Version Space correctly describes the target concept, provided :
 - There are no errors in the training examples.
 - There is some hypothesis that correctly describes the target concept.
- If S and G converge to a single hypothesis, the concept is exactly learned.
- In case of errors in the training, useful hypothesis are discarded, no recovery possible.
- An empty version space means no hypothesis in H is consistent with training examples.
- The learned version space doesn't change with different orderings of training examples.
- optimal strategy (if you are allowed to choose) :
 - Generate instances that satisfy half the hypotheses in current version space.
 - Ex: $\langle \text{sunny}, \text{warm}, \text{Normal}, \text{Light}, \text{Warm}, \text{Same} \rangle$
 - satisfies 3/6 hyp.
 - Ideally the VS can be reduced by half at each experiment.
 - Correct target found in $\lceil \log_2 |VS| \rceil$ experiments.

* INDUCTIVE BIAS :

- Remarks on Candidate Elimination & Version Space Algorithm
- 1) Will the Candidate Elimination algorithm give us correct hypothesis? whether it will give us correct hypothesis or not
 - 2) What training example should the learner request next?
- Inductive learning : From examples we derive rules.
 - Deductive learning - Existing rules are applied to our examples.
 - Biased hypothesis space : Doesn't consider all types of training examples. (Solution → include all hypothesis)
 - Unbiased hypothesis space : Provided a hypothesis capable of representing set of all examples.
- ex: sunny & warm & normal & strong & cool & change = Yes
- Possible Instances = $3 \times 2 \times 2 \times 2 \times 2 \times 2 = 96$
- Target Concept = 2 (huge) [Practically not possible]
- Idea of Inductive Bias : Learner generalizes beyond the observed training examples to infer new examples.
- '>' → inductively inferred from
- $x > y \rightarrow y$ is "inductively inferred from" x .

Example: Learning from a set of training instances.

Suppose, learning algorithm L takes
training data $D_{\text{train}} = \{x, c(x)\}$
New instance = x_i (to be classified)

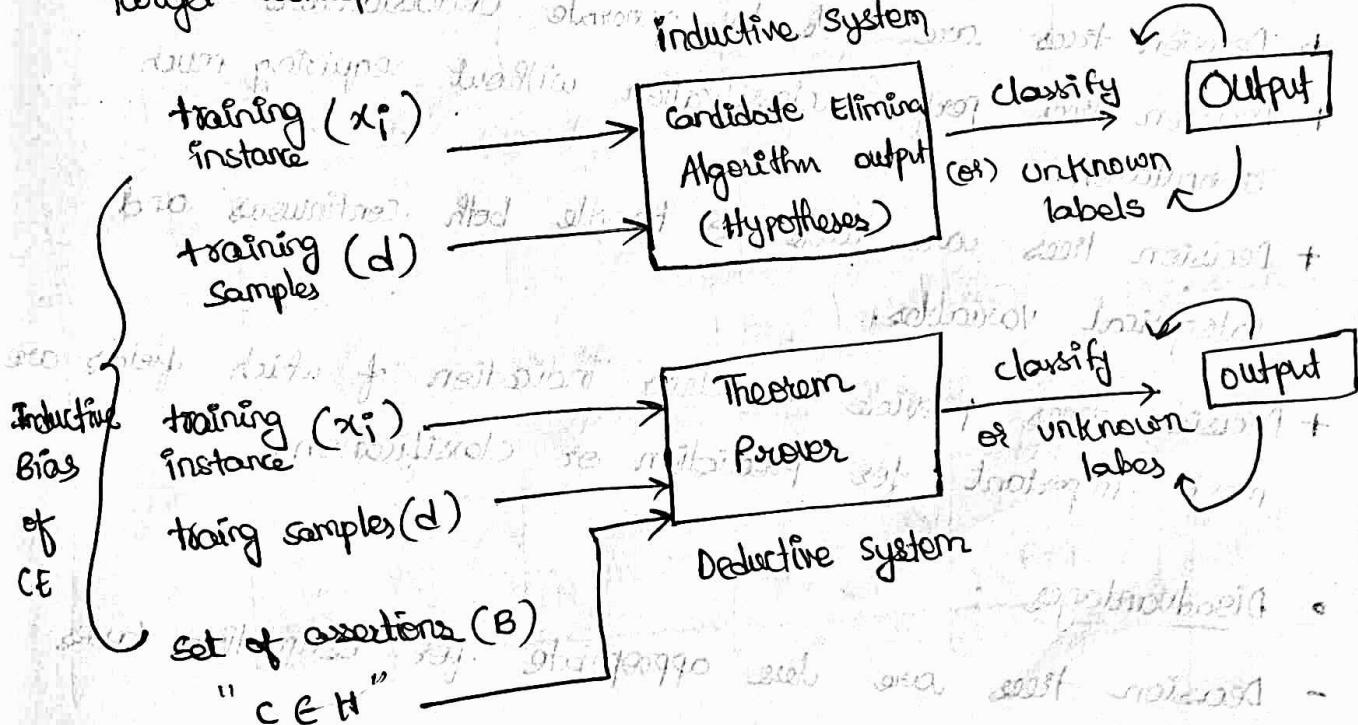
Represented as $L(x_i, d_c)$

$$\rightarrow (d_c \wedge x_i) > L(x_i, d_c) \text{ then } L(x_i, d_c) \text{ is inductively inferred.}$$

definition: For all instances, which belongs to x_i , it will hold the assertions (B) & has training instances (d) with target concept (c) for all values of x it implies, Generalizing (L) that has values of $x \in$ training instance with target concept.

$$\forall x_i \in X [(B \wedge d_c \wedge x_i) \vdash L(x_i, d_c)] * \text{IMP}$$

- * Case of Candidate Elimination Algorithm says that your "target concept is consistent with the hypothesis."



- Benefit: It essentially follows non-procedural way of classifying instances.
- ② → It will give a comparative strength of different classification algorithms.
- Rule Learner: No Bias
- Find-S Algo: most specific hypotheses (+ve examples)
- Candidate Elimination: VS H, d = Classification

Rank these algorithms first:
based on Inductive Bias

Rule learner	<	Candidate elimination	<	Find-S
--------------	---	-----------------------	---	--------

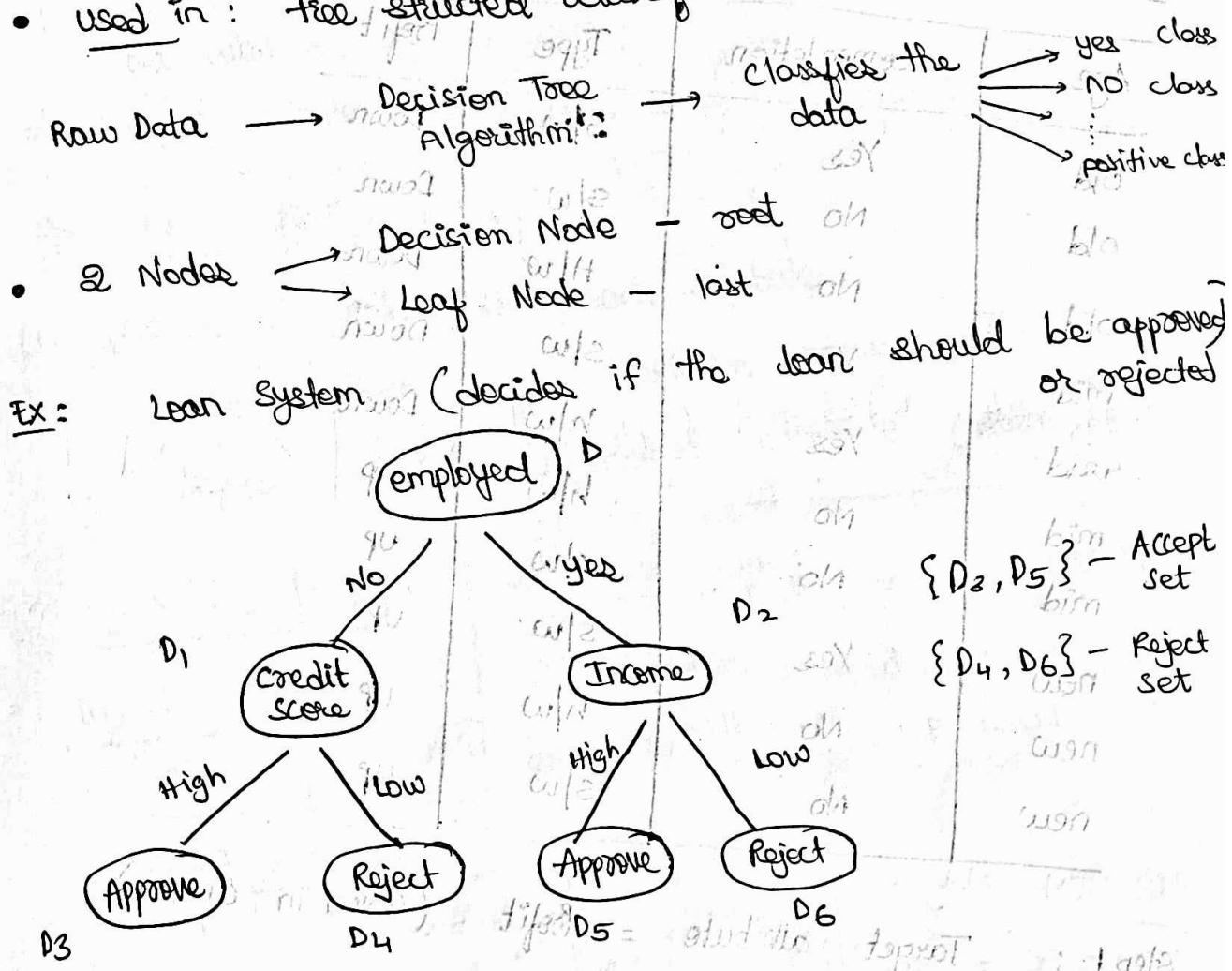
* DECISION TREE LEARNING :

- It is a flow chart like structure, where each internal node denotes a test on an attribute, each branch represents an output / outcome of the test, & each leaf node (terminal node) holds a class label.
- Advantages:
 - + Decision trees are able to generate understandable rules.
 - + Decision trees perform classification without requiring much computation.
 - + Decision trees are able to handle both continuous and categorical variables.
 - + Decision trees provide a clear indication of which fields are most important for prediction or classification.
- Disadvantages:
 - Decision trees are less appropriate for estimation tasks

GUDI VARAPRASAD

where the goal is to predict the value of a continuous attribute.

- Decision Trees are prone to errors in classification problems with many classes and relatively small no. of training examples.
- Decision Tree can be computationally expensive to train.
- The process of growing a decision tree is expensive.
- used in: tree structured classification & Regression



- Algorithm (ID3):
- In the given dataset, choose a target attribute.
- Calculate Information gain of target attribute.

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

GUDI VARAPRASAD

3. For remaining attributes, find entropy

$$\text{Entropy} = \text{IG} \times \text{probability}$$

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

$$4. \text{ Calculate } \text{Gain} = \text{IG} - E(A)$$

Example:

Age	Competition	Type	Profit
old	Yes	s/w	Down
old	No	s/w	Down
old	No	h/w	Down
old	Yes	s/w	Down
mid	Yes	h/w	Down
mid	No	h/w	Up
mid	No	s/w	Up
new	Yes	s/w	Up
new	No	h/w	Up
new	No	s/w	Up

Step 1 : Target attribute = Profit

(Given in question)

Step 2 : Information Gain (IG)

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

Here

$$P = \text{Count of down class} = 5 \rightarrow \text{how many times down occurred}$$

$$N = \text{Count of up class} = 5 \rightarrow \text{how many times up occurred}$$

$$= \left(-\frac{5}{10} \log_2 \left(\frac{5}{10} \right) - \frac{5}{10} \log_2 \left(\frac{5}{10} \right) \right) = 1 \text{ (Ans) DT}$$

so, Information Gain, IG = 1

Step 3: calculate entropy for remaining attributes
 only for non-target attributes

$$[\text{Entropy} = \text{IG} \times \text{Probability}] \rightarrow$$

attributes $\rightarrow \{\text{age, competition, Type, profit}\}$

Target attributes $\rightarrow \{\text{profit}\}$ (step 1)

Non-Target attributes $\rightarrow \{\text{age, competition, Type}\}$

For Age; finding Entropy :

(1) Prepare a table for each attribute

rows - values of undertaken attribute (old, mid, new)
 columns - values of target attribute (down, up)

Age values	down	up	$\text{Entropy} = \text{IG} \times \text{probability}$
old	3	0	$P_{\text{down}} = \frac{3}{10}$, down count
mid	2	2	$N = \text{up count}$
new	0	3	$N = \text{total pft.}$

$$\text{IG (old)} = - \left(\frac{3}{3} \log \left(\frac{3}{3} \right) + \frac{0}{3} \log \left(\frac{0}{3} \right) \right) = 0$$

$$\text{Probability} = \frac{3}{10} \left[\frac{\text{No. of times old}}{\text{Total values}} \right]$$

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

$$\text{IG (mid)} = - \left(\frac{2}{4} \log \left(\frac{2}{4} \right) + \frac{2}{4} \log \left(\frac{2}{4} \right) \right) = 1$$

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

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

$$IG(\text{new}) = - \left[\left(\frac{0}{3} \log \left(\frac{0}{3} \right) + \frac{3}{3} \log \left(\frac{3}{3} \right) \right] = 0$$

$$\text{Probability} = 3/10 \quad (= P.I. \text{ and no change})$$

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

$$\begin{aligned}\text{Entropy (age)} &= \text{Entropy (old)} + \text{Entropy (mid)} + \text{Entropy (new)} \\ &= 0 + 0.4 + 0 = 0.4\end{aligned}$$

Step 4 : Gain = $IG - E(A)$

$$= 1 - 0.4 = 0.6$$

Repeat the same calculation of Gain of Competition, Type

So, we obtain:

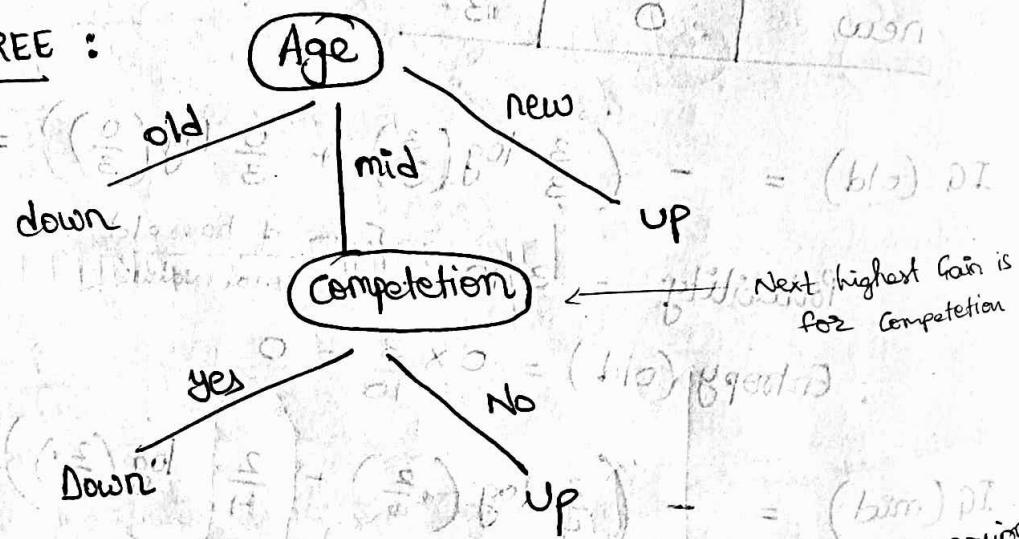
$$\text{Gain (competition)} = 0.124$$

$$\text{Gain (Age)} = 0.6$$

$$\text{Gain (Type)} = 0$$

Highest Gain \rightarrow decision node (root) \rightarrow Age bin

DECISION TREE :



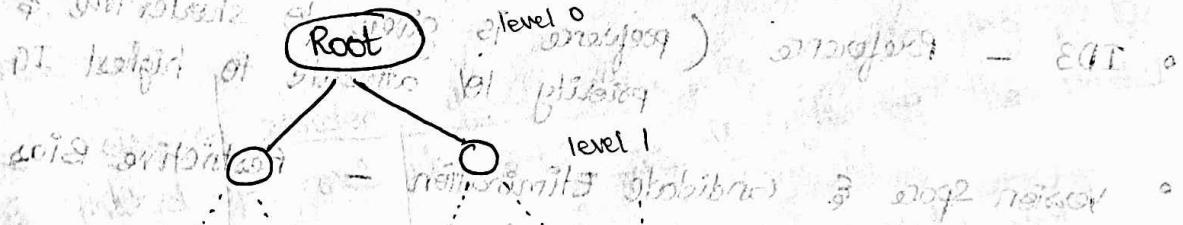
Type \rightarrow gain = 0 & even in Decision Tree, no requirement

- 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 o/p values.
3. Training data can have errors.
4. May contain missing attribute values also.

*. HYPOTHESIS SPACE SEARCH IN ID3: (ID3 is DTL)

- ID3 can be characterized as searching a space of hypotheses for one that fits the training examples.
- ID3 will search set of possible decision trees from available hypothesis.
- ID3 performs simple to complex searching.



- *. First start with empty tree & keep on adding nodes.
- Every discrete valued (finite) function can be described by some decision tree.
- Avoids major risk of incomplete hypothesis.
- Has only single current hypothesis.
- Cannot determine alternative decision trees.
- Backtracking is not possible.
- Can be extended easily to noisy data also.

* INDUCTIVE BIAS IN ID3 | DTL : {Preference Bias}

Set of assumptions

- Inductive Bias of ID3 consists of describing the basis by which ID3 chooses one consistent decision tree over all the possible ID3 DT's.
- ID3 Search strategy: Selects shortest tree over longer ones.
- 1) Select in favour of shorter trees over longer ones.
- 2) selects element with highest IG as root attribute over lowest IG ones.
- Types of Inductive Bias:

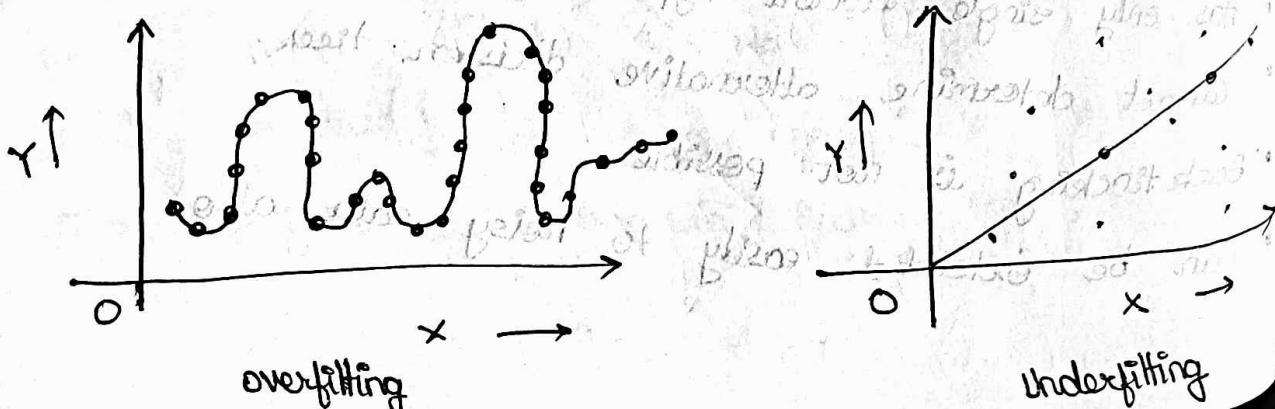
 - Restrictive bias - based on conditions.
 - Preference bias - based on priorities.

