

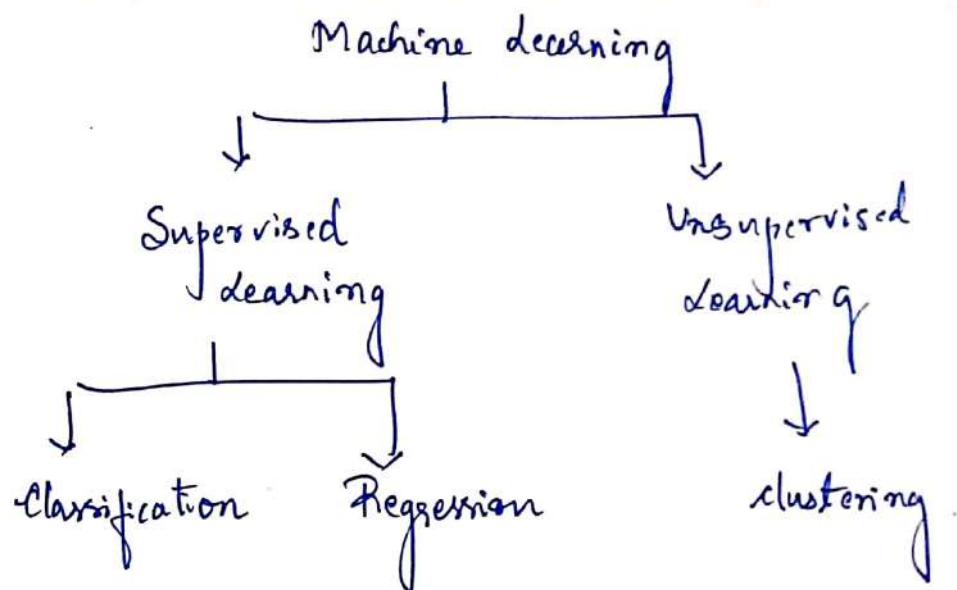
Computational Intelligence

↳

Machine Learning.

Basics :-

- The ability of machine to categorize correctly new Examples that differ from those used for training is known as Generalization.
- Machine learning is done in 2 phases
 - Training phase
 - Testing phase
- * In training phase we train the machine with N variants of a random continuously varying data say x .
- * In testing phase a new variant of the continuously varying quantity x is given to test the prediction of the Machine learning Algorithm.
- There can be 2 ways of machine learning
 - ↳ Supervised
 - ↳ un-Supervised.



Linear Basis Function Models :- [Linear Regression]

- Regression analysis is a form of predictive modeling technique which gives relationship between dependent (target) & independent variables (predictor)
- linear Regression is a statistical Model which gives linear relationship between two or more variables.

Note :- Regression algorithms are used to predict the continuous values such as price, salary, age etc.

- Simple Linear Model of regression is given as

$$y(x, w) = w_0 + w_1 x_1 + \dots + w_D x_D$$

where $x = (x_1, \dots, x_D)^T$

$x \rightarrow$ independent variable

$y \rightarrow$ dependent variable.

→ To Extend the class of models to include non-linearity of g/p variables Equation is modified as

$$y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x)$$

↓
Basis function
Total m parameters

* $w_0 \rightarrow$ Bias parameter

$$y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x) = w^T \phi(x)$$

where $w = (w_0, \dots, w_{M-1})^T$ & $\phi = (\phi_0, \dots, \phi_{M-1})^T$

Note :- Here original variables are vector x , features of x is denoted as $\{\phi_j(x)\}$

→ By using non-linear Basis fn we allow $y(x, w)$ to non-linear fn of g/p vector x .

→ Basis function

- polynomial
- Gaussian
- Sigmoidal

* polynomial $\phi_j(x) = x^j$

* Gaussian $\phi_j(x) = \exp\left\{ -\frac{(x-\mu)^2}{2s^2} \right\}$

* Sigmoidal $\phi_j(x) = \sigma\left(\frac{x - \mu}{s}\right)$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

Note :- Other basis func is Fourier Basis, wavelets.
We can also assume $\phi(x)$ to simply equal to x .
 $\phi(x) = x$.

Maximum Likelihood & Least Squares :-

* Basics to be known

- Polynomial curve fitting
- Gaussian function
- Bayes theorem.

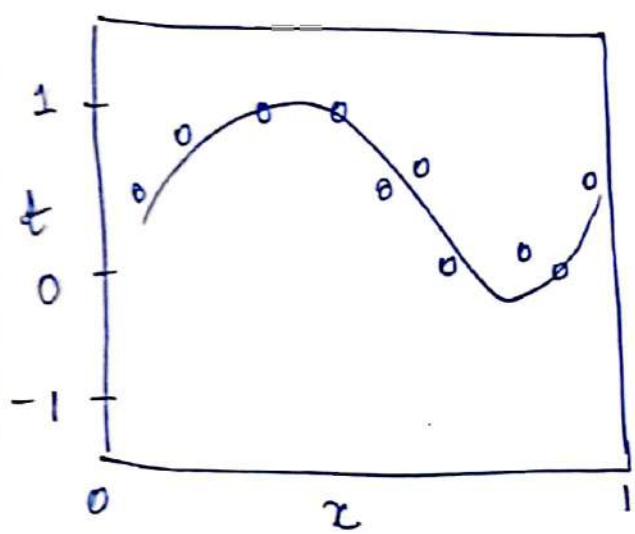
→ Applications in which training data comprises Examples of input vectors along with their corresponding target vectors are known as Supervised learning problems. If the desired output consists of one or more continuous variables, then task is called regression.

→ If training data consists of a set of input vector x without any corresponding target values. This is unsupervised learning where we need to discover groups of similar examples in the data. called as clustering.

→ The technique of reinforcement learning is concerned with the problem of suitable actions to be taken in a given situation in order to maximize a reward. Here learning is an error process.

Ex:- polynomial curve fitting

→ Given a training set with N observations of x
 $x = (x_1, \dots, x_N)^T$ with $t = (t_1, \dots, t_N)^T$



* x is uniformly spaced in range $[0, 1]$
* t is calculated by taking an example function $\sin(2\pi x)$ with small level of random noise.

→ Our goal is to exploit this training set in order to make predictions of \hat{t} of the target for some new value \hat{x} of the input variable. & discover the underlying function $\sin(2\pi x)$ through Machine Learning Algorithm.

→ Polynomial function is given as (6)

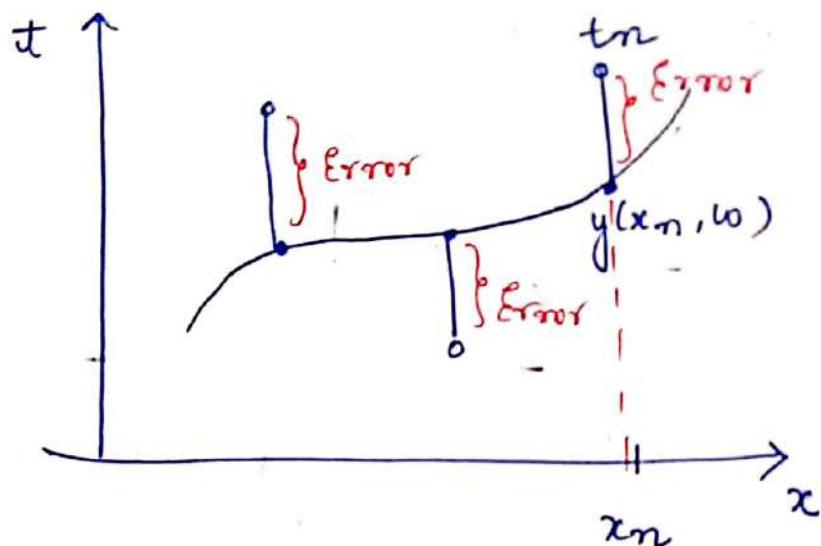
$$y(x, w) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M$$

$$= \sum_{j=0}^M w_j x^j$$

- * $M \rightarrow$ Order of polynomial
- * polynomial co-efficients w_0, \dots, w_M are denoted as vector w .
- * polynomial function $y(x, w)$ is a non-linear function of x , and linear function of w .
- The values of co-efficients will be determined by fitting the polynomial to the training data
- In fitting the polynomial curve to the training data we occur with error or misfit ~~is~~ between $y(x, w)$ for any given w & x .
- To minimize error widely used error function is square of errors between $y(x_n, w)$ ^{for} data point x_n & corresponding target t_n .

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2$$

→ fig below shows the geometrical interpretation of sum of square errors. (7)

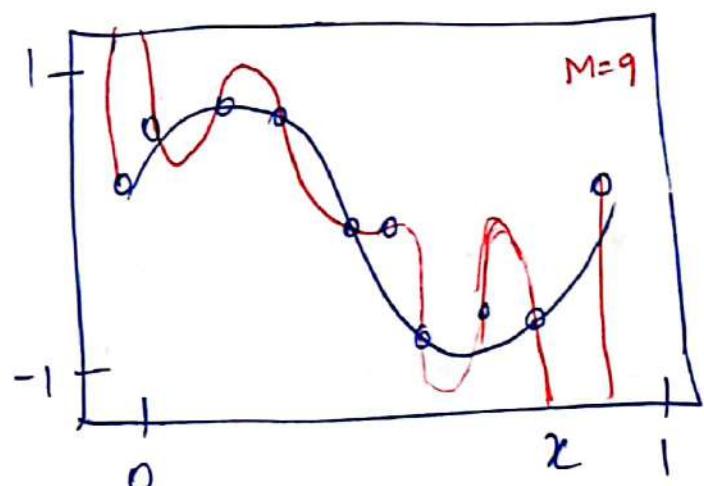
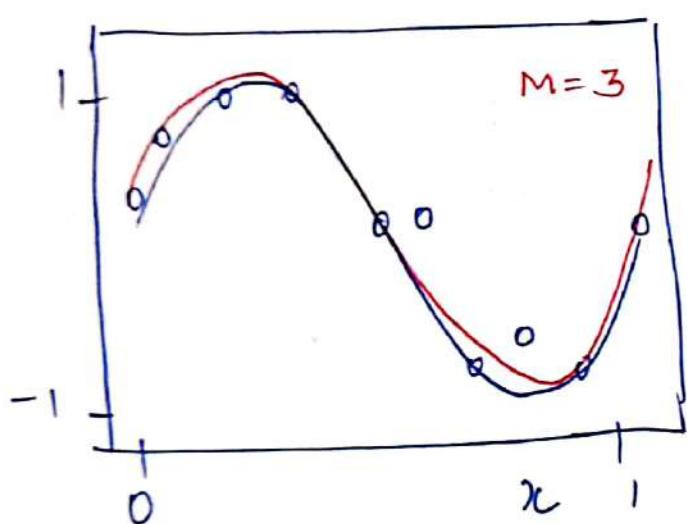
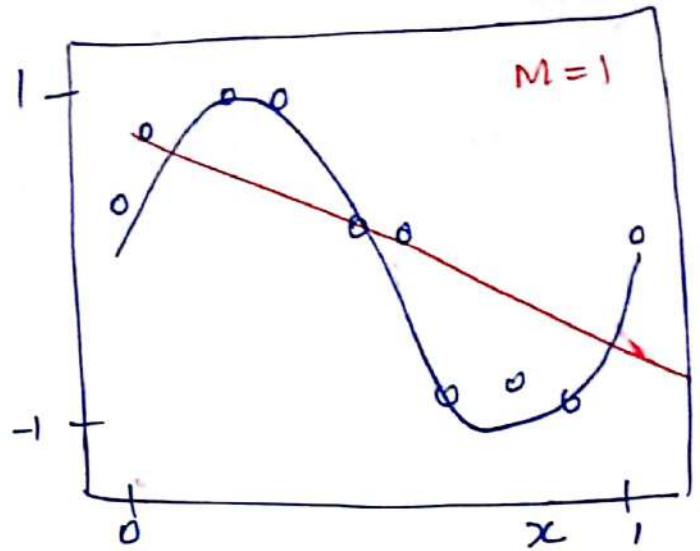
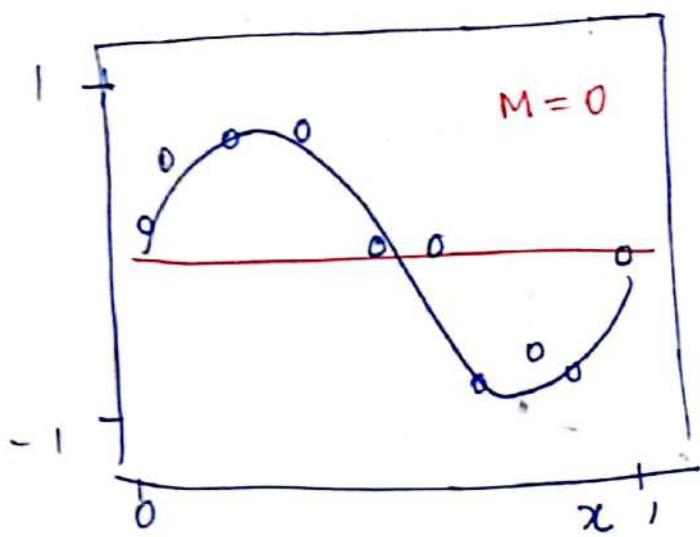


* Here we can solve the curve fitting problem by choosing the value of w . $E(w)$ is as small as possible.

* Resulting polynomial with changed w is given as $y(x, w^*)$

→ Other than reducing $E(w)$, there remains the problem of choosing order of M of the polynomial.

→ -



→ we notice that with $M=0$ & $M=1$ polynomials
we get poor fits of data [we need a $\sin(2\pi x)$
function fit]

Note :- Blue colour curve → Actual fit needed
Red colour curve → fit got with the order
of polynomial.

→ with $M=3$ we get Best fit to the function
 $\sin(2\pi x)$

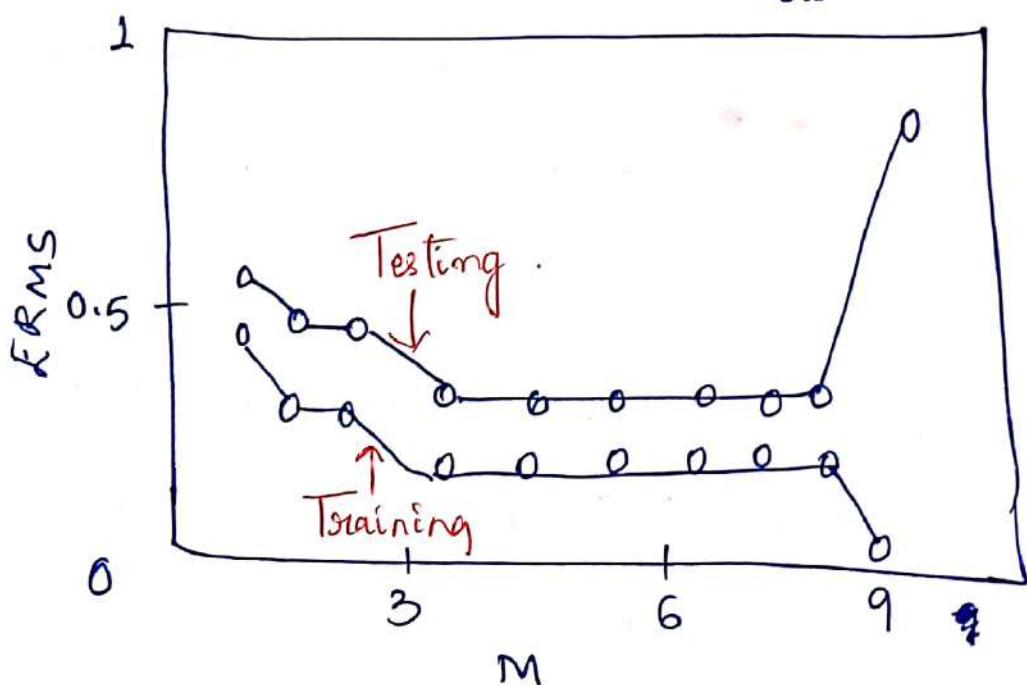
→ with $M=9$ we get Excellent fit to the
training data .

→ In fact polynomial passes exactly through each data point & $E(w^*) = 0$; But curve oscillates wildly & gives poor representation of the function $\sin(2\pi x)$. This is called over-fitting.

→ It is sometimes more convenient to use root-mean-square of error function

$$E_{RMS} = \sqrt{2E(w^*)/N}$$

↓
different sizes of data sets.



* We note in above fig that for $3 \leq M \leq 8$ we get small values of E_{RMS} (error).

* For $M=9$, training set error goes to zero

(10)

but test-set error becomes very large, \therefore giving wild oscillations in $y(x, w^*)$ & the predictions.

- The over-fitting problem becomes less severe as the size of data-set increases.
- Over-fitting problem is a general property of maximum likelihood, By adapting Bayesian approach, overfitting problem can be avoided.

Note :- In Bayesian Model the effective number of parameters adapts automatically to the size of the data set.

- One more approach to control the over-fitting phenomenon is regularization. Regularization involves adding a penalty term to the error function.

$$\tilde{E}(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \frac{\lambda}{2} \|w\|^2$$

where $\|w\|^2 = w^T w = w_0^2 + w_1^2 + \dots + w_M^2$ ⑩

- the coefficient λ governs the relative importance of the regularization term.

Probability theory :-

→ Probability theory provides a consistent framework for the quantification & manipulation of uncertainty, when combined with decision theory. will allow us to make optimal predictions given all the information available to us, even though that information may be incomplete or ambiguous.

Rules of Probability :-

Sum rule :- $p(x) = \sum_y p(x, y)$

Product rule :- $p(x, y) = p(y|x)p(x)$

Note :- $p(x, y) \rightarrow$ joint probability {probability of $x \& y$ }

$p(y|x) \rightarrow$ probability of y given x .

$p(x) \rightarrow$ marginal probability ./ probability of x

→ From product rule, together with the symmetry property $p(x, y) = p(y, x)$

(15)

we obtain the following relationship between
Conditional probabilities .

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

Bayes theorem

* Baye's theorem plays a central role in pattern
recognition & machine learning.

* Denominator of Baye's theorem can be written

as $P(x) = \sum_y P(x,y) \rightarrow \text{sum rule}$

$$P(x_{ref}) = \sum_y P(x|y) P(y)$$

$$(\because P(x,y) = P(y|x) P(x))$$

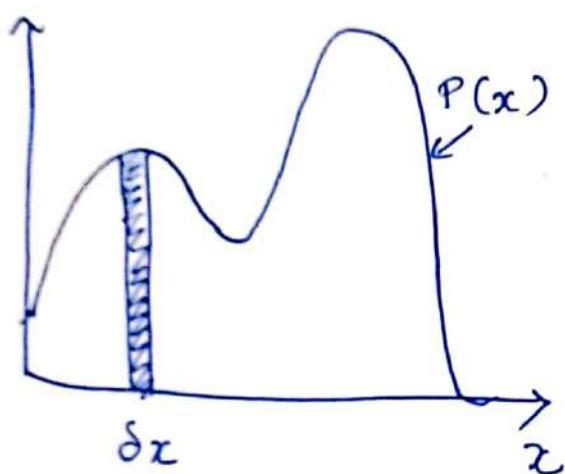
$$\therefore P(x) = \sum_y P(x|y) P(y)$$

↓
normalization constant-

Probability densities :-

→ we wish to consider probabilities with continuous variables rather than just discrete set of events .

→ If we have to take probability of continuously varying real-valued variable x falling in the interval $(x, x + \delta x)$ is given as $P(x) \delta(x)$ for $\delta x \rightarrow 0$, then $P(x)$ is called ~~Prob~~ Probability density function of ~~over~~ over x .



→ The probability that x will lie in an interval (a, b) is given by

$$P(x \in (a, b)) = \int_a^b P(x) dx$$

→ Probabilities are nonnegative

$$\boxed{\begin{aligned} & \therefore P(x) \geq 0 \\ & \int_{-\infty}^{\infty} P(x) dx = 1 \end{aligned}}$$

→ $P(x)$ must satisfy these two conditions.

→ Most important operations involves probabilities is that of finding weighted averages of functions
 → The average value of some function $f(x)$ under a probability distribution $p(x)$ is called Expectation of $f(x)$. It is denoted by $E[f]$.

Discrete Distribution

$$E[f] = \sum_x p(x) f(x)$$

Continuous variables

$$E[f] = \int p(x) f(x) dx$$

→ Applying / Re-writing concept of Bayes theorem with respect to linear Basis function model

$$y(w|t) = \frac{P(t|w) P(w)}{P(t)}$$

↓
Posterior Probability

Likelihood

Prior Probability Info.

* w → weights associated with Input

t → target value achieved when weight
weights are set

→ Bayes's theorem is used to convert a prior probability into a posterior probability by incorporating the Evidence provided by the observed data.

$\text{Posterior} \propto \text{likelihood} \times \text{prior}$

Note :- all of the quantities in Bayes's probability theorem are viewed as function of w .

* The denominator is normalization constant which

Ensures that posterior distribution on d-H-S is a valid probability density & integrates to one.

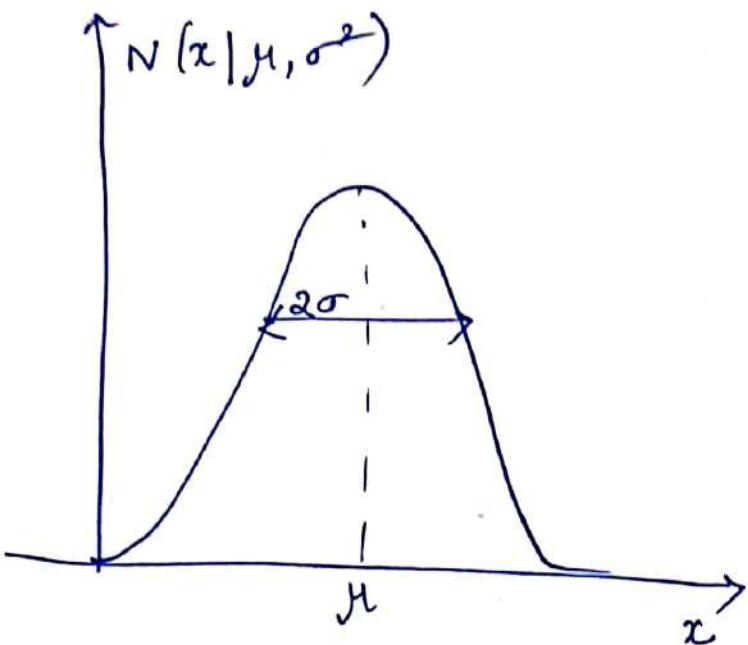
- In Baye's theorem, the likelihood function plays a central role.
- A widely used estimator of w is maximum likelihood in which w is set to a value that maximizes likelihood function $P(t|w)$ $P(t|w)$.

