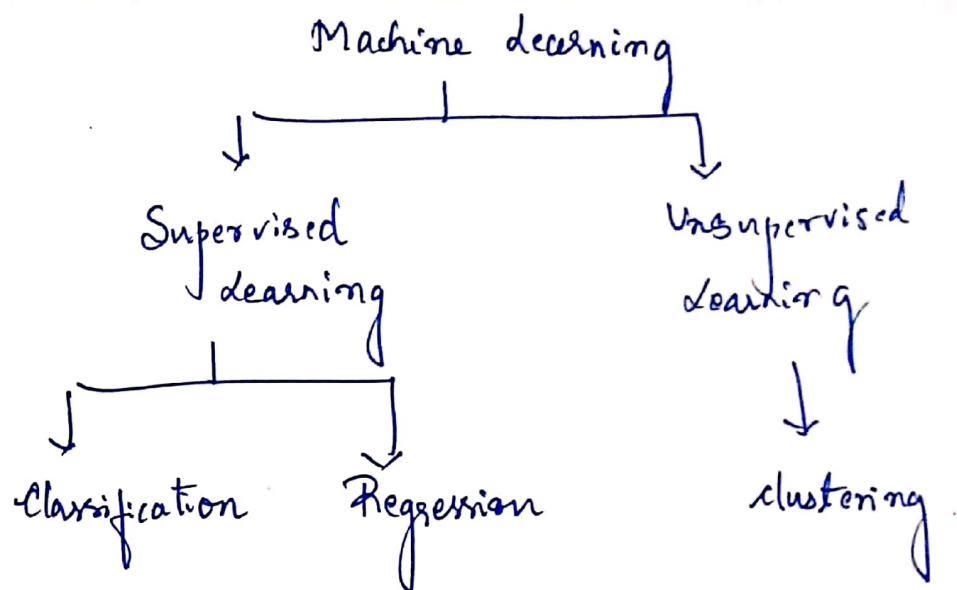


Computational Intelligence

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

* 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.

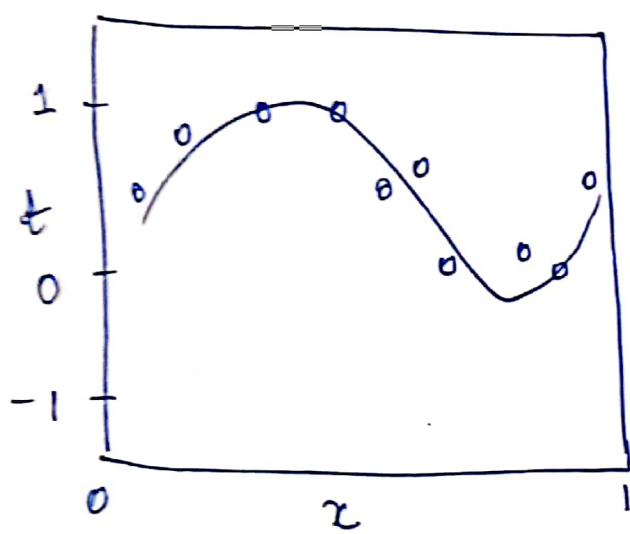
→ Applications in which training data comprises Examples of input vectors along with their corresponding target vectors are known as Supervised learning problems. If the desired output consists of one or more continuous variables, then task is called regression.

→ 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 trail 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

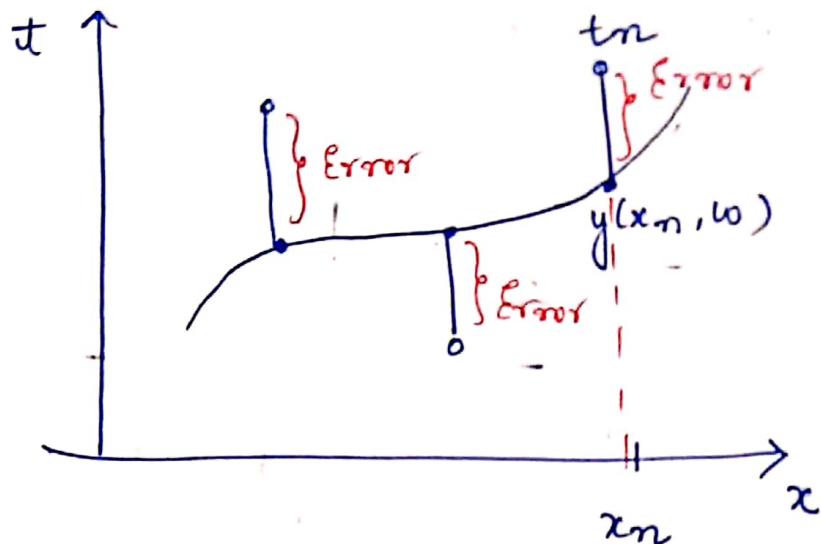
→ In fitting the polynomial curve to the training data we occur with Error or misfit ~~is~~ between

