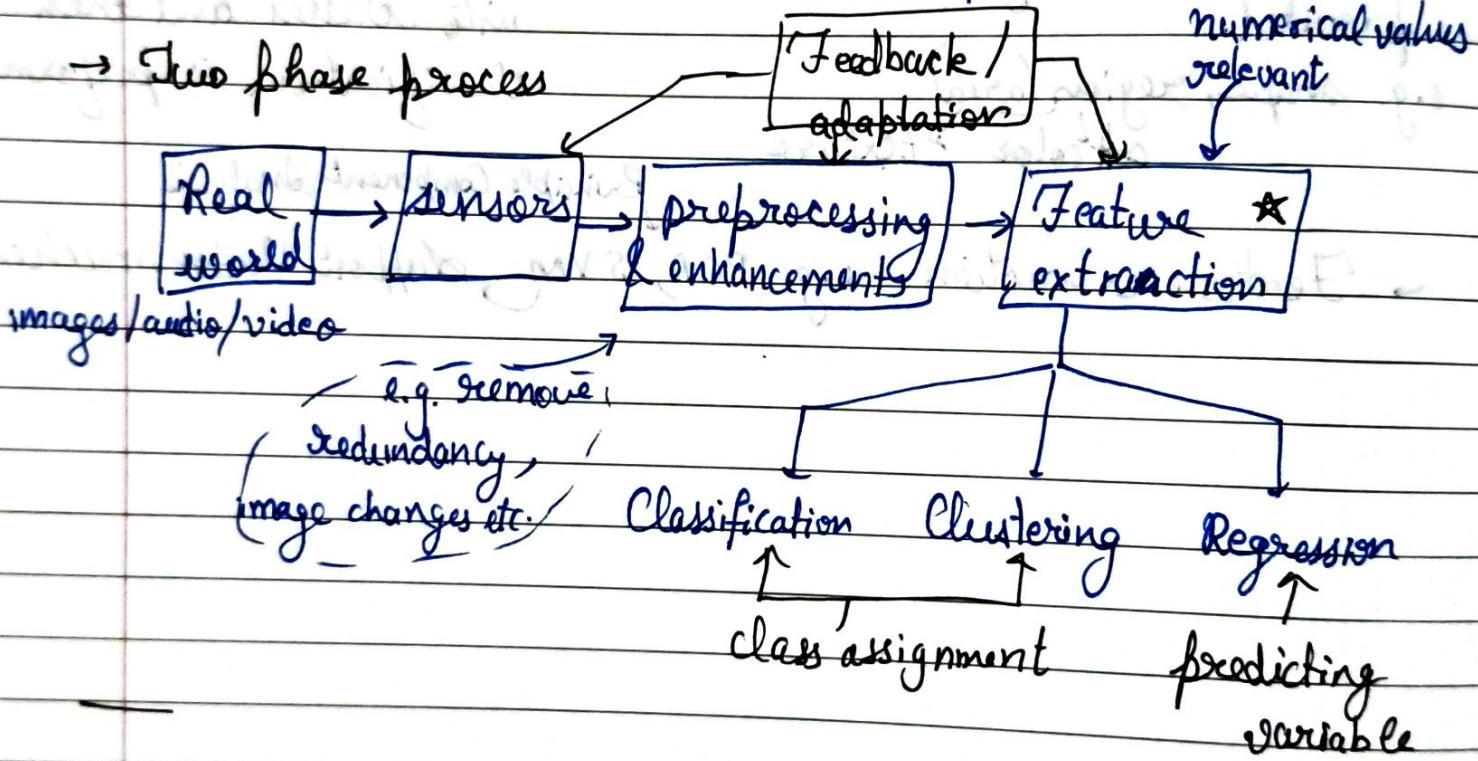


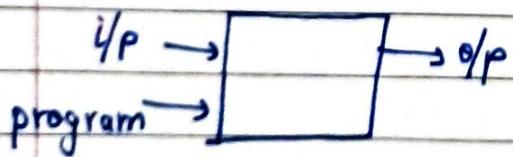
Introduction

- Pattern Recognition: identifying and analyzing patterns in data is called pattern recognition.
- Machine Learning is a type of pattern recognition done automatically by machine.
- Pattern: similarities which are not identical but can help in classification.
- There are several standard models for PR

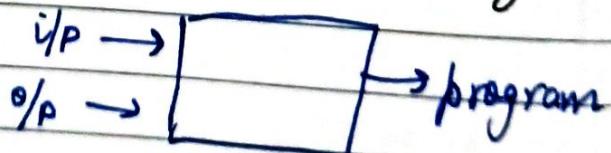
- Two phase process

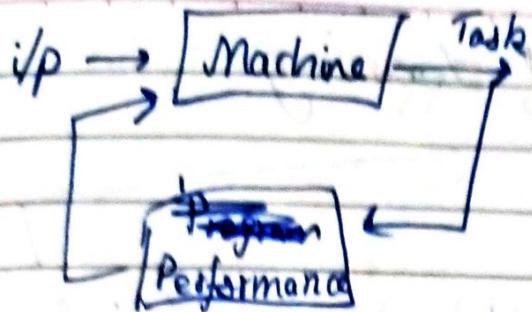


Traditional programming



Machine Learning





⇒ Advantages :-

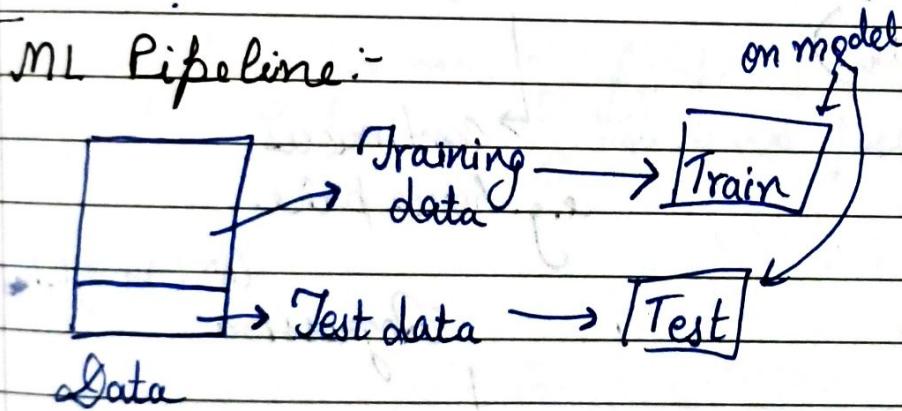
- ↳ (i) We are not explicitly programming
- ↳ (ii) Handle large data
- ↳ (iii) Make human being independent

Feature matrix :-

Features →

Samples		F_1	F_2	F_3	F_4	...
s_1						
s_2						
s_3						
:						

ML Pipeline :-



⇒ Training data : test data is generally 70:30, 80:20, 90:10

$50:50 \rightarrow 60:40 \rightarrow 70:40$ | The efficiency keeps improving. We stop when saturation is reached.

