

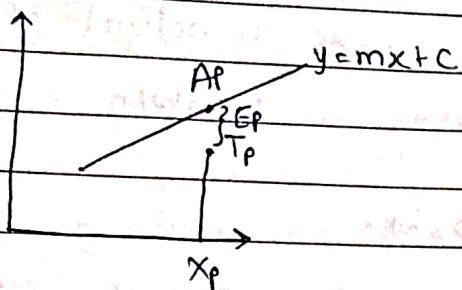
* Supervised learning algorithm

In this algorithm, input & output is known.
Data is labelled (with teacher)

* Unsupervised learning algorithm

In this algorithm input & output are not known.
Data is unlabelled (without teacher)

Eg:



T_p = Target point

A_p = Actual point

E_p = Error at point P.

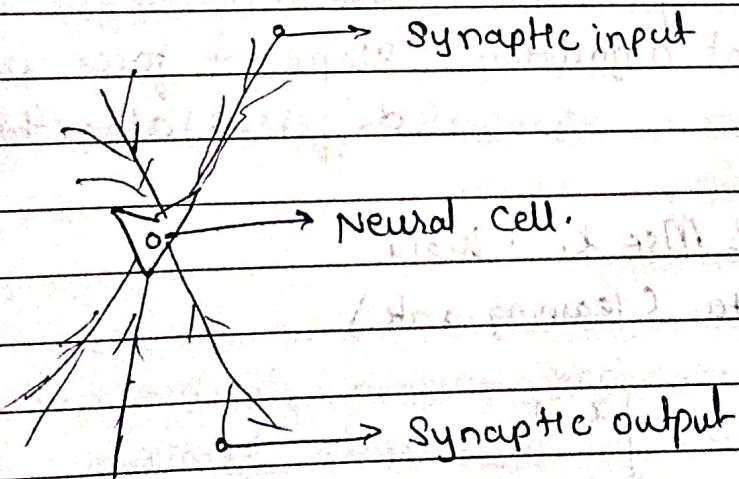
$$E_p = T_p - A_p$$

$$\text{MSE} = \frac{\sum_{i=0}^p (T_p - A_p)^2}{p}$$

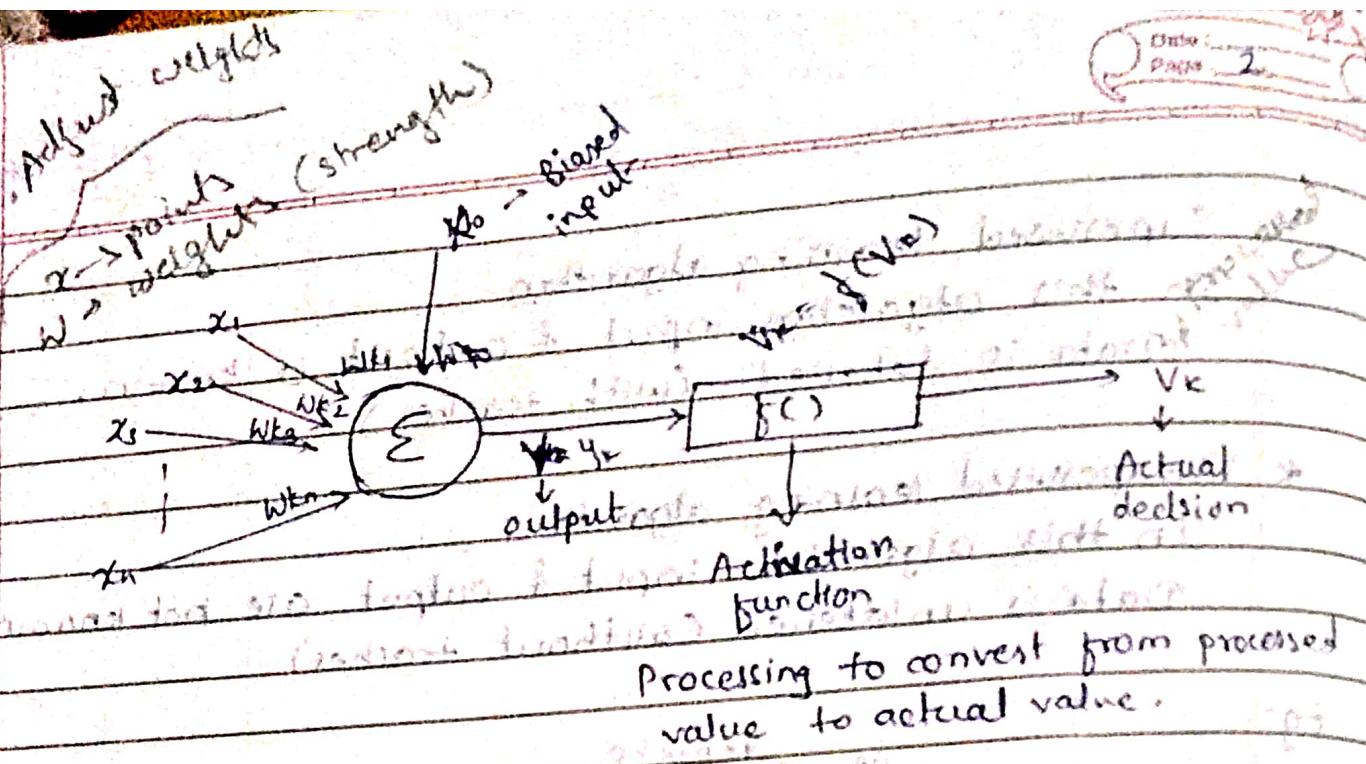
Mean square error

To minimize MSE we have to change m & c - intercept
Here m & c are free parameters.

*



In our body,
we adjust
strength of
connection by
free parameter



$$y_r = x_1 w_{ki} + x_2 w_{k2} + \dots + x_n w_{kn}$$

$$y_r = \sum_{i=1}^n x_i w_{ki} + x_0 w_{k0}$$

$$y_r = \sum_{i=0}^n x_i w_{ki}$$

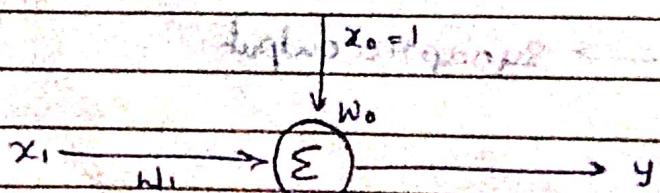
$$y_r = \begin{cases} 1 & y_r > 0 \\ 0 & \text{else} \end{cases}$$

- Gradient descent algorithm: Slope & force are in same direction.

- Steepest descent algorithm: Slope & force are in different directions.

$$W_{\text{new}} = \eta Mse x_i + W_{\text{old}}$$

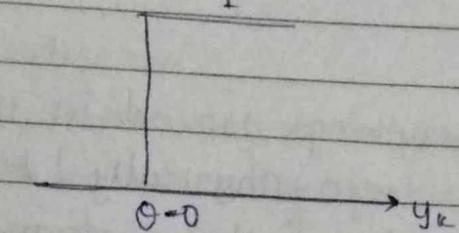
\downarrow
netta (learning rate)



$$y = x_1 w_1 + x_0 w_0$$

$$y = x_1 w_1 + w_0 \quad \{ x_0 = 1 \}$$

- * Types of activa" func"
- i) Threshold activa" func"



Types of learning

- i) Supervised learning: (inductive): Training data includes desired output.

Regression Classification,

linear non linear
The algorithm learns relationship bet" specific inputs & outputs based on training data & human feedback

- 2] Unsupervised learning: Training data does not include desired clustering association output.

The algo. analyzes data for trends & patterns without being given a specific output variable / human feedback.

- 3] Semi supervised: Training data includes few outputs.

- 4] Reinforcement:- Rein Rewards from sequence of action (eg. game)

The problem of problem recognition is to identify the undeline structure within data.

Intelligence is defined as ability of to comprehend to understand & profit from experience

- Plato said "Concept of abstract ideas are known to us a prior knowledge, through a mystic connection with world."
- Pattern - Pattern is everything around in this digital world.
A pattern can either be seen physically / it can be observed mathematically by applying algorithm.
Eg. Colors on clothes, speech pattern, etc
In CS, pattern is represented using vector feature values.
- Pattern recognition -
PR is process of recognizing patterns by using machine learning algorithm.
PR can be defined by classification of data based on knowledge already gained or on statistical information extracted from patterns and / or their representation.
- Feature
Feature may be represented as continuous, discrete or discrete binary variables.
- Numerical value given info. about object signal etc is called feature.
- Broadly speaking feature are any extractable measure used.
- Features may be :
 - 1) Symbolic (color)
 - 2) Numeric (Weight)
 - 3) Both.

- feature when it is put into certain order is called feature recognition.

- feature extraction.

feature extrac" refers to process of transforming raw data into numerical feature that can be processed while preserving the info. in the original data set.

- It yields better results than applying machine learning directly into data

- Feature extrac" can be accomplished manually or automatically.