The Gaussian Distribution:-

- Let us assume that the random variable which is continuously varying & is used to train the Machine has a gaussian distribution instead of polynomial distribution, discussed earlier.
- For a real-valued variable x , Gaussian distribution is defined as .

$$N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x-\mu)^2 \right\}$$

Gaussian is also called normal distribution.



$\mu \rightarrow$ mean

$\sigma^2 \rightarrow$ variance

$\sqrt{\sigma^2}$ / Variance = σ = std deviation

$\frac{1}{\text{Variance}} = \frac{1}{\sigma^2} = \beta$ [precision]

* Above Gaussian distribution satisfies the below Condition

$$\boxed{N(x | \mu, \sigma^2) > 0}$$

$$\boxed{\int_{-\infty}^{\infty} N(x | \mu, \sigma^2) dx = 1}$$

* Expectation of x under this Gaussian distribution is given as

$$\boxed{E[x] = \int_{-\infty}^{\infty} N(x | \mu, \sigma^2) x dx = \mu}$$

average value
of x

→ We can obtain or Maximize likelihood function by writing the probability of data set, given μ, σ^2

Gaussian distribution is given as (7)

$N(x|\mu, \sigma^2) \rightarrow$ taking probability of this distribution w.r.t μ, σ^2

w.r.t μ

$$P(x|\mu, \sigma^2) = \prod_{n=1}^N N(x_n|\mu, \sigma^2)$$

* It is more convenient

to maximize log of likelihood function.

∴ Above Equation becomes

$$\ln p(x|\mu, \sigma^2) = -\frac{1}{2\sigma^2}$$

$$\sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

∴ ∵ Likelihood or w.r.t μ is maximizing the above expression. we get

$$\boxed{\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n}$$

max likelihood

w.r.t β

$$P(t|x, w, \beta) = N(t|y(x, w), \beta^{-1})$$

$$P(t|x, w, \beta) = \prod_{n=1}^N N(t_n|x_n, w, \beta^{-1})$$

$$\boxed{\beta = \frac{1}{\sigma^2} \quad \therefore \beta^{-1} = \sigma^2}$$

∴ taking log for above Expression we get

$$\ln p(t|x, w, \beta) = -\frac{\beta}{2}$$

$$\sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \frac{N}{2} \ln(2\pi)$$

$$+ \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

$$\boxed{\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, w_{ML}) - t_n\}^2}$$

→ Now returning back to maximum likelihood

& least square concept

→ Let us consider target t with Gaussian noise model

$$\begin{aligned} t &= y(x, w) + \epsilon \leftarrow \text{Gaussian noise} \\ y(x, w) &= w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x) \quad \begin{array}{l} \text{linear} \\ \text{Basis for} \\ \text{model} \end{array} \end{aligned}$$

→ Taking probability of the above Expression with precision parameter from Gaussian distribution

$$p(t|x, w, \beta) = N(t | y(x, w), \beta^{-1})$$

↑ normal distribution ↑ inverse variance
precision

$$\left[\frac{1}{\sigma^2} = \beta \quad \frac{1}{\beta^{-1}} = \sigma^2 \right]$$

↓ Variance

* Here Expectation of new x to a target t

is given as

$$E[t|x] = \int t \cdot p(t|x) dt = y(x, w)$$

→ Now consider a data set of inputs

$$x = \{x_1, \dots, x_N\} \text{ with } t_1, \dots, t_N$$

→ We can obtain the expression for the likelihood function as a function of adjustable parameter w, β (19)

$$\therefore p(t|x, w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$$

$y(x, w) = w^T \phi(x_n)$

→ To maximize this likelihood function

$$P(w|t) = \frac{p(t|w) \cdot P(w)}{\text{J. likelihood } P(t)}$$

→ Bayes theorem

↳ Keeping notation of LHS dropping x we get

$$P(t|w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$$

→ Taking log on both sides to maximize the above function $p(t|w, \beta)$ [refer gaussian distribution] we get

$$\ln p(t|w, \beta) = \sum_{n=1}^N \ln N(t_n | w^T \phi(x_n), \beta^{-1})$$

$$= \left[\frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \right] - \beta^{-1} E_D(w)$$

↑ Likelihood ↓ Error to Likelihood

→ (1)

where where $E_D \rightarrow$ Sum of Square Error. fn⁽²⁰⁾

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2$$

→ We can maximize the likelihood function under a conditional Gaussian noise distribution by minimizing a sum-of-squares Error given by $E_D(w)$:

→ Minimize Minimize by taking gradient of $E_D(w)$

$$\nabla \ln p(t|w, \beta) = \beta \frac{d}{dw} \left[\frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2 \right]$$

~~* Simply differentiate by w~~

βE_D .

$$0 = \frac{d}{dw} \left[\frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2 \right]$$

Assume this

[const's = 0]

Note :- above Expression is $\frac{dc}{dx}$ to

$$\frac{d}{dx} x^2 = 2x \frac{d}{dx} x$$

$$\therefore 0 = \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\} \cdot \frac{d}{dw} \{t_n - w^T \phi(x_n)\}$$

(21)

$$\therefore 0 = \sum_{n=1}^N \cancel{x} \cdot \frac{1}{\cancel{x}} (t_n - w^T \phi(x_n)) \cdot [0 - \phi^T(x_n)]$$

$$0 = \sum_{n=1}^N (t_n - w^T \phi(x_n)) \cdot \phi^T(x_n)$$

$$0 = \sum_{n=1}^N t_n \phi(x_n)^T - \sum_{n=1}^N w^T \phi(x_n) \cdot \phi^T(x_n)$$

∴ For making the equation in generalized form
we will write the above equation as

$$\rightarrow 0 = t \cdot \phi(x_n)^T - w^T \phi(x_n) \phi(x_n)^T$$

$$\rightarrow t \cdot \phi(x_n)^T = w^T \phi(x_n) \phi(x_n)^T$$

or

$$\rightarrow w \phi(x_n)^T \phi(x_n) = t \cdot \phi(x_n)^T$$

$$\rightarrow \therefore w \phi(x_n)^T \phi(x_n) \cdot \left[\phi^T(x_n) \cdot \phi(x_n) \right]^{-1} = \left[\phi^T(x_n) \cdot \phi(x_n) \right]^{-1} \cdot t \cdot \phi(x_n)^T$$

$\left[\because A \cdot A^{-1} = I / \phi^T(x_n) \cdot \phi(x_n)^{-1} = I \right]$

$$\therefore \underbrace{w \phi(x_n)^T \phi(x_n)}_{w_{ML}} \cdot \left[\phi^T(x_n) \cdot \phi(x_n) \right]^{-1} = \left[\phi^T(x_n) \cdot \phi(x_n) \right]^{-1} \cdot t \cdot \phi(x_n)^T$$

$$\boxed{\therefore w_{ML} = [\phi^T \phi]^{-1} \cdot \phi^T \cdot t}$$

→ The quantity $w_{ML} = (\phi^T \phi)^{-1} \phi^T t$

is

↓
normal Equations for
least squares problem

$\phi \rightarrow n \times m$ matrix called Design matrix

$$\phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{m-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{m-1}(x_2) \\ \vdots & & & \\ \phi_0(x_n) & \phi_1(x_n) & \dots & \phi_{m-1}(x_n) \end{bmatrix}$$

The quantity $\phi = (\phi^T \cdot \phi)^{-1} \phi^T$ is known as Moore - Penrose pseudo inverse of matrix ϕ . [ϕ is a $n \times m$ matrix]

Note :- $w_{ML} = [\phi^T \phi]^{-1} \cdot \phi^T \cdot t$

\downarrow \downarrow \downarrow \downarrow
 $(m \times n)(n \times m)$ $n \times m$ $n \times 1$

$m \times m$

$w_{ML} = (m \times 1) \rightarrow$ Row vector
gs a.

→ 11th we can maximize the log likelihood function

$$\ln p(t|\omega, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\omega)$$

wrt to precision parameter β .

$$\begin{aligned} \frac{d}{d\beta} [\ln p(t|\omega, \beta)] &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\omega) \\ &= \frac{N}{2} \cdot \frac{1}{\beta} - 0 - E_D(\omega) \end{aligned}$$

Equating above equation to zero in order to minimize it

$$\frac{N}{2\beta} - E_D(\omega) = 0$$

$$\frac{N}{2\beta} = E_D(\omega)$$

$$\frac{1}{\beta} = \frac{2}{N} \cdot E_D(\omega)$$

$\left. \frac{1}{\beta} \right|_{\text{subst } E_D(\omega)}$

$$\therefore \frac{1}{\beta_{ML}} = \frac{2}{N} \cdot \frac{1}{\sum_{n=1}^N (t_n - \omega^T \phi(x))^2}$$

$$\therefore \frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N (t_n - \omega^T \phi(x))^2$$

∴ we can conclude that to maximize the likelihood function by 2 parameters w, β . (24)

wrt w :-

$$w_{ML} = (\phi^T \phi)^{-1} \phi^T t$$

wrt β :-

$$\beta_{ML} = \frac{1}{N} \sum_{n=1}^N \{t_n - w_{ML}^T \phi(x_n)\}^2$$

(24)

∴ we can conclude that to maximize the likelihood function by 2 parameters w, β .

wrt w :-

$$w_{ML} = (\phi^T \phi)^{-1} \phi^T t$$

wrt β :-

$$\underset{\beta_{ML}}{\perp} = \frac{1}{N} \sum_{n=1}^N \{t_n - w_{ML}^T \phi(x_n)\}^2$$

Sequential learning :-

- When we have input data set sufficiently large & we use the complete set or entire batch for training the machine [Batch processing] the computationally costly
- ∴ we can use sequential algorithms as an alternative which is also known as on-line algorithms.
- Here data sets of the input are considered one at a time & parameters are updated each time.

→ We can obtain a Sequential learning algorithm by applying the technique of stochastic gradient descent ./ sequential gradient descent.

→ Error function is given for all data points as summation function.

$$\mathcal{E} = \sum_n E_n$$

→ The stochastic gradient descent algorithm updates the parameter w using

$$w^{(T+1)} = w^T - \eta \nabla E_n$$

$T \rightarrow$ Iteration number

$\eta \rightarrow$ Learning rate

→ we Initialize value of ' w ' to some starting vector $w^{(0)}$

→ For the case of Sum of square of error function $\mathcal{E}(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2$

$$w^{(T+1)} = w^T + \eta (t_n - w^{(T)} \phi_n) \phi_n$$

(26)

where $\phi_n = \phi(x_n)$. This is known as Least - mean - Squares or LMS algorithm.

Regularized Least Squares :-

→ Earlier concept of adding a regularization-term to an error function in order to control over-fitting to control total error function & minimize it.

$$E_D(w) + \lambda E_N(w)$$

\downarrow
regularization term

$E_D(w) \rightarrow$ data dependent Error

$E_N(w) \rightarrow$ regularization term

→ Simple regularizer is given by sum-of-squares of weight vector Elements.

$$E_N(w) = \frac{1}{2} w^T w$$

$$\text{also } E(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2$$

∴ Total error fn is

[Subst done for $E_D(w) + \lambda E_N(w)$]

$$\frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2 + \frac{\lambda}{2} w^T w \rightarrow (1)$$

$\underbrace{\sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2}_{E(w)} + \lambda \underbrace{w^T w}_{E_N(w)}$

→ This choice of regularizer is known as in machine learning as weight decay as in Sequential learning algorithms, it encourages weight-vectors to decay towards zero, unless supported by the data.

→ Setting gradient of Eqn (1) wrt to ' w ' to zero & solving for ' w ' we get

Eqn (1) is

$$\frac{1}{2} \sum_{n=1}^N \{ t_n - w^T \phi(x_n) \}^2 + \frac{\lambda}{2} w^T \cdot w$$

$$\therefore \frac{1}{2} \sum_{n=1}^N \{ t_n - w^T \phi(x_n) \}^2 + \frac{\lambda}{2} w^2$$

Differentiating above Expression wrt ' w '

$$= \frac{1}{2} \cdot 2 \cdot \{ t_n - w^T \phi \} \phi^T + \frac{\lambda}{2} \cdot 2w$$

$$= t\phi^T - w^T \phi \phi^T + \lambda \cdot w$$

$$= t\phi^T - w^T (\phi \phi^T - (-)\lambda I)$$

$$0 = t\phi^T - w^T (\phi \phi^T + \lambda I)$$

$$\therefore w(\phi^T \phi + \lambda I) = t\phi^T$$

or

$$\boxed{w = (\lambda I + \phi^T \phi)^{-1} \phi^T t}$$

↑ regularization factor

Simple Extension of Least Square
with regularization factor.

$$w_{ML} = (\phi^T \phi)^{-1} \phi^T t$$

Multiple outputs :-

→ If we wish to predict more target variables 't' denoted collectively by target vector 't', then

$$y(x, w) = w^T \phi(x)$$

↓

$k \rightarrow$ dimensional column vector

$w \rightarrow m \times k$ matrix

$\phi(x) \rightarrow M$ dimensional column vector

* Gaussian distribution

$$p(t|x, w, \beta) = N(t | w^T \phi(x), \beta^{-1} I)$$

$$\therefore \ln p(T|x, w, \beta) = \sum_{n=1}^N \ln N(t_n | w^T \phi(x_n), \beta^{-1} I)$$

$\left. \begin{matrix} \uparrow \\ \{t_1 \dots t_n\} \end{matrix} \right\}$
 $\left. \begin{matrix} \downarrow \\ \{x_1 \dots x_n\} \end{matrix} \right\}$

$$\rightarrow \text{III}^* = \frac{NK}{2} \ln \left(\frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \| t_n - w^T \phi(z_n) \|^2$$

$\left[\text{III}^* \rightarrow \frac{N}{2} \ln \beta - \frac{N}{2} \ln (2\pi) - \beta E_D(w) \right]$ derived

Earlier, the above expression can also be written as

$$\frac{N}{2} \ln \left(\frac{\beta}{2\pi} \right) - \beta E_D(w) \left[\log A - \log B = \log \frac{A}{B} \right]$$

\therefore maximizing above fm wrt w

$$w_{ML} = (\phi^T \phi)^{-1} \phi^T \cdot T \quad \xrightarrow{\text{Target vector}}$$

$$\uparrow \text{II}^* \rightarrow w_{ML} = (\phi^T \phi)^{-1} \phi^T \cdot t$$

\rightarrow If we examine result for target variable t_K

$$w_K = (\phi^T \phi)^{-1} \phi^T t_K = \phi^T t_K.$$

single target

The Bias-variance decomposition :-

\rightarrow Use of maximum likelihood or least squares can lead to severe over-fitting if complex data sets are taken.

→ We can limit the no. of basis function in order to avoid over-fitting but this will limit the flexibility of the model to capture important trends in data.

→ Here adding a regularization term λ [>] can control over-fitting for models with many parameters.

→ Bayesian treatment of linear regression will avoid the over-fitting problem of maximum likelihood.

→ Note:- Bias is used to simplify the assumptions made by the model to make the target function easier to approximate.

2) Variance is the amount that the estimate of the target will change given different training data.

3) Bias-variance tradeoff is the property of a model that the variance of the parameters estimated across samples by increasing the bias in the estimated parameters.

→ Bias - variance decomposition or tradeoff is the conflict in trying to simultaneously minimize these two sources of error that prevent supervised learning algorithms from generalizing beyond the training set.

Bias Error - erroneous assumptions in the learning algorithm.

Variance - Error from sensitivity to small fluctuations in the training set.

* Bias - Variance decomposition is a way of analyzing a learning algorithm's expected generalization error with respect to three terms bias, variance, irreducible error. resulting from noise itself.

→ Assuming $g(x)$ as generalization error which is taken for predicting a new target t for a new input x

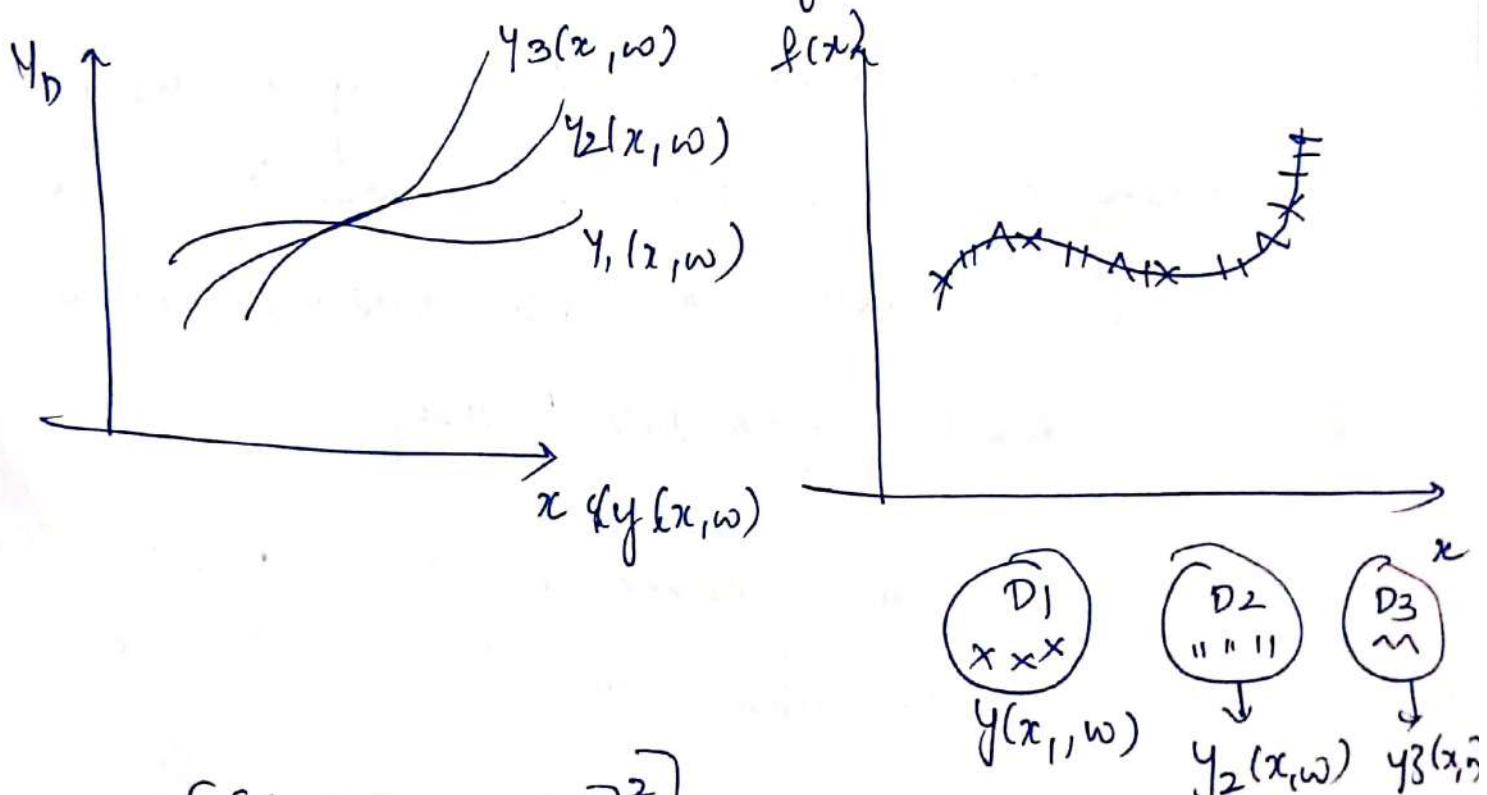
$$E [(t - g(x))^2]$$

→ Error in Estimation of t , can be written as

(32)

$$\begin{aligned} \mathbb{E} [(t - y(x))^2] &= \mathbb{E} [(t - g(x))^2] \\ &= \mathbb{E} \left[(\underbrace{t - g(x)}_a + \underbrace{g(x) - y(x)}_b)^2 \right] \\ &= \mathbb{E} \left[(\underbrace{t - g(x)}_a)^2 \right] + \mathbb{E} \left[(\underbrace{g(x) - y(x)}_b)^2 \right] \\ &\quad + 2 \cdot \mathbb{E} \left[(\underbrace{t - g(x)}_a) \underbrace{g(x) - y(x)}_b \right]. \end{aligned}$$

→ Here $y(x)$ is obtained from data set-



$$\rightarrow \mathbb{E} [(y(x) - g(x))^2] \xrightarrow{(1)}$$

we can write the above eqn (4) as

$$\rightarrow E_D [(g(x) - y(x, D))^2].$$

$$= \mathbb{E}_D \left[\underbrace{(y(x, D) - \mathbb{E}_D[y(x, D)])}_{P} + \underbrace{\mathbb{E}_D[y(x, D)] - g(x)}_{Q} \right]^2$$

$$= \mathbb{E}_D [P^2 + Q^2 + 2PQ]$$

$$= \mathbb{E}_D \left[(y(x, D) - \mathbb{E}_D[y(x, D)]) \right]^2 + \mathbb{E}_D \left[\mathbb{E}_D[y(x, D)] - g(x) \right]^2 \\ + 2 \mathbb{E}_D \left[(y(x, D) - \mathbb{E}_D[y(x, D)]) \left(\mathbb{E}_D[y(x, D)] - g(x) \right) \right]$$

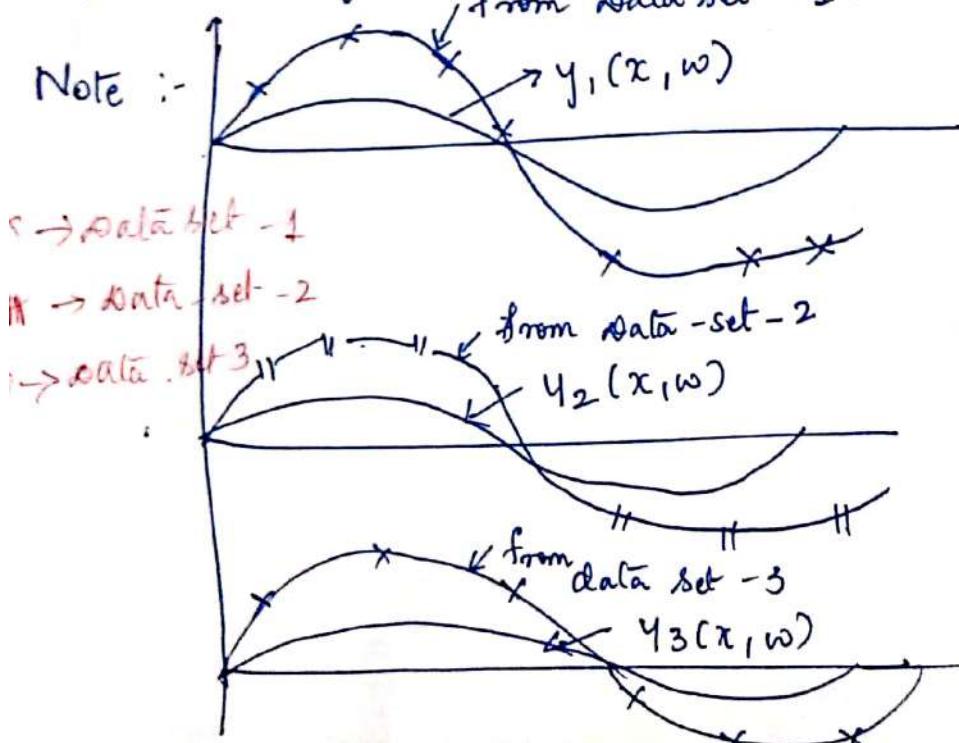
Variance Bias²

Noise.

