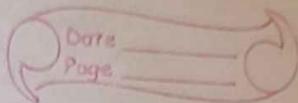
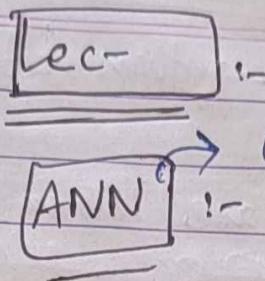


WJMK
= $\frac{S}{M}$
= MM

WJMK
= $\frac{S}{M}$
= MM



30/09/2024
Monday



artificial Neural n/w,

Binary class $\rightarrow \{0, 1\} \rightarrow \frac{1}{1+e^{-x}} \rightarrow \text{sigmoid}$

For multiclass

we have to modify sigmoid, can't use as it is.

How can change classification to regression model?

→ why classification can not be applied by regression ??

→ we want to normalize it, so that "scaling factor" can remain same!

ONE-HOT-CODING

say 10-class

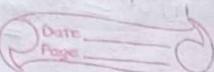
↪ $\{1, 2, 3, \dots, 10\}$

\downarrow
 $(1, 0, 0, \dots, 0)$

$(0, 1, 0, \dots, 0)$

so $y_i = (0, 0, \dots, \underset{i^{\text{th position}}}{1}, 0, 0, \dots, 0)$

WORK
30/30



→ if i^{th} position = 1 → then belongs to i^{th} class.

* Sigmoid in this case known as "SOFT MAX"

$$S = (y_1, y_2, \dots, y_n) \\ y_i \text{ in class } i$$

y_i in class - 2

In soft max, for each position, compute!

$$\sigma(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\sigma(z_1) = \frac{e^4}{e^4 + e^6 + e^7 + e^8 + \dots}$$

formula not changing, representation we are changing, so that same formula can be used.

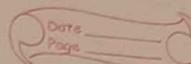
99% data in 1 class → Data Imbalance

G-them model with accuracy = 99%

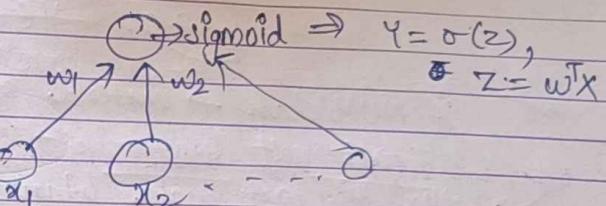
class kernels

↓
but get bad model

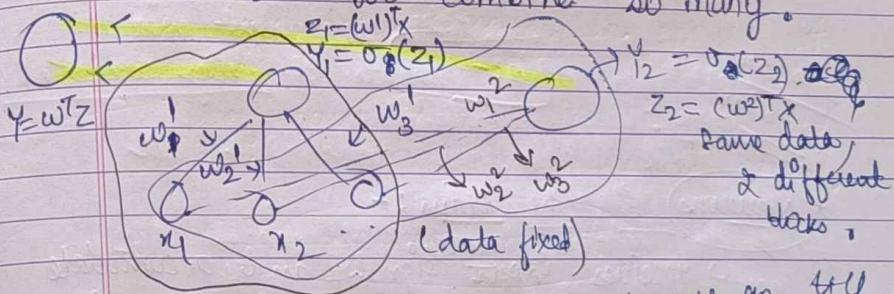
WORK
30/30



If data biased then model will always be biased choose kuch bhi karo.



ANN → also known as multi-perceptron, as we combine so many.



we go till
 $z_2 \dots$
can go till
 $z_{100} \text{ etc.}$

single hidden
layer and 2 nodes
in it

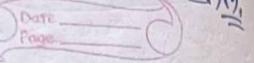
we can ↑ hidden layers as
well nodes in
each hidden layer.

→ What change in ANN classification is that
from regression to gaaye?
→ sigmoid itado!

WOMK
= 30
MM

WOMK
= 30
MM

WOMK
= 30
MM



saare htaayein? yaa sirf last wala?

↑
Dono kr skta hain,
as then ~~ye~~ 10. 4 E 10.

If multiclass → instead of sigmoid,
put soft max.

→ better ki last wala htaayein! → as
andar say w_1 jyada; then
wka weightage jyada ho jaayega, so
andar sigmoid alao!

7/10/24
Monday

lec

* Neurons in the brain :-

↳ receives input signals and accumulate
voltage. After some threshold, they will
fire spiking responses.

* A simpler Neuron :- for neural nets, we use
simpler model for neuron :-

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

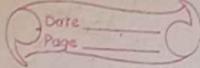
output .

inputs

weights

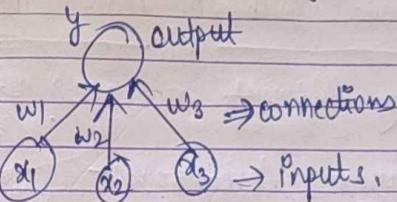
activation function

WOMK
= 30
MM



similar to logistic regression ⇒

$$y = \sigma(w^T x + b),$$



y output

w_1 w_2 w_3 → connections

x_1 x_2 x_3 → inputs.