Types of Learning

Supervised

(Classification)

Unsupervised

(Clustering)

Reinforcement

→ Supervised Learning :-

↳ A teacher provides labelled training sets to train a classifier.

Supervised learning

classification
(binary/multiclass)

↳ discrete

e.g. Yes/No

regression

↳ continuous
e.g. house price

↳ Classification

↳ Logistic regression

↳ Bayesian classification

↳ KNN

↳ Decision tree

↳ Random forest

↳ Neural Network

↳ Support Vector Machine

Regression (Not in syllabus)

↳ Linear regression

↳ KNN regression

↳ Decision tree regression

↳ SVM regression

↳ Radial basis network for regression.

→ Unsupervised learning

↳ Clustering

↳ k-means clustering

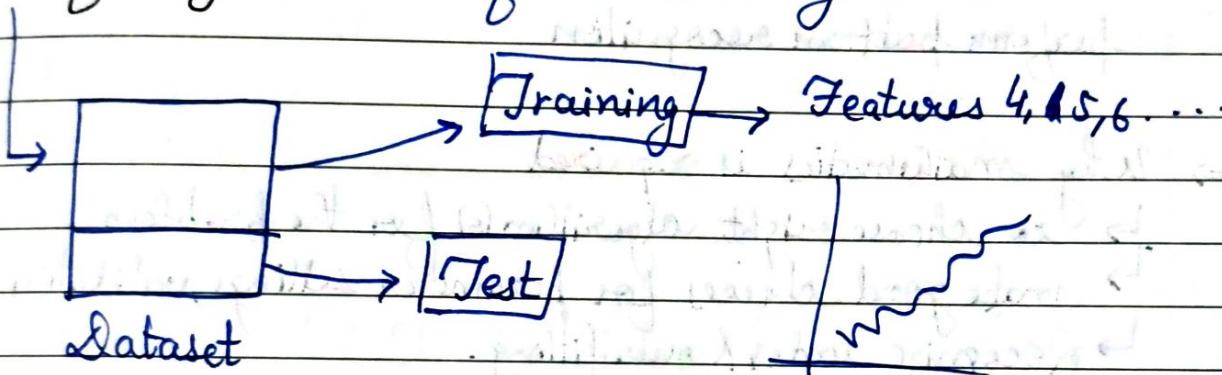
↳ k-medoid K medoids clustering

↳ Hierarchical clustering

↳ DB scan

↳ Fuzzy means clustering

→ Overfitting (Curse of dimensionality)



→ Underfitting → no. of features is too less

↳ generalization is not possible.

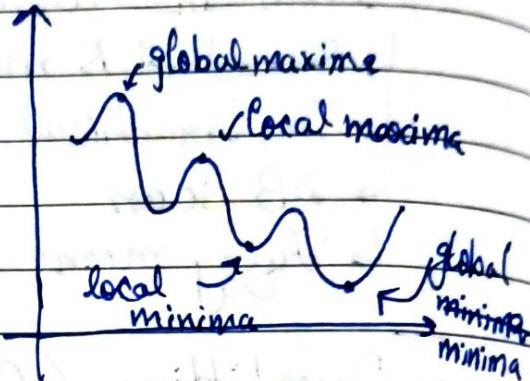
⇒ Optimization

→ Optimization is the process of training the model iteratively to maximize/minimize a value.

→ Common optimization methods are:

- Maximum likelihood
- Expectation maximization
- Gradient descent

∴ We need mathematics in order to perform pattern recognition



→ Why mathematics is required

- ↳ to choose right algorithm(s) for the problem
- ↳ make good choices for parameter settings, validation strategies
- ↳ Recognize underfitting overfitting.

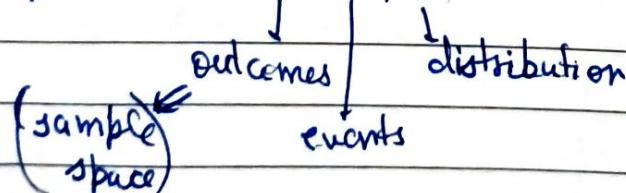
→ Mathematical notation

$|B| \rightarrow$ cardinality

$\|\cdot\| \rightarrow$ norm

$\mathbb{R}^n \rightarrow$ dimensionality

→ Probability space: It is a random process/experiment with 3 components (Ω, E, P)



Discrete space

- $|\Omega|$ is finite
- Σ is used

Continuous space

- $|\Omega|$ is infinite
- \int is used

\Rightarrow In PR, data must be unbiased / unskewed.

\rightarrow pdf : probability density function \rightarrow continuous
{ pmf : probability mass function \rightarrow discrete

\Rightarrow probability function \rightarrow maps probability to sample space
 \hookrightarrow area is always 1.

\rightarrow cdf : cumulative density function.

\hookrightarrow probability till a value e.g. $P(x \leq 3) = P(x=1) + P(x=2) + P(x=3)$

\Rightarrow Random variable

\hookrightarrow It is a function which associates number x_i with each outcome ω of a process.

\rightarrow Multivariate Probability Distribution

\hookrightarrow several random processes occur (in parallel or in sequence)

\hookrightarrow can be described as joint prob. of several random variables.

e.g. $y = a_1x_1 + a_2x_2 + \dots + a_nx_n + c$, we want to find

(i) change in y w.r.t. x_1 , change in y w.r.t. x_2

(ii) relationship b/w x_1, x_2, \dots, x_n

\rightarrow prob. given one event occurs: $P(X=x_1 | Y=y) = P(X=x_1, Y=y) / P(Y=y)$

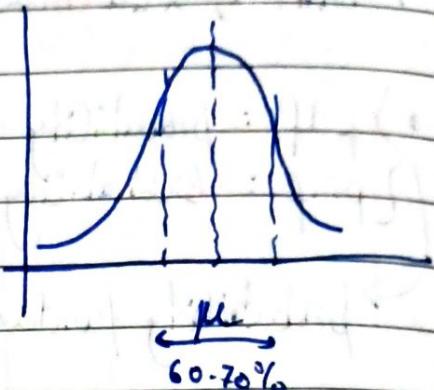
\Rightarrow same concepts of joint, marginal & conditional prob. applies for continuous \Rightarrow multivariate distribution.

