

SEM-2

Machine Learning for Big Data

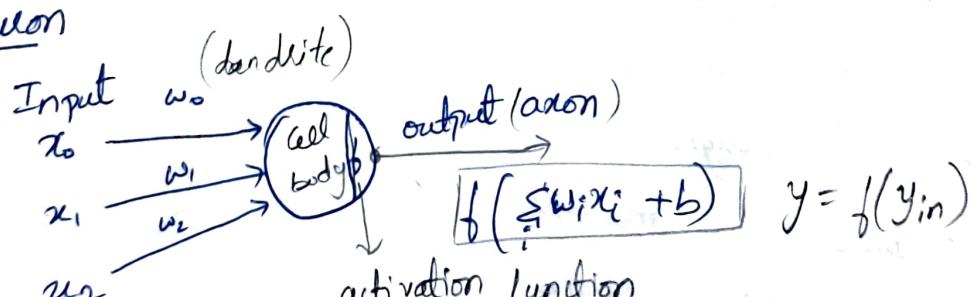
Topics

1. Artificial Neural Network
2. Clustering
3. Support vector Machines & kernel methods
4. Deep Learning \rightarrow ANN with extra hidden layer
5. Reinforcement Learning

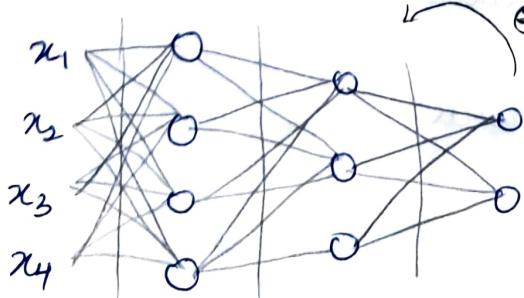
Linear \rightarrow can be classified with single classifier

Non-linear \rightarrow take the data to higher dimension and then classify

ANN \rightarrow Neuron



$$y_{in} = x_0w_0 + x_1w_1 + x_2w_2 + \dots + x_nw_n + b \quad (\text{Net input})$$



② Recurrent Neural Network

\rightarrow Feed forward Neural network

(Backpropagating) Neural Network is Feed Forward Neural Network
Back propagates the error.

What are all the different types of errors in ML algorithms?
How gradient descent approaches helps to reduce errors?
gradient descent is 'log-loss function'

Neural net



Multiple classifiers for non linear data -
Neural Network.

ANN applications - (Previously matlab)

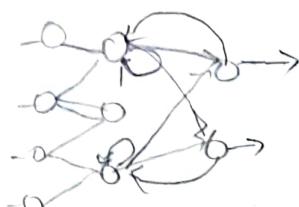
- pattern recognition - image analysis - shifted to neural network
- prediction
- optimization
- associative memory

ANN in medicine -

- Image analysis
- Drug design
- Biochemical analysis
- Diagnostic systems

Network Architecture - 2 types

- ① Feed Forward Neural Network - Deep learning, Backpropagating Neural network
- ② Recurrent neural network



Weights - settings the values for weights - enable learning/training
weight is information to feed to neural network to train the model

Activation function - linear $f(x) = x$

non-linear (used in multilayer)

$$f(x) = \begin{cases} 1 & \text{if } f(x) > 0 \\ 0 & \text{if } f(x) \leq 0 \end{cases}$$

