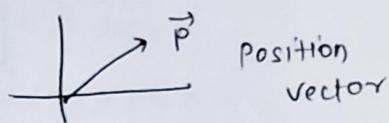


* Essential Maths for ML *

Scalar \rightarrow only magnitude

Vector \rightarrow magnitude with direction
Column vector

$\begin{bmatrix} \quad \end{bmatrix} \downarrow 1 \text{ column } (\text{for } n \text{ dimension} \rightarrow n \text{ row})$



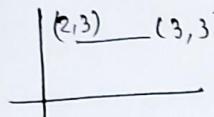
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\rightarrow \mathbf{x} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad a = 1$$

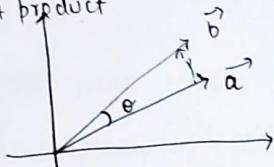
$x + a \times$ not in math

in Python (broadcast) \rightarrow convert in same dim

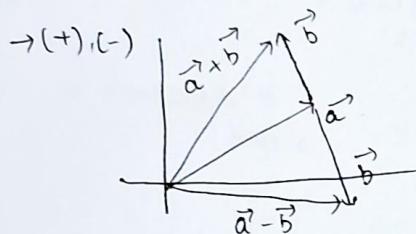
$$\text{translation } \rightarrow (2, 3) \rightarrow (3, 3) \quad \begin{bmatrix} 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$



\rightarrow Dot product



$$\vec{a} \cdot \vec{b} = ab \cos \theta = b(a \cos \theta) = a(b \cos \theta)$$



\rightarrow equation

$$\begin{bmatrix} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 = b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 = b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 = b_3 \end{bmatrix}$$

$$\Rightarrow A(\mathbf{x}) = \mathbf{b} \quad \text{vector} \rightarrow \text{vector}$$

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \rightarrow \text{transformations}$$

Dimension $\Rightarrow \mathbb{R}^n$

$$y = f(x)$$

$\hookrightarrow b = A(x) \rightarrow$ assume it transforming vector A

$$A: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$[1 \times 3] \times [3 \times 1] = [1 \times 1]$$

$w: \mathbb{R}^n \rightarrow \mathbb{R}$ for Linear classifier

$$\hat{y} = f(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$$

$$y \approx f(x) = w^T x \quad (x^T w) \rightarrow \text{not followed}$$

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}_{4 \times 1} \quad x = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{1500 \times 1}$$

$$y = \begin{bmatrix} \dots \end{bmatrix}_{1500 \times 1}$$

$$(y - \hat{y}_i)^2 = (y - w^T x)(y - w^T x)$$

$$\cancel{\frac{\partial w_i}{\partial w_i}} = -2(y - \hat{y}_i)x^T \cancel{x} = \cancel{1 \times 1500} \cancel{1 \times 1500} \cancel{x^T x} = (1500 \times 1)$$

$$\frac{\partial w_i}{\partial w_i} = -2(y - \hat{y})x_i^T \quad (1 \times 1500) \quad (500 \times 1) \rightarrow 1 \times 1 \checkmark$$

Normalisation

1) Remove Duplicate (Min max scaling)

2) Handle missing value (put 0 or other technique)

$$\text{min max scaling} = \frac{x_i - \text{min}}{\text{max} - \text{min}} \quad R \in (-1, 0)$$

$$\text{another way} \quad \frac{x_i - \text{min}}{\text{max} - \text{min}} \quad R \in (0, 1)$$

$$\text{Standard} \rightarrow \frac{x - \mu}{\sigma} \quad \text{usually } (-2, +2)$$

Replace with mean value (ignore the missing value)

(or make them 0) by $n \rightarrow \text{increase}$

Preprocessing step in ML :-

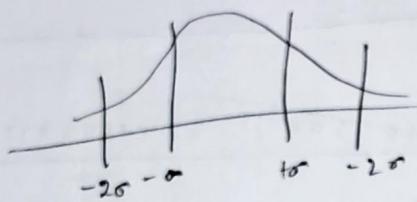
2) Normalisation

- Min-Max scaling (on feature) $\rightarrow \frac{x - \min}{\max - \min}$
- standard scaling $\rightarrow u, \sigma$ $\frac{x - \mu}{\sigma} \rightarrow (-3, 3)$

In error we have to compute $(y - (\alpha x + \gamma))^2$
→ original error
→ normalized error $\rightarrow (y_{\text{norm}} - \hat{y}_{\text{norm}})^2$ both check

→ we have to handle outlier

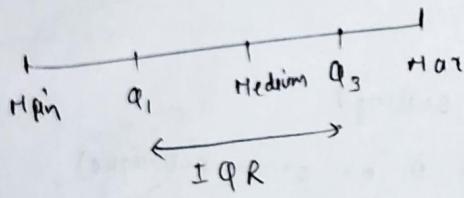
①



$$-3 > \frac{x - \mu}{\sigma} > 3$$

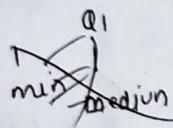
consider outlier

②



Inter quantile range
 $= Q_3 - Q_1$

if value $< Q_1 - 1.5 \text{ IQR}$
or value $> Q_3 + 1.5 \text{ IQR}$



Percentile $\rightarrow P(X \leq x)$

How many point's less than x

$$\frac{n}{N} \times 100$$

→ How to handle string.