- feature selection imp features rakhna bat sab nikal do
feature selec" is the process of reducing no. of input variables when developing a predictive model

It is desirable to reduce no. of input variables to both reduce the computational cost of modeling & in some cases, to improve performance of model.

feature selection is primarily focused on removing non-informative / redundant predictors from the model.

- Difference bet" feature vector & feature space.

- A feature is a numerical / symbolic property of an aspect of an object.

- A feature vector is a vector containing multiple elements about an object.

- Putting feature vector for objects together can make up a feature space.

Scatter plot

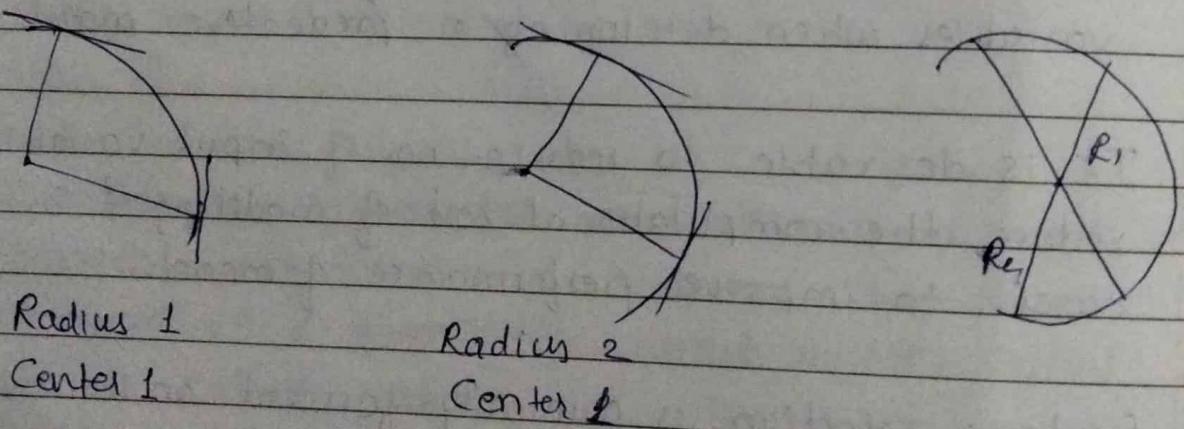
Scatter plot are graphs that represent the relationship between 2 variables in a data set.

It represents data points on 2-D plane / cartesian system.

The independent variable or attribute is plotted on x-axis while dependent variable is plotted on y-axis.

Boundary based → Chain Coding Technique.
Region based.

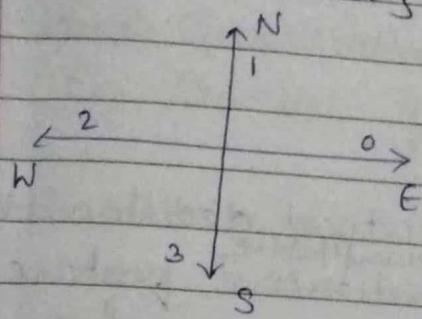
* Curve Pattern Recognition.



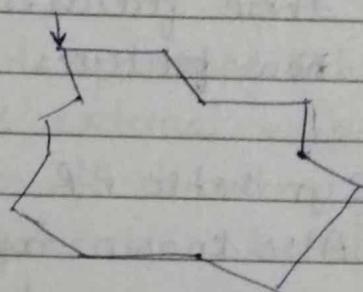
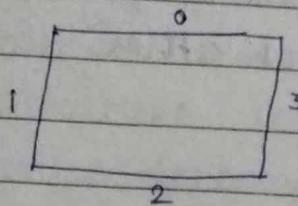
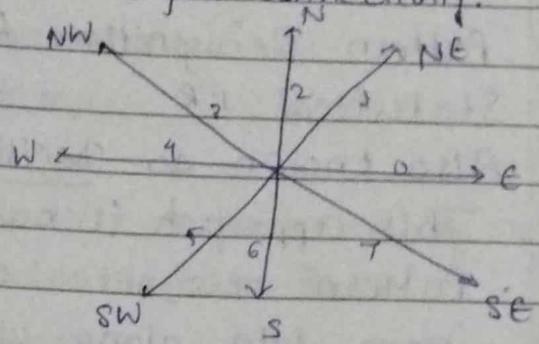
Chain Coding Technique.

- 1) Four Connectivity
- 2) Eight Connectivity

Four connectivity



Eight connectivity



(0, 3, 2, 1)

Chain code

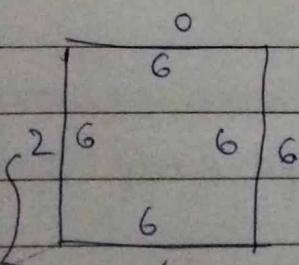
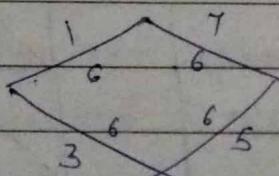
(0, 1, 0, 6, 7, 5, 3, 4, 3, 1, 2, 1)

Chain-code

Translational invariant (Problem of origin rotation)
problem

Differential Chain Coding.

To overcome translational invariant problem.

If pattern is rotated then also it gives same difference
(Advantage)Bet 2
& 4 there
are 6 lines

II Pattern Recognition Approaches:

1) Statistical PR:

Also known as Decision Theory.

- This approach is based on statistical decision theory.
- Pattern recognizer extract quantitative features from data along with multiple samples & compare those features. However it does not touch upon how features are related to each other.

2) Synthetic PR

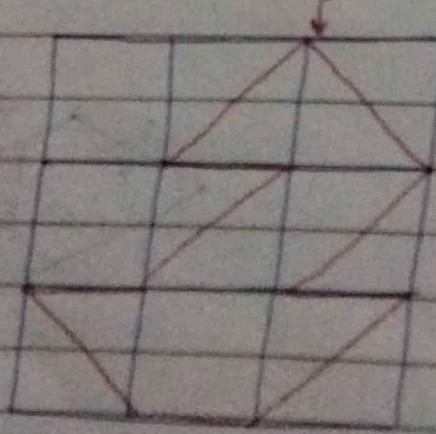
Also known as Structural

- This approach is closest to how human perceps works. It extract morphological features from one data sample & checks how these are connected & related.

3) Neural PR:

In this approach artificial neural network are utilized.

Compared to ones mentioned above, it allows more flexibility in learning & is the closest to naturally learning.



→ Chain code (8-connectivity)

7 5 0 5 4 3 0 1 4 1

Features.

- 1) Quantitative data is numerical
Eg. Amount of rainfall / temp.
- 2) Qualitative data is discrete, with a finite no. of well-known values like weather descrip? in table
- 3) Morphological features.

This includes aspects of outward appearance (shape, structure, color, pattern, size) ie external morphology (or eidonomy) as well as form & structure of internal parts like bones & organs ie internal morphology (or anatomy).

Minimum perimeter length polygonization.

- Boundary is converted into set of "concaved cells".
- Concatenated → Join | combine.
- Cord try to shrink towards outward direction | inward direction until & unless reaches the limiting point.
- Moves only in one direction
- Moves towards resulting direction upto limiting points given by inward & outward direction.

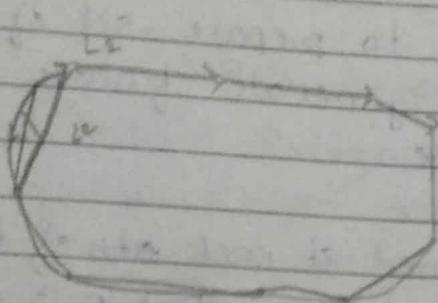
Splitting of boundary.

- Identify the set of points on arbitrary boundary along with this boundary can be split.
- The boundary points becomes vertexes of the polygon.
- Criterion funcⁿ: helps to decide part of point on arbitrary boundary along with which boundary can be splitted.
- Approximation 1: Algorithm

- s1: Initially decide 2 points along which we can divide the boundary.
- s2: Draw straight line passing through these 2 points.
- s3: Find max. distance on boundary point from the line in both segments.
- s4: If max. distance is greater than threshold, then take maximum point as a next split point.
- s5: Again find max. distance from the boundary & if distance is greater than threshold value then take as next split point else stop the algorithm.

Approximation 2: What is the area of portion enclosed within straight line & original boundary.

Approximation 3:



Select an origin point.

Start drawing lines from that point
(L₁, L₂, L₃, ...)

Suppose in the abv fig;
L₂ > 0 (threshold val.)

then L₂ will be the final boundary.

