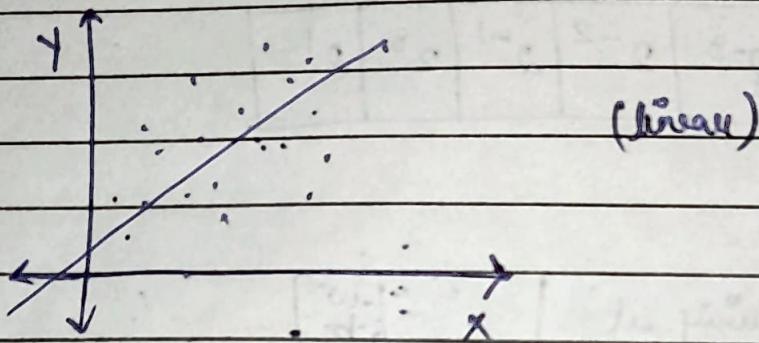


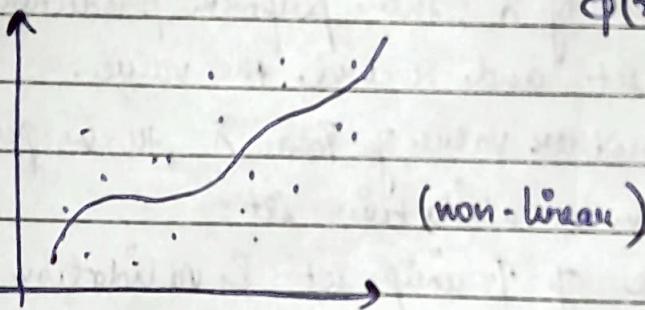
23/2/23

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$f(x) = w^T x + b, \quad x \in \mathbb{R}^n, \quad w \in \mathbb{R}^n, \quad b \in \mathbb{R}.$$



$$f(x) = w^T \phi(x) + b, \quad w \in \mathbb{R}^n$$

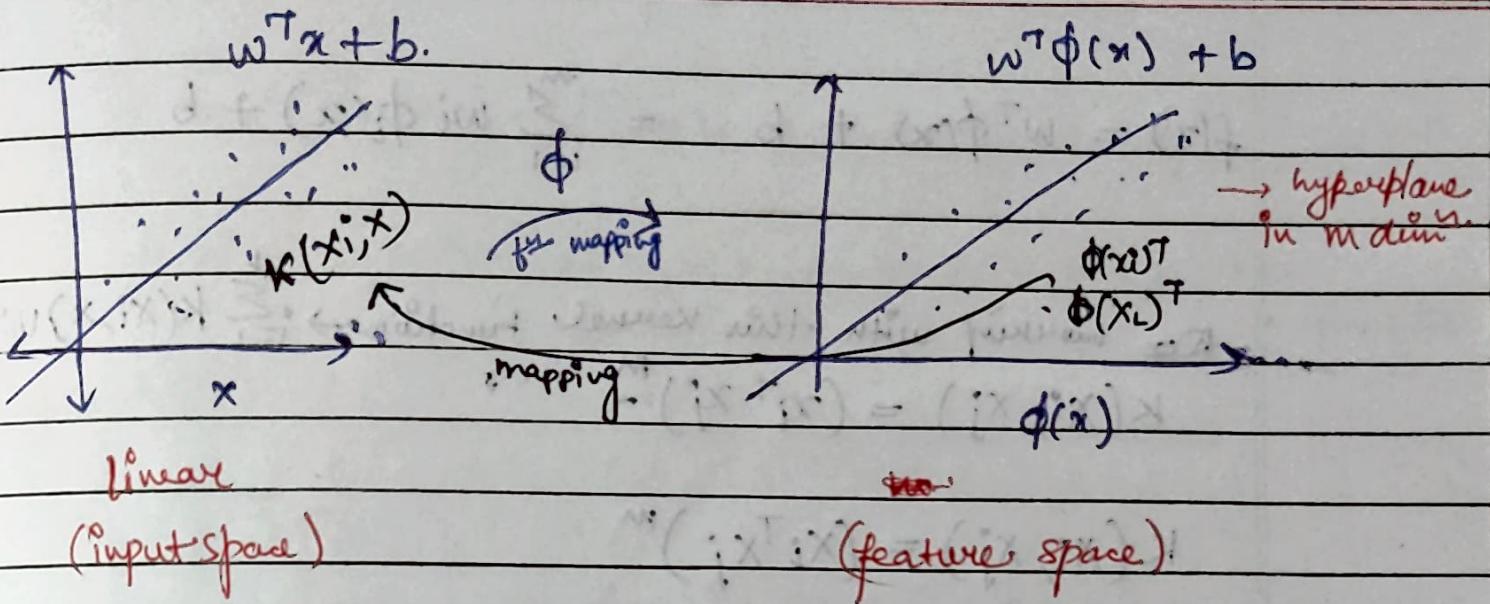


$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_m(x) \end{bmatrix}, \quad m \rightarrow \infty$$

\* mapping  $x$  to  $\phi(x)$  when we work with non-linear  $f(x)$ .  
So, at last, we are working with linear function  $(w^T \phi(x) + b)$ .

$w^T \phi(x) + b$  is a line in  $m$ -dimension. for  $\phi(x) = \phi_1(x), \dots, \phi_m(x)$

\* if the no. of pts in  $\phi(x)$ , we ke dimension hoga.