- ID3 - Preference (preference is given to shorter tree & priority to attribute to highest IG)
- Voronoi Space & Candidate Elimination - Restrictive Bias

Why short hypothesis? (why short tree?)

→ According to "Occam's Razor" - Prefer simplest hypothesis that fits the data (into table).

* OVERFITTING : [IMPORTANT]



- Given a hypothesis space H , a hypothesis $h \in H$ is said to OVERFIT the training data if there exists some alternative hypothesis $h' \in H$, such that h has smaller error than h' over the training examples but h' has smaller error than h over the entire distribution of instances (including the instances beyond the training set)

- Reasons for overfitting:
 - Errors and noise in training examples.
 - small no. of training samples to produce a representative procedure.
 - overfitting happens when a model learns the details of noise in the training data to the extent that it negatively impacts the performance of the model on new data.
 - The noise or random fluctuations in the training data is picked up and learned as concepts by the model.
 - These concepts are not applicable to new data and negatively impact the models ability to generalize.
- Avoid overfitting?
 - stop growing when data split is not statistically significant
 - Grow full tree, then post-prune.
 - Use separate set of examples, distinct from training set.
 - Reduced Error Pruning.

* Fails to classify test data (new data)

* REDUCED ERROR PRUNING

- The process of adjusting decision tree to minimize the misclassification error is, pruning.
- Aim is to minimize the overfitting of the model.
- Pre-pruning / early stopping rule:
 - Halt subtree construction @ some node after checking some measure.
 - Stops the tree growth pre-maturely.
 - Avoids generating overly complex sub-trees which overfits the training data.
- Post Pruning:
 - Deletes DT entirely.
 - Trims the nodes of DT in bottom-up fashion.
 - Replace nodes by leaf.
 - If error increases after trimming replace sub-tree by a leaf node.
- REP holds out some available instances called Pruning set.
- Pruning set provides an estimate of error rate on future instances than of training data.


```

if (error (parent) >= error (child))
    {
        prune the tree & replace node with leaf
    }
    else
        dont prune, repeat the process.
            
```

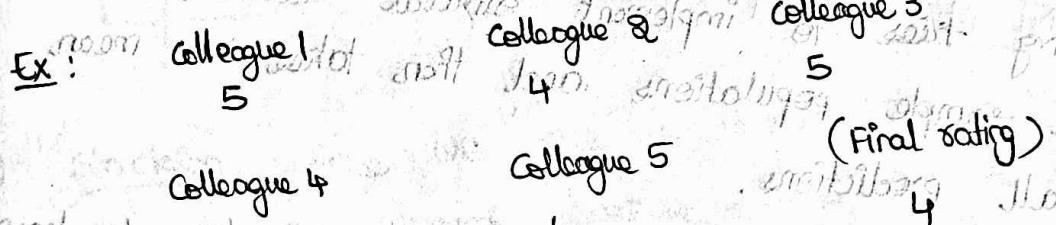
*. ENSEMBLE LEARNING

- It is a process by which multiple models, such as classifiers or experts, are "strategically" generated & combined to solve a problem.

- Ensemble learning is primarily used to improve the (classification, prediction, function approximation) performance of a model.

*. Basic Ensemble Techniques :

- ① Max Voting : The max voting method is generally used for classification problems. In this technique, multiple models are used to make predictions for each data point.



- ② Averaging : Multiple predictions are made for each data point in averaging. In this method, we take average of predictions from all the models and use it to make predictions (final).

$$\text{Final rating} = \frac{5 + 4 + 5 + 4 + 4}{5} = 4.4$$

- ③ Weighted Average: "All models are assigned different weights defining the importance of each model for prediction."

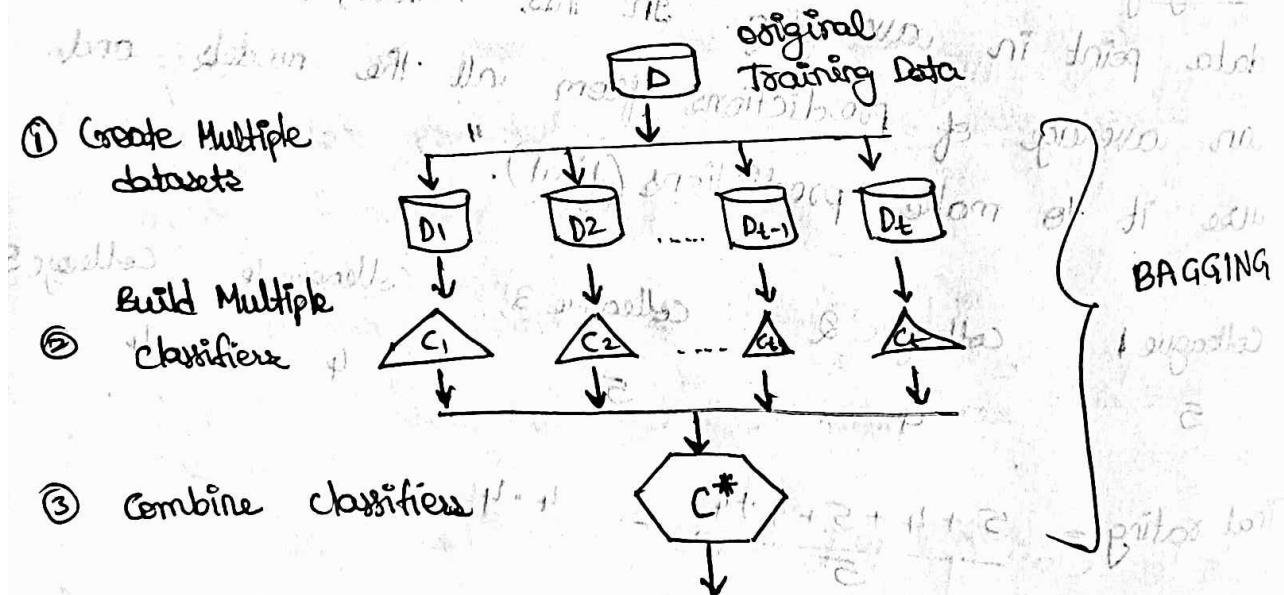
	Colleague 1	Colleague 2	Colleague 3	Colleague 4
Weight	0.23	0.23	0.18	0.18
Rating	5	4	5	4

$$\text{Final rating} = (5 \times 0.23) + (4 \times 0.23) + (5 \times 0.18) + \\ (4 \times 0.18) + (4 \times 0.18) = 4.41$$

* - ADVANCED ENSEMBLE LEARNING

① Boggling :

- Bagging tries to implement similar learners on small sample populations and then takes mean of all predictions.
 - In generalized bagging, you can use different learners on different population. This helps to reduce variance.



② Boosting:

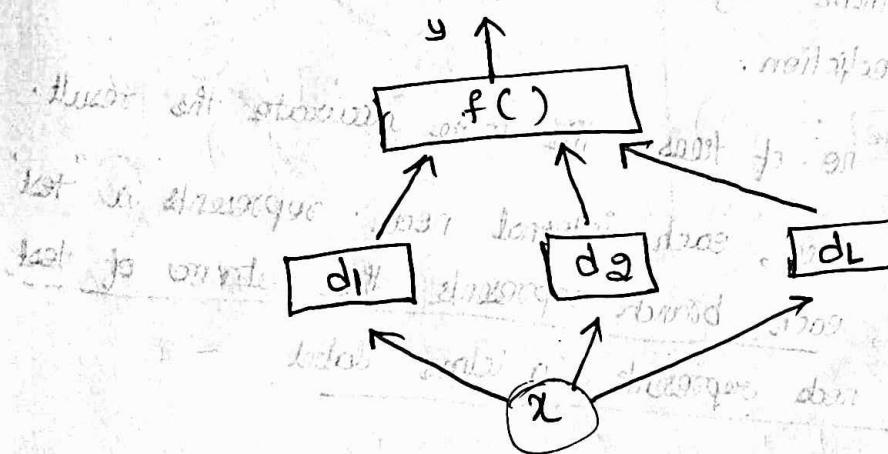
- Boosting is an iterative technique which adjust the weight of an observation based on the last classification.
- If an observation was classified incorrectly, it tries to increase the weight of this observation and vice-versa.
- In general, it decreases the bias error and builds strong predictive models. However, they may sometimes over fit on the training data.

+	+	-
+	-	-
+	-	-

+	+	-
+	0	-
0	0	-
0	+	-

+	+	-
+	-	-
-	-	-
-	-	+

- ③ Stacking: Here we use a learner to combine output from different learners. This can lead to decrease in either bias or variance depending on the combining learner we use.



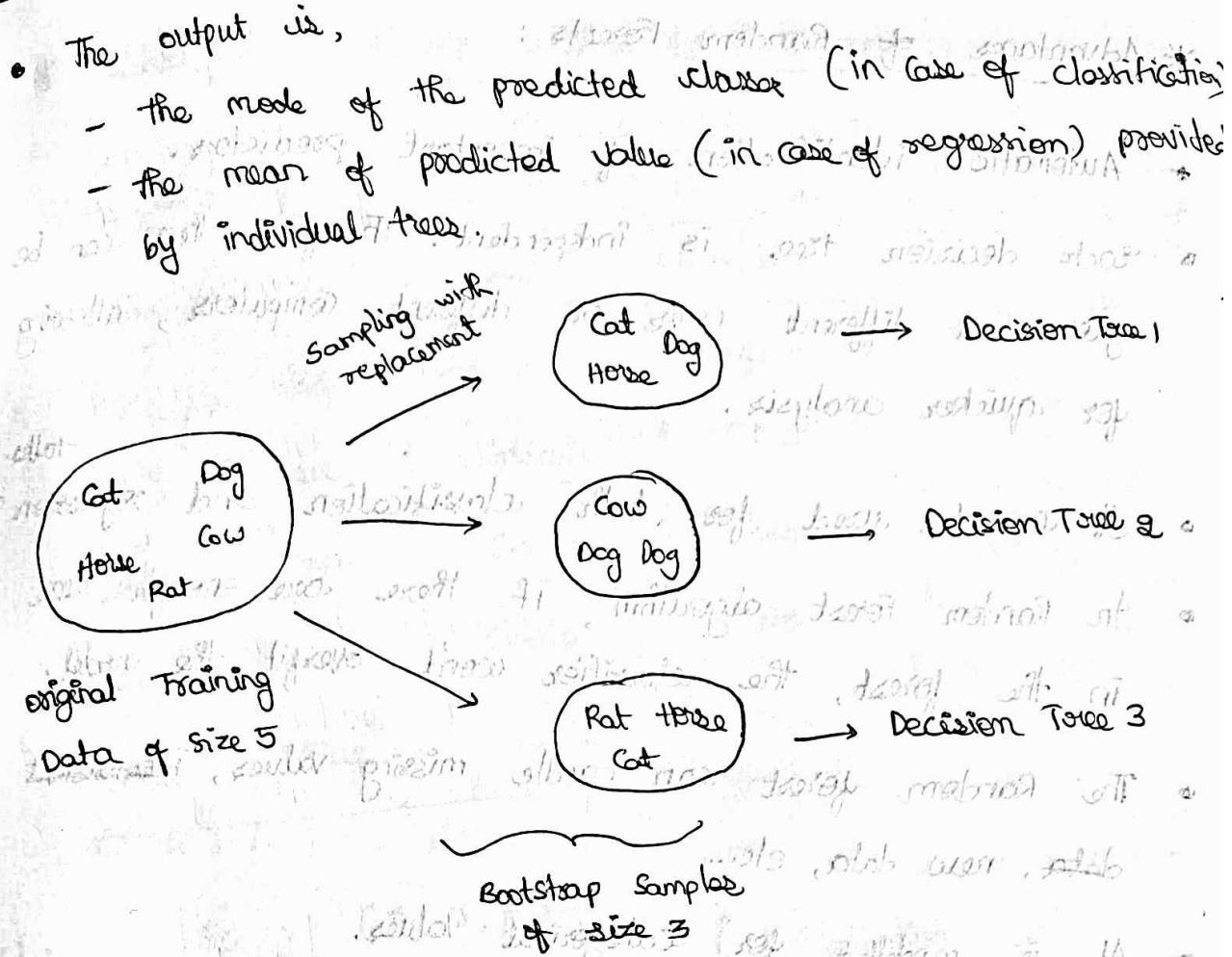
*- TYPES OF SAMPLING :-

- ① Simple Random sampling: There is an equal probability of selecting any particular item.
- ② Sampling without replacement: Once an object is selected, it is removed from the population.
- ③ sampling with replacement: A selected object is not removed from population.
- ④ stratified sampling: Partition the data set, and draw samples from each partition (proportionally i.e., approximately the same percentage of the data).

*- RANDOM FOREST :

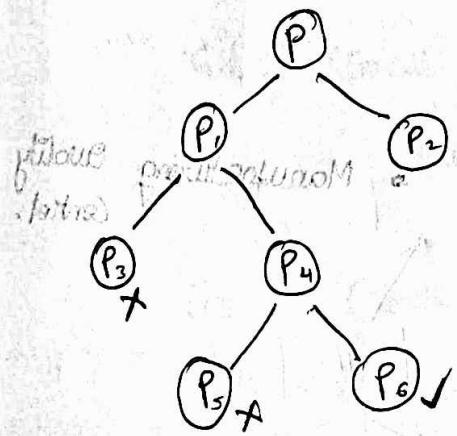
- Random Forest is a supervised learning algorithm.
- The forest it builds, is an ensemble of decision trees, usually trained with the bagging method.
- Working Rule: Random forest builds multiple decision trees and merges them together to get a more accurate and stable prediction.
- The larger the no. of trees, the more accurate the result.
- In a decision tree, each internal node represents a "test" on a attribute, each branch represents the outcome of test, and each leaf node represents a class label.

GUDI VARAPRASAD



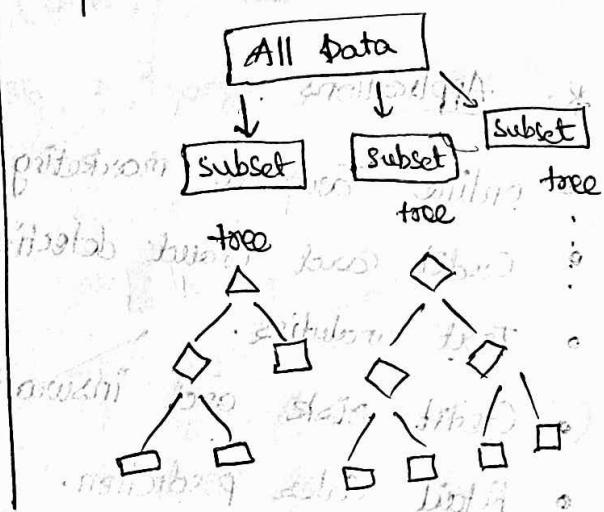
Decision Tree

- Decision tree builds classification models in the form of tree.
- It breaks down a dataset into smaller & smaller subsets.



Random Forest

- Random Forest is an ensemble classifier made of many decision trees.
- Ensemble models combine the results from different models.



* Advantages of Random Forests :

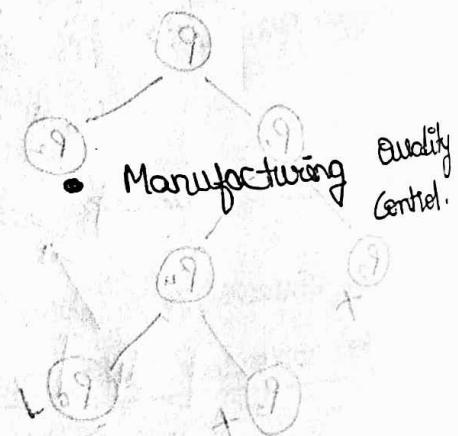
- Automatic identification of important predictors.
- Each decision tree is independent. Therefore trees can be grown on different cores or different computers, allowing for quicker analysis.
- It can be used for both classification and regression.
- In Random Forest algorithm, if there are enough trees in the forest, the classifier won't overfit the model.
- The Random forest can handle missing values, new data, etc..
- It is modeled for categorical values.

* Disadvantages :

- Suited for wide datasets with only a moderate no. of features. Breiman recommends the use of other tools for larger datasets.
- Large memory needed to store built models.
- Overfitting might be used/seen with noisy data.

