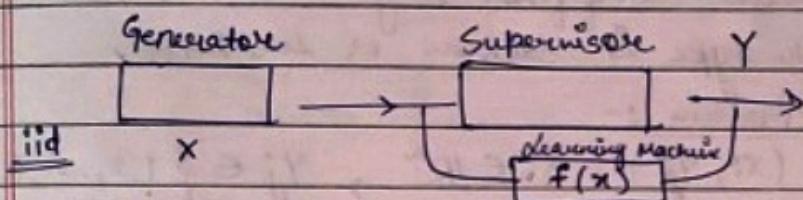


9/1/23

IT - 60

Machine Learning

- Assignments
- Surprise Quiz
- In-sem
- End-Sem



Relationship b/w X & Y is random.

Chaotic behaviour \rightarrow random nature.

(label present) \rightarrow supervised learning \rightarrow supervisor respond is must.

(label not present) \rightarrow Unsupervised " \rightarrow " " " not must.

(label not present) \rightarrow semi-supervised learning \rightarrow using few labels to get info.

Problem of Classification

How to classify different groups

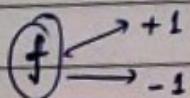
feature used to differentiate b/w diamond or gold (eg.)

X_1	X_2	X_3	f	Y	(label data)
color	light intensity	Density			
2,1,0	5	9.2		-1	(diamond)
1,2,1	6	1.6		1	(gold)
1,0,1	3	1.2		1	
1,1,1	5	0.7		1	
1,1,1	4	0.2		-1	(diamond)
1,0,1,	5	0.1			

Now this will be using ML model to predict which is gold / diamond.

Classification Problem

We need variables or labels data to classify types of fish.



e.g. - length, height, no. of fins, weight, color. to decide whether fish type is Salmon or sea bass.

Multiclassification Problem :-

\rightarrow Training set $\{ (x_i, y_j) : x_i \in \mathbb{R}^n, y_j \in \{1, 2, \dots, k\} \}$
 $i = 1, 2, \dots, l \}$ which can efficiently predict

* For binary classification $\rightarrow 1$ or 2 , $+1$ or -1 etc.

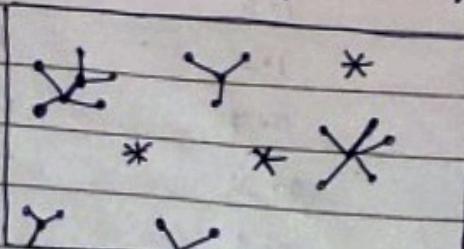
For every data pt., we have vectors in \mathbb{R}^n .

10-class classification (\because 10 numbers)

0 1 2 3 4

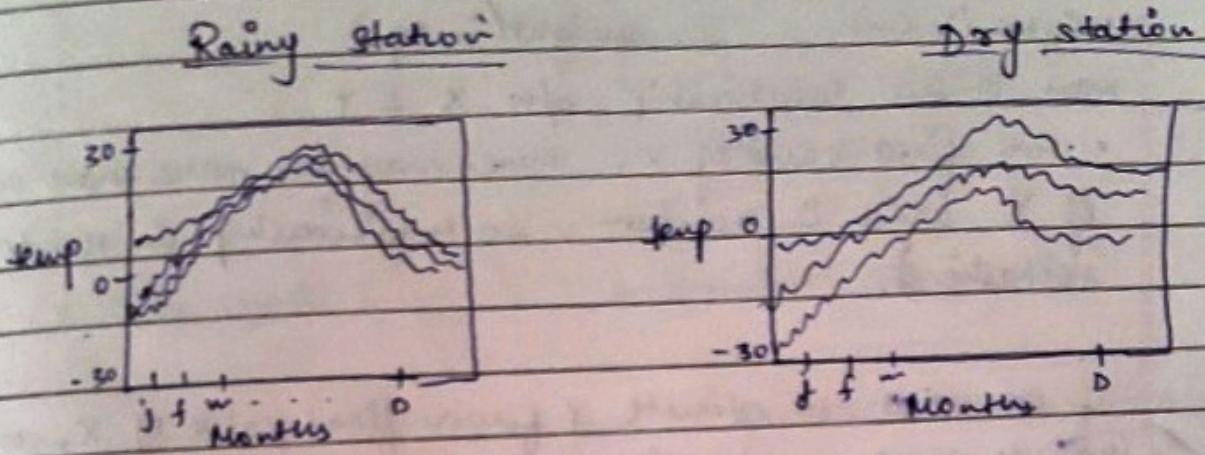
5 6 7 8 9

Graph-Data : \rightarrow e.g. Instagram data



To refer which people are students, which are not, for this we need graph data.

Functional Data : eg. - weather data , EEG data (signal)



Functional data \rightarrow when x is a function not a vector.
 e.g. ~~car~~ traffic \rightarrow which roads are prone to accidents.

Multi-class and Multi-label classification.

Multi-class			Multi-label		
$c=3$			Samples		
(100)	(010)	(001)			
labels (\pm)			Labels (\pm)		
[0 0 1]	[1 0 0]	[0 1 0]	[1 0 1]	[0 1 0]	[1 1 1]

Multi-label \rightarrow movie (components - genre, actor, director etc.)
 \rightarrow wikipedia page (components \rightarrow academic, author, entertainment etc.)

- Y is a vector
- categories are present (more than 1, so multi-label).

Regression Problem (one parameter)

Y

X

Height (in cm) , weight (in kg)

want to see relationship b/w X & Y

- For fixed value of X, there may be more than one value of Y, which is random, so relationship is random or stochastic.

If we want to estimate Y for a fixed value of X, we should check central tendency (mean, median, mode).

$$f(x) = E(Y/x)$$

more than one parameter.

we have to estimate income of a person.

having gender, education, seniority, age, work class, income

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_e, y_e)\}$$

$x_i \in R^n$ same

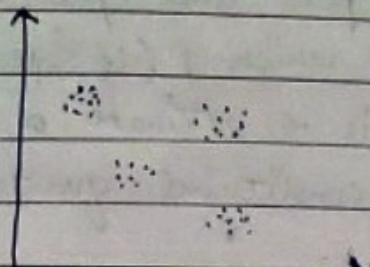
$y_i \in R$.

$$f(x) \rightarrow Y$$

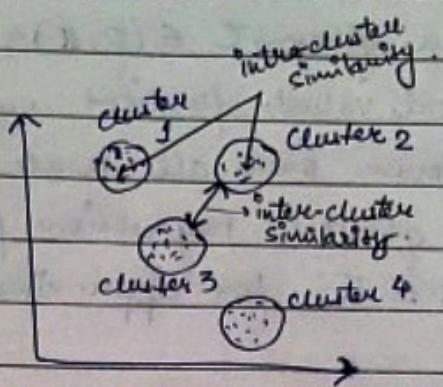
problem \rightarrow to estimate f in X of data to give approx value of Y.

If we want to predict, both income & saving, then this is the case of multi-output regression problem.

Clustering Problem:



a. Data objects



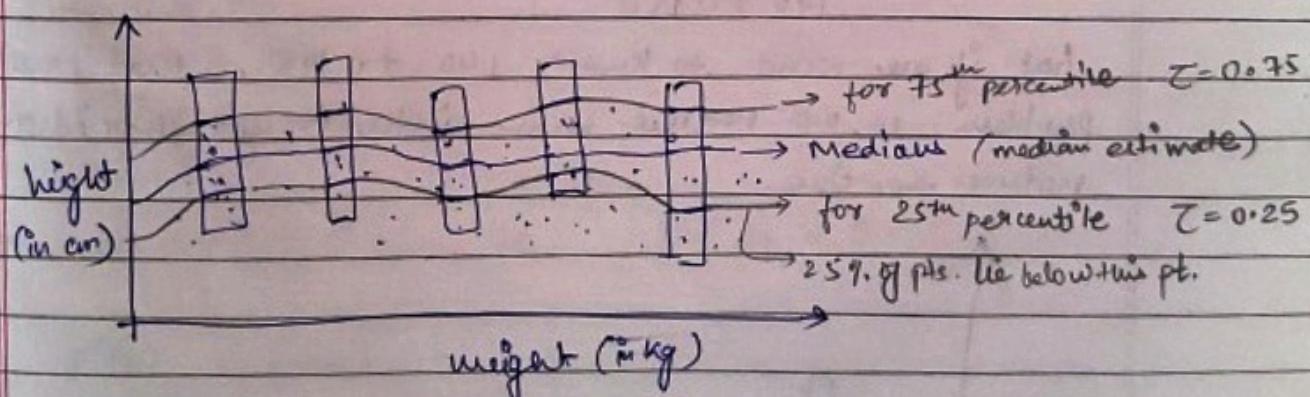
b. clustered data objects.

16/1/23 Most common problem is regression problem as any classification problem can be classified as a regression problem.

Measures of central tendency. \rightarrow Mean, Median, Mode

Mean is affected by outliers.

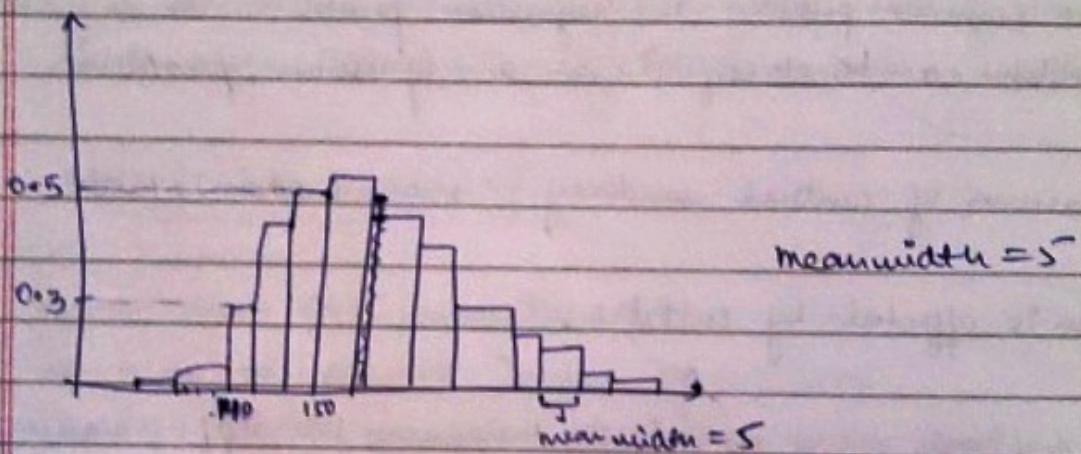
For a fixed value of weight, there may be diff. values of height.



$$f(x) = \text{median}(Y/x)$$

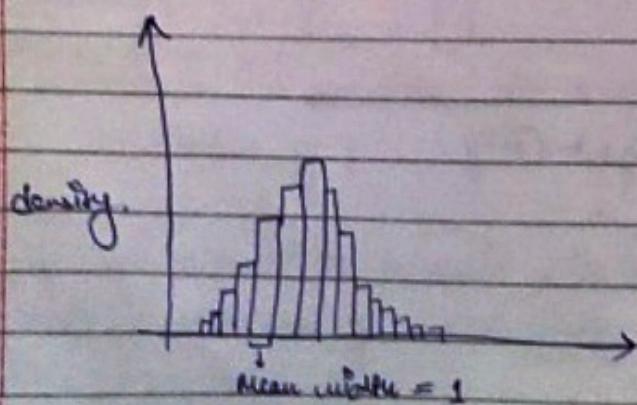
for a given $\tau \in (0,1)$, the conditional quantile μ_{τ} is a real valued function such that point wise $\mu_{\tau}(x)$ is the minimum over all real y for which $P(Y \leq \mu_{\tau}(x)) = \tau$.
 The quantile regression problem is to estimate a function $f_{\tau} \in F$ for approximating a conditional quantile $\mu_{\tau}(x)$.