$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\}$$

$x_i \in \mathbb{R}^n$

$y_i \in \mathbb{R}$

$$w = \sum_{i=1}^L v_i \phi(x_i)$$

$$= v_1 \phi(x_1) + v_2 \phi(x_2) + \dots + v_L \phi(x_L)$$

\* In higher dimension space,  $w$  can be obtained by linear comb<sup>n</sup> of  $L$  training pts. in feature space.

$$\begin{aligned} & \rightarrow w^T \phi(x) + b \\ & \Rightarrow \left( \sum_{i=1}^L v_i \phi(x_i) \right)^T \phi(x) + b \\ & \Rightarrow \sum_{i=1}^L v_i \phi(x_i)^T \phi(x) + b \end{aligned}$$

### Kernel Generated Function

$$\sum v_i K(x_i, x) + b$$

↳ this will tell in higher dimension how  $\phi(x)$  looks like.

$$f(x) = w^T \phi(x) + b = \sum_{i=1}^m w_i \phi_i(x) + b$$

\* working with this Kernel function  $\rightarrow \sum_{i=1}^n k(x_i, x) v_i + b$

$$k(x_i, x_j) = (x_i^T x_j)^m$$

$$k(x_i, x_j) = (x_i^T x_j)^m$$

\* If we are working with Gaussian basis fn & have taken m no. of basis, then we use this type of Kernel fn.

$$k(x_i, x_j) = \exp \frac{-\|x_i - x_j\|^2}{\sigma^2}$$

$$= \phi(x_i)^T \phi(x_j)$$

→ Kernel fn. for Gaussian basis fn.

7/2/23

Fitting the linear function

$$f(x) = w_1 x_1 + w_2 x_2 + b \rightarrow \text{linear}$$

$$f(x) = w_1 x_1^2 + w_2 x_2^2 + w_3 x_1 x_2 + w_4 x_1 + w_5 x_2 + b$$

Kernel function

$x_1$	$x_2$	$\phi$	$x_1$	$x_2$	$x_3$
-	-		-	-	-
-	-		-	-	-
-	-		-	-	-

$$\phi(x) = \phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \end{bmatrix}$$

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

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$w_1 x_1^2 + w_2 x_2^2 + \sqrt{2} w_3 x_1 x_2 + b. \rightarrow \text{linear fn. in 3-D.}$$

$$w_1 x_1^2 + w_2 x_2^2 + \sqrt{2} w_3 x_1 x_2 + w_4 x_1 + w_5 x_2 + b \hookrightarrow \text{Linear fn. in 2D}$$

$$\phi(x)^T \phi(y) = ?$$

$$\begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \\ x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} y_1^2 \\ y_2^2 \\ \sqrt{2} y_1 y_2 \\ y_1 \\ y_2 \end{bmatrix} = (x_1 y_1)^2 + (x_2 y_2)^2 + 2 x_1 x_2 y_1 y_2 + x_1 y_1 + x_2 y_2 = (x^T y)^2 + (x^T y)$$

$$\therefore \phi(x)^T \phi(y) = (x^T y)^2 + (x^T y)$$

$$K(x, y) = (x^T y)^2 + x^T y \rightarrow \text{Kernel Generated fn.}$$

$$\sum_{i=1}^5 K(x_i^o, x) u_i^o + b$$

$x_1$	$x_2$
5.9	3
6.9	3.1
6.6	2.9
4.6	3.2
6	2.2

$$\sum_{i=1}^5 K(x_i^o, x) u_i^o + b \quad u_i^o = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ b \end{bmatrix}$$

$$K\left(\begin{bmatrix} 5.9 \\ 3 \end{bmatrix}, \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = ((5.9 x_1 + 3 x_2)^2 + (5.9 x_1 + 3 x_2) u_1 + b$$

$u_i \rightarrow$  we need to find.  $\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ b \end{bmatrix}$

$b \rightarrow$  bias term

$$K(x_1, x_1) = K\left(\begin{bmatrix} 5 & 9 \\ 3 & 7 \end{bmatrix}; \begin{bmatrix} 5 & 9 \\ 3 & 7 \end{bmatrix}\right)$$

$$= [5 \cdot 9 \ 3] \begin{bmatrix} 5 & 9 \\ 3 & 7 \end{bmatrix} \Rightarrow \therefore K(x_1, x_2) = \phi(x_1)^T \phi(x_2)$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

\* For every kernel function, we have a Kernel matrix.

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\} \quad x_i \in \mathbb{R}^n \\ y_i \in \mathbb{R}$$

Kernel Matrix:

$$K = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_L) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_L) \\ \vdots & \vdots & \ddots & \vdots \\ K(x_L, x_1) & K(x_L, x_2) & \dots & K(x_L, x_L) \end{bmatrix}$$

- always +ve semi definite
- symmetric matrix

→ Quadratic Kernel.

$$K(x, y) = (c + x^T y)^2$$

$c \rightarrow$  constant (user defined)

$$K(x, y) = (c + x^T y)^m \rightarrow \text{feature polynomial of order } m.$$

$$K(x, y) = \exp^{-\frac{\|x - y\|^2}{2\sigma^2}}$$

$\sigma^2 \rightarrow$  user defined

This Kernel function can be used to learn any type of function in any dimension.