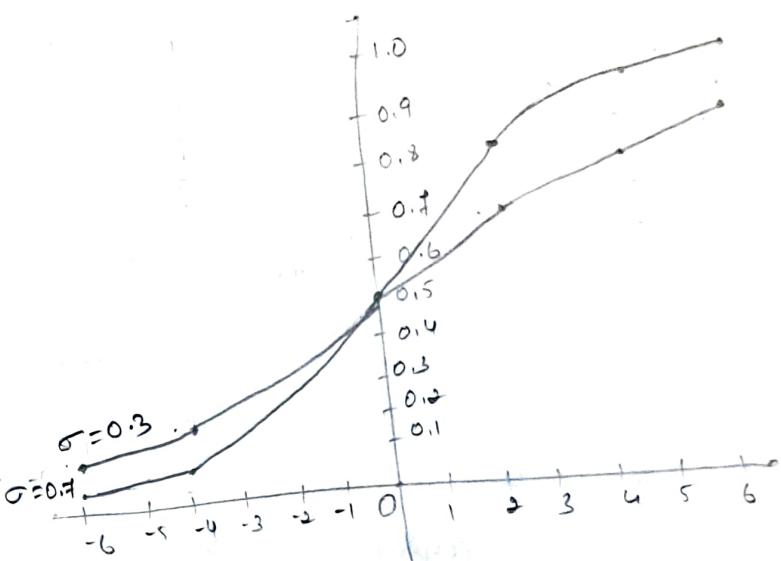


-2	1.0
-1	0.5
0	0.5
1	0.5
2	1.0

* Sigmoid function

$$f(x) = \frac{1}{1 + e^{-\sigma x}}$$

x	$f(x) = \frac{1}{1 + e^{-\sigma x}}$ $\sigma=0.3$	$f(x) = \frac{1}{1 + e^{-\sigma x}}$ $\sigma=0.7$
-6	0.14	0.014
-4	0.23	0.05
0	0.5	0.5
2	0.64	0.80
4	0.76	0.94
6	0.85	0.98



Bipolar sigmoid function

$$b(x) = 2 \left(\frac{1}{1 + e^{-\sigma x}} \right) - 1$$

$$b(x) = \frac{e^{\sigma x}}{1 + e^{\sigma x}} - 1$$

$$b(x) = \frac{2 - 1 - e^{-\sigma x}}{1 + e^{-\sigma x}}$$

$$b(x) = \frac{1 - e^{-\sigma x}}{1 + e^{-\sigma x}}$$

$$f(x) = \begin{cases} 1 & \text{if } x \geq 1 \\ 0 & \text{if } x < 1 \end{cases} \quad b(x) = 2f(x) - 1$$

vanishing gradient

Sigmoid activation function has a problem of vanishing gradient after large no of inputs are sent to hidden layers.

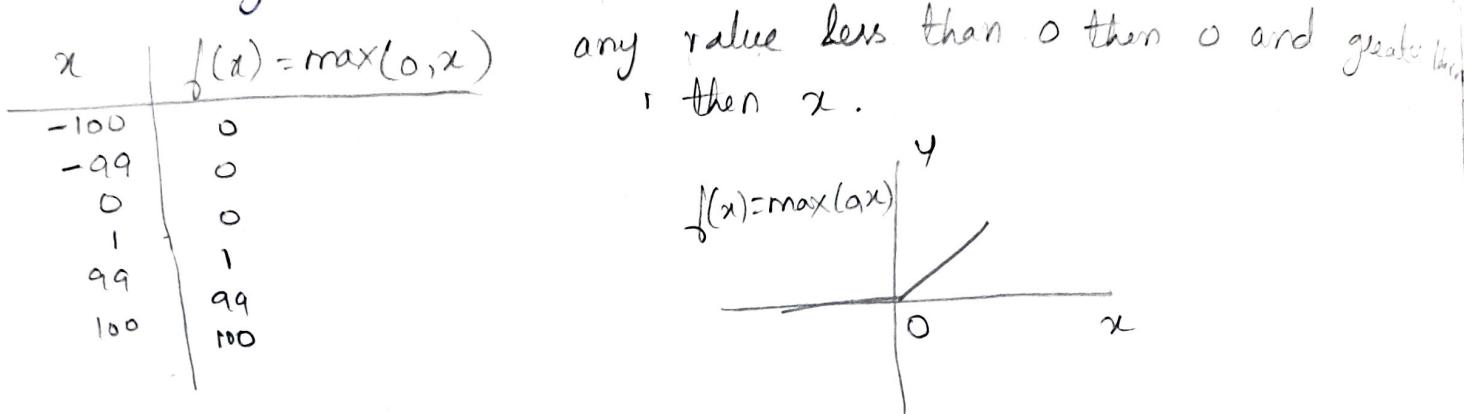
tanh → non-linear & differentiable function similar to sigmoid function, but output values ranges from -1 to +1.

in deep NN, prefer ReLU for hidden layers as activation function.
if multiple output then softmax

ReLU (Rectified linear unit)

commonly used activation function in CNN.