Again repeat the same process of drawing lines from end point of L₂.

$$\theta_2 = \sum_{i=1}^k x_i \theta_{i-1}$$

$$= x_n \theta_{n-1} + x_{n-1} \theta_{n-2} + \dots + x_k \theta_{k-1}$$

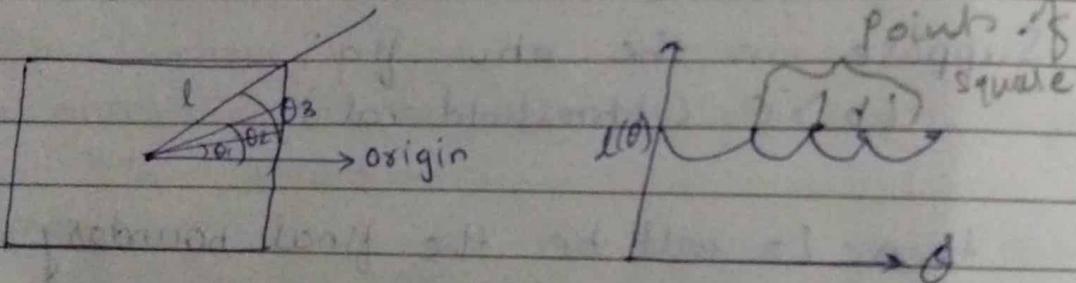
Feature vector.

- Polygon approximation of arbitrary boundary such polygon is not much helpful for PR purpose
- Therefore, approximate with piecewise linear segments.

Finally, we have to prepare set of feature vector.
One approach to generate feature vector is
On n^{th} angle.

- Inner angle sustained at each n^{th} of the vertices.
- Autoregressive model, such that n^{th} angle can be represented by linear combination of k no. of previous angles.

Signature.

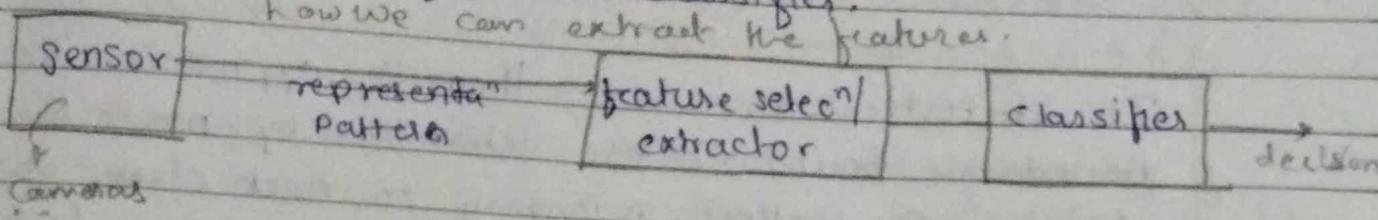


Main purpose of feature vector - multiple features are kept in series

Statistical Pattern Recognition

WORLD STAR™
Date: 3/8/22
Page: 13

- # Block diagram of pattern classifier.



- # Pattern classifier..

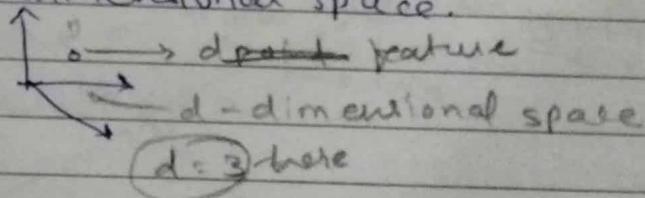
Statistical Pattern Recognition.

Eg of student in class
60 boys
20 girls
probability of new student coming will be as per prior knowledge

Statistical pattern recogni" refers to use of statistics to learn from examples.

It means to collect observa", study & digest them in order to infer general rules or concept that can be applied to new unseen observation.

- In statistical approach, each pattern is represented in terms of d features or measurements & is viewed as point in d-dimensional space.



- The goal is to choose those features that allows pattern vector belonging to different categories to occupy compact & disjoint region in d-dimensional feature space.
- The effectiveness of representa" space (feature set) is determined by how well patterns from different classes can be separated.
- Given a set of training patterns from each class, objective is to establish decision boundaries in feature space

which separates pattern belonging to same different classes.

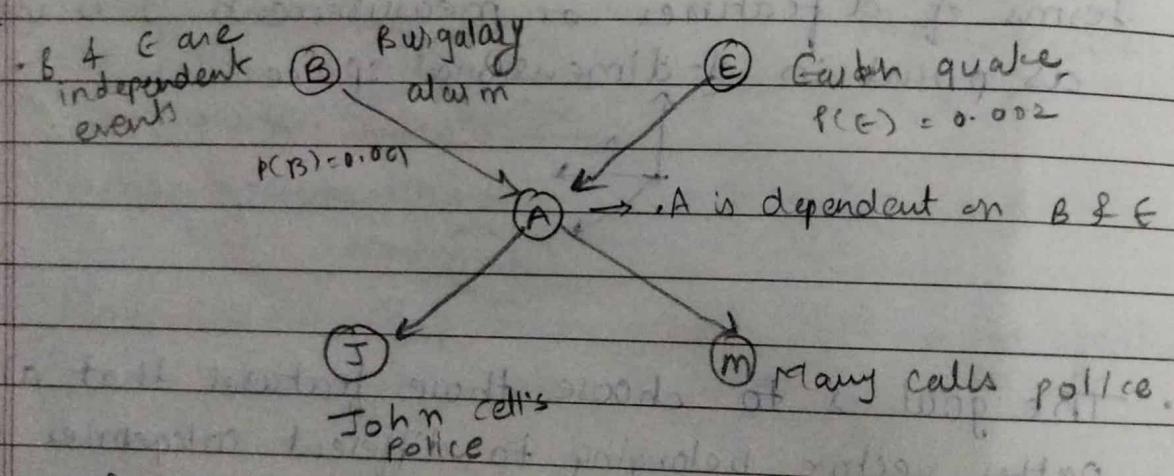
In statistical decision, theoretic approach, decision boundaries are determined by probability distribution of pattern belonging to each class, which must either be either be specified or learned [41], [44].

Gaussian case & class dependence

Bayesian Network

- Directed Acyclic Graph (DAG)
- Nodes : Set of random variables / Events
- Directed link : $X \rightarrow Y$
 X has direct influence on Y
- Conditional Probability Table (CPT)

Eg:-



For A

B	E	$P(A)$
T	T	0.95
T	F	0.95
F	T	0.29
F	F	0.001

for J

A	$P(J)$
T	0.9
F	0.05

for M

A	$P(M)$
T	0.7
F	0.9

Here probability depends on previous event hence sum of probability is not one.

q-not

T-product

Baye's rule:

Probability of occurrence of event A & event B is given by probability of A $p(A \cap B)$ [probability of A & B]

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

Already event B is occurred & the probability of occurrence A is $p(A|B)$. This is known as conditional probability.

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

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

Joint probability distribution.

$$P(x_1, x_2, x_3, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{Prob. of parent})$$

In Above ex.

$$\begin{aligned} P(A, J, M) &= P(A) * P(J|A) * P(M|A) \\ &= P(A \cap J \cap M) \end{aligned}$$

$$Q.1 \quad P(J \cap M \cap A \cap \neg B \cap \neg E) \cdot ?$$

$$= P(J|A) * P(M|A) * P(A|\neg B, \neg E) * P(\neg B) * P(\neg E)$$

$$= 0.9 * 0.7 * 0.001 * (1-0.001) * (1-0.002)$$

$$= 0.999 \quad 0.998$$

$$= 0.0006281$$

$$Q.2 \quad P(J) \cdot ?$$

$$P(J) = P(J|A) * P(A) + P(J|\bar{A}) * P(\bar{A})$$

$$P(A) = P(A|B, E) * P(B \cap E) +$$

$$P(A|\bar{B}, E) * P(\bar{B} \cap E) +$$

$$P(A|B\bar{E}) * P(B \cap \bar{E}) +$$

$$P(A|\bar{B}\bar{E}) * P(\bar{B} \cap \bar{E})$$

$$\begin{aligned}
 &= 0.9 * 0.001 * 0.002 + \\
 &\quad 0.29 * (1 - 0.001) * 0.002 + \\
 &\quad 0.95 * 0.001 * 0.998 + \\
 &\quad 0.001 * 0.999 * 0.998 \\
 &= 0.0000018 + 0.000879 + 0.0009481 \\
 &\quad 0.00093 \\
 &= 0.002518
 \end{aligned}$$

$P(\bar{A})$
 $1 - P(A)$

$$\begin{aligned}
 P(\bar{A}) &= P(\bar{A} | B, E) * P(B \wedge E) + \\
 &\quad P(\bar{A} | \bar{B}, E) * P(\bar{B} \wedge E) + \\
 &\quad P(\bar{A} | B, \bar{E}) * P(B \wedge \bar{E}) + \\
 &\quad P(\bar{A} | \bar{B}, \bar{E}) * P(\bar{B} \wedge \bar{E})
 \end{aligned}$$

$$\begin{aligned}
 &= 0.05 * 0.001 * 0.002 + \\
 &\quad 0.71 * 0.999 * 0.002 + \\
 &\quad 0.05 * 0.001 * 0.998 + \\
 &\quad 0.999 * 0.999 * 0.998
 \end{aligned}$$