$\mathcal{L}^{(\text{loss})} = \text{variance} + \text{Bias}^2 + \text{loss Noise}$

* $g(x)$ is from Expectation

$y(x)$ is from parametric approach.



* $g(x)$ is reference
 $y_1(x, w)$, $y_2(x, w)$
 are different solution

* how far $g(x)$ is
 deviated from $y_1(x, w)$
 $y_2(x, w)$ is known
 as Bias.

* Estimation within $y_1(x, w)$, $y_2(x, w)$, $y_3(x, w)$
or deviation within $y_1(x, w)$, $y_2(x, w)$, $y_3(x, w)$
is called as variance.

-: Module 2 :-

Introduction:-

- The goal in classification is to take an ^{Input} vector $x \in \mathbb{R}^D$ to assign it to one of K discrete classes C_k .
- The Input space is divided into decision regions whose boundaries are called as decision boundaries.
- As we are discussing linear classification, it means that the decision surfaces are linear functions of Input vector $x \in \mathbb{R}^D$ are defined by $(D-1)$ dimensional hyperplanes within the D -dimensional Input space.
- There are 3 distinct approaches to the classification problems & the simplest involves constructing a discriminant function which directly assigns each vector ' x ' to a specific class.

Discriminant functions:-

- A discriminant function takes an \mathbb{R}^p vector ' x ' and assigns it to one of K classes denoted as C_k .

→ Here we assume linear discriminations in which decision surfaces are hyperplanes.

Two classes :-

Simplest representation of a linear discriminant function is

$$y(x) = w^T x + w_0$$

w → weight vector

w₀ → Bias

→ If vector x is assigned to class C₁ if y(x) ≥ 0
C₂ if ~~y(x)~~ otherwise.

→ The decision boundary is corresponding to a (D-1) dimensional hyperplane within the D-dimensional input space.

→ Consider two points x_A & x_B both lie on the decision surfaces. Because y(x_A) = y(x_B) = 0 we have $w^T(x_A - x_B) = 0$ & hence 'w' vector is orthogonal to every vector lying on the decision surface, ∴ 'w' determines the orientation of the decision surface.

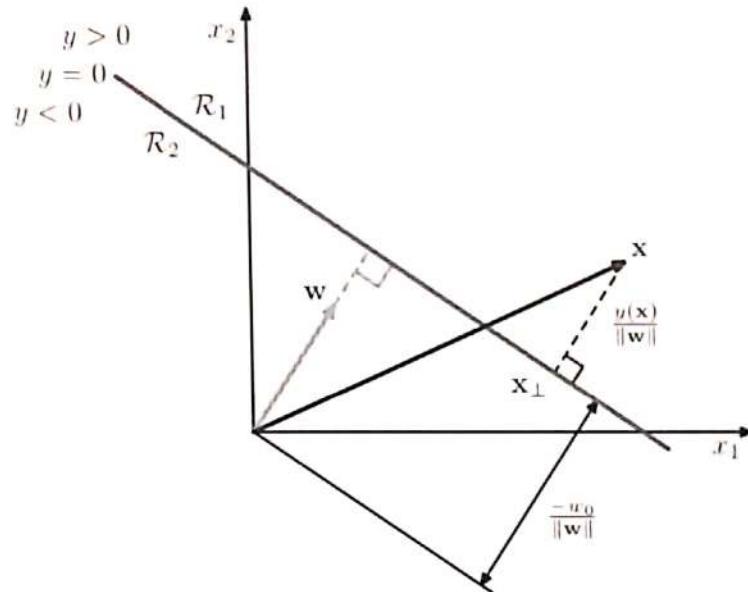
→ The normal distance from the Origin to
the decision surface is given by ③

$$\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$$

→ Consider the diagram shown below.

182 4. LINEAR MODELS FOR CLASSIFICATION

Figure 4.1 Illustration of the geometry of a linear discriminant function in two dimensions. The decision surface, shown in red, is perpendicular to \mathbf{w} , and its displacement from the origin is controlled by the bias parameter w_0 . Also, the signed orthogonal distance of a general point \mathbf{x} from the decision surface is given by $y(\mathbf{x})/\|\mathbf{w}\|$.



→ Here $y(\mathbf{x})$ gives a signed measure of the perpendicular distance r of the point \mathbf{x} from the decision surface.

→ For Example let \mathbf{x} be the arbitrary point if \mathbf{x}_\perp be its orthogonal projection onto the decision surface.

$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

→ multiplying both sides with w^T & adding
 w_0 we get - (4)

$$x \cdot w^T + w_0 = \left[x \perp + r \frac{w}{\|w\|} \right] w^T + w_0$$

* using $y(x) = w^T x + w_0$ ~~zero~~

$$y(x \perp) = w^T x \perp + w_0 = 0$$

we get

$$y(x) = r \boxed{r = \frac{y(x)}{\|w\|}}$$

$$\boxed{y(x) = \tilde{w}^T \tilde{x}}$$

ii) Multiple classes :-

→ Consider Extension of linear discriminations
to $k > 2$ classes.

→ Consider the use of $k-1$ classifiers, each
of which solves a two class problem of
separating points into a particular class
 C_k from points not in that class.

(5)

→ This is known as a one-versus-the-rest classifier.

→ Fig below shows an example of 3 classes with ambiguous classification.

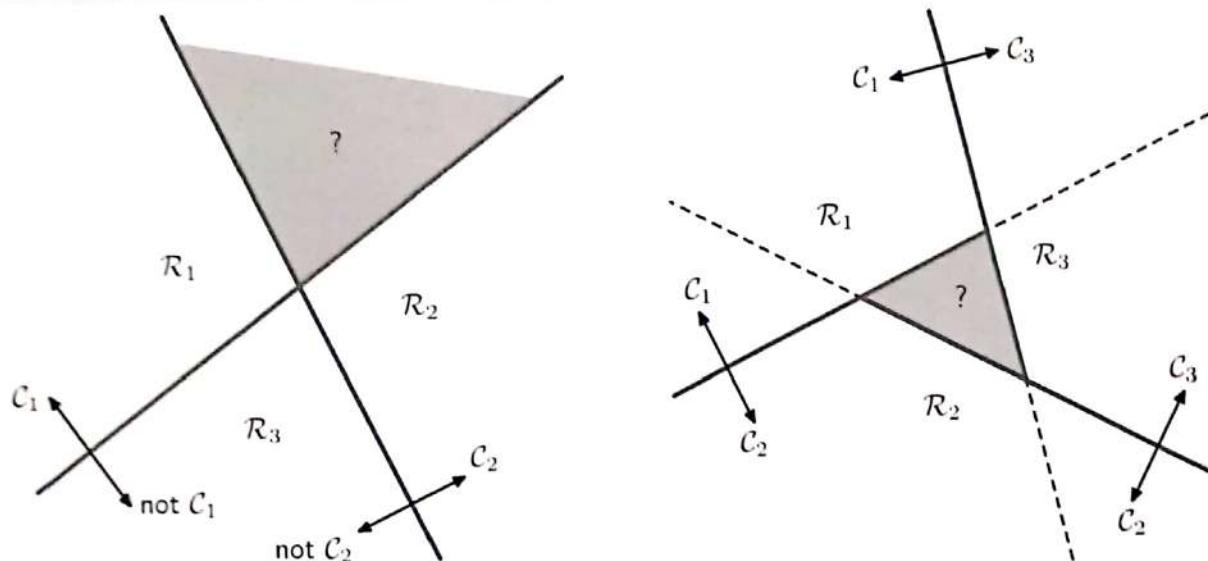


Figure 4.2 Attempting to construct a K class discriminant from a set of two class discriminants leads to ambiguous regions, shown in green. On the left is an example involving the use of two discriminants designed to distinguish points in class C_k from points not in class C_k . On the right is an example involving three discriminant functions each of which is used to separate a pair of classes C_k and C_j .

→ An alternate is to introduce $\frac{K(K-1)}{2}$ binary discriminant functions, one for every possible pair of classes. This is known as a one-versus-one classifier.

→ Further illustrating in Equations.

$$y_K(x) = w_K^T x + w_K^0$$

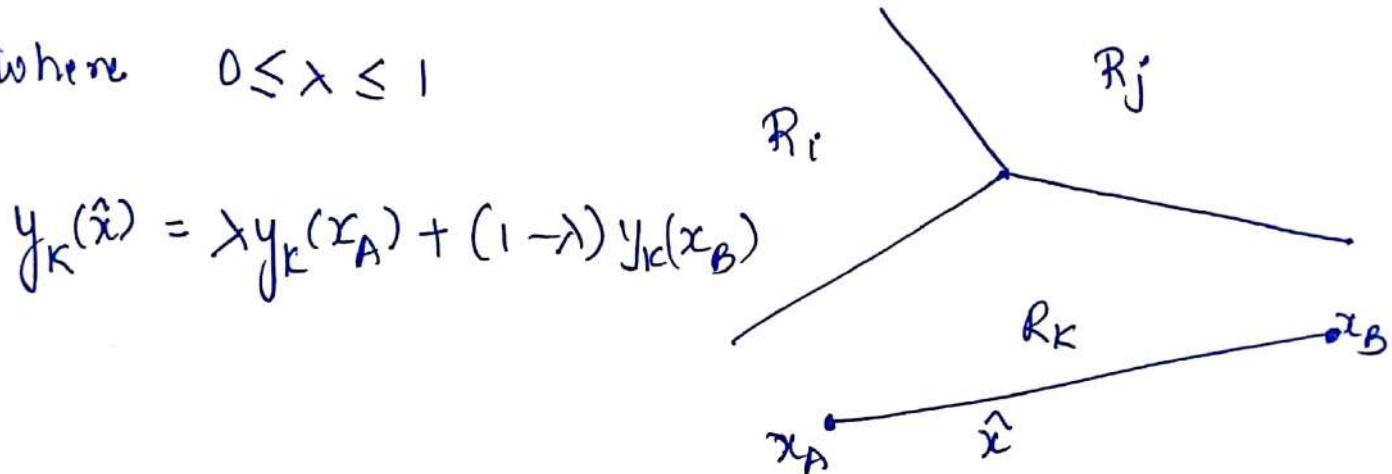
$$\{ (w_K - w_j)^T x + (w_{K0} - w_{j0}) = 0$$

where k & j are two classes C_k & C_j . (6)

→ consider two points $x_A \in R_k$ both of which lie inside decision region R_k . Any point \hat{x} that lies on the line connecting $x_A \in R_k$ can be expressed in the form,

$$\hat{x} = \lambda x_A + (1-\lambda)x_B.$$

where $0 \leq \lambda \leq 1$



Note :- Both $x_A \in R_k$ & $x_B \in R_k$. R_k is singly connected & convex.

Least Squares for classification :-

→ we know that minimization of sum-of-squares error function led to a simple closed-form solution for the parameter values. in regression, so that can be applied for classification.

→ Each class C_k described by its linear model

$$y_k(x) = \omega_k^T x + \omega_{k0}$$

where $k = 1, \dots, K$, grouping these together using vector notation.

$$y(x) = \tilde{\omega}^T \tilde{x}$$

→ we now determine parameter matrix $\tilde{\omega}$ by minimizing sum-of-squares error function.

Error function is defined as

$$E_D(\tilde{\omega}) = \frac{1}{2} \text{Tr} \{ (\tilde{x}\tilde{\omega} - T)^T (\tilde{x}\tilde{\omega} - T) \}$$

* Setting the derivative with respect to $\tilde{\omega}$ to zero. & rearranging we obtain

$$\tilde{\omega} = (\tilde{x}^T \tilde{x})^{-1} \tilde{x}^T T = \tilde{x}^+ T$$

\downarrow
Pseudo-Inverse of
the matrix \tilde{x} .

Note :- Sum-of-squares error function penalizes predictions that are 'too correct' & which lie along the decision boundary.

Principle Component Analysis :-

- It is a way of Identifying pattern in data and Expressing the data in such a way to highlight their similarities & differences
- It is also way to reduce the dimension [Dimensionality Reduction] to reduce the complexity of the Analysis .

Example problem :-

- Given the following data use PCA to reduce the dimension from 2 to 1 .

Feature	1	2	3	4
x	4	8	13	7
y	11	4	5	14

Step 1 :- From the data set find n & N

$$n = \text{no. of features} = 2$$

$$N = \text{no. of samples} = 4$$

Step 2 :- Computation of mean of variables

$$\bar{x} = \frac{4+8+13+7}{4} = 8$$

$$\bar{y} = \frac{11 + 4 + 5 + 14}{4} = 8.5$$

Step 3 :- Computation of Covariance matrix

i> find ordered pairs $\rightarrow x, y$

(x, x) (x, y) (y, x) (y, y)

ii) find covariance of all ordered pairs

$$\text{Cov}(z, x) = \frac{1}{N-1} \sum_{k=1}^N (x_{ik} - \bar{x}_i)(z_{jk} - \bar{z}_j)$$

07

$$\text{a) } \text{Ges}(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^2$$

$$= \frac{1}{d_4 - 1} \left[(4 - 8)^2 + (8 - 8)^2 + (13 - 8)^2 + (7 - 8)^2 \right]$$

三 14

$$4) \text{cov}(x,y) = \frac{1}{4-1} \left[(4-8)(11-8.5) + (8-8)(4-8.5) + (18-8)(5-8.5) + (7-8)(14-8.5) \right] \\ = -11$$

$$\text{cov}(y, x) = \text{cov}(x, y) = -11$$

$$\begin{aligned}\text{cov}(y, y) &= \frac{1}{4-1} \left[(11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2 \right] \\ &= 23\end{aligned}$$

Covariance matrix :-

$$\begin{aligned}S &= \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix} \\ &= \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}\end{aligned}$$

Step 4 :- Eigen value , Eigen vector , normalized eigen vector .

i) Eigen value :-

$$\det(S - \lambda I) = 0$$

$$\text{Identity matrix } I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \times \lambda = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$S = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

$$\therefore \det(S - \lambda I) = 0$$

$$\det \begin{bmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{bmatrix} = 0$$

$$(14-\lambda)(23-\lambda) - (-11 \times -11) = 0$$

$$\lambda^2 - 37\lambda + 201 = 0$$

$$\therefore \lambda_1 = 30.3849 \quad \left[\lambda_1 > \lambda_2 \right]$$

$$\lambda_2 = 6.6151$$

* pick largest value which becomes first principal component.

ii) Eigen vector of λ_1

$$(S - \lambda_1 I) U_1 = 0$$

\downarrow \hookrightarrow Eigen vector of λ_1

$$\begin{bmatrix} 14-\lambda_1 & -11 \\ -11 & 23-\lambda_1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 0$$

$$\begin{bmatrix} (14-\lambda_1)u_1 - 11u_2 \\ -11u_1 + (23-\lambda_1)u_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(14-\lambda_1)u_1 - 11u_2 = 0 \rightarrow (1)$$

$$-11u_1 + (23-\lambda_1)u_2 = 0 \rightarrow (2)$$

* use any one linear equation to find values of eigen vector

∴ Using Eqn (1)

$$(14 - \lambda_1) u_1 - 11 u_2 = 0$$

$$\frac{u_1}{11} = \frac{u_2}{14 - \lambda_1} = t$$

$$\text{when } t = 1$$

$$u_1 = 11, \quad u_2 = 14 - \lambda_1$$

$$\text{Eigen vector } u_1 \text{ of } \lambda_1 = \begin{bmatrix} 11 \\ 14 - \lambda_1 \end{bmatrix}$$

$$= \begin{bmatrix} 11 \\ 14 - 30.3849 \end{bmatrix} = \begin{bmatrix} 11 \\ -16.3849 \end{bmatrix}$$

* Normalizing Eigen vector u_1

$$e_1 = \begin{bmatrix} \frac{11}{\sqrt{11^2 + (-16.3849)^2}} \\ \frac{-16.3849}{\sqrt{11^2 + (-16.3849)^2}} \end{bmatrix} = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$$

$$\text{Hence } e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$



Step 5 :- Define new dataset-

Frist principle component	1	2	3	4
	P_{11}	P_{12}	P_{13}	P_{14}

$$P_{11} = e_1^T \begin{bmatrix} x_i - \bar{x} \\ y_i - \bar{y} \end{bmatrix} \quad e_1 = e_1^T \begin{bmatrix} 4 - 8 \\ 11 - 8.5 \end{bmatrix}$$

$$= \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \begin{bmatrix} -4 \\ 2.5 \end{bmatrix} = -4.3052$$

$$P_{12} = \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \begin{bmatrix} 8 & -8 \\ 4 & -8.5 \end{bmatrix} = 3.73$$

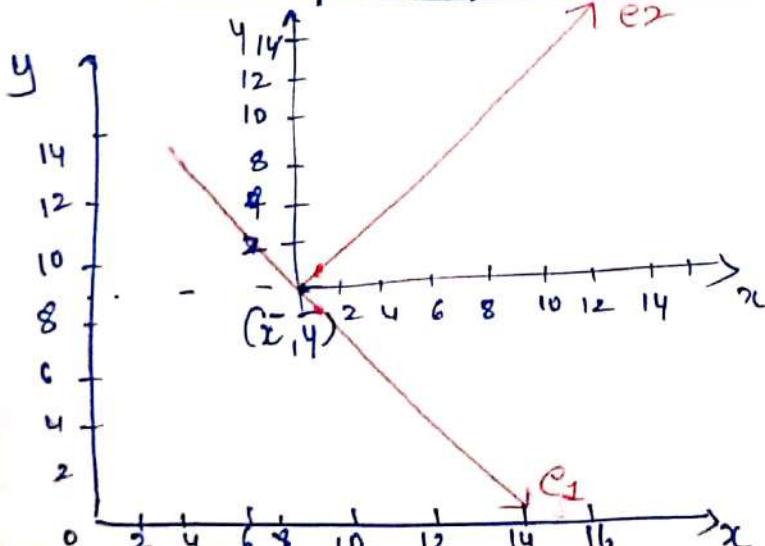
111th

$$P_{13} = 5.69$$

$$P_{14} = -5.12$$

PC1	1	2	3	4
	-4.3	3.7	5.6	-5.12

⇒ new data set
with dimension 1



Linear discriminant analysis :-

(14)

- LDA is a dimensionality reduction technique used as a pre-processing step for pattern-classification & Machine learning applications.
- LDA is similar to PCA but LDA in addition finds the axes that maximizes the separation between multiple classes.