conditional quantile \rightarrow $\tau\%$ of data points should be lying below $f_{\tau}(x)$.



$$0.3 + 0.5 = 0.8 \text{ % of prop. of students lies under } 140 \text{ to } 150.$$

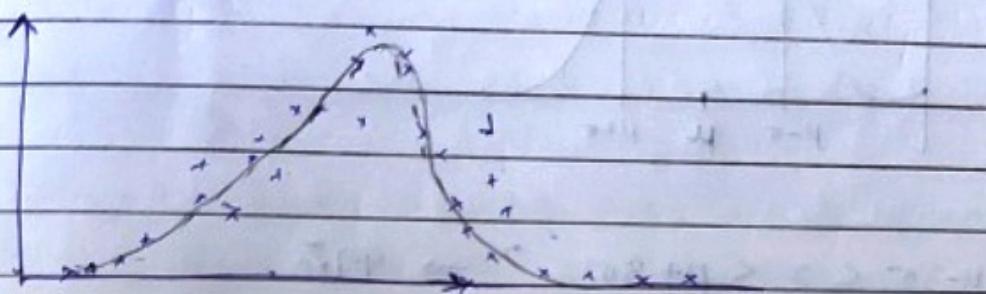
but if we want to know 140 to 148, that creates problem. So, we reduce mean width to get clear idea, i.e., reduce the size.



If width is unique, then Y axis gives relative frequency.

$$\text{Density} = \text{prob.} \times \text{width}$$

Density curve se prob. nikalna h to add kar daage dono ka prob.



$$f(x), \quad [a, b]$$

$$a f(a) + (a+h) f(a+h) + \dots$$

$$\int_a^b f(x) dx.$$

Properties:

$$\rightarrow f(x) \geq 0$$

$$\rightarrow \int_{-\infty}^{\infty} f(x) dx = 1$$

$$\Rightarrow P(x=a) = 0 \quad \text{ie} \quad \int_a^a f(x) dx = 0$$

$$\Rightarrow F'(x) = f(x)$$

$$\Rightarrow P(a \leq x \leq b) = F(b) - F(a)$$

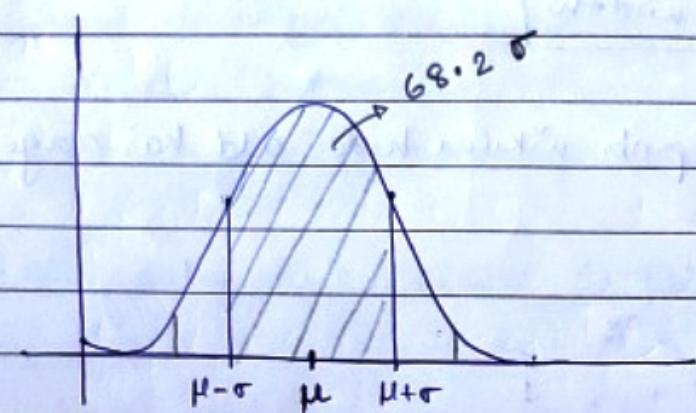
$$\Rightarrow F(a) = P(X \leq a)$$

$$\Rightarrow F(b) = P(X \leq b)$$

18/1/23



Normal Distribution



$$\mu - 3\sigma < x < \mu + 3\sigma \rightarrow 99\%$$

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}, \quad \sigma > 0, \quad -\infty < \mu < \infty$$

$$\min_c \sum_{i=1}^n (x_i - c)^2, \quad x_i \in \mathbb{R} \quad \text{find } c \rightarrow \text{centroid}$$

Mean \rightarrow sum of squares of deviation is minimum.

$$\partial \sum_{i=1}^n (x_i - c)^2 = 0$$

$$\Rightarrow \boxed{c = \frac{1}{n} \sum_{i=1}^n x_i}$$

Manhattan distance:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2}$$

$$d(x, y) = \sum_{i=1}^n |x_i - y_i| + |x_2 - y_2| + \dots + |x_n - y_n|$$

$$d(x, y) = \left[\sum_{i=1}^n (x_i - y_i)^p + (x_2 - y_2)^p + \dots + (x_n - y_n)^p \right]^{1/p}$$

(A, d)
↳ distance

- (i) → distance is always symmetric, ie, $d(x, y) = d(y, x) \forall x, y \in A$.
- (ii) → distance from $(x, y) = 0$
ie, $d(x, y) = 0 \Leftrightarrow x = y$.
- (iii) $d(x, y) + d(y, z) \geq d(x, z), x, y, z \in A$
ie, sum of 2 sides of distance is always greater than the third side.

→ Minimizing sum of Euclidean distance

$$\min_{C \in \mathbb{R}^2} \sum_{i=1}^n (x_i - c)^T (x_i - c), x_i \in \mathbb{R}^2$$

$$c = \frac{1}{n} \sum_{i=1}^n x_i$$

$$[c_1, c_2] = \frac{1}{n} \sum_{i=1}^n (x_i, y_i)$$

$$x_i = [x_i, y_i]$$

$$\min_{C \in \mathbb{R}^2} \sum_{i=1}^n (x_i, y_i)$$

$$\min_{\substack{C \in \mathbb{R}^2 \\ (c_1, c_2)}} (x_1 - c_1)^2 + (y_1 - c_2)^2 + (x_2 - c_1)^2 + (y_2 - c_2)^2 + (x_3 - c_1)^2 + (y_3 - c_2)^2 + \dots + (x_n - c_1)^2 + (y_n - c_2)^2$$

$$\text{Let } \frac{\partial L}{\partial c_1} = 0 \Rightarrow c_1 = \frac{1}{n} \sum_{i=1}^n x_i$$

var. of projected data? pt. on d.
median, mode of n data pts.



$x^T x \rightarrow$ dot product of x^T and x .

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n x_i^T d \\ &= \frac{1}{n} (x_1^T d + x_2^T d + \dots + x_n^T d) \\ &= \frac{1}{n} (x_1 + x_2 + \dots + x_n)^T d \\ &= \underbrace{\bar{x}^T d}_{\text{mean of actual data}} \quad \rightarrow (\text{projected in } d.) \end{aligned}$$

* $\bar{x}^T d$: mean of projected data in d .

Take the mean of actual data & project it on d .
(d is unit vector).

Same data points ko ek unit vector pe project kia to
bivariate or n-variate data ke scalar points milne waala
mean nikala. (formula se pata chala mean waala pt. ko
unit vector pe project karne se bivariate data ka
mean milta h).

19/1/23

Date _____
Page _____
 $x_1 \in \mathbb{R}$
 $x_2 \in \mathbb{R}$

$\bar{x} = \frac{1}{n} \sum x_i$

 \vdots
 $x_n \in \mathbb{R}$

$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

$E(\bar{x}) = \mu$

 \bar{x} is unbiased estimator of μ .

 $x \uparrow \downarrow y \rightarrow$ covariance -ve.

 $x \uparrow, y \uparrow \rightarrow$ covariance +ve.

$\text{Cov}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$

-+	++
--	+-

most
of data pts. lie b/w 1st & 3rd quadrant \rightarrow +ve covariance
 " " " " 2nd & 4th " \rightarrow -ve covariance.

$\text{Correl}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$

 $\gamma(X, Y)$

$-1 \leq \gamma \leq 1.$

$\sigma_x \sigma_y$

$x_1 \quad y_1 \rightarrow x_1 \rightarrow d^T x_1 \rightarrow (d^T x_1 - d^T \mu)^2$

$x_2 \quad y_2 \rightarrow x_2 \rightarrow d^T x_2 \rightarrow (d^T x_2 - d^T \mu)^2$

$x_L \quad y_L \rightarrow x_L \rightarrow \frac{d^T x_L}{d^T \mu} \quad (d^T x_L - d^T \mu)^2$

$\underbrace{\mu_x \quad \mu_y}_{\mu}$

Var. of projected matrix:

$\frac{1}{n-1} \sum_{i=1}^n (d^T(x_i - \mu))^2$

$= \frac{1}{L-1} (d^T(x_i - \mu))^T (d^T(x_i - \mu))$

- inner product
- norm
- dist.
- dot product

Date _____
Page _____

$x \in \mathbb{R}^n$ length of vector (Norm)

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$* \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad (\text{L}_2 \text{ Norm of } x) \quad ; \text{ distance, so tve}$$

$$* \|x\|_1 = \sum_{i=1}^n |x_i| \quad (\text{L}_1 \text{ Norm of } x)$$

$$* \|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \quad (p^{\text{th}} \text{ norm of } x)$$

$$* \|x\|_\infty = \max_{i=1,2,\dots,n} (x_i) \quad (\infty \text{ norm})$$

e.g. $x = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$

$$\|x\|_1 = 1+2+3 = 6$$

$$\|x\|_2 = \sqrt{1^2+2^2+3^2} = \sqrt{14}$$

$$* \|x+y\| \leq \|x\| + \|y\|$$

$$\|x\|_2^2 = x^T x$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$x^T x = [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$x^T x = x_1^2 + x_2^2 + \dots + x_n^2$$

$$= \frac{1}{L-1} \sum_{i=1}^L d^T (x_i - \mu)^2$$

$$= \frac{1}{L-1} (d^T (x - \mu))^2 \quad (d^T (x - \mu))^T$$

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_L \\ y_1 & y_2 & \dots & y_L \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_L \end{bmatrix}$$

$$(x - \mu) = \begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_2 & \dots & x_L - \mu_L \\ y_1 - \mu_2 & y_2 - \mu_2 & \dots & y_L - \mu_2 \end{bmatrix}$$

$$= \frac{1}{L-1} d^T (x - \mu) (x - \mu)^T d$$

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_L \\ y_1 & y_2 & \dots & y_L \end{bmatrix}$$

then projection, $x^T d$.

$$= \frac{1}{L-1} \sum_{i=1}^L [(x^T d) - (\mu^T d)]^2$$

mean
 $\mu^T d$

$$= \frac{1}{L-1} \sum_{i=1}^L [(x^T d) - (\mu^T d)]^2$$

$d^T \sum d$

$$= \frac{1}{L-1} (x^T d - \mu^T d)^T (x^T d - \mu^T d)$$

$$= \frac{1}{L-1} d^T (x^T - \mu^T)^T (x^T - \mu^T) d$$

$$= \frac{1}{L-1} d^T (x - \mu) (x - \mu)^T d.$$

Variance in the direction :

$$= \frac{1}{L-1} d^T (x-\mu) (x-\mu)^T d$$

$$\frac{1}{L} (x-\mu) (x-\mu)^T = \begin{bmatrix} x_1 & x_2 & \dots & x_L \\ y_1 & y_2 & \dots & y_L \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

$$(x-\mu) (x-\mu)^T = \frac{1}{L} \begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_1 & \dots & x_L - \mu_1 \\ y_1 - \mu_2 & y_2 - \mu_2 & \dots & y_L - \mu_2 \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 & y_1 - \mu_2 \\ x_2 - \mu_1 & y_2 - \mu_2 \\ \vdots & \vdots \\ x_L - \mu_1 & y_L - \mu_2 \end{bmatrix}$$

$$= \frac{1}{L-1} \begin{bmatrix} \sum_{i=1}^L (x_i - \mu_1)^2 & \sum_{i=1}^L (x_i - \mu_1)(y_i - \mu_2) \\ \sum_{i=1}^L (x_i - \mu_1)(y_i - \mu_2) & \sum_{i=1}^L (y_i - \mu_2)^2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix}$$

\downarrow
Variance-Covariance Matrix.

$$\begin{bmatrix} x_1 & x_2 & \dots & x_L \\ y_1 & y_2 & \dots & y_L \end{bmatrix} \in \mathbb{R}^2$$

$$\begin{bmatrix} x_1 & x_2 & \dots & x_L \\ y_1 & y_2 & \dots & y_L \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix}$$

$$\text{Variance} = d^T \Sigma d.$$

$$\begin{bmatrix} X_{11} & X_{21} & \cdots & X_{L1} \\ X_{12} & X_{22} & \cdots & X_{L2} \\ \vdots & \vdots & & \vdots \\ X_{1n} & X_{2n} & \cdots & X_{Ln} \end{bmatrix}$$

$$d = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix}$$

$$\|d\|_2 = 1.$$

$$\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22}^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_{nn}^2 \end{bmatrix}$$

→ Var Cov Matrix for
n components.