→ isome map  $K(x, y) \rightarrow$  feature space.

$$\min_{(w,b)} \frac{\lambda}{2} w^T w + \frac{1}{L} \sum_{i=1}^L (y_i - (w^T x_i + b))^2$$

$$\underbrace{\begin{bmatrix} w \\ b \end{bmatrix}}_{\text{solution}} = (A^T A + \lambda I)^{-1} A^T Y$$

$$A = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_L) & 1 \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_L) & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ K(x_L, x_1) & K(x_L, x_2) & \dots & K(x_L, x_L) & 1 \end{bmatrix}$$

16/23

Given the training set,

$$T = \{(x_i^*, y_i^*) : x_i \in \mathbb{R}^n, y_i \in \mathbb{R}\}$$

We are solving

$$\min_{w,b} \frac{\lambda}{2} w^T w + \frac{1}{L} \sum_{i=1}^L (y_i^* - (w^T x_i + b))^2$$

Gradient Descent Method

$$\frac{\partial}{\partial w} J(w, b) = \nabla_w J(w, b) = \lambda w - \frac{2}{L} \sum_{i=1}^L x_i^* (y_i^* - (w^T x_i + b))$$

$$\frac{\partial}{\partial b} J(w, b) = -\frac{2}{L} \sum_{i=1}^L (y_i^* - (w^T x_i + b))$$

$$\begin{bmatrix} w^0 \\ b^0 \end{bmatrix} \in \mathbb{R}^{n+1} \rightarrow \text{Initialise}$$

Repeat

$$\begin{bmatrix} w^{k+1} \\ b^{k+1} \end{bmatrix} = \begin{bmatrix} w^k \\ b^k \end{bmatrix} - \eta \underbrace{\begin{bmatrix} \nabla_w J(w^k, b^k) \\ \nabla_b J(w^k, b^k) \end{bmatrix}}_{\text{Gradient}}$$

$$f(x) = w^T x + b$$

WEIR  
b ∈ ℝ.

$$\min_{(w,b)} J(w,b)$$

$$\cdot \frac{d(x^T A x)}{dx} = 2Ax$$

$$\frac{d(w^T w)}{dw} = \frac{\partial w^T w}{\partial w} = \frac{\partial w^T}{\partial w} w = 2w$$

Until

$$\left\| \begin{bmatrix} \nabla_w J(w^k, b^k) \\ \nabla_b J(w^k, b^k) \end{bmatrix} \right\|_2 \leq \epsilon.$$

→ \* For any basis function :-

$$\min_{w, b} \frac{\lambda}{2} w^T w + \sum_{i=1}^L (\gamma_i - (w^T \phi(x_i) + b))^2 = \min_{(w, b)} J(w, b)$$

$$\nabla_w J(w, b) = \lambda w - \frac{2}{L} \sum_{i=1}^L (\gamma_i - (w^T \phi(x_i) + b)) \phi(x_i)$$

$$\nabla_b J(w, b) = - \frac{2}{L} \sum_{i=1}^L (\gamma_i - (w^T \phi(x_i) + b))$$

\* For Kernel generated function :-

$$\min_{u, b} \frac{\lambda}{2} u^T u + \frac{1}{L} \sum_{i=1}^L (\gamma_i - \left( \sum_{j=1}^L K(x_j, x_i) u_j + b \right))^2$$

$$f(x) = \sum_{j=1}^L K(x_j, x) u_j + b$$

→ ←  $\nabla_u J(u, b) = \lambda u - \frac{2}{L} \sum_{i=1}^L (\gamma_i - \left( \sum_{j=1}^L K(x_j, x_i) u_j + b \right)) \sum_{j=1}^L K(x_j, x_i)$

→ ←  $\nabla_b J(u, b) = - \frac{2}{L} \sum_{i=1}^L (\gamma_i - \left( \sum_{j=1}^L K(x_j, x_i) u_j + b \right))$

$$W = \begin{bmatrix} w \\ b \end{bmatrix} =$$

$$k = [ ]_{100 \times 100}$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$K(u, v) = (u^T v)^2$$

$\uparrow$  kernel fn

	$x_1$	$x_2$	$y$
$x_1$	0	0	0
$x_2$	0	1	0
$x_3$	1	1	0
$x_4$	1	0	1

$$K(x_1, x_2) = (x_1^T x_2)^2$$

write Gradient Descent Algo.

→ Kernel matrix,  $K = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_L) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_L) \\ \vdots & \vdots & \ddots & \vdots \\ K(x_L, x_1) & K(x_L, x_2) & \dots & K(x_L, x_L) \end{bmatrix}$

$$K = \begin{bmatrix} K(0, 0) & 0 & K(0, 1) & 0 & K(0, 2) & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix}_{4 \times 4}$$

$$y_1 = \left( \sum_{j=1}^4 K(x_j, x_1) u_j + b \right) \stackrel{?}{=} K(x_1, x_1)$$

$$(x_1^T x_1)^2 / \sqrt{(x_2^T x_2)^2 / \sqrt{(x_3^T x_3)^2 / \sqrt{(x_4^T x_4)^2}}}$$

$$y_2 = \left( \sum_{j=1}^4 K(x_j, x_2) u_j + b \right)$$

$$x_1^T = [0 \ 0]$$

$$x_2^T = [0 \ 1]$$

$$x_3^T = [1 \ 1]$$

$$x_4^T = [1 \ 0]$$

$$0 - \left( (K(x_1, x_1)u_1 + K(x_2, x_1)u_2 + K(x_3, x_1)u_3 + K(x_4, x_1)u_4) + b \right) (K(x_1, x_1) \cdot K(x_2, x_1) \cdot K(x_3, x_1) \cdot K(x_4, x_1))$$

$$(x_2^T x_1)^2 = [0 \ 1] \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0^2 = 0$$