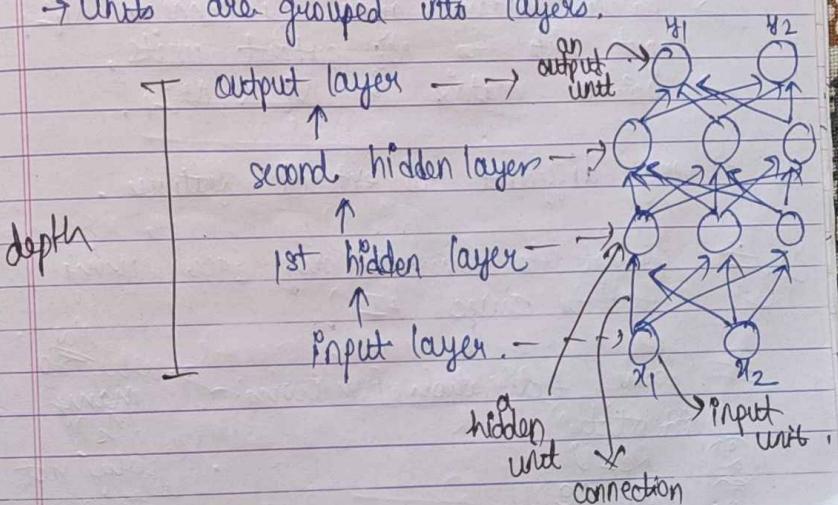
→ By throwing together lot of these simple
neuron-like processing units, we can do
some powerful computations!

→ sigmoid function is non-linear.

* A feed-forward Neural N/w

→ a DAG (Directed acyclic graph)

→ Units are grouped into layers.



layers left
brain
I don't know
WORKS
Hyperparameter, what else have to learn
It through algo.

size/only no. of nodes in each layer

* layers ↑ have same no. of nodes
kya difference in model?

one model with 'n' nodes and 'l' layers,
2nd — (n) More

Ques Are they different or same?
one model with 3 nodes in each layer (2 layers)
one model with all 6 nodes in 1 layer

Multilayer Perceptrons:-

→ outputs are function of input units :-
 $y = f(x) = \phi(Wx + b)$

→ multilayer net consists of fully connected layers.
 ↳ applied component wise

↳ all inputs are connected to all output units.

→ Each hidden layer i connects N_{i-1} input units to N_i output units.

→ weight matrix is $N_i \times N_{i-1}$.

→ very difficult to extract features.

Ques How to generate my face, using "text to Image" AI?

- # Some Activation Functions:-
- (1) Identity
 - (2) ReLU
 - (3) soft ReLU
 - (4) Sigmoid
- many fns?
why not only sigmoid?

as each layer can have different no. of nodes.

WORKS

Date _____
Page _____

when sigmoid is powerful, then why so many?

Reason:- Non-Convex loss fn's.

Difference?

here gradient = 0

soh stop ho jayega.

Gradient vanishing problem.
to remove this problem

(1) Hard threshold

(2) logistic

(3) ?

Computation in Each Layer:-

Each layer computes a fn.

$$h^{(1)} = f^{(1)}(x) = \phi(W^{(1)}x + b^{(1)})$$

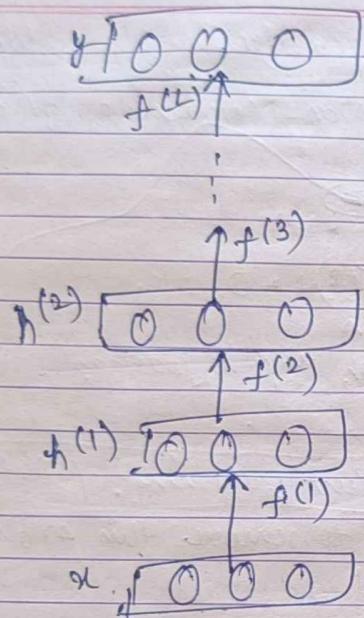
$$h^{(2)} = f^{(2)}(h^{(1)}) = \phi(W^{(2)}h^{(1)} + b^{(2)})$$

$$y = f^{(L)}(h^{(L-1)})$$

$$\text{If regression} \Rightarrow y = (w^{(L)})^T(h^{(L-1)}) + b^{(L)}$$

$$\text{If classification} \Rightarrow y = \sigma((w^{(L)})^T(h^{(L-1)}) + b^{(L)})$$

Sigmoid (if binary), soft max (if multiclass)



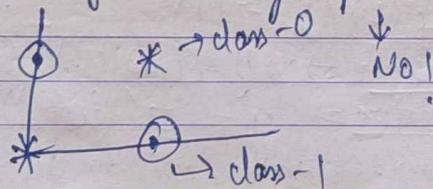
video is sequence of image
Image is not sequence

Q. # Truth table for XOR:-

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

Now classify this!
 ① using linear regression
 ② using ANN

→ will they be linearly separable?