$\begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} = A$

$$A = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$\mu_1 = \frac{2}{4} = \frac{1}{2}$$

$$\mu_2 = \frac{1}{2}$$

$$(A - \mu)^T = \begin{bmatrix} -0.5 & -0.5 & 0.5 & 0.5 \\ -0.5 & 0.5 & -0.5 & 0.5 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

$$(A - \mu)^T = \begin{bmatrix} -0.5 & -0.5 \\ -0.5 & 0.5 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$\textcircled{2}$ $(A - \mu)(A - \mu)^T = \begin{bmatrix} \dots \end{bmatrix}$

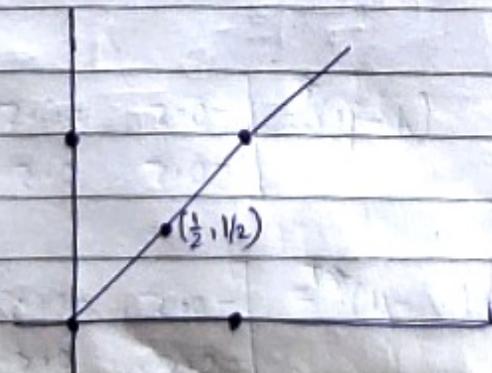
$$(A - \mu I)(A - \mu I)^T = \begin{bmatrix} -0.5 & -0.5 & -0.5 & 0.5 \\ 0.5 & 0.5 & 0.5 & -0.5 \\ 0.5 & -0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0.5 & 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma = \frac{1}{L-1} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad L=4$$

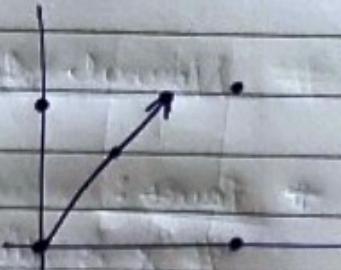
$$= \frac{1}{3} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

\textcircled{B} $\Sigma = \begin{bmatrix} 1/3 & 0 \\ 0 & 1/3 \end{bmatrix}$



$$\begin{bmatrix} -0.5 & -0.5 \\ -0.5 & 0.5 \\ 0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$d = \begin{bmatrix} 3/5 \\ 4/5 \end{bmatrix}$$



$$d^T \Sigma d = \begin{bmatrix} 3/5 & 4/5 \end{bmatrix} \begin{bmatrix} 1/3 & 0 \\ 0 & 1/3 \end{bmatrix} \begin{bmatrix} 3/5 \\ 4/5 \end{bmatrix}$$

$$= \begin{bmatrix} 3/5 & 4/5 \end{bmatrix} \begin{bmatrix} 3/15 \\ 4/15 \end{bmatrix} \stackrel{1 \times 2}{=} \frac{9}{75} + \frac{16}{75} = \frac{25}{75} = \frac{1}{3}$$

$$= \frac{9}{75} + \frac{16}{75} = \frac{25}{75} = \frac{1}{3}$$

$$d^T \Sigma d = \frac{1}{3} \rightarrow \text{Variance of projected data pt. in } d.$$

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}_{4 \times 2} \begin{bmatrix} 3/5 \\ 4/5 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 0 \\ 4/5 \\ 3/5 \\ 7/5 \end{bmatrix}$$

↓

projectionen in d dir^n.

25/1/23

Page

Normal distribution,

* Result:

If X is distributed as $N_p(\mu, \Sigma)$ then any linear combination of variables $a'x = a_1x_1 + a_2x_2 + \dots + a_px_p$ is distributed as $N(a'\mu, a'\Sigma a)$.

Also if $a'x$ is distributed as $N(a'\mu, a'\Sigma a)$ for every a , then X must be $N_p(\mu, \Sigma)$.

* If X is distributed as $N_p(\mu, \Sigma)$, then q linear combinations,

$$A_{(q \times p)} X_{(p \times 1)} = \begin{bmatrix} a_{11}x_1 + \dots + a_{1p}x_p \\ a_{21}x_1 + \dots + a_{2p}x_p \\ \vdots \\ a_{q1}x_1 + \dots + a_{qp}x_p \end{bmatrix}$$

In 1-D

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$\hat{y} \rightarrow \begin{bmatrix} 8.0 & 0.2 \end{bmatrix} \rightarrow \text{var. cov. matrix. (let)}$$

var
(matrix)
var
(y dir*)

$$\text{Mean} = 0, 0$$

$$(x - \mu)^T \Sigma^{-1} (x - \mu)$$

* dist. of $\frac{(x-\mu)}{\sigma}$ for pt \rightarrow Mahalanobis distance.

Date _____
Page _____

$$\text{Variance} \rightarrow (x-\mu)^T \Sigma^{-1} (x-\mu)$$

x mean se kitne s.d. dene h $\xrightarrow{\text{given}}$ $(x-\mu)^T \Sigma^{-1} (x-\mu)$

For multi-Dimension, : $X \sim N_p(\mu, \Sigma)$

$$f(x) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$
 (in p-dimension)

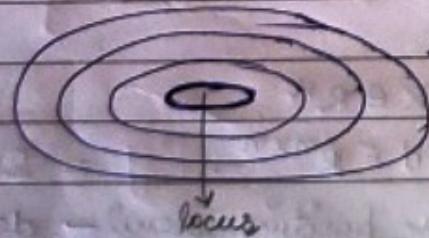
$$-\infty < x_i < \infty, i=1, 2, \dots, p.$$

where, $\Sigma \rightarrow$ positive definite

6/2/23

Plotting multivariate Normal Density.

$$f = (x-\mu)^T \Sigma^{-1} (x-\mu)$$



contour plot

for symmetric distribution \rightarrow contour is circle otherwise,
it is ellipse.

Bivariate Normal Density:

Evaluating for $p = 2$ variate normal density in terms of individual parameters,

$$\mu_1 = E(X_1)$$

$$\mu_2 = E(X_2)$$

$$\text{Var}(X_1) = \sigma_{11}$$

$$\text{Var}(X_2) = \sigma_{22}$$

$$\text{Cov}(X_1, X_2) = \rho_{12} = \frac{\sigma_{12}}{\sqrt{\sigma_{11}} \sqrt{\sigma_{22}}}$$

Result: If Σ is positive definite, so that Σ^{-1} exists.
then

$$\Sigma e = \lambda e$$

$$\Rightarrow \Sigma^{-1} e = \frac{1}{\lambda} e.$$

so (λ, e) is an eigenvalue - eigenvector pair for Σ corresponding to the pair $(\frac{1}{\lambda}, e)$ for Σ^{-1} .

Also Σ^{-1} is positive definite.

If $X^T H X \geq 0$, $X \in \mathbb{R}^n$,
 $H \in \mathbb{R}^{n \times n}$.

then H is said to be a positive semi-definite matrix.

Constant probability density contours :-

= {all x such that $(x-\mu)^T \Sigma^{-1} (x-\mu) = c^2$ }

= surface of an ellipsoid centered at μ .

Contours of constant density for the p -dimensional normal distib. are ellipsoids defined by x such that

$$(x-\mu)^T \Sigma^{-1} (x-\mu) = c^2$$

These ellipsoids are centered at μ & have axes $\pm c\sqrt{\lambda_i}e_i$
where $e_i = \Sigma e_i = \lambda_i$ for $i=1,2,\dots,p$.

Convex sets :

$C \subseteq \mathbb{R}^n$ is convex if. $tx + (1-t)y \in C$ for any $x, y \in C$ & $0 \leq t \leq 1$.



convex set



not convex set.

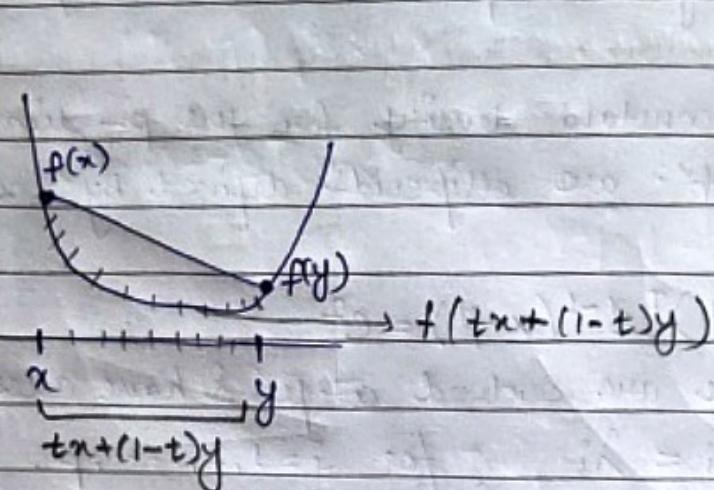
Convex function :

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex

if $\text{dom}(f)$ (the domain of f) is a convex set,
& if $f(tx + (1-t)y) \leq t f(x) + (1-t)f(y)$ for any $x, y \in \text{dom}(f)$ & $0 \leq t \leq 1$.

$$tx + (1-t)y \rightarrow \text{pt. on axis.}$$

$$\Rightarrow f(tx + (1-t)y) \leq t f(x) + (1-t)f(y)$$



Convex function:

$f(x)$ is convex fn.

then, $f(x)$ would be concave fn.

also, $\alpha f(x)$ where $\alpha \geq 0$ is convex func.

$f_1(x) \cdot f_2(x)$ is convex then,

$f_1(x) + f_2(x)$ would be convex fn.

then, $\max(f_1(x) + f_2(x))$ is convex fn.

$$\Rightarrow y = x^2 \rightarrow \text{convex fn. } (\because \frac{d^2y}{dx^2} = 2 > 0)$$

$$\Rightarrow y = -x^2 \rightarrow \text{concave fn.}$$

$$\Rightarrow y = \sin x \rightarrow \text{concave on interval } [0, \pi].$$

5/2 | 29

strictly convex function :

$$f(\lambda x + (1-\lambda)y) < \lambda f(x) + (1-\lambda)f(y)$$

where $0 \leq \lambda \leq 1$.