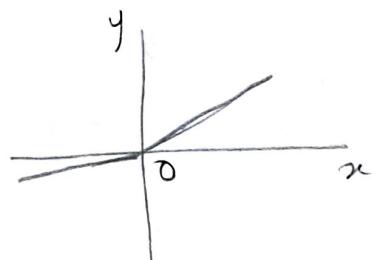
$$f(x) = \max(0, x)$$



Leaky ReLU

it is an improved version of ReLU function, instead of defining ReLU function as 0 for x less than 0, we define it as a small linear component of x .

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



Softmax function

large no of outputs at the last layer to normalize the output.

softmax → probabilistic form

multiclass classification, probabilities sum is 1.

range is between 0 & 1.

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum e^{x_j}} \text{ for all } j$$

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix} \rightarrow \begin{bmatrix} e^{x_i} \\ \vdots \\ e^{x_j} \end{bmatrix} \rightarrow \begin{bmatrix} 0.02 \\ 0.9 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$

* Difference between PyTorch & Tensorflow?

① McCulloch-Pitts Neuron Model

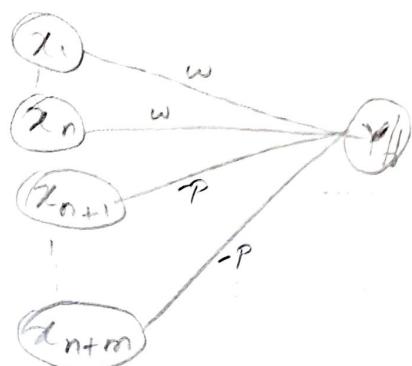
The most fundamental unit of deep learning neural network is a perceptron.

But the very first step towards the perceptron we use today was taken in 1943 by McCulloch and Pitts, by mimicking the functionality of a biological neuron.

find weights = net input
⇒ activation function
⇒ output

M-P performs two tasks

- ① Task-1 : performs an aggregation of all inputs
- ② Task-2 : based on the aggregated values the M-P net makes the decision.



$$y_{in} = x_1 w_1 + x_2 w_2 + \dots + x_n w_n \\ = \sum_{i=1}^n x_i w_i$$

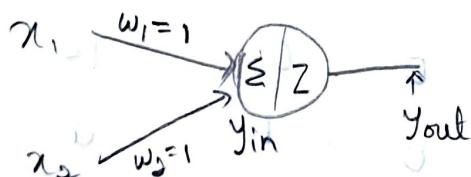
θ = threshold

$\theta = 1$

$$y_{out} = f(y_{in}) = \begin{cases} 1 & \text{if } y_{in} \geq 0 \\ 0 & \text{if } y_{in} < 0 \end{cases}$$

- ① Generate the output of logic AND function by McCulloch-Pitts neuron model.

AND - 2 inputs



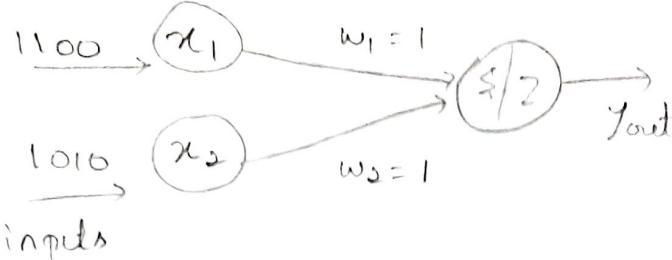
$$\text{net input } y_{in} = x_1 w_1 + x_2 w_2 = x_1 + x_2$$

Activation function $y_{out} = f(y_{in})$

$$\text{activation } f(y_{in}) = \begin{cases} 1 & \text{if } y_{in} \geq 2 \\ 0 & \text{if } y_{in} < 2 \end{cases}$$

x_1	x_2	$x_1 \text{ AND } x_2$	$y_{in} = x_1 + x_2$
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	2