before experiment after experiment

→ Expected value ~~mean~~

$$\Rightarrow E(f) = \sum_i p(x_i) \cdot f(x_i)$$

→ generally, missing values are filled with the average value



$$0.4 \times 2 + 0.6 \times 0.5 = 0.8 + 0.3 = 1.1$$

$$0.3 \times 3 + 0.9 \times 0.2 = 0.9 + 0.18 = 1.08$$

→ Expected value = mean in ideal case

→ Variance = mean value - expected value

$$\Rightarrow \text{mean} = \frac{\sum_i (x_i)}{N}$$

the spread

$$\Rightarrow \text{Var}(x) = \frac{1}{N} \sum_i (x_i - \mu)^2 = \sum_i (x_i - \mu)^2 \cdot p(x_i); \mu = \sum_i (x_i)$$

$$\Rightarrow \text{standard deviation} = \sqrt{\text{Var}(x)}$$

→ Covariance

↳ tendency for x & y to deviate from their means in same (or opposite) directions

$$\star \star \quad \frac{1}{N} \rightarrow \text{population}$$

$\frac{1}{N-1} \rightarrow \text{sample (a part of population)}$

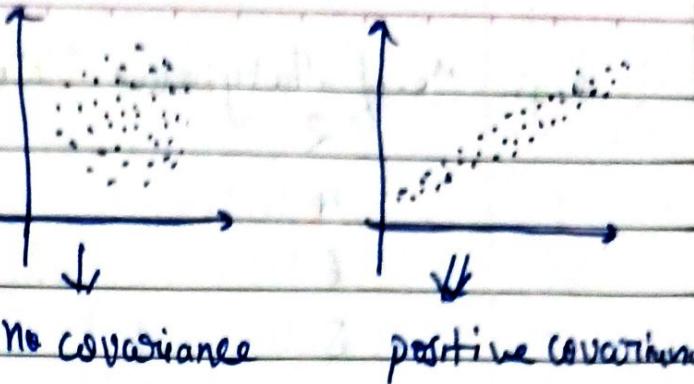
$$\rightarrow \text{cov}(x, y) = \frac{1}{N-1} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

e.g. only an
small region

$$\rightarrow f(x_i) = (x_i - \mu_x)$$

$$g(y_i) = (y_i - \mu_y)$$

$$\text{cov} = \frac{1}{N-1} \sum_i f(x_i) g(y_i)$$



→ Correlation

↳ we will only study Pearson's correlation coeff.

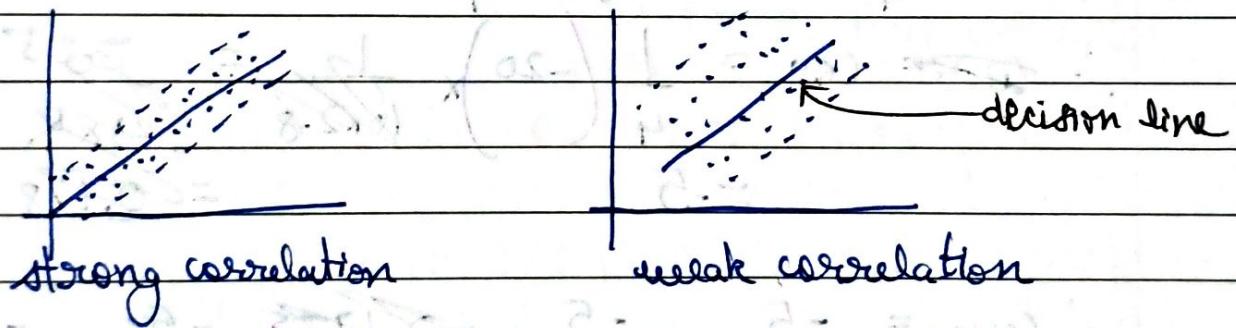
$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

↳ $-1 \leq \text{corr} \leq 1$

- || if $-1 \rightarrow 0$: negatively correlated
- if at 0 : no relation b/w features
- if 0 to 1 : positively correlated

↳ It only reflects linear dependence b/w variables.

↳ we will study techniques to convert non linear to linear relationship



$$\therefore \text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{where } \sigma_x = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2}, \quad \sigma_y = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2}$$

$$\text{and } \bar{x} = \frac{\sum_{i=1}^N x_i}{N}, \quad \bar{y} = \frac{\sum_{i=1}^N y_i}{N}$$

e.g. No. of study hours No. of sleeping hours

2	10
4	9
6	8
8	7
10	6

x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	y_i	$y_i - \bar{y}$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
2	-4	16	10	2	4	8
4	-2	4	9	1	1	2
6	0	0	8	0	0	0
8	2	4	7	-1	1	-2
10	4	16	6	-2	4	-8
<u>30</u>	<u>4</u>	<u>40</u>	<u>40</u>	<u>0</u>	<u>10</u>	<u>-20</u>

$$\therefore \bar{x} = \frac{40}{4} = 10, \quad \bar{y} = \frac{10}{4} = 2.5$$

$$\therefore \text{Cov} = \frac{1}{4} (-20) \times 2.5 = -5$$

$$\therefore \text{Corr} = \frac{-5}{\sqrt{10} \times 2.5} = \frac{-5}{\sqrt{25}} = -1$$

- Applications of linear algebra
- ↳ operations between vectors & matrices
 - ↳ coordinate transformations
 - ↳ dimensionality reduction
 - ↳ linear regression
 - ↳ solution of linear systems of equations etc.

e.g. If two features are highly correlated, we can remove one of them.

→ SVD → singular value decomposition

↳ it is used to reduce the no. of dimensions

→ Linear regression → reln

↳ Linear eqn

Homogeneous

$$AX = 0$$

↳ matrix

Non-Homogeneous

$$AX = B$$

↳ Gauss elimination

consistent

inconsistent

$$\frac{a_1}{a_2} = \frac{b_1}{b_2} \neq \frac{c_1}{c_2}$$

infinite many

unique solution

$$a_1 - b_1 = c_1$$

$$a_1 \neq b_1$$



Axioms of Probability

$$\textcircled{1} \quad 0 \leq P(A) \leq 1$$

$$\textcircled{2} \quad P(S) = 1$$

\textcircled{3} If A_1, A_2, \dots, A_n are mutually exclusive

$$P(A_i \cap A_j) = 0 \text{ then}$$

$$P(A_1 \cup A_2 \cup A_3 \dots \cup A_n) = \sum_{i=1}^n P(A_i)$$

→ Conditional probability

→ Calculating probability of an event A given that event B has already happened.

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

→ Bayes' rule :-

→ Priori probability : based on observations

→ Decision rule : rule which the model forms to classify the data.

→ It is never based on priori probability
Raw data : data captured through sensors

↳ our model cannot understand raw data

e.g. Feature vector for salmon & seabass

→ Posterior probability : $P(w_i/x_i)$

↑ ↑
class feature
vector

how much this feature is likely if output is the given class.