e.g. size → large $\rightarrow 0$ } label encoding
medium $\rightarrow 1$
small $\rightarrow 2$ } cardinal data
(medium - small = large - med)

why not $\rightarrow (100 - 0)^2$ } Difference
's long \rightarrow

one hot encoding \rightarrow Sparse Representation

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

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

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

Maths \rightarrow x^2 $\frac{d}{dx} x^2 = 2x$ \downarrow scalar

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad x^T x = [x_1 \ x_2 \ x_3] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = x_1^2 + x_2^2 + x_3^2$$

$$y = x^T x \quad \downarrow x \rightarrow \text{vector}$$

$\frac{dy}{dx} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \frac{\partial y}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 2x_1 \\ 2x_2 \\ 2x_3 \end{bmatrix} = 2 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 2x$

vector derivative with
scalar \rightarrow vector

$\xrightarrow{\text{vector}} \text{3 component}$

$$\rightarrow x^T A x = y$$

$$1 \times 2 \quad 2 \times 2 \quad 2 \times 1 \quad y \rightarrow 1 \times 1 \quad \downarrow \text{scalar}$$

$$x = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \quad A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\frac{dy}{dx} = (A + A^T) x$$

\hookrightarrow if A is not symmetric

$2Ax \rightarrow A$ is symmetric

$$\rightarrow x^T A x$$

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \xrightarrow{1 \times 2} \begin{bmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{bmatrix}$$

$$a_{11}x_1^2 + a_{12}x_1x_2$$

$$2x_1 = + a_{21}x_1x_2 + a_{22}x_2^2$$

$$y = a_{11}x_1^2 + (a_{12} + a_{21})x_1x_2 + a_{22}x_2^2$$

$$\frac{\partial y}{\partial x_1} = 2a_{11}x_1 + (a_{12} + a_{21})x_2$$

$$\frac{\partial y}{\partial x_2} = (a_{12} + a_{21})x_1 + 2a_{22}x_2$$

$$\begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \end{bmatrix} = 2 \begin{bmatrix} \text{fill} & 2AX \end{bmatrix}$$

Essential math in ML

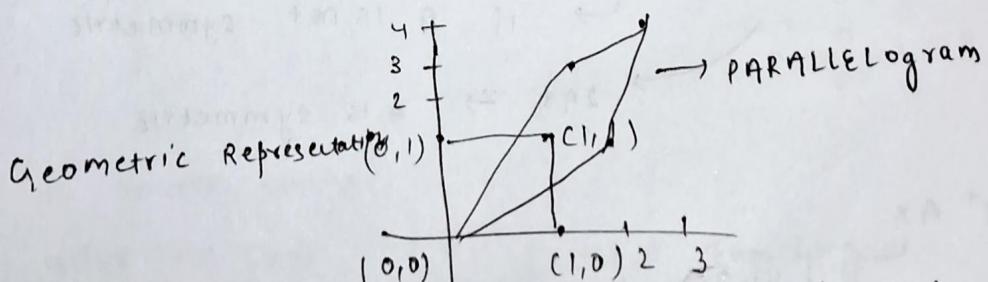
$$\begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}_{2 \times 4} = \begin{bmatrix} 0 & 0 & 2 & 2 \\ 0 & 3 & 0 & 3 \end{bmatrix}$$

2×2

square transformed to rectangle

$AX \rightarrow$ Translation (if non-diagonal is $\neq 0$)

$$\begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 3 & 1 & 4 \end{bmatrix}$$



The above operation is shear operation / rotation

The basis (axes change)

W^X
Transformations for linear regression

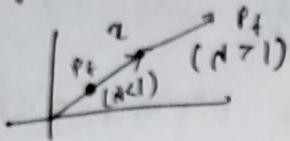
→ Gradient Descent is slow

A transformation

$\mathbf{z} \rightarrow$ eigen vector

$d \rightarrow$ eigen value

$$A \mathbf{z} = d\mathbf{z}$$



How to solve $\rightarrow (A - dI)\mathbf{z} = 0$ or $(dI - A)\mathbf{z} = 0$

$\mathbf{x} = 0$ (Trivial case)

$\det(A - dI) \rightarrow$ represent Area

$$A = \begin{bmatrix} 1 & 3 \\ 2 & 2 \end{bmatrix} \Rightarrow A - dI = \begin{bmatrix} 1-d & 3 \\ 2 & 2-d \end{bmatrix}_{2 \times 2} \text{ (so } 2d)$$

$$(1-d)(2-d) - 6 = 0$$

No, for rectangular matrix

$$2 - d - 2d + d^2 - 6 = 0 \Rightarrow d^2 - 3d - 4 = 0$$

$$d = 4, -1$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$x + 3y = -2$$

$$2x + 2y = -4$$

$$\begin{pmatrix} 2x = -3y \\ -3, 2 \end{pmatrix}$$

→ collection of vectors

$$\mathbf{X} \in \mathbb{R}^{m \times n}$$

$$\mathbf{S} = \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \\ \vdots \\ \mathbf{s}_n \end{bmatrix}$$

$\mathbf{s}_1 \rightarrow$ column vector

$$\mathbb{R}^{n \times m}$$

(column as vector)

$$[s_1 \ s_2 \ s_3 \ \dots \ s_m]$$

$$\mathbb{R}^{m \times n}$$

(row as vector)

$$\mathbf{X}^T \mathbf{X} \rightarrow \sum_{i=1}^m \sum_{j=1}^n x_i x_j = x_1 x_2 + x_1 x_3 + x_1 x_4 - x_1 x_n$$

→ Identify the eigen value & vector for this

Reduce vector:

$X^T X \rightarrow$ Eigen vector are actually

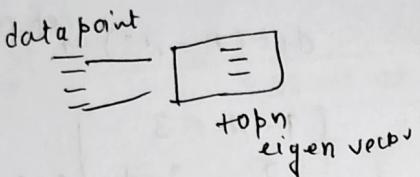
corresponding to top eigen value are actually encoding maximum information

→ Sort d_i value
→ take top s value

$$0.4 d_1 0.3 d_2 0.2 d_3 0.2 d_4 0.1 d_5$$

$$x_{\text{new}} = 0.4 d_1 + 0.3 d_2 + 0.2 d_3 + 0.1 d_4$$

compute the loss



vector Norm → for x a vector of n dimensions

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

$$\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{\frac{1}{2}}$$

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

norm ≥ 0

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

$\Rightarrow A \in \mathbb{R}^{m \times n}$

vector norm

1) Positive $|x| \geq 0$
(non-negative)

2) Homogeneous $|\alpha x| = \alpha |x|$

3) Triangle Inequality $|x+y| \leq |x| + |y|$

Matrix form:
1st norm

$$\|A\|_1 = \left(\max_{j=1 \dots n} \sum_{i=1}^m |a_{ij}| \right) \quad \begin{array}{l} \text{max from} \\ \text{column sum} \end{array}$$

$$\|A\|_\infty = \left(\max_{i=1 \dots m} \sum_{j=1}^n |a_{ij}| \right) \rightarrow \begin{array}{l} \text{max from} \\ \text{row sum.} \end{array}$$

Frobenius norm or f-norm

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

$$\|A\| \geq 0$$

$$\alpha \|A\| = \|\alpha A\|$$

$$\|A+B\| \leq (\|A\| + \|B\|)$$

$$\Rightarrow X$$

$\mu \rightarrow$ mean

$$A = X - \mu$$

$A^T A \rightarrow$ find eigen value
evector

\rightarrow sort & find top eigen val.

$$\begin{cases} x_1, x_2, \dots, x_{10} \\ a_i^\circ = x_i^\circ - \mu \end{cases}$$

find the projection of a_i on eigen vector.
 \rightarrow Dot product

$K \rightarrow$ eigen vector
 μK

$$a_{\text{const}, i} = \sum_{i=1}^K w_i x_i^\circ$$

$$a_{\text{reconst}}^{\text{update}} = a_{\text{reconst}} + u$$

Now find

$$\left\| (x_i^\circ - a_{\text{reconst}, \text{update}})^2 \right\| \quad \begin{array}{l} \text{Make Norm} \\ \text{or} \end{array}$$

EIGEN FACES

(Improved) Gradient descent:

→ Momentum based GD

(We accumulate the gradient)

Previous formula

$$w_{t+1} = w_t - \eta \nabla w_t$$

$\eta \rightarrow$ less slow & stable

$\eta \rightarrow$ more fast & unstable

$v_t \rightarrow$ Velocity

$$v_0 = 0$$

normally $\gamma = 0.9$

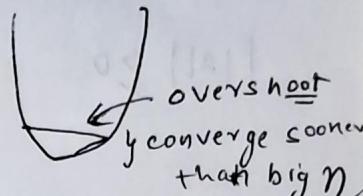
$$v_{t+1} = \gamma v_t + \eta \nabla w_t$$

$$w_{t+1} = w_t - v_{t+1}$$

$$\Rightarrow (w_{\text{prev}} - w_{\text{curr}}) \ll 0.001 \text{ stop}$$

→ Another improved version

NAG → (Nesterov Accelerated gradient)



so, they check if it overshoot if overshoot
make some change & don't overshoot

$$v_{t+1} = \gamma v_t + \eta \nabla (w_t + \gamma v_t)$$

↓
lookahead position
gradient

$$\gamma = \begin{cases} 0.9 \\ 1 \end{cases}$$

$$(0.7, 0.8, 0.9 \dots 1.0)$$

$$\eta = [0.01, 0.001 \dots 0.005]$$

$$w_{t+1} = w_t - v_{t+1}$$

⇒

Probability & statistics for ~~ML~~

Random Variable: outcome space \rightarrow all possible outcome

↳ showing result of random experiment

$$\sum_{i=1}^n p_r(x_i) = 1$$

$$P_r(e) = \frac{|E|}{|S|} = \frac{n(E)}{n(S)}$$

→ Probability Mass fn

$$P(X=1) = 1/6$$

$P(X=2) = 1/6$ by equally likely event

$\rightarrow N, 2, \mathbb{Q}, R, C$

$N \subseteq \mathbb{Z} \subseteq \mathbb{Q} \subseteq R \subseteq C \rightarrow$ non countable

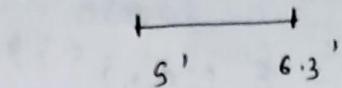
$$f(x) = \{\text{real } | x \in N\} \rightarrow \text{countable infinite}$$

\rightarrow Let have height
continuous random experiment.