$y(x, w)$ 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)

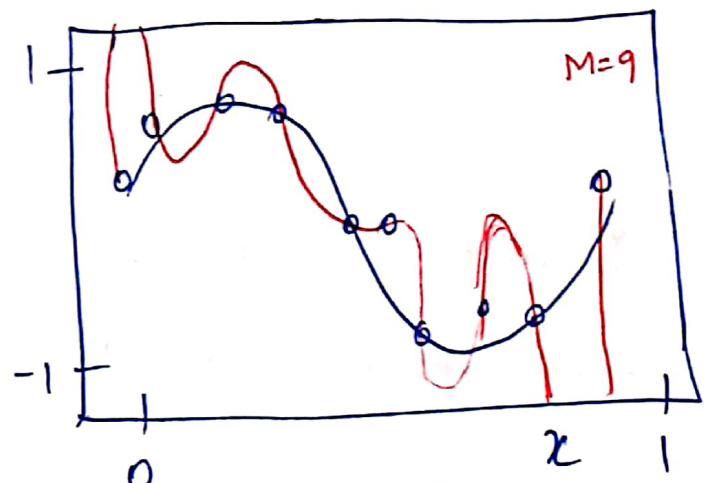
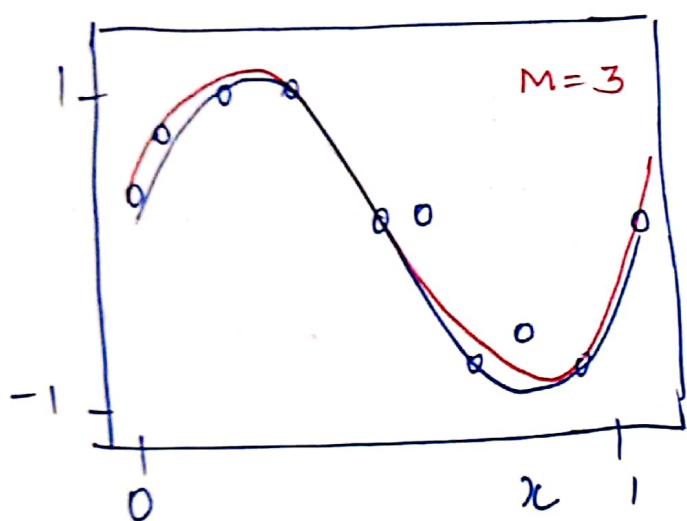
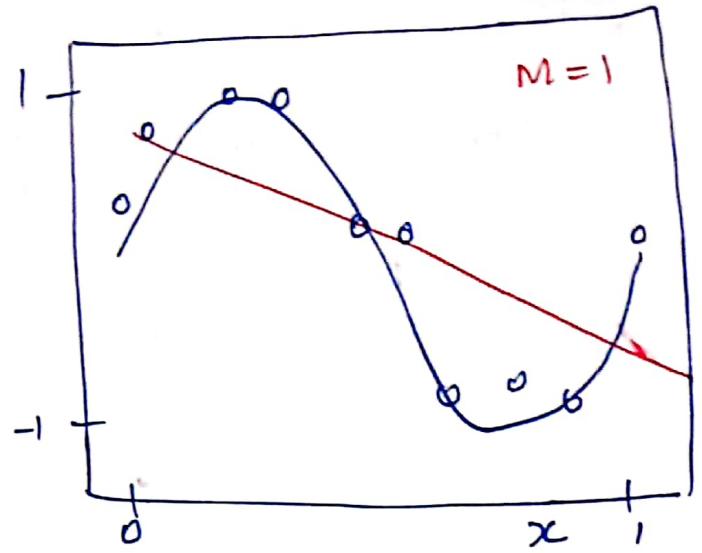
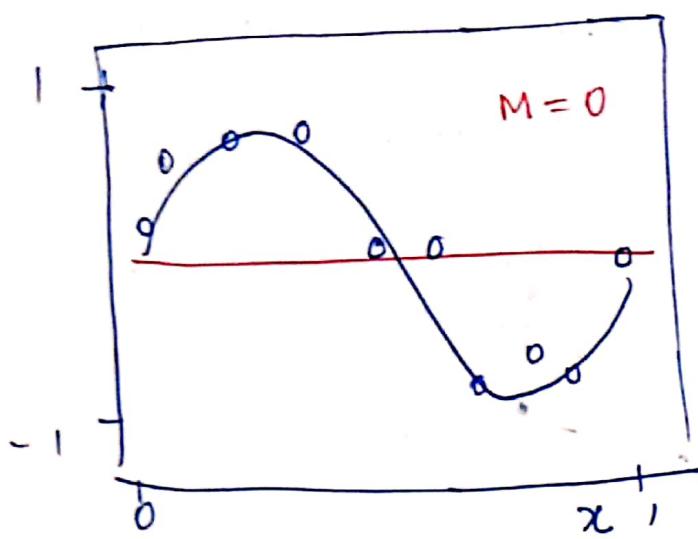