$$= 0 - [0 + 0 + 0 + 0 + b](0)$$

$$(x_2^T x_2)^2 = [0 \ 1] \begin{bmatrix} 0 \\ 1 \end{bmatrix} = 1^2 = 1$$

$$= 0$$

13/3/23

## Least Square Kernel Regression.

$$\min_{(u,b)} \frac{\lambda}{2} u^T u + \sum_{i=1}^l \left( y_i - \left( \sum_{j=1}^l K(x_j, x_i) u_j + b \right) \right)^2$$

$$\begin{bmatrix} u \\ b \end{bmatrix} = (A^T A)^{-1} A^T Y$$

$$A = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_l) & 1 \\ \vdots & \ddots & & & \\ K(x_l, x_1) & K(x_l, x_2) & \dots & K(x_l, x_l) & 1 \end{bmatrix}$$

first let us consider

$$J(u, b, x_k, y_k) = \left( y_k - \left( \sum_{j=1}^l K(x_j, x_k) u_j + b \right) \right)^2$$

$$\frac{\partial J_u(u, b, x_k, y_k)}{\partial u} = -2 \left( y_k - \left( \sum_{j=1}^l K(x_j, x_k) u_j + b \right) \right) \sum_{j=1}^l K(x_j, x_k)$$

Differentiating this term w.r.t.  $u$ .

$$\sum_{j=1}^l K(x_j, x_k) u_j = K(x_1, x_k) u_1 + K(x_2, x_k) u_2 + \dots + K(x_l, x_k) u_l + b$$

$$\begin{bmatrix} K(x_1, x_k) \\ K(x_2, x_k) \\ \vdots \\ K(x_l, x_k) \end{bmatrix}$$

$$\frac{\partial J(u, b, x_k, y_k)}{\partial u} = -2 \left( y_k - \left( \sum_{j=1}^l K(x_j, x_k) u_j + b \right) \right) \begin{bmatrix} K(x_1, x_k) \\ K(x_2, x_k) \\ \vdots \\ K(x_l, x_k) \end{bmatrix}$$

now w.r.t  $b$ ,

$$\frac{\partial J(u, b, x_k, y_k)}{\partial b} = -2 \left( y_k - \left( \sum_{j=1}^l K(x_j, x_k) u_j + b \right) \right)$$

## Gradient Descent least square Kernel Regression

Algorithm : gradient descent method :-

Initialize  $u^0 = u^{start} \in R^l$  and  $b \in R$

Repeat

$$u^{(j+1)} = u^{(j)} - \eta_k \left( \lambda u + \sum_{i=1}^l \frac{\partial J(u, b, x_k, y_k)}{\partial u} \right)$$

$\eta_k = \text{step size}$

$$b^{(j+1)} = b^{(j)} - \eta_k \left( \sum_{i=1}^l \frac{\partial J(u, b, x_k, y_k)}{\partial b} \right)$$

Until  $\left\| \lambda u + \sum_{i=1}^l \frac{\partial J(u, b, x_k, y_k)}{\partial u} \right\| \leq \epsilon$

## Stochastic Gradient Descent Least Square Regression

$$\min_{(w,b)} \frac{\lambda}{2} \cdot w^T w + \sum_{i=1}^l (y_i - (w^T x_i + b))^2$$

$$\frac{\partial}{\partial w} = \lambda w - 2 \sum_{i=1}^l (y_i - (w^T x_i + b)) x_i$$

$$J(w, b, x_k, y_k) = (y_k - (w^T x_k + b))^2$$

$$\nabla J(w, b, x_k, y_k) = -2(y_k - (w^T x_k + b)) x_k$$

$$\frac{\partial J(w, b, x_k, y_k)}{\partial b} = -2(y_k - (w^T x_k + b))$$

### Algorithm: Stochastic Gradient Descent Method

\* At every iteration, value of  $\eta_k$  should be changed.

Initialise  $w^0 = w^{start} \in R^l$  and  $b^0 \in R$

Repeat

Randomly select subset  $B$  from Training set  $T$ .

$$w^{(j+1)} = w^{(j)} - \eta_k \left( \lambda w_j + \sum_{(x_k, y_k) \in B} \frac{\partial J(w^{(j)}, b^{(j)}, x_k, y_k)}{\partial w} \right)$$

$$b^{(j+1)} = b^{(j)} - \eta_k \left( \sum_{(x_k, y_k) \in B} \frac{\partial J(w^{(j)}, b^{(j)}, x_k, y_k)}{\partial b} \right)$$

Until  $\| \begin{bmatrix} w^{(j+1)} \\ b^{(j+1)} \end{bmatrix} - \begin{bmatrix} w^{(j)} \\ b^{(j)} \end{bmatrix} \| \leq \epsilon$

# One of the method to reduce the value of  $\eta$ .

$$\eta = 0.2$$

for  $i = 1; 1000$

$$\eta_i = \frac{\eta}{(1 + 0.0008)}$$

and.