→ $P(x_i | w_i)$ → likelihood probability

⇒ Conditional probability =

$$P(A|B) = P(A \cap B) / P(B) \quad \text{(i)}$$

⇒ Bayes rule

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

from (i) and (ii), $P(A \cap B) = P(A) \cdot P(B|A) + P(B) \cdot P(A|B)$

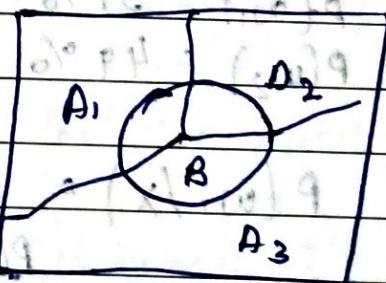
$$\therefore P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)} \quad \begin{array}{l} \xrightarrow{\text{prior}} \\ \xleftarrow{\text{likelihood}} \end{array}$$

\uparrow
posterior
probability

\nwarrow
(marginalization or
normalization)

→ Intersection of events

$$P(B) = P(A_1 \cap B) + P(A_2 \cap B) + P(A_3 \cap B)$$



$$= P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3)$$

$$\therefore P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

→ Bayes's Theorem :-

$$P(A_i | B) = \frac{P(B|A_i) \cdot P(A_i)}{P(B)}$$

where, $P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$

→ In terms of PR,

$$P(w_i | x) = \frac{P(x|w_i)P(w_i)}{\sum_{k=1}^n P(x|w_k)P(w_k)}$$

↑
Class Feature

e.g. In a country, there are 51% of adults are male, 9.5% males smoke cigar whereas 1.7% of females smoke cigars. Find probability that person smokes cigar is a male.

Ans. $P(w_1) = 51\%$

$P(w_2) = 49\%$

$P(x|w_1) = 9.5\%$

$P(x|w_2) = 1.7\%$

$$P(w_1 | x) = \frac{P(x|w_1)P(w_1)}{P(x|w_1)P(w_1) + P(x|w_2)P(w_2)}$$

$$= \frac{9.5 \times 51}{9.5 \times 51 + 49 \times 1.7}$$

$$= \frac{484.5}{567.8} = 0.85$$

→ Vector :-

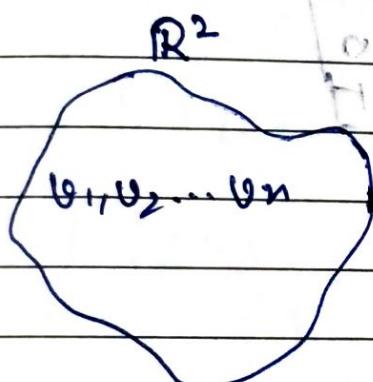
→ addition, subtraction, dot product, angle b/w 2 vectors, length, magnitude

→ Linear combination of vectors.

Consider a set $S = \{v_1, v_2, \dots, v_n\}$; then, a new vector of same space is a linear combination of vectors from same space

$$v = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n$$

e.g.



$\alpha_1, \alpha_2, \dots, \alpha_n \rightarrow$ scalars

→ Basis : dimension of vector space which can span the whole vector space

→ they must be linearly independent i.e. must be orthonormal

↳ orthogonal + same unit length

→ Cosine similarity.

→ span: Let V be a ~~is~~ vector space and let S be a subset of V , $S = \{v_1, v_2, \dots, v_n\}$. Then, span of S is the set of all linear combination of vectors v_1, v_2, \dots, v_n denoted by $\text{Span}(v)$

$\rightarrow \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n$ belongs to
span of V

e.g. A vector $\vec{v} = \begin{bmatrix} 19 \\ 10 \\ -1 \end{bmatrix}$ in $\text{span}(S)$ where $S = \left\{ \begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix}, \begin{bmatrix} -5 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 7 \\ -4 \end{bmatrix} \right\}$

$S = \left\{ \begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix}, \begin{bmatrix} -5 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 7 \\ -4 \end{bmatrix} \right\}$. Find value of α_1, α_2 and α_3 .

[Ans: $\begin{bmatrix} 4 \\ -1 \\ 2 \end{bmatrix}$]

Ans

$$\begin{bmatrix} 3 \\ -1 \\ 2 \end{bmatrix} \alpha_1 + \begin{bmatrix} -5 \\ 0 \\ 1 \end{bmatrix} \alpha_2 + \begin{bmatrix} 1 \\ 7 \\ -4 \end{bmatrix} \alpha_3 = \begin{bmatrix} 19 \\ 10 \\ -1 \end{bmatrix}$$

$$\therefore \begin{bmatrix} 3 & -5 & 1 \\ -1 & 0 & 7 \\ 2 & 1 & -4 \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} = \begin{bmatrix} 19 \\ 10 \\ -1 \end{bmatrix}$$

$$\therefore \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} = \begin{bmatrix} 19 \\ 10 \\ -1 \end{bmatrix} \times \begin{bmatrix} 3 & -5 & 1 \\ -1 & 0 & 7 \\ 2 & 1 & -4 \end{bmatrix}^{-1}$$

$$\begin{bmatrix} 3 & -5 & 1 \\ -1 & 0 & 7 \\ 2 & 1 & -4 \end{bmatrix}^{-1} = \begin{bmatrix} -7 & 10 & -1 \\ -19 & 14 & -13 \\ -35 & -22 & -5 \end{bmatrix} \times \frac{1}{(-72)}$$

$$\therefore \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} = \begin{bmatrix} 19 \\ 10 \\ -1 \end{bmatrix} \times \begin{bmatrix} -7 & 10 & -1 \\ 10 & -4 & -22 \\ -1 & -13 & -5 \end{bmatrix} \times \frac{1}{-72} = \begin{bmatrix} 4 \\ -1 \\ 2 \end{bmatrix}$$

$$\therefore \alpha_1 = 4, \alpha_2 = -1, \alpha_3 = 2$$

→ Norm of a vector

- (i) $\|x\| \geq 0$ & $\|x\| = 0$ if and only if $x = 0$
- (ii) $\|x+y\| = \|x\| + \|y\|$ → linear transformation
- (iii) $\|\alpha x\| = |\alpha| \|x\|$

$$\text{Pnorm} = \|x\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}$$

$\therefore l_1$ norm, $\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|$

ℓ_2 norm, $\|x\|_2 = \sqrt{(x_1)^2 + (x_2)^2 + \dots + (x_n)^2}$

~~l_p norm~~

ℓ_∞ norm, $\|x\|_\infty = \max(|x_1|, |x_2|, \dots, |x_n|)$

→ Similarity measures

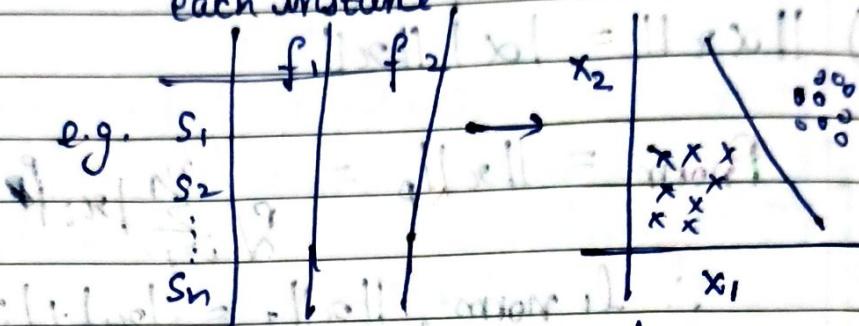
Supremum

- ① Supremum distance
- ② Manhattan distance
- ③ Euclidean distance
- ④ Minkowski distance

Bayesian Decision Theory

- Prior probability: $P(w_1), P(w_2)$ → calculate using N_1/N , N_2/N
- Likelihood probability: $P(X|w_i)$

↳ we will calculate feature vector for each instance



decision boundary
is a straight line
we can calculate likelihood.