→ Not linearly separable, so go to higher domain / dimension → where they are linearly separable.

feature learning

Note :-
Using ANN, can learn any fn,

why? Sigmoid / activation fn is making it so much powerful.

Expressive power of linear N/w's.

Designing an N/w to classify XOR
 N/w → or two
 one hidden layer,

again
 solo
 same
 weight
 deliver
 kya
 hogा?

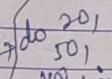
2-3 nodes in that
 hard threshold activation fn / w/o
 activation fn.

- ① Manually
- ② code

Back propagation of error
 9/10/24

Back propagation :-

getting ' ψ ' in feature learning is a difficult task.

WORK 
50
Date _____
Page _____

WIMK
Date _____
Page _____

Expressivity of logistic Activation functions :-

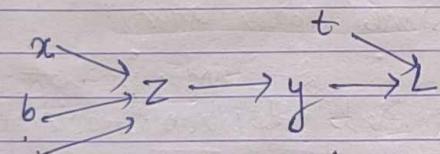
→ at last layer → we ~~can't~~ sigmoid or soft max function as activation function, don't use ReLU etc.
they can be used in intermediate layers only,
↳ WHY?

gradient descent

$$w_{i+1} = w_i - \eta \left(\frac{\partial L}{\partial w_i} \right)$$

Hyperparameter
(we've to fix)
usually keep small to let train

Logistic least squares:-



$$z = wx + b$$

$$y = \sigma(z)$$

$$L = \frac{1}{2} (y - t)^2$$

$$\frac{dz}{dx} = \left(\frac{dz}{dy} \right) \left(\frac{dy}{dx} \right) \text{ [Univariate chain rule]}$$

$$\# \text{ as: } L = \frac{1}{2} (\sigma(wx+b) - t)^2$$

$$\frac{\partial L}{\partial w} = \frac{\partial}{\partial w} \left(\frac{1}{2} (\sigma(wx+b) - t)^2 \right)$$

WIMK
Date _____
Page _____

gradient for b

$$= (\sigma(wx+b) - t) \sigma'(wx+b)x$$

$$\frac{\partial L}{\partial b} = (\sigma(wx+b) - t) \sigma'(wx+b)$$

→ Disadvantages of this approach?

Now gradient for 'w' →

$$\frac{\partial h}{\partial w} = \left(\frac{\partial h}{\partial y} \right) \left(\frac{\partial y}{\partial w} \right)$$

$$= \left(\frac{\partial L}{\partial y} \right) \left(\frac{\partial y}{\partial z} \right) \left(\frac{\partial z}{\partial w} \right)$$

$$= (y - t) \sigma'(z)(x)$$

$$= (\sigma(wx+b) - t) \sigma'(wx+b)x$$

Computing loss →

$$z = wx + b$$

$$y = \sigma(z)$$
$$L = \frac{1}{2} (y - t)^2$$

gradient for b:-

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial b}$$

$$= (y - t) \sigma'(z)(1)$$

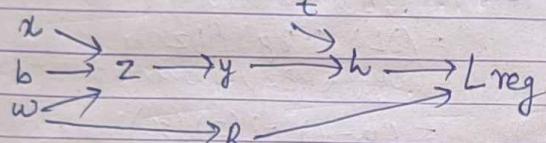
→ if layer $T_1 \rightarrow z, z', z''$... so on, derivative will also change accordingly. (abhi jitni computation kri, that is only for one hidden layer)

Computing derivation :- $\bar{y} = (y - t) = \frac{\partial L}{\partial y}$

$$\bar{z} = \bar{y}(\sigma'(z)) = \frac{\partial L}{\partial z} = \frac{\partial L}{\partial y} \left(\frac{\partial y}{\partial z} \right)$$

$$\left(\frac{\partial L}{\partial y} \right) \left(\frac{\partial y}{\partial z} \right) \left(\frac{\partial z}{\partial w} \right) = \frac{\partial L}{\partial w} = \bar{w} = \bar{z}x, \quad b = \bar{z}$$

L_2 - Regularized Regression:-



$$R = \frac{1}{2} w^2$$

$$h_{reg} = h + \lambda R$$

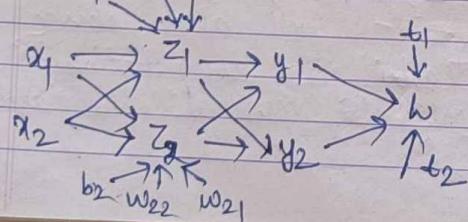
Softmax Regression:-

$$z_e = \sum_j (w_j x_j + b_e)$$

$$y_k = \frac{e^{z_k}}{\sum_l e^{z_l}}$$

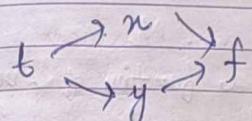
Cross-Entropy Loss

$$l = - \sum_k t_k \log(y_k)$$



Multi-Variate chain Rule:-

$$\frac{d}{dt} (f(x(t), y(t))) = \frac{\partial f}{\partial x} \left(\frac{dx}{dt} \right) + \frac{\partial f}{\partial y} \left(\frac{dy}{dt} \right)$$