* Applications :

- Online Targeted marketing.
- Credit Card Fraud detection.
- Text analytics.
- Credit risk and insurance risk.
- Retail Sales prediction.
- Biological & Medical Research



GUDI VARAPRASAD

* Candidate Elimination Algorithm Example :

Given Dataset,

Size	Color	Shape	Class / Label
Big	Red	Circle	No
Small	Red	Triangle	No
Small	Red	Circle	Yes
Big	Blue	Circle	No
Small	Blue	Circle	Yes

Set:

so : $(\emptyset, \emptyset, \emptyset)$ } Initialization

Go : $(?, ?, ?)$

S1 : $(\emptyset, \emptyset, \emptyset)$

G1 : $(\text{Small}, ?, ?)$ $(?, \text{Blue}, ?)$ $(?, ?, \text{Triangle})$

S2 : $(\emptyset, \emptyset, \emptyset)$

G2 : $(\text{Small}, \text{Blue}, ?)$ $(\text{Small}, ?, \text{Circle})$

$(?, \text{Blue}, ?)$ $(\text{Big}, ?, \text{Triangle})$ $(?, \text{Blue}, \text{Triangle})$

G3 : $(\text{Small}, \text{Red}, \text{Circle})$

G4 : $(\text{Small}, ?, \text{Circle})$

S4 : $(\text{Small}, \text{Red}, \text{Circle})$

G5 : $(\text{Small}, ?, \text{Circle})$

S5 : $(\text{Small}, ?, \text{Circle})$

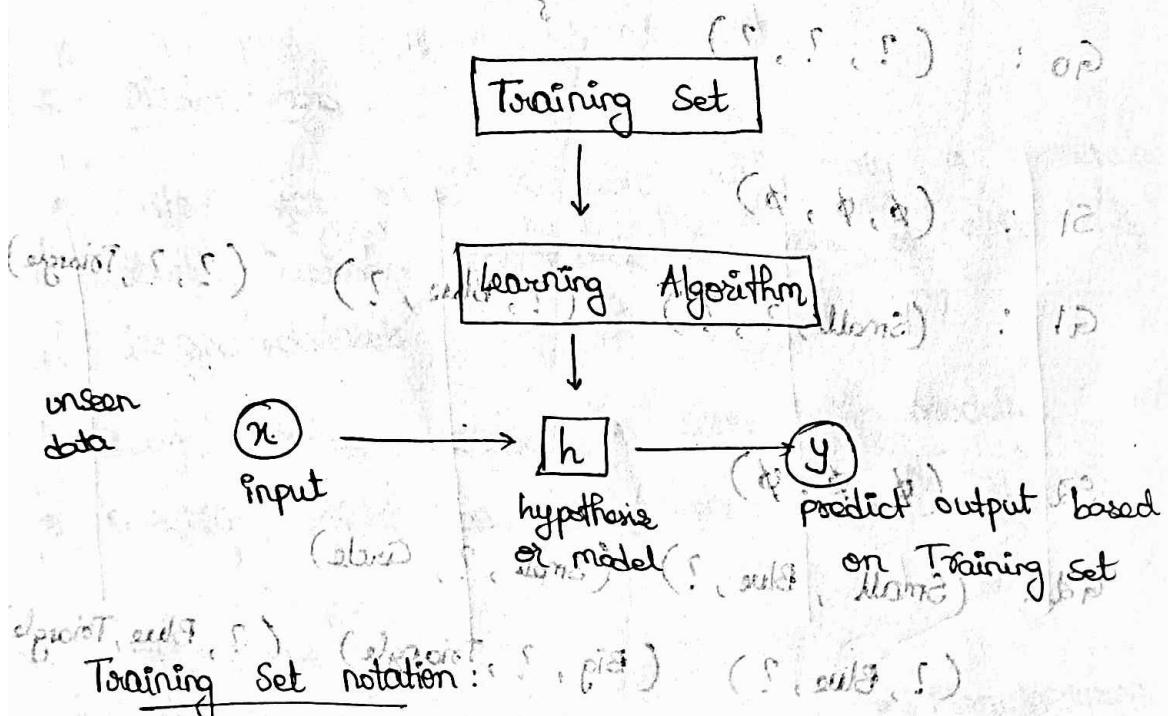
G6 : $(\text{Small}, ?, \text{Circle})$

* MODULE - 2

* LINEAR REGRESSION :

Labelled Data	Supervised Learning		Unsupervised learning
	Discrete	Continuous	clustering
	Classification	Regression	Dimensionality Reduction
	(x)	(y)	(x)

- In a simple way, regression is estimating unknown values from set of known values.



Training Set notation:

x = Input Variable / features

y = Output Variable / target Variable

(x, y) = One training example

$(x^{(i)}, y^{(i)})$ = i^{th} training example

$$y = h_{\theta}(x) = \theta_0 + \theta_1 x$$

θ_0, θ_1 are parameters/weights that control hypothesis

$h_{\theta}(x)$ - predicted output from Training set

y - actual output

- error = $|y - h_{\theta}(x)|$ → this error should be very low
Minimizing error is our goal & that is the best fit for

* Cost Function:

Idea: choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training example (x, y)

$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = J(\theta_0, \theta_1)$$

Simplified version: (1 parameter)

1. Hypothesis: $h_{\theta}(x) = \theta_1 x$, $\theta_0 = 0$

2. Parameters: θ_1

3. Cost Function: $J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

4. Goal: minimize $J(\theta_1)$ → we use method called Gradient Descent

* Gradient Descent:

repeat until convergence {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j=0, j=1$)

α : learning rate (step size)

GUDI VARAPRASAD

- Gradient descent can converge to a local minimum even with the learning rate α fixed.
- As we approach a local minimum, gradient descent automatically takes smaller steps. So, no need to decrease α over time.
- Linear regression model:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$
- "Batch": Each step of gradient descent uses all the training examples to make progress.

* MULTIPLE LINEAR REGRESSION

Ex:

$x^{(i)}$	Size in ft ² (x_1)	No. of Bedrooms (x_2)	No. of floors (x_3)	Age of home (x_4)	Price (y)
x^1	2104	5	1	45	460
x^2	1416	3	2	40	232
x^3	1534	3	2	30	315
x^4	852	2	1	36	178
:	:	:	:	:	:

Notation:

n = No. of features

$x^{(i)}$ = input features of i th training example

$x_j^{(i)}$ = value of feature j in i th training example

$$x_3^{(2)} = 2$$

$$x_3^{(4)} = 1$$

$$n = 4$$

$$x_4^{(3)} = 30$$

Now, Hypothesis $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

$$\Rightarrow x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n = \theta^T \cdot x$$

$$h_{\theta}(x) = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix} = [H]_{1 \times n} \cdot x_{n \times 1}$$

Now, Gradient Descent will be:

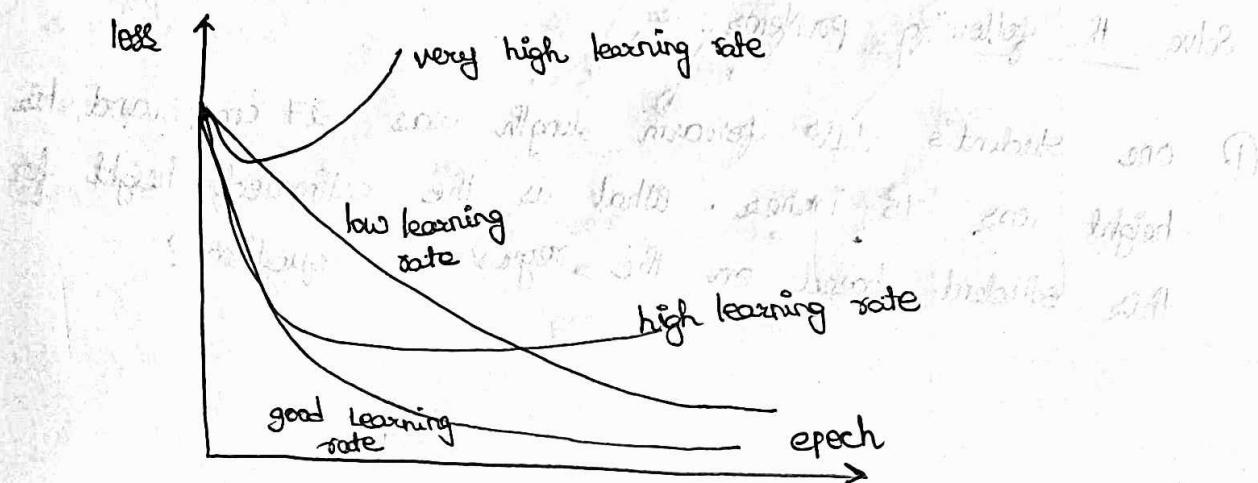
Repeat until convergence {

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x_j^i$$

{ } $n \geq 1$

* Gradient Descent : Learning rate

- Automatic convergence test
- α too small : slow convergence
- α too large : may not converge



*. POLYNOMIAL REGRESSION :

→ House price prediction :

- Features : frontage, depth.

- $h_{\theta}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$

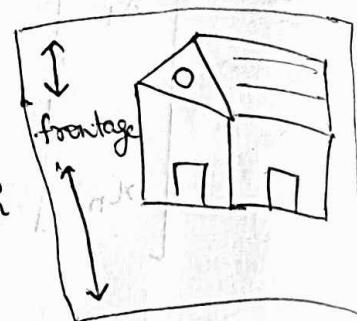
- From given features we can generate new feature called Area.

- Area, $x = \text{frontage} \times \text{depth}$

- $h_{\theta}(x) = \theta_0 + \theta_1 x$

- $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \dots$

$$\begin{aligned} &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3 \dots \\ (\text{a}) &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\sqrt{\text{size}}) + \theta_3 (\sqrt[3]{\text{size}}) \dots \end{aligned}$$



*. Example Problems : (LINEAR REGRESSION)

The question is based on a regression equation for 55 college students with $x = \text{left forearm length (cm)}$ & $y = \text{height}$. The forearm lengths ranged from 22 cm to 31 cm. The regression equation is $\hat{y} = 30.3 + 1.49x$.

Solve the following problems :

- ① One student's left forearm length was 27 cm, and his height was 75 inches. What is the estimated height for this student, based on the regression equation?

GUDI VARAPRASAD

Q. One student's left forearm length was 22 cm, and her height was 63 inches. What is the estimated height for this student, based on the regression equation?

③ What is the residual (Error) in problem 1, problem 2?

$$\Delta y = \text{residual} = |\text{given} - \text{estimated}|$$

① $\hat{y} = 30.3 + 1.49x$

Given, forearm length = 27 cm
height of student = 75 inch = 190.5 cm

height = $30.3 + 1.49$ (forearm length)

height = $30.3 + 1.49(27) = 68.1$ cm

$\hat{y} - y = \Delta y$
 $\Delta y = 190.5 - 68.1 = 122.4$ cm (residual error)

② $x = 22$ cm
 $y = 63$ inch = 160.02 cm

According to our given regression

$$\begin{aligned}\hat{y} &= 30.3 + 1.4 \times 22 \\ &= 61.1 \text{ cm}\end{aligned}$$

$\Delta y = 160.02 - 61.1 = 98.92$ (residual error)

③ Problem 1 residual error = 122.4 cm
Problem 2 residual error = 98.92 cm

* LOGISTIC REGRESSION :

- Regression used to fit a curve to data in which the dependent variable is binary, or dichotomous.
 - Ex: Coronary Heart Disease (CD) and age : In this study sampled individuals were examined for signs of CD (present = 1 / absent = 0) and the potential relationship between this outcome & their ages (yrs.) was considered.
 - Logistic Regression is basically a Binary Classification problem statement.
 - Sometimes, a (Linear Regression) straight line or best fit line can sometimes classify the data (yes/no). But if there is any outlier point that completely changes the best fit line, then the "Linear Regression" concept fails to solve binary classification problem. (high error rate)
- Hence, "LOGISTIC REGRESSION" comes into picture



$$x \rightarrow y_i: +ve = 1$$

$$\bullet \rightarrow y_i: -ve = -1$$

above the plane > 0

below the plane < 0

GUDI VARAPRASAD

case ①: $y_1 = +1$ (because it is x)
 $w^T x_1 > 0$ (above the plane)

$(y_1 \cdot w^T x_1) > 0 \Rightarrow$ This is easily classified.

case ②:
 $y_2 = -1$ (because it is \circ)

$w^T x_2 < 0$ (below the plane)

$(y_2 \cdot w^T x_2) > 0 \Rightarrow$ This is getting classified properly

case ③: $y_3 = -1$ (because it is \circ)

$w^T x_3 > 0$ (above the plane)

$(y_3 \cdot w^T x_3) < 0 \Rightarrow$ This (x_3, y_3) is incorrectly classified.

[OR] $y_4 = +1$ (because it is x)

$w^T x_4 < 0$ (below the plane)

$(y_4 \cdot w^T x_4) < 0 \Rightarrow$ This (x_4, y_4) point is not properly classified.

The best fit line that linearly separates (x, \circ) points make sure that the sum of all the points along with the distance should be maximum (cost function must be maximum)

$$\left\{ = \sum_{i=1}^n y_i w_i^T x_i \right\} \rightarrow \text{this should be maximum.}$$

Sigmoid Function: $F(z) = \frac{1}{1+e^{-z}}$

$$\text{where } z = y_i \cdot w^T x_i$$

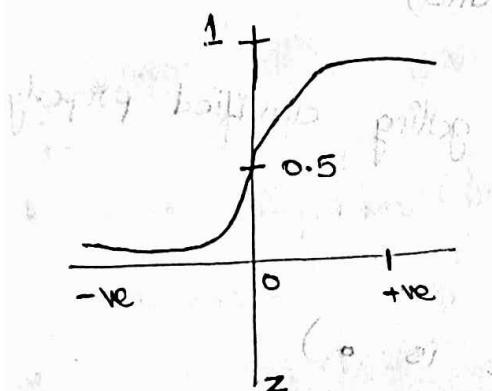
Here, to prevent chance of best fit line, (outliers)

$$\left\{ = \max \sum_{i=1}^n f(y_i \cdot w^T x_i) \right\} \text{ between 0 to 1}$$

any no. of outliers

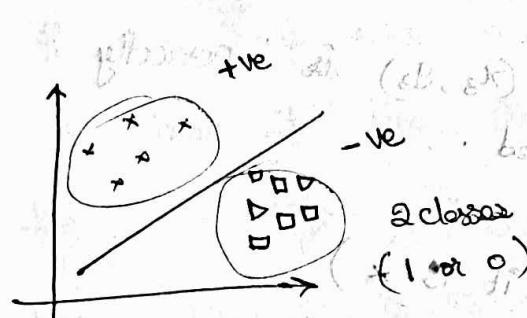
that may interrupt the

possibility of best fit line

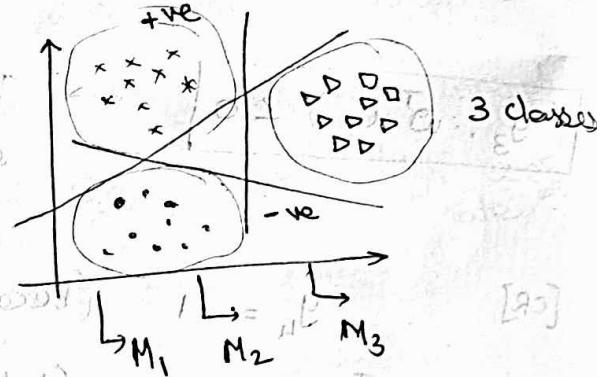


\Rightarrow The value comes out to be between 0 & 1.

*. Multiclass classification (one vs Rest) :



Binary class classification



Multiclass classification

Suppose, Example :

Inputs	f_1	f_2	f_3	Output		O_1	O_2	O_3
I_1		I_2	I_3	O_1		+1	-1	-1
I_4		I_5	I_6	O_2		-1	+1	-1
I_7		I_8	I_9	O_3		-1	-1	+1
I_{10}		I_{11}	I_{12}	O_1		+1	-1	-1
I_{13}		I_{14}	I_{15}	O_2		-1	+1	-1

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x

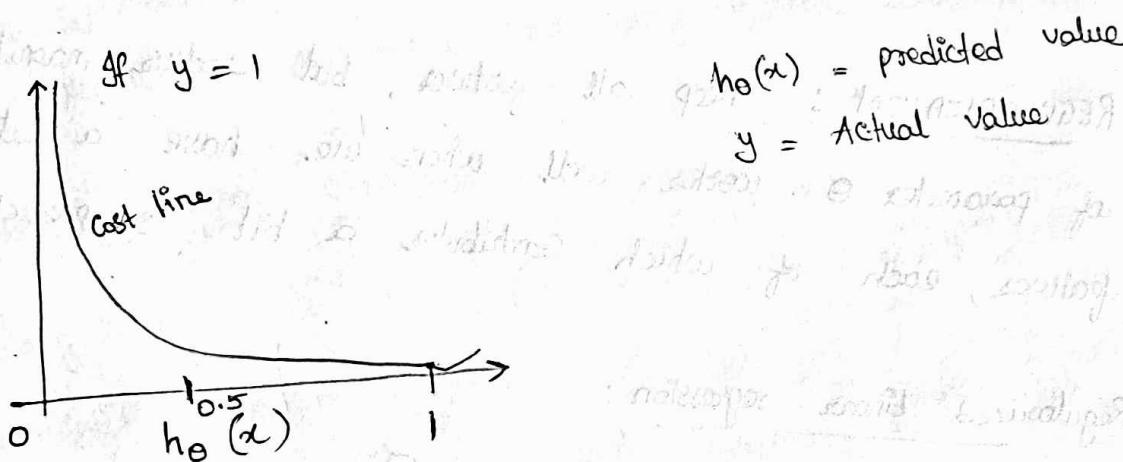
$$\text{Ex: If } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorsize} \end{bmatrix} \Rightarrow h_{\theta}(x) = 0.7$$

Tell the patient that 70% chance of tumor being malignant

Probability that $y=1$, given

* Logistic Regression Cost Function:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



* REGULARIZATION:

- Overfitting: If we have too many features, the learned hypothesis may fit the training set very well, but fail to generalize to new examples (predict prices on new examples).

$$\circ J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$$

underfitting - high bias, low variance.

overfitting - high variance, low bias.

- Bias - Variance Trade off:

- 1) Bias - difference between what you expect to learn and truth. It decreases with more complex model.
 - 2) Variance - difference between what you expect to learn and what you learn from a particular dataset. It measures how sensitive learner is to specific dataset. Increases with more complex model.
- 1. Reduce no. of features (not so easy)
- 2. Regularization. (better sometimes)

- REDUCTION: Keep all features, but reduce magnitude of parameters θ . Works well when we have a lot of features, each of which contributes a bit to predicting y .

- Regularized Linear regression:

$$\theta_j = \theta_j \left(1 - \alpha \frac{1}{m}\right) - \alpha \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

- Unregularized Linear regression:

$$\theta_j = \theta_j - \alpha \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

Here $1 - \alpha \frac{1}{m} < 1$: weight decay $(\text{as}) \quad \alpha \left(\frac{1}{m}\right) > 0$

* INSTANCE BASED LEARNING :

- Instance based learning is often termed "Lazy learning", as there is typically no "transformation" of training instances into more general "statements".
- Instead, the presented training data is simply stored and, when a new query instance is encountered, a set of similar, related instances is retrieved from memory and used to classify new query instance.
- Hence, instance-based learners never form an explicit general hypothesis regarding the target function. They simply compute the classification of each new query instance as needed.

* Nearest neighbor classifier:

- Training Data : $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(n)}, y^{(n)})$
- Learning : Do nothing
- Testing : $h(x) = y^k$, where $K = \arg \min_{\{1, 2, \dots, n\}} D(x, x^{(i)})$. take min distance
- k-NN Approach : (k must be odd)
- The simplest, most used instance-based learning algorithm is the K-NN algorithm.
- K-NN assumes that all instances are points in some n -dimensional space and defines neighbors in terms of distance (usually Euclidean in R -space).
- k is the no. of neighbors considered.
- Voronoi cell: Graphic depiction of k -nearest neighboring points.

GUDI VARAPRASAD

Algorithm :

1. Load Data.
2. Initialize K . (Rule of thumb)
3. For each sample in the training data,
 - 3.1 Calculate distance between query point \mathbf{E} current point.
 - 3.2 Add the distance \mathbf{E} index of example to an ordered collection.
4. Sort \rightarrow ordered collection of distances \mathbf{E} indexes from small to large.
5. Pick first K entries from sorted collection.
6. Get the labels of selected K entries.
7. If regression \rightarrow return mean of K labels.
 If classification \rightarrow return mode of K labels.