\rightarrow we cannot enumerate

\rightarrow so, use probability density
function

$\int_a^b f(x) dx \rightarrow$ interval
PDF probability when $dx \rightarrow 0$



$$\text{for } \rightarrow \int_a^b f(x) dx$$

function determine the height of person

$$\text{CDF} = \Pr(X \leq x)$$

cumulative

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

$$\text{CDF} = \int_{-\infty}^x f(x) dx$$

for discrete

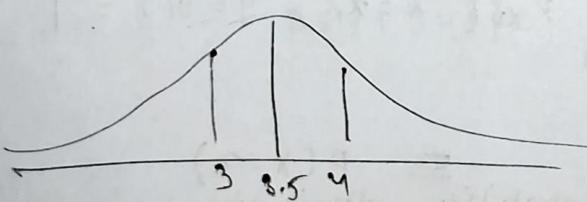
$$\Pr(X < x_i) = \sum_{i=1}^c p_i(x_i)$$

Expectation

$$E(x) = \sum_{i=1}^N x_i p(x_i) = \frac{1}{6} + \frac{2}{6} + \frac{3}{6} + \frac{4}{6} + \frac{5}{6} + \frac{6}{6} = 3.5 \quad (\text{mean})$$

Law of Large Number : \rightarrow if a experiment infinite time,
the $E(x)$ is happen

Centre limit theorem:



$$E[(x - \bar{x})]^2 = E[x^2] - (E[x])^2$$

↓
variance

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

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

$$P(A) = \sum_B P(A, B) \rightarrow \text{Marginal prob}$$

$$P(B) = \sum_A P(A, B)$$

entropy, cross entropy / information
 \Rightarrow Joint, conditional and Marginal Distribution

$$\underbrace{P_{x,y}(x=\alpha, y=\beta)}_{\downarrow \text{PMF}}$$

$P_{x,y} \rightarrow$ Joint distribution

A then B

$$\text{Rule of multiplication } P(A, B) = P(A) P(B|A)$$

$$\text{B then A } P(B|A) = P(B) P(A|B)$$

Rule of addition

$$\sum_{i=1}^n P(x_i) = 1 \quad P(A \cup B) = P(A) + P(B)$$

$$P(A \cup B) = P(A) + P(B)$$

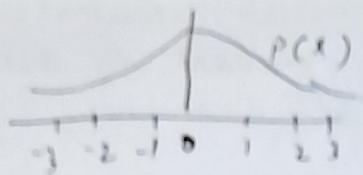
A & B are mutually exclusive

$$\begin{aligned} \Rightarrow P(x_1 x_2 x_3) &= P(x_1) P(x_2|x_1) P(x_3|x_1 x_2) \\ &= P(x_1) \frac{P(x_2|x_1)}{P(x_1)} \frac{P(x_1 x_2 x_3)}{P(x_1 x_2)} = P(x_1 | x_2 x_3) \end{aligned}$$

$$\rightarrow \underbrace{P_{x,y}(x=\alpha, y=\beta)}_{\substack{\text{sample}}}$$

$$\sum_{i=1}^N \sum_{j=1}^M P_{xy}(x=x_i, y=y_j) = 1$$

$x \sim p(x)$ \rightarrow probability \equiv distribution
 \downarrow
 sample



$$\mu = 0 \quad \sigma = 1 \quad P(X) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

PDF

(ii) $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p_{X,Y}(x,y) dx dy = 1$

X	$y=0$	$y=1$	$y=2$	
$x=0$	0.2	0.1	0.2	0.5
$x=1$	0	0.2	0.1	0.3
$x=2$	0.1	0	0.1	0.2
	0.3	0.3	0.4	1

$P(y=0)P(y=1)P(y=2)$

$$P(X=0, Y=0) = 0.2$$

$$P(X=0, Y=0) = 0.2$$

$$P(Y=0) = P(X=0, Y=0) + P(1, 0) + P(2, 0)$$

$$= 0.3$$

$$P(Y=y_j) = \sum_{x_i \in X} P(X=x_i, Y=y_j)$$

$$P(X=x_j) = \sum_{y_i \in Y} P(X=x_j, Y=y_i)$$

conditional in two's:

$$P(Y=1 | X=0) \rightarrow \text{represent row } x=0$$

$$P(X=1 | Y=1) \rightarrow \text{represent column } y=1$$

conditional distribution

$$= \frac{P(X, Y=1)}{P(Y)} \rightarrow \text{depend on } X$$

$$P(X=1 | Y=1) = \frac{P(X=1, Y=1)}{P(Y=1)} = \frac{0.2}{0.3} = \frac{2}{3}$$

for continuous:

$$P(X+Y=2) \quad P(X|Y=1)$$

\Rightarrow Information theory \rightarrow information encoded in distribution
 i) measure in bits

$$\text{Surprise} = \log\left(\frac{1}{P}\right) \quad P \approx 0$$

$$= -\log P$$

$$\text{Entropy} = E[\text{surprise}]$$

$$= \sum_{x_i \in X} \text{surprise}_i P_i (\text{surprise}_i)$$

$$= \sum_i P_i (-\log P_i)$$

$$\text{entropy} = \sum_i P_i \log P_i \quad \text{by } E[f(x)] = E[x]$$

CROSS ENTROPY \rightarrow
 KL-Divergence

Entropy:-

Represents no. of bits

$$I(x) = \log_2\left(\frac{1}{P(x)}\right) \quad (H(P) = f[I(x)])$$

$$= E[-\log(P(x))]$$

$$E[f(x)] \leq f[E(x)] \quad = E[-\log(P(x)) P(x)]$$

Jensen's Inequality

$$= -\sum_{x \in X} P(x) \log P(x)$$

coin toss

P	H	T
y_2	y_2	
1	0	
0	1	

$$\text{Entropy} = -\sum \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 1$$

$$\log 0 = \infty \Rightarrow$$

$$\text{Entropy} = -\sum 0 \log 0 + 1 \log_2 1$$

$$= 0$$

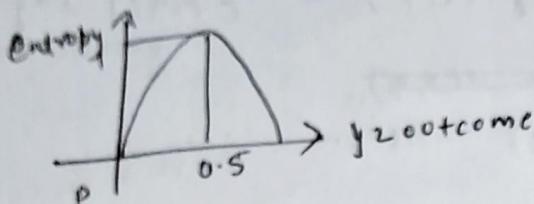
Entropy → calculate randomness

$$P(H) = 0.2 \quad P(T) = 0.8$$

$$\begin{aligned} \text{Entropy} &= - [(0.2 \log 0.2 + 0.8 \log 0.8) \\ &= - (-2.3 \times 0.2 + 0.8 \times 0.3) \\ &= 0.46 + 0.24 = 0.7 \end{aligned}$$