OR function



x_1	x_2	$x_1 \text{ or } x_2$	$y_{in} = x_1 + x_2$
0	0	0	0
0	1	1	1
1	0	1	1
1	1	1	2

$$y_{in} = x_1 w_1 + x_2 w_2 = x_i w_i \\ = x_1 + x_2$$

$$\text{Activation function} = Y_{out} = f(y_{in}) = \begin{cases} 1 & \text{if } f(y_{in}) \leq 1 \\ 0 & \text{if } f(y_{in}) > 1 \end{cases}$$

Assignment

NOT function

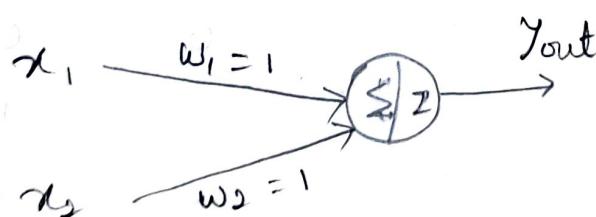


$$y_{in} = x_1 w_1 + b$$

x_1	$x_1 \text{ NOT}$	$y_{in} = xw_1 + b$
0	1	0
1	0	1

$$\text{Activation function} = Y_{out} = f(y_{in}) = \begin{cases} 0 & \text{if } y_{in} \geq 1 \\ 1 & \text{if } y_{in} < 1 \end{cases}$$

NAND function



$$y_{in} = w_1 x_1 + x_2 w_2 + b$$

$$w_1 = 1, w_2 = 1, b = 0, \theta = 1$$

x_1	x_2	$\text{NAND}(x_1, x_2)$	y_{in}
0	0	1	0
0	1	0	1
1	0	0	1
1	1	0	2

$$\text{Activation function} = Y_{out} = f(y_{in}) = \begin{cases} 0 & \text{if } y_{in} \geq 1 \\ 1 & \text{if } y_{in} < 1 \end{cases}$$

NOR Function

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$w_1 = 1 \quad w_2 = 1 \quad b = 0 \quad 0 = 1$$

Activation function

$$y_{out} = f(y_{in})$$

$$= \begin{cases} 0 & \text{if } y_{in} \geq 1 \\ 1 & \text{if } y_{in} < 1 \end{cases}$$

Limitations of M-P Neuron

- * Applies only to Boolean inputs
- * we always need to hard code the threshold
- * All inputs are holding equal priority, what if we want to assign more importance to single layer. & XOR doesn't work.

② Perception Networks

Frank Rosenblatt - 1962

Iterative weight adjustment (powerful technique) - learning rule

$\Delta t x_i \rightarrow$ input
learning rate \rightarrow target

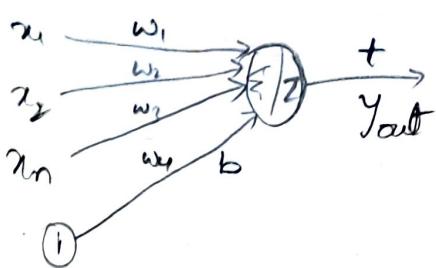
Purpose of training - to reduce error

Classification purpose

Types of perception Networks. - Single layer perceptron

Multilayer perceptron

(Back propagation neural network)



Bias is always to balance the network
- no error back propagating.

$$y_{in} = x_1 w_1 + x_2 w_2 + \dots + x_n w_n + b$$

steps

$$y_{in} = b + \sum_{i=1}^n x_i w_i$$

① Inputs

② initialize

w_1, w_2, x_1, b

activation $y_{out} = f(y_{in}) = \begin{cases} 1 & \text{if } y_{in} \geq 0 \\ -1 & \text{if } y_{in} < 0 \end{cases}$

$$y_{out} = t \rightarrow \text{stop}$$

update weights

$$w_i(\text{new}) = w_i(\text{old}) + \alpha t x_i$$

update bias

$$w_1(\text{new}) = w_1(\text{old}) + \alpha t x_1$$

$$b(\text{new}) = b(\text{old}) + \alpha t$$

$$w_2(\text{new}) = w_2(\text{old}) + \alpha t x_2$$

$$\leftarrow x_1 = 1$$

Example

- Develop a perception for AND function with bipolar inputs and targets. (initial weights are 0; learning rate = 1; and threshold = 1)

$$w_1 = w_2 = b = 0, \alpha = 1, \theta = 1$$

x_1	x_2	$t = x_1 \text{ AND } x_2$	y_{out}
1	1	1	1
1	-1	-1	-1
-1	1	-1	-1
-1	-1	-1	-1