How to choose K ?

- Rule of thumb, $K = \sqrt{n}$, n is no. of examples in data
- In practice, $K=1$ is often used for efficiency, but can be sensitive to "noise"

CAT

Ex : We have data from the questionnaire survey and objective testing with two attributes (acid durability and strength) to classify whether a special paper tissue is good or not. Consider below 4 training samples.

$X_1 = \text{Acid Durability}$	$X_2 = \text{Strength}$	$Y = \text{Classification}$
1	7	BAD
1	4	BAD
3	4	GOOD
1	4	GOOD

GUDI VARAPRASAD

Now the factory produces a new paper tissue that passes the laboratory test with $x_1 = 3$, $x_2 = 7$. Guess the classification of new tissue.

Sol: Step 1: Initialize & define k (assume $k=3$)

Step 2: Compute distance between input sample and training sample

- Coordinate of the input sample is $(3, +)$.
- Instead of calculating the Euclidean distance, we calculate the squared Euclidean distance.

x_1	x_2	squared Euclidean distance	Rank	?
7	7	$(7-3)^2 + (7-7)^2 = 16$	3	Yes
7	4	$(7-3)^2 + (4-7)^2 = 25$	4	No
3	4	$(3-3)^2 + (4-7)^2 = 13$	1	Yes
1	4	$(1-3)^2 + (4-7)^2 = 13$	2	Yes

Total set of points

7	7	Bad
7	4	X
3	4	Good
1	4	Good

Majority = Good $\Rightarrow (3, +) \Rightarrow$ GOOD

\therefore The classification of new tissue = GOOD $(3, +)$

- Conclusions on KNN:
- KNN works well on many practical problems and is fairly noise tolerant (depending on value of k)

- K-NN is subject to the curse of dimensionality (i.e., presence of many irrelevant attributes).
- K-NN needs adequate distance measure.
- K-NN relies on efficient indexing.

Applications of KNN :

- Used in classification.
- Used to get missing values.
- Used in pattern recognition.
- Used in gene expression.
- Used to get 3D structure.
- Used to measure document similarity.

* DISTANCE - WEIGHTED KNN :

$$\text{Weights} = f(\text{distance}) = \frac{1}{\text{distance}}$$

Give weights based on distance & Predict class labels based on the weighted sum not on majority vote.

Ex:

Neighbors	True label	Distance	Weight	Sum of Weights
x_1	Positive	0.1	10	$10 + 3.3 = 13.3$
x_2	Positive	0.3	3.3	
x_3	Negative	1	1	$1 + 0.5 + 0.3$
x_4	Negative	2	0.5	
x_5	Negative	3	0.3	$= 1.8$

- KNN is a lazy algorithm because it doesn't learn a discriminative function from the training data but memorizes the training dataset instead. There is no training time in K-NN.

- * LOCALLY - WEIGHTED REGRESSION:
 - To overcome the problem of non-linearly separable data, Local Weighted Regression assigns weights to data to overcome the ~~parallel~~ problem.
 - This is computationally more expensive.
 - Finding / Assigning weights \rightarrow Kernel Smoothing [Weight $\propto \frac{1}{\text{dist}}$]

$D = a e^{-\frac{\|x - x_0\|^2}{2c^2}}$

x — each training input.
 x_0 — value we are predicting.

We construct a weight matrix (w)

- for each training input (x) and for the value we are trying to predict (x_0)

Weight matrix - diagonal matrix

$$B = (x^T w x)^{-1} \cdot x^T \cdot w \cdot y$$

B = Model Parameter

then prediction can be defined as

$$y = B x_0$$

$y \rightarrow$ prediction

Drawbacks:

1. Need to evaluate whole dataset everytime.
 2. Computation cost is more.
 3. Memory requirement is more.
- * RADIAL BASIS FUNCTIONS (Model - 5 with Neural Networks)

- * CASE BASED LEARNING : All instance based learners have 3 properties.
 - 1) They are lazy learners.
 - 2) classification is different from each instance.
 - 3) Instances are represented with n dimensional euclidean space.

- In case based reasoning, everything is considered as case and based on previous cases - we propose a solution.

- Instances are represented as symbols (not values).

CBR has 3 components :

1. Similarity functions or distance measures.
 2. Approximation / Adjustment of instances.
 3. Symbolic representation of instances.
- For modelling Case based Reasoning, we use Case based design tool (CADET) system that has 75 predefined libraries.

* EXAMPLES : (KNN)

The table shows lists of dataset that was used to create a nearest neighbour model that predicts whether it will be a good day to go surfing.

GUDI VARAPRASAD

ID	Wave Size	Wave Period	Wind Speed	Good Surf
1	6	15	5	yes
2	1	6	9	no
3	7	10	4	yes
4	7	12	3	yes
5	2	2	10	no
6	10	2	20	no

Assuming that the model uses Euclidean distance to find the nearest neighbour, what prediction will the model return for each of the following query instances.

ID	Wave Size	Wave Period	Wind Speed	Good Surf
a1	8	15	2	?
a2	8	2	18	?
a3	6	11	4	?

Sol: Given query instance a1 has
 Wave size = 8, Wave Period = 15
 Wind speed = 2

Finding the Euclidean distance between a1 & the given lists in dataset for ID1, ID2, ID3, ..., ID6.

GUDI VARAPRASAD

ID1 :

$$\begin{aligned} \text{wave size} &\Rightarrow (6-8)^2 \\ \text{wave Period} &\Rightarrow (15-15)^2 \\ \text{Wind speed} &\Rightarrow (5-2)^2 \end{aligned} \quad \left\{ \begin{array}{l} d = \sqrt{(-2)^2 + (0)^2 + (3)^2} \\ d = \sqrt{4+0+9} = \sqrt{13} \\ d = 3.605 \end{array} \right.$$

ID2 :

$$\begin{aligned} \text{wave size} &\Rightarrow (1-8)^2 \\ \text{wave Period} &\Rightarrow (6-15)^2 \\ \text{Wind speed} &\Rightarrow (9-2)^2 \end{aligned} \quad \left\{ \begin{array}{l} d = \sqrt{(-7)^2 + (-4)^2 + (7)^2} \\ d = \sqrt{49+81+49} = \sqrt{179} \\ d = 13.379 \end{array} \right.$$

ID3 :

$$\begin{aligned} \text{wave size} &\Rightarrow (7-8)^2 \\ \text{wave Period} &\Rightarrow (10-15)^2 \\ \text{Wind speed} &\Rightarrow (4-2)^2 \end{aligned} \quad \left\{ \begin{array}{l} d = \sqrt{(-1)^2 + (-5)^2 + (2)^2} \\ d = \sqrt{1+25+4} = \sqrt{30} \\ d = 5.477 \end{array} \right.$$

ID4 :

$$\begin{aligned} \text{wave size} &\Rightarrow (7-8)^2 \\ \text{wave Period} &\Rightarrow (12-15)^2 \\ \text{Wind speed} &\Rightarrow (3-2)^2 \end{aligned} \quad \left\{ \begin{array}{l} d = \sqrt{(-1)^2 + (-3)^2 + (1)^2} \\ d = \sqrt{1+9+1} = \sqrt{11} \\ d = 3.316 \end{array} \right.$$

ID5 :

$$\begin{aligned} \text{wave size} &\Rightarrow (2-8)^2 \\ \text{wave Period} &\Rightarrow (2-15)^2 \\ \text{Wind speed} &\Rightarrow (10-2)^2 \end{aligned} \quad \left\{ \begin{array}{l} d = \sqrt{(-6)^2 + (-13)^2 + (8)^2} \\ d = \sqrt{36+169+64} = 16.401 \end{array} \right.$$

ID6 :

$$\begin{aligned} \text{wave size} &\Rightarrow (10-8)^2 \\ \text{wave Period} &\Rightarrow (2-15)^2 \end{aligned} \quad \text{Wind speed} \Rightarrow (20-2)^2 \quad d = 22.293$$

GUDI VARAPRASAD

ID	Euclidean distance	Rank	Output
1	3.605	2	yes
2	13.379	4	no
3	5.477	3	yes
4	3.316	1	yes
5	16.401	5	no
6	22.293	6	no

Assuming $k = 3$, considering 3 nearest neighbors

$ID = 4 \rightarrow \text{yes}$
 $ID = 1 \rightarrow \text{yes}$
 $ID = 3 \rightarrow \text{yes}$

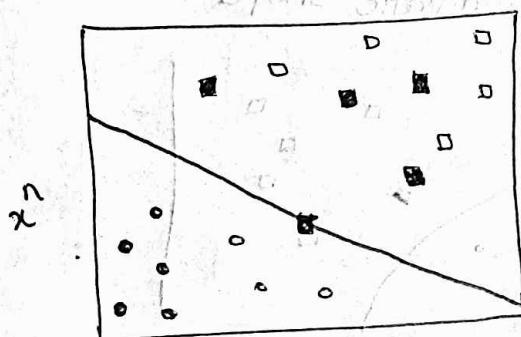
} majority = YES

\therefore New Query Instance at with Wave size = 8 FT
 Wave Period = 15 sec and Wind Speed = 2 MPH
 It will be a GOOD SURF (YES).

$$\textcircled{1} \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{i=1}^n \xi_i \quad \text{s.t. } \forall i = 1, \dots, n$$

$$\xi_i \geq 0 \quad \& \quad (\mathbf{w} \cdot \mathbf{x}_i + b) y_i - (1 - \xi_i) \geq 0$$

- The smaller error penalty C means that the decision boundary will have a larger margin



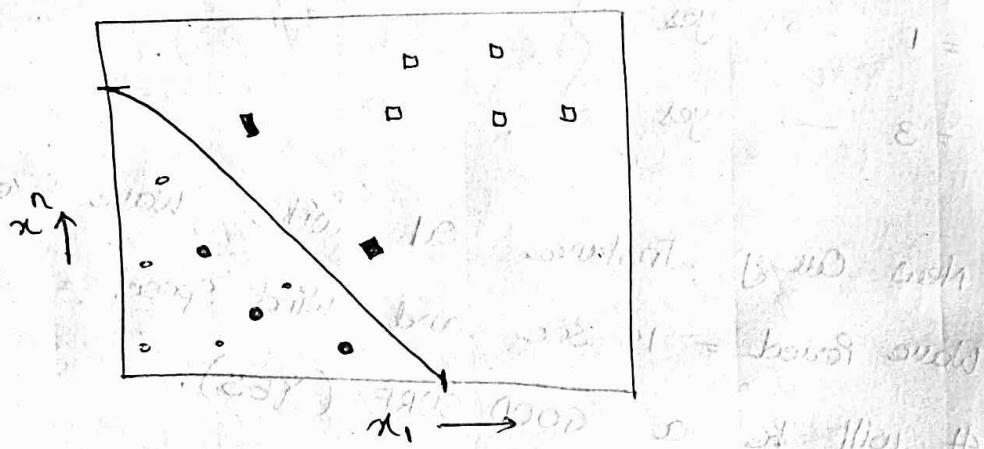
GUDI VARAPRASAD

$$② \min \frac{1}{2} w \cdot w + C \sum_{i=1}^n \xi_i \text{ s.t. } \forall i = 1, \dots, n$$

$$\xi_i \geq 0, \quad (w \cdot x_i + b) y_i - (1 - \xi_i) \geq 0$$

$C=1$

- From the distance of the support vectors to the decision boundary, we can infer that margin here is small because the primal optimization reduces the width of the margin as the penalty of margin as the penalty increases, it is Linear Kernel with the larger error penalty C .

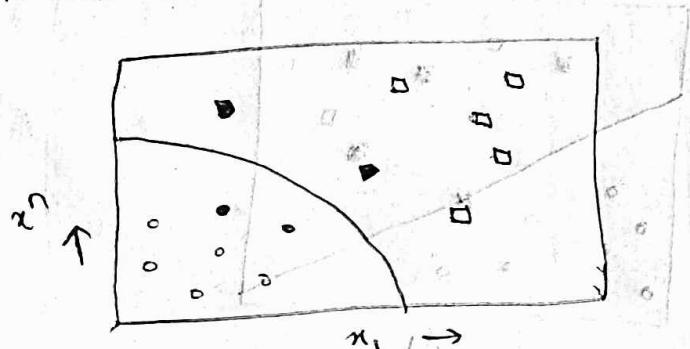


$$③ \max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

such that $\sum_{i=1}^n \alpha_i y_i = 0; \quad \alpha_i \geq 0 \quad \forall i = 1, \dots, n$

$$\text{where } K(u, v) = u \cdot v + (u \cdot v)^2$$

- It is sum of Linear Kernel & Quadratic Kernel so, it will be in quadratic shape



GUDI VARAPRASAD

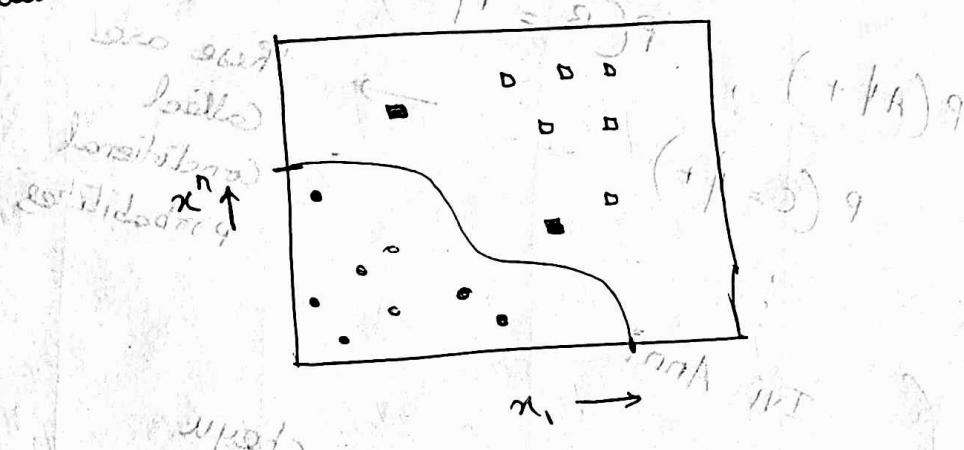
$$\textcircled{4} \quad \max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

such that $\sum_{i=1}^n \alpha_i y_i = 0$; $\alpha_i \geq 0 \quad \forall i = 1, \dots, n$

$$\text{where } K(u, v) = \exp\left(-\frac{\|u-v\|^2}{2}\right)$$

Gaussian RBF kernel, $K(u, v) = \exp(-\gamma \|u-v\|^2)$,

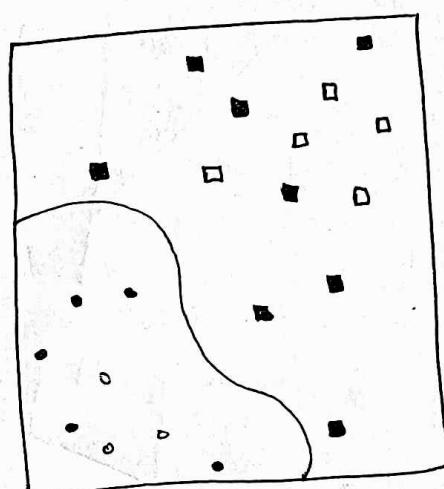
a smaller γ results a decision boundary with smoother curvature.



$$\textcircled{5} \quad \max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

such that, $\sum_{i=1}^n \alpha_i y_i = 0$; $\alpha_i \geq 0 \quad \forall i = 1, \dots, n$

$$\text{where } K(u, v) = \exp(-\gamma \|u-v\|^2)$$



- Larger γ will result in gaussian RBF kernel to have more distinguished curvature.

- A larger γ means Gaussian distribution has smaller variance & more distinguished curve.

- What is Bagging? (decrease variance)
 - An ensemble learning method that is commonly used to reduce variance within a noisy dataset. A random sample of data in a training set are selected with replacement.

- What is Boosting? (decrease bias)
 - Boosting creates a collection of predictors. It refers to a family of algorithms which converts weak learners to strong learner. It is an ensemble method for improving model predictions of given learning algorithm.

- What is pruning?
 - A data compression technique in ML that reduces the size of decision tree by removing sections of tree which are not required to classify instances.
- What is pre-pruning?
 - Early stopping the growth of tree before it classifies completely the training set.
- What is post-pruning?
 - Remove the sections of tree after it has completely grown.
- What is version space?
 - Intermediate space between hypothesis space & specific hypothesis. It is also hierarchical representation of knowledge.

- Disadvantages of Candidate Elimination Algo?
 - Fails for noisy / inconsistent data
 - Uses partially learned, sometimes fails to predict right hypothesis for new training sample.

GUDI VARAPRASAD

- Why we only consider particular rules for decision trees?
- A small change in data can cause a large change in the structure of the decision tree causing instability.
- For Decision Trees, sometimes calculations can go for more complex compared to other algorithms.

Because of certain rules, decision trees requires less effort for data preparation during pre-processing.

- Applications of Machine Learning?
 - Product Recommendation
 - Speech Recognition
 - Self driving Cars, Medical diagnosis
- What is Regression?
 - A statistical method to predict continuous outcomes based on one or more predictor variables.
- What is Linear Regression?
 - A supervised ML model in which it finds the best fit linear line between the independent and dependent variable.
- What is Cost function?

- Cost function (J) of Linear Regression is the Root Mean Squared error between predicted y values & true y value.

$$J = \frac{1}{n} \sum_{i=1}^n (pred_i - y_i)^2$$

θ_0 - Intercept

- What is hypothesis of Linear Regression?

$$h_\theta(x) = y = \theta_0 + \theta_1 \cdot x$$

(θ_0, θ_1 are parameters that control the hypothesis y)

GUDI VARAPRASAD

- What is gradient descent?
 - To update θ_1 & θ_2 values in order to minimize cost function (i.e. minimizing RMSE value) and achieving best fit line the model uses Gradient Descent.
- Here, the idea is to start with random θ_1, θ_2 values and then iteratively update them, reaching minimum cost.
- Why logistic Regression needed?
 - When our data has outliers, the best fit line predicted using linear regression may deviate and incorrectly predicts the output values.
 - The output in linear regression is sometimes > 1 or < 0 .
- What is sigmoid function?
 - $F(z) = \frac{1}{1 + e^{-z}}$; to prevent the output of linear regression from being > 1 or < 0 . This sigmoid function is an activation function that has range $(0, 1)$. Since we have to predict probability (lies between $(0, 1)$) this sigmoid function is right choice.
- What is Bias?
 - The difference between what you expect to learn & truth.
- What is variance?
 - The difference between what you expect to learn & what you learnt.

Small variance, high bias - underfitting (increase features)
High variance, small bias - overfitting (decision trees)
→ Regularization

GUDI VARAPRASAD

- What is Regularization?
 - Keep all features, but reduce parameter's magnitude. It discourages learning more complex or flexible model, to prevent overfitting.
- Define Logistic Regression?
 - Regression used to fit a curve of data in which the dependent variable is binary values or dichotomous.
- What is Confusion Matrix?
 - A summarized table of no. of correct & incorrect predictions obtained from by a classifier or classification model. It is a performance measurement for ML algorithm.