* Here we can solve the curve fitting problem by choosing the value of w . $E(w)$ is as small as possible.

* 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
 $\sin(2\pi x)$

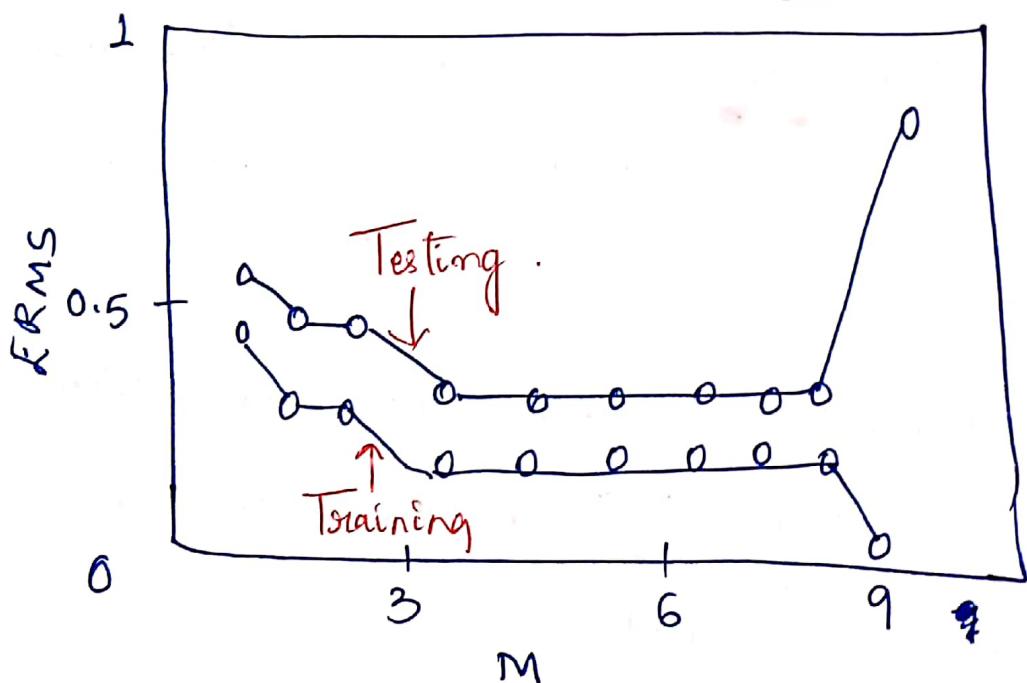
→ with $M=9$ we get Excellent fit to the
training data .

→ 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 $y(x, w^t)$ & 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.