Concave function :

$$f(\lambda x + (1-\lambda)y) \geq \lambda f(x) + (1-\lambda)f(y)$$

where $0 \leq \lambda \leq 1$.

strictly concave function :

$$f(\lambda x + (1-\lambda)y) > \lambda f(x) + (1-\lambda)f(y)$$

where $0 \leq \lambda \leq 1$.

Gradient and derivatives :-

- for a fn. $f(x) = f(x_1, x_2, \dots, x_n)$, and a unit vector $u = (u_1, \dots, u_n)$, the directional derivative is

$$\text{as } \nabla_u f(x) = \lim_{h \rightarrow 0} \left(\frac{f(x+hu) - f(x)}{h} \right)$$

$f(x)$ is differentiable implies that $\nabla_u f(x)$ is well defined for all x & u can be obtained as

gradient
(rate of change)

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

& the directional derivative of f at pt. x ,
for any dir. (unit vector) u can
be obtained as,

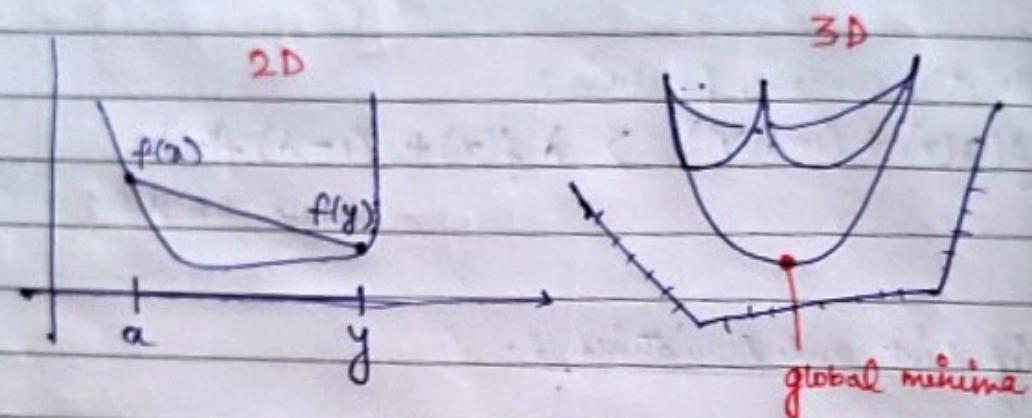
$$\nabla_u f(x) = u^T \nabla f$$

Unconstrained Convex Optimization

Consider the optimization problem

$$\min_{\mathbf{x}} f(\mathbf{x}) ; \quad f: \mathbb{R}^n \rightarrow \mathbb{R} \text{ is convex \& smooth,}$$

then the necessary \& sufficient condition for optimal sol. \mathbf{x}_0
 is $\nabla f(\mathbf{x}) = 0$ at $\mathbf{x} = \mathbf{x}_0$



$$(\nabla f(\mathbf{x}))^\top = \left[\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2}, \dots, \frac{\partial L}{\partial x_n} \right]$$

If d dirⁿ. we pta karne h, then,

$$d^\top \nabla f(\mathbf{x}), \quad \text{where } d \text{ is unit vector} \\ \& \|d\|_2 = 1$$

Find the dirⁿ. u such that, $d^\top \nabla f(\mathbf{x})$ is max^m(u)
 such that $\|u\|_2 = ?$

* if max_b a^T b (we have to ~~max~~ max^m b, then a^T b
 should be in dirⁿ. of a^T, then a^T b
 will be max^m).)

→ So, in dirⁿ. of $\nabla f(\mathbf{x})$, then $d^\top \nabla f(\mathbf{x})$ will be
 sharply increasing, ie, in dirⁿ. of gradient.
 $\Rightarrow \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}$

* for minimum, the dirn. of gradient will be -ve.

i.e if $\min_b a^T b$, then, $b = -a^T$.

$$\Rightarrow \frac{\nabla f(x)}{\|\nabla f(x)\|} = -\frac{\nabla f(x)}{\|\nabla f(x)\|}$$

Gradient-Descent Algorithm :

- An iterative algorithm.
- The -ve gradient dirn. is dirn. of steepest descent.

* Method of Gradient Descent :

(i) Initialise $x_0 = x_{\text{start}} \in \mathbb{R}^n$

$K = \text{no. of steps}$

(ii) Repeat $x^{(k+1)} = x^{(k)} - \gamma_k \nabla f(x^{(k)})$

$\gamma_k = \text{eta/gamma}$

until $\|\nabla f(x^{(k)})\| \leq \epsilon$.

\downarrow

tells how much
we are moving
in a dirn.

(Step length)

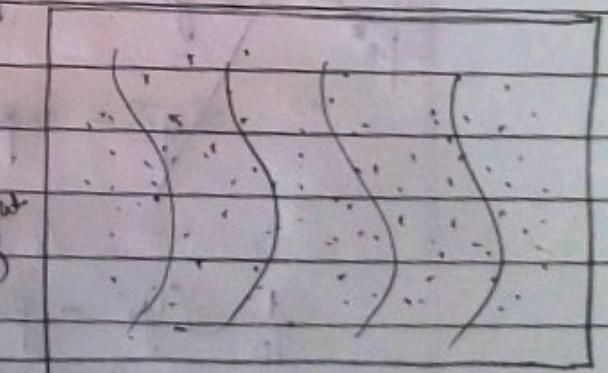
Regression :

$$Y_i = E(Y_i/x) + \varepsilon_i$$

for mean,

$$f(Y_i) = E(Y/x)$$

Weight
(mm)



Assumptions :

(i) all X_i & Y_i are i.i.d.

Weight (kg)

$$(ii) Y_i = E(Y_i/x) + \varepsilon_i$$

\downarrow

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

(i.e., mean of error term = 0).

random
in nature

Task: predict balance using income.

- Training set.
- Testing set

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

$$\varepsilon = [y - (\beta_0 + \beta_1 x_1)]^2 \quad \text{we have to minimize error.}$$

$$\min_{\beta_0, \beta_1} \frac{1}{10} \sum_{i=1}^{10} (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$\begin{aligned} \frac{\partial}{\partial \beta_0} &= 0 \\ \Rightarrow \frac{2}{10} \cdot \sum (y_i - (\beta_0 + \beta_1 x_i)) &= 0 \\ \frac{\partial^2}{\partial \beta_0^2} &\Rightarrow \sum y_i = \sum (\beta_0 + \beta_1 x_i) = 0 \\ \Rightarrow \sum y_i &= n \beta_0 \end{aligned}$$

$$\left. \begin{aligned} \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{(x_i - \bar{x})(y_i - \bar{y})}{(x_i - \bar{x})^2} \end{aligned} \right\}$$

$$\frac{1}{10} \sum (y_i - \beta_0 - \beta_1 x_i)^2 = 0$$

$$\frac{\partial}{\partial \beta_1} = \frac{2}{10} \sum (y_i - \beta_0 - \beta_1 x_i) (-1) = 0$$

$$\Rightarrow -\sum y_i + \beta_0 + \beta_1 \bar{x} = 0$$

$$\Rightarrow \sum y_i = \beta_0 + \beta_1 \bar{x} \quad \Rightarrow \hat{\beta}_0 = \frac{\sum (y_i - \hat{\beta}_1 x_i)}{n} \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\frac{\partial}{\partial \beta_1} = \frac{2}{n} \sum (y_i - \beta_0 - \beta_1 x_i) (-x_i)$$

$$= -\sum y_i x_i + \beta_0 n + \beta_1 n x_i^2 = 0$$

$$\Rightarrow \sum y_i x_i = \beta_0 n + \beta_1 n x_i^2$$

$$\hat{\beta}_1 = \frac{\sum y_i x_i - \beta_0 n}{n x_i^2}$$

$$\hat{\beta}_1 = \frac{\sum y_i x_i}{n x_i^2}$$

$$\hat{\beta}_1 = \frac{\sum y_i^2}{n x_i^2} - \frac{\beta_0}{n}$$

$$= \frac{\sum y_i^2}{n^2} - \frac{1}{n^2} (\sum y_i^2 - \hat{\beta}_1 \sum x_i^2)$$

$$= \frac{\sum y_i^2}{n^2} - \frac{\sum y_i^2}{n^2} + \hat{\beta}_1$$

$$(3+2)(4-2)$$

$$(3+2)(4-2)$$

$$(3+2)(4-2)$$

$$(3+2)(4-2) = (3A-1)^T (3A-1)$$

$$(3+2)(4-2)$$

$$(3+2)(4-2)$$

$$(3+2)(4-2)$$

$$(3+2)(4-2)$$

$$(3+2)(4-2) + (3+2)(4-2) = (3A-1)^T (3A-1)$$

$$-2(3+2)(4-2) +$$

Fitting least Sq. Regression polynomial

Date _____
Page _____

our estimate is $wx + b \rightarrow f(x)$

$$\min_{w, b \in \mathbb{R}} \sum_{i=1}^L (y_i - (wx_i + b))^2$$

$$\min_{u \in \mathbb{R}^2} (Y - Au)^T (Y - Au)$$

$$Au = \begin{bmatrix} wx_1 + b \\ wx_2 + b \\ \vdots \\ wx_L + b \end{bmatrix}$$

$$= \begin{bmatrix} y_1 - (wx_1 + b) \\ y_2 - (wx_2 + b) \\ \vdots \\ y_L - (wx_L + b) \end{bmatrix}$$

$$A = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_L & 1 \end{bmatrix}; u = \begin{bmatrix} w \\ b \end{bmatrix}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_L \end{bmatrix}$$

$$\min_{u \in \mathbb{R}^2} (Y - Au)^T (Y - Au) = \begin{bmatrix} y_1 - (wx_1 + b) \\ y_2 - (wx_2 + b) \\ \vdots \\ y_L - (wx_L + b) \end{bmatrix}^T \begin{bmatrix} y_1 - (wx_1 + b) \\ y_2 - (wx_2 + b) \\ \vdots \\ y_L - (wx_L + b) \end{bmatrix}$$

$$\min_{u \in \mathbb{R}^2} (Y - Au)^T (Y - Au) = (y_1 - (wx_1 + b))^2 + (y_2 - (wx_2 + b))^2 + \dots + (y_L - (wx_L + b))^2$$

$$\min_u (Y - Ay)^T (Y - Ay) \simeq J(u)$$

$$\nabla_u J(u) = 2A^T(Y - Ay) = 0$$

$$\Rightarrow A^T(Y - Ay) = 0$$

$$\Rightarrow u = (A^T A)^{-1} A^T y$$

$$\frac{\nabla J(u)}{\nabla u} = \frac{(Y - Ay)}{u}$$

$$\min x^T x = F(x), x \in \mathbb{R}^n$$

$$\min \sum_{i=1}^n (x_1^2 + x_2^2 + \dots + x_n^2)$$

$$\frac{\partial}{\partial x_1} = 2x_1, \frac{\partial}{\partial x_2} = 2x_2$$

:

$$\# \frac{\partial (x^T x)}{\partial x} = 2Ax$$

↓
matrix

$$\# \nabla_x Ax = AT$$

$$\text{or } \frac{\partial}{\partial x} (Ax) = AT$$

$$\Rightarrow \frac{\partial \min_u (Y - Ay)^T (Y - Ay)}{\partial u}$$

$$= -2(Y - Ay) = 0$$

$$\Rightarrow A^T(Y - Ay) = 0$$

$$u = (A^T A)^{-1} A^T y$$

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = (x^T x)^{-1} x^T y$$

$$u = \begin{bmatrix} \omega \\ b \end{bmatrix} = \begin{bmatrix} \beta_1 \\ \beta_0 \end{bmatrix}$$

$$A = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_L & 1 \end{bmatrix}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_L \end{bmatrix}$$

$$f(x) = 2 \cdot 6 + 6 \cdot 9x$$

$$\begin{array}{c} A^T \\ \hline X & Y \\ \begin{matrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 0 & 1 \end{matrix} \end{array}$$