		1	0
Predicted values	1	TP	FP
	0	FN	TN

Actual values

TP - True Positive

FP - False Positive

FN - False Negative

TN - True Negative

TP - You predicted positive & it is True.

True Negative - You predicted negative & it is True.

False Negative - You predicted negative & it is False. (Type II error)

False Positive - You predicted positive & it is False. (Type I error)

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (\text{as high as possible}) \quad \text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{Total}} \quad (\text{high as possible})$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (\text{as high as possible})$$

GUDI VARAPRASAD

- What is K-NN? (An instance based learner)
 - A supervised ML algorithm that is used to solve both classification & regression problems.
 - It works by finding the distances between a query instance & all the data (training examples) nearest to query & then takes k-nearest examples & votes for most frequent output.
- Why is KNN called lazy Algorithm?
 - Because it does not learn a discriminative function from training data rather memorizes the training data & stores it to classify new training example.
- What is difference between KNN & weighted-KNN?
 - Weighted KNN is modified version of KNN. While taking majority vote, the nearest neighbors vary widely in their distance & closest neighbors more reliably indicate the class of object.
- What is Support Vector Machine?
 - The objective of SVM algorithm is to find a hyperplane in an N-dimensional space (N - the no. of features) that distinctly classifies the data points.
- What is hyperplane?
 - In an n-dimensional space that divides the space into two disconnected parts. called hyperplane has a flat, $n-1$ dimensional subset of
- What is margin?
 - The distance between the different classes. the hyperplanes drawn to differentiate

GUDI VARAPRASAD

- What is dimension of Hyperplane?
 - Hyperplanes are decision boundaries that help classify the data points. Data points falling on either side of the hyperplane can be attributed to different classes.
 - The dimension of hyperplane depends upon the no. of features. If the no. of input features is 2, then hyperplane is just a line. If the no. of input features is 3, then the hyperplane becomes 2-D plane.
- What are support vectors?
 - The data points that are closer to the hyperplane and influence the position & orientation of hyperplane. Using these, we can maximize the margin of the classifier.
- Explain Large Margin Intuition?
 - We take the output of the linear function and if that output is greater than 1, identified is in other class. The reinforcement range $[-1, 1]$ acts as margin (also called threshold value).
- What is hinge loss?
 - The loss function that helps maximize the margin i.e hinge loss.
- What is Kernel-trick?
 - A method of using linear classifier to solve a non-linear problem by mapping non-linear data into higher dimensional space.

- Why do we need Kernel trick?
- If the data is not linearly separable in 2-dimensional space, to build a linear classifier, we have to transform our data into 3D space while dealing with 2-D data.
- Types of Kernels?
 - 1) Linear Kernel -
 - 2) Gaussian Kernel -
 - 3) Polynomial Kernel -
 - 4) String Kernel -
 - 5) Chi-square Kernel -
- What is Linear Kernel?
 - Used when data is linearly separable, i.e. it can be separated using single line.
- What is Gaussian Kernel?
 - Gaussian is one such kernel that gives good linear separation in higher dimension for many non-linear problems.
- Logistic Regression vs. SVM?
 - $n = \text{no. of features}$
 - $m = \text{no. of training examples}$
 - n is large than m : use logistic regression or SVM without linear kernels.
 - 2: n is small, m is intermediate: use SVM with Gaussian kernel.
 - 3: n is small, m is large: Create/add more features, then use logistic regression (or linear SVM).

• Geometrical interpretation of margins -

• Margin = $x_1 - x_2$ - needed constraints - $\|x_1 - x_2\|$

GUDI VARAPRASAD

- SVM multiclass classification?
 - Use one vs all method. Train K SVMs, one to distinguish $y = i$ from the rest, get $\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(K)}$. Then pick class i with the largest $[\Theta^{(i)T} \cdot x]$
- SVM Regression?
 - Support Vector Regression is supervised learning algorithm that is used to predict discrete values. The basic idea behind SVR is to find best fit line.
 - In SVR, the best fit line is the hyperplane that has the maximum no. of points. Here, SVR tries to fit the best line within the distance between hyperplane & boundary line.
- Disadvantages of SVM?
 - 1) Not suitable for large data.
 - 2) Doesn't fit for data with noise (target class are overlapping).
 - 3) no. of features $>>$ no. of training samples, SVM fails.
 - 4) Long training time on datasets.
 - 5) choosing a "good" Kernel is not easy.
- What is Radial Basis?

$$K(x_1, x_2) = e^{-\frac{\|x_1 - x_2\|^2}{2\sigma^2}}$$

{ Find a non linear classifier or regression line.

6 - Variance of hyperparameter.

$\|x_1 - x_2\|$ - Euclidean distance between x_1, x_2 points.

MODULE - 4 :

* - NAIVE BAYES CLASSIFIER :

Ex : Given Data below. Classify this New Instance

(color = Green , Legs = 2 , Height = Tall , Smelly = No) Using

Naive Bayes classifier.

No	Color	Legs	Height	Smelly	Species
1	White	3	Short	Yes	M
2	Green	2	Tall	No	M
3	Green	3	Short	Yes	M
4	White	3	Short	Yes	M
5	Green	2	Short	No	H
6	White	2	Tall	No	H
7	White	2	Tall	No	H
8	White	2	Short	Yes	H

Sol:

$$P(M) = \frac{4}{8} = 0.5$$

$$P(H) = \frac{4}{8} = 0.5$$

Color	M	H
White	$\frac{2}{4}$	$\frac{3}{4}$
Green	$\frac{2}{4}$	$\frac{1}{4}$

Legs	M	H
2	$\frac{1}{4}$	$\frac{4}{4}$
3	$\frac{3}{4}$	$\frac{0}{4}$

Height	M	H
Tall	$\frac{3}{4}$	$\frac{1}{4}$
Short	$\frac{1}{4}$	$\frac{3}{4}$

GUDI VARAPRASAD

Smelly	M	H
Yes	$\frac{3}{4}$	$\frac{1}{4}$
No	$\frac{1}{4}$	$\frac{3}{4}$

$$g_{NB}^1 = \underset{v_j \in \{\text{yes, No}\}}{\operatorname{argmax}} P(v_j) \prod_i P(a_i | v_j)$$

$$g_{NB}^1(\text{yes}) = \frac{g_{NB}^1(\text{yes})}{g_{NB}^1(\text{yes}) + g_{NB}^1(\text{No})}$$

$$g_{NB}^1 = \underset{v_j \in \{\text{yes, No}\}}{\operatorname{argmax}} P(v_j)$$

$$g_{NB}^1(\text{No}) = \frac{g_{NB}^1(\text{No})}{g_{NB}^1(\text{yes}) + g_{NB}^1(\text{No})}$$

$$P(M | \text{New Instance}) = P(M) * P(\text{Color} = \text{Green} | M) *$$

$$P(\text{legs} = 2 | M) * P(\text{height} = \text{tall} | M) *$$

$$P(\text{smelly} = \text{No} | M)$$

$$P(M | \text{New Instance}) = 0.5 * \frac{2}{4} * \frac{1}{4} * \frac{3}{4} * \frac{1}{4} = 0.0117$$

$$P(H | \text{New Instance}) = P(H) * P(\text{Color} = \text{Green} | H) *$$

$$P(\text{legs} = 2 | H) * P(\text{height} = \text{tall} | H) *$$

$$P(\text{smelly} = \text{No} | H)$$

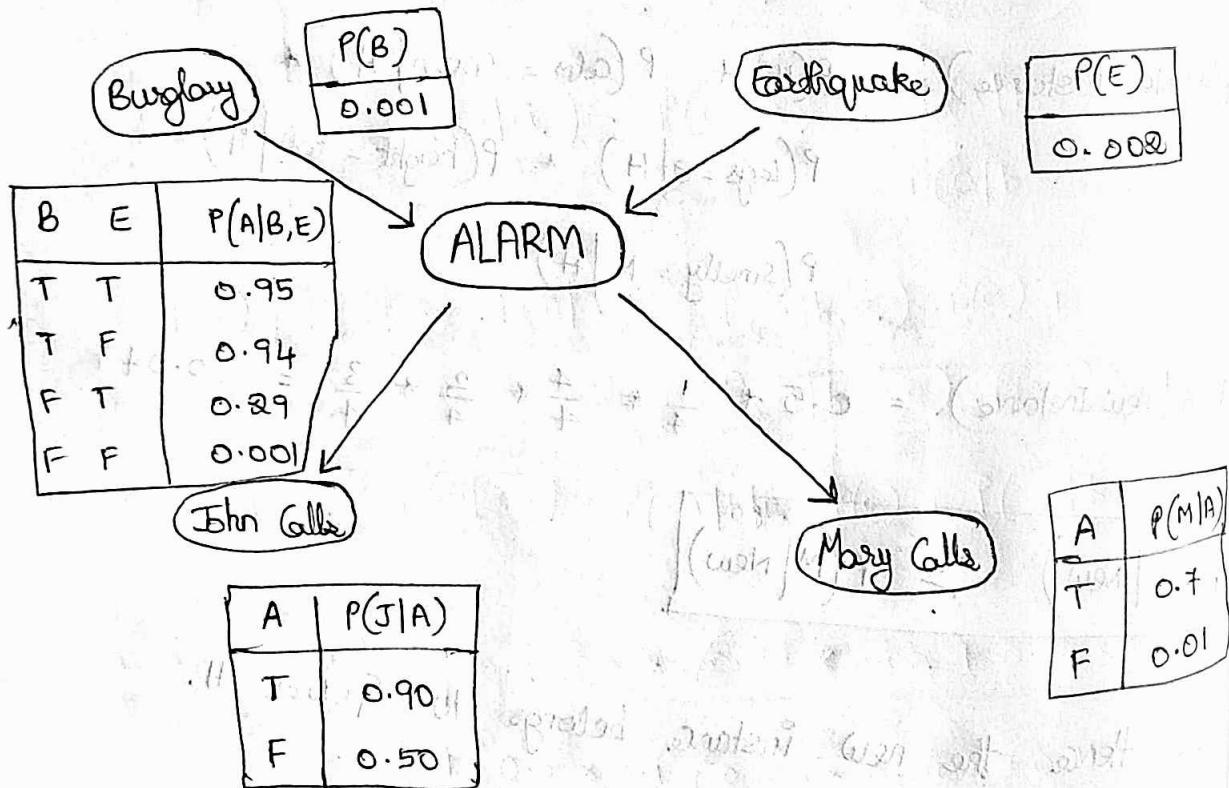
$$P(H | \text{New Instance}) = 0.5 * \frac{1}{4} * \frac{4}{4} * \frac{2}{4} * \frac{3}{4} = 0.047$$

$$P(H | \text{New}) > P(M | \text{New})$$

\therefore Hence the new instance belongs to Species H.

*. BAYESIAN BELIEF NETWORKS :

- You have a new burglar alarm installed at home.
- It is fairly reliable at detecting burglary, but also sometimes responds to minor earthquakes.
- You have two neighbors, John & Merry, who promised to call you at work when they hear the alarm.
- John always calls when he hears the alarm, but sometimes confuses telephone ringing with the alarm and calls too.
- Merry likes loud music and sometimes misses the alarm. Given the evidence of who has or has not called, we would like to estimate the probability of a burglary.



① What is the probability that the alarm has sounded that but neither a burglary nor an earthquake has occurred, and both John & Merry call?

Sol:

$$P(j \wedge m \wedge a \wedge \sim b \wedge \sim e) = P(j|a) \cdot P(m|a) \cdot P(a|\sim b, \sim e) \cdot P(\sim b) \cdot P(\sim e)$$

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

② What is the probability that John call?

Sol:

$$P(j) = P(j|a) \cdot P(a) + P(j|\sim a) \cdot P(\sim a)$$

$$= P(j|a) \cdot \{ P(a|b,e) * P(b,e) + P(a|\sim b,e) * P(\sim b,e) + P(a|b,\sim e) * P(b,\sim e) + P(a|\sim b,\sim e) * P(\sim b,\sim e) \}$$

$$+ P(j|\sim a) \cdot \{ P(\sim a|b,e) * P(b,e) + P(\sim a|\sim b,e) * P(\sim b,e) + P(\sim a|b,\sim e) * P(b,\sim e) + P(\sim a|\sim b,\sim e) * P(\sim b,\sim e) \}$$

$$= 0.90 * 0.00252 + 0.05 * 0.9974 = 0.0521$$

③ What is the probability that there is a burglary given that John and Merry calls?

Sol:

Suppose we have given for the evidence variables E_1, E_2, \dots, E_m their values e_1, \dots, e_m , and we want

GUDI VARAPRASAD

to predict whether the query variable X has the value x or not.

For this, we compute and compare the following:

$$P(x | e_1, \dots, e_m) = \frac{P(x, e_1, \dots, e_m)}{P(e_1, \dots, e_m)} = \alpha \cdot P(x, e_1, \dots, e_m)$$

$$P(\sim x | e_1, \dots, e_m) = \frac{P(\sim x, e_1, \dots, e_m)}{P(e_1, \dots, e_m)} = \alpha \cdot P(\sim x, e_1, \dots, e_m)$$

$$\Rightarrow \alpha = \frac{1}{P(x, e_1, \dots, e_m) + P(\sim x, e_1, \dots, e_m)}$$

$$P(b | j, m) = \alpha \cdot P(b) \leq_a P(j | a)$$

$$P(m | a) \leq_e P(a | b, e) \cdot P(e)$$

$$= \alpha \cdot P(b) \leq_a P(j | a) \cdot P(m | a) \cdot \{ P(a | b, e) \cdot P(e) + P(a | b, \sim e) \cdot P(\sim e) \}$$

$$= \alpha \cdot P(b) [P(j | a) \cdot P(m | a) \{ P(a | b, e) \cdot P(e) + P(a | b, \sim e) \cdot P(\sim e) \}]$$

$$+ P(j | \sim a) \cdot P(m | \sim a) \{ P(\sim a | b, e) \cdot P(e) + P(\sim a | b, \sim e) \cdot P(\sim e) \}$$

$$= \alpha * (0.001) * [0.9 \times 0.7 \times (0.95 * 0.02 + 0.94 * 0.98) + 0.5 \times 0.1 (0.05 \times 0.002 + 0.71 \times 0.98)]$$

$$= \alpha * 0.00059$$

GUDI VARAPRASAD

similarly,

$$\begin{aligned}
 P(\sim b | j, m) &= \alpha \cdot P(\sim b) \leq_a P(j|a) \cdot P(m|a) \leq_e P(a|\sim b, e) \\
 &= \alpha \cdot P(\sim b) \leq_a P(j|a) \cdot P(m|a) \left\{ P(a|\sim b, e) \cdot P(e) + \right. \\
 &\quad \left. P(a|\sim b, \sim e) \cdot P(\sim e) \right\} \\
 &= \alpha \cdot P(\sim b) \left[P(j|a) \cdot P(m|a) \left\{ P(a|\sim b, e) \cdot P(e) + \right. \right. \\
 &\quad \left. \left. P(a|\sim b) \cdot P(\sim e) \right\} \right] \\
 &\quad + P(j|a) \cdot P(m|\sim a) \left\{ P(\sim a|\sim b, e) \cdot P(e) + \right. \\
 &\quad \left. P(\sim a|\sim b, \sim e) \cdot P(\sim e) \right\} \\
 &= \alpha * 0.999 * \left\{ (0.9 * 0.7 (0.29 * 0.02 + 0.01 * 0.998) \right. \\
 &\quad \left. + 0.05 * 0.01 (0.71 * 0.02 + 0.999 * 0.998) \right\} \\
 &= \alpha * 0.0015
 \end{aligned}$$

We get, $P(b | j, m) = \alpha \cdot P(b, j, m)$

~~we get~~, $P(\sim b | j, m) = \alpha \cdot P(\sim b, j, m)$

And, $P(b, j, m) = 0.00059$

~~from~~ $P(\sim b, j, m) = 0.0015$

We know that, $\alpha = \frac{1}{P(b, j, m) + P(\sim b, j, m)} = \frac{1}{(0.00059 + 0.0015)}$

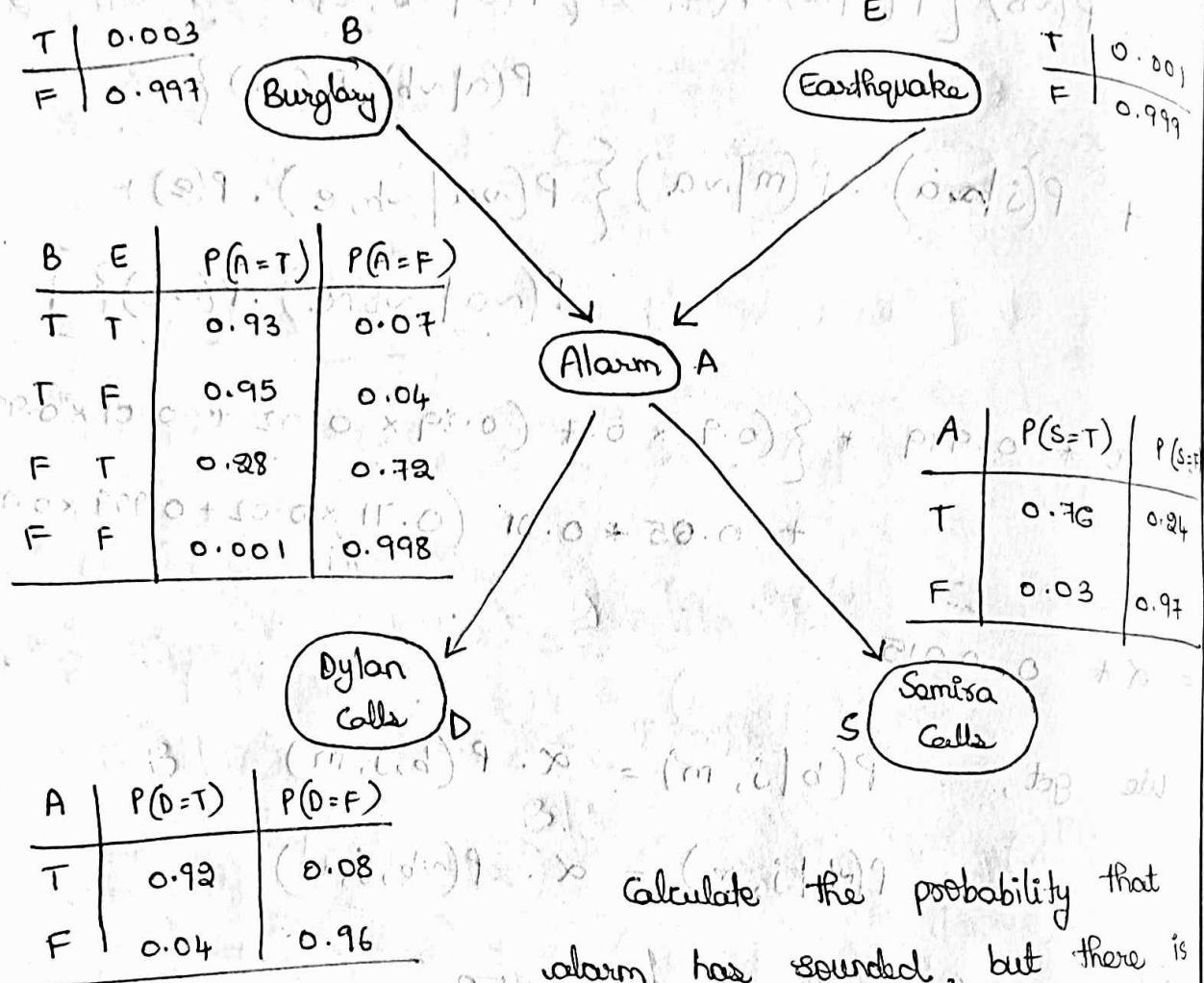
$\alpha = \frac{1}{0.00209} = 478.4688 \approx 478.5$

GUDI VARAPRASAD

$$\rightarrow P(b|j, m) = 478.5 \times 0.00059 = 0.28$$

$$P(\neg b|j, m) = 478.5 \times 0.0015 = 0.72$$

Example 2 : (Consider same Burglar Alarm Case before)



Calculate the probability that alarm has sounded, but there is neither a burglary, nor an earthquake occurred, and Dylan & Samira both called Munish.