Activation function

$$y_{out} = f(y_{in}) = \begin{cases} 1, & \text{if } y_{in} \geq 1 \\ -1, & \text{if } y_{in} < 1 \end{cases}$$

$$\textcircled{1} \quad x_1 = 1 \quad x_2 = 1 \quad t = 1$$

$$w_1 = 0 \quad w_2 = 0 \quad b = 0$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$y_{in} = 0 + 0 + 0 = 0$$

$$y_{out} = f(y_{in}) = f(0) = -1$$

$$\boxed{\text{target} = 1 \neq y_{out} = -1}$$

update the weights

y_{in} = net input

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$w_1(n) = w_{1(0)} + \alpha t x_1$$

$$w_2(n) = w_{2(0)} + \alpha t x_2$$

$$w_1(n) = 0 + 1 \times 1 \times 1 = 1$$

$$w_2(n) = 0 + 1 \times 1 \times 1 = 1$$

$$b(n) = b_{(0)} + \alpha t$$
$$= 0 + 1 \times 1 = 1$$

$$y_{in} = w_1 x_1 + w_2 x_2 + b$$

$$y_{in} = 1 \times 1 + 1 \times 1 + 1 = \underline{\underline{3}}$$

$$y_{out} = f(y_{in}) = f(3) = 1$$

(target = +1) == (y_{out} = 1)

stop
the learning

③ Next pattern

$$x_1 = 1 \quad x_2 = 1 \quad t = -1$$

$$w_1 = 1 \quad w_2 = 1 \quad b = 1$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$y_{in} = 1 - 1 + 1 = \underline{\underline{1}}$$

$$y_{out} = f(y_{in}) = f(1) = \underline{\underline{1}}$$

target = -1 ≠ y_{out} = 1

update weights

$$w_1(n) = w_{1(0)} + \alpha t x_1$$

$$w_1(n) = 1 + 1 \times -1 \times 1 = 0$$

$$w_2(n) = w_{2(0)} + \alpha t x_2$$

$$w_2(n) = 1 + 1 \times -1 \times -1 = 2$$

$$b(n) = b_{(0)} + \alpha t$$

$$b(n) = 1 + 1 \times -1 = 0$$

~~target = -1~~
 $y_{in} = x_1 w_1 + x_2 w_2 + b$
 $y_{in} = 0 + 2 + 0 = -2$
 $y_{out} = f(y_{in}) = -1$

④ Next pattern

$$x_1 = -1 \quad x_2 = 1 \quad t = -1$$

$$w_1 = 0 \quad w_2 = 2 \quad b = 0$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$y_{in} = 0 + 2 + 0 = 2$$

$$y_{out} = f(y_{in}) = 1$$

target = -1 ≠ y_{out} = 1

update weights

$$w_1(n) = w_{1(0)} + \alpha t x_1$$

$$= 0 + 1 \times -1 \times -1 = 1$$

$$w_2(n) = w_{2(0)} + \alpha t x_2$$

$$= 2 + 1 \times -1 \times 1 = 1$$

$$b(n) = b_{(0)} + \alpha t$$

$$= 0 + 1 \times -1 = -1$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$y_{in} = -1 + 1 - 1 = -1$$

check the activation $y_{out} = f(y_{in}) = -1$

$$\boxed{(target = -1) == (y_{out} = -1)}$$

— Stop

④ next pattern

$$x_1 = -1 \quad x_2 = -1 \quad t = -1 \\ w_1 = 1 \quad w_2 = 1 \quad b = -1$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$y_{in} = -1 + 1 - 1 = -3$$