Goal :- To project a feature space (n -dimensional data) onto a smaller subspace (K ($K \leq n-1$) while maintaining the class discriminating information

Example :- Let's take a 2-D dataset -

$$C_1 \Rightarrow x_1 = (x_1, x_2) = \{(4,1), (2,4), (2,3), (3,6), (4,4)\}$$

$$C_2 \Rightarrow x_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$$

Step 1 :- Compute within-class scatter matrix (S_W)

$$S_W = S_1 + S_2$$

$S_1 \Rightarrow$ Covariance matrix for class C_1

$S_2 \Rightarrow$ Covariance matrix for class C_2 .

i) finding S_2

$$S_2 = \sum_{x \in C_2} (x - \mu_2)(x - \mu_2)^T$$

$\mu_2 \rightarrow$ mean of class C_2

$$\mu_2 = \left\{ \frac{4+2+2+3+4}{5}, \frac{1+4+3+6+4}{5} \right\}$$

$$\mu_2 = [3.00 \quad 3.60]$$

Similarly $\mu_2 = [8.4 \quad 7.60]$

$$\therefore (x - \mu_1) = \begin{bmatrix} 1 & -1 & -1 & 0 & 1 \\ -2.6 & 0.4 & -0.6 & 2.4 & 0.4 \end{bmatrix}$$

Now, for each x we calculate $(x - \mu_1)(x - \mu_1)^T$.

$$\therefore \begin{bmatrix} 1 \\ -2.6 \end{bmatrix} \begin{bmatrix} 1 & -2.6 \end{bmatrix} = \begin{bmatrix} 1 & -2.6 \\ -2.6 & 6.76 \end{bmatrix} \rightarrow (1)$$

$$\begin{bmatrix} -1 \\ 0.4 \end{bmatrix} \begin{bmatrix} -1 & 0.4 \end{bmatrix} = \begin{bmatrix} 1 & -0.4 \\ -0.4 & 0.16 \end{bmatrix} \rightarrow (2)$$

$$\begin{bmatrix} -1 \\ -0.6 \end{bmatrix} \begin{bmatrix} -1 & -0.6 \end{bmatrix} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 0.36 \end{bmatrix} \rightarrow (3)$$

$$\begin{bmatrix} 0 \\ 2.4 \end{bmatrix} \begin{bmatrix} 0 & 2.4 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 5.76 \end{bmatrix} \rightarrow (4)$$

$$\begin{bmatrix} 1 \\ 0.4 \end{bmatrix} \begin{bmatrix} 1 & 0.4 \end{bmatrix} = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 0.16 \end{bmatrix} \rightarrow (5)$$

(14)

Adding (1) + (2) + (3) + (4) + (5) & taking average
we get covariance matrix S_1

$$S_1 = \begin{bmatrix} \frac{4}{5} & -\frac{2}{5} \\ -\frac{2}{5} & \frac{13.2}{5} \end{bmatrix} = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 2.64 \end{bmatrix}$$

Similarly for the class 2, the covariance matrix is given by

$$S_2 = \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix} \text{ & } \mu_2 = \begin{bmatrix} 8.4 & 7.6 \end{bmatrix}$$

$$S_W = S_1 + S_2$$

$$S_W = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

Step 2 :- Compute Between class scatter matrix (S_B)

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$\mu_1 = \begin{pmatrix} 3 & 3.6 \end{pmatrix} \quad \mu_2 = \begin{pmatrix} 8.4 & 7.6 \end{pmatrix}$$

$$S_B = \begin{bmatrix} -5.4 \\ -4 \end{bmatrix} \begin{bmatrix} -5.4 & -4 \end{bmatrix} = \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16.00 \end{bmatrix}$$

Step 3:- Find the best LDA projection vector

* Similar to principal component analysis we find this using eigen vectors using largest Eigen value.

$$S_w^{-1} S_B \cdot w = \lambda \cdot w \quad \xrightarrow{\text{Projection vector}} \quad (1)$$

i) Solving the eigen value

$$|S_w^{-1} S_B - \lambda I| = 0$$

$$S_w^{-1} = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}^{-1} = \begin{bmatrix} 0.38 & 0.03 \\ 0.03 & 0.19 \end{bmatrix}$$

$$S_B = \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16.00 \end{bmatrix}$$

$$\lambda I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \times \lambda = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$S_w^{-1} \cdot S_B = \begin{bmatrix} 11.72 & 8.6 \\ 4.9 & 3.6 \end{bmatrix}$$

$$\therefore |S_w^{-1} S_B - \lambda I| = \begin{vmatrix} 11.72 - \lambda & 8.6 \\ 4.9 & 3.6 - \lambda \end{vmatrix} = 0$$

Solving for λ

$$(11.72 - \lambda)(3.6 - \lambda) - (4.9 \times 8.6) = 0$$

$$42.19 - 11.72\lambda - 3.6\lambda + \lambda^2 - 42.14 = 0$$

$$\lambda^2 - 15.32\lambda + 0.05 = 0.$$

$$\therefore \lambda_1 = 15.31, \lambda_2 = 3.26 \times 10^{-3}$$

Substituting highest Eigen value in Eqn (1)

$$\begin{bmatrix} 11.72 & 8.6 \\ 4.9 & 3.6 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 15.31 \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

~~$$11.72w_1 + 8.6w_2 = 15.31w_1 \rightarrow (2)$$~~

$$4.9w_1 + 3.6w_2 = 15.31w_2 \rightarrow (3)$$

$$\begin{cases} 11.72w_1 - 15.31w_1 + 8.6w_2 = 0 \\ -3.59w_1 + 8.6w_2 = 0 \end{cases} \rightarrow (4)$$

$$4.9w_1 + 3.6w_2 - 15.31w_2 = 0$$

$$4.9w_1 - 11.7w_2 = 0$$

$$\begin{bmatrix} -3.59 & 8.6 \\ 4.9 & -11.7 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$-3.59w_1 + 8.6w_2 = 0$$

$$w_1 = \frac{8.6}{3.59}w_2 = 2.39w_2$$

$$4.9w_1 + 11.7w_2 = 0$$

$$\text{Let } w_2 = 1$$

$$\therefore w_1 = \frac{11.7}{4.9} = 2.38$$

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 2.38 \\ 1 \end{bmatrix}$$

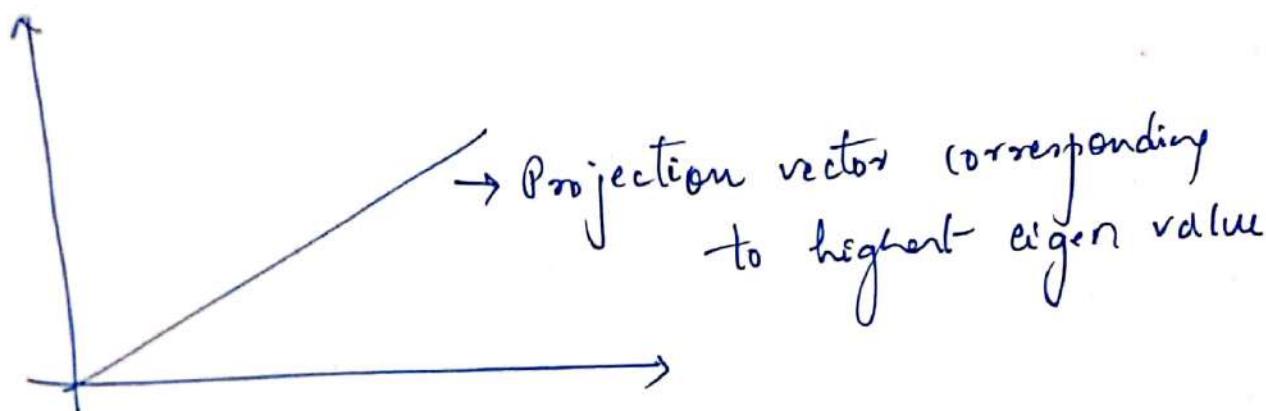
normalizing

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 2.38 \\ \sqrt{1^2 + 2.38^2} \end{bmatrix}$$

$$\therefore \text{we get } \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0.91 \\ 0.39 \end{bmatrix}$$

Step 4 :- Dimension Reduction

$$y = w^T x \rightarrow \begin{array}{l} \text{Input data sample} \\ \downarrow \\ \text{Projection vector} \end{array}$$



The perceptron algorithm :-

→ Here let us consider two-class model in which input vector x is first transformed using a fixed nonlinear transformation to give a feature vector $\phi(x)$ & then construct a generalized linear model

$$y(x) = f(w^T \phi(x)).$$

where the non-linear activation function $f(\cdot)$ is a step function of form

$$f(a) = \begin{cases} +1 & a > 0 \\ -1 & a < 0 \end{cases}$$

- In Perception algorithm it convenient to use $t = +1$ for c_1 & $t = -1$ for c_2 as target values matching the choice of activation function.
- To determine the parameters w correctly we can minimize error function but alternatively we can also use perception criterion.

- To derive perception criterion based Error we 1st consider

pattern $x_n \rightarrow c_1 \Rightarrow$ when $w^T \phi(x_n) > 0$

$x_n \rightarrow c_2 \Rightarrow$ when $w^T \phi(x_n) < 0$

$$t \in \{-1, +1\}$$

note :- v satisfy condition $w^T \phi(x_n) t_n > 0$
All patterns must

→ perception criterion is given by

$$E_p(w) = - \sum_{n \in M} w^T \phi_n t_n$$

$M \rightarrow$ denotes the set of misclassified patterns

→ we apply gradient-descent algorithm to this error function. The change in weight is given by

$$\begin{cases} w^{(t+1)} = w^{(t)} - \eta \nabla E_p(w) \\ = w^{(t)} + \eta \phi_n t_n. \end{cases}$$

where $\eta \rightarrow$ learning rate

$T \rightarrow$ integer values of steps of algorithm.

Note :- As the weight are change during training the misclassified pattern will change.
set of

Probabilistic Discriminative Models :-

→ There are 2 Models

- i) Generative
- ii) discriminative.

Analogy :-

Task \rightarrow Determine the language that someone is speaking

- i) Generative approach :- Learn Each language & determine as to which language the speech belongs to.
- ii) Discriminative approach :- Determine the linguistic differences without learning any language

- \Rightarrow We wish to learn $f: x \rightarrow y$ eg., $P(y|x)$
- i) In generative we model the joint distribution of all data
 - ii) In discriminative we model only points at the boundary.

Note :- Discriminative models make predictions on the unseen data based on conditional probability

\rightarrow Discriminative classification is also called as "informative".

Probability Basics :-

(23)

* Prior, conditional & joint probability for random variables.

1) Prior Probability - $P(x)$

2) Conditional probability - $P(x_1|x_2), P(x_2|x_1)$

3) Joint probability :- $x = (x_1, x_2) \quad P(x) = P(x_1, x_2)$

4) Relationship $P(x_1, x_2) = P(x_2|x_1) P(x_1)$
 $= P(x_1|x_2) P(x_2)$

5) Independence

$$P(x_2|x_1) = P(x_2), P(x_1|x_2) = P(x_1), P(x_1, x_2) = P(x_1)P(x_2)$$

Bayes rule :-

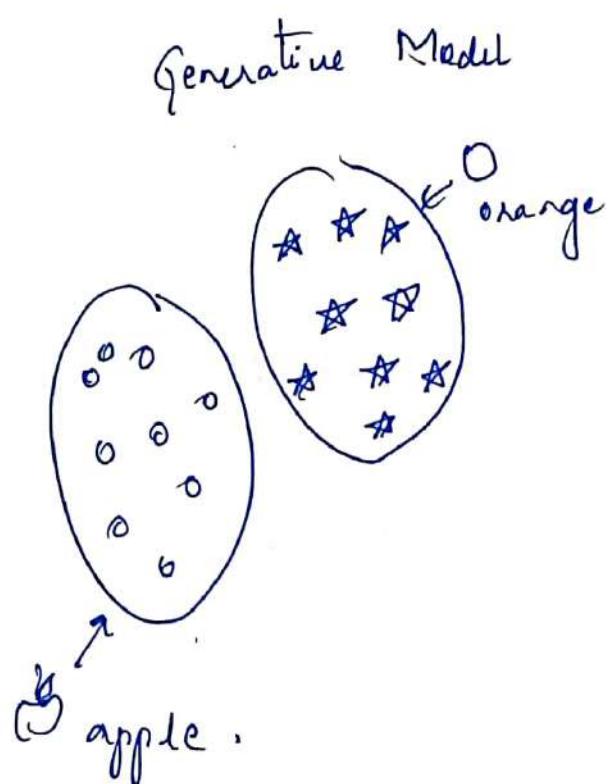
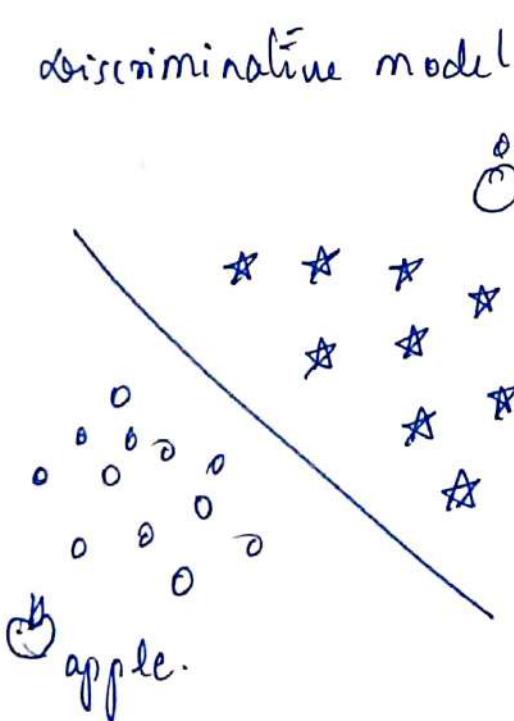
Generative .

$$P(c|x) = \frac{P(x|c) P(c)}{P(x)}$$

discriminative

$$\text{Posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{Evidence.}}$$

- Given some data points, the discriminative model (24) learns to classify the data points into their respective classes by learning the decision boundary that separates the classes.
- The generative models can also classify the given data points, but instead of learning the decision boundary, they learn the characteristics of the classes.



- * Above shows example of image classification & whether an input image is an apple/orange.
- i) Discriminative Model learns the optimal

boundary that separates the apple & orange classes

ii) Generative model learn their distribution by learning the characteristics of apple & orange classes.

→ while both types of models are used to predict $P(y|x)$. Calculation methods for both Models are different.

i) Generative Model → first find joint probability distribution $P(x,y)$, then use Bayes' theorem to calculate $P(y|x)$

ii) discriminative Model → will learn the conditional probability $P(y|x)$ directly.

Ex :- Data set $(0,0)$ $(0,1)$ $(1,0)$ $(1,1)$

$$\begin{array}{ccc}
 P(x,y) & & \\
 \begin{array}{c} x=0 \\ x=1 \end{array} &
 \begin{array}{c} y=0 \\ \hline y_1 \\ y_2 \end{array} &
 \begin{array}{c} y=1 \\ \hline y_3 \\ y_4 \end{array} \\
 & &
 \begin{array}{c} y_1 \Rightarrow P(x_1) \\ y_2 \Rightarrow P(x_2) \end{array} \\
 & &
 \hline
 p(y) = 4/8 & & p(y) = 1/8
 \end{array}$$

(25)

$$\text{and } P(y|x) = \frac{P(x,y)}{P(x)} \text{ or } \frac{P(x,y)}{P(y)}$$

$$P(x_1) = \frac{1}{2} \quad P(x_2) = \frac{1}{2}$$

$$\therefore P(y|x) \quad \begin{array}{c} y=0 \\ x=0 \end{array} \quad \begin{array}{c} y=1 \\ x=1 \end{array} \quad \frac{\text{prior})}{\text{post}} = \frac{y_2}{y_1}$$

$$x=1 \quad \frac{1}{2} \quad \frac{1}{2}$$

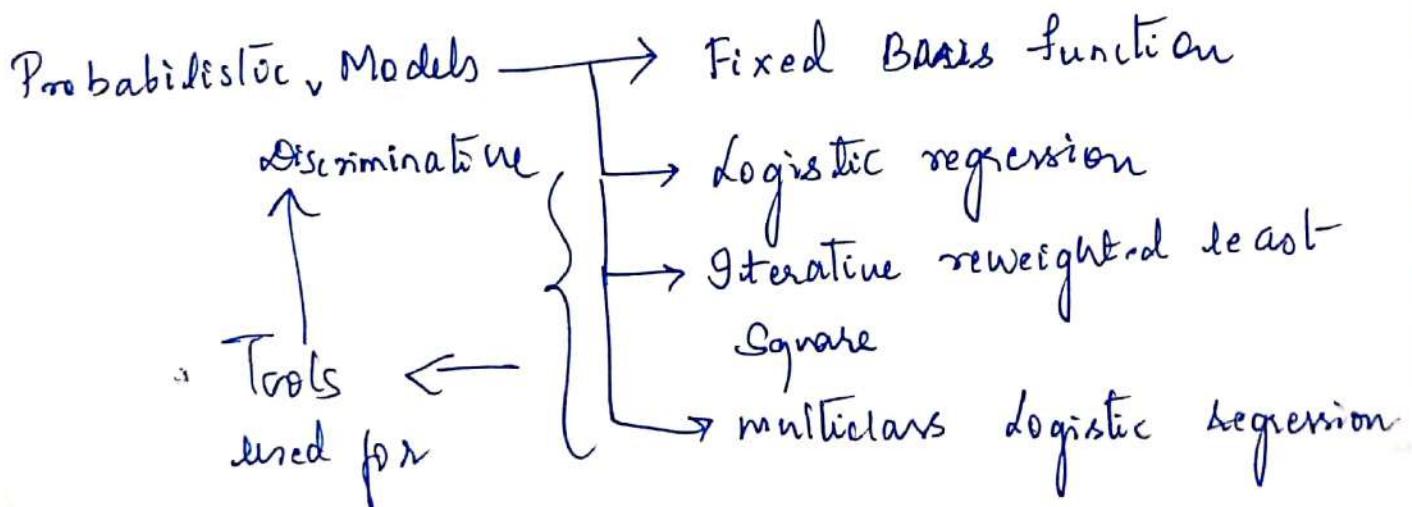
Summary:-

Generative Model

→ works good for small data & delivers high accuracy.

Discriminative Model

- Accuracy is high
- Resource saving
- fewer parameters to be determined.



Logistic Regression :-

(27)

→ Here we model the posterior probabilities directly assuming that they have a sigmoid-shaped distribution [without Modeling class prior & class conditional densities]

→ The sigmoid-shaped function (σ) is a model function logistic regression.

→ 1st non-linear transformation of inputs using a vector of basis function $\phi(x)$

$$P(c_1|\phi) = y(\phi) = \sigma(w^T\phi)$$

$$P(c_2|\phi) = \cancel{y(\phi)} \quad 1 - P(c_1|\phi)$$

∴ if we have M -dimensional feature space ϕ , we use $M+1$ adjustable parameters to represent w directly.

→ Here $\sigma(\cdot) \rightarrow$ logistic sigmoid function or logistic regression.

→ We can use maximum likelihood to determine the parameters of the logistic regression model.

$$\frac{d\sigma}{da} = \sigma(1-\sigma)$$

→ For a data set $\{\phi_n, t_n\}$, $t_n \in \{0, 1\}$ & $\phi_n = \phi(x_n)$ with $n = 1 \dots N$.

Likelihood function can be written as

$$P(t|w) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

where $t = [t_1, \dots, t_N]^T$

$$y_n = P(c_1 | \phi_n) \quad \& \text{Error Equation}$$

$$\nabla E(w) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

Iterative reweighted least squares :-

→ Logistic regression has no closed-form solution of $\nabla E(w) = 0$. due to non-linearity of logistic sigmoid function.

→ ∴ Error function can be minimized by an efficient iterative technique.

$$\text{ie., } w^{(\text{new})} = w^{(\text{old})} - H^{-1} \nabla E(w)$$

(20) 9

where $H \rightarrow$ Hessian matrix $\rightarrow 2^{\text{nd}}$ derivatives of $E(w)$
wrt 'w'

$$H = \nabla^2 E(w)$$

$$\nabla E(w) = \phi^T \phi w - \phi^T t.$$

$$\therefore H = \nabla \nabla E(w) = \sum_{n=1}^N \phi_n \phi_n^T = \underbrace{\phi^T \phi}_{\text{Diagonal matrix}}$$

Module 3 :-

K-Means algorithm :-

Clustering :-

→ clustering is the classification of objects into different groups or partitioning of a data set into subsets or clusters. Each subset of data share some common trait

Types of clustering are

- 1) Hierarchical algorithms
- 2) Partitional clustering

* K means algorithm comes under partitional clustering type of clustering.

Common distance Measures :-

→ Distance measure will determine how the similarity of two elements is calculated & will influence the shape of the clusters.

Euclidean distance (2-norm distance) is given by

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

numbers mean "stop"

→ k means clustering is an algorithm to cluster n objects based on attributes into k partitions. where $k < n$.

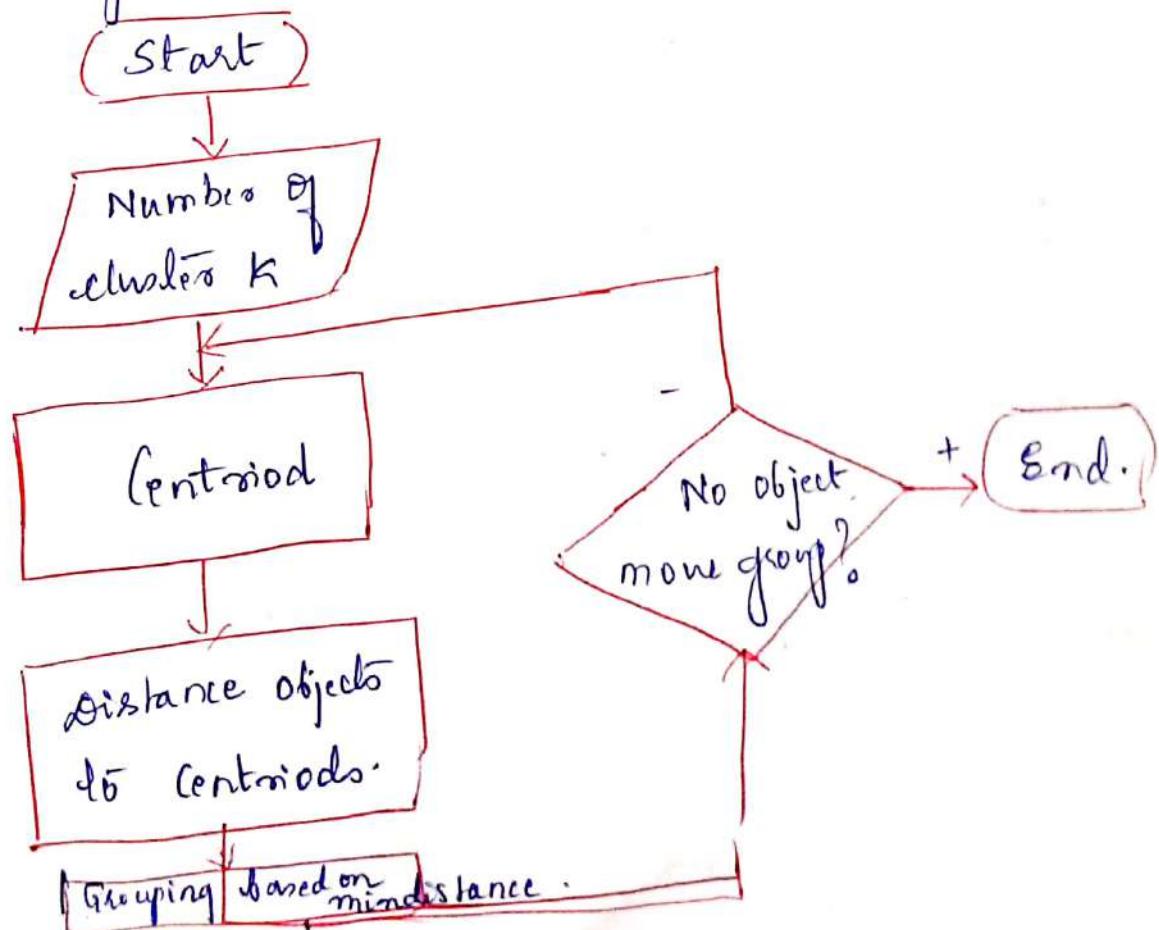
→ It assumes that the object attributes form a vector space.

note :- Supervised algorithms → Target variables are known

Unsupervised Algorithms → set of features are known, target or label are not known

clustering → find clusters of data

How the Algorithm works :-



(3)

Step 1 :- For each data point, place it in a cluster whose current centroid to which it is nearest.