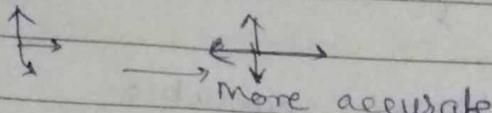
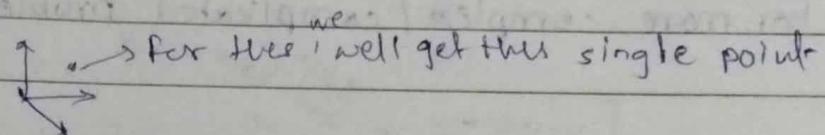
$$\begin{aligned}
 &= 0.0000001 + 0.001418 + 0.0000499 + \\
 &\quad 0.99600499 - 0.99600499
 \end{aligned}$$

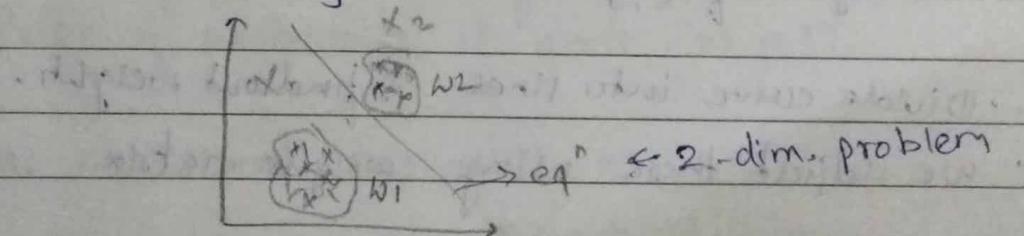
$$= 0.9974729$$

$$\begin{aligned}
 P(J) &= P(J|A) * P(A) + P(J|\bar{A}) * P(\bar{A}) \\
 &= 0.9 * 0.002518 + 0.05 * 0.9974729
 \end{aligned}$$

$$= 0.022662 + 0.0498273$$

$$= 0.07253564$$

- 1) Dimensionality of feature vector dependent on problem under consideration.
- 2) Accuracy of descrip" increased as we increase dimensionality of feature vector.

- 3) Increase in dimensionality increases complexity of problem
- 4) Dimension of PR problem based on the complexity of the problem.
 - Extract more no. of features
- 5) Once we represent pattern by pattern feature then this entire feature point is single point in dimensional feature space.

- 6) In PR problem we have to find out boundaries for patterns. Designing / training of classifier means to find out equation of decision boundary for given training samples. Therefore this is called supervised learning.



- 7) In above example, training samples are used to separate using boundaries & it is possible to separate using single line. Hence it is called as linearly separable problem.

3-dim → Plane
multidimensional → Hyperplane

2 dim - line

features \rightarrow Input of classifier
(Neural network always in bet' n 0 to 1)

WORLD STAR™

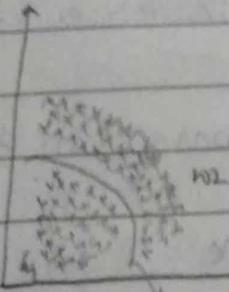
Date _____
Page 18

- e) In case of 3-D problem classes can be sp separated using plane. Hyperspace is required in multidimensional problem.

In all cases, eqⁿ of decision boundary is linear.

labelled
data is
for given
input output
there

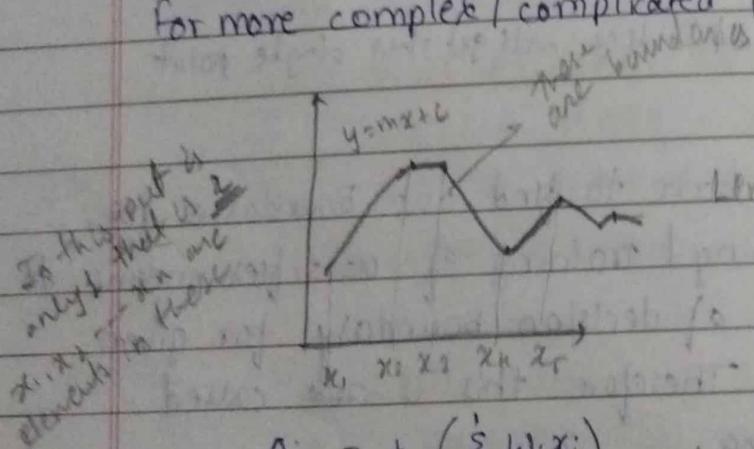
* If classes are not linearly separable



Non linear decision

boundary.

For more complex / complicated problems



$$o_j = f \left(\sum_{i=0}^n w_i x_i \right)$$

• Divide curve into lines & find out weight.

• we adjust these using weight matrix

No W

• $x_0 = 1$ is always present in front of weight

• $y_p = y_t$ after adjustment of weight

Weight is balanced such that all signs are

• weight after each iteration

• weight after each iteration

* Joint probability.

- Q. The height of students are given below. Classify new students having height 3 feet.

$P(B)$	$P(a)$
6	6
7	5
5	3
7	2
6	3
5	3
4	4
3	5
5	2
6	1

$$P(x|G) = P(a)$$

Prior of probability of boys $\Rightarrow 0.5$

Prior of probability of girls $\Rightarrow 0.5$

Class conditional probability \Rightarrow

$$P(3|B) = P(3)/10 = 1/10 = 0.1$$

$$P(3|G) = P(3)/10 = 3/10 = 0.3$$

Joint prob:-

$$P(G|3) = \frac{P(3|G) \cdot P(a)}{P(3)}$$

$$= \frac{0.3 \times 0.5}{0.2} = 0.75$$

$$P(3) = 0.1 \times 0.5 + 0.3 \times 0.5 \\ = 0.2$$

$$P(B, 3) = \frac{P(3|B) \times P(B)}{P(3)}$$

$$= \frac{0.1 \times 0.5}{0.2} = 0.25$$

$$P(4, 3) > P(B, 3)$$

\therefore It is a girl.

- Q. Consider following dataset are features & t is the class dataset = a, b, c.

a	b	c	t
1	0	1	1
1	1	1	1
0	1	1	0
1	1	0	0
1	0	1	0
0	0	0	1

Classify the test instance given below into class 1 or 0 using Bayes classifier.
 $(a, b, c) = (0, 0, 1)$ $(a, b, c) = (1, 0, 0)$

	k			value of k=1 & a=0
a -	0	1	1	
	1	2	2	value of k=1 & a=1
b -	0	1	2	
	1	2	1	
c -	0	1	1	
	1	2	2	

Since we have 3 IIP we divide each element by 3 to create likelihood table.

	k			
a -	0	1/3	1/3	1/3
	1	2/3	2/3	2/3
		3/3	3/3	3/3
b -	0	1/3	2/3	3/3
	1	2/3	1/3	1/3
		3/3	3/3	3/3
c -	0	1/3	1/3	2/3
	1	2/3	2/3	4/3
		3/3	3/3	3/3

$$(a, b, c) = (0, 0, 1)$$

$$P(1|x) = P(a_0|1) \times P(b_0|1) \times P(c_1|1) \times P(0)$$

$$= \frac{1}{3} \times \frac{2}{3} \times \frac{2}{3} \times \frac{1}{2}$$

Joint
prob.

$$= \frac{4}{27 \times 2}$$

$$= \frac{2}{27}$$

$$P(0|x) = P(a_0|0) \times P(b_0|0) \times P(c_0|0) \times P(0)$$

$$= \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{1}{2}$$

conditional

$$= \frac{1}{27}$$

$$P(x|1) = \frac{P(1|x)}{P(0|x) + P(1|x)}$$

$$= \frac{2}{27} \times \frac{27}{3}$$

$$= \frac{2}{3}$$

$$P(x|0) = \frac{P(0|x)}{P(0|x) + P(1|x)}$$

$$= \frac{1}{27} \times \frac{27}{3}$$

$$= \frac{1}{3}$$

$$P(x|1) > P(x|0)$$

∴ It belongs to class 1

$$(a_1, b_1, c_1) = (1, 0, 0)$$

$$\begin{aligned}
 P(1|x) &= P(a_1|1) \times P(b_0|1) \times P(c_0|1) \times P(1) \\
 &= \frac{2}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{1}{2} \\
 &= \frac{4}{27} \\
 &= \frac{2}{27}
 \end{aligned}$$

$$\begin{aligned}
 P(0|x) &= P(a_1|0) \times P(b_0|0) \times P(c_0|0) \times P(0) \\
 &= \frac{2}{3} \times \frac{1}{3} \times \frac{1}{3} \times \frac{1}{2} \\
 &= \frac{1}{27}
 \end{aligned}$$

$$\begin{aligned}
 P(x|1) &= \frac{P(1|x)}{P(1|x) + P(0|x)} \\
 &= \frac{\frac{2}{27}}{\frac{2}{27} + \frac{1}{27}} \\
 &= \frac{2}{3}
 \end{aligned}$$

$$\begin{aligned}
 P(x|0) &= \frac{P(0|x)}{P(0|x) + P(1|x)} \\
 &= \frac{\frac{1}{27}}{\frac{1}{27} + \frac{2}{27}} \\
 &= \frac{1}{3}
 \end{aligned}$$