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

$$\text{where } \|w\|^2 = w^T w = w_0^2 + w_1^2 + \dots + w_M^2 \quad \textcircled{11}$$

- the co-efficient 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 :-

$$\text{Sum rule : } p(x) = \sum_y p(x, y)$$

$$\text{Product rule : } p(x, y) = p(y|x) p(x)$$

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

$$P(x_{\text{ref}}) = \sum_y P(x|y) P(y) \quad (\because P(x,y) = P(y|x) P(x))$$

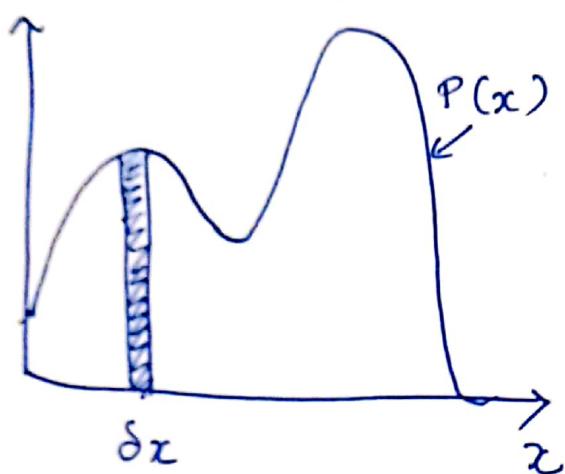
$$\therefore P(x) = \sum_y P(x|y) P(y)$$

↓
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

likelihood

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

(posterior Probability)

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.

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

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.

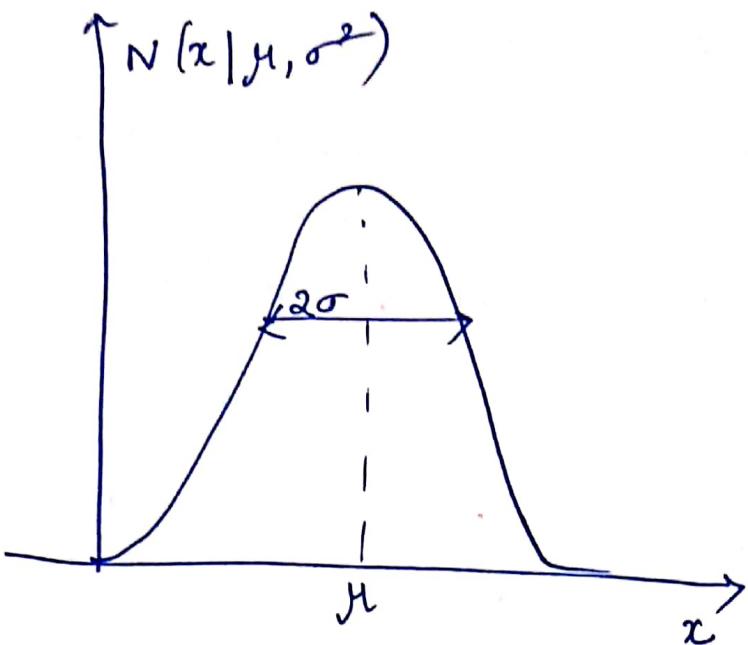
- In Baye's theorem , the likelihood function plays a central role .
- A widely used Estimator of w is maximum likelihood in which w is set to a value that maximizes likelihood function $P(t|w)$ $P(t|w)$.

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

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

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

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

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

w.r.t β

w.r.t β

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

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

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

By taking log for above expression we get

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

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

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

∴ 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

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

→ Now returning back to maximum likelihood

& least square concept

→ Let us consider target t with Gaussian noise model

$$\therefore t = y(x, w) + \epsilon \leftarrow \text{Gaussian noise}$$

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

(linear
Basis for
model)

→ 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

$$\begin{cases} \frac{1}{\sigma^2} = \beta & \frac{1}{\beta} = \sigma^2 \\ \beta^{-1} = \sigma^2 & \end{cases}$$

↓ Variance

* Here Expectation of new x to a target t

is given as

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

→ 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

$$P(w|t) = \frac{p(t|w) \cdot P(w)}{\text{Likelihood } P(t)}$$

→ 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

$$\begin{aligned} \ln p(t|w, \beta) &= \sum_{n=1}^N \ln N(t_n | w^T \phi(x_n), \beta^{-1}) \\ &= \left[\frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \right] - \beta E_D(w) \end{aligned}$$