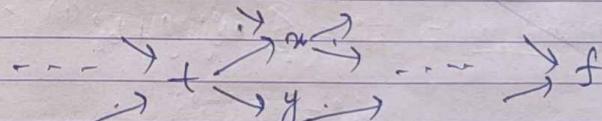


① ex:- $f(x, y) = y + e^{xy}, \quad x(t) = \cos(t), \quad y(t) = t^2$

$$\begin{aligned} \frac{d}{dt} (f(x, y)) &= \frac{\partial f}{\partial x} \left(\frac{dx}{dt} \right) + \frac{\partial f}{\partial y} \left(\frac{dy}{dt} \right) \\ &= (y e^{xy})(-\sin t) + (1 + x e^{xy})(2t) \end{aligned}$$

$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \rightarrow$ values already computed by our program.
 $\frac{dx}{dt}, \frac{dy}{dt} \rightarrow$ expressions to be evaluated.

Notation $\rightarrow \bar{x} = \bar{x} \frac{dx}{dt} + \bar{y} \frac{dy}{dt}$



black box \rightarrow as hidden layer k has node ka kitne-ka contribution in final output?
kitne sleeping nodes? \rightarrow jinko no contribution

WJMK
Date _____
Page _____

→ so node active hi nahi hai!

→ kya unko kisi aur kaam pe lga sake hai,
w/o affecting others.

VERY DIFFICULT 😐

Backpropagation Algorithm :-

① forward pass → for $i = 1, \dots, N$

compute o^i as a function of parents (w^i)

② backward pass → for $i = N-1, \dots, 1$

Backpropagation for Regularized Logistic Least Square :-

forward pass :- $z = wx + b$, $y = o(z)$

$$L = \frac{1}{2} (y - t)^2, R = \frac{1}{2} w^2$$

$$L_{reg} = L + \alpha R$$

backward pass :- $L_{reg} = 1$

$$R = \frac{\partial L_{reg}}{\partial R} = (\text{L}_{reg})(\alpha)$$

Backpropagation for Two-layer Neural Net :-

Forward pass :- $z_i^{(1)} = \sum_j w_{ij}^{(1)} x_j + b_i^{(1)}$

$$h_i^{(1)} = o(z_i^{(1)})$$

$$L = \frac{1}{2} \sum_k (y_k - t_k)^2$$

$$y_k = \sum_i w_{ki}^{(2)} h_i^{(1)} + b_k^{(2)}$$

WJMK
Date _____
Page _____

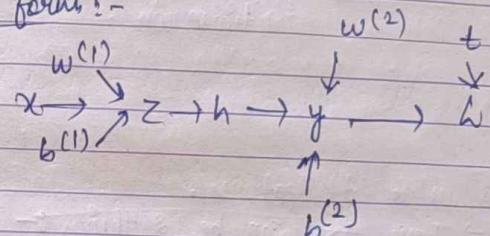
Backward pass :- $t = 1, \bar{y}_k = t(y_k - t_k)$

$$\bar{w}_{ki}^{(2)} = \bar{y}_k h_i^{(1)}$$

$$\bar{b}_k^{(2)} = \bar{y}_k, \bar{h}_i^{(1)} = \sum_k (\bar{y}_k)(\bar{w}_{ki}^{(2)})$$

$$\bar{z}_i^{(1)} = \bar{h}_i^{(1)} (\sigma'(z))$$

③ In vectorised form :-



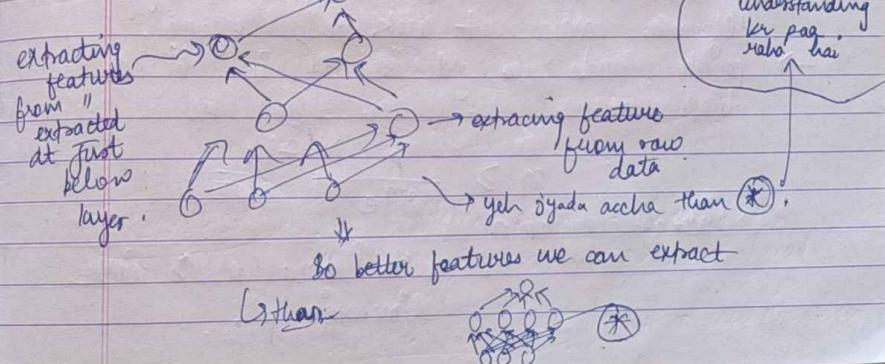
Forward pass :-

$$z = w^{(1)}x + b^{(1)}, h = o(z), y = w^{(2)}h + b^{(2)}$$

$$L = \frac{1}{2} \|t - y\|^2$$

Backward pass :-

$$\begin{aligned} t &= 1, \bar{y} = \bar{t}(y - t), \bar{w}^{(2)} = \bar{y} h^T, \\ \bar{b}^{(2)} &= \bar{y}, \bar{h} = (w^{(2)})^T \bar{y}, \bar{z} = \bar{h} o \sigma'(z), \\ \bar{w}^{(1)} &= \bar{z} x^T, \bar{b}^{(1)} = \bar{z} \end{aligned}$$



WJMK
Date _____
Page _____

WJMK
Date _____
Page _____

14 Oct 2024
Monday

(lec) :-

$$\{x_i, y_i\}_{i=1}^n$$

$n \rightarrow$ total no. of training data

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

$d \rightarrow$ attribute / quality / feature.