Step 2 :- After all data points are assigned, repeat update the locations of centroid of the k clusters.

Step 3 :- Reassign all points to their closest centroid

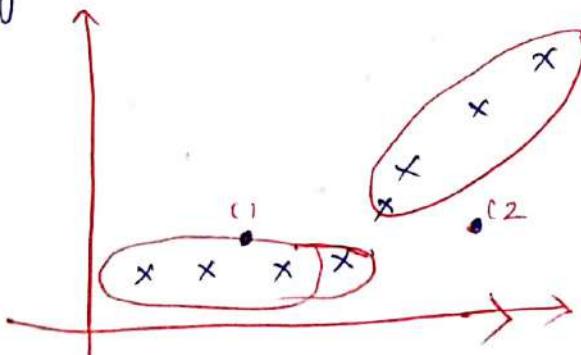
Note :- Sometimes data points changes its cluster centroids.

Step 4 :- Repeat step 2 & Step 3 until no data points change clusters & until all converge.

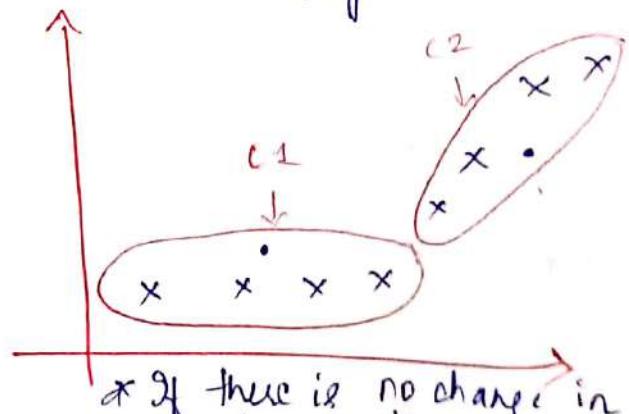
Con note :- Convergence - Points do not move between clusters & centroids are stable.

Example:-

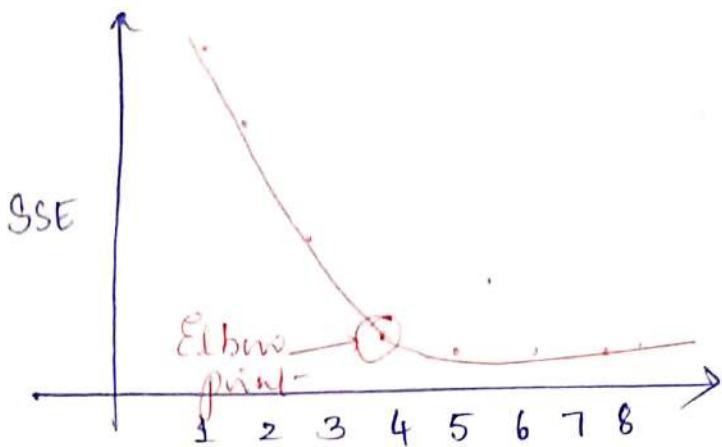
Step 1



Step 2



→ How to determine correct number of clusters (K) (4)



* Elbow Technique.

So value of $\boxed{K = 4}$ for a set of data points.

Note :- SSE [sum of Square Error]

$$SSE = \sum_{i=0}^n \text{dist}(x_i - \underbrace{c_i}_{\downarrow \text{cluster number}})^2$$

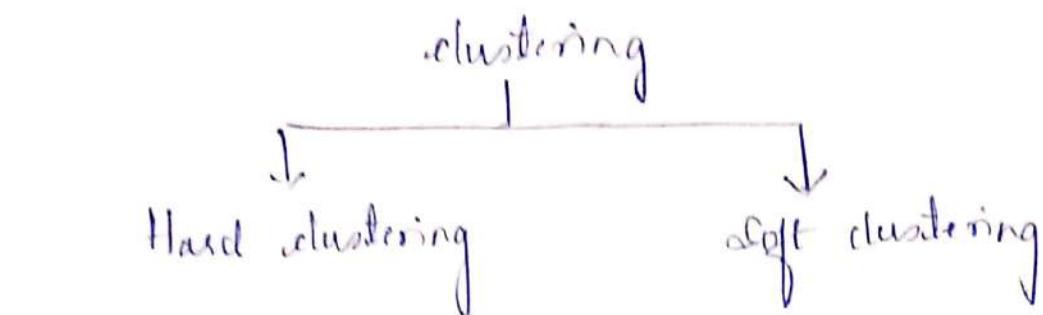
$$\therefore SSE = SSE_1 + SSE_2 + \dots + SSE_K$$

\uparrow for cluster 1 \downarrow cluster 2 \downarrow for cluster K.

Conclusion:-

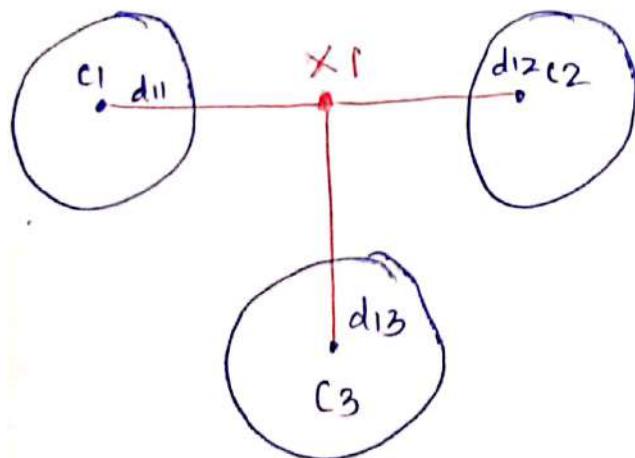
→ K-means algorithm is useful for undirected knowledge discovery & is relatively simple.

Fuzzy K-means Algorithm :-



Note :- K-means is a hard-clustering approach.

- In soft or fuzzy clustering, instead of putting each data points into separate clusters, a probability of that point to be in that cluster assigned.
- In soft clustering or fuzzy clustering, each data point can belong to multiple clusters along with its probability score or likelihood.
- Here Fuzzy set theory is used to obtain the optimal values of the centroid.



- In this technique the particular vector (x_1) belongs to all the clusters with different membership value.

(6)

note :- sum of membership value is '1'.

→ Consider vector x_1 belongs to the cluster 1 whose centroid is c_1 then $(x_1 - c_1)^2$ must be minimized. But if x_1 belongs to cluster 1 with membership value m_1 then $m_1(x_1 - c_1)^2$ has to be minimized.

→ The fuzzy K-means problem is treated as the optimization technique for obtaining the membership values along with the centroids of individual clusters such that

$$\sum_{i=1}^N \sum_{k=1}^K m_{ik}^2 (x_i - c_k)^2 \text{ is minimized}$$

$m_{ik} \rightarrow$ Membership value of vector x_i belongs to cluster k

$x_i \rightarrow i^{\text{th}}$ vector

$c_k \rightarrow c_k$ is k^{th} centroid.

note :- Here we use converging the membership values to find an optimal solution.

Algorithm :-

(7)

Step 1 :- Initialize the membership values (M_{ij}) randomly.

Step 2 :- Compute the centroids of the clusters

$$C_1 = \underbrace{\left[\text{vector } j * (\text{membership values of vector } j \text{ belonging to cluster 1})^2 \right]^2}_{\text{e.g., } M_{11}^2 + M_{12}^2 + \dots},$$

Sum of the squared values of membership values belonging to the cluster 1.

III) all centroids (e.g.: - C_2, C_3, \dots, C_6) are obtained.

Step 3 :- Update the membership values M using new centroids.

$$M_{11} = \frac{1}{\sum_{i=1}^6 \left[\frac{(x_1 - c_i)^2}{(x_1 - c_1)^2} \right]^{1/2}} = \frac{1}{\sqrt{\left[\frac{(x_1 - c_1)^2}{(x_1 - c_1)^2} \right] + \left[\frac{(x_1 - c_2)^2}{(x_1 - c_2)^2} \right] + \dots + \left[\frac{(x_1 - c_6)^2}{(x_1 - c_6)^2} \right]}}$$

Step 4 :- Compute the sum of squared difference between the previous membership value & current membership value.

$$[m_{11}^{\text{old}} - m_{11}^{\text{new}}]^2 + [m_{12}^{\text{old}} - M_{12}^{\text{new}}]^2 + \dots$$

- * If computed value is not less than the threshold then go to step 2.
- * If computed value is less than the threshold then stop the iteration.

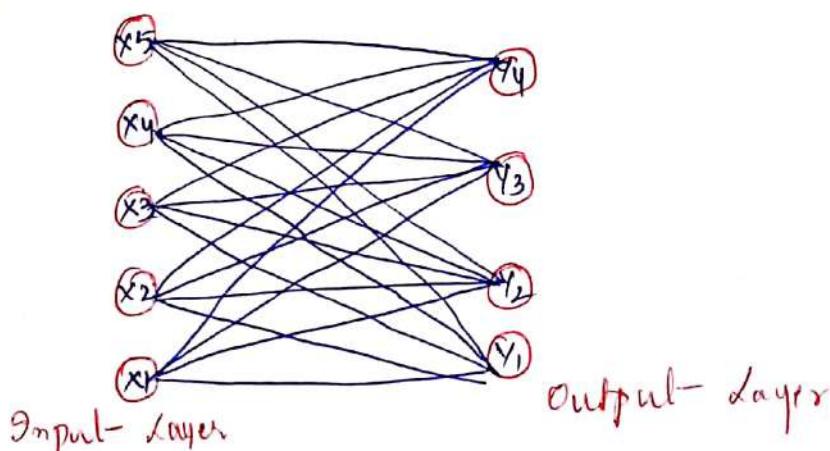
Self Organizing Maps :-

- Supervised training techniques, where there is a target output for each input pattern & the network learns to produce the required output.
- Unsupervised training - is where the network learn to form their own classifications of the training data without external help.
- One particularly interesting class of unsupervised based on competitive learning in which the neurons compete amongst themselves to be activated & this activated neuron is called a winner-takes-all neuron or winning neuron.

→ Such competition can be implemented by having lateral inhibition connections between the neurons resulting in neurons organizing themselves. Such a network is called as self Organizing Maps (SOM)

→ The goal of SOM is to transform an incoming signal pattern of arbitrary dimension into a 2 dimensional discrete map, ∴ setting up a SOM.

→ We shall concentrate on a particular kind of SOM known as Kohonen network. This is a feed forward structure with a single computational layer arranged in rows & columns. Each neuron is fully connected to all source nodes in the input layer.



→ The aim is to learn a "feature map" from the spatially "continuous Input Space" to a low dimensionality spatially discrete output space which is formed by arranging the computational neurons into a grid.

Algorithm :-

Step 1 :- Each node weights are initialized
Step 2 :- Every node in the network is examined to calculate which one's weight are most like the input vector. The winning node is commonly known as the "Best matching unit"
Step 3 :- Update the weight after finding the winning neuron

Step 4 :- Training will stop if input & weight vectors are identical otherwise goto step - 2

$$w^{new} = w^{old} + \alpha \eta (g_{ip} - \text{weight vector})$$

K-means algorithm :-

1

* Assume Euclidean space / distance

* start by picking k , the number of clusters

* Initialize clusters by picking one point per

cluster
note: for this moment, assume we pick the k points at random.

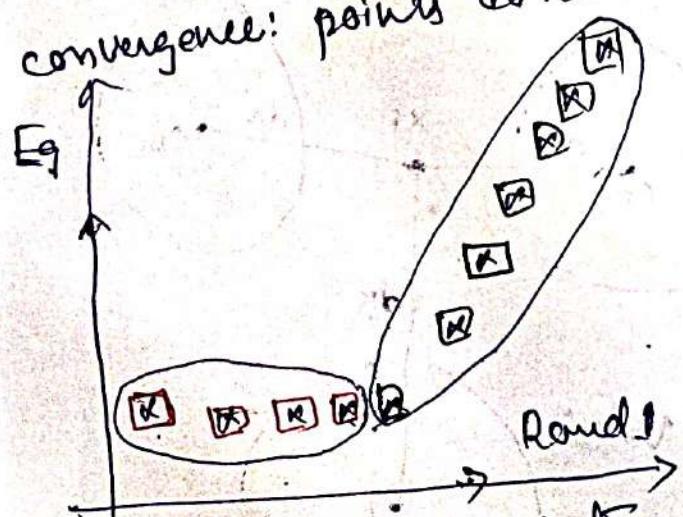
Step 1: for each point, place it in the cluster whose current centroid it is nearest.

Step 2: after all points are assigned, update the locations of centroid of the k clusters

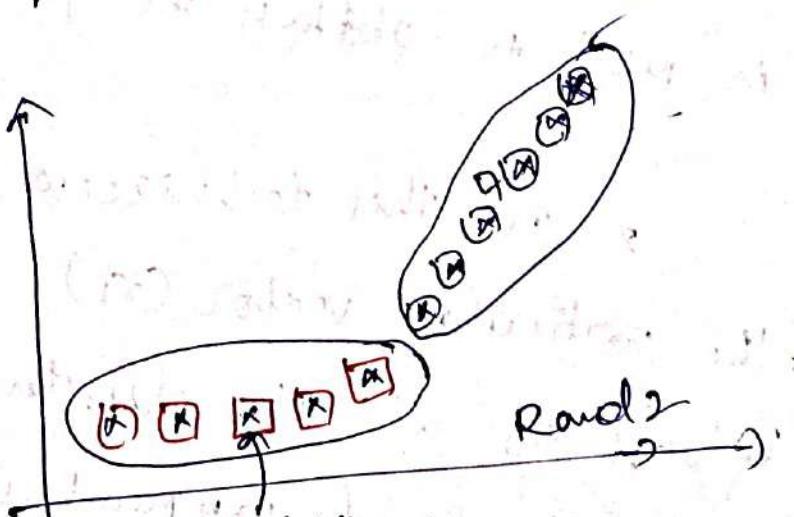
Step 3: Reassign all points to their closest centroid
note: some times moves points b/w clusters

Step 4: Repeat step 2 and step 3 until convergence

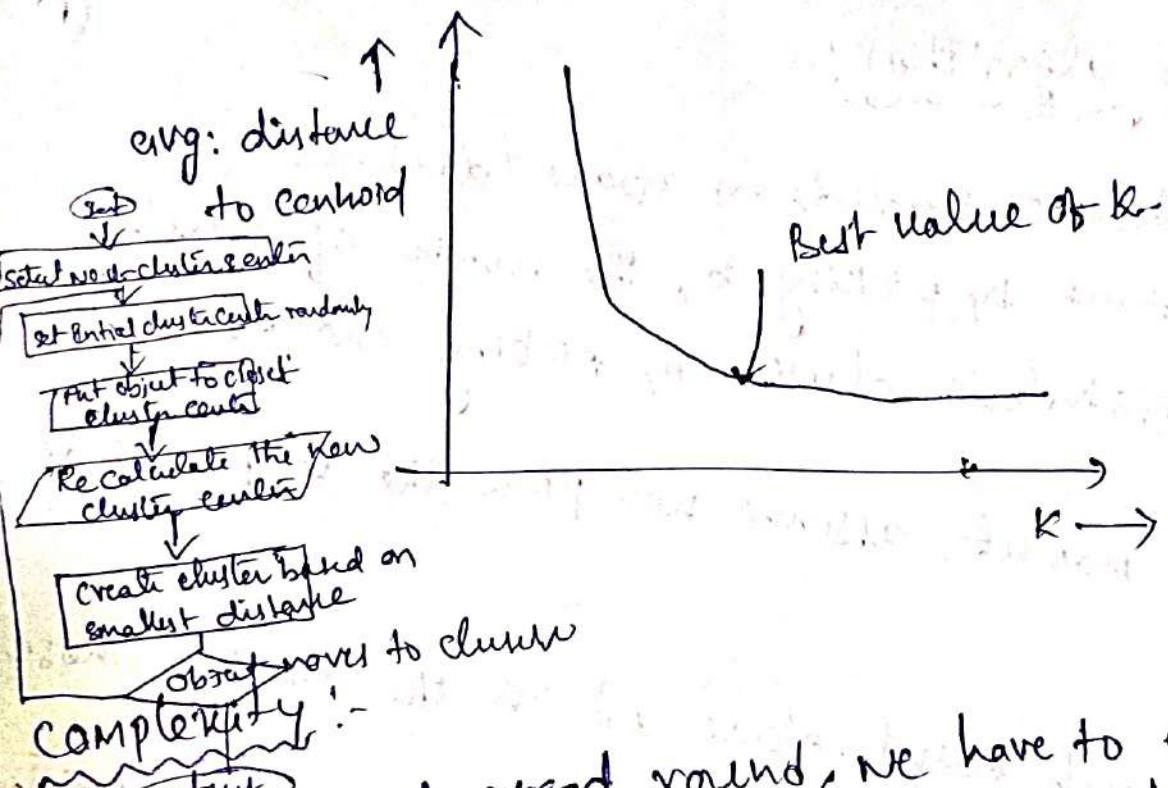
convergence: points do not move b/w clusters and centroids stable



Because of new data
centroids change



in round-'3' no change in
centroid position.



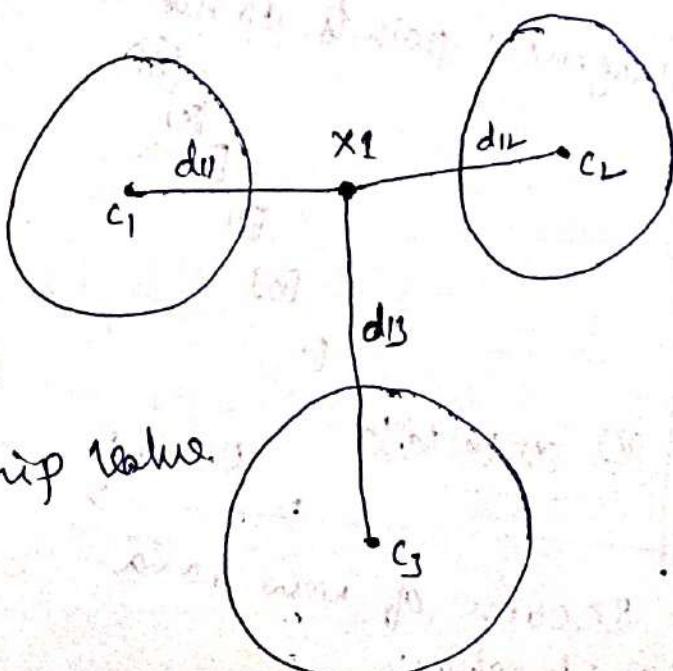
Complexity :-

- 1) In each round, we have to examine each input point exactly once to find closest centroid.
- 2) If data is too large then computation of ~~all~~ iterations are required.

2) Fuzzy K-means algorithm

1. In fuzzy - k means techniques, fuzzy set theory is used to obtain the optimal values of the centroid

2. In this technique the particular vector (x_1) belongs to all the clusters with different membership value.



Note: Sum of membership value is '1'

3. Consider the vector x_1 belongs to the cluster 1.

whose centroid is c_1 , then $(x_1 - c_1)^2$ must be minimized.
But if x_1 belongs to cluster one with membership value m_1 then $m_1(x_1 - c_1)^2$ has to be minimized. similarly

~~x_1 belongs to cluster 2 with membership value m_2~~

~~x_1 belongs to cluster 2 with membership value m_2~~

then $(x_1 - c_2)^2 \cdot m_2$ has to be minimized.

The fuzzy k -means problem is treated as the

optimization technique for obtaining the membership values along with the centroids of the individual clusters

such that

$$\sum_{i=1}^N \sum_{k=1}^K m_{ik}^2 (x_i - c_k)^2 \text{ is minimized}$$

$m_{ik} \rightarrow$ membership value of vector x_i belongs to the cluster k .

$x_i \rightarrow$ i th vector

$c_k \rightarrow$ c_k is k th centroid

~~Note:~~ Here we need converging the membership values.

Fuzzy K-means Algorithm:

Step 1: Initialize the membership values (M_{ij}) randomly.