Sol: Given Probabilities are:

$$P(B = \text{True}) = 0.003$$

$$P(B = \text{False}) = 0.997$$

$$P(E = \text{True}) = 0.001$$

$$P(E = \text{False}) = 0.999$$

GUDI VARAPRASAD

Probability of Alarm : (A)

Burglary (B)	Earthquake (E)	$P(A=T)$	$P(A=F)$
T	T	0.93	0.07
T	F	0.95	0.04
F	T	0.28	0.72
F	F	0.001	0.998

Probability of Dylan Calls : (D)

Alarm (A)	$D(D=True)$	$\sim D(D=False)$
T	0.92	0.08
F	0.04	0.96

Probability of Samira Calls : (S)

Alarm (A)	$S(S=True)$	$\sim S(S=False)$
T	0.76	0.24
F	0.03	0.97

Now, Alarm sounded $\Rightarrow A = \text{True} \Rightarrow P(A=F)$
 No Burglary $\Rightarrow B = \text{False} \Rightarrow P(B=F)$
 No Earthquake $\Rightarrow E = \text{False} \Rightarrow P(E=F)$

Dylan called $\Rightarrow D = \text{True} \Rightarrow P(D=T)$

Samira called $\Rightarrow S = \text{True} \Rightarrow P(S=T)$

Now $P(A, \sim B, \sim E, D, S) = ?$

$$(280.0 \times 0.76) + (280.0 \times 0.03) + (280.0 \times 0.001)$$

$$\Rightarrow = P(D|A) \cdot P(S|A) \cdot P(A|\sim B, \sim E) \cdot P(\sim B) \cdot P(\sim E)$$

$$= 0.92 \times 0.76 \times 0.001 \times 0.997 \times 0.999$$

$$= 0.0006964 \approx \underline{\underline{0.0007}}$$

* In an oximeter manufacturing factory, Wing A, B, C produce 25%, 36% and 42% oximeters respectively. Out of total 5.2%, 3.5%, 2.5% are defective ones. An oximeter is drawn at random from product. If oximeter drawn is found to be defective. What is the probability it is manufactured by wing B machine.

Sol: Let A_1 be event that oximeter manufactured by wing A. & B_2 for wing B, C be for wing C

$$P(B|E) = ?$$

E : Defective oximeter (A) manuf.

$B|E$: Defective oximeter given that is manufactured by wing B.

Given,

$$P(A) = 25\% = 0.25$$

$$P(E|A_1) = 0.052$$

$$P(B) = 36\% = 0.36$$

$$P(E|B) = 0.035$$

$$P(C) = 42\% = 0.42$$

$$P(E|C) = 0.025$$

$$P(B|E) = \frac{P(B) \times P(E|B)}{P(A) \times P(E|A) + P(B) \times P(E|B) + P(C) \times P(E|C)}$$

$$= \frac{0.36 \times 0.035}{(0.25 \times 0.052) + (0.36 \times 0.035) + (0.42 \times 0.025)} = \frac{126}{361} = 0.349$$

GUDI VARAPRASAD

* Rishav has joined a new company and his age is 48 years and he is getting a salary amount of Rs. 142,000/- In that company, the salary amount is labeled as High (H) or Low (L) based on the employee's age as shown below.

S.No. (ID)	Employee Name	Age	Salary	Salary Type (High/Low)
1	Amit	33	150000	H
2	Bikash	48	220000	H
3	Seema	60	100000	H
4	Pratik	40	62000	L
5	Arbab	33	95000	H
6	Ridhi	52	18000	L
7	Biju	35	120000	H
8	Rakesh	20	20000	L
9	Kiran	45	80000	H
10	Jhuma	35	60000	L
11	Piyus	25	40000	L

Use K-NN algorithm and find the Rishav's salary

label for, $K=5$ & $K=7$?

New Instance $\Rightarrow (12, \text{Rishav}, 48, 142000, ?)$

GUDI VARAPRASAD

For ID 1 :

$$\left. \begin{array}{l} \text{Age} \Rightarrow (33 - 48)^2 \\ \text{Salary} \Rightarrow (150000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{225 + 64000000}$$

$$d = 8000.014069$$

For ID 2 :

$$\left. \begin{array}{l} \text{Age} \Rightarrow (48 - 48)^2 \\ \text{Salary} \Rightarrow (220000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{0 + 78000000}$$

$$d = 28000.007800$$

For ID 3 :

$$\left. \begin{array}{l} \text{Age} \Rightarrow (60 - 48)^2 \\ \text{Salary} \Rightarrow (60000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{144 + 1764000000}$$

$$d = 42000.00171$$

For ID 4 :

$$\left. \begin{array}{l} \text{Age} \Rightarrow (40 - 48)^2 \\ \text{Salary} \Rightarrow (62000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{64 + 6400000000}$$

$$d = 80000.0004$$

For ID 5 :

$$\left. \begin{array}{l} \text{Age} \Rightarrow (23 - 48)^2 \\ \text{Salary} \Rightarrow (95000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{625 + 2209000000}$$

$$d = 47000.00665$$

For ID 6 :

$$\left. \begin{array}{l} \text{Age} \Rightarrow (52 - 48)^2 \\ \text{Salary} \Rightarrow (18000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{16 + 15376000000}$$

$$d = 124000.0001$$

GUDI VARAPRASAD

For ID 7:

$$\left. \begin{array}{l} \text{Age} \Rightarrow (35-48)^2 \\ \text{Salary} \Rightarrow (120000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{169 + 169000000}$$

$$d = 13000.0065$$

For ID 8:

$$\left. \begin{array}{l} \text{Age} \Rightarrow (20-48)^2 \\ \text{Salary} \Rightarrow (20000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{784 + 14884000000}$$

$$100d = 122000.0032$$

For ID 9:

$$\left. \begin{array}{l} \text{Age} \Rightarrow (45-48)^2 \\ \text{Salary} \Rightarrow (80000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{9 + 3844000000}$$

$$100d = 62000.00007$$

For ID 10:

$$\left. \begin{array}{l} \text{Age} \Rightarrow (35-48)^2 \\ \text{Salary} \Rightarrow (60000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{169 + 6724000000}$$

$$d = 82000.00103$$

For ID 11:

$$\left. \begin{array}{l} \text{Age} \Rightarrow (25-48)^2 \\ \text{Salary} \Rightarrow (40000 - 142000)^2 \end{array} \right\} \Rightarrow d = \sqrt{529 + 10404000000}$$

$$d = 102000.0026$$

GUDI VARAPRASAD

ID	Euclidean distance	Rank	Output
1	8000.014062	1	H
2	78000	6	H
3	42000.00171	3	H
4	80000.0004	7	L
5	47000.00665	4	H
6	124000.0001	11	L
7	13000.0065	2	H
8	122000.0032	10	L
9	62000.00007	5	H
10	22000.00103	8	L
11	102000.0026	9	L

i. For $K=5 \Rightarrow \{ID = 1, 7, 3, 5, 9\}$

Corresponding outputs are H, H, H, H, H

Majority vote = H for $K=5$

∴ New Query Instance : (12, Rishav, 48, 142000, H)

So, Rishav gets High Salary

ii. For $K=7 \Rightarrow \{ID = 1, 7, 3, 5, 9, 2, 4\}$

Corresponding outputs are H, H, H, H, H, H, L

Majority vote = H for $K=7$

∴ New Query Instance : (12, Rishav, 48, 142000, H)

So, Rishav gets High Salary.

* EM Algorithm : (Expectation - Maximization)

- used to find latent variable (not directly derivable)
- basic for many unsupervised clustering algorithm.

① steps involved in EM-algorithm:

① Initially, a set of initial values are considered.

② A set of incomplete data is given to system.

③ Next step - expectation step (E-step). Here we use

observed data to estimate or guess the values of missing/incomplete data.

④ Maximisation step or M-step. Here we use the complete data generated in preceding E-step to update the values.

⑤ We check if values are converging or not.
If converging - stop
Else, Repeat step 2 & step 3 till the convergence occurs.

② Usage:

- 1) used to fill missing data.
- 2) used for unsupervised clustering.
- 3) used to discover values of latent variables.

③ Advantages:

- + With each iteration, likelihood increases.
- + E-step & M-step are easy to implement.

④ Disadvantages:

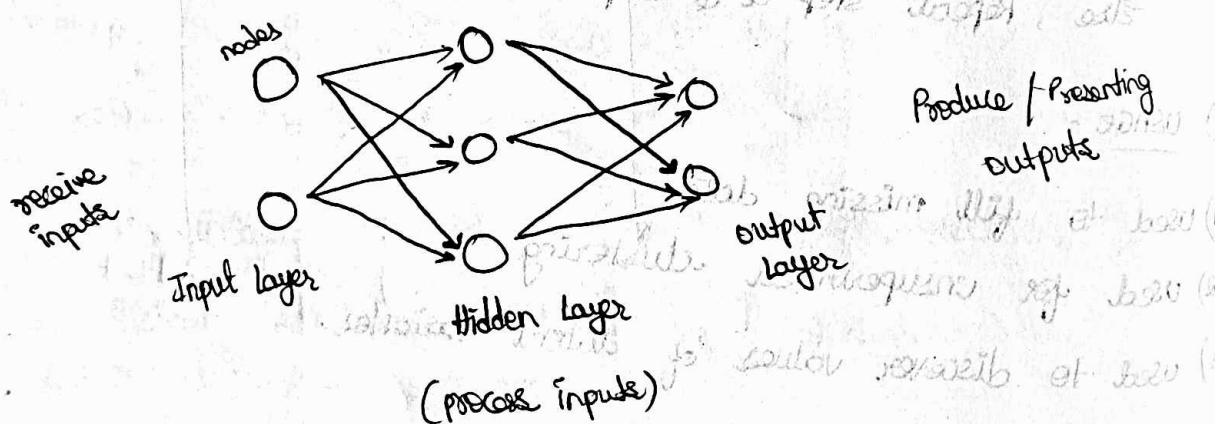
- slow convergence.
- makes convergence to local optimal only.

MODULE - 5

- Human Brain - millions of neurons
 - Biological - Neurons, formal term \rightarrow Nodes
 - Natural functions are done artificially.
- Diagram of a node:
-
- $$y = f(\sum x_i w_i)$$
- Labels for diagram:
- x_i : Input signal
 - w_i : weight of corresponding input signal
 - Σ : summation function
 - f : Activation function
 - y : output function
- A node has 2 parts:

Representation of Artificial Neural Networks: (3 parts)

1. Input layer 2. Hidden Layer 3. Output Layer

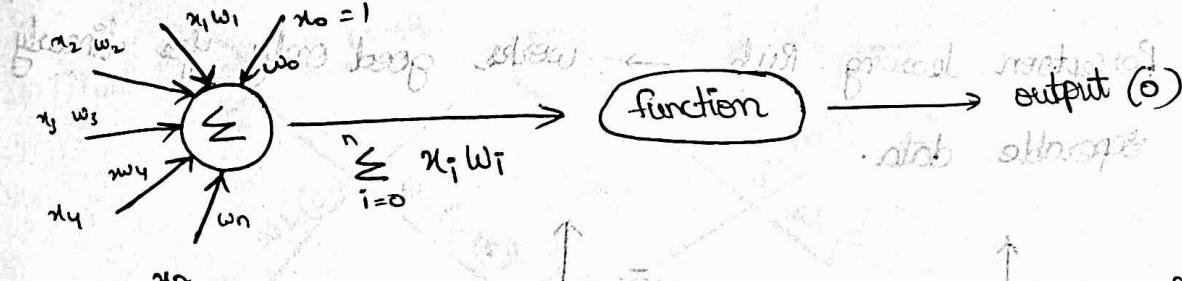


- Appropriate problems for Neural networks Learning:
- 1) Instances have many attribute value pairs.
 - 2) Target function has discrete values, continuous values or combination of both.

3. Training examples with errors or missing values.
4. Long training times are acceptable.
5. Fast evaluation of the target function learnt.
6. Ability for humans to understand the target function learnt by machine is not important.

* PERCEPTION

- basic unit used to build ANN.
 - takes real valued input, calculates linear combination of these inputs and generates output.
- output = 1 if the result \geq threshold
 -1 otherwise result \leq threshold



o : actual output

t : target output

if $actual = target \rightarrow$ weights are fixed
 otherwise \rightarrow changed.

How to change weights?: Initially, random weights. Later, keep on applying iterations and check if $(o = t)$.

Formula: $w_i \leftarrow w_i + \Delta w_i$

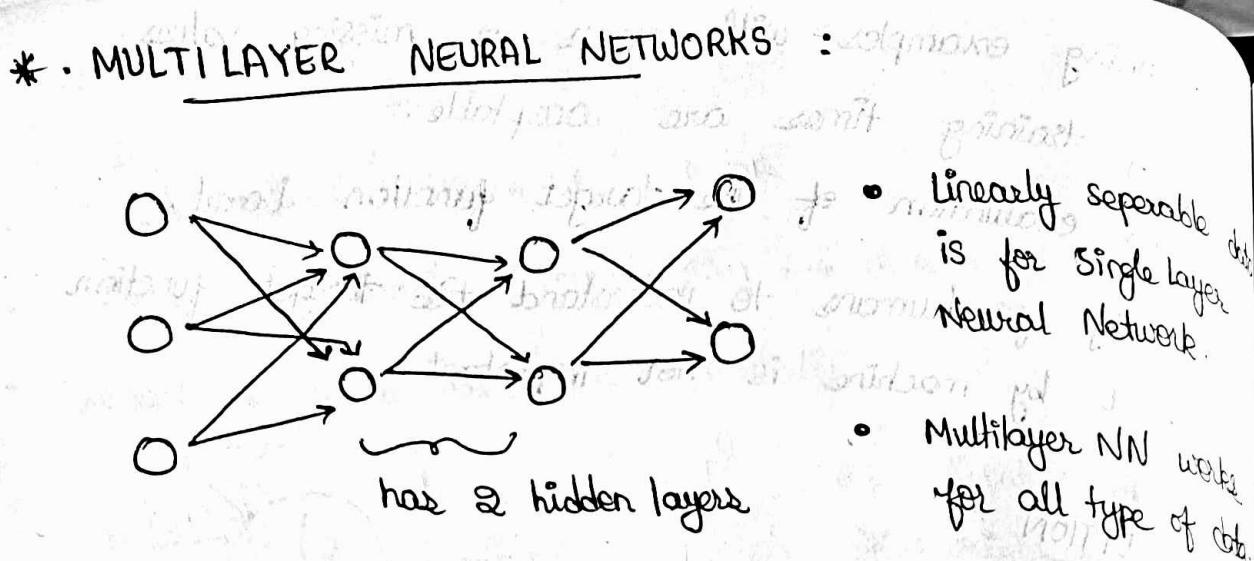
$$\Delta w_i = n \cdot (t - o) \cdot x_i$$

n: learning rate
 t: target output
 o: actual output

x_i : input associated with w_i

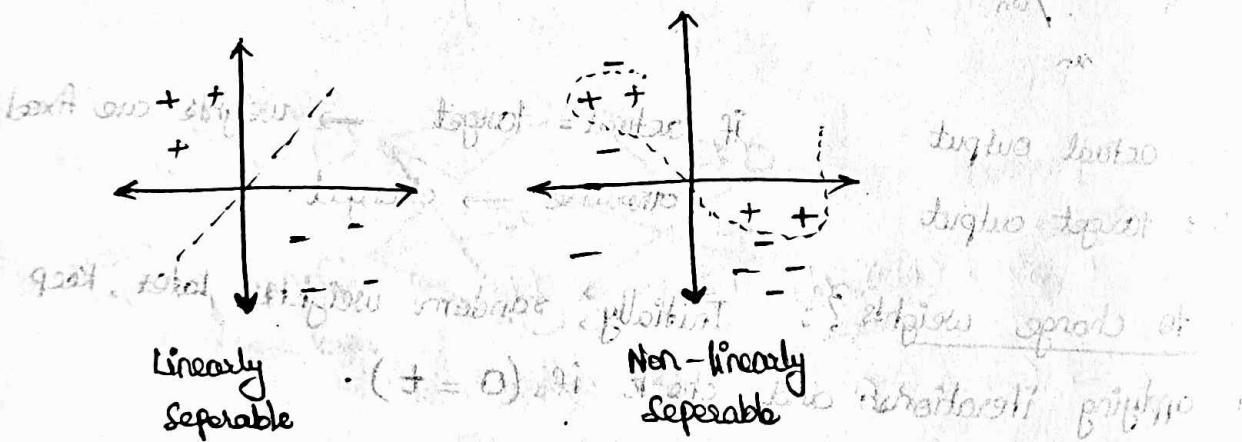


GUDI VARAPRASAD

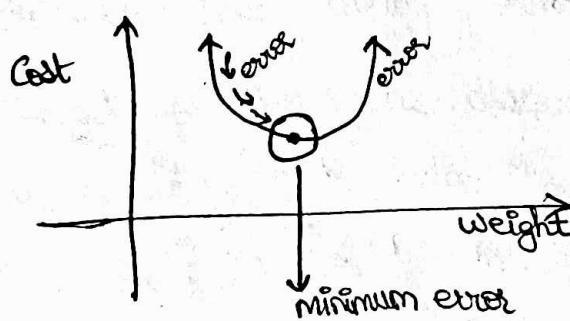


- output is not matching \rightarrow go back & modify weights.
- Back propagation is possible.