$P(x|1) > P(x|0)$
 So it belongs to class 1

$$(a, b, c) = (1, 0, 1)$$

$$\begin{aligned} P(1|x) &= p(a|1) \times p(b|1) \times p(c|x_1) \\ &= \frac{2}{3} \times \frac{2}{3} \times \frac{2}{3} \times \frac{1}{2} \\ &= \frac{4}{27} \end{aligned}$$

$$\begin{aligned} P(0|x) &= p(a|0) \times p(b|0) \times p(c|x_0) \\ &= \frac{2}{3} \times \frac{1}{3} \times \frac{7}{3} \times \frac{1}{2} \\ &= \frac{2}{27} \end{aligned}$$

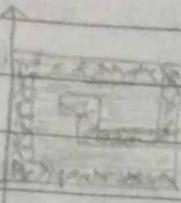
$$\begin{aligned} P(x|1) &= \frac{P(1|x)}{P(1|x) + P(0|x)} \\ &= \frac{\frac{4}{27}}{\frac{4}{27} + \frac{27}{27}} = \frac{4}{31} = \frac{2}{3} \end{aligned}$$

$$\begin{aligned} P(x|0) &= \frac{P(0|x)}{P(0|x) + P(1|x)} \\ &= \frac{\frac{2}{27}}{\frac{2}{27} + \frac{27}{27}} = \frac{1}{3} \end{aligned}$$

$$P(x|1) > P(x|0)$$

∴ It belongs to class 1.

Bayes decision Theory. (Boxer)



Linear $y = mx + c$.
Non linear quadratic.

Let $w_1 \rightarrow$ Probability of acceptance $P(w_1)$
 $w_2 \rightarrow$ Probability of rejection $P(w_2)$

Rule 1 :-

If $P(w_1) > P(w_2)$ then class (w_1)

Rule 2:-

If $P(w_2) > P(w_1)$ then class (w_2)

→ New input

A priori probability (APP)

It is completely based on History

[APP + some features (x)]

$x \rightarrow$ can be have various values.

$P(x|w_1) \rightarrow$ Probability density func" (PDF) of variable x for the object that belongs to w_1 .

Class conditional
PDF

$P(x|w_2) \rightarrow$ Probability density func" (PDF) of variable x for the object that belongs to w_2 .

→ NO of classes (number of regions)

Joint probability

$$P(w_i, x) = P(x|w_i) \cdot P(w_i)$$

→ By probability theory

$$P(w_i|x) \cdot P(x) = P(x|w_i) \cdot P(w_i)$$

$$P(w_i|x) = \frac{P(x|w_i) \cdot P(w_i)}{P(x)}$$

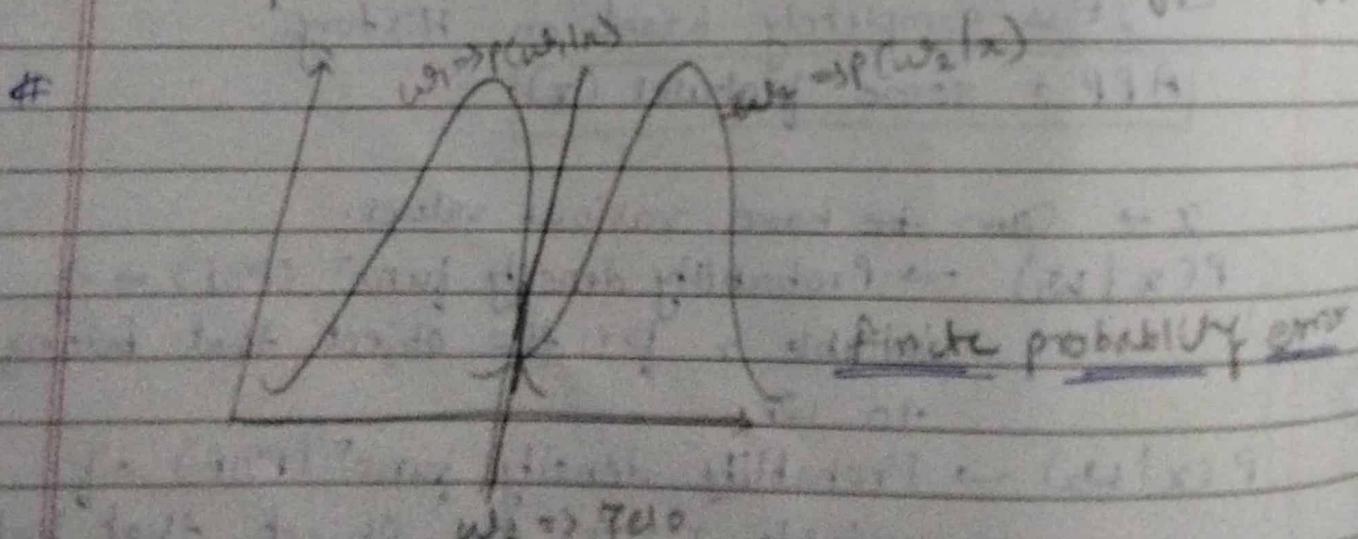
$$P(x) = \sum_{i=1}^n P(x|w_i) \cdot P(w_i)$$

Probability $P(w_1|x) > P(w_2|x) \rightarrow x \text{ belongs to } w_1$
if $P(w_2|x) < P(w_1|x) \rightarrow x \text{ belongs to } w_2$

$$P(x|w_1) \cdot P(w_1) > P(x|w_2) \cdot P(w_2) \rightarrow w_1$$

$$P(x|w_1) \cdot P(w_1) < P(x|w_2) \cdot P(w_2) \rightarrow w_2$$

If $P(w_1) \neq P(w_2)$ is same, the decision depends on class conditional probability density.



- $p(w_1|x)$ should be zero (left side)
If it has some value, then the finite probability error occurs. & vice versa.

Generalize \neq Baye's theory.

- Use more than 2 states of nature (classes)
- In this classification, use more than 1 feature.
If there are more than 1 feature then it is feature vector.
- Allow other actions other than merely deciding steps of nature.
- Introduce a loss function which is more general than probability vector.

Let us consider.

C no. of classes $\{w_1, w_2, w_3, \dots, w_c\}$

a - No. of actions $\{x_1, x_2, x_3, \dots, x_a\}$

$\lambda(x_i|w_j) \Rightarrow$ loss function.

{ Loss incurred for taking action x_i when true state nature is w_j }

$x \rightarrow d$ dimensional feature vector.

for x feature vector, we consider $a^m x_i$

$$R(x_i|z) = \sum_{j=1}^c \lambda(x_i|w_j) \cdot p(w_j|z)$$

(conditional Risk / Risk factor / Accepted loss .

Classifying z by taking $a^m x_i$

Minimum risk classifier.

• Minimum risk classifier classifier.

$w_1, w_2 \rightarrow$ Events

$x_1, x_2 \rightarrow$ Actions

$$\Rightarrow (x_1 | w_1) \rightarrow x_1$$

Action
class

$$R(x_i | z) = \lambda_{11} P(w_1 | z) + \lambda_{12} P(w_2 | z)$$

$$R(x_2 | z) = \lambda_{21} P(w_1 | z) + \lambda_{22} P(w_2 | z)$$

$w_i \Rightarrow$

$$R(x_1 | z) < R(x_2 | z)$$

$$[\lambda_{11} P(w_1 | z) + \lambda_{12} P(w_2 | z)] < [\lambda_{21} P(w_1 | z) + \lambda_{22} P(w_2 | z)]$$

$$(\lambda_{21} - \lambda_{11}) P(w_1 | z) > (\lambda_{12} - \lambda_{22}) P(w_2 | z)$$

Wrong class implied
action to make correct decision

$$\lambda_{21} P(w_1 | z) > \lambda_{12} P(w_2 | z)$$

Minimum Error rate Classifier.

We have to maximize the probability to decrease the risk.

$x_i \Rightarrow$ True State of nature is w_i .

$$(\lambda(x_i | w_i)) \quad \left\{ \begin{array}{ll} 0 & i=j \text{ if true condition} \\ 1 & i \neq j \text{ (High loss)} \end{array} \right.$$

Loss func

$$R(x_i | z) = \sum_{j=1}^2 \lambda(x_i | w_j) \cdot P(w_j | z)$$

$$\sum_{i \neq j} \lambda(x_i | w_j) \cdot P(w_j | z) \quad \begin{array}{l} \text{Not possible for one value} \\ \text{not possible condition} \end{array}$$

$\frac{2}{M} P(w_j | x)$ (Not possible condition)

$$R(x_i | x) = 1 - P(w_j | x)$$

This will minimize the risk.

Discriminant function, writing func. to every class

This classifier will compute the number equals to classes i.e. if no. of classes are c then we have to calculate c no. of risk + accepted risk / loss function / risk function.

or

$$\begin{cases} g_1(x) = P(w_1 | x) \\ g_2(x) = P(w_2 | x) \\ \vdots \\ g_c(x) = P(w_c | x) \end{cases}$$

are called
as
discrimi-
nant
function

Whichever discriminant func. gives max. val. we will put x into that class.