Step 2: Compute the centroids of the clusters (Eg: 6 clusters)

$$\text{Hint: } C_1 = \underline{\underline{x}_{1:4}}$$

Hint

$C_1 = [\text{vector}_1 * (\text{membership value of vector}_1 \text{ belongs to cluster}_1) \text{ i.e. } M_{11} + \text{vector}_2 * (\text{membership value of vector}_2 \text{ belongs to cluster}_1) \text{ i.e. } M_{12} + \dots + \text{vector}_{100} * (\text{membership value of vector}_{100} \text{ belongs to cluster}_1) \text{ i.e. } M_{100}] / \text{sum of the squared values of the membership values belonging to the cluster}_1.$

values belonging to the cluster-1 by all centroids (Eg: c_1, c_2, \dots, c_6) are obtained

Step 3: Update the membership values M using new (current) values of centroids.

$$M_{ii} = \frac{(x_i - c_1)^v}{(x_i - c_1)^v} + \frac{(x_i - c_2)^v}{(x_i - c_2)^v} + \dots + \frac{(x_i - c_6)^v}{(x_i - c_6)^v}$$

Similarity
Other membership
values are computed.

$$M_{45} = \frac{(x_4 - c_1)^v}{(x_4 - c_1)^v} + \frac{(x_4 - c_2)^v}{(x_4 - c_2)^v} + \dots + \frac{(x_4 - c_6)^v}{(x_4 - c_6)^v}$$

(3)

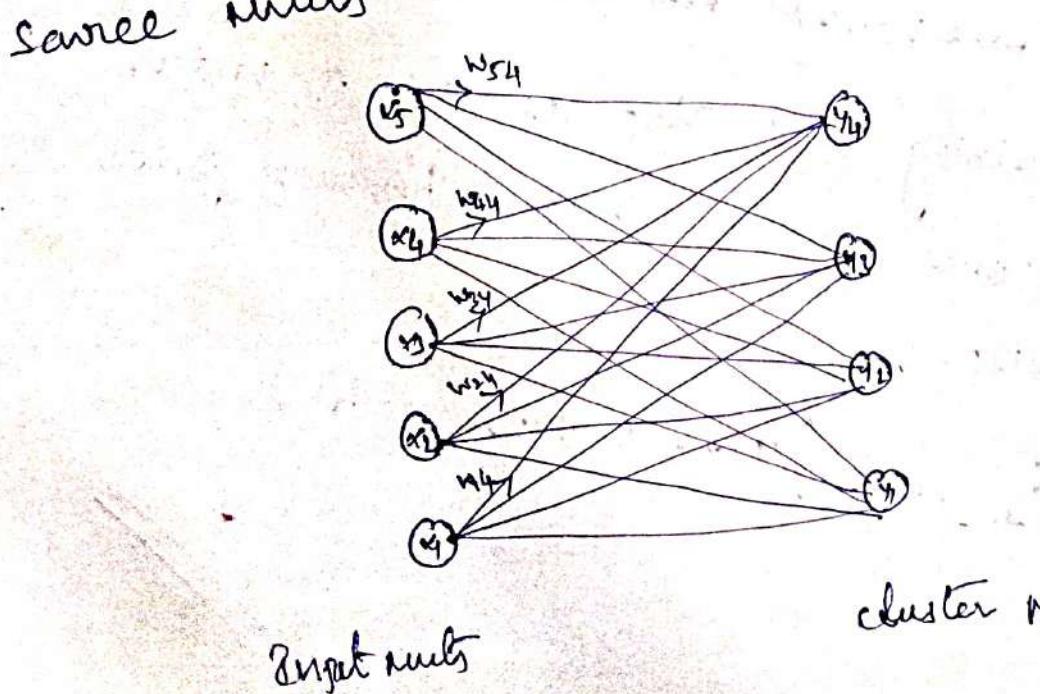
Step 4 : compute the sum of squared difference between the previous membership value and the current membership value.

If computed value is not less than the threshold then go to step-2.
 If computed value is less than the threshold then stop the iteration.

Eg : $(M_{11}^{old} - M_{11}^{new})^2 + (M_{12}^{old} - M_{12}^{new})^2 + \dots$

KSOM (Kohonen self-organizing maps)

This has a feed-forward structure with a single computational layer of neurons, arranged in a 2D grid. Each neuron is fully connected to all the input units in the input layer.



input units

cluster unit

The aim is to learn a "feature map" from the specially "continuous input space", in which our I/p vectors live, to the low dimensional spatially discrete output space which is formed by arranging the computational neurons into a grid.

KSOM (Algorithm)

Step 1: Each node's weights are initialized
(all I/p & O/p are normalized).

Step 2: Every node in the network is examined to calculate which ones' weights are most like the I/p vector.
The winning node is commonly known as the "Best matching unit".

Step 3: Update the weight after find the winner vector
(more importance to the neuron which is closer to winner, less importance to the other neurons which is far from winner)

Step 4: Training will stop if I/p & weight vectors are identical. Otherwise go to step-2.

$$w_{\text{new}} = w_{\text{old}} + \alpha^n (I/p - \text{weight vector})$$

Particle swarm Algorithm :-

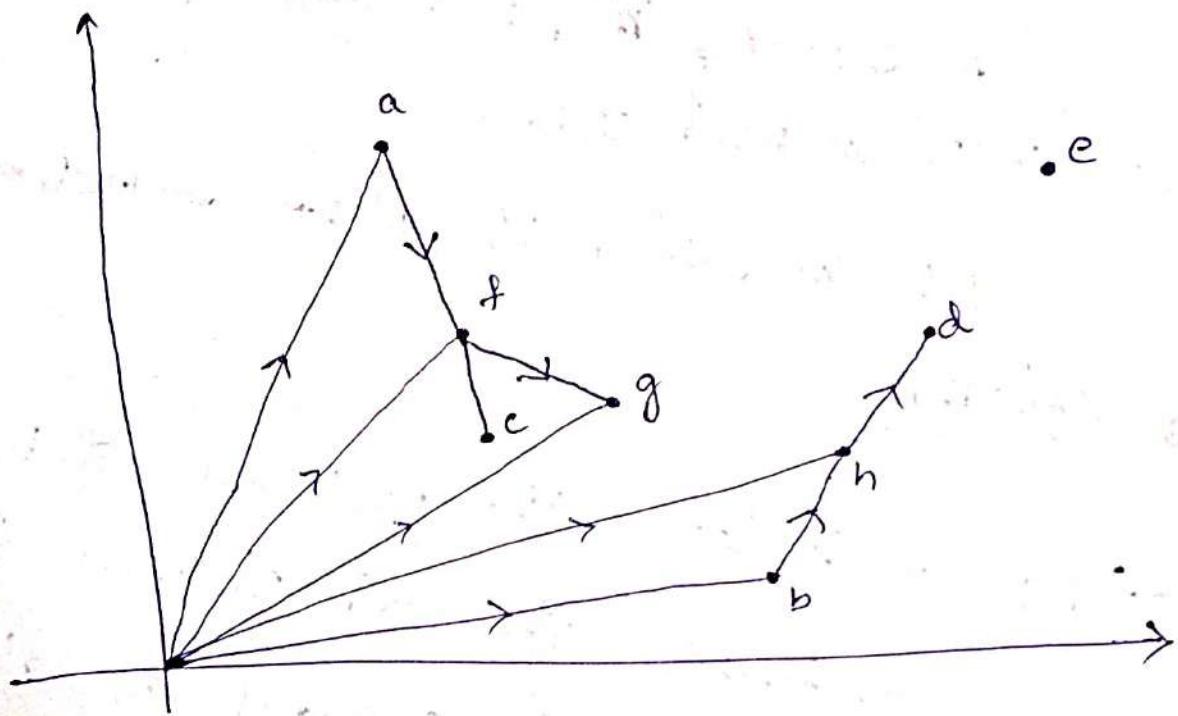


fig: vector representation of PSO Algorithm.

- Both are trying to reach 'e'!
- 1) Let the position of the swarms be at 'a' and 'b'. Both are trying to reach 'e'!
 - 2) Let 'a' decided to move 'c' and 'b' is decided to move 'd'. 'c' is the distance b/w position 'c' and 'e' is greater than the distance b/w 'd' and 'e'.
 - 3) So, 'd' is a common position decided by the a and b
 - 4) The position 'c' is the individual decision taken by 'a'
 - 5) The position 'd' is the "common position" to take by both 'a' and 'b'. The position 'd' is common position to take by both 'a' and 'b'.

Now let us consider three swarms (A, B, C), are trying to reach the particular destination point. ~~Book~~
A desired "A", B decides B' , C decides "C'"

which is its next position.

Let the distance b/w A & ~~C~~ be d . ~~and~~ what

PSO algorithm

Step 1 : Initialize the positions a, b, c, d, e

Step 2 : Initialize the next positions decided by the individual swarms as a', b', c', d', e'

Step 3 : Global decision regarding the next position is computed as follows.

compute $f(a', b, c, d, e)$, $f(a, b', c, d, e)$, $f(a, b, c', d, e)$

$f(a, b, c, d', e)$ and $f(a, b, c, d, e')$. find the minimum among the computed values.

If $f(a, b, c, d, e)$ is minimum among all, the global position decided regarding the next position is (a) .

If $f(a, b', c, d, e)$ is minimum the global position is b' and so on.

Let

Let the selected global position is represented as "global"

step 4: Net value for 'a' is computed as the linear combination of c_1 , $(a'-a)$, $(\text{global} - a)$. i.e

$$\text{next } a = a + \frac{\text{weight of individual decision}}{\text{weight of random decision}} c_1 * \text{RAND} * (a' - a) + \frac{\text{weight of random decision}}{\text{weight of individual decision}} c_2 * \text{RAND} * (\text{global} - a)$$

$$\text{next } b = b + \frac{\text{weight of individual decision}}{\text{weight of random decision}} c_1 * \text{RAND} * (b' - b) + \frac{\text{weight of random decision}}{\text{weight of individual decision}} c_2 * \text{RAND} * (\text{global} - b)$$

$$\text{next } c = c + \frac{\text{weight of individual decision}}{\text{weight of random decision}} c_1 * \text{RAND} * (c' - c) + \frac{\text{weight of random decision}}{\text{weight of individual decision}} c_2 * \text{RAND} * (\text{global} - c)$$

step 5 : change the current value for a', b', c', d', e'

as $\text{next } a, \text{next } b, \text{next } c, \text{next } d, \text{next } e$.

step 6 : If $f(\text{next } a, b', c', d', e')$ is less than $f(a', b', c', d', e')$ then update the value for $f(a', b', c', d', e')$ otherwise 'a' is not chosen "a" as "next a" otherwise 'a' is not chosen "a" as "next a" (do this for b', c', d', e' also.)

step 7 : Repeat the steps 3 to 6 for much iteration to reach the final decision.

The values c_1 & c_2 are decided based on the weightage given $\frac{w_1}{w_2}$ to individual decision and global decision respectively.

Let $\Delta a(t)$ is the change in the value for updating the value for 'a' in t^{th} iteration, then next at $(t+1)^{th}$ iteration can be computed using the following formula.

This is considered as the velocity for updating the position of the swarm in every iteration

$$next\ a(t+1) = a(t) + \Delta a(t+1)$$

where

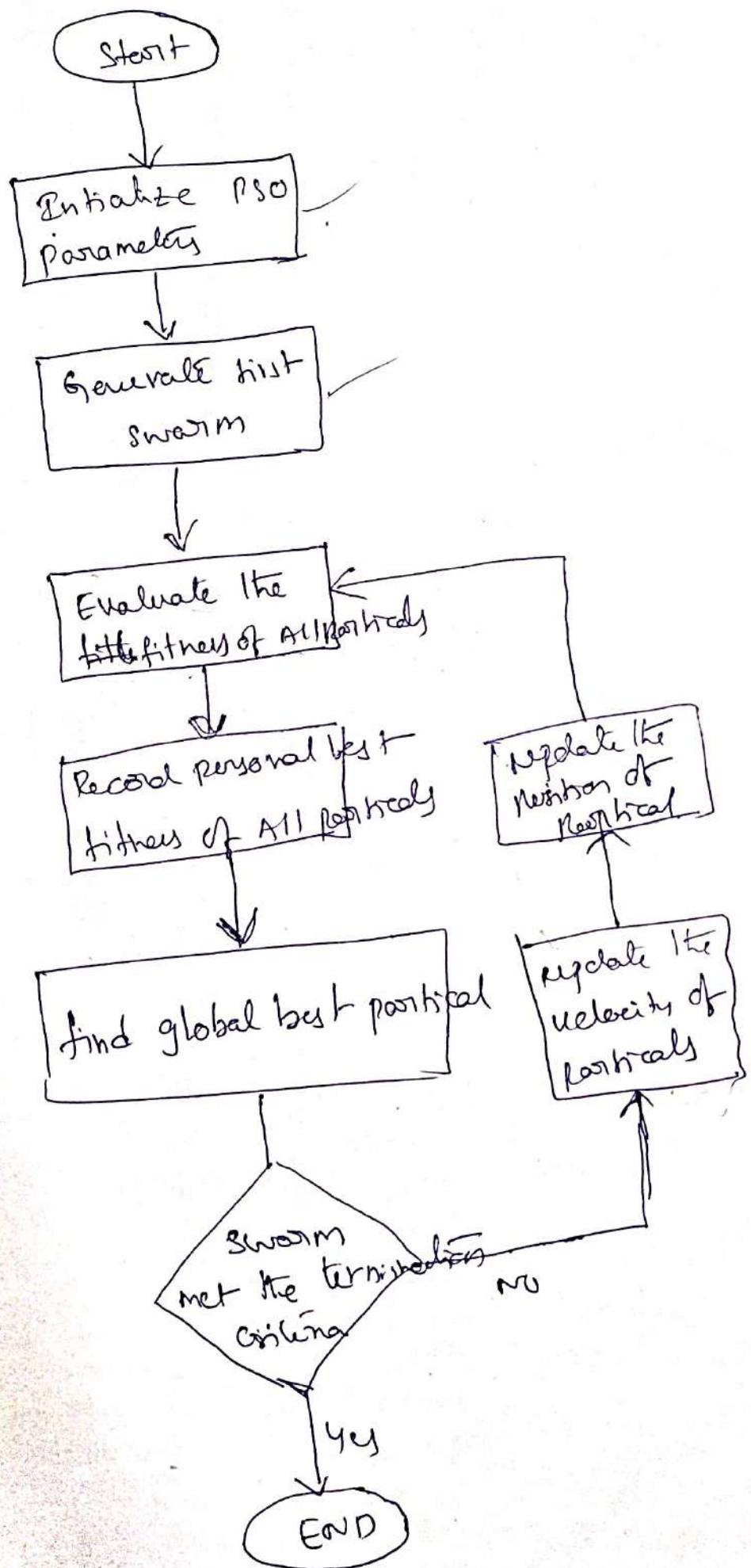
$$\Delta a(t+1) = c_1 * rand * (a' - a) + c_2 * rand * (\text{Global}\ a) + w(t) * 4a(t)$$

$w(t)$ is the weight at t^{th} iteration. The value for 'w' is adjusted at every iteration as given below, where 'iter' is total no. of iterations

used.

$$w(t+1) = \frac{w(t) - t * w(t)}{iter}$$

PSO Algorithm flow chart



(b)

Ant colony optimization

Ant colony optimization technique exploit the natural behavior of ants in finding the shortest path to reach destination from the source.

Algorithm

Consider the problem of finding the optimum order in which the numbers from 1 to 8 are arranged so that the cost of the order is maximized. Cost of the particular order is computed using the two matrices A and B as described below

a ₁								
a ₂	a ₃							
a ₄	a ₅	a ₆						
a ₇	a ₈	a ₉	a ₁₀					
a ₁₁	a ₁₂	a ₁₃	a ₁₄	a ₁₅				
a ₁₆	a ₁₇	a ₁₈	a ₁₉	a ₂₀	a ₁			
a ₂₁	a ₂₂	a ₂₃	a ₂₄	a ₂₅	a ₂₆	a ₂₇		
a ₂₉	a ₃₀	a ₃₁	a ₃₂	a ₃₃	a ₃₄	a ₃₅	a ₃₆	

matrix A

b ₁								
b ₂	b ₃							
b ₄	b ₅	b ₆						
b ₇	b ₈	b ₉	b ₁₀					
b ₁₁	b ₁₂	b ₁₃	b ₁₄	b ₁₅				
b ₁₆	b ₁₇	b ₁₈	b ₁₉	b ₂₀	b ₂₁			
b ₂₂	b ₂₃	b ₂₄	b ₂₅	b ₂₆	b ₂₇	b ₂₈		
b ₂₉	b ₃₀	b ₃₁	b ₃₂	b ₃₃	b ₃₄	b ₃₅	b ₃₆	

matrix B

The cost of the order [3 2 1 5 7 2 4 6 8] is computed as the sum of the products of values obtained from the positions (3,1) $\delta(1,3)$, (2,2), (3,1), (4,5) $\delta(5,4)$, (5,7) $\delta(7,5)$, (6,2), (7,4), (8,6) of matrix A & B
 $(a_6 \times b_6) + (a_3 \times b_3) + (a_6 \times b_6) + (a_{14} \times b_{14}) + (a_{28} \times b_{28}) + (a_{25} \times b_{25}) + (a_{36} \times b_{36})$

The problem of finding the optimal order is achieved using Ant colony optimization technique as described below.

Step 1: Generate the five sets of orders randomly which are treated as the initial values selected by the 5 ants respectively.

The following are the orders selected by the 5 ants along with the corresponding cost as given below

<u>ANT</u>	<u>ORDERS</u>	<u>COST</u>
ANT 1	7 1 3 4 6 2 5 8	c1
ANT 2	4 2 8 5 7 1 3 6	c2
ANT 3	5 1 3 2 7 6 8	c3
ANT 4	3 7 4 1 8 6 5 2	c4
ANT 5	8 1 2 4 3 5 7 6	c5

Step 2: Pheromone is the value assigned for i^{th} number for the j^{th} position of the sequence selected by all the ants. Initially all the values of matrix are completely filled rep with ones. The values of pheromone matrix are updated in every iteration.

$$\text{pheromone matn}(n+1) = \text{pheromone matn}(n) + \text{pheromone matn}.$$

(7)

	1	2	3	4	5	6	7	8
1	0	$\frac{1}{c_1} + \frac{1}{c_3} + \frac{1}{c_5}$	0	$\frac{1}{c_4}$	0	$\frac{1}{c_2}$	0	0
2	0	$\frac{1}{c_2}$	$\frac{1}{c_5}$	$\frac{1}{c_3}$	0	$\frac{1}{c_1}$	0	$\frac{1}{c_4}$
3	$\frac{1}{c_4}$	0	$\frac{1}{c_1} + \frac{1}{c_3}$	0	$\frac{1}{c_5}$	0	$\frac{1}{c_2}$	0
4	$\frac{1}{c_2}$	0	$\frac{1}{c_4}$	$\frac{1}{c_5} + \frac{1}{c_1}$	$\frac{1}{c_3}$	0	$\frac{1}{c_5}$	$\frac{1}{c_4} + \frac{1}{c_3}$
5	$\frac{1}{c_3}$	0	0	$\frac{1}{c_2}$	0	$\frac{1}{c_4}$	$\frac{1}{c_1}$	$\frac{1}{c_5} + \frac{1}{c_3}$
6	0	0	0	0	0	$\frac{1}{c_1}$	$\frac{1}{c_4}$	$\frac{1}{c_3} + \frac{1}{c_5}$
7	$\frac{1}{c_1}$	$\frac{1}{c_4}$	0	0	0	$\frac{1}{c_2}$	$\frac{1}{c_3}$	$\frac{1}{c_5}$
8	$\frac{1}{c_5}$	0	$\frac{1}{c_2}$	0	$\frac{1}{c_4}$	0	0	$\frac{1}{c_1} + \frac{1}{c_3}$

Step 3: Neat set of orders selected by the ants
 as described below:

Eg: The random position sequence is $7\ 8\ 2\ 3\ 6\ 4\ 15$ (say)

Let the order decided by the ant 1 be represented as $[u_1, u_2, u_3, u_4, u_5, u_6]$ [u_7, u_8]

The value for the u_7 in the order as mentioned above is selected first as the first number of the randomly generated position sequence is 7.

* To select the value for u_7 , 7th column of the pheromone matrix is used as below,

$$\text{column } 7 = [0 \ 0 \ \frac{1}{c_2} \ 0 \ \frac{1}{c_4} + \frac{1}{c_5} \ \frac{1}{c_3} \ \frac{1}{c_5} \ 0]$$

* Now calculate the normalized value of column 7.

This is also called as the probability vector used for selecting the value for u_7 in the order

$$S = [\frac{1}{c_2} + \frac{1}{c_4} + \frac{1}{c_5} + \frac{1}{c_3} + 1] = \text{value}$$

$$\text{Normalized vector } t = \frac{\text{column } 7}{S} = \left[\frac{0}{S} \ \frac{0}{S} \ \frac{\frac{1}{c_2}}{S} \ \frac{0}{S} \ \frac{\frac{1}{c_4} + \frac{1}{c_5}}{S} \right]$$

$$= [p_1 \ p_2 \ p_3 \ p_4 \ p_5 \ p_6 \ p_7 \ p_8] \left(\frac{1}{c_3} + \frac{1}{c_5} \ \frac{0}{S} \right)$$

* find the lowest value position, for example

* find the lowest value the 5 is assigned to u_7 .

If p_5 is the lowest value the 5 is assigned to u_7 .

Note: position no. '5' is not used again.

* The next value in position sequence is 8,

so column '8' is selected from pheromone matrix, and select

the normalized value

$$\text{column } 8 = [0 \ \frac{1}{c_4} \ 0 \ 0 \ \frac{1}{c_2} + \frac{1}{c_3} \ 0 \ \frac{1}{c_4} + \frac{1}{c_3}]$$

$$S = \left[\frac{1}{c_4} + \left(\frac{1}{c_2} + \frac{1}{c_1} \right) + \left(\frac{1}{c_1} + \frac{1}{c_3} \right) \right]$$

Normalized vector $\bar{v} = \frac{\text{column } 8}{S} = (a_1, a_2, a_3, a_4, a_5, a_6, a_7)$

find the lowest value in normalized vector '8' and replaced with v8. Examp if a_3 is lowest then number '3' is assigned to v8.

* This process is repeated to obtain the order selected by the ~~ant~~ Ant 1, Ant 2, Ant 3, Ant 4, Ant 5.