- * GRADIENT DESCENT AND DELTA RULE
- Delta Rule \rightarrow used for non-linearly separable data
 - Perceptron Learning Rule \rightarrow works good only for linearly separable data.



- Main idea behind delta rule is, it uses gradient descent rule and finds out the best weight.



Modifying weights: $w_i = w_i + \Delta w_i$ Here, $\Delta \Sigma(w)$ is derivative of error w.r.t weights. Also $\Delta w_i = -\eta \Delta \Sigma(w)$ called as gradient.

$$\Delta w_i = -n \sum (t_d - o_d) x_d \quad \partial \text{PE} \cdot o = (o_d)_d$$

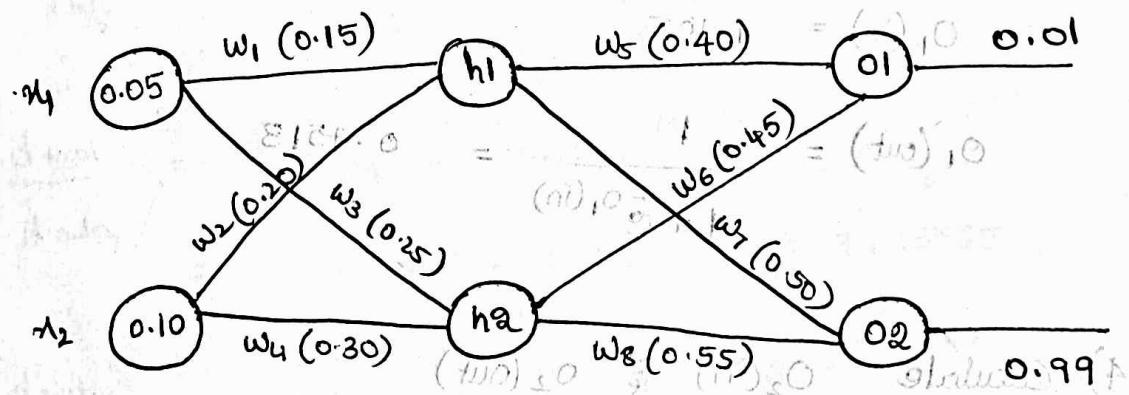
*. BACK PROPAGATION ALGORITHM:

- When error occurs, we go in backward direction & fix the error. output \rightarrow hidden \rightarrow input layer

Part 1: Calculate forward propagation (error) $= (o_i)_d$

$$= 0.05 \times (0.15 \times 0.20 \times 0.40) + (0.10 \times 0.20 \times 0.50) =$$

Ex:



$$b_1 = 0.35, \quad b_2 = 0.60 \quad (\text{Bias factors})$$

1) Calculate $h_1 (\text{in}) \text{out} + (\text{e} \cdot \text{c} \times \text{EP} \cdot o)$

$$h_1 (\text{in}) = w_1 x_1 + w_2 x_2 + b_1 = 0.15 \times 0.005 + 0.20 \times 0.10 + 0.35 =$$

$$= (0.15 \times 0.005) + (0.20 \times 0.10) + 0.35 = 0.377$$

$$h_1 (\text{out}) = \frac{1}{1 + e^{-h_1 (\text{in})}} = \frac{1}{1 + e^{-0.377}} = 0.5932$$

$(\text{Intrinsic} = 0.5932)$

GUDI VARAPRASAD

8) Calculate $h_2(\text{in})$ & $h_2(\text{out})$

$$h_2(\text{in}) = x_1w_3 + x_2w_4 + b_1$$

$$= (0.05 \times 0.25) + (0.1 \times 0.3) + 0.35 = 0.7025$$

$$h_2(\text{in}) = 0.7025$$

$$h_2(\text{out}) = \frac{1}{1 + e^{-h_2(\text{in})}} = \frac{1}{1 + e^{-0.7025}} = 0.5968$$

3) Calculate $O_1(\text{in})$ & $O_1(\text{out})$

$$O_1(\text{in}) = h_1(\text{out}) * w_5 + h_2(\text{out}) * w_6 + b_2$$

$$= (0.593 \times 0.4) + (0.5968 \times 0.45) + 0.6$$

$$O_1(\text{in}) = 1.105$$

$$O_1(\text{out}) = \frac{1}{1 + e^{-O_1(\text{in})}} = 0.7513$$

4) Calculate $O_2(\text{in})$ & $O_2(\text{out})$

$$O_2(\text{in}) = h_1(\text{out}) \times w_7 + h_2(\text{out}) \times w_8 + b_2$$

$$= (0.593 \times 0.5) + (0.5968 \times 0.45) + 0.6$$

$$O_2(\text{in}) = 1.22484$$

$$O_2(\text{out}) = \frac{1}{1 + e^{-O_2(\text{in})}} = \frac{1}{1 + e^{-1.22484}} = 0.7729$$

5) Calculate E_{Total}

$$E_{\text{Total}} = \sum \frac{1}{2} (\text{target} - \text{output})^2$$

$$E_{\text{total}} = \epsilon_{01} + \epsilon_{02}$$

$$= \frac{1}{2} (0.01 - 0.7513)^2 + \frac{1}{2} (0.99 - 0.7729)^2$$

$$\epsilon_{01} = 0.274 + 0.0235 = 0.29837$$

Part 2 : Calculating Backward Propagation Error

(Output layer \rightarrow Hidden Layer)

weights are w_5, w_6, w_1, w_8

$$w_5^* = w_5 - \eta \frac{dE_{\text{total}}}{dw_5}$$

$$\frac{dE_{\text{total}}}{dw_5} = \frac{dE_{\text{total}}}{dout_0} \times \frac{dout_0}{dnet_0} \times \frac{dnet_0}{dw_5}$$

$$\frac{dE_{\text{total}}}{dout_0} = \text{Out}_0 - \text{target}_0$$

$$= 0.751365 - 0.01 = 0.7413565$$

$$\frac{dout_0}{dnet_0} = \text{out}_0 (1 - \text{out}_0)$$

$$= 0.751365 (1 - 0.751365) = 0.186815602$$

$$\frac{dnet_0}{dw_5} = \text{out}_{h_1} = 0.59326992 = h_1(\text{out})$$

$$\frac{dE_{\text{total}}}{dw_5} = 0.7413565 \times 0.186815602 \times 0.59326992$$

$$= 0.08216704$$

$$w_5^* = w_5 - \eta \frac{dE_{\text{total}}}{dw_5} = 0.4 - 0.6 \times 0.08216704$$

$$= 0.350699776$$

GUDI VARAPRASAD

$$w_8^* = w_8 - \eta \frac{d E_{\text{total}}}{d w_8} = 0.5 + 0.03 = 0.53$$

$$(0.5 + 0.03) \frac{1}{2} + (0.5 - 0.03) \frac{1}{2} =$$

$$\frac{d E_{\text{total}}}{d w_8} = \frac{d E_{\text{total}}}{d \text{out}_{O_2}} \times \frac{d \text{out}_{O_2}}{d \text{net}_{O_2}} \times \frac{d \text{net}_{O_2}}{d w_8}$$

$$\frac{d E_{\text{total}}}{d \text{out}_{O_2}} = \frac{0.03 \text{out}_{O_2} - 0.99 \text{target}_{O_2}}{\text{out}_{O_2} (\text{target}_{O_2} - \text{out}_{O_2})} = \frac{0.03 \cdot 0.7129 - 0.99}{0.7129 (1 - 0.7129)} = -0.2171$$

$$\frac{d \text{out}_{O_2}}{d \text{net}_{O_2}} = \frac{\text{out}_{O_2} (1 - \text{out}_{O_2})}{\text{net}_{O_2} (1 - \text{net}_{O_2})} = \frac{0.7129 (1 - 0.7129)}{0.17552559} = \frac{0.17552559}{0.17552559}$$

$$\frac{d \text{net}_{O_2}}{d w_8} = \frac{\text{out}_{h_2}}{\text{out}_{h_2}} \times \frac{\text{out}_{h_2}}{\text{out}_{h_2}} \times \frac{\text{out}_{h_2}}{\text{out}_{h_2}} = \frac{0.5968}{0.5968} = \frac{0.5968}{0.5968}$$

$$\frac{d E_{\text{total}}}{d w_8} = -0.2171 \times 0.17552559 \times 0.5968 = -0.022742$$

$$w_8^* = 0.55 - 0.6 \times (-0.022742) = 0.563645$$

Similarly calculating w_6^* and w_7^*

$$w_6^* = \frac{w_6}{(1 - w_6)} \times \frac{w_7^*}{(1 - w_7)} = \frac{0.5}{(1 - 0.5)} \times \frac{0.563645}{(1 - 0.563645)} = \frac{0.5}{0.5} \times \frac{0.563645}{0.4363545} = \frac{0.563645}{0.4363545}$$

Part 3: Calculating Backward Propagation error

Suppose $\{w_1, w_2, w_3, w_4\}$ (hidden layer \rightarrow input layer)

$$w_1^* = w_1 - \eta \frac{d E_{\text{total}}}{d w_1} = 0.5 + 0.03 = 0.53$$

$$\frac{d E_{\text{total}}}{d w_1} = \frac{d E_{\text{total}}}{d \text{out}_{h_1}} \times \frac{d \text{out}_{h_1}}{d \text{net}_{h_1}} \times \frac{d \text{net}_{h_1}}{d w_1} = \frac{0.03 \text{out}_{h_1} - 0.99 \text{target}_{h_1}}{\text{out}_{h_1} (\text{target}_{h_1} - \text{out}_{h_1})} \times \frac{\text{out}_{h_1} (1 - \text{out}_{h_1})}{\text{net}_{h_1} (1 - \text{net}_{h_1})} \times \frac{\text{net}_{h_1}}{w_1} = \frac{0.03 \cdot 0.5 - 0.99}{0.5 (1 - 0.5)} \times \frac{0.5 (1 - 0.5)}{0.17552559} \times \frac{0.17552559}{w_1} = \frac{0.03 \cdot 0.5 - 0.99}{0.5} \times \frac{0.5}{0.17552559} \times \frac{0.17552559}{w_1} = \frac{0.03 \cdot 0.5 - 0.99}{0.5} \times \frac{0.5}{w_1} = \frac{0.03 \cdot 0.5 - 0.99}{0.5} \times \frac{0.5}{0.5} = \frac{0.03 \cdot 0.5 - 0.99}{0.5} = -1.92$$

GUDI VARAPRASAD

$$\frac{\partial \Sigma_{\text{total}}}{\partial \text{outh}_1} = (\text{out}_1 \frac{\partial \Sigma_{\text{out}_1}}{\partial \text{outh}_1}) + (\text{out}_2 \frac{\partial \Sigma_{\text{out}_2}}{\partial \text{outh}_1})$$

\downarrow \downarrow

$\frac{\partial \Sigma_{\text{out}_1}}{\partial \text{outh}_1}$ $\frac{\partial \Sigma_{\text{out}_2}}{\partial \text{outh}_1}$

$$\frac{\partial \Sigma_{\text{out}_1}}{\partial \text{net}_{\text{out}_1}} \times \frac{\partial \text{net}_{\text{out}_1}}{\partial \text{outh}_1}$$

\downarrow

$$\frac{\partial \Sigma_{\text{out}_1}}{\partial \text{out}_{\text{out}_1}} \times \frac{\partial \text{out}_{\text{out}_1}}{\partial \text{net}_{\text{out}_1}}$$

$$(\text{out}_2 \frac{\partial \Sigma_{\text{out}_2}}{\partial \text{net}_{\text{out}_2}}) \times \frac{\partial \text{net}_{\text{out}_2}}{\partial \text{outh}_1}$$

\downarrow

$$\frac{\partial \Sigma_{\text{out}_2}}{\partial \text{out}_{\text{out}_2}} \times \frac{\partial \text{out}_{\text{out}_2}}{\partial \text{net}_{\text{out}_2}}$$

$$\frac{\partial \Sigma_{\text{out}_2}}{\partial \text{out}_{\text{out}_2}} = \text{out}_{\text{out}_2} - \text{target}_{\text{out}_2}$$

$$= 0.7129284165 - 0.99$$

$$= -0.217071535$$

$$\frac{\partial \text{out}_{\text{out}_2}}{\partial \text{net}_{\text{out}_2}} = \text{out}_{\text{out}_2} (1 - \text{out}_{\text{out}_2})$$

$$= 0.7129 (1 - 0.7129) = 0.175510652$$

$$\frac{\partial \text{net}_{\text{out}_1}}{\partial \text{out}_{\text{out}_1}} = \text{from } h_1 \text{ from } h_1 = w_5 = 0.4$$

$$\frac{\partial \Sigma_{\text{out}_2}}{\partial \text{net}_{\text{out}_2}} = -0.21707 \times 0.1755 = -0.038$$

$$\frac{\partial \text{net}_{\text{out}_2}}{\partial \text{out}_{\text{out}_2}} = \text{from } h_1 \text{ from } h_1 = w_1 = 0.5$$

$$\frac{\partial \Sigma_{\text{out}_2}}{\partial \text{outh}_1} = -0.038 \times 0.5 = -0.01904$$

$$\frac{\partial \Sigma_{\text{out}_1}}{\partial \text{net}_{\text{out}_1}} = 0.13849856 \text{ from } \frac{\partial \Sigma_{\text{out}_1}}{\partial \text{outh}_1} = 0.1384 \times 0.4$$

$$= 0.0554$$

GUDI VARAPRASAD

$$\frac{\partial E_{\text{total}}}{\partial w_1} = 0.055396 + (-0.01904) \quad \frac{\partial E_{\text{total}}}{\partial w_1}$$

$$= 0.036350306 \quad \frac{\partial E_{\text{total}}}{\partial w_1}$$

$$\frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_{\text{out}_1}}{\partial w_1} \times \frac{\partial E_{\text{out}_1}}{\partial w_1} = \frac{\partial E_{\text{out}_1}}{\partial w_1} \times \frac{\partial E_{\text{out}_1}}{\partial w_1}$$

$$\frac{\partial E_{\text{out}_1}}{\partial w_1} = \frac{\partial E_{\text{out}_1}}{\partial w_1} \times \frac{\partial E_{\text{out}_1}}{\partial w_1} = \frac{\partial E_{\text{out}_1}}{\partial w_1} \times \frac{\partial E_{\text{out}_1}}{\partial w_1}$$

$$\frac{\partial E_{\text{out}_1}}{\partial w_1} = \frac{\partial}{\partial w_1} (w_1 x_1 + w_2 x_2 + b_1) = \frac{\partial}{\partial w_1} (0.15 \times 1.0 + 0.15 \times 0.5 + 0.15) = 0.05$$

$$\frac{\partial E_{\text{total}}}{\partial w_1} = 0.036350306 \times 0.241300709 \times 0.05 = 0.00438568$$

$$w_1^* = w_1 - n \frac{\partial E_{\text{total}}}{\partial w_1} = 0.15 - (0.6) \times 0.00438568 = 0.1497368592$$

Similarly,

$$w_2^* = \frac{\partial E_{\text{total}}}{\partial w_2} = \frac{\partial E_{\text{out}_1}}{\partial w_2} = \frac{\partial E_{\text{out}_1}}{\partial w_2}$$

$$w_3^* = \frac{\partial E_{\text{total}}}{\partial w_3} = \frac{\partial E_{\text{out}_1}}{\partial w_3} = \frac{\partial E_{\text{out}_1}}{\partial w_3}$$

Repeat this all 3 parts until you get minimum error.

- Remarks on Backpropagation Algorithm :

1. Convergence & local minima

2. Representational power of feed forward networks

3. Hypothesis space search and inductive Bias

4. Hidden layer Representation

5. Generalization, overfitting & stopping criteria.

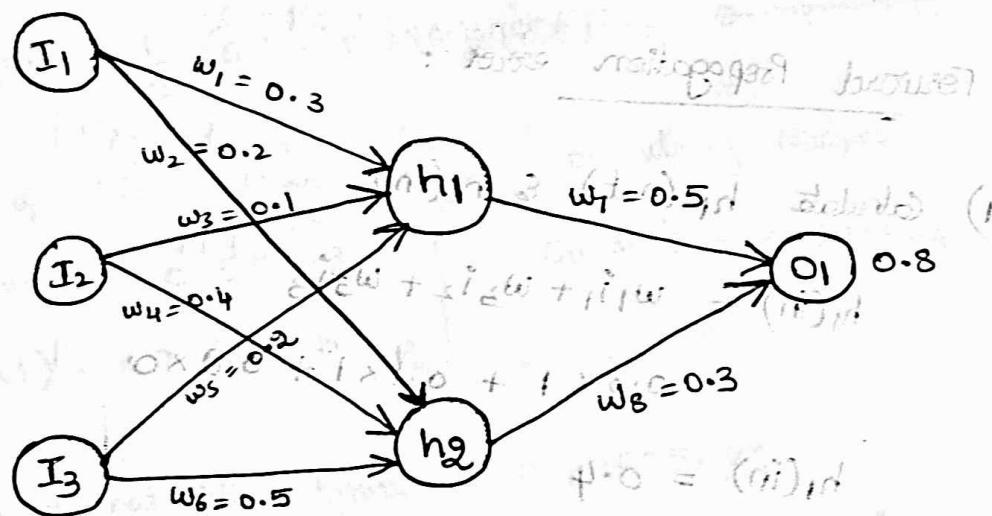
functions: Boolean, Continuous, Arbitrary

Ex: The following is a network where the output of each neuron in hidden and output layer is computed by applying the hidden layer and output layer neurons for the given:

1. Compute the output of hidden layer and output layer neurons for the given input $(1, 1, 0)$

2. What is the error of this network with the given weights.

3. Update the weights of w_7, w_8 by applying the backpropagation algorithm.



Sol: It computes the gradient of the loss function

for a single weight by the chain rule.

1. Inputs x , arrive through the preconnected graph-path.

2. Input is modeled using real weights w . The weights

are usually randomly selected to train a neural network.