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} \quad (\because A = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix})$$

$$u = (A^T A)^{-1} A^T Y$$

$$A^T A = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix}_{2 \times 4} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ 0 & 1 \end{bmatrix}_{4 \times 2}$$

$$= \begin{bmatrix} 2 & 2 \\ 2 & 4 \end{bmatrix}$$

$$(A^T A)^{-1} = \frac{1}{4} \begin{bmatrix} 4 & -2 \\ -2 & 2 \end{bmatrix} = \begin{bmatrix} 1 & -1/2 \\ -1/2 & 1/2 \end{bmatrix}$$

$$(A^T A)^{-1} A^T Y = \begin{bmatrix} 1 & -1/2 \\ -1/2 & 1/2 \end{bmatrix}_{2 \times 2} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}_{4 \times 4} \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}_{4 \times 1}$$

$$= \begin{bmatrix} 1 & -1/2 \\ -1/2 & 1/2 \end{bmatrix}_{2 \times 2} \begin{bmatrix} 1 \\ 3 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 1 - 3/2 \\ -1/2 + 3/2 \end{bmatrix}_{2 \times 1}$$

$$(A^T A)^{-1} A^T Y = \begin{bmatrix} -1/2 \\ 1 \end{bmatrix}$$

RMSE:

RMSE → Root Mean Square Error

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

* RMSE = $\sqrt{\frac{1}{K} \sum_{i=1}^K (y_i - (\beta_0 + \beta_1 x_i))^2}$

Quadratic Fitting.

$$f(x) = \beta_0 + \beta_1 x^1 + \beta_2 x^2$$

$$f(x) = w_2 x^2 + w_1 x + b$$

$$\min_{(w_2, w_1, b)} \sum_{i=1}^n (y_i - (w_2 x_i^2 + w_1 x_i + b))^2$$

$$A = \begin{bmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots \\ x_L^2 & x_L & 1 \end{bmatrix} \quad u = \begin{bmatrix} w_2 \\ w_1 \\ b \end{bmatrix} \xrightarrow{\text{u}} \begin{bmatrix} \beta_2 \\ \beta_1 \\ \beta_0 \end{bmatrix}$$

$$\min_u (Y - Au)^T (Y - Au)$$

Cubic Fitting

$$f(x) = w_3 x^3 + w_2 x^2 + w_1 x + b$$

$$\min_{(w_3, w_2, w_1, b)} \sum_{i=1}^n (y_i - (w_3 x_i^3 + w_2 x_i^2 + w_1 x_i + b))^2$$

$$A = \begin{bmatrix} x_1^3 & x_1^2 & x_1 & 1 \\ x_2^3 & x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ x_L^3 & x_L^2 & x_L & 1 \end{bmatrix}_{L \times 4} \quad u = \begin{bmatrix} w_3 \\ w_2 \\ w_1 \\ b \end{bmatrix}$$

$$f(x) = w_1 x_1 + w_2 x_2 + b$$

~~for quad-fitting in n dimension, A=? u=? , w=?~~



Fitting with fifth order polynomial

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4$$

$$\min_{(w_4, w_3, w_2, w_1, b)} \sum_{i=1}^L \left(y_i - (w_4 x^4 + w_3 x^3 + w_2 x^2 + w_1 x + b) \right)^2$$

$$A = \begin{bmatrix} x_1^4 & x_1^3 & x_1^2 & x_1 & 1 \\ x_2^4 & x_2^3 & x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_L^4 & x_L^3 & x_L^2 & x_L & 1 \end{bmatrix}_{L \times 5}$$

$$u = \begin{bmatrix} w_4 \\ w_3 \\ w_2 \\ w_1 \\ b \end{bmatrix}$$

$$u = (A^T A)^{-1} A^T y$$

when $f(x) = w_1 x_1 + w_2 x_2 + b$ (ie, $x_1 \& x_2$)

$$f(x) = w_1 x_1 + w_2 x_2 + b$$

$$\min_{(w_2, w_1, b)} \sum_{i=1}^L \left(y_i - (w_1 x_1 + w_2 x_2 + b) \right)^2$$

$$\text{if then, } X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix},$$

$$\text{then, } w^T X + b = f(x)$$

i.e,

$$f(x) = w^T X + b = w_1 x_1 + w_2 x_2 + b$$

$$\text{or, } \min_{\substack{(w \in \mathbb{R}^2) \\ (b \in \mathbb{R}^2)}} \sum_{i=1}^L \left[y_i - (w^T X + b) \right]^2$$

$$\min_{\substack{w \in \mathbb{R}^2 \\ b \in \mathbb{R}}} \sum_{i=1}^n (y_i - (w^T x_i + b))^2$$

$$\min_{\substack{w \in \mathbb{R}^2 \\ u}} (Y - Au)^T (Y - Au)$$

where, $A = \begin{bmatrix} x_{11} & x_{12} & 1 \\ x_{21} & x_{22} & 1 \\ \vdots & \vdots & \vdots \\ x_{L1} & x_{L2} & 1 \end{bmatrix}$ $u = \begin{bmatrix} w_2 \\ w_1 \\ \vdots \\ b \end{bmatrix}$, $b \in \mathbb{R}$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_L \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \in \mathbb{R}^2 \quad X_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iL} \end{bmatrix} \in \mathbb{R}^2$$

$$u = (A^T A)^{-1} A^T Y$$

$$\# \quad u = (A^T A)^{-1} A^T Y = \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} = \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$$x_1 = 5.6, \quad x_2 = 3.7 \quad \text{Predict balance}$$

$$\therefore f(x) = w_1 x_1 + w_2 x_2 + b \\ = 0.7 x_1 + 0.2 x_2 + 0.1.$$

i.e., x_1 contributes most in balance.

Improving the prediction for $M=8$

$$\# \quad \min \frac{1}{10} \sum_{i=1}^{10} [y_i - (\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5 + \beta_6 x_i^6 + \beta_7 x_i^7)]^2$$

$$\# \quad x^T H x + f^T x + b$$

$$\rightarrow x^T H x + f^T x + b \quad ; \quad w^T x + b ; \quad w \in \mathbb{R}^n \\ b \in \mathbb{R}$$

$$H = \text{sq. matrix}, \quad = \begin{bmatrix} w_5 & w_3/2 \\ \frac{w_3}{2} & w_4 \end{bmatrix} = H$$

$$f = \begin{bmatrix} w_2 \\ w_1 \end{bmatrix}$$

$$f(x) = f \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = w_5 x_1^2 + w_4 x_2^2 + w_3 x_2 x_1 + w_2 x_2 + w_1 x_1 + b$$

$$\Rightarrow \min_{(w_5, w_4, \dots, w_1, b)} \sum_{i=1}^l \left[y_i - (w_5 x_{i1}^2 + w_4 x_{i2}^2 + w_3 x_{i1} x_{i2} + w_2 x_{i2} + w_1 x_{i1} + b) \right]^2$$

$$A = \begin{vmatrix} x_{11}^2 & x_{12}^2 & x_{11}x_{12} & x_{12} & x_{12} & 1 \\ x_{21}^2 & x_{22}^2 & x_{21}x_{22} & x_{21} & x_{22} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ x_{l1}^2 & x_{l2}^2 & x_{l1}x_{l2} & x_{l1} & x_{l2} & 1 \end{vmatrix}$$

$$u = \begin{bmatrix} w_5 \\ w_4 \\ w_3 \\ w_2 \\ w_1 \\ b \end{bmatrix} \quad \text{or} \quad u = (A^T A)^{-1} A^T Y$$

when $x \in \mathbb{R}^n$

i.e. $x \in \mathbb{R}^n$, $y \in \mathbb{R}$

$$A = \begin{bmatrix} x_{11}^n & x_{12}^n & \dots & x_{1n}^n & x_{11} x_{21}^{n-1} & x_{11}^m x_{12}^{n-2} & \dots & 1 \\ x_{21}^n & x_{22}^n & \dots & x_{2n}^n & x_{21} x_{22}^{n-1} & x_{21}^2 x_{22}^{n-2} & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{L1}^n & x_{L2}^n & \dots & x_{Ln}^n & x_{L1} x_{L2}^{n-1} & x_{L1}^2 x_{L2}^{n-2} & \dots & 1 \end{bmatrix}$$

9/2/23

#1. $T = \{(x_i, y_i) : x_i \in \mathbb{R}^n, y_i \in \mathbb{R}, i=1, 2, \dots, L\}$

$$f(x) \sim Y$$

The regression problem is to find $f(x)$ which can approximate y values well. For this we need to minimize error.

$$\min_f \sum_{i=1}^L (y_i - f(x_i))^2.$$

$$L(x, y, f)$$

$$\min_f \sum_{i=1}^L L(x_i, y_i, f)$$

$$\sum_{i=1}^L (y_i - f(x_i))^2$$

One of the type of loss fun. :-

- Least square loss function :-

$$\begin{aligned} L(x_i, y_i, f) \\ = (y - f(x_i))^2 \end{aligned}$$

2. Find a $f \in F$ such that $\sum_{i=1}^L (y_i - f(x_i))^2$.

→ Set of linear fun. in \mathbb{R}^n

$$F = \{w^T x + b : w \in \mathbb{R}^n, b \in \mathbb{R}\}$$

$$\min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \sum_{i=1}^L (y_i - (w^T x_i + b))^2.$$

$$[w] = (A^T A)^{-1} A^T y$$

$$B = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} & 1 \end{bmatrix}$$

$$x_{21} & x_{22} & \dots & x_{2n} & 1$$

$$A = \begin{bmatrix} \vdots \\ x_{L1} & x_{L2} & \dots & x_{Ln} & 1 \end{bmatrix}$$

\Rightarrow

* In n dimension, for m no. of parameter, no. of possible terms:
 $= \frac{(m+n)!}{m! n!}$

For quadratic,

$$\frac{(2+n)!}{2! n!}$$

$$= \frac{(2+n)(1+n)n!}{2! n!} = \frac{(n+1)(n+2)}{2}$$

for cubic, $\frac{(3+n)!}{3! n!}$

$$A = \begin{bmatrix} X_{11}^m & X_{12}^m & \dots & X_{1n}^m & \dots & \dots & 1 \\ X_{21}^m & X_{22}^m & \dots & X_{2n}^m & \dots & \dots & 1 \\ \vdots & & & & & & \\ X_{L1}^m & X_{L2}^m & \dots & X_{Ln}^m & \dots & \dots & 1 \end{bmatrix}$$

$$\begin{bmatrix} w_k \\ w_{k-1} \\ \vdots \\ w_1 \\ b \end{bmatrix} = (A^T A)^{-1} A^T y$$

$K = \frac{(m+n)!}{m! n!} - 1$

Training set

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\} \quad x_i \in \mathbb{R}^n, \quad y_i \in \mathbb{R}, \quad i=1, 2, \dots, L$$