More random, $P(H) = 0.4, 0.6$

$$-1.3 \times 0.4 + 0.7 \times 0.6 = 0.52 + 0.42 = 0.94$$



→ Maximum entropy is possible when the distribution is uniform

For dice $\rightarrow 1/6$

Joint Entropy: $f_1(x,y) = - \sum_{x,y} P(x,y) \log P(x,y) \quad H(X|Y)$

$$H(X|Y) = \sum_{x,y} P(x,y) \log \left(\frac{P(x,y)}{P(y)} \right) \leftarrow \log P(X|Y)$$

$H \geq 0$ (non-negative)

→ Cross Entropy

→ Relative Entropy (KL Divergence)

$$CE = - \sum_x P(x) \log q(x) \quad \begin{array}{l} P, q \rightarrow \text{distribution} \\ \text{True Distribution} \quad \text{Estimated / Modeled Distribution} \end{array}$$

$$x, y, \hat{y} \quad \hat{y} = f(x) \quad \rightarrow \text{Loss fn dependent on CE}$$

Logistic Regression

true distribution $\rightarrow y$ estimated $\rightarrow \hat{y}$

$$CE = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

$$\frac{1}{1+e^{-x}} \quad \underbrace{\nabla x}_{\rightarrow w^T x}$$

KL Divergence $\rightarrow P, q$ show different 2 distribution are

$$D(P||q) = \int p(x) \log \frac{p(x)}{q(x)} dx \text{ for continuous}$$

$P \rightarrow$ true

$q \rightarrow$ estimated

$$= \sum_{x \in X} p(x_i) \log \frac{p(x_i)}{q(x_i)}$$

$$\Rightarrow -(-\sum p(x_i) \log p(x_i)) - \sum p(x_i) \log q(x_i)$$

$$= H(P) - H(P, q)$$

$$(H(P, q) - H(P))$$

Mutual Information

$$I(X; Y) = \sum_Y \sum_x p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

when independent $I(X; Y) = 0$

encode mutual info

Vocabulary of ML

T, P, E

$P_1 \rightarrow \epsilon_{t_1}$

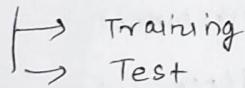
$P_2 > P_1 \rightarrow \epsilon_{t_2}$

(training experience)

$t_2 > t_1$

Testing phase

Two phase



ML algo

→ Parametric

→ Non-Parametric \rightarrow RNN

Types of ML training style)

→ Supervised \rightarrow Semi-Supervised

→ Unsupervised \rightarrow Reinforcement Learning

→ clustering

→ k-means

→ Agglomerative Clustering

→ DBSCAN (Density base)

(KNN \rightarrow reverse KNN)