$\rightarrow g_i(x) > g_j(x)$ This condition is true only except $\forall j, i$ All values of i should not be equal to j .

If this condition is true then given feature vector belong to class w_i .

• $\sum_{i=1}^C P(w_i|x)$ (Not possible condition)

$$R(x|w_i) = 1 - P(w_i|x)$$

This will minimize the risk.

iii. Discriminant function. "assigning func" to every class

This classifier will compute the number equals to classes i.e. if no. of classes are C then we have to calculate C no. of risk + accepted risk / loss function / risk function.

or

$$\begin{cases} g_1(x) = P(w_1|x) \\ g_2(x) = P(w_2|x) \\ \vdots \\ g_C(x) = P(w_C|x) \end{cases}$$

are called
as
discrimi-
nant
function

Whichever discriminant func gives max. val. we will put x into that class.

$\rightarrow g_i(x) > g_j(x)$ This condition is true only except $\forall j \neq i$ All values of j should not be equal to i .

If this condition is true then given feature vector belong to class w_i .

ii) Minimum Error Classifier

$$g_1(x) = P(w_1|x)$$

for minimum error classifier, here require to find
risk func. to minimum if it is possible by many
 $P(w_1|x)$

As a result,

$(R(\alpha, 1)) = 1 - P(w_1|x))$, risk factor will be
reduced.

6. $g_1(x) = g_2(x)$

monotonically increasing function.

$$g_1(x) = P(w_1|x)$$

$$\cdot g_0(x) = \ln P(w_0|x) + \ln P(w_1|x)$$

- Q. The height of students are given below. Classify
new student having height 3 feet

Boys

6

7

5

7

6

5

4

3

5

6

Girls.

6

5

3

2

3

4

5

1

2

Prior probability of boy's $\Rightarrow 0.5$
 Prior probability of girl's $\Rightarrow 0.5$

Class conditional probability \Rightarrow

$$P(3|B) = P(3) / 10 = 1/10 = 0.1$$

$$P(3|G) = P(3) / 10 = 3/10 = 0.3$$

Joint prob:-

$$P(G|3) = \frac{P(3|G) \cdot P(G)}{P(3)}$$

$$= \frac{0.3 \times 0.5}{0.2} = 0.75$$

$$\begin{aligned} P(3) &= 0.1 + P(3|B) \times P(B) + P(3|G) \times P(G) \\ &= 0.1 \times 0.5 + 0.3 \times 0.5 \\ &= 0.2 \end{aligned}$$

$$P(B, 3) = \frac{P(3|B) \cdot P(B)}{P(3)}$$

$$\frac{0.1 \times 0.5}{0.2}$$

$$= 0.25$$

$$\therefore P(G, 3) > P(B, 3)$$

\therefore It is a girl.

- # Gaussian Parameter Estimation of Prob. Density Function.
- MLE (maximum likelihood estimation)

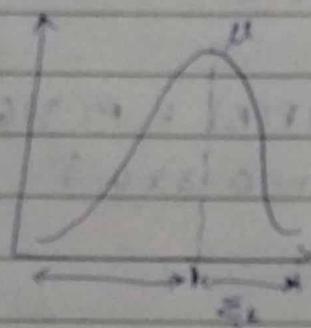
Gaussian distribution

$$w_0 = w_1, \dots, w_n$$

Gaussian distn

$$P(x|w_i) \sim N(\mu_i, \sigma_i^2), P(x|w_i)$$

Arg : Standard deviation (σ)
 If feature vector belongs to w_i



$$P(x|w_i) = \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x - \mu_i)^2}{2\sigma_i^2}}$$

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

$$\hat{\sigma}_i^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu}) \cdot (x_k - \hat{\mu})^T$$

Reinforcement
At every stage
will get some reward.

WORLD STARTM

Date : 38

Page : 7/9/22

Learning Rules.

- 1] Supervised - Labeled data / with teacher - Input & output are known
 - 2] Unsupervised - Unlabelled data / without teacher. - Input is known
Output is not known
 - 3] Reinforcement - Rewards.
-
- 1) Delta learning Rule Rule.
 - 2) Hebbian learning Rule
 - 3) Memory based learning Rule
 - 4) Competitive learning Rule
 - 5) Outstar learning Rule
 - 6) Boltzman learning rule.

① Delta learning rule:

- Widrow Hoff
- Error correctn LR
- ~~Mean~~ Mean square learning.
- Adaptive learning Elements (Adaline)

Let x_1, x_2, \dots, x_n are inputs connected with the strength of connections also called as weights w_1, w_2, \dots, w_m

w_{kj} is weight of connection from j^{th} input to k^{th} processing unit.

$$V_k = \sum_{i=1}^n x_i w_{ki} + x_0 w_{ko}$$

for p^{th} input

$$E_{kp} = T_{kp} - Y_{kp}$$

Error at Targeted Actual
point P point point
(swaragte) (satil bang)

$$E_{\text{tot}} = \frac{1}{n} \sum_{p=1}^n (T_{kp} - Y_{kp})^2$$

↓
Total error

$$\Delta W_{kj(p)} = n E_{kp} \cdot x_{j(p)}$$

↓ Change in weight learning rule ↓

$$W_{kj} = W_{kj(0)} + \Delta W_{kj(p)}$$

↓ Old weight ↓

② Hebbian Learning Rule / Activity Product Rule.

If one neuron simulates second neuron & frequently contributes in its activation, the synaptic association bet' these 2 neurons is strengthened & second neuron become extra sensitive to stimuli from 1st neuron.

$$\Delta W_{kj(p)} = f(x_{j(p)} \cdot y_{kp})$$

sum function of product

$$\Delta W_{kj(p)} = n (x_{j(p)} \cdot y_{kp})$$

$$\Delta W_{kj(p)} = n (x_{j(p)} \cdot y_{kp}) - \phi n (x_{j(p)} \cdot y_{kp})$$

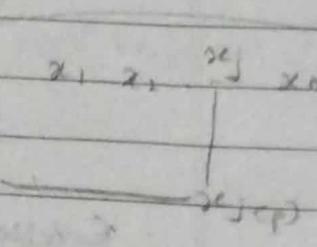
[0-1 range
[0.1 - 0.001]]

③ Memory based learning / Instance based learning
Unsupervised learning

$$X = \{x_1, x_2, x_3, \dots, x_m\} \text{ in dimension } d,$$

$m \Rightarrow$ dimension

$d =$ training patterns / no. of feature vectors.



↓
input vector (x_1, x_2, \dots, x_n)
pth pattern ($x_{j(p)} - x_j$)