$d \rightarrow$ dimension

↓
may be
independent /
dependent
or anything.

→ we've ' d ' features.

→ There, ~~are~~ can be 2 things :-

- (1) n ^{very} large, or,
- (2) d very large.

not curse
1st curse
sparse data
dimension increases

Curse of dimensionality

→ As you ↑ dimension, points tend to stay in outer region / periphery.

→ (distance) → as measure of similarity / difference may not work any longer!
points tends to get isolated in higher dimensions.

hardly any m
difference in
similarity

$$n \geq \left(\frac{p}{d \pi e} \right)^{p/2} (\text{SPTT})$$

↑ how ' n ' is with dimension ' p '.

WJMK
Date _____
Page _____

WJMK
Date _____
Page _____

Fluctuations accumulate :-

$$E [\| f(x_1) - f(0_1) \|^2] \leq E [|e_1|^2] = \sigma^2$$

where,

$\sigma^2 \rightarrow$ variance of noise

→ In ' p ' dimension, this error becomes ' $p\sigma^2$ '.

→ This error can be very large in high-dimensional settings, even if σ^2 is small.

→ point

→ every data is a vector (from origin)

↳ as surf tip make hair vector k_i ,
is like 'point' bola.

→ In higher dimensions, numerical computations can become very intensive.

Circumventing the curse of dimensionality :-

Assumptions on which techniques are based :-

Need to go to lower dimensions.

Features are not INDEPENDENT,

find synthetic features which are independent, low in number.

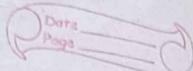
① a) points lying in high dimension are actually embedded in low dimension.
known as b) points locally lie in Euclidean dimension space.

Riemannian
Manifold
Assumption!

(low dimensional space)

like our earth is 'flat' but locally 'flat' but

WEEK 5



- ① An unsupervised learning task :-
Dimensionality Reduction

$$x \in \mathbb{R}^M \quad M < D$$

↑ f_o

Useful for:- $x \in \mathbb{R}^D$

- 1) Data compression (lossy)
- 2) Dataset visualization (2D or 3D)
- 3) discovering most important features.

② PCA (Principal components Analysis) :- (statistical interpretation)

Want to
data → project to lower dimension
(Want features)

↳ statistical independent
↳ orthogonal features

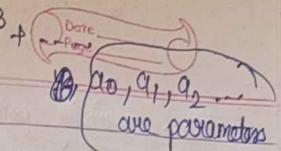
Assumption of PCA :-

- 1) features in high 'd' are independent, so
(a) can project to low dimensions.
(b) features are independent
(c) features (m) = linear combination of features in higher dimension

* Linear regression.

↳ linear combination of non-linear
↳ linear in vector
↳ linear in y and parameters.

WEEK 5



→ not linear in y and x .

like, every signal is linear combination of sinusoidal, f_n 's.
↳ non-linear.

↳ orthogonal

* PCA = basis set of features.

- PCA projects data in least square sense.
↳ It captures big (principal) variability in the data and ignores small variability.
- PCA is way to reduce data dimensionality.
- PCA projects high dimensional data to a lower dimension.

→ say $x^i, i=1, 2, \dots, N$ data points in ' p ' dimensions (' p ' is large)

$$\bar{x}_0 = \bar{m} = \frac{1}{N} \sum_{i=1}^N x_i^i$$

derivation about mean is minimum

but quite uninformative.

thus mean like to represent

→ so represent data by straight line :-

$$x = m + ae$$

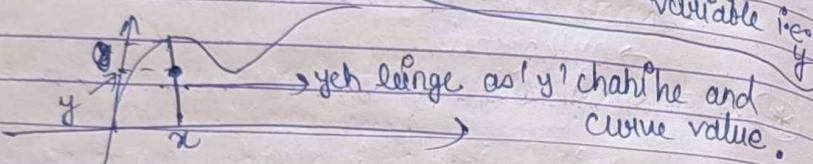
e → unit vector along straight line

a → signed distance of x from m

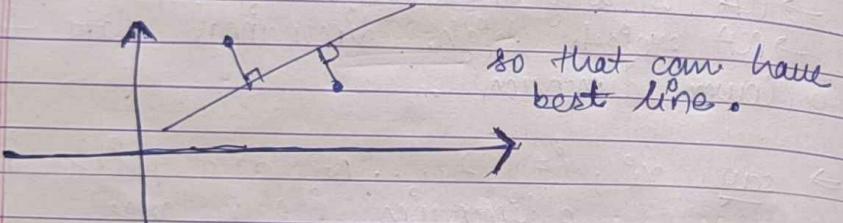
WJMK
= 30mm

Thus, training points projected on straight line \rightarrow
 $x_i^o = m + a_i e, i=1, 2, \dots, N$

In linear regression:- parallel to dependent variable $i.e.$



but in PCA \rightarrow perpendicular fitting!