$(x_1, x_2, \dots, x_n) [y]$

$$L_{\text{loss}} = \sum_i (y_i - \hat{y}_i)^2$$

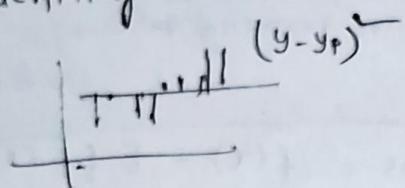
minimize over all loss

3 Data
 ↳ Training → testing → unseen
 ↳ validation → come from the same distribution

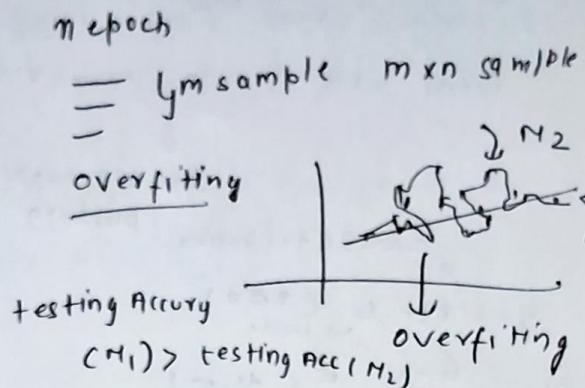
why validation →

→ overfitting

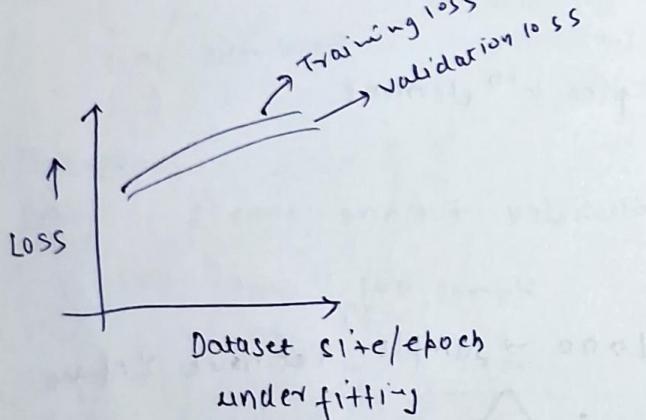
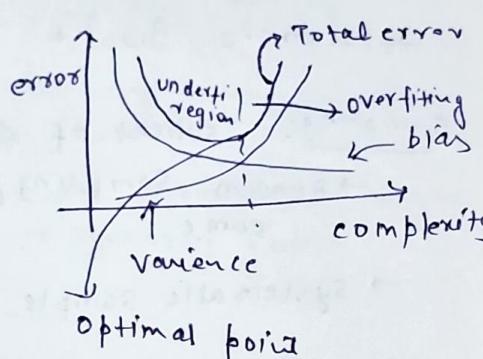
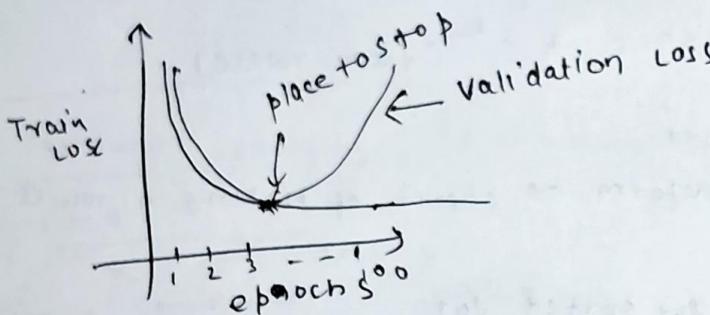
→ underfitting



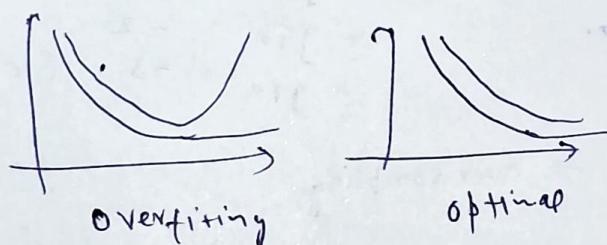
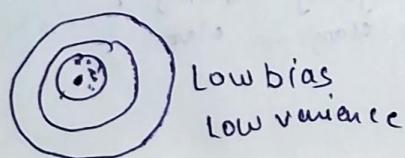
$$MSE = \frac{1}{N} \sum (y_i - \hat{y}_i)^2$$



$$MSE = \frac{1}{N} \sum (y_i - \hat{y}_i)^2$$



Low-bias & low variance



high bias low variance



Low bias High variance



High bias / high variance



Bias vs variance Trade off

Inverse proportional

Total error = Bias error + variance error + irreducible error

↓
does not capture pattern

↓
very sensitive to training data
noise

True Distribution

$$y = f(x) + \epsilon \quad (\epsilon \approx 0)$$

Noise

$$\Rightarrow \mathcal{L} \sim \mathcal{N}(u, \sigma^2)$$

$$\text{Bias} = f(x) - E[\hat{f}(x)]$$

$$\text{variance} = E[(x - u)^2] = E[x^2] - (E[x])^2$$

$$\text{variance} = E[(f(x) - E[\hat{f}(x)])^2]$$

$$\text{total error} = \text{Bias}^2 + \text{variance} + \sigma^2 \quad (\text{from noise})$$

sampling: subset of dataset.

→ Random Sampling / uniform → chance of picking a row is same

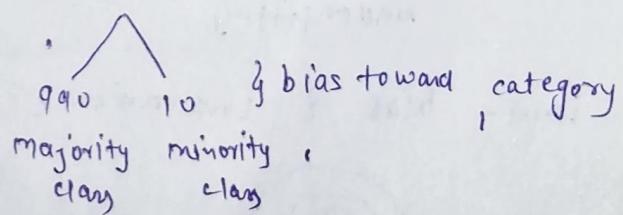
→ Systematic Sample → for sorted data

$$\begin{array}{c} \equiv \\ \equiv \\ \equiv \end{array} \begin{array}{l} \text{category-1} \\ \text{category-2} \\ \text{category-3} \end{array} \quad \begin{array}{l} \text{pick } k^{\text{th}} \text{ element} \\ \dots \end{array}$$

→ oversampling

→ repeat minority class
to make 990

1000 → sample let have 2 classes



Undersampling : undersample the majority class

SMOTE :- Synthetic Minority OverSampling Technique

$$d x_1 + (1-d)x_2 = 0$$

$$\begin{matrix} & \xrightarrow{\quad} \\ x_1 & & x_2 \end{matrix}$$

→ Take point from minority class & create new point

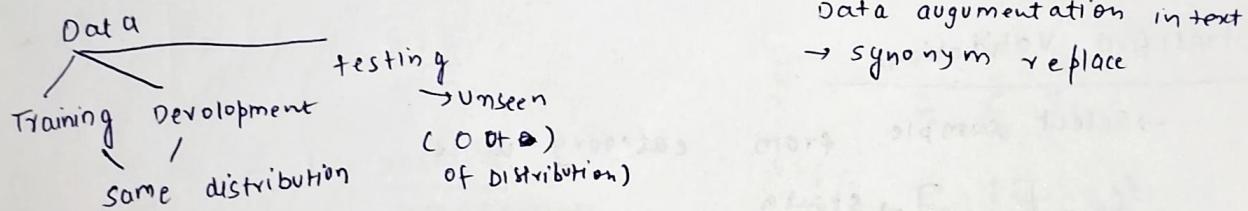
process:

take a point & check nearest and find new point using different d

⇒ Label has 3 class

c_1 } how many model take:
 c_2 1st Model (binary) → c_1 , Not c_1
 c_3 2nd Model → c_2 , Not c_2
 3rd Model → c_3 , Not c_3