# other methods to reduce the value of  $\eta$  (eta)

- Line search method
- Momentum concept

Stochastic Gradient Descent, Kernel Method Regression  
Least Square

Algorithm: Stochastic gradient descent method

Initialize  $x^0 = u^{start} \in R^l$  and  $b \in R$

repeat

Randomly select subset  $B$  from Training set  $T$ .

$$u^{(j+1)} = u^{(j)} - \eta_k \left( \lambda u + \sum_{(x_k, y_k) \in B} \frac{\partial J(u, b, x_k, y_k)}{\partial u} \right)$$

$$b^{(j+1)} = b^{(j)} - \eta_k \left( \sum_{(x_k, y_k) \in B} \frac{\partial J(u, b, x_k, y_k)}{\partial b} \right)$$

Until  $\| \begin{bmatrix} u^{(j+1)} \\ b^{(j+1)} \end{bmatrix} - \begin{bmatrix} u^{(j)} \\ b^{(j)} \end{bmatrix} \| \leq \epsilon$

$$(x_1, y_1) \quad (x_2, y_2) \quad \dots \quad (x_k, y_k)$$

$$f(x_1) \quad f(x_2) \quad \dots \quad f(x_k)$$

- RMSE =  $\sqrt{\frac{1}{K} \sum_{i=1}^k (y_i - f(x_i))^2}$

- MAE =  $\frac{1}{K} \sum_{i=1}^k |y_i - f(x_i)|$

- NMSE =  $\frac{SSE}{SST} = \frac{\sum_{i=1}^k (y_i - f(x_i))^2}{\sum_{i=1}^k (y_i - \bar{y})^2}$  → Normalised Mean Sq. Error

SST → Sum of sq. Training Deviation

Variance of predicted  $f(\bar{x})$

$$\frac{1}{K^2} \left( \sum_{i=1}^k f(x_i) - f(\bar{x}) \right)^2$$

- $R^2 = \frac{\sum_{i=1}^k (y_i - \bar{y})^2}{\sum_{i=1}^k (f(x_i) - f(\bar{x}))^2}$

- Sparsity =  $\frac{\#(w_i = 0)}{\#(w_i)}$   $\Rightarrow \frac{\text{no. of } w_i = 0}{\text{no. of } w_i}$

Sparsity Calculates how many  $w_i = 0$ , ie it calculates simplicity of the model and how many total  $w_i$ .

def,

$$w_1 = 2.8$$

$$w_2 = 0$$

$$w_3 = 0$$

$$b = 0.4$$

$$\text{sparsity} = \frac{2}{3} \quad (\text{here}) = \frac{\#\{w_i=0\}}{\#\{w_i\}} \quad \begin{matrix} \# w_i = 0 \\ \# \text{total } w_i = 3 \end{matrix}$$

\* Sparsity should be high is preferable.

2 if  $w_1 = 2.8$

$$w_2 = 1.2$$

$$w_3 = 0.8$$

$$b = 0.4$$

then sparsity = 0, because no  $w_i = 0$ .

15/3/23

## MLE

if  $\{X_1, X_2, \dots, X_n\} \sim D(\theta)$   
(any distnb.)

$\{39.5, 68, 78, 105, 101, 82, \dots\} \sim N(\mu, \sigma^2)$

we select parameter  $\theta$ , for which has max<sup>m</sup> prob. given sample D.

$$\text{i.e., } \max_{\theta} P(D|\theta)$$

$$= \max_{\theta} P(D|\theta) \cdot P(\theta)$$

Likelihood expect value (prior info)

sometimes not given,

so we maximise only  $P(D|\theta)$

$$= \max_{\theta} P(D|\theta)$$

pts. in D are i.i.d.

$$\max_{\theta} P(D|\theta) = \max_{\theta} P(X_1, X_2, \dots, X_n|\theta)$$

$$\max_{\theta} P(D|\theta) = \max_{\theta} \prod_{i=1}^n P(X_i|\theta) = \max_{\theta} \sum_{i=1}^n \log P(X_i|\theta)$$

\* if  $\{X_1, X_2, \dots\} \sim N(\mu, \sigma^2)$

$\bar{\sigma}$  is any const. then estimate  $\mu$ .

$$\max_{\theta=\mu} \sum_{i=1}^n P(X_i|\mu, \bar{\sigma}) \rightarrow \text{in case of normal distribution,}$$

$$J(u) = \max_{\mu} \sum_{i=1}^n \log \frac{1}{\bar{\sigma} \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{(X_i - \mu)}{\bar{\sigma}} \right)^2}$$

$$J(u) = \max_{\mu} \sum_{i=1}^n \log \left( \frac{1}{\bar{\sigma} \sqrt{2\pi}} \right) - \frac{1}{2} \sum_{i=1}^n \left( \frac{(X_i - \mu)}{\bar{\sigma}} \right)^2$$

$$\text{derivative} = 0, \text{ i.e. } \frac{\partial J(u)}{\partial u} = 0$$

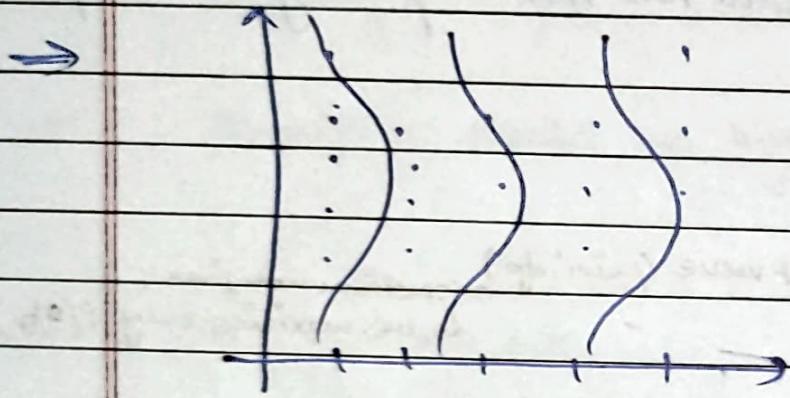
$$\frac{\partial J(u)}{\partial \mu} \frac{2}{2} \sum_{i=1}^n (X_i - \mu) = 0$$

$$\mu = \frac{1}{n} \sum x_i$$

$$\boxed{\mu = \bar{x}}$$

\* If we take  $\mu$  as const. & want to estimate  $\sigma$  & maximize  $\sigma$ , then following same step we get as,

then 
$$\hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$



$y/x$ .