nearest neighbour } k-nearest neighbour
algo. used in ml for

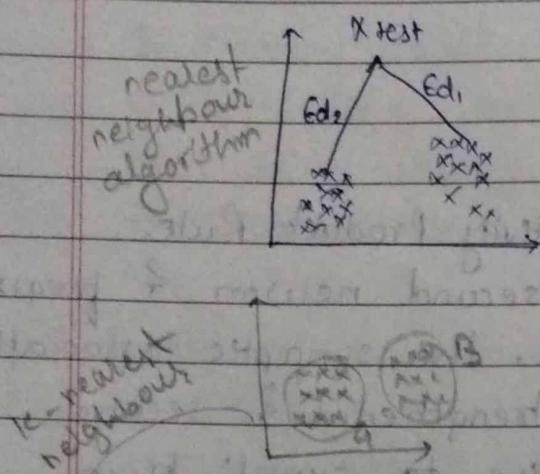
WORLD STAR
Date: 9/8/20
Page: 40

$$\{x_i\}_{i=1}^N$$

$$\{d_i\}_{i=1}^N$$

Euclidean $\{\sqrt{x_i - d_i}\}_{i=1}^N$

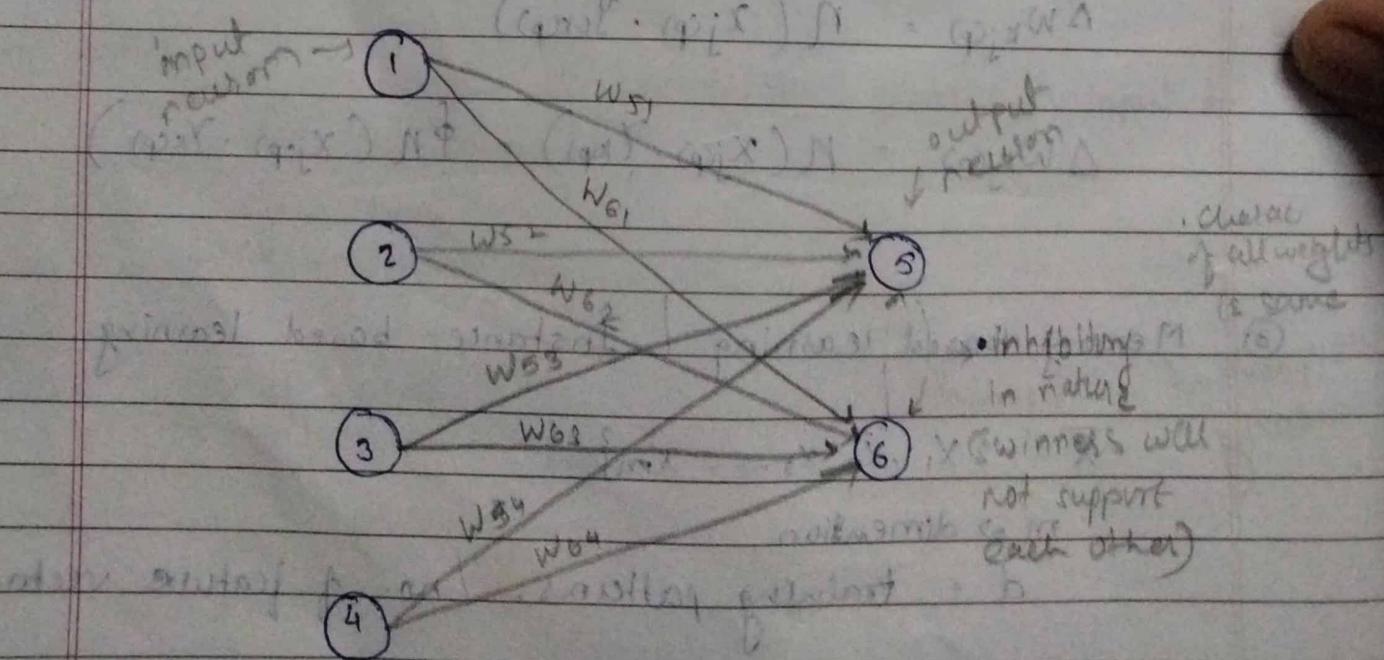
Euclidean distance.



$Ed_1 < Ed_2 \rightarrow x \text{ test will classify into}$
 Ed_2

To minimize the error, we have to find out mean of the inputs. Then this is known as k-nearest neighbour algo.

(4) Competitive learning Rule.



- Condition:-

$$w_{S1} + w_{G1} = 1$$

$$\text{i.e. } \sum_{j=1}^n w_{kj} = 1$$

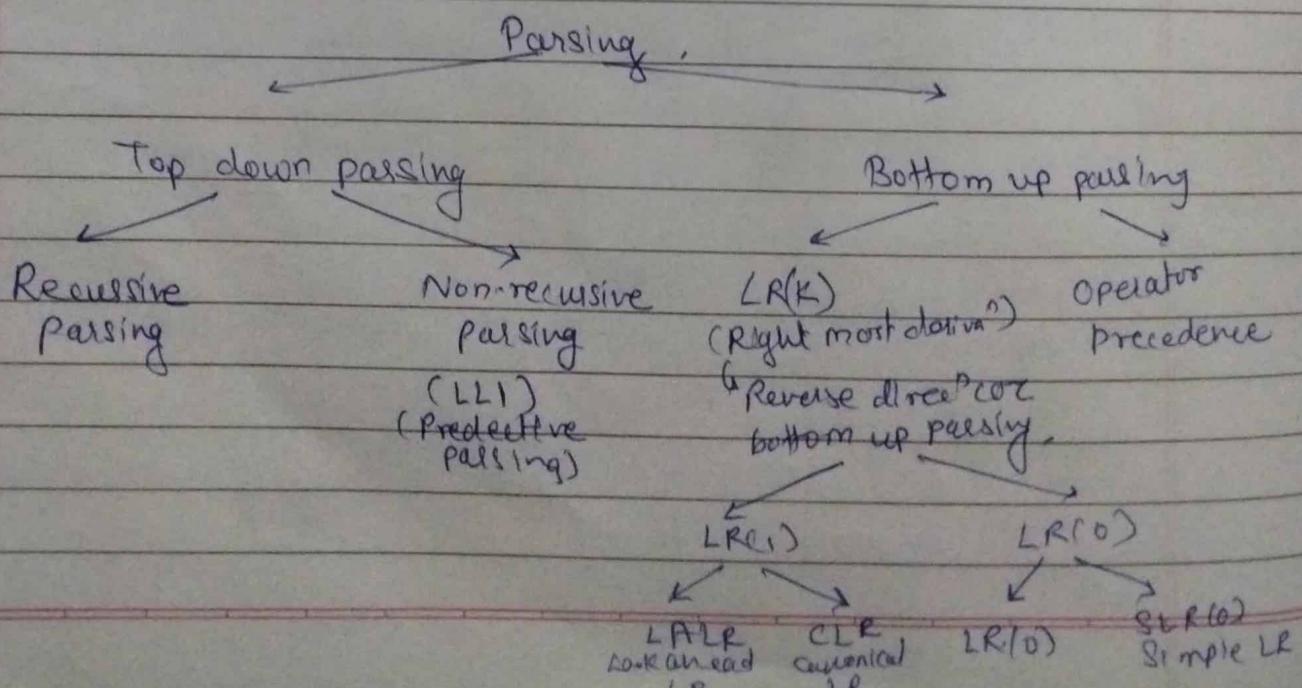
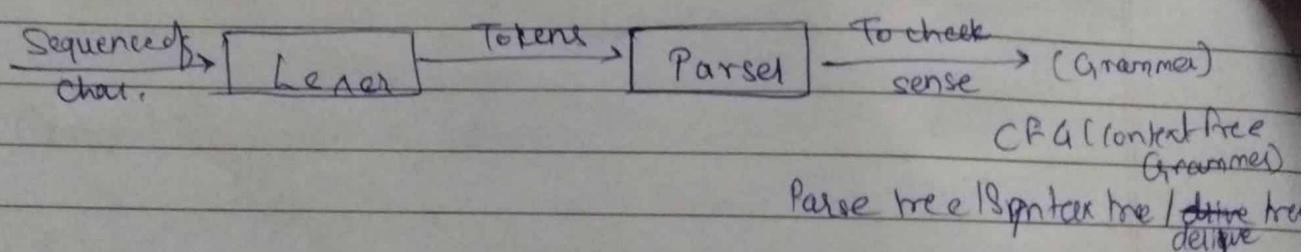
- Winner will be favour.

- In this algo. Winner takes all neurons.

Parsing

- Parsing is a process of deriving string from the given grammar.
- Grammar: Study of the way words are used to make sentences.
- Rule of language governing the sounds, words, sentences & other elements as well as their combination of interpretation (grammar).

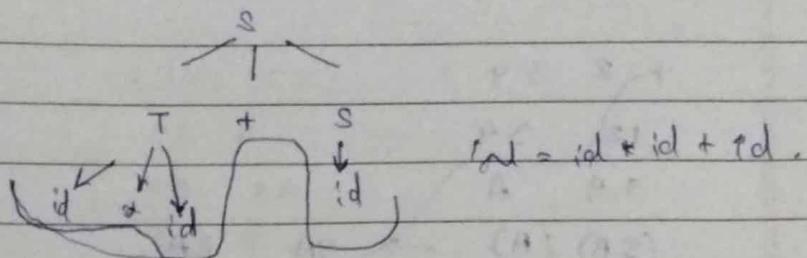
- Lexical Analysis.
- It is a process of converting sequence of characters into sequence of steps.
- In computer science, lexical Analysis is process of converting a sequence of characters (such as in computer pages / web page), into a sequence of lexical tokens (string with page).
- Lexer, Tokenizer or scanner



* Ambiguous grammar: A CFG is said to be ambiguous if there exists more than one derivation tree for the given input string ie more than one LADT) more than one RMDT.

$$S \rightarrow T + S \mid T$$

$$T \rightarrow id \cdot T \mid id \mid s$$



* Unambiguous grammar: CFG (Context Free Grammar)

CYK (Parsing) Algorithm.

Code Younger Karmi

$$S \rightarrow AB$$

$$A \rightarrow BB \mid a$$

$$B \rightarrow AB \mid b$$

- CFG $\xrightarrow{\text{form}} \text{CNF} \mid \text{GNF}$

Chomsky Normal form

Arebach Normal form.

- CYK only support CNF.