- Posterior probability: $P(w_i|x) = \frac{P(x|w_i) P(w_i)}{P(x)}$

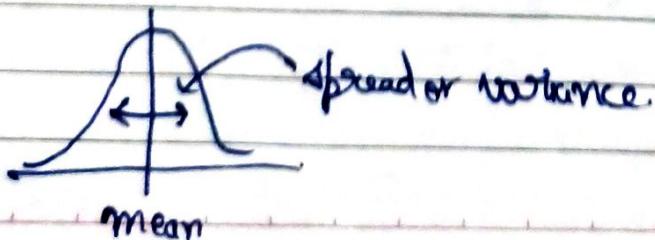
$$\text{where } P(x) = \sum_{i=1}^n P(x|w_i) P(w_i)$$

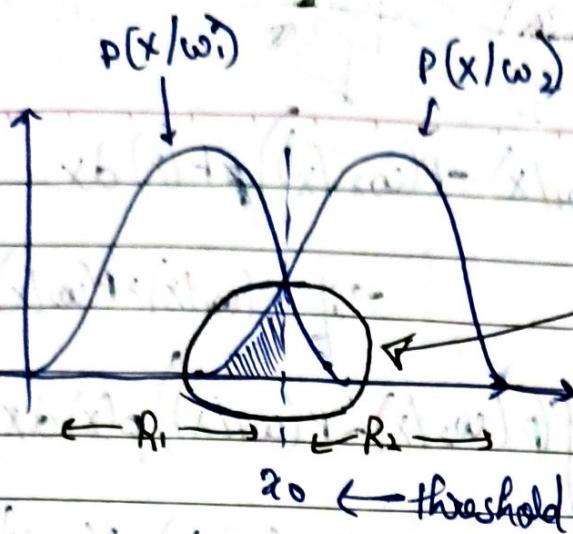
- Bayes classification: if $P(w_1|x) > P(w_2|x)$
then new sample is w_1 , else w_2

$$P(x|w_1) \cdot P(w_1) > P(x|w_2) \cdot P(w_2)$$

↳ if $P(w_1) = P(w_2) = 1/2$, $P(x|w_1) > P(x|w_2) \Rightarrow w_1$

- Most of examples in real world follow normal distribution





this portion has collision
 i.e. there is error: sample belonging to w_2 is classified into w_1 .

If $x < x_0$, it will classify in w_1 ,

$x > x_0$, it will classify in w_2

$$\therefore \text{probability error, } P_e = \left(\int_{-\infty}^{x_0} P(x|w_2) dx + \int_{x_0}^{\infty} P(x|w_1) dx \right) / 2$$

→ How we can minimise this error probability

↳ Bayesian classifier is an optimal classifier.

$$P_e = P(x \in R_2, w_1) + P(x \in R_1, w_2)$$

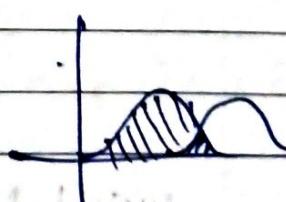
↳ i.e. $P(x \in R_2 \text{ and } x \in w_1)$ or $P(x \in R_2 \cap w_1)$

$$\text{From here, } P_e = P(x \in R_2 | w_1) P(w_1) + P(x \in R_1 | w_2) P(w_2)$$

$$\rightarrow P(x \in R_2 | x \in w_1) \cdot P(x \in w_1)$$

$$\therefore P_e = P(w_1) \int_{R_2} P(x|w_1) dx + P(w_2) \int_{R_1} P(x|w_2) dx$$

[Now, for R_1 , $P(w_1/x) > P(w_2/x)$
 R_2 : $P(w_2/x) > P(w_1/x)$]



$$\therefore P_e = \int_{R_1} P(w_1/x) P(x) dx + \int_{R_2} P(w_2/x) P(x) dx$$

[Using Bayes rule]

priori prob
↓



$$\therefore P_e = P(w_1) - \int_{R_1} (P(w_1/x) - P(w_2/x)) p(x) dx$$

→ if $P(w_1/x) > P(w_2/x) \Rightarrow$ error min
↓

$$\text{Similarly, } P_e = P(w_2) - \int_{R_2} (P(w_2/x) - P(w_1/x)) p(x) dx$$

→ if $P(w_2/x) > P(w_1/x) \Rightarrow$ error min

Therefore, Bayes decision rule \Rightarrow minimize the errors.

→ w = state of nature = class = random variable

→ $p(x)$ = probability density function = evidence

→ $p(x/w_i)$ = conditional probability density = likelihood

→ $P(\text{error}) = \min(P(w_1), P(w_2)) \leftarrow$ prior probability decision rule

Risk :-

↳ we also consider $\lambda_i (\alpha_i / w_j)$ i.e. weight function.
where, $\alpha_i \rightarrow$ action

∴ It tells how much risk is there e.g. tumor is benign or malignant

$$\rightarrow \text{risk} = R = \lambda_{12} P(w_1) \int_{R_2} P(x/w_1) dx + \lambda_{21} P(w_2) \int_{R_1} P(x/w_2) dx$$

factor

original

factor

original

weight of class 1
but assigning it to
class 2.

weight of class 2
but assigning it to
class 1.

e.g. if tumor $w_1 \rightarrow$ malignant
 $w_2 \rightarrow$ benign
 then $\lambda_{12} > \lambda_{21}$

→ In general for m classes,

$$\delta_{ik} = \sum_{i=1}^m \lambda_{ki} \int_{R_i} P(x/w_k) dx \quad \begin{matrix} \leftarrow \text{belonging to class} \\ k \text{ but classified} \\ \text{to other classes} \end{matrix}$$

$$\therefore \text{total risk}, \delta = \sum_{k=1}^m \delta_{ik} P(w_k)$$

$$\rightarrow \therefore \forall x \in R_i, \text{ if } \sum_{k=1}^m \lambda_{ki} P(x/w_k) P(c_{ik}) < \sum_{k=1}^m \lambda_{kj} P(x/w_k) P(c_{jk})$$

~~if~~ $j \neq i$
 because in that case the risk is minimum

$$\rightarrow \therefore l_1 = \lambda_{11} P(x/w_1) P(w_1) + \lambda_{21} P(x/w_2) P(w_2) \quad \} \text{for 2 classes}$$

$$l_2 = \lambda_{12} P(x/w_1) P(w_1) + \lambda_{22} P(x/w_2) P(w_2) \quad \} \text{classifier}$$

↳ 4 weights needed, where, $\lambda_{11} = \lambda_{22} = 0$.