→ to fit a line/hyperplane

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

$$= w_K x_1^m + w_{K-1} x_2^m + \dots + w_{k-m} x_n^m + w_{k-(m+1)} b$$

x_n^m

$$+ w_{k-n} x_1^{m-1} x_2 + \dots + \dots + b.$$

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

(Non-linear) $= w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_m \phi_m(x)$

$$\text{where, } \phi_1(x) = x_1^m$$

$$\phi_2(x) = x_2^m$$

$$\phi_{m-1}(x) = x_n$$

$$\phi_m(x) = 1 \quad (\text{may be } b)$$

$$\phi(x) = [x_1^T, x_2^T, x_3^T]^T; \text{ s.t. } x_1 + x_2 + \dots + x_n \leq m$$

(Non-linear) $\Rightarrow f(x) = w^T \phi(x)$

$$= \begin{bmatrix} w_m \\ w_{m-1} \\ \vdots \\ w_3 \\ w_2 \\ b \end{bmatrix}^T \begin{bmatrix} \phi_m(x) \\ \phi_{m-1}(x) \\ \vdots \\ \phi_3(x) \\ \phi_2(x) \\ 1 \end{bmatrix} \rightarrow x^m$$

$$= \begin{bmatrix} w_m \\ w_{m-1} \\ \vdots \\ w_2 \\ b \end{bmatrix}^T \begin{bmatrix} \phi_m(x) \\ \phi_{m-1}(x) \\ \vdots \\ 1 \end{bmatrix}$$

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

$$\text{for Quadratic} \rightarrow \begin{bmatrix} w_2 \\ w_1 \\ \vdots \\ b \end{bmatrix}, \begin{bmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots \\ x_L^2 & x_L & 1 \end{bmatrix}$$

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

$\downarrow \phi_1(x)$ $\downarrow \phi_2(x)$ $\downarrow \phi_3(x)$ $\downarrow \phi_4(x)$ $\downarrow \phi_5(x)$ $\downarrow \phi_6(x)$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \\ b \end{bmatrix}, \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \phi_3(x) \\ \phi_4(x) \\ \phi_5(x) \\ 1 \end{bmatrix}$$

$$\min_{w \in \mathbb{R}^n} \sum_{i=1}^L (y_i - (w^T x_i + b))^2 \quad \xrightarrow{\text{linear}}$$

$$b \in \mathbb{R}$$

$$\min_{w \in \mathbb{R}^m} \sum_{i=1}^L (y_i - (w^T \phi(x_i) + b))^2$$

$$b \in \mathbb{R}$$

take derivative
of square it to 0.

$$\nabla_w J(w, b) = 0$$

$$\text{where } \phi(x^*) = \begin{bmatrix} \phi_m(x) \\ \phi_{m-1}(x) \\ \vdots \\ \phi_1(x) \end{bmatrix}$$

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

$$A = \begin{bmatrix} \phi_m(x_1) & \phi_{m-1}(x_1) & \cdots & \phi_1(x_1) & 1 \\ \phi_m(x_2) & \phi_{m-1}(x_2) & \cdots & \phi_1(x_2) & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \phi_m(x_4) & \phi_{m-1}(x_4) & \cdots & \phi_1(x_4) & 1 \end{bmatrix}$$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \cdots & x_{11}x_{12} & x_{12}x_{13} & \cdots \\ x_{21} & x_{22} & x_{23} & \cdots & x_{21}x_{22} & \cdots & \cdots \\ x_{31} & x_{32} & x_{33} & \cdots & x_{31}x_{32} & \cdots & \cdots \end{bmatrix}$$

$\phi_1(x_1)$ $\phi_2(x_1)$ $\phi_3(x_1)$

$\phi_1(x_2)$

$$x_1 \in \mathbb{R}^3, x_2 \in \mathbb{R}^3, x_3 \in \mathbb{R}^3$$

$$x_1 = \begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \end{bmatrix} \quad x_2 = \begin{bmatrix} x_{21} \\ x_{22} \\ x_{23} \end{bmatrix}$$

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

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

$b \in \mathbb{R}$

Gaussian function :-

$$\phi_m(x) = e^{-\frac{1}{2} \left(\frac{(x - \mu_m)}{\sigma_m} \right)^2}$$



→ diff.
basis
functions

$$\phi_{m-1}(x) = e^{-\frac{1}{2} \left(\frac{(x - \mu_{m-1})}{\sigma_{m-1}} \right)^2}$$



$$\phi_1(x) = e^{-\frac{1}{2} \left(\frac{(x - \mu_1)}{\sigma_1} \right)^2}$$



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

where,

$$A = \begin{bmatrix} \phi_m(x_1) & \phi_{m-1}(x_1) & \cdots & \phi_1(x_1) & 1 \\ \phi_m(x_2) & \phi_{m-1}(x_2) & \cdots & \phi_1(x_2) & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \phi_m(x_L) & \phi_{m-1}(x_L) & \cdots & \phi_1(x_L) & 1 \end{bmatrix}_{L \times (m+1)}$$

$$\text{where, } \phi_m(x_1) = e^{-\frac{1}{2} \left(\frac{(x_1 - \mu_m)}{\sigma_m} \right)^2}$$

$$\phi_1(x_L) = e^{-\frac{1}{2} \left(\frac{(x_L - \mu_1)}{\sigma_1} \right)^2}$$

Basis

Sigmoidal function:

$$\phi_m(x) = \frac{1}{1 + e^{-(\beta_m^T x + \beta_m^0)}}$$

β_m^0 is not power 0,
it is a constant.

$$\phi_{m-1}(x) = \frac{1}{1 + e^{-(\beta_{m-1}^T x + \beta_{m-1}^0)}}$$

$$\phi_1(x) = \frac{1}{1 + e^{-(\beta_1^T x + \beta_1^0)}}$$

β_1^0 is not power 0,
it is a constant.

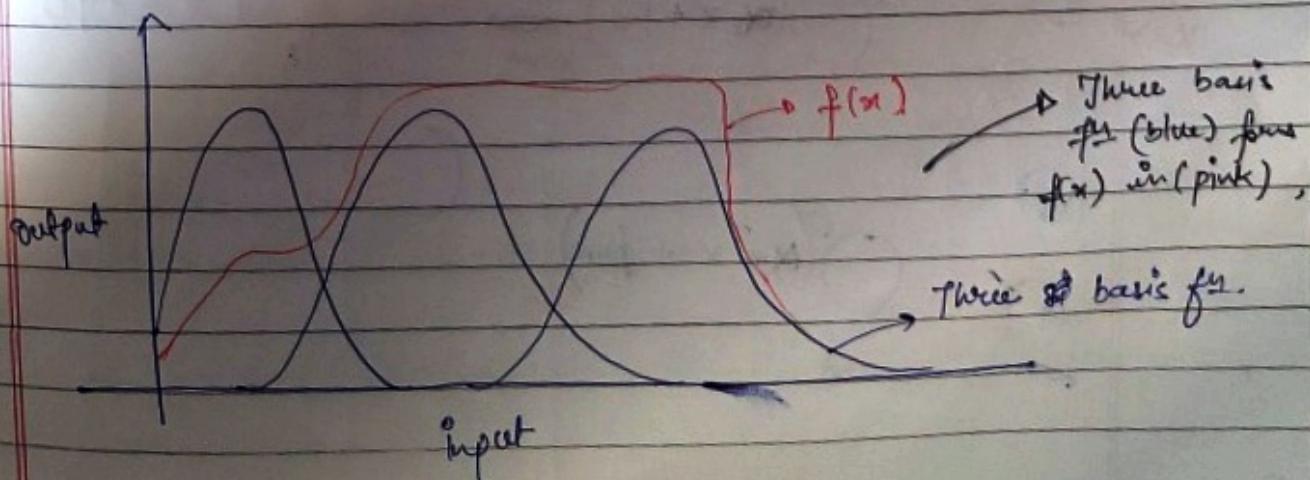
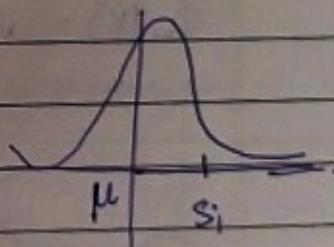
m=2

m-type of Gaussian functions

$$\phi_j(x) = e^{-\frac{1}{2s_j^2} \|x - \mu_j\|^2}$$

$$\Rightarrow \phi_1(x) = e^{-\frac{1}{2s_1^2} \|x - \mu_1\|^2}$$

$$\phi_m(x) = e^{-\frac{1}{2s_m^2} \|x - \mu_m\|^2}$$



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

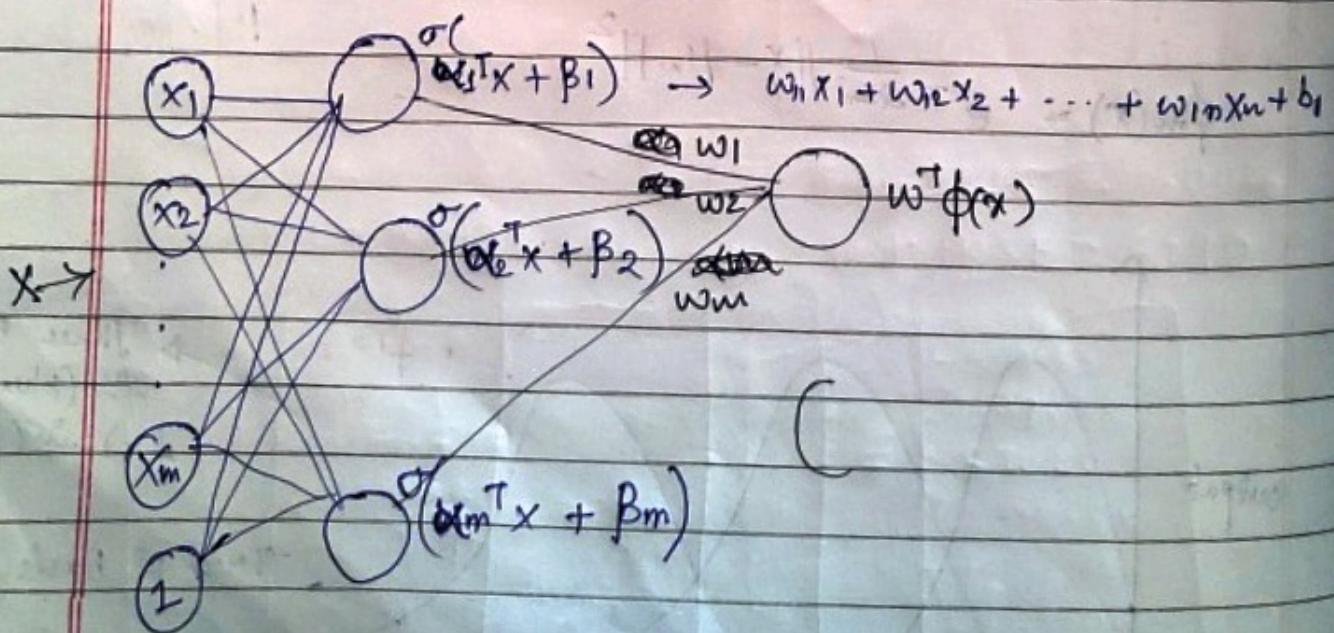
can be obtained by

Sigmoidal function:

$$\phi_1(x) = \sigma(\alpha_1, \beta_1, x) = \frac{1}{1 + e^{-(\alpha_1^T x + \beta_1)}}$$

$$\phi_2(x) = \sigma(\alpha_2, \beta_2, x) = \frac{1}{1 + e^{-(\alpha_2^T x + \beta_2)}}$$

$$\phi_m(x) = \sigma(\alpha_m, \beta_m, x) = \frac{1}{1 + e^{-(\alpha_m^T x + \beta_m)}}$$



$$w_1(\sigma(\alpha_1^T x + \beta_1)) + w_2(\sigma(\alpha_2^T x + \beta_2)) + \dots + w_m(\sigma(\alpha_m^T x + \beta_m)) + b$$