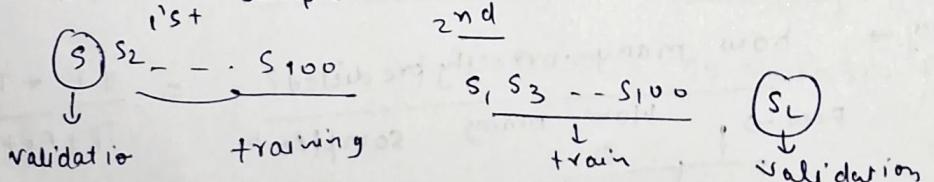
⇒ Validation Strategies ⇒ to validate model



Technique:

→ LOOV :- (Leave one out validation)

Let have 100 sample



→ 100 train & test and take average

K-fold cross validation:-

100 sample \rightarrow divide in K part

Let's 4-fold \rightarrow $\frac{P_1}{25} | \frac{P_2}{15} | \frac{P_3}{25} | \frac{P_4}{25}$

First time $\rightarrow P_1 \rightarrow \text{test} \rightarrow 25$ } Fold-1
 $P_2 P_3 P_4 \rightarrow \text{train} \rightarrow 75$

Fold two $\rightarrow P_2 \rightarrow \text{test}$
 $P_1, P_3, P_4 \rightarrow \text{train}$

Fold 3 $\rightarrow P_3 \rightarrow \text{test}$
 $P_1, P_2, P_4 \rightarrow \text{train}$

Fold 4 \rightarrow _____

\rightarrow Take average of all accuracy

dis

$\rightarrow C_1 N/K$ if data is in order
 $\rightarrow C_2 N/K$ fold $\rightarrow C_1$
 $\rightarrow C_3$ training $\rightarrow C_2 - C_1$
 \vdots
 $\rightarrow C_K \Rightarrow \text{accuracy} = 0$

Stratified Validation:

\rightarrow Select sample from category wise

$C_1 [] \rightarrow \text{strata}$

$C_2 [] \rightarrow \text{strata} \dots C_K [] \rightarrow \text{strata}$

Evaluation Technique:

\rightarrow Accuracy \rightarrow how many correctly predicted? $= \frac{TP + TN}{TP + FP + FN + TN}$

pred \rightarrow How many samples?

		0	1
0	TP	FN	
1	FP	TN	

Gold
actual)

True positive $\rightarrow TP$

false Negative $\rightarrow FN$

False Positive $\rightarrow FP$

True Negative $\rightarrow TN$

Total positive sample = TP + FN

Total negative sample = FP + TN

Total model predict positive = TP + FP

$$\text{Precision (P)} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall (R)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$F_1 \text{ score} = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2 \cdot P \cdot R}{P + R}$$

$$F_\beta \text{ score} = (\beta + 1) \frac{P \cdot R}{P + \beta^2 R} \quad \left. \begin{array}{l} \text{weight for recall more} \\ \text{than precision} \end{array} \right\}$$

For Regression

$$\text{MSE} \rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \underbrace{\hat{y}}_{\text{Mean Residual score}}$$

$$\text{Sum of Squared error} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{Mean Absolute error} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\rightarrow R^2 = 1 - \frac{SSE}{SST}$$

$$= 1 - \frac{\text{residual error}}{\text{sum of average error}}$$

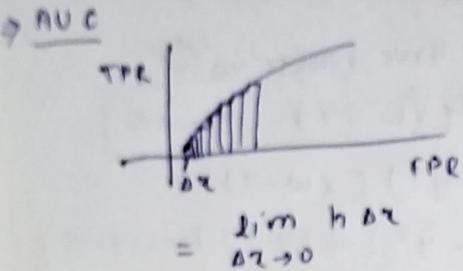
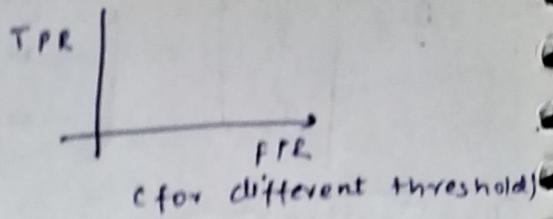
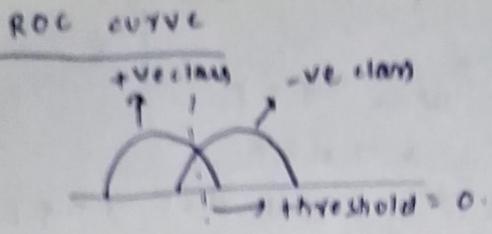
$$\bar{y} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$0 \leq R^2 \leq 1 \quad (\text{very close to 1 is good})$$

$R^2 \approx 1 \rightarrow \text{variability}$

$R^2 \approx 0 \rightarrow \text{no variability}$



TPR & FPR
→ atleast 100 points
(bins)

Supervised Learning

→ require a dataset

- feature / data
- label / classes / target / value
- classification
- Regression

ith sample $\rightarrow (x_i, y_i)$

vector | collection
of feature

Linear Regression :-

⇒ Mini Batch : medium update
Batch size $K=16$ (normally take in in power of 2)

for batch : error is

$$L = \frac{1}{2K} \sum_{i=1}^K (y_i - \hat{y}_i)^2$$

$$\Delta w = 2x \frac{1}{2} \sum_{i=1}^K (y_i - \hat{y}_i) x_i$$

Another way:



Normal equation

$$\frac{\partial L}{\partial w} = 0$$

$$\hat{y} = w^T x$$

$$\hat{y}_1 = w^T x_1$$

$$\hat{y}_2 = w^T x_2$$

$$\hat{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} w^T x_1 \\ w^T x_2 \\ \vdots \\ w^T x_n \end{bmatrix} \xrightarrow{\text{w} \rightarrow d \times 1, \text{x} \rightarrow d \times n} n \times d \rightarrow n \times 1$$

$$\boxed{(A+B)^T = A^T + B^T}$$

$$(AB)^T = B^T A^T$$

$$\frac{\partial}{\partial w} w^T a = \frac{\partial a^T w}{\partial w} = a$$

$$\frac{\partial L}{\partial w} = 0$$

If w & a has same dim(shap)

$$w = [w_1, w_2, w_3]$$

$$a = [a_1, a_2, a_3]$$

$$w^T a = w_1 a_1 + w_2 a_2 + w_3 a_3$$

$$\frac{\partial x}{\partial w} = \begin{bmatrix} \frac{\partial x}{\partial w_1} \\ \frac{\partial x}{\partial w_2} \\ \vdots \\ \frac{\partial x}{\partial w_d} \end{bmatrix}$$

$$\frac{\partial w^T a}{\partial w} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = a^T$$

$$\frac{\partial a^T w}{\partial w} = a^T$$

$$\frac{\partial}{\partial w} w^T x^T y = x^T y$$

$$\hat{y} = x w$$

$$L = (x w - y)^2 \rightarrow \text{true target value}$$

$$\begin{aligned} L &= (x w - y)^T (x w - y) \\ &= ((x w)^T - y^T) (x w - y) \\ &= (w^T x^T - y^T) (x w - y) \\ &= w^T x^T x w - w^T x^T y - y^T x w + y^T y \end{aligned}$$

$$\frac{\partial L}{\partial w} = 0$$

$$\boxed{w = (x^T x)^{-1} x^T y}$$

$$\frac{\partial}{\partial w} (w^T x^T x w - w^T x^T y - y^T x w + y^T y) = 0$$

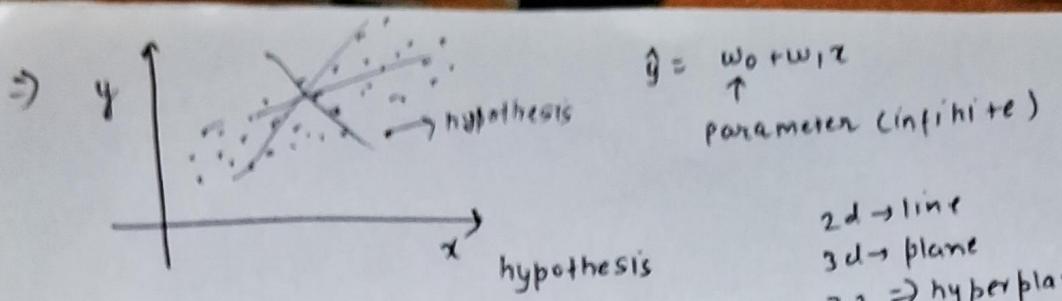
$$\Rightarrow \frac{\partial}{\partial w} (y^T x w) = \frac{\partial}{\partial w} (x^T y)^T w \\ = x^T y$$

$$\Rightarrow \frac{\partial}{\partial w} \underbrace{w^T x^T x w}_{u} \underbrace{\frac{\partial}{\partial w}}_{v}$$

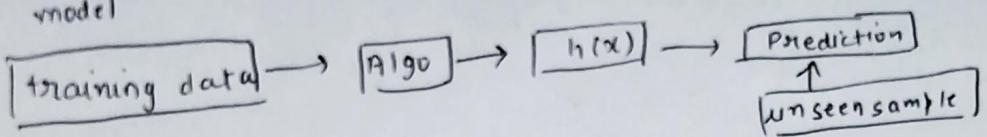
$$= w^T \frac{\partial}{\partial w} x^T x w + x^T x w \frac{\partial}{\partial w} w^T$$

$$= w^T x^T x + x^T x w I$$

$$= w^T x^T x + x^T x w$$



\Rightarrow hypothesis is a function which your algo consider while creating a model



\Rightarrow Hypothesis space:-
 \rightarrow Best fit line (Mean squared error) \rightarrow Regression

$$\Rightarrow \hat{y} = w_0 + w_1 x \quad (x_i, y_i) \Rightarrow L = (y - \hat{y})^2$$

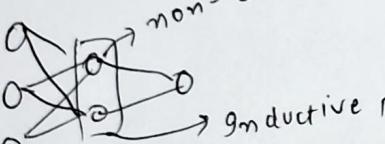
$$\Rightarrow (x_i^*, y_i^*, L_i^*) \Rightarrow \text{if plot elliptical (parabola)}$$



\Rightarrow if in 2D it get ellipse

 \rightarrow contours
 \sim in contour line gradient are same

Inductive Bias: Bias \rightarrow set of assumptions


non-linear activ fn
 \Rightarrow Inductive Bias \rightarrow generalise \rightarrow predicting label
discrete obj on unseen ex.

\Rightarrow Occam's razor \rightarrow choose the simplest hypothesis