$$y_{out} = f(y_{in}) = -1$$

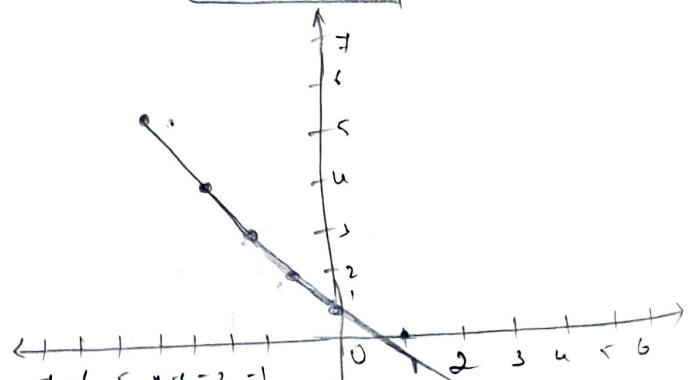
$$\boxed{(target = -1) == (y_{out} = -1)}$$

For the unseen class : $w_1 x_1 + w_2 x_2 + b = 0$

$$w_1 = 1 \quad w_2 = 1 \quad b = -1 \quad \text{find weight}$$

$$x_1 + x_2 - 1 = 0$$

$$\boxed{x_2 = 1 - x_1}$$



perception
classified

x_1	$x_2 = 1 - x_1$
-4	$1 - (-4) = 5$
-3	4
-2	3
-1	2
0	1
1	0
2	-1
3	-2
4	-3

• Develop a perceptron for OR function with bipolar inputs and targets. (initial weights are 0; learning rate = 1; and threshold = 1)

$$w_1 = w_2 = b = 0, \alpha = 1, \theta = 1$$

x_1	x_2	$t = x_1 \text{ or } x_2$	y_{out}
1	1	1	1
1	-1	1	1
-1	1	1	1
-1	-1	-1	-1

Activation function

$$y_{out} = f(y_{in}) = \begin{cases} 1, & \text{if } y_{in} \geq 1 \\ -1, & \text{if } y_{in} < 1 \end{cases}$$

② Next pattern

$$x_1 = 1 \quad x_2 = -1 \quad t = 1$$

$$w_1 = 1 \quad w_2 = 1 \quad b = 1$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$y_{in} = 1 + -1 \times 1 + 1 = 1$$

$$y_{out} = f(y_{in}) = f(1) = 1$$

$$\boxed{\text{target} = 1 == y_{out} = 1}$$

① First pattern

$$x_1 = 1 \quad x_2 = 1 \quad t = 1$$

$$w_1 = 0 \quad w_2 = 0 \quad b = 0$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$y_{in} = 0 + 0 + 0 = 0$$