↑ bias point

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

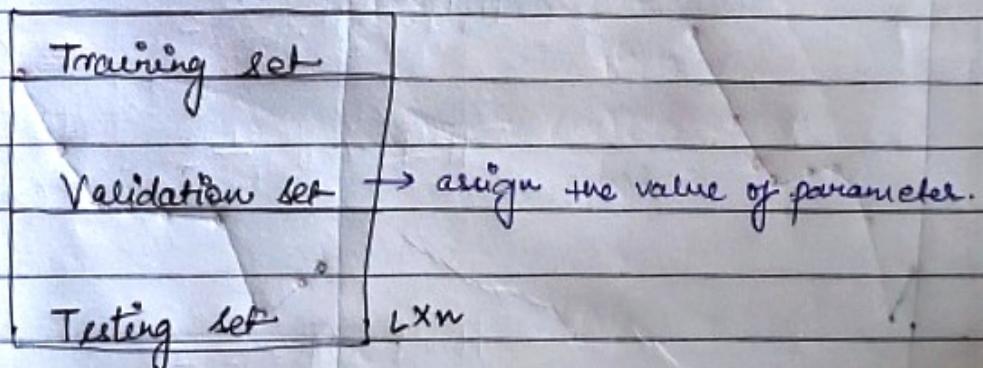
for m data pts $\rightarrow w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_m \phi_m(x) + b$

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

for $(m+1)$ data pts. $\rightarrow w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_m \phi_m(x) +$

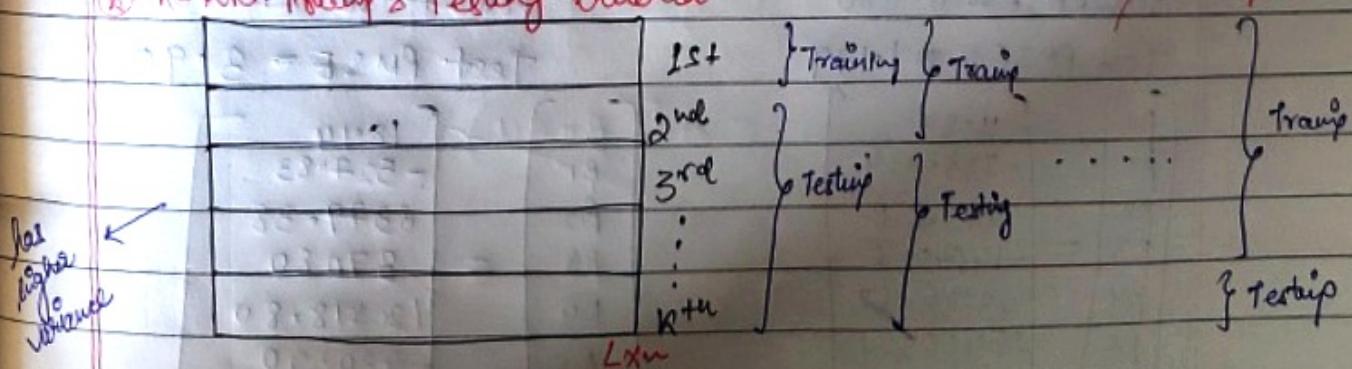
$$w_{m+1} \phi_{m+1}(x) + b$$

* if we consider $w_{m+1} \phi_{m+1}(x)$ as 1, then, for $(m+1)$ data pts., $f(x)$ will be equal to for m data pts.



K-fold Training & Testing criteria

leave one point

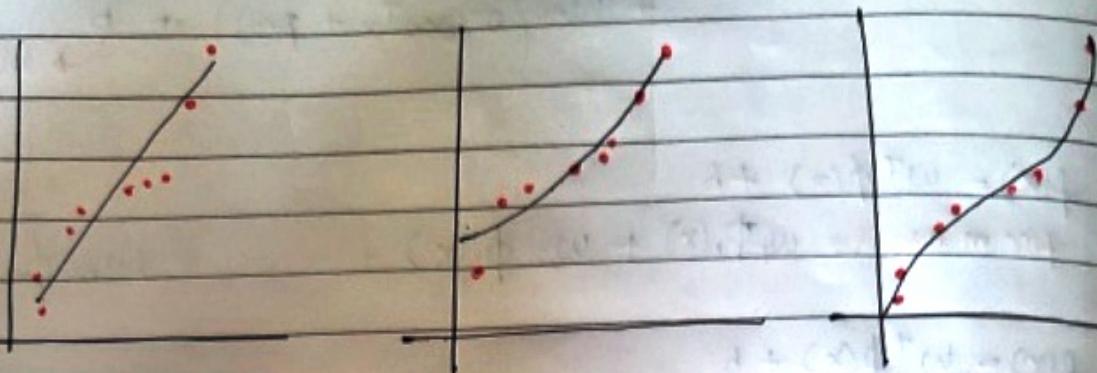


So, we have (k) RMSE, since we have K data pts.

22/2/23

Date _____
Page _____

M = 2



$$\text{Train RMSE} = 0.5947$$

$$\text{Test RMSE} = 0.9426$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} 2.6 \\ 6.9 \end{bmatrix}$$

M = 3



$$\text{Train RMSE} = 0.4980$$

$$\text{Test RMSE} = 0.7711$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} 3.36 \\ 2.51 \\ 4.58 \end{bmatrix}$$

M = 5

$$\text{Train RMSE} = 0.4997$$

$$\text{Test RMSE} = 0.9811$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix} = \begin{bmatrix} 2.62 \\ 11.14 \\ 15.58 \\ 9.17 \\ 4.26 \end{bmatrix}$$

M = 8

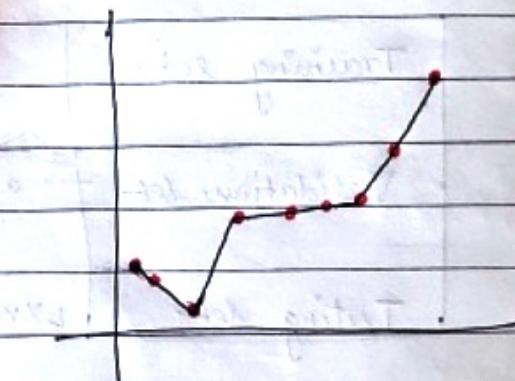


$$\text{Train RMSE} = 0.1186$$

$$\text{Test RMSE} = 1.0179$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \\ \beta_7 \end{bmatrix} = \begin{bmatrix} 11.89 \\ -283.034 \\ 3015.61 \\ -14643.7 \\ 38006.62 \\ -54565.9 \\ 40844.45 \\ -12458.6 \end{bmatrix}$$

M = 9



$$\text{Train RMSE} = 0.0026$$

$$\text{Test RMSE} = 3.90$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \\ \beta_7 \\ \beta_8 \end{bmatrix} = \begin{bmatrix} 18.14 \\ -527.83 \\ 6379.38 \\ 370.80 \\ 120518.80 \\ -230390 \\ 256860.2 \\ 154208 \\ 308542.75 \end{bmatrix}$$

* we can see that value β_i is very large.

* problem of overfitting increases exponentially with $\# \text{ of data}$.



- * higher order of β , has more or very large values.

for $M=8$,

$$\min \frac{1}{10} \left\{ \sum_{i=1}^{10} (y_i - (\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_7 x_i^7))^2 \right\} + \lambda \left[\frac{1}{2} (\beta_0^2 + \beta_1^2 + \dots + \beta_7^2) \right]$$

$\lambda = \text{regularisation parameter}$.

minimising eq. of second norm.

(*) value of λ , $f(x)$ will tend to constant.

* we \downarrow the value of $\beta_0, \beta_1, \beta_2, \dots$ by λ regularisation technique.

* we use regularisation technique.

- to control flexibility.

- to decrease value of β_i , to control overfitting.

* one of method to find $\lambda \rightarrow$ grid search method.

$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$

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

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

diff. basis function

* Optimization Problem \rightarrow

$$\min_{w, b} \underbrace{\frac{1}{2} w^T w}_{w_1^2 + w_2^2 + \dots} + \frac{1}{2} \sum_{i=1}^n (y_i - (w^T \phi(x_i) + b))^2$$

least sq.

$$\frac{d}{du} (u^T u) = 2u$$

$$\min_{w, b} \frac{1}{2} w^T I_0 w + \frac{1}{2} ((y - Aw)^T (y - Aw))$$

$$I_0 = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & & & \ddots & 0 \\ 0 & 0 & \dots & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix}$$

derivative wrt w

$$= -A^T (y - Aw)$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \\ b \end{bmatrix} \quad (m+1) \times 1$$

$$I_0 u^T u = \begin{bmatrix} w_1^2 \\ w_2^2 \\ \vdots \\ b^2 \end{bmatrix}$$

~~$$\frac{\partial (u^T A x)}{\partial x} = 2Ax$$~~

Date _____
Page _____

$$u^T \otimes I_0 u = \begin{bmatrix} w_1^2 & 0 & \cdots & 0 \\ 0 & w_2^2 & \cdots & 0 \\ \vdots & \ddots & \cdots & 0 \\ 0 & 0 & \cdots & b \\ 0 & 0 & \cdots & 0 \end{bmatrix}$$

$$u = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \\ b \end{bmatrix} = (\lambda I_0 + A^T A)^{-1} A^T \gamma$$

$$J(w) = w_1^2 + w_2^2$$

$$\nabla w = \frac{\partial J(w)}{\partial w_1} = 2w_1 = 0$$

$$\nabla w = \frac{\partial J(w)}{\partial w_2} = 2w_2 = 0$$

$$\checkmark * \quad \frac{\partial}{\partial x} (x^T A x) = 2Ax.$$

$$\checkmark * \quad \frac{\partial}{\partial u} \left(\frac{1}{2} [(\gamma - Au)^T (\gamma - Au)] \right) = \frac{1}{2} 2 \cdot (\gamma - Au)(-A)^T$$

$$= -A^T(\gamma - Au)$$

$$\checkmark * \quad \frac{\partial}{\partial u} (u^T u) = 2u$$

$$\checkmark * \quad \frac{\partial}{\partial u} \left(\frac{\lambda}{2} u^T I_0 u \right) = \frac{\lambda}{2} 2u I_0 = \lambda I_0 u$$

* Regularized Least Square Regression Model.