* This is called one Iteration.

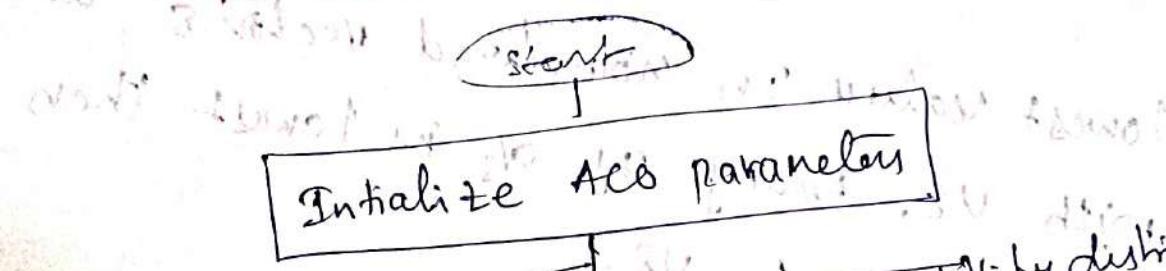
Step 4 : Updating the pheromone matrix for the second iteration is computed as described in step 2. This is called updating pheromone matrix (2)

$$\therefore \text{pheromone matrix}(2) = \text{pheromone matrix}(1) + \text{updating pheromone matrix}$$

Step 5 : Next set the orders selected by the five ants are computed as described in the step 3

Step 6 : Step 3 ~~and~~ to step 5 is ~~repeated~~ for the particular iterations. Thus set of orders are obtained using Ant colony algorithm

Note: Best among the best set of orders selected by the five ants in the last iteration is selected as the global best order.



Initialize ACO parameters

construct solution using the probability distribution
(Pheromone trail and randomization)

Local updating of pheromone

NO

All ants visited
through all cities?

Yes

congest the length of the optimal paths
and update only the amount of the
pheromone on the optimal paths

NO

The termination condition satisfied

Yes

output values contain the
maximum pheromones

①

~~local foraging~~ ^{also see} ~~some~~

1- Bacterial tries to move towards the region where more nutrients are available (i.e) moving ~~to~~ towards increasing concentrations of nutrients

$\nabla I(\Theta) \downarrow \downarrow \rightarrow$ objective function

Θ = the position of the Bacteria

$I(\Theta)$ = is the nutrient profile

$I(\Theta) < 0 \rightarrow$ presence of nutrients

$I(\Theta) = 0 \rightarrow$ neutral

$I(\Theta) > 0 \rightarrow$ noxious substance

1. Chemo taxis: tendency to move ~~to~~ towards distant sources of nutrients

a) Tumbling [change in direction]:

$$\Theta_{\text{new}} = \Theta^i + C \phi$$

→ run walk in random directions

$$\phi = \frac{\Delta}{\Delta t \alpha}$$

chemotactic step size

b) Swimming:

movement in same direction as the previous tumbling.

Note: Bacterium will tend to keep moving if it is headed in the direction of an increasingly favorable environment.

If at θ , then η is lower (Bacter) than the cost at θ_i , another swimming step is taken and is continued as long as it contains maximum no. of swimming

steps Ns.

2. Reproduction

After Ns chemoatetic steps, reproduction

steps is taken

- a) Least healthy bacteria die
- b) " or healthiest bacteria split into two bacteria which are then placed in the same location.

Note: Health of the bacteria are measured by

Objective function J.

Higher cost represents that the bacterium did not get many nutrients (not as healthy) and hence unlikely to reproduce.

3. Elimination - Dispersion

In real world process, bacteria can be dispersed to new locations, for example wind dispersal.

Individual bacterium is stochastically selected for elimination from the population and is replaced by a new bacterium located at a random new location within the optimization domain, according to the present probability P_{ed} .

4. social communication

Bacterium that has perceived good source of nutrients can release a "chemical attract signal" which attracts the other bacteria to converge to its current location. (This help in making all the bacterium to come closer)

Also "chemical repellent signal" is released to ensure that the bacteria do not get too close to each other.

Let Θ^i be the position of the i^{th} bacteria

then the objective function J is modified as,

$$J_{sw}^i = J^i + J_{cc}(\Theta^i, \Theta)$$

$F_{cc}(\theta^i, \theta)$ is chosen such that the distance b/w it with bacteria and all other bacteria are maximized (for repellent signal) and distance b/w it's bacteria and all other bacteria are minimized (for attractant signal).

Note: The width of the attractant and repellent signals are used band as the requirement.

$$F_{cc}(\theta^i, \theta) = -M \left(\sum_{k=1}^{s_{\text{attraction}}} e^{-\frac{\theta^i - \theta_k}{w_{\text{attraction}}}} - \sum_{k=1}^{s_{\text{repellent}}} e^{-\frac{\theta^i - \theta_k}{w_{\text{repellent}}}} \right)$$

width of the
attraction

width of the
repellent

(a)

Let P = Dimension of the search space

P = Dimension of the population

N_b = Total no. of bacteria in the population

N_c = Total no. of chemotactic steps

N_s = The swimming length

N_{re} = No. of reproduction steps

N_{ed} = No. of elimination-dispersal events

P_{ed} = Elimination-dispersal probability

$c(i)$ = The site of the step taken in the random

direction specified by the tumble.

(3)

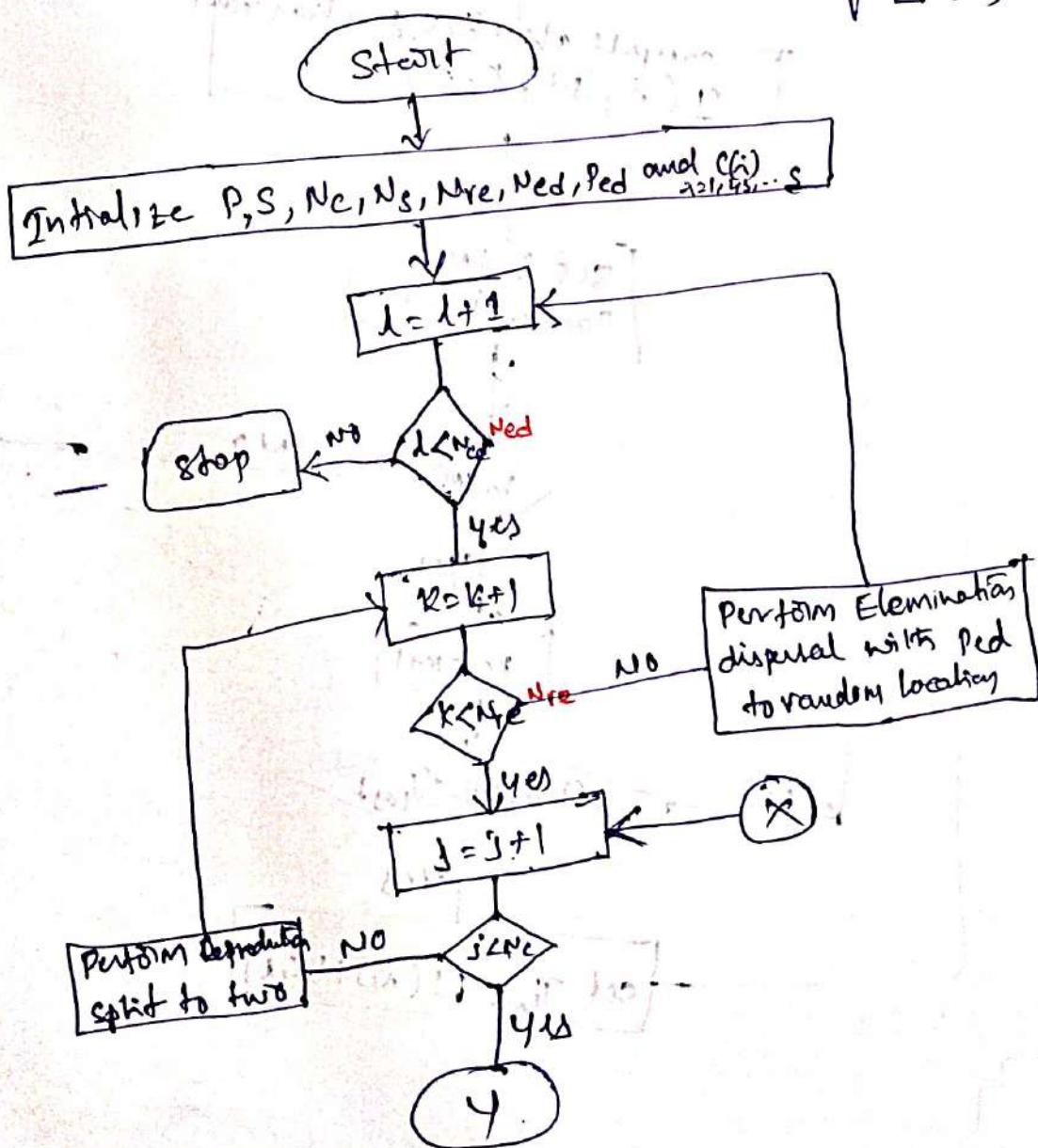
Let : $\Theta^i(j, k, l) =$ i^{th} bacterium at j^{th} chemoattractant
 k^{th} reproductive and l^{th} elimination

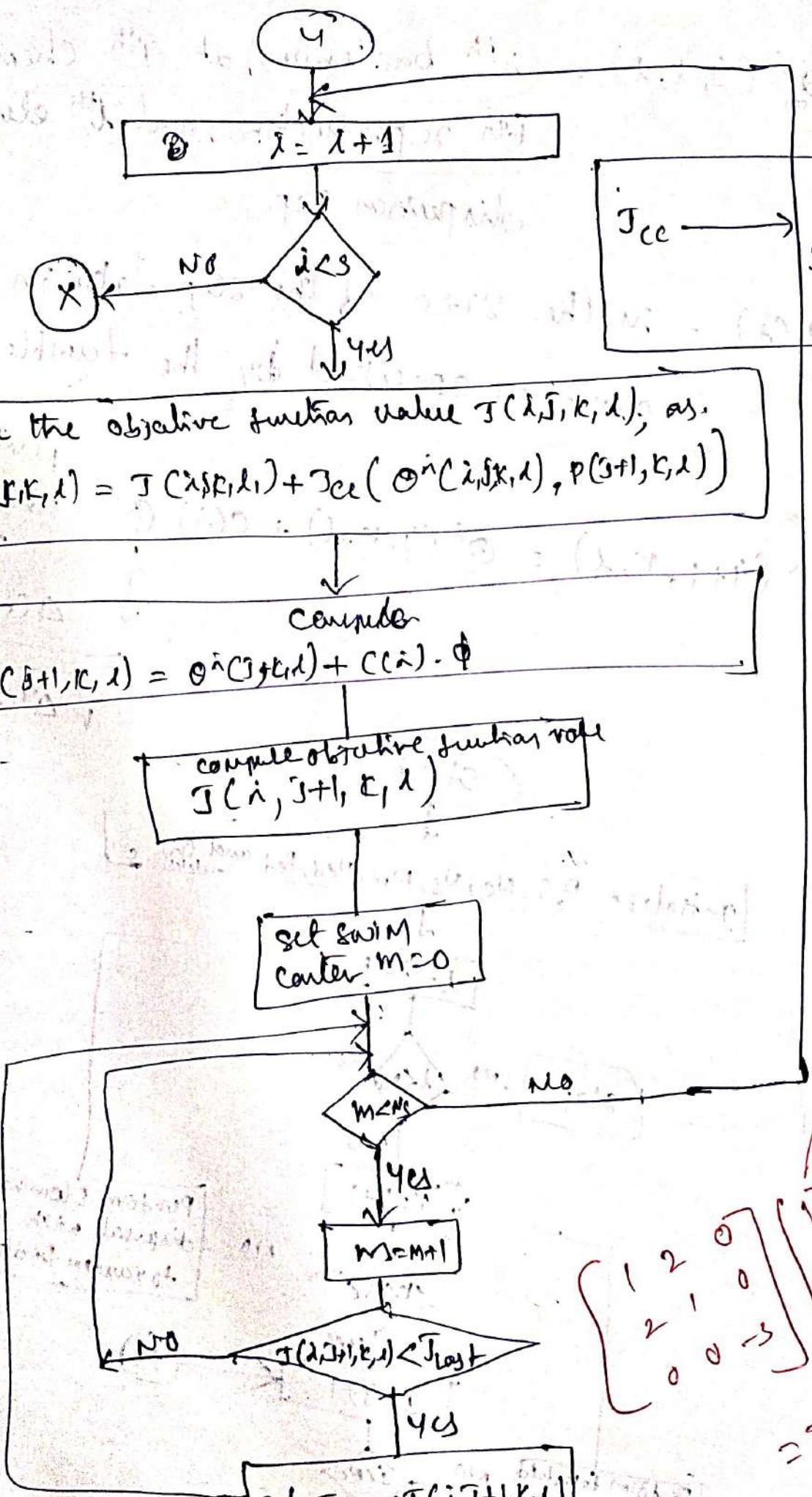
dispersion step.

$c(i)$ = is the size of the step taken in the random direction, specified by the tumble.

$$\Theta^i(j+1, k, l) = \Theta^i(j, k, l) + c(i) \phi$$

↗ unit walk
 in random
 directions
 $\Delta(i)$
 $\sqrt{\Delta T(i) \cdot \Delta(i)}$





Q.1 Use k-means clustering.
To divide data into 2 clusters.

x_1	1	2	2	3	4	5
x_2	1	1	3	2	3	5

Step 1: choose 2 random cluster centres.

Say $v_1 = (2, 1)$ $v_2 = (2, 3)$

Data point	Distance from v_1	Distance from v_2
a ₁ (1, 1)	1	2.24
a ₂ (2, 1)	0	2
a ₃ (2, 3)	2	0
a ₄ (3, 2)	1.41	1.41
a ₅ (4, 3)	2.83	2
a ₆ (5, 5)	5	3.61

$$(x_1, x_2) \\ (1, 1)$$

$$(y_1, y_2) \\ (2, 1)$$

$$\sqrt{(1-2)^2 + (1-1)^2} = \sqrt{1+0} \\ = 1$$

Step 3 :- Cluster 1 of V_1

$$: \{ \cancel{a_1, a_2} \{ a_1, a_2, a_4, \} \}$$

cluster 2 of V_2 ; $\{ a_3, a_5, a_6 \}$

Step 4 :- Recalculate cluster centers.

$$V_1 = \frac{1}{3} [(1, 1) + (2, 1) + (3, 2)]$$

$$\text{and} \\ \text{where} \\ V_1 = \frac{1}{3} [6, 4]$$

$$V_1 = (2, 1.33)$$

$$V_2 = \frac{1}{3} [a_3 + a_5 + a_6]$$

$$V_2 = \frac{1}{3} [(2, 3)(4, 3)(5, 5)]$$

$$= \frac{1}{3} [(11, 11)]$$

$$= (3.67, 3.67).$$

Step 5 :- Repeat from step 5
until we get same cluster centers
as same cluster elements.

Distance table :-

Data point -	$\text{dis. from } v_1(2, 1.33)$	$v_2(3.67, 3.67)$	$v_1 \text{ cen}$
a ₁	1.05	3.78	<u>v₁</u>
a ₂	0.33	3.15	<u>v₁</u>
a ₃	1.67	1.8	<u>v₁</u>
a ₄	1.204	1.8	<u>v₁</u>
a ₅	2.605	0.75	<u>v₂</u>
a ₆	4.74	1.88	<u>v₂</u>

cluster 1 of v_1 : {a₁, a₂, a₃, a₄}

v_2 : {a₅, a₆}

Recalculating cluster centers

$$v_1 = \frac{1}{4} [a_1 + a_2 + a_3 + a_4]$$

$$= \frac{1}{4} [(1, 1) + (2, 1) + (2, 3) + (3, 2)]$$

$$= \frac{1}{4} (8, 7) = (2, 1.75)$$

$$v_2 = \frac{1}{2} [a_5 + a_6]$$

$$= \frac{1}{2} [(4,3) + (5,5)]$$

$$= \frac{1}{2} (9,8) = 4.5, 4$$

* centers not same.

Again distance table.

	v_1	v_2	Assigned v_1 center
a_1	1.25	4.61	v_1
a_2	0.75	3.9	v_1
a_3	1.25	2.69	v_1
a_4	1.03	2.5	v_1
a_5	2.36	1.12	v_2
a_6	4.42	1.12	v_2

cluster 1 of $v_1 = \{a_1, a_2, a_3, a_4\}$

of $v_2 = \{a_5, a_6\}$

* cluster elements are same as in previous iteration.

cluster 1 : { (1, 1), (2, 1), (2, 3)
 (3, 2) }

cluster 2 : { (4, 3), (5, 5) }

$$[m_{11}^{\text{old}} - m_{11}^{\text{new}}]^2 + [m_{12}^{\text{old}} - m_{12}^{\text{new}}]^2 + \dots$$

- * If computed value is not less than the threshold then goto step 2.
- * If computed value is less than the threshold then stop the iteration.

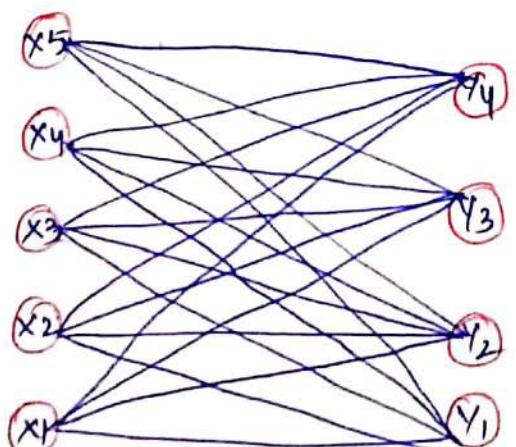
Self Organizing Maps :-

- Supervised training techniques, where there is a target output for each input pattern & the network learns to produce the required output.
- Unsupervised training - is where the network learn to form their own classifications of the training data without external help.
- One particularly interesting class of unsupervised based on competitive learning in which the o/p neurons compete amongst themselves to be activated & this activated neuron is called a winner-takes.

→ Such competition can be implemented by having lateral inhibition connections between the neurons resulting in neurons organizing themselves. Such a network is called as self Organizing Maps (SOM)

→ The goal of SOM is to transform an incoming signal pattern of arbitrary dimension into 1 or 2 dimensional discrete map, ∴ setting up a SOM.

→ We shall concentrate on a particular kind of SOM known as Kohonen network. This is a feed forward structure with a single computational layer arranged in rows & columns. Each neuron is fully connected to all source nodes in the input layer.



→ The aim is to learn a "feature map" from the spatially "continuous input space" to a low dimensionality spatially discrete output space, which is formed by arranging the computational neurons into a grid.

Algorithm :-

Step 1 :- Each node weights are initialized

Step 2 :- Every node in the network is examined to calculate which one's weight are most like the input vector. The winning node is commonly known as the "Best matching unit"

Step 3 :- Update the weight after finding the winning neuron

Step 4 :- Training will stop if input & weight vectors are identical otherwise goto step - 2

$$w^{new} = w^{old} + \alpha_n (s_ip - \text{weight vector})$$

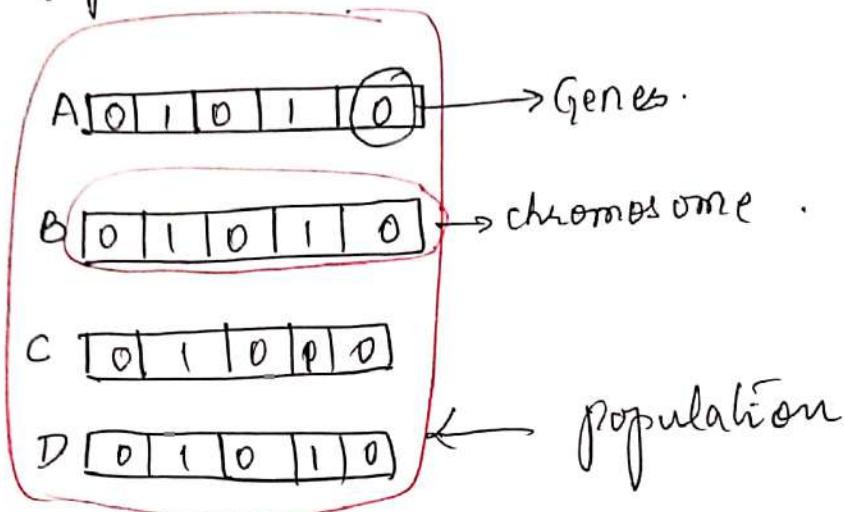
Genetic Algorithm :-

- The Genetic Algorithm was developed by John Holland & is a model or abstraction of biological Evolution based on Charles Darwin theory of natural selection.
- Charles Darwin theory of natural selection Explains survival of fittest overtime, change within species leads to replacement of old species by new species as less successful species becomes Extinct.
- Genetic Algorithms are a class of probabilistic optimization algorithms inspired by biological Evolution process. These algorithms use "natural Selection" & "genetic Inheritance".
- Phases of Genetic Algorithm :-
 - GE Begins with a large population & gradually reduces it using some heuristics until only the best is left -

➤ Initial population :-

(11)

→ Process starts with a set of individuals known as a population.



note :- Genes \rightarrow chromosomes are divided into several parts called genes .

chromosome - All genetic information of a species is stored in chromosomes .