\rightarrow points have normal distribution

↳ independent

↳ and sum them, total variance = $n o^2$

Standard deviation = $(n)^{1/2} o$

so any like $\rightarrow o \propto k \sqrt{n}$. (???)
(factor)

$$\text{Now } \rightarrow J(a_1, \dots, a_N, e) = \sum_{i=1}^N \|m + a_i e - x_i^o\|^2$$

so parabolic (convex)

WJMK
= 30mm

unit vector
thus = 1

$$J(a_1, \dots, a_N, e) = \sum_{i=1}^N a_i^2 \|e\|^2 - 2 \sum_{i=1}^N a_i e^T (x_i^o - m) + \sum_{i=1}^N (x_i^o - m)^2$$

* differentiate w.r.t a_i and equating to zero

$$a_i = e^T (x_i^o - m)$$

thus, plugging these values $\rightarrow N$

$$J(e) = -e^T Se + \sum_{i=1}^N (x_i^o - m)^2$$

where,

$$S = \sum_{i=1}^N (x_i^o - m)(x_i^o - m)^T$$

(scatter matrix)

\rightarrow minimizing $J_e \equiv$ maximize $e^T Se$.

constraint $\rightarrow e^T e = 1$ (as e is unit vector)

By Lagrange multiplier method \Rightarrow

$$(say) J_p = e^T Se + \lambda (e^T e - 1)$$

\uparrow objective fn

\uparrow differentiate w.r.t (e) and equate to 0

$$2Se + 2\lambda e = 0$$

$$Se = \lambda e$$

so, e is eigenvector of S corresponding to largest eigenvalue.

so reduces to eigen value problem

WJMK
Date _____
Page _____

* Eigen value problem :-

U have vector \vec{x} , multiply by a, s.t. it got either enhanced or diminished w/o changing direction.

so,

Dimensionality reduction problem reduces to Eigen value problem.

→ 1st eigen vector is 1st principal component.
→ 2nd

WJMK
Date _____
Page _____
[lec.] -

* Dimensionality Reduction

↳ curse of dimensionality
↳ more data (features) do not translate into better understanding

↳ Linear Dimensionality Reduction (PCA).

→ features not independent.

smaller no. of independent features
(dim - k)

→ linear features in d-dimensions gives independent features in m-dimensions ($m < d$).

PCA → going to lower dimensions

↳ Eigen value Problem

WJMK
Date _____
Page _____

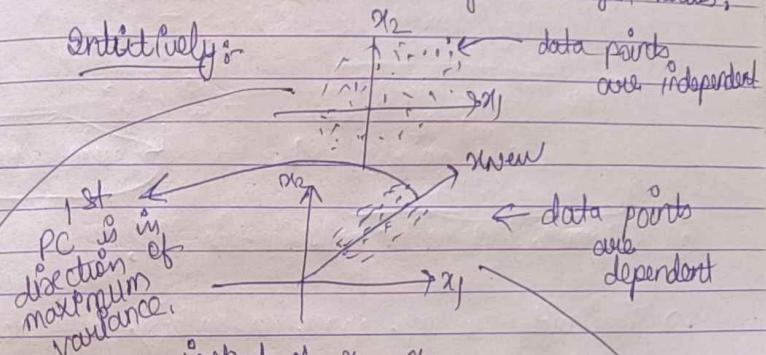
covariance matrix

XX^T

↳ Centred → Eigen value dec (SVD).
→ sort eigen values in descending order.

↳ Take largest ('k') eigen values.

Intuitively :-



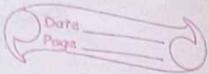
Instead of x_1, x_2 ,
take x_2 & x_{new} ,
as then can
work in 1 dimension only !!
PC → principal component .

$$\text{covariance matrix} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

$$\text{covariance matrix} = \begin{bmatrix} \sigma_1^2 & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \sigma_2^2 \end{bmatrix}$$

$$\text{and } \text{cov}(x_1, x_2) = \text{cov}(x_2, x_1)$$

WJMK
= MM.



→ says matrix given telling similarity b/w different weights.

Say → position : x, y

distance : $d(x, y) \leftarrow L_p$ Norm

similarity (kernel) $K = \langle x, y \rangle = x^T y$

Here (Linear kernel),
but it
can be any
kernel

$$= \sum x_i y_i$$

↳ cosine similarity

① difference b/w vector and function ????

→ Similarity in any space, is nothing but (dot product).