∴ we assign $x \rightarrow w$, if $l_1 < l_2$ and vice versa

$$\therefore \text{loss matrix, } L = \begin{bmatrix} 0 & \lambda_{12} \\ \lambda_{21} & 0 \end{bmatrix} \quad \begin{matrix} \leftarrow \text{in our case } \lambda_{12} > \lambda_{21} \end{matrix}$$

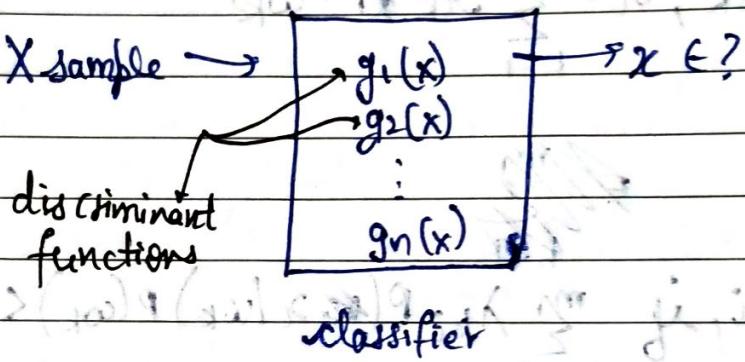
∴ $P(x/w_2) \lambda_{21} < P(x/w_1) \lambda_{12}$ in our case.

Confusion matrix :

		actual prediction	
		T	F
T	T.P.	F.P.	true positive (type-1 error)
	F	F.N.	false negative (type-2 error)
F		T.N.	true negative false positive

more dangerous

→ discriminant function



→ It is opposite to risk function

→ It is monotonically increasing function

∴ if $g_i(x) > g_j(x), x \in w_i$

$$g_i(x) = -R(\alpha_i | x) \quad \text{where } R(x | t) = 1 - p_{\text{out}}(x)$$

Discriminant function

→ It is a function which is used to create decision boundary / decision surface.

$$\text{as we know, risk, } R(x_i/x) = \sum_{j=1}^c \lambda_j (x_i/w_j) P(w_j/x)$$

$$\approx 1 - P(w_i/x)$$

$$\rightarrow \therefore \text{discriminant function } g_i(x) = \max_c \{R(x_i/x)\}$$

$$= -R(x_i/x)$$

$$= P(w_i/x)$$

→ density function

↳ it can be any F s.t. $F(\cdot)$ is monotonically increasing function then $F(g_i(x))$ is also monotonically increasing function.

$$\rightarrow g_i(x) = P(w_i/x)$$

take $F = \log$ function

$$\therefore F(g_i(x)) = \ln P(w_i/x)$$

$$\text{now, } P(w_i/x) = \frac{P(x/w_i) P(w_i)}{\underbrace{P(x)}_{\text{constant}}}$$

$$\therefore F(g_i(x)) = \ln [P(x/w_i) P(w_i)]$$

$$= \ln \underbrace{P(x/w_i)}_{\substack{\uparrow \\ \text{prob. dist. functions}}} + \ln P(w_i) \leftarrow \begin{array}{l} \text{discriminant} \\ \text{function} \end{array}$$

Prob. density function → Parametric distribution curve known.

Discrete PDF :-

① Bernoulli distribution

→ Denoted by $x \sim \text{Ber}(p)$ ↑ known parameter

→ There is one element which can have only 2 outputs

→ Prob. of success = p

Prob. of failure = $1-p$

② Binomial distribution → n independent elements
or n identical Bernoulli trials

$$P(X=x) = {}^n C_x p^x q^{n-x}$$

③ Poisson distribution

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad \left. \begin{array}{l} \text{when } n \rightarrow \infty \\ p \text{ is small} \\ np = \text{constant } (1) \end{array} \right.$$

④ Geometric distribution

→ keep trying until you get success.

$$P(X=x) = pq^{x-1}$$

e.g. if x = toss of coin & head is success

T → TT → TTT → TTH.

Success occurs only one time in all the trials.

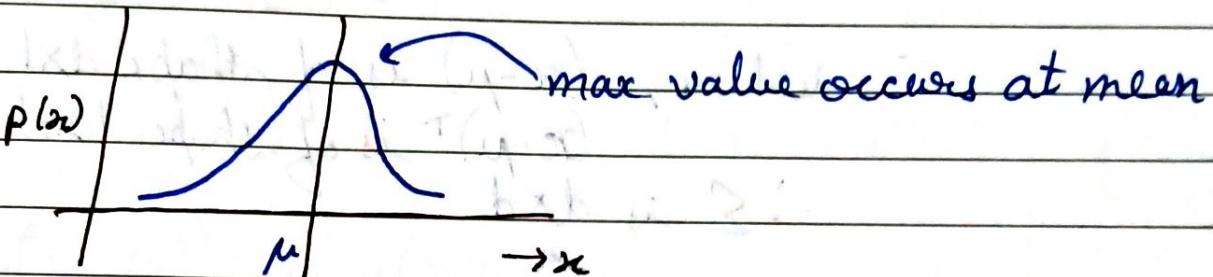
Normal distribution / Gaussian distribution

→ Represented as $X \sim N(\mu, \sigma^2)$

\uparrow mean \uparrow variance

$$\rightarrow \text{PDF}, p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

→ Curve :



$$\begin{aligned} \rightarrow \mu &= \int_{-\infty}^{\infty} x p(x) dx, \quad \sigma^2 = E[(x-\mu)^2] \\ &= \int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx \end{aligned}$$

Multivariate Normal Density

⇒ univariate $\Rightarrow x$ i.e. there is only one feature

→ But we cannot take only single feature to decide which feature belongs to which class.

→ Instead we have feature vector.

$$p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \times e^{-\frac{1}{2} \{(x-x_\mu)^T \Sigma^{-1} (x-x_\mu)\}}$$

where, μ = expected value = $E[\mathbf{x}] = \int_{\mathbb{R}^d} \mathbf{x} p(\mathbf{x}) d\mathbf{x}$

\mathbf{x} is d -dimensional feature vector

Σ is covariance matrix

$|\Sigma|$ = determinant of covariance matrix

$$\Rightarrow \Sigma = E[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T]$$

where, $(\mathbf{x} - \mu)$ is of shape $d \times 1$

$(\mathbf{x} - \mu)^T$ is of shape $1 \times d$

$\therefore \Sigma$ is $d \times d$

\Rightarrow Inner product = scalar = $(\mathbf{x} - \mu)^T (\mathbf{x} - \mu)$

$$\therefore \Sigma = \int_{-\infty}^{\infty} (\mathbf{x} - \mu)(\mathbf{x} - \mu)^T p(\mathbf{x}) d\mathbf{x}$$

$$\sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)] \rightarrow \text{covariance}$$

$$\sigma_{ii} = E[(x - \mu_i)^2] \rightarrow \text{variance}$$

Diagonals represent variance and rest are covariance

e.g. Bivariate, $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

Lets assume x_1 and x_2 are statistically independent

$\therefore \sigma_1^2, \sigma_2^2$ only exist while $\sigma_{12} = \sigma_{21} = 0$

$$\therefore \Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \wedge p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} x^T \Sigma^{-1} x}$$

$$\left(-\frac{1}{2} \left(\left(\frac{x_1 - \mu_1}{\sigma_1} \right)^2 + \left(\frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right) \right)$$

~~Final pdf~~

$$\rightarrow p(x/w_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2} \left\{ (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right\}}$$

$$g_i(x) = p(x/w_i)$$

$$\therefore f(g_i(x)) = \ln p(x/w_i)$$

$$\text{Final discriminant function} = \frac{-1}{2} \left[(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right] - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_i| + \ln p(w_i)$$

K-NN Classifier

Value nearest neighbour

→ This is the oldest method.