Natural Evolution

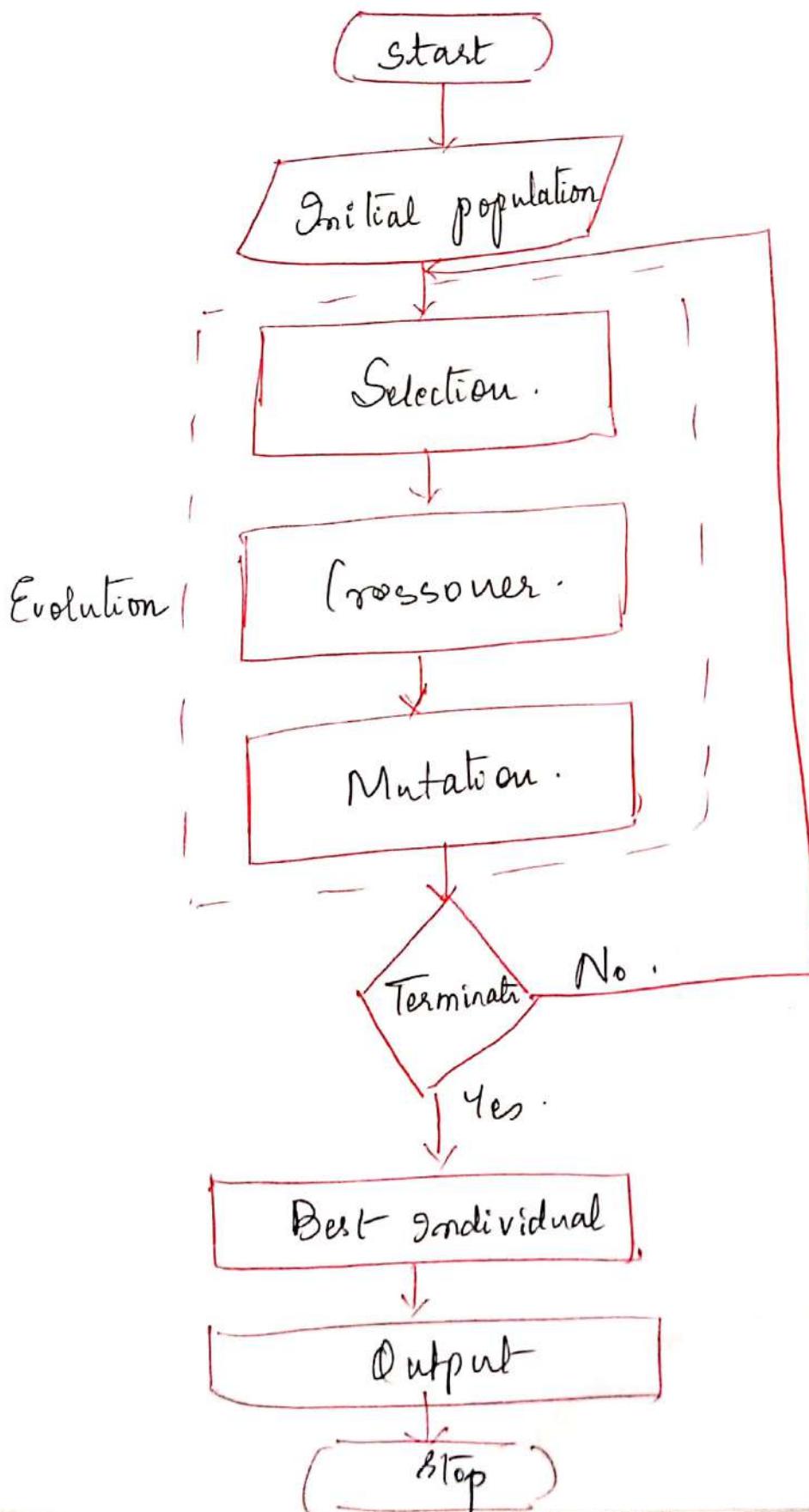
- 1> Gene
- 2> Chromosome
- 3> Genotype
- 4> Phenotype

Genetic Algorithm

- } feature
string .
coded string
Decoded structure

Applications of Genetic Algorithm :-

- 1) Traveling salesman problem.
- 2) Robotics.
- 3) Shipment routing.



2) Fitness Function :-

→ Determines how fit an individual member of the population is.

→ A fitness score is associated with each individual & higher the fitness score , higher the probability that individual will be selected for reproduction to next generation .

3) Selection :-

→ This phase main goal is to find the region where getting the best solution is more .

4) Crossover :-

→ Also called as recombination , where genetic merging takes place to produce new offspring .

5) Convergence .

When there is no improvement in solution quality .

Bacterial Foraging Optimization Algorithm (BFO) :-

(14)

- Foraging means searching for food & Exploiting food resources. It affects an animal's fitness because it plays an important role in an animal's ability to survive & reproduce.
- Natural selection tends to Eliminate those poor "foraging strategies" & favor the reproduction of those animals that have successful foraging strategies.
- Such Evolutionary principles have led Scientists in the field of "foraging theory" to hypothesize.
- Here the algorithm is based on a bacteria E. coli ∵ this optimization is named as Bacterial Foraging Optimization.
- "E. coli" [*Escherichia coli*] is a type of bacteria that normally lives in our intestines of healthy people & animals.

It helps digest the food we eat.

(15)

→ During foraging a real bacteria does two basic operations.

1) Swim

2) Tumble.

Algorithm:-

i) CHEMOTAXIS

→ a) When a bacterium meets a favorable Environment (rich in nutrients)

i) It will continuously swim in the same direction

tumble → run → run

b) When it meets an unfavorable Environment

i) It will tumble → change the direction of swim.

tumble → run → tumble.

2) SWARMING

→ E. coli bacterium has a specific sensing activation & decision making mechanism.

→ On each move it releases attract to signal

(16)

other bacteria towards it.

3) Reproduction :-

→ After calculating fitness value of each bacteria reproduction allows the bacteria to survive & reproduce.

4) Elimination & dispersal

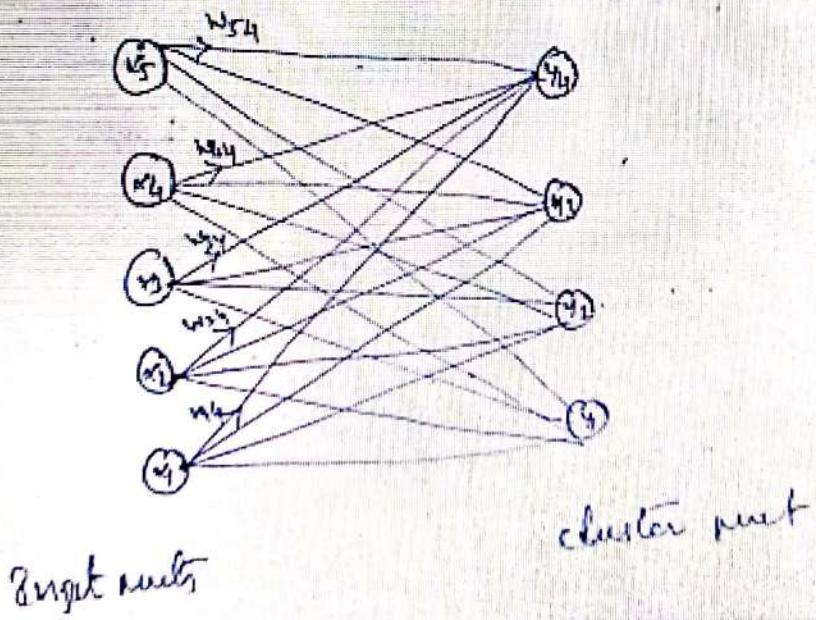
→ Unfit bacterias are eliminated

Eg: SOM

KSOM (Kohonen self-organizing maps)

This has a feed-forward structure with a single computational layer of neurons. Each neuron is fully connected to all the inputs in the input layers.

Source



Step 4: Compute the sum of squared difference between the previous membership value and the current membership value.

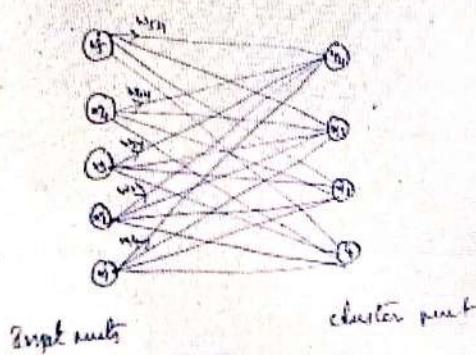
If computed value is not less than the threshold then go to step-2.

If computed value is less than the threshold then stop the iteration.

$$\text{Eg: } (m_{11}^{\text{old}} - m_{11}^{\text{new}})^2 + (m_{12}^{\text{old}} - m_{12}^{\text{new}})^2 + \dots$$

KSOM (Kohonen self-organizing maps)

This has a feed-forward structure with a single computational layer of neurons. ~~connected~~
Each neuron is fully connected to all the source units in the I/p layers.



$$[(x_1 - c_1)^2] [(x_2 - c_2)^2]$$

The aim is to learn a "feature map" from the spatially continuous input space, in which one or 2D vectors live, to the low dimensional spatially discrete output space which is formed by arranging the computational neurons into a grid.

KSOM (Algorithm)

Step 1: Each node weights are initialized (all 2D & 0/p are normalized).

Step 2: Every node in the network is examined to calculate which ones' weights are most like the I/p vector. The winning node is commonly known as the "Best matching unit".

Step 3: Update the weight after find the winner neuron (more importance to the neuron which is closer to winner, less importance to the next neurons which is far from winner)

Step 4: Training will stop if 2/p & weight vectors are identical. otherwise go to step-2.

$$w_i^{\text{new}} = w_i^{\text{old}} + \alpha^{(t)} (I/p - \text{weight vector})$$

Particle Swarm Algorithm:-

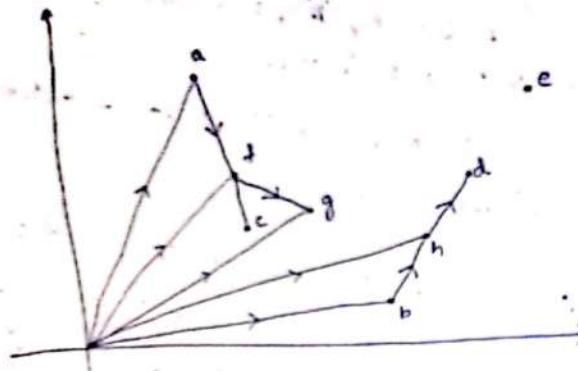


Fig: vector representation of pso algorithm

- 1) Let the position of the swarms be at 'a' and 'b'. Both are trying to reach 'e'.
- 2) Let 'b' decided to move 'c' and 'b' is decided to move 'd'.
- 3) The distance b/w position 'c' and 'e' is greater than the distance b/w 'd' and 'e'.
- 4) So, 'd' is a common position decided by the 'a' and 'b'
- 5) The position 'c' is the individual decision taken by 'a'; The position 'd' is the common position decided by 'a' and 'b'; The position 'd' is the common position to take by 'a' and 'b'.

④

Now, let us consider that swarm (a,b,c) are trying to reach the particular destination point 'e'. Let 'a' decided 'd', 'b' decides 'e', 'c' decided 'd' which is the next position.

Let's discuss this in more detail about pso algorithm.

- Step 1: Initialize the positions a, b, c, d, e
- Step 2: Initialize the next positions decided by the individual swarms as a', b', c', d', e'
- Step 3: Global decision regarding the next position is computed as follows:

compute $f(a', b', c', d', e)$, $f(a, b', c, d, e)$, $f(a, b, c', d, e)$, $f(a, b, c, d', e)$ and $f(a, b, c, d, e)$. find the minimum among the computed values.

If $f(a, b, c, d, e)$ is minimum among all, the global position decided regarding the next position is (a, b, c, d, e) .

If $f(b, c, d, e)$ is minimum the global position is b', c', d', e' and so on.

Let the selected global position is (5) represented as "global".

Step 4: Net value for 'a' is computed as the linear combination of 'a', $(a-a)$, $(\text{global}-a)$. i.e.

$$\text{next}a = a + C_1 \times \text{rand} \times (a-a) + C_2 \times \text{rand} \times (\text{global}-a)$$

$$\text{next}b = b + C_1 \times \text{rand} \times (b-b) + C_2 \times \text{rand} \times (\text{global}-b)$$

$$\text{next}c = c + C_1 \times \text{rand} \times (c-c) + C_2 \times \text{rand} \times (\text{global}-c)$$

Step 5: change the current value for a, b, c, d, e as $\text{next}a, \text{next}b, \text{next}c, \text{next}d, \text{next}e$.

Step 6: If $f(\text{next}a, \text{next}b, \text{next}c, \text{next}d, \text{next}e)$ is less than $f(a, b, c, d, e)$ then update the value for 'a' as "next a" otherwise 'a' is not change (do the for b, c, d, e also)

Step 7: Repeat the steps 3 to 6 for much iteration to reach the final decision.

The values C_1 & C_2 are decided based on the weightage given to individual decision and global decision respectively.

PSO algorithm flow chart

Let $\Delta a(t)$ is the change in the values from step 4. The value for 'a' in 't' iteration, thus next at $(t+1)$ in iteration, can be computed using the following formula:

This is considered as the velocity for updating the position of the swarms in every iteration

$$\text{next}a(t+1) = a(t) + \Delta a(t+1)$$

where

$$\Delta a(t+1) = C_1 \times \text{rand} \times (a-a) + C_2 \times \text{rand} \times (\text{global}-a) + w(t) \times \Delta a(t)$$

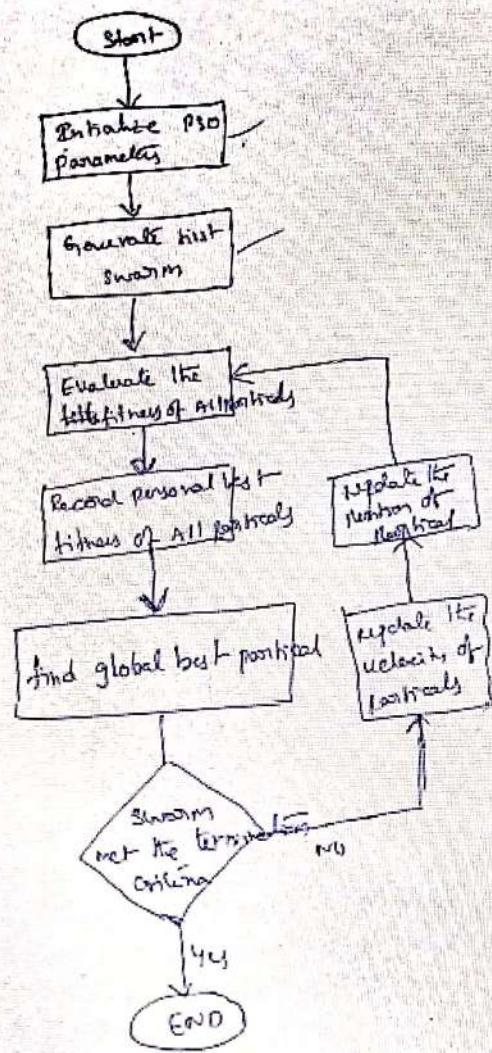
$w(t)$ is the weight at t^{th} iteration. The value for 'w' is adjusted at every iteration as given below, where 'iter' is total no. of iteration.

W.t.d.

$$w(t+1) = \frac{w(t) - t \times w(t)}{\text{iter}}$$

Ant colony optimization

PSO algorithm flow chart



Ant colony optimization technique exploit the natural behavior of ants in finding the shortest path to reach destination from the source.

Algorithm

Consider the problem of finding the opinion order in which the numbers from 1 to 8 are arranged so that the cost of the order is minimized. Cost of the particular particular order is computed using the two matrices A and B as described below

a ₁₁									
a ₁₂	a ₂₃								
a ₂₄	a ₃₅	a ₄₆							
a ₃₇	a ₈₈	a ₉₁₀							
a ₄₁₁	a ₁₂	a ₁₃	a ₁₄	a ₁₅					
a ₁₆	a ₁₇	a ₁₈	a ₁₉	a ₂₀	a ₂₁	a ₂₂			
a ₂₃	a ₂₄	a ₂₅	a ₂₆	a ₂₇	a ₂₈				
a ₁₉	a ₂₀	a ₂₁	a ₂₂	a ₂₃	a ₂₄	a ₂₅	a ₂₆	a ₂₇	a ₂₈

Matrix A

b ₁₁									
b ₁₂	b ₂₃								
b ₂₄	b ₃₅	b ₄₆							
b ₃₇	b ₈₈	b ₉₁₀							
b ₄₁₁	b ₁₂	b ₁₃	b ₁₄	b ₁₅					
b ₁₆	b ₁₇	b ₁₈	b ₁₉	b ₂₀	b ₂₁	b ₂₂			
b ₂₃	b ₂₄	b ₂₅	b ₂₆	b ₂₇	b ₂₈				
b ₁₉	b ₂₀	b ₂₁	b ₂₂	b ₂₃	b ₂₄	b ₂₅	b ₂₆	b ₂₇	b ₂₈

Matrix B

The cost of the order [3 2 1 5 7 2 4 6 8] is computed as the sum of the products of values obtained from the positions (1,1), (1,2), (2,1), (3,1), (4,1), (5,1), (6,1), (7,1), (8,1) of matrix A & B

$$(0.6 \times b_{11} + 0.3 \times b_{12} + 0.1 \times b_{13} + 0.6 \times b_{14}) + (0.1 \times b_{21} + 0.6 \times b_{22} + 0.3 \times b_{23} + 0.1 \times b_{24}) + \dots + (0.1 \times b_{81} + 0.6 \times b_{82} + 0.3 \times b_{83} + 0.1 \times b_{84})$$

The process of solving the Traveling Salesman Problem using Ant Colony Optimization technique is described below:

Step 1: generates five sets of orders randomly which are treated as the initial values selected by the 5 ants respectively.

The following are the values selected by the 5 ants along with the corresponding cost as given below:

ANT	ORDERS	COST
ANT 1	8 7 3 4 6 5 8	c1
ANT 2	4 2 8 5 3 1 6	c2
ANT 3	5 1 3 2 4 6 8	c3
ANT 4	3 4 1 8 6 5 2	c4
ANT 5	8 1 2 4 5 3 6	c5

Step 2: Pheromone is the value assigned for i^{th} number in the j^{th} position of the sequence selected by all the ants. Initially all the values of matrix are completely filled up with zero. The values of pheromone matrix are updated in every iteration.

Iteration matrix (initial) = pheromone matrix after initializing pheromone matrix.

To select the value for $(i, j)^{th}$ column of

1	0	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	0	0
2	$\frac{1}{4}$	0	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	0	$\frac{1}{4}$
3	$\frac{1}{4}$	0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	0	$\frac{1}{4}$
4	$\frac{1}{4}$	0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	0	$\frac{1}{4}$
5	$\frac{1}{4}$	0	0	0	0	$\frac{1}{4}$	$\frac{1}{4}$
6	0	0	0	0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
7	$\frac{1}{4}$	$\frac{1}{4}$	0	0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
8	$\frac{1}{4}$	0	$\frac{1}{2}$	0	$\frac{1}{4}$	0	$\frac{1}{4}$

Step 3: Next set of cities selected by the ants is as described below:

Eg: The random position of 8 2 3 6 4 15 (day) let the cities decided by the ant 1 be separated at (1, 1), (2, 2) and (3, 3) as $\boxed{(1, 2)}$ are the value for the set in the order as numbered above is selected not as the first number of the generated position sequence as 7.

* To select the value for a_{ij} , j^{th} column of the pheromone matrix is used as below,

$$\text{column } 7 = [0 \ 0 \ \frac{1}{c_2} \ 0 \ \frac{1}{c_4} \ \frac{1}{c_3} \ \frac{1}{c_5} \ 0]$$

* Now calculate the normalized value of column 7. This is also called as the probability vector and for selecting the values in a_{ij} in the order

$$S = [\frac{1}{c_2} + \frac{1}{c_4} + \frac{1}{c_3} + \frac{1}{c_5} + 1]^{-1} = \text{value}$$

$$\begin{aligned} \text{Normalized vector } 7 &= \frac{\text{column } 7}{S} = \left[\frac{0}{S} \ \frac{0}{S} \ \frac{\frac{1}{c_2}}{S} \ \frac{0}{S} \ \frac{\frac{1}{c_4}}{S} \ \frac{0}{S} \ \frac{\frac{1}{c_3}}{S} \ \frac{0}{S} \right] \\ &= [p_1 \ p_2 \ p_3 \ p_4 \ p_5 \ p_6 \ p_7 \ p_8] \xrightarrow{\frac{1}{c_2}} \left[\frac{0}{S} \ \frac{0}{S} \ \frac{1}{S} \ \frac{0}{S} \ \frac{1}{S} \ \frac{0}{S} \ \frac{0}{S} \right] \end{aligned}$$

* Find the lowest value position, for example

If p_5 is the lowest value then 5 is assigned to a_{ij} .

Note: position no. '5' is not used again.

* The next value in position sequence is 8, so column '8' is selected from pheromone matrix, and select the normalized value

$$\text{column } 8 = [0 \ \frac{1}{c_4} \ 0 \ 0 \ \frac{1}{c_2} + \frac{1}{c_3} \ 0 \ \frac{1}{c_4} + \frac{1}{c_2}]$$

$$8 = \left[\frac{1}{c_4} + \left(\frac{1}{c_2} + \frac{1}{c_3} \right) + \left(\frac{1}{c_4} + \frac{1}{c_2} \right) \right]$$

$$\text{Normalized vector } 8 = \frac{\text{column } 8}{S} = (a_1 \ a_2 \ a_3 \ a_4 \ a_5 \ a_6 \ a_7)$$

Find the lowest value in normalized vector '8' and replace with U8. Example if a_3 is lowest then number '3' is assigned to U8.

* This process is repeated to obtain the order selected by the Ant 1, Ant 2, Ant 3, Ant 4, Ant 5.

* This is called one Iteration.

Step 4: updating the pheromone matrix for the second iteration is completed as described in step 2. This is called updating pheromone matrix (2)

$$\therefore \text{pheromone matrix}(2) = \text{pheromone matrix}(1) + \text{updating pheromone matrix}$$

Step 5: Next set the orders selected by the five ants are computed as described in the step 3.

Step 6: steps 3 and to step 5 is repeated for the particular iterations. Thus best set of orders are obtained using Ant colony algorithm

Note: Best among the best set of orders selected by the five sets ants in the last iteration is selected as the global best order.