Pythagoras Theorem :

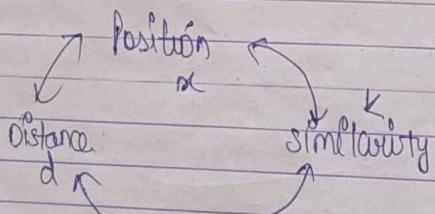
↳ Iska Reverse

↳ given distance, how to find best points $\in K$ in k -dimension?

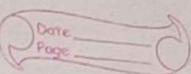
2, 3, ...
etc.

so can be used for dimension reduction.

PCA
 $x \rightarrow K'$
main
concept
kernels
val



WJMK
= MM.



XX^T

X

Position - Similarity :-

$$K = XX^T$$

→ what is Distance Matrix in Neural Networks?

Similarity-Distance :-

$$d(x_i^o, x_j^o)^2 = (x_i^o - x_j^o)^2$$

$$= \langle x_i^o, x_i^o \rangle + \langle x_j^o, x_j^o \rangle - 2 \langle x_i^o, x_j^o \rangle$$

$$= K_{ii}^o + K_{jj}^o - 2 K_{ij}^o$$

$$= D_{S, i, j}$$

we can determine D_S from K .

(Kernel) SVM mein pdha tha,

why so special?

↓ basis of (dot product). as all, decisions are made on

$w \cdot x < 0 \rightarrow$ then line w neeche, o/p if > 0 ,
↑ toh upar, if $=$ then on line.
(Inclination)

→ thus classification done on basis of dot product.
with kernel trick \rightarrow w/o going to high dimension, bmean \rightarrow millions of

$\langle w, x \rangle$ dot product

works on 3 assumption:-

i) all decision require $\langle w, x \rangle$

ii) $w \leftarrow f^n$ of x 's.

iii) we can calculate $\langle w, x \rangle$ in high-d,
w/o going to high-d.

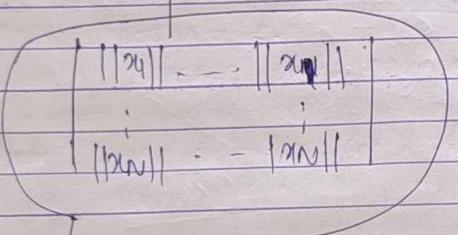
that's
gives us
(kernel)
version

Distance to position = MDS (multi dimensional scaling) :-

$$m = (\text{distance}(x_i^o, x_j^o))^2$$

$$= \|x_i^o - x_j^o\|^2$$

$$= \|x_i^o\|^2 + \|x_j^o\|^2 - 2 \langle x_i^o, x_j^o \rangle$$



want to get rid of this.

$$J_{MM} = -2X^T X = -2B$$

↳ Double-centring

The manifold hypothesis :-

Assumption:-

→ Natural data in high dimensional spaces concentrates close to lower dimensional manifolds.

→ probability density decreases very rapidly when moving away from the supporting manifold.

→ probability distribution is nothing but a histogram.

Stand with prob = 1/2, what will he do?

↳ Take a coin, if head → stand (say), if tail of no sit.

WJMK
Date _____
Page _____

WJMK
Date _____
Page _____

→ agar 100 baar bola,
tak 50 race set, and other
50 times stand.
→ manifold assumption holds for smaller distances.
↳ thus flat)
euclidean
↳ like earth flat for
smaller distances.

4/11
Monday

[GMM]

[E-M]

→ expectation - Maximization algorithm,

→ neither of model we learned till now like SUM, KNN etc. is a generative model, as we are not learning distributions in them, from pdf (any)

↳ (so) data can be generated!
Thus Hallucination is one of major challenge till now in gen AI.

product rule :-

$$P(x, y) = P(x) \cdot P(y|x)$$

$$= P(y) \cdot P(x|y)$$

which to choose?

depends on availability

of data set.

→ Thus direction matters:

W.M.K
Date _____
Page _____

M.W [Naive Bayes Model] → extension of Bayes Thm)

If directed, Bayesian belief Net.
↳ dependency graph.

A depends on B or B depends on A ... etc.

K-Means

- 1. fold
- 2. compute mean
- 3. assign nodes to cluster

↳ unsupervised (labels not available)
no. of clusters = K.

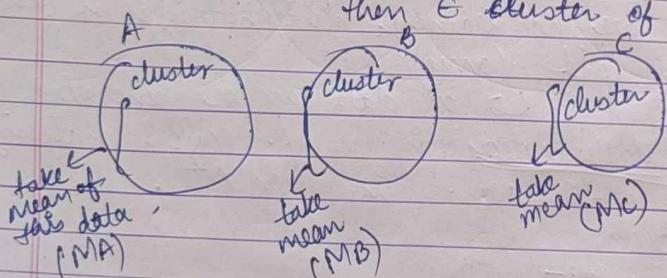
If $K=3 \rightarrow$ choose any random 3 pts. from data set (A, B, C) .

then har pt. ka sun teeno se distance lo, unka minimum lo, teeno ka.

↳ uss cluster ko that

pt. will e.

say if pt. D closest to A,
then e cluster of A.