→ It is non-parametric i.e. doesn't assume any distribution.

→ Also called lazy learner bcz. it is slow & prepares everything only when new data comes.

e.g.	S. No	Maths	English	Result
	1	4	3	Fail
	2	6	7	Pass
	3	7	8	Pass
	4	5	5	Fail
	5	8	8	Pass
	6	6	8	?

Assume $k = 3$, 6th is the testing data.

Ans	1	4	3	$\sqrt{(4-6)^2 + (3-8)^2} = \sqrt{29}$	F
	2	6	7	$\sqrt{(6-6)^2 + (7-8)^2} = \sqrt{1}$	P
	3	7	8	$\sqrt{(7-6)^2 + (8-8)^2} = \sqrt{1}$	P
	4	5	5	$\sqrt{(5-6)^2 + (8-5)^2} = \sqrt{10}$	F
	5	8	8	$\sqrt{(8-6)^2 + (8-8)^2} = \sqrt{2}$	P
	6	6	8		

∴ Result is Pass since it is 3rd closest

Q

Height Weight Class

167	51	UW
182	62	N
176	69	N
173	64	N
172	65	N
174	56	UW
169	58	N
173	57	N
170	55	N

Test

170	57
-----	----

or?

← test data

distance

10.8

8.87

1.674

4.35

8.52

6.65

4.74

2.7

5.8

Distance

6.7

13

13.41

7.61

8.24

9.05

1.41

7

2

4.85

i.e. normal

Ans ~~mean H = 173.44, mean W = 58.67~~

For k=3, Nearest values are : N, N, UW

∴ Using k-NN, the class of testing data is N.

flamming code example

	Pepper	Ginger	Chilly	Lites
A	T	T	T	F
B	T	F	F	T
C	F	T	T	F
D	F	T	F	T
E	T	F	F	T

flamming distance
 $1+0+0=1 \leftarrow 2nd$

$1+1+1=3$

$0+0+0=0 \leftarrow \text{closest}$

$0+0+1=1 \leftarrow 2nd$

$1+1+1=3$

Test → F T T ?

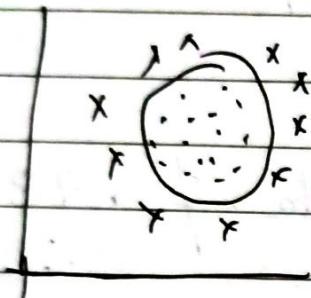
$k=3$

if $k=3$, nearest are A B C ~~and~~ and an average answer
 is False

Artificial Neural Network

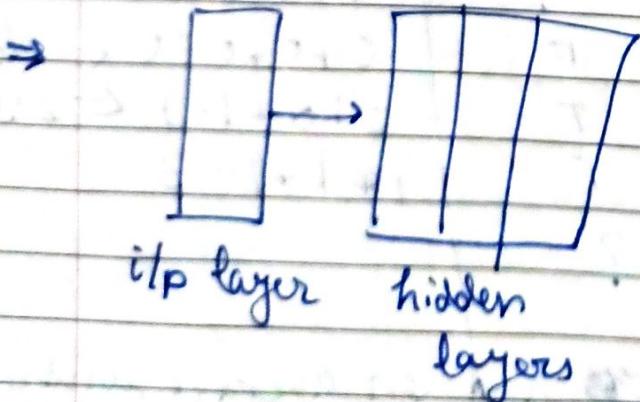
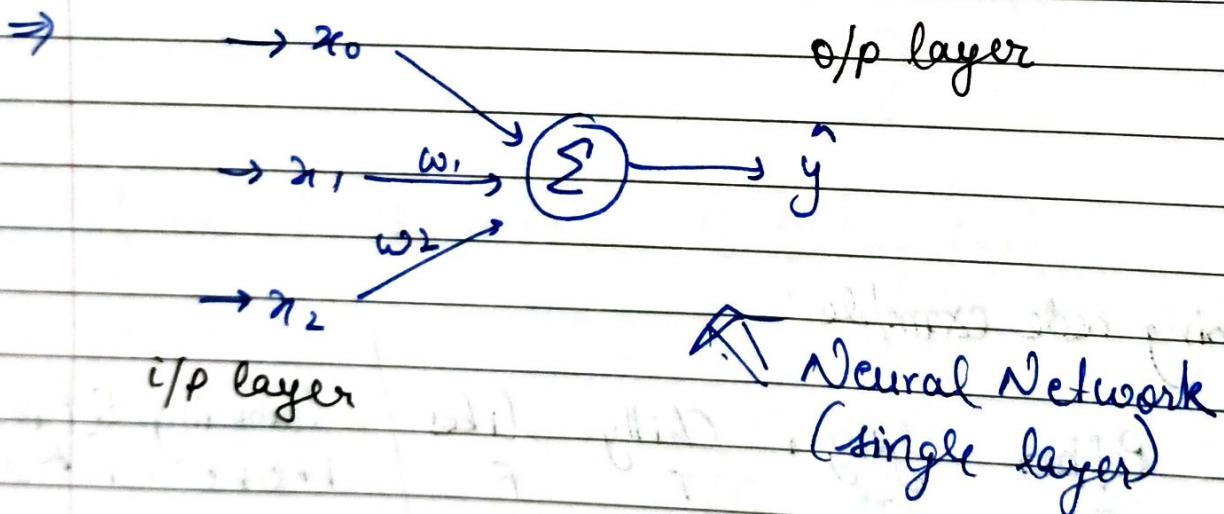


← Linear boundary can be solved using mathematics



← Non-linear boundary which is very difficult mathematics
∴ DL is used.

gn
⇒ Deep Learning, the model itself performs feature extraction by applying algo unknown to us



machine

part of human

brain

Artificial Neural Network

⇒ Why?

↳ if data is too complex, it is difficult to extract the features.

e.g.



we cannot easily determine the decision boundary.