$$y_i^o = f(x_i) + \varepsilon_i^o$$

$$\mathbb{E}(y_i^o | x) \rightarrow w^T x_i + b$$

$$w^T \phi(x_i) + b \quad (\text{for basis } \phi)$$

$$y_i^o | x \sim N(w^T \phi(x_i) + b, \sigma)$$

$$\varepsilon_i^o = y_i - (w^T x_i + b) \sim N(0, \sigma)$$

In regression problem, we have,

$$T = \left\{ (x_1, y_1), (x_2, y_2), \dots, (x_L, y_L) \right\} \\ x_i \in \mathbb{R}^n, y_i \in \mathbb{R}$$

$$y_i/x_i \sim N(w^T \phi(x_i), \sigma)$$

$$\max P(Y_1, Y_2, \dots, Y_L | X_1, X_2, \dots, X_L)$$

$$\max P(Y_1/X_1, Y_2/X_2, \dots, Y_L/X_L)$$

$$\max \prod_{i=1}^L P(Y_i/x_i)$$

$$= \max \prod_{i=1}^L \frac{1}{\sqrt{2\pi}} e^{-\frac{(Y_i - w^T \phi(x_i))^2}{2\sigma^2}}$$

$$\because Y_i/x \sim N(w^T \phi(x_i), \sigma)$$

$$\approx \max \sum_{i=1}^L \log \left( \frac{1}{\sqrt{2\pi}} e^{-\frac{(Y_i - w^T \phi(x_i))^2}{2\sigma^2}} \right)$$

$$= \max \sum_{i=1}^L \log \left( \frac{1}{\sqrt{2\pi}} \right) - \frac{(Y_i - w^T \phi(x_i))^2}{2\sigma^2}$$

$$\max - \sum_{i=1}^L (Y_i - w^T \phi(x_i))^2$$

$$\min_w \sum_{i=1}^L (Y_i - w^T \phi(x_i))^2$$

→ ~~dead~~ Squared Loss fn.

### Advantages :

This is sq. fn. of  $w$  &  $b$ ,

smooth fn., solving system of equation is easy.

For solving  $Y/X \sim N()$ , then, only this loss fn. will work.

### Disadvantages

It is sensitive to outliers.



To solve this, we need to solve or go for weighted least square regression model.

20/3/23

# Bias - Variance

$$\mathcal{T} = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\}$$

$$y \rightarrow f(x) \rightarrow x.$$

$$\min_f \frac{1}{L} \sum_{i=1}^L (y_i - f(x_i))^2 \rightarrow \text{Empirical risk}$$

We should minimise below one instead of above because, only  $L$  data points are considered above.

$$\iint (y_i - f(x_i)) P(x, y) d(x, y)$$

$$\underbrace{\left[ E_{\mathcal{T}}(f_T(x_i)) - \underbrace{f_0(x_i)}_{\substack{\text{Training} \\ \text{set}}} \right]}_{\substack{\text{Bias} \\ \text{actual} \\ \text{estimate}}}^2$$

$$\text{Var}(f_T(x))$$

→ When we increase the degree of polynomial / we use more complex function.

- Bias decreasing
- Variance increasing

$$\rightarrow y_i = f_0(x_i) + \varepsilon_i$$

Assumptions:

$$(i) E(\varepsilon_i) = 0$$

(ii)  $x_i$ 's are iid.

$$\text{minimise } E_T(y_i - f_T(x_i))^2$$

can be written as,

$$\min E_T(y_i - f_0(x_i) + f_0(x_i) - f_T(x_i))^2$$

$$\Rightarrow E_T(y_i - f_0(x_i))^2 + E_T(f_0(x_i) - f_T(x_i))^2 + 2 E_T(y_i - f_0(x_i))(f_0(x_i) - f_T(x_i))$$

Expanding this term below

$$= E_T[(y_i - f_0(x_i))(f_0(x_i) - f_T(x_i))]$$

$$= E_T(y_i f_0(x_i)) - E_T(f_0(x_i) f_0(x_i)) - E_T(y_i f_T(x_i)) + E_T(f_0(x_i) f_T(x_i))$$

~~$E_T(y_i f_0(x_i))$~~  Writing  $y_i = f_0(x_i) + \varepsilon_i$ , we get,

$$= E_T[(f_0(x_i) + \varepsilon_i) f_0(x_i)] - E_T[f_0(x_i) f_0(x_i)] - E_T[(f_0(x_i) + \varepsilon_i) f_T(x_i)] + E_T[f_0(x_i) f_T(x_i)]$$

$$= E_T(f_0(x_i) f_0(x_i)) + E_T(\varepsilon_i f_0(x_i)) - E_T(f_0(x_i) f_0(x_i)) - E_T(f_0(x_i) f_T(x_i)) - E_T(\varepsilon_i f_T(x_i)) + E_T(f_0(x_i) f_T(x_i))$$