Now take each pt. → take distance from MA,
MB, MC, jis minimum, uss cluster mein daalte jaao... fise mean lo ...

Jab tak mean change hota jaaye!!

* Drawback of K-mean! -

① Complexity of this algo?????

W.M.K
Date _____
Page _____

Hidden Markov model

N → data

K → clusters

M → Iterations

$TC = O(N \times K \times M)$ → very high!

② using only euclidean distance, what will happen? → will create circle (shape of cluster).

↓
↳ so behaves like circular or spherical cluster about mean.

③ outliers will misguide

→ If know hidden variables, can predict output in more better way. (like coin biased or not)

* Can u fit "gaussian curve" to get cluster?

→ Data set hogya

↳ Uska μ and σ^2
vekaab and

gaussian curve bharlo simple! (ii)

→ Say want to fit 2 gaussian! how?

↳ ditne dia khe can see from bar chart (har bar ~~bar~~ chart) bhaba not easy).

→ How to mix 2 gaussian?

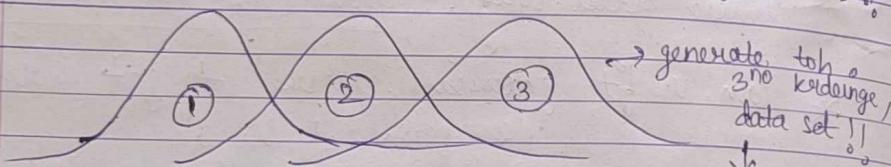
6/11/2024
Wednesday

WDMK
SMM

(lec 1)

GMM (Gaussian Mixture Model) :-

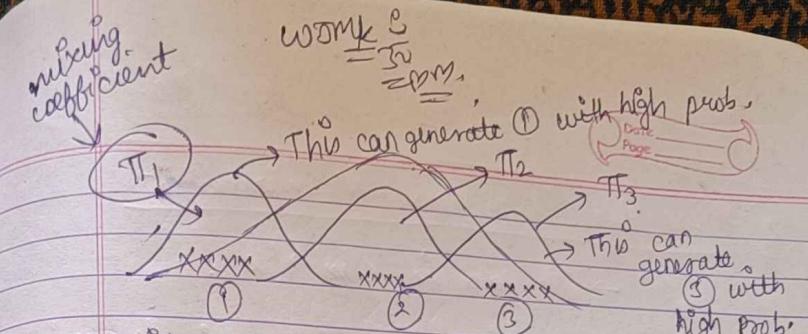
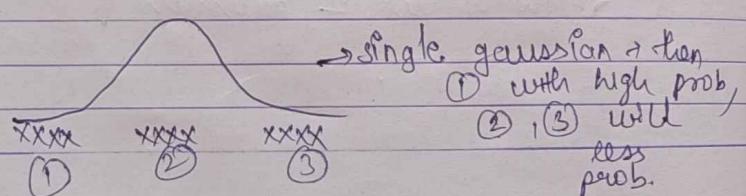
- as gaussian curve is from $-\infty$ to ∞ , thus can generate same dataset from ∞ distribution.
- any distribution (normal) can generate data set (any).
- if can generate a particular data set from any normal distribution (any μ, σ), then why we want to learn particular normal distribution (just take random μ, σ), then where is difference ???



but only one will give high probability (more likely) to generate that data set.

- If have 100 data, learn μ, σ , then can generate ' ∞ ' data, not only 100. (Yehi hota hai in Gen AI)
- But kya several data valid hai? need validation!
→ some may be fake, thus not acceptable.

e.g.:-



- It does not mean we can not generate ② or ① by this.
That's why we have to use multiple gaussian curves sometimes.

$$\sum_{k=1}^K \pi_k T_k \rightarrow \text{weighted sum / convex combination}$$

$$\pi_k > 0 \quad \forall k$$

Σ_k covariance (thus multidimensional data).

$$p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

$$p(x) = N(\mu, \Sigma)$$

$\hookrightarrow x \rightarrow$ iid type data.

$$p(x_1, x_2, \dots, x_n) = \prod N(x_i | \mu_i, \Sigma_i)$$

$$\Rightarrow \log(p(x_1, x_2, \dots, x_n)) = \sum \log(N(x_i | \mu_i, \Sigma_i))$$

\hookrightarrow then differentiate w.r.t Σ .
(???)

$$\Rightarrow p(x_1, x_2, \dots, x_n) = \prod \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k)$$

$$\log(p(x_1, x_2, \dots, x_n)) = \sum \log \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k)$$

complex, no close form exist to solve.

WJMK
2 Step algo.

H/No. 5



E-M algo (Expectation - Maximization) :-

$r_k^{(i)}$ \rightarrow i^{th} data in class \rightarrow 'k'.

* firstly learn (z) \rightarrow means kaun kis class mein
hai usko udhar bhejdo, then har class ka
 μ and σ nikalo,

\hookrightarrow 2-step process.

\Rightarrow can't learn z, μ, σ simultaneously,
almost impossible.

* using ANN, can approximate any f^n ,
 \hookrightarrow so powerful.