↑ likelihood ↓ Error to likelihood

→ (1)

where where $E_D \rightarrow$ Sum of Square Error. for $\text{Eqn } 20$

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

→ 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 $E_D(w)$

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

~~if~~ $\left[\begin{array}{l} \text{Simply differentiate} \\ \text{by } w \end{array} \right]$

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

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

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

(21)

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

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

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

∴ 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

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

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

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

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

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

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

⊗

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

$\underbrace{\hspace{10em}}$
 $m \times m$

$\overbrace{\hspace{10em}}$
 $w_{ML} = (m \times 1) \rightarrow \text{Row vector}$ as a

→ 11th we can maximize the log likelihood function

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

w.r.t to precision parameter β .

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

Equating above equation to zero in order to minimize it

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

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

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

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

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

∴ 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 :-

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

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.

$$E = \sum_n E_n$$

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

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

(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.

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

\downarrow
regularization term

$E_D(w) \rightarrow$ data dependent error

$E_w(w) \rightarrow$ regularization term

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

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

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

∴ Total error fn is

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

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

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

→ 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

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

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

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

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

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

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

Earlier, the above expansion can also be written

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

\therefore maximizing above fm wrt w

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

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

\rightarrow If we examine result for target variable t_K single target

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

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 parameter 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

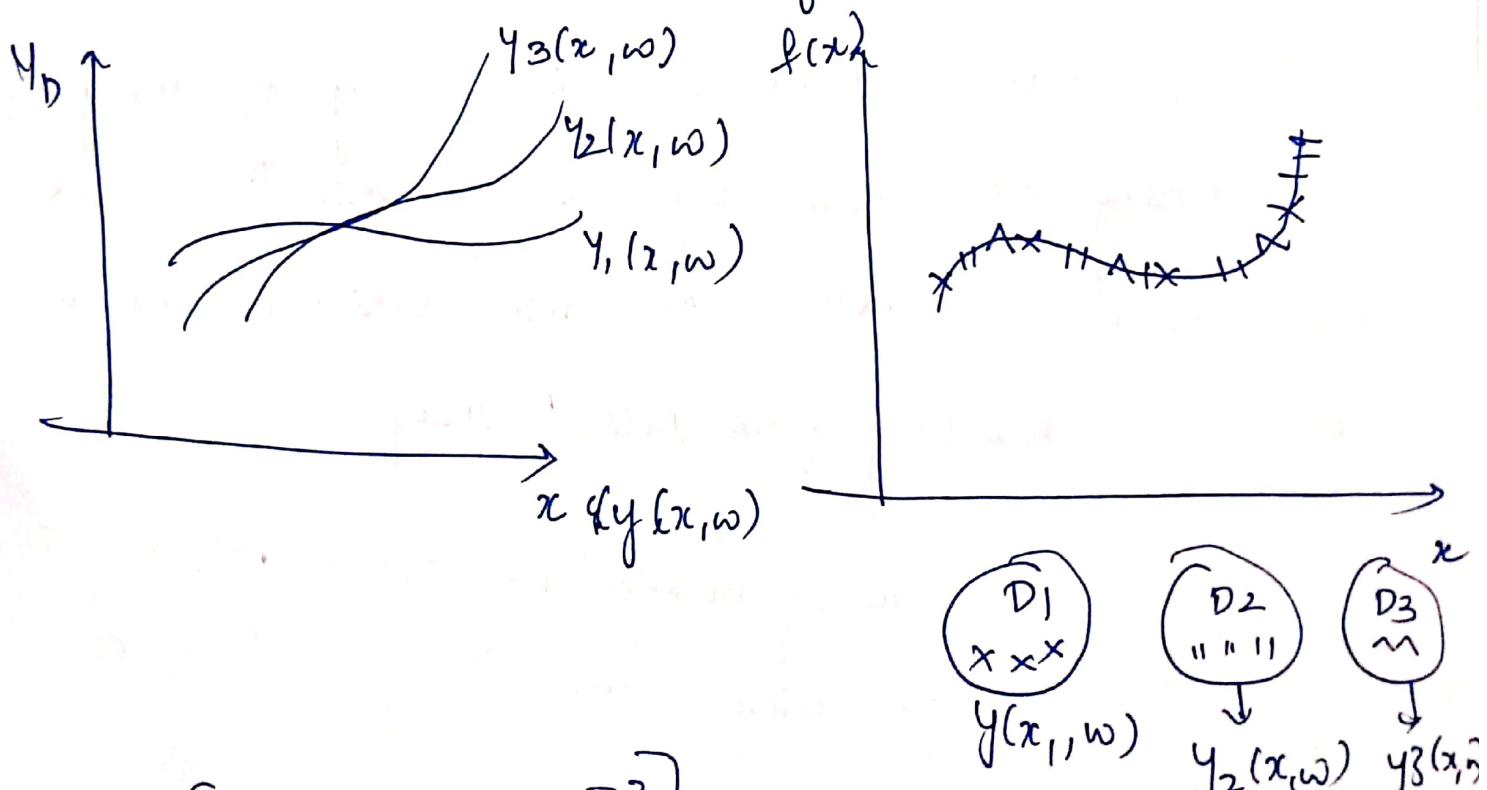
$$E [(t - g(x))^2]$$

→ 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

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

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

$$= E_D [P^2 + Q^2 + 2PQ]$$

$$= E_D \left[(y(x, D) - E_D[y(x, D)]) \right]^2 + E_D \left[E_D[y(x, D)] - g(x) \right]^2 \\ + 2 E_D \left[(y(x, D) - E_D[y(x, D)]) \left(E_D[y(x, D)] - g(x) \right) \right]$$

Variance Bias²

+ Noise.

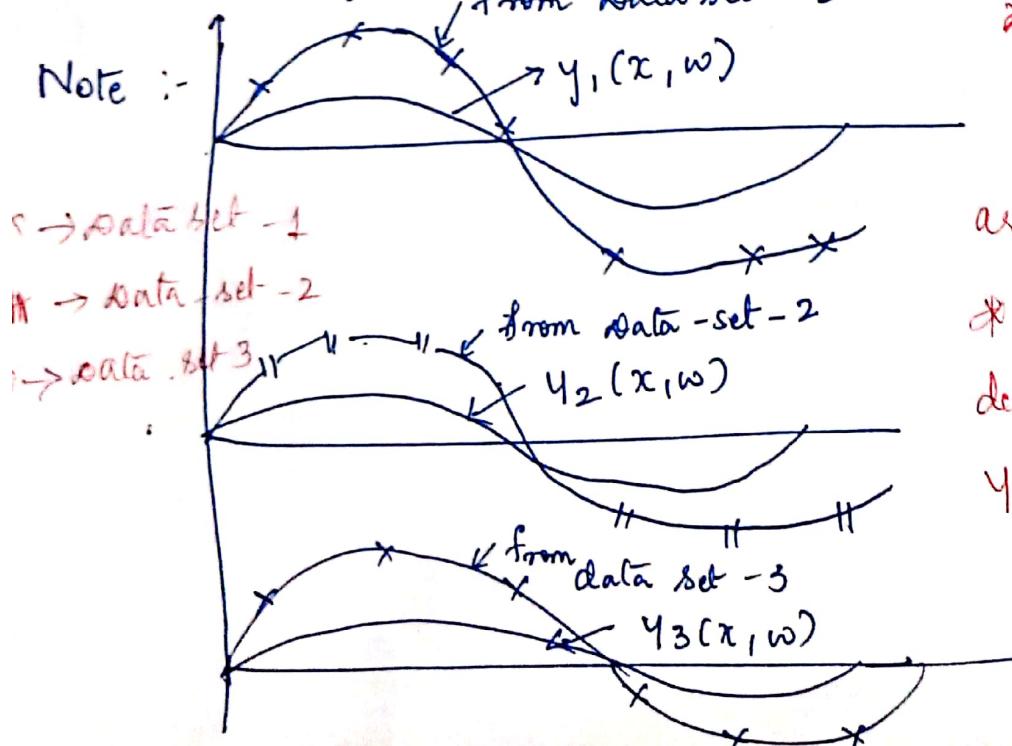
$E^{(\text{loss})} = \text{Variance} + \text{Bias}^2 + \text{Loss} \cdot \text{Noise}$

* $g(x)$ is from Expectation

$y(x)$ is from parametric approach.

from data set - 1.

Note :-



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

are different solutions

* 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.