$$= E_T(\varepsilon_i f_0(x_i)) - E_T(\varepsilon_i f_T(x_i))$$

$$\because E(\varepsilon_i) = 0$$

$$\therefore = 0$$

Hence,  $\min E_T(y_i - f_T(x_i))^2 =$

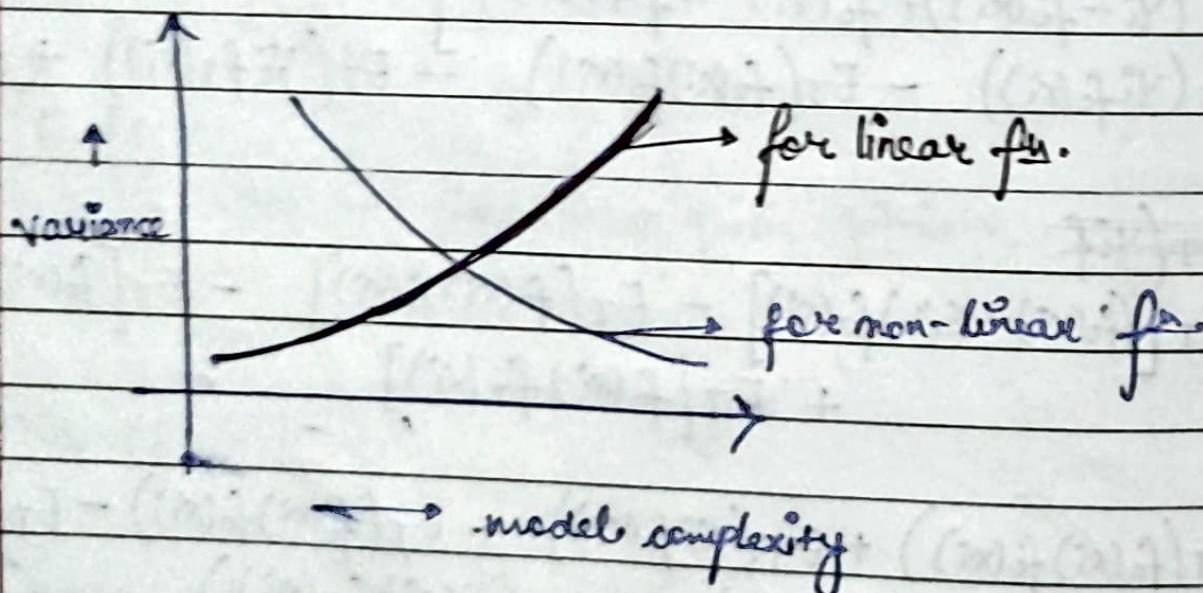
~~$E_T(y_i - f_0(x_i))^2$~~   $+ E_T(f_0(x_i) - f_T(x_i))^2 + 0$ .

$$= E(\varepsilon_i^2) + E_T(f_0(x_i) - f_T(x_i))^2$$

$$\begin{aligned}
 \text{diff } \Sigma &= f_i(x_i) - f_T(x_i) \\
 \mathbb{E}(\Sigma^2) &= \text{Var}(\Sigma) + [\mathbb{E}(\Sigma)]^2 \quad (\because \mathbb{V}(z) = \mathbb{E}(z^2) - (\mathbb{E}(z))^2) \\
 &= \mathbb{E}(S^2) + \text{Var}_{\mathbb{T}}[f_i(x_i) - f_T(x_i)] + \mathbb{E}_{\mathbb{T}}[(f_i(x_i) - f_T(x_i))^2] \\
 &= \mathbb{E}(S^2) + \mathbb{V}_{\mathbb{T}}(f_T(x_i)) + \mathbb{E}_{\mathbb{T}}[f_i(x_i) - f_T(x_i)]^2 \\
 &\quad (\because \text{Var}(kz) = k^2 \text{Var}(z))
 \end{aligned}$$

$$\mathbb{E}_{\mathbb{T}}[(y_i - f_T(x_i))^2] = [\mathbb{E}_{\mathbb{T}}(f_i(x_i) - f_T(x_i))]^2 + \text{Var}_{\mathbb{T}}(f_T(x_i)) + \mathbb{E}(S^2)$$

$$\Rightarrow \mathbb{E}_{\mathbb{T}}(y_i - f_i(x_i))^2 = \underbrace{[\mathbb{E}_{\mathbb{T}}(f_i(x_i) - f_T(x_i))]^2}_{\text{Bias}} + \underbrace{\mathbb{V}(f_T(x_i))}_{\text{variance}} + \underbrace{\mathbb{E}(S^2)}_{\text{irreducible error}}$$



- for linear  $\rightarrow$  as <sup>model</sup> complexity  $\uparrow$ , variance  $\uparrow$
- for non-linear  $\rightarrow$  as complexity  $\uparrow$ , variance  $\downarrow$

02/3/23

$$\frac{(m+n)!}{m! n!}$$

m = features  
n = degree

n = no. of pts. Date \_\_\_\_\_  
Page \_\_\_\_\_

LASSO  $\rightarrow$  least Absolute Shrinkage & Selection Operator.

Lasso Regression Model -

$$f(x) = w^T x + b \quad w \in \mathbb{R}^n, b \in \mathbb{R}$$

$$:= \bar{w}^T \bar{x} \quad \text{where, } \bar{w} \in \mathbb{R}^{n+1}$$

$$\bar{x} = [x, 1]$$

$$f(x) = w^T \phi(x) + b$$

$$:= \bar{w}^T \bar{\phi}(x)$$

$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_{m-1}(x) \end{bmatrix}, \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{m-1} \end{bmatrix}$$

$$\& \quad \bar{\phi}(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_{m-1}(x) \end{bmatrix} \quad \bar{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ b \end{bmatrix}$$

$$f(x) = w^T \phi(x), \quad w \in \mathbb{R}^m$$

$$\phi(x) \in \mathbb{R}^m$$

- 'm' should be chosen as how many features we have.

- Problem with Least Sq. Regression Model: -

$\rightarrow$  sensitive with outliers

$\rightarrow$  overfitting, as no. of pts. is very less and features is very high.