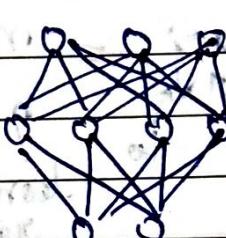


→

Input layer

Hidden layer

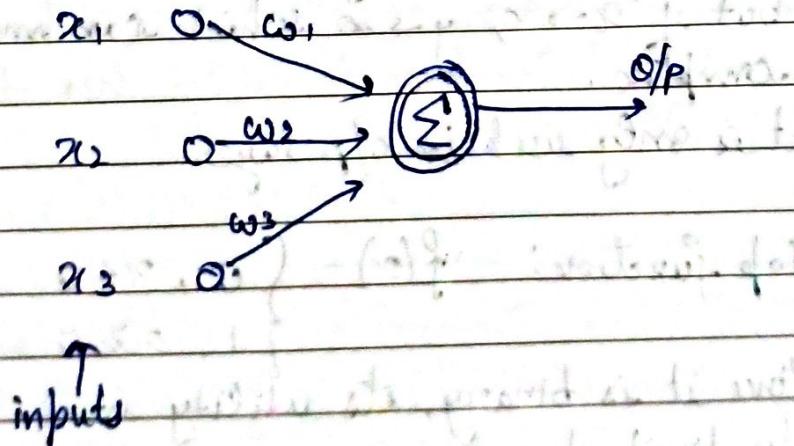
Output layer



We have linear natured lines which are interconnected

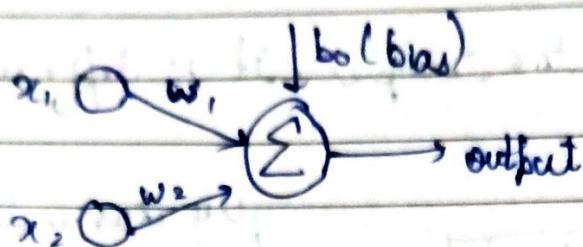
INN → 0 hidden layer

2NN → 1 hidden layer



$$\Sigma = x_1 w_1 + x_2 w_2 + x_3 w_3$$

- We want to convert linear output to non-linear output with the help of activation function.



Forward pass: i/p layer \rightarrow hidden layer \rightarrow o/p layer

Backward pass: Actual o/p is compared with predicted propagation output & we go backward to update the weight.

$$\Delta w_{ij} = \eta \delta_j O_j$$

from $\frac{\partial J}{\partial w_{ij}}$ learning rate error

\therefore o/p layer \rightarrow hidden layer \rightarrow i/p layer

Activation function :-

① Linear activation function: $y = mx + c$

but if $x \rightarrow \infty$, $y \rightarrow \infty$ which can make the system complex.

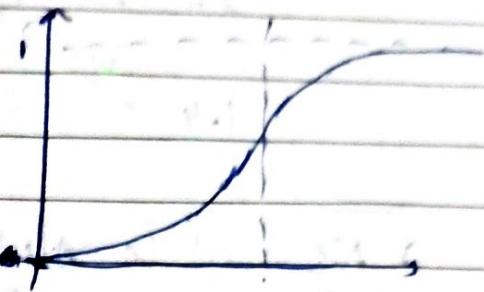
\therefore it is only used in o/p layer

② Binary Step function: $f(z) = \begin{cases} 0, z < 0 \\ 1, z \geq 0 \end{cases}$

\rightarrow since it is binary, its utility is also limited

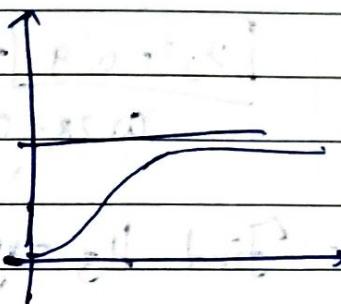
③ Sigmoid function : $A = \frac{1}{1+e^{-x}}$

- Non-linear activation function
- as $x \rightarrow \infty$, slope becomes constant
∴ we get optimized solution



④ tanh function : It is updated & more efficient than sigmoid

$$F(x) = \tanh(x) = \left(\frac{e^x - 1}{e^x + 1} \right)$$



⑤ ReLU (Rectified Linear Unit) : most popular

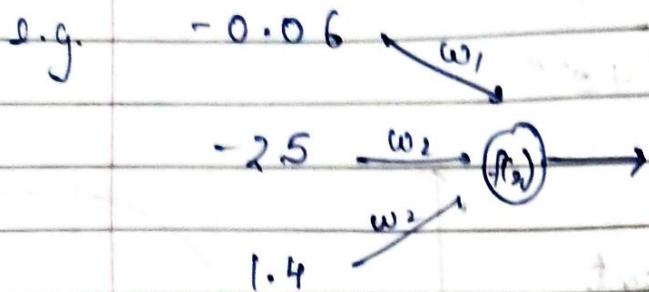
$$f(x) = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$$

Feedforward NN is unidirectional

Recurrent NN : they have directed cycles

→ Perceptron = Neural Network

→ If data has variation then learning rate is usually kept low so that it fits.

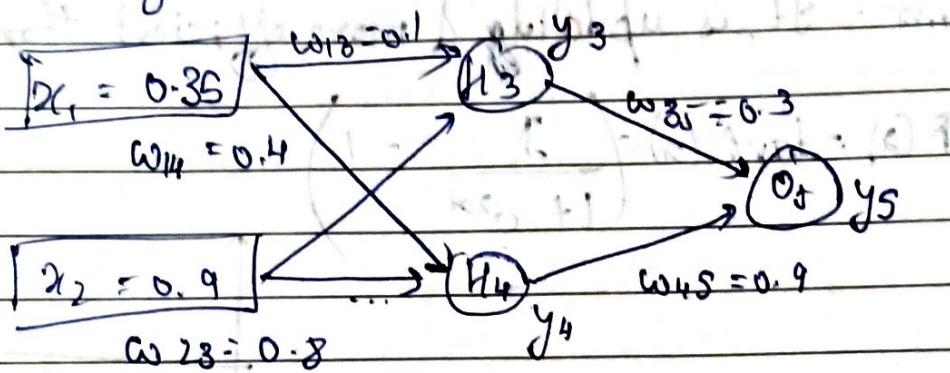


$$\text{Output} = f(y) = \frac{1}{1 + e^{-y}}$$

$$\text{where } y = -0.06w_1 + (-2.5)w_2 + 1.4w_3$$

→ No need to study how training is performed in Neural Network (G.P.L, softmax etc.)

Backprop
e.g.



Find $h_3 \rightarrow y_3$, $h_4 \rightarrow y_4$, $o_5 \rightarrow y_5$

$$w_{ji}^{\text{new}} = w_{ji}^{\text{old}} + \Delta w_{ji}$$

where $\Delta w_{ji} = \alpha \delta_j o_j$ target output

$$\Rightarrow \underline{\delta_j (\text{o/p unit})} = y_j (1 - y_j) (t_o - y_o)$$

$$\text{e.g. } \delta_5 = y_5 (1 - y_5) (t_5 - y_5)$$

$$\Rightarrow \underline{\delta_j (\text{hidden layer})} = y_j (1 - y_j) \sum w_{kj} \delta_k$$

$$\text{e.g. } \delta_3 = y_3 (1 - y_3) (w_{35} \delta_5)$$

if o_6 as well then
add $w_{36} \delta_6$ as well