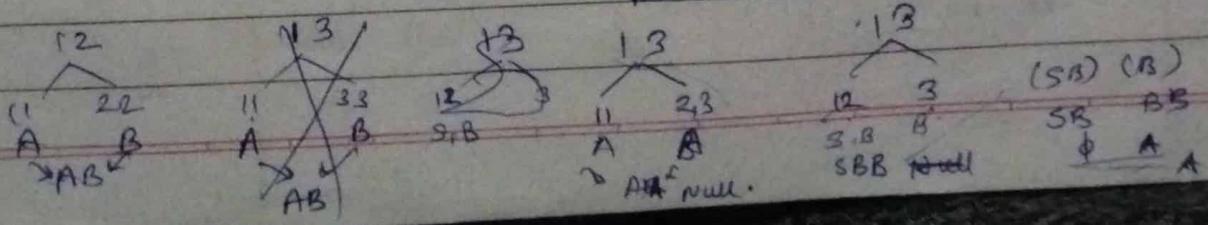
- It is universal (applicable for all CNF)

Eg: Pattern: $\rightarrow abbb$.

a b b b
1 2 3 4

L to L
(11) (22)
 $\begin{matrix} a \\ \downarrow \\ A \end{matrix}$ $\begin{matrix} b \\ \downarrow \\ B \end{matrix}$

	4	3	2	1
1	S, B	S, B A	S, B	A
2	S, B A	A	B	
3	A	B		
4	B			



~~2, 3, 4.~~

23 4	2 3 4
23 44	22 < 34
A B	B A
AB	BA
S, B	Null. (\emptyset)

~~1 2 3 4~~

12 34	SA	BA
S, B A	\emptyset	\emptyset
(SB) (A)		

~~1 2 3 4~~

(11) 234	(125) 44
A S, B	A B
(AS) (AB)	\rightarrow S, B
\emptyset <u>S, B</u>	

~~a b a b a~~

~~a b a b a~~

1	2	3	4	5
---	---	---	---	---

$S \rightarrow AB$

$A \rightarrow BB/a$

$B \rightarrow AB/b$

	5	4	3	2	1
1	Null	S, B	A	S, B	A
2	Null	AB	Null	B	
3	A	S, B	A		4, 5
4	Null	B		44 55	23 44 BA BB 8 A
5	A			B A	

1, 2, 3	12, 3	2, 3, 4	3, 4, 5
12 33	11, 23	23 44 22 34	33 45 34 55
S, B A	A \emptyset	\emptyset B B, S, B	A \emptyset S, B A
SA BA	A	<u>B</u> (BS) (BB)	— A SA BA
\emptyset			\emptyset \emptyset

$$\begin{array}{c} 1234 \\ (11) \quad (234) \\ A \quad AB \\ AB \quad AB \\ \phi \quad \underline{S,B} \end{array} \qquad \begin{array}{c} (123)(44) \\ A \quad B \\ AB \\ \underline{S,B} \end{array}$$

$$\begin{array}{c} 2345 \\ (22) \quad (3,45) \\ B \quad A \\ \phi \end{array} \qquad \begin{array}{c} (234)(55) \\ AB \quad A \\ AA \quad BA \\ \phi \quad \phi \end{array}$$

$$\begin{array}{c} 12345 \\ (11) \quad (2345) \\ A \quad \phi \\ A \end{array} \qquad \begin{array}{c} (1234)(55) \\ S,B \quad A \\ \underline{SA} \quad BA \quad \phi \end{array}$$

#1 Syntactic Analysis System

- ↳ Phrase structure grammar
- ↳ The transformational grammar
- ↳ Transition network grammar.

1) Phrase structure grammar.

Purely syntactic analysis of language does not do very much in language processing.

- Syntax adjust to a more complete analysis process.
- Syntax is better understood than most other linguistic procedure.
- Linguists are noted for their work on generative grammars.
- Generative grammars are used to generate grammatically correct statements.
- Generative grammars have many property of simplicity.
- Particular type of grammar is known as phrase structure.
- Simple & usable both for sentence generation & for sentence analysis.

Consider the rewrite rules of the form

$$S \rightarrow A + B$$

↳ statement

S followed by A & B

- Variable S can be rewritten as A followed by B
- The symbol plus separates the variables.
- When presenting rewrite rules capital letters denote non terminal elements. i.e elements that can be rewritten further, by appearing on the left side of some rewrite rule.

Where the bar stands for logical connective OR.
 The result would be huge grammar which would not be simple.

Eg: the man hit the ball.

the man hit the man.

the ball hit the ball.

the ball hit the ball.

the man took the man.

the man took the ball.

the ball took the man.

the ball took the ball.

Derivation tree.

~~Statement~~ S → NP + VP

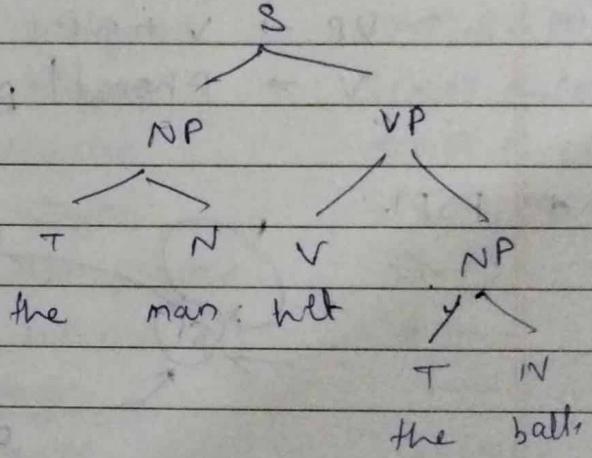
* NP → T + N.

T → the

N → man, ball, etc. →

* VP → V + NP

V → hit / took.



Context free

Grammars such that if expression 1 is known as context free grammar phrase structure grammar bcoz non terminal symbols of LHS of the rewrite rules can be replaced by the right sides of the rules regardless of the context in which these symbols may appear, ie no restrictions apply to any rewrite rule.

John phoned Mary.

Mary phoned John.

John phoned John.

Mary phoned Mary.

John phoned up Mary

John phoned up John.

Mary phoned up John.

Mary phoned up Mary

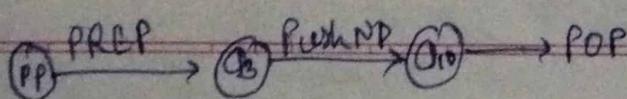
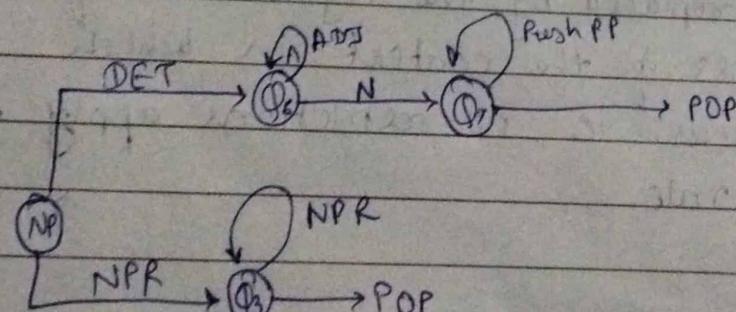
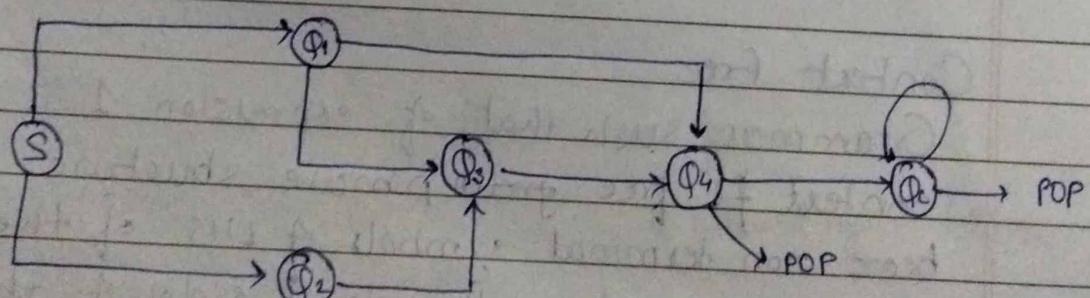
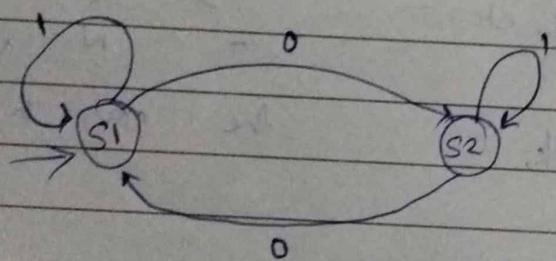
S → NP + VP.

NP → John | Mary

VP → V + NP

V → phoned | phoned up.

* FSM.



The tall man in the Stetson is John Wayne.

	Current State	IIP String rec.	Stack.
1	S		--
2	NP	The tall man --	
3	Q6	The tall man --	S(NP)
4	Q6	tall man in --	S(NP)
5	Q7	man in the --	S(NP)
6	PP	in the Stetson --	S(NP), Q7(PP)
7	Q9	the Stetson is John --	S(NP), Q7(PP)
8	NP	the Stetson is John --	S(NP), Q7(PP), Q9(NP)
9	Q6	Stetson is John Wayne	S(NP), Q7(PP), Q9(NP)
10	Q8	is John Wayne	S(NP), Q7(PP), Q9(NP)
11	POP	is John Wayne	S(NP), Q7(PP), Q9(NP)
12	Q10	is John Wayne	S(NP), Q7(PP)
13	POP	is John Wayne	S(NP), Q7(PP)
14	POP	is John Wayne	S(NP)
15	Q1	is John Wayne	S --
16	Q4	John Wayne	←
17	NP	John Wayne	Q4(NP)
18	Q8	Wayne	Q4(NP)
19	Q8	—	Q4(NP)
20	POP	—	—
21	Q4	—	—
22	Pop	Accept	—