Date _____
Page _____

* Grid Search Method

$$\lambda^{-5} \quad \lambda^{-4} \quad \lambda^{-3} \quad \lambda^{-2} \quad \lambda^{-1} \quad \lambda^0 \quad \lambda^1$$

λ

Training set $[w]$
 b

Validation set

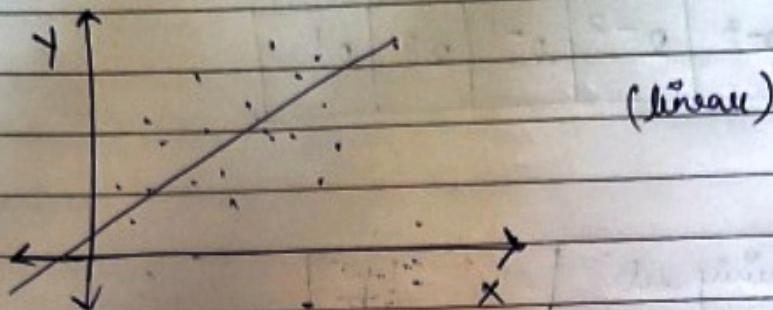
* we choose one value of λ , then perform prediction on training set & then validation set and record the value.

we then choose another value of λ , then perform same steps on Training & then on validation set.

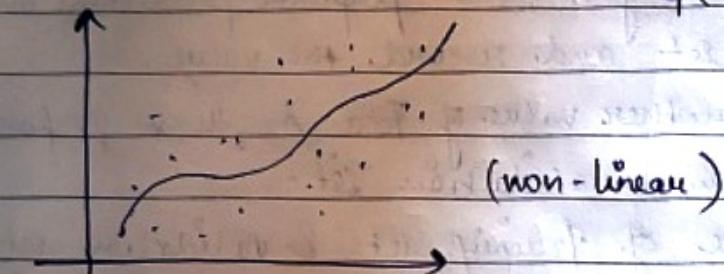
for best value of Training set & validation set ; we choose that λ for our regularised least square Regression model.

23/2/23

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



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



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

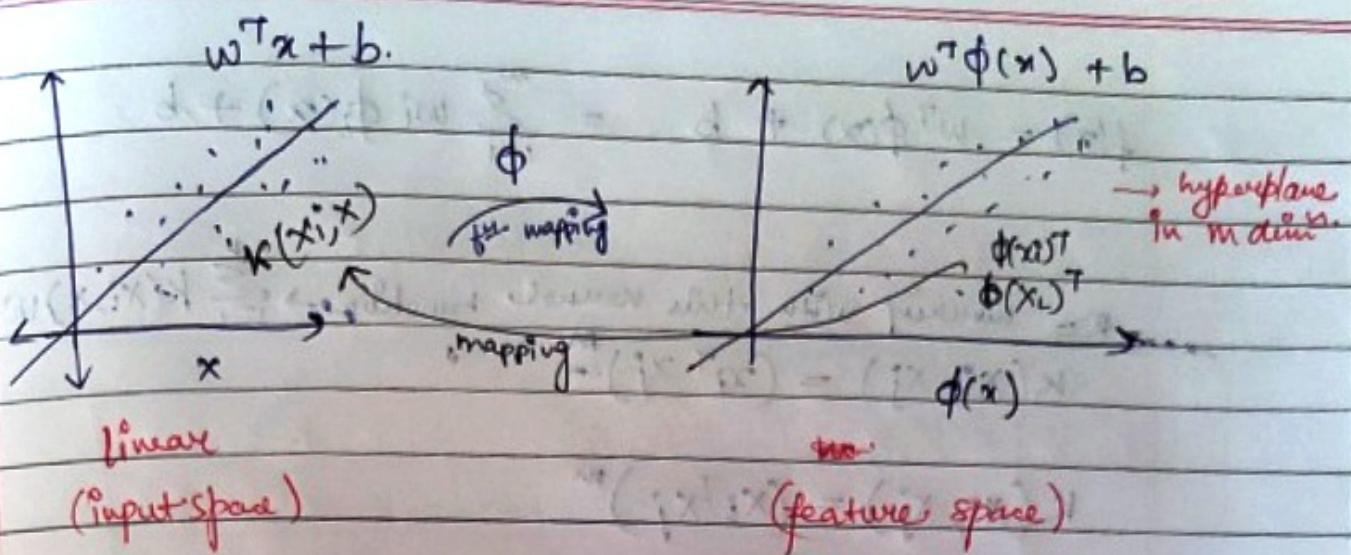
$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_m(x) \end{bmatrix}$$

m > 0

- * mapping x to $\phi(x)$ when we work with non-linear f .
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 $\phi(x)$, then k 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 combⁿ 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 dim how $\phi(x)$ looks like.

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

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

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

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

* 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

27/2/23 Fitting the linear function

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

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

x_1	x_2
-	-
-	-
-	-

ϕ

x_1	x_2	x_3
-	-	-
-	-	-
-	-	-

Kernel function

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

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

$$w_1 x_1^2 + w_2 x_2^2 + \sqrt{2} w_3 x_1 x_2 + w_4 x_1 + w_5 x_2 + b$$

↳ Linear fn.
in 5D

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

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

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

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

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

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

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

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

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

$b \rightarrow$ bias term

* For every kernel fn., we have a kernel matrix.

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

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

$$K(x, y) = (c + x^T y)^2 \quad c \rightarrow \text{constant (user defined)}$$

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



$$K(x, y) = \exp^{-\frac{\|x - y\|^2}{2\sigma^2}} \quad \sigma^2 \rightarrow \text{user defined}$$

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

→ isomorphism $\text{Kerel Space} \rightarrow \text{feature space}$.

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

$$\begin{bmatrix} w \\ b \end{bmatrix} = (A^T A + \lambda I)^{-1} A^T y$$

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

1b/23

Given the training set,

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

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

$w \in \mathbb{R}^n$
 $b \in \mathbb{R}$

We are solving.

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

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

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

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

Gradient Descent Method.

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

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

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

Repeat

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

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 basic 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 :-

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

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

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

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

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

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

Date _____
Page _____

	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

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

Kernel fn

~~K(x1, x2) = (x1^T x2)^2~~

write Gradient Descent Algo.

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

$$K = \begin{bmatrix} K(0,0) 0 & K(0,1) 0 & K(0,2) 0 & K(0,3) 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 4 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}_{4 \times 4}$$

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

$$\frac{(x_1^T x_1)^2}{\sqrt{[0]_1^2}} / \sqrt{\frac{(x_2^T x_2)^2}{\sqrt{[0]_1^2}}} / \sqrt{\frac{(x_3^T x_3)^2}{\sqrt{[0]_1^2}}} / \sqrt{\frac{(x_4^T x_4)^2}{\sqrt{[0]_1^2}}}$$

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

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

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

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

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

$$0 = (K(x_1, x_1)u_1 + K(x_2, x_1)u_2 + K(x_3, x_1)u_3 + K(x_4, x_1)u_4) + b$$

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

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

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

13/3/23

Date _____
Page _____Least Square Kernel Regression.

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

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

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

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 wst b,

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

Gradient Descent Least Square Kernel Regression

Algorithm : gradient descent method :-

Initialise $u^0 = u^{start} \in R^k$ and $b \in R$

repeat

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

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

until $\left\| \begin{bmatrix} \lambda u + \sum_{i=1}^k \frac{\partial J(u, b, x_k, y_k)}{\partial u} \\ \sum_{i=1}^k \frac{\partial J(u, b, x_k, y_k)}{\partial b} \end{bmatrix} \right\| \leq \epsilon$

Stochastic Gradient Descent Least Square Regression

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

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

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

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

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

Algorithm : Stochastic Gradient Descent Method

* At every iteration, value of η_k should be changed.

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

Repeat

Randomly select subset B from Training set T .

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

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

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 η .

$$\eta = 0.2$$

for $i = 1 : 1000$

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

end.

other methods to reduce the value of η (Ceta)

- Line search method
- Momentum concept

Stochastic Gradient Descent, Kernel Method Regression

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 .

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

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

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

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

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

- 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

$SST \rightarrow$ Sum of Sq. Training Deviation

Variance of predicted $f(x)$

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

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

- 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,

$$w_1 = 2.8$$

$$w_2 = 0$$

$$w_3 = 0$$

$$b = 0.4$$

$$\text{sparsity} = \frac{2}{3} \quad (\text{here}) = \frac{\#(w_i^o)}{\# w_i}$$

$\sum w_i^o = 0$
 $\sum \text{total } w_i = 3$

* Sparsity should be high is preferable.

& if $w_1 = 2.8$

$$w_2 = 1.2$$

$$w_3 = 0.8$$

$$b = 0.4,$$

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

* MLE

def $\{x_1, x_2, \dots, x_n\} \sim D(\theta)$
(any distibⁿ)

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

We select parameter θ , for which has max^m. prob. given sample D.
ie, $\max_{\theta} P(D|\theta)$

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

Likelihood expert value (prior info)

sometimes not given,

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

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

pts. in D are i.i.d.

$$\max_{\theta} P(D|\theta) = \max_{\theta} P(x_1, x_2, \dots, x_n|\theta)$$

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

* def $\{x_1, x_2, \dots, x_n\} \sim N(\mu, \sigma^2)$

\bar{x} is any const. then estimate μ .

$$\max_{\theta} \sum_{i=1}^n P(x_i|\mu, \sigma^2) \rightarrow \text{in case of normal distribution,}$$

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

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

derivative = 0, ie $\frac{\partial J(\mu)}{\partial \mu} = 0$

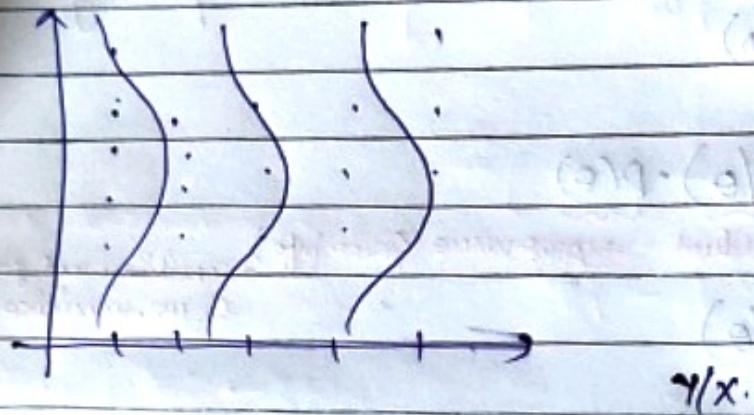
$$\frac{\partial J(\mu)}{\partial \mu} = \frac{2}{2} \sum_{i=1}^n (x_i - \mu) = 0$$

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

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

* if we take μ as const. & want to estimate σ^2 & maximize σ , then following same step we get as,

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



$$Y_i = f(x_i) + \varepsilon_i$$

$$E(Y_i | x) \rightarrow w^T x_i + b$$

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

$$Y_i | x \sim N(w^T \phi(x_i) + b, \sigma^2)$$

$$\varepsilon_i = Y_i - (w^T x_i + b) \sim N(0, \sigma^2)$$

In regression problem, we have,

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

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

$$\max P(y_1, y_2, \dots, y_L | (x_1, x_2, \dots, x_L))$$

$$\max P(y_1/x_1, y_2/x_2, \dots, y_L/x_L)$$

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

$$= \max \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(y_i - w^T \phi(x_i))^2}{2\sigma^2}} \quad (\because y_i/x_i \sim N(w^T \phi(x_i), \sigma))$$

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

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

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

$$\min_w \sum_{i=1}^L (y_i - w^T \phi(x_i))^2 \rightarrow \text{Least Squared Loss fn.}$$

Advantages :

This is eq. fn. of w & b ,

smooth fn., solving system of equation is easy.

for solving $y/x \sim N(\cdot)$, 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.