To overcome these problem, we use Lasso Regression Model.

$$\|w\|_1 = \sum_{i=1}^n |w_i|$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$f(x) = w^T x$$

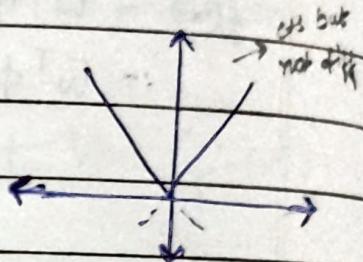
### LASSO Regression Model :

$$J(w) = \min_w \frac{\lambda}{2} \|w\|_1 + \frac{1}{2} \sum_{i=1}^L (y_i - (w^T \phi(x_i)))^2$$

cts but not differentiable

differentiable

\*  $\|w\|_1 \rightarrow$  cts but not differentiable.



$\therefore \|w\|_1$  is not differentiable, so we work with sub-gradient.

### Sub-Gradient :- (S)

A sub-gradient of a convex fn. f at  $x \in \mathbb{R}^m$  is any g, such that

$$f(y) \geq f(x) + g^T(y-x) \quad \forall y \in \mathbb{R}^m$$

all fn. which satisfies above property, is called the subgradient of that fn.

#  $g(\|x\|_1) = \begin{cases} 1, & \text{if } x \geq 0 \\ -1, & \text{if } x < 0 \\ \text{any element in } [-1, 1], & x = 0 \end{cases} \quad x \in \mathbb{R}$

↓  
sub gradient  
of  $\|x\|_1$

# Subgradient of  $\|w\|_1$ , where  $w \in \mathbb{R}^m$ .

$$\|w\|_1 = \sum_{i=1}^m |w_i|$$

$$g(i) = \begin{cases} \text{sign}(w_i) \\ \text{any element in } [-1, 1] \end{cases}, \text{ if } w(i) = 0$$

~~partial  
subgradient  
w.r.t.  
component~~

dimension of  $g(i)$  is  $m$  i.e.,  $g(i) \in \mathbb{R}^m$  ( $i$  here  $i = 1 \dots m$ )

$$S(\mathcal{J}(w)) = \frac{\lambda}{2} g - \frac{2}{L} \sum_{i=1}^L (y_i - (w^T \phi(x_i))) \phi(x_i)$$

Lasso Reg. model.

Algorithm: Subgradient of Lasso Regression Model

Input Training set  $T$ , tol

Initialize  $w^0 \in \mathbb{R}^m$

Repeat

$$w^{k+1} = w^k - \eta_k S(\mathcal{J}(w^k))$$

until

$$\|S(\mathcal{J}(w^k))\| \leq \text{tol}$$

↳ tolerance.

## Stochastic Subgradient Method :-

Input :- Training set  $T$ , max\_iter

Initialize  $w^0 \in \mathbb{R}^m$

for  $i = 1$  to max\_iter

Draw a random subset  $B$  of size  $\gamma$  from Training set  $T$

$$w^{k+1} = w^k - \eta_k \nabla J_B(w^k)$$

end

where,

$$J_B(w^k) = \frac{\lambda}{2} \|w\|_2 + \sum_{(x_i, y_i \in B)} (y_i - w^T \phi(x_i))^2$$

\* Regularised least sq. Regression Model is also called Ridge Regression Model.

$$\min_{w,b} \frac{\lambda}{2} w^T w + \frac{1}{L} \sum_{i=1}^L (y_i - (w^T x_i + b))^2 \rightarrow \text{Ridge}$$

$$\min_{w,b} \frac{\lambda}{2} \|w\|_2 + \frac{1}{L} \sum_{i=1}^L (y_i - (w^T x_i + b))^2 \rightarrow \text{Lasso}$$

$$\min_{w,b} \frac{\lambda_1}{2} (w^T w) + \frac{\lambda_2}{2} \|w\|_2 + \frac{1}{L} \sum_{i=1}^L (y_i - (w^T x_i + b))^2 \rightarrow \text{Elastic-Net}$$

$$\min_{w,b} \frac{\lambda}{2} w^T w + \frac{1}{L} \sum_{i=1}^L \|y_i - (w^T x_i + b)\|_1$$

$\downarrow$   
This loss func is optimal when noise is uniformly distributed.

## L1 - Norm Regression Model :-

$$\min_{w,b} \frac{\lambda}{2} \|w\|_1 + \frac{1}{L} \sum_{i=1}^L \|y_i - (w^T x_i + b)\|_1$$

$$\min_{w,b} \text{Regularisation} + \frac{1}{L} \sum_{i=1}^L L(y_i, w, b)$$

23/3/23

## CLASSIFICATION.

If any of the data pts. is very large & all others are very less w.r.t to the largest, then it will effect Euclidean distance & every other operation, so we normalised the data. as,

### Normalisation :-

For every variable of transaction  $x = \frac{\text{Value of var.} - \min^m \text{ value of var.}}{\max^m \text{ value of var.} - \min^m \text{ value of var.}}$

\* Relative frequency converge to probability for large number of data points.

### # Bayes formula.

$$P(w_i/x) = \frac{P(x/w_i) P(w_i)}{P(x)}$$

↑ likelihood      ↑ prior prob.  
evidence

{ if  $P(w_1/x) \geq P(w_2/x)$  ; decide  $w_1$   
else, decide  $w_2$ .

{ if  $P(x/w_1) P(w_1) \geq P(x/w_2) P(w_2)$  ; decide  $w_1$   
else, decide  $w_2$