3. Calculate the output for every neuron from the input layer, to the hidden layers, to the output layer.
4. Calculate the error in the outputs.
5. Travel back from output layer to the hidden layer to adjust the weights such that error is minimized.
6. Keep repeating this process until the desired output is achieved.

$$w_1 = 0.3$$

$$w_2 = 0.2$$

$$w_3 = 0.1$$

$$w_4 = 0.4$$

$$w_5 = 0.2$$

$$w_6 = 0.5$$

$$w_7 = 0.5$$

$$w_8 = 0.3$$

$$o_1 = 0.8$$

Forward Propagation error :

- 1) Calculate $h_1(\text{out})$ & $h_1(\text{in})$

$$h_1(\text{in}) = w_1 i_1 + w_3 i_2 + w_5 i_3$$

$$= 0.3 \times 1 + 0.1 \times 1 + 0.2 \times 0$$

$$h_1(\text{in}) = 0.4$$

$$h_1(\text{out}) = \frac{1}{1 + e^{-h_1(\text{in})}} = \frac{1}{1 + e^{-0.4}} = 0.5986876601$$

GUDI VARAPRASAD

$$2) h_2(\text{in}) = w_2 i_1 + w_4 i_2 + \underline{w_6 i_3} \quad \text{E.O.} - \text{pw} = ?_{\text{pw}}$$

$$= 0.2 \times 1 + 0.4 \times 1 + 0.5 \times 0$$

$$h_2(\text{in}) = \frac{0.6}{\text{total}} \times \frac{\text{total}}{\text{total}} = \frac{0.6}{\text{total}}$$

$$h_2(\text{out}) = \frac{1}{1 + e^{-h_2(\text{in})}} = \frac{1}{1 + e^{-0.6}} = 0.6456563062$$

3) Calculate $O_1(\text{in}) \& O_1(\text{out})$

$$O_1(\text{in}) = h_1(\text{out}) * w_7 + h_2(\text{out}) * w_8$$

$$= 0.5986876601 \times 0.5 + 0.6456563062 \times 0.3$$

$$= 0.2993438301 + 0.1936968919$$

$$O_1(\text{in}) = 0.493040782$$

$$O_1(\text{out}) = \frac{1}{1 + e^{-0.493040782}} = 0.6208224$$

4) Calculate E_{total}

$$E_{\text{total}} = \sum \frac{1}{2} (\text{target} - \text{output})^2 = E_{O1}$$

$$E_{O1} = \frac{1}{2} (0.8 - 0.6208224)^2 = 0.01605229064$$

① output of hidden layer :

$$h_1(\text{out}) = 0.598688 \quad | \quad O_1(\text{out}) = 0.620822$$

$$h_2(\text{out}) = 0.645656$$

② error of network, $E_{\text{total}} = 0.01605229$

③ updating weights w_7, w_8 using Backpropagation:

$$w_7^* = w_7 - \eta \frac{dE_{\text{total}}}{dw_7} \quad \begin{matrix} 0.01605229 \\ \text{total} \\ \text{total} \end{matrix} \quad h = \text{learning rate} = 0.5$$

$$w_8^* = w_8 - \eta \frac{dE_{\text{total}}}{dw_8} \quad \begin{matrix} 0.01605229 \\ \text{total} \\ \text{total} \end{matrix} \quad (\text{assume})$$

GUDI VARAPRASAD

$$w_7^* = w_7 - 0.5 \frac{\partial E_{\text{total}}}{\partial w_7} \quad \text{target}_7 + \text{err}_7 = (\text{err})_{\text{ad}}$$

$$0.235402044316 \times 1.0 + 1 \times 6.0$$

$$\frac{\partial E_{\text{total}}}{\partial w_7} = \frac{\partial E_{\text{total}}}{\partial \text{out}_{01}} \times \frac{\partial \text{out}_{01}}{\partial \text{net}_{01}} \times \frac{\partial \text{net}_{01}}{\partial w_7} \quad (\text{err})_{\text{ad}}$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out}_{01}} = \text{out}_{01} - \text{target}_{01}$$

$$\frac{\partial \text{out}_{01}}{\partial \text{net}_{01}} = 0.620822 - 0.8 = -0.179178$$

$$\frac{\partial \text{net}_{01}}{\partial w_7} = (\text{err})_{\text{ad}} + \text{err}_7 = (\text{err})_{\text{ad}}$$

$$\frac{\partial \text{out}_{01}}{\partial \text{net}_{01}} = \text{out}_{01} (1 - \text{out}_{01}) = 0.235402044316$$

$$= \frac{0.620822}{(-0.179178)} \frac{1 - 0.620822}{(1 - 0.179178)} = -0.2112827557$$

$$\frac{\partial \text{net}_{01}}{\partial w_7} = \frac{\text{out}_{01}}{1} = 0.598688 \quad (\text{err})_{\text{ad}}$$

$$\frac{\partial E_{\text{total}}}{\partial w_7} = -0.179178 \times -0.2112827557 \times 0.598688$$

$$= 0.235402044316 \times 0.02525198182$$

$$w_7^* = w_7 - 0.5 (0.02525198182)$$

$$= 0.5 + 0.02525198182 = 0.5126259909 \quad (\text{err})_{\text{ad}}$$

$$w_7^* = 0.5126259909 \quad (\text{err})_{\text{ad}}$$

$$w_8^* = w_8 - 0.5 \frac{\partial E_{\text{total}}}{\partial w_8} \quad \text{target}_8 + \text{err}_8$$

$$\frac{\partial E_{\text{total}}}{\partial w_8} = \frac{\partial E_{\text{total}}}{\partial \text{out}_{02}} \times \frac{\partial \text{out}_{02}}{\partial \text{net}_{02}} \times \frac{\partial \text{net}_{02}}{\partial w_8}$$

$$\frac{\partial \text{out}_{02}}{\partial \text{net}_{02}} = \frac{\text{out}_{02}}{1 - \text{out}_{02}} = \frac{0.235402044316}{1 - 0.235402044316} = 0.235402044316$$

GUDI VARAPRASAD

$$\frac{dE_{\text{total}}}{dw_1} = -0.179178$$

$$\frac{dout_1}{dw_1} = 0.2354020443$$

$$dout_1 =$$

$$\frac{dout_1}{dw_2} = out_2 = 0.645656$$

$$dw_2$$

$$\begin{aligned} \frac{dE_{\text{total}}}{dw_2} &= -0.179178 \times 0.235402044316 \times 0.645656 \\ &= -0.0272330574 \end{aligned}$$

$$w_g^* = w_g - (0.5)(-0.0272330574)$$

$$w_g^* = 0.3 + 0.0136165287 = 0.3136165287$$

updated

weights are:

$$w_7^* = 0.5126259909 \approx 0.5126$$

$$w_8^* = 0.3136165287 \approx 0.3136$$

GUDI VARAPRASAD

*. FACE RECOGNITION

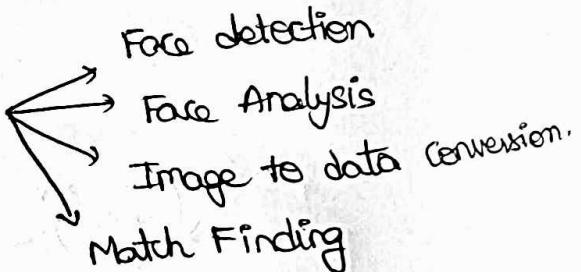
- Face Recognition is a category of biometric software that maps individuals facial features and stores the data as a face point.
- Machine learning and Face Recognition.

*. IMAGE PROCESSING

- Involves the process of Computer vision.
- Image Reading : Reads any image between the range of 0 - 255.
- For any image, there are 3 colors ^{RGB}.
- OpenCV : Python Library to solve computer vision.

*. Machine Learning :

- Takes dataset as input and learn from that data
- identifies the pattern (height, width, color, etc...)
- ML does 3 major functions in Face Recognition.
 - 1) Deriving the Feature Vector
 - 2) Matching Algorithms
 - 3) Face Recognition operations



- Face Recognition softwares : DeepVision AI, SenseTime, Amazon Rekognition, faceFirst, TrueFace, Cognitec...

• Applications :

- 1) Hospital Security
- 2) Mobile Devices

2 - 3 JULY 2014

3) Airliner Industry

4) Banking sectors

* Advancements (in Neural Networks) : (3 advanced Topics)

① Alternative error functions :

- Adding an extra penalty term to $E(w)$ can reduce the problem of overfitting.
- each weight is multiplied with a constant in each iteration. so, we can redefine E as

$$E(\bar{w}) = \frac{1}{2} \sum_{d=1}^D \sum_{k=1}^{K_d} (t_{kd} - o_{kd})^2 + \gamma \sum_{j=1}^J w_{ji}^2$$

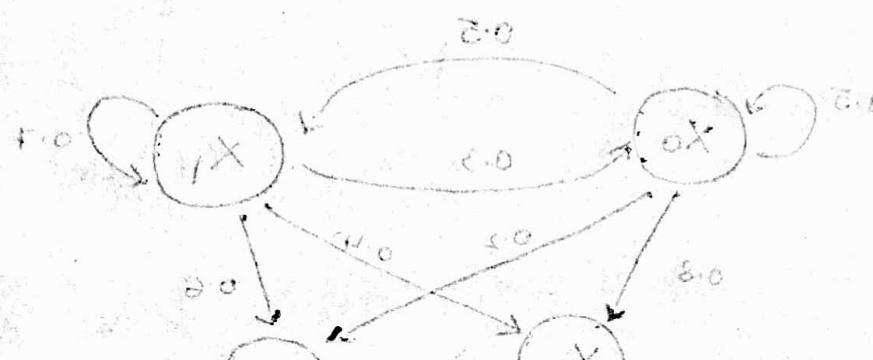
② Alternative error Minimization procedures :

A) Line search

B) Conjugate Gradient Descent method.

③ Recurrent Neural Networks (RNN) :

- ANN that can be applied to time series data and uses the output of network units at time t as input to other units at time $t+1$.



MODULE - 6

*. MARKOV MODEL

$$\bullet \quad P(v_1, v_2, \dots, v_n) = \prod_{i=1}^n P(v_i | v_{i-1})$$

Ex 3: What is the probability of given sequence S, R, R, R, C, C
 Initial probabilities given as $\pi = \{0.7, 0.25, 0.05\}$

Sol: $P(S) = 0.7$

$$P(R|S) = 0.05$$

$$P(R|R) = 0.6$$

$$P(C|R) = 0.2$$

$$P(C|C) = 0.5$$

$$P(\text{Sequence}) = P(S, R, R, R, C, C)$$

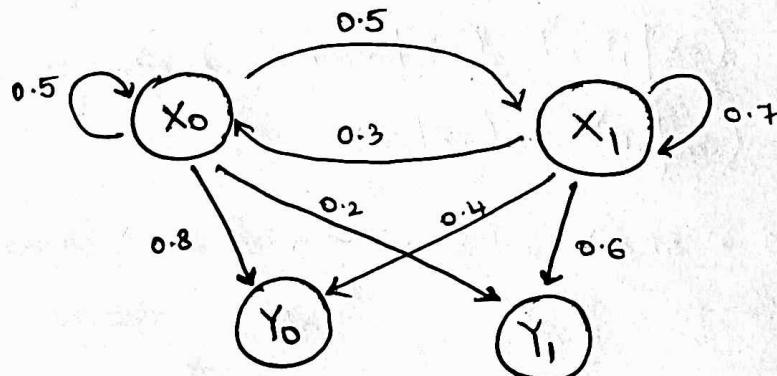
$$= P(S) \times P(R|S) \times P(R|R) \times P(R|R) \times P(C|R) \times P(C|C)$$

$$= 0.7 \times 0.05 \times 0.6 \times 0.6 \times 0.2 \times 0.5$$

$$P(S, R, R, R, C, C) = 0.00126$$

*. FORWARD ALGORITHM :

Ex:



What is the probability of $\{X_0, Y_0, Y_1, Y_1\}$ sequence given the parameters of our HMM?

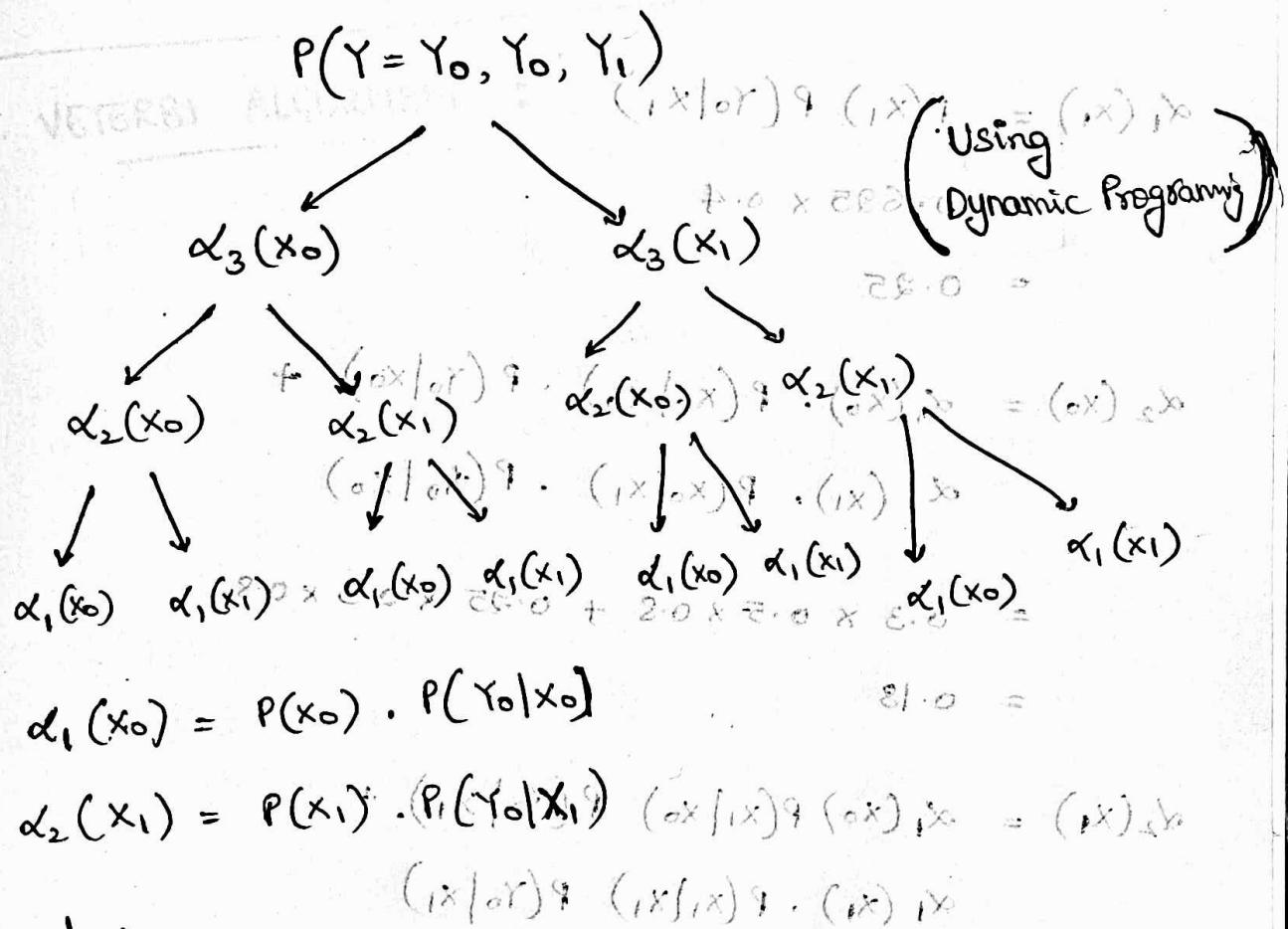
GUDI VARAPRASAD

$$A = \begin{bmatrix} x_0 & x_1 \\ x_0 & 0.5 & 0.5 \\ x_1 & 0.3 & 0.7 \end{bmatrix} \quad (ix|or) \quad B = \begin{bmatrix} y_0 & y_1 \\ y_0 & 0.8 & 0.2 \\ y_1 & 0.4 & 0.6 \end{bmatrix} \quad (ix) \quad (ix) \quad (ix) \quad (ix) \quad (ix)$$

Transition Matrix Emission Matrix

stationary distribution for hidden states : $\pi \cdot A = \pi \Rightarrow \pi = [0.375 \ 0.625]$

Here no. of Hidden states = 2



Formula :

$$\alpha_t(x_i) = \sum_{j=0}^{n-1} \alpha_{t-1}(x_j) \cdot P(x_i|x_j) \cdot P(Y^t|x_i)$$

For $j=0, n=1$

$$\alpha_t(x_i) = \alpha_{t-1}(x_0) \cdot P(x_i|x_0) \cdot P(Y^t|x_i)$$

$$\alpha_t(x_i) = \alpha_{t-1}(x_1) \cdot P(x_i|x_1) \cdot P(Y^t|x_i)$$

Here, n is no. of Hidden states $= 2$

GUDI VARAPRASAD

$$\alpha_i(x_i) = \prod_{j=1}^n P[j] \cdot P(Y_j | x_i)$$

$$P(Y_1 Y_2 \dots Y_t) = \sum_{i=0}^{n-1} \alpha_{t-1}(x_i)$$

initial state
transition

$$\alpha_0(x_0) = P(x_0) \cdot P(Y_0 | x_0)$$

$$= 0.375 \times 0.8$$

$$= 0.3$$

initial state
transition

$$\alpha_1(x_1) = P(x_1) \cdot P(Y_1 | x_1)$$

$$= 0.625 \times 0.4$$

$$= 0.25$$

(ix) \rightarrow (ok) \rightarrow (ox) \rightarrow

$$\alpha_2(x_0) = \alpha_1(x_0) \cdot P(x_0 | x_0) \cdot P(Y_0 | x_0) +$$

$$\alpha_1(x_1) \cdot P(x_0 | x_1) \cdot P(Y_0 | x_0)$$

$$= 0.3 \times 0.5 \times 0.8 + 0.25 \times 0.3 \times 0.8$$

$$= 0.18$$

(ox) \rightarrow (ox) \rightarrow (ox) \rightarrow (ox) \rightarrow (ox) \rightarrow

$$\alpha_2(x_1) = \alpha_1(x_0) P(x_1 | x_0) (P(Y_0 | x_1) + P(Y_1 | x_1))$$

$$\alpha_1(x_1) \cdot P(x_1 | x_1) P(Y_0 | x_1)$$

$$(ix) \rightarrow (0.3 \times 0.5 \times 0.4) + (0.25 \times 0.7 \times 0.4)$$

$$= (ix) \rightarrow$$

(ix) \rightarrow (ok) \rightarrow (ox) \rightarrow (ox) \rightarrow (ox) \rightarrow

$$\alpha_3(x_0) = \alpha_2(x_0) P(x_0 | x_0) \cdot P(Y_1 | x_0) +$$

$$\alpha_2(x_1) \cdot P(x_0 | x_1) \cdot P(Y_1 | x_0)$$

$$= (0.18 \times 0.5 \times 0.2) + (0.13 \times 0.3 \times 0.2)$$

$$= 0.0258$$

(ix) \rightarrow (ok) \rightarrow (ox) \rightarrow (ox) \rightarrow (ox) \rightarrow

$$\begin{aligned}\alpha_3(x_1) &= \alpha_2(x_0) \cdot P(x_1|x_0) \cdot P(y_1|x_1) + \\ &\quad \alpha_2(x_1) \cdot P(x_1|x_1) \cdot P(y_1|x_1) \\ &= 0.18 \times 0.5 \times 0.6 + 0.13 \times 0.7 \times 0.6 \\ &= 0.1086\end{aligned}$$

$$\begin{aligned}P(Y=y_0, y_1) &= \alpha_3(x_0) + \alpha_3(x_1) \\ &= 0.0258 + 0.1086 = 0.1344\end{aligned}$$

* VITERBI ALGORITHM :