$$y_{out} = f(y_{in}) = f(0) = -1$$

$$\boxed{\text{target} = 1 \neq y_{out} = -1}$$

update the weights

$$w_1(n) = w_1(0) + \alpha t x_1$$

$$w_1(n) = 0 + 1 \times 1 \times 1 = 1$$

$$w_2(n) = w_2(0) + \alpha t x_2$$

$$w_2(n) = 0 + 1 \times 1 \times 1 = 1$$

$$y_{in}^{(n)} = 0 + 1 \times 1 = 1$$

$$y_{in} = 1 \times 1 + 1 \times 1 + 0 = 2$$

$$y_{out} = f(y_{in}) = 1$$

$$\boxed{(\text{target} = 1) == (y_{out} = 1)}$$

stop

③ next pattern

$$x_1 = -1 \quad x_2 = 1 \quad t = 1$$

$$w_1 = 1 \quad w_2 = 1 \quad b = 1$$

$$y_{in} = x_1 w_1 + x_2 w_2 + b$$

$$= -1 + 1 + 1$$

$$= 1$$

$$\checkmark \boxed{(\text{target} = 1) == (y_{out} = 1)}$$

$$④ x_1 = -1 \quad x_2 = -1 \quad t = -1$$

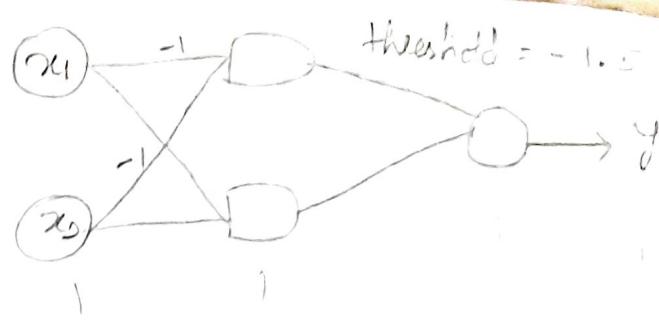
$$w_1 = 1 \quad w_2 = 1 \quad b = 1$$

$$y_{in} = -1 - 1 + 1 = -1$$

$$f(-1) = -1$$

$$\boxed{(\text{target} = -1) == (y_{out} = -1)}$$

XOR Problem



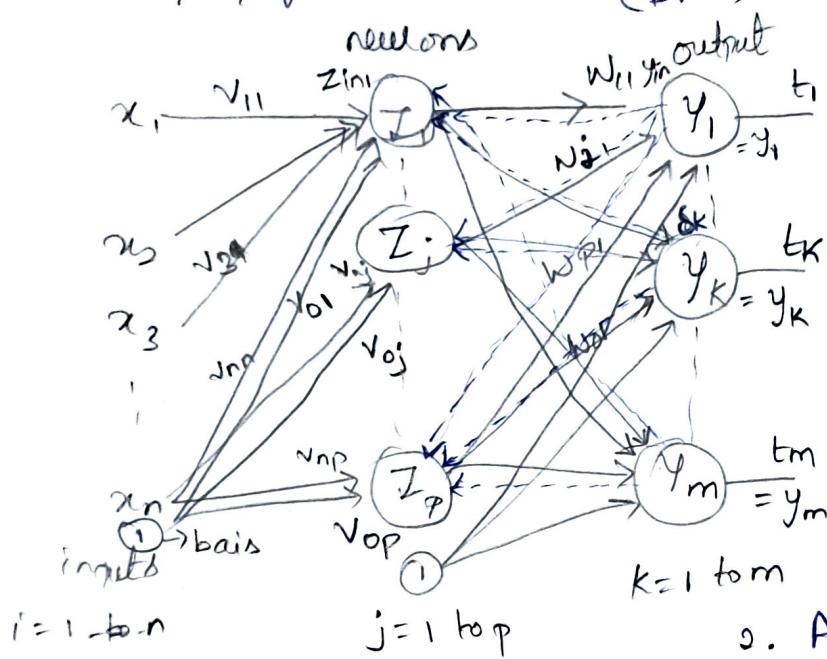
Multilayer Networks

① Feed forward networks

no feedbacks, output is calculated for every input & hidden.
eg - Backpropagation nn

Extension of XOR, first is input - last is output, middle are hidden layer - deprn

Back propagation Network (BPNN)



feed forward calculations is done
by keras. sequential

Input is net neuron

input = x_1, x_2, x_3

Input layer & hidden layer = z_1, z_2, z_m

Output layer = y_1, y_2, \dots, y_m

① Feed forward Step

1. net input at Z_i

$$Z_{inj} = v_{0j} + v_{1j}x_1 + v_{2j}x_2 + \dots + v_{nj}x_n + v_{nj}$$

$$Z_{inj} = v_{0j} + \sum_{i=1}^n x_i v_{ij}$$

2. Activation function at hidden layer

$$z_j = f(Z_{inj}) = \frac{1}{1 + e^{-Z_{inj}}}$$

if sigmoidal

3. Feed forward

$$y_{inj} = Z_i w_{i1} + Z_j w_{j1} + \dots + Z_m w_{m1} + b_0$$

$$y_{ink} = w_{0k} + \sum_{j=1}^m Z_j w_{jk}$$

4. Activation

$$y_{ik} = f(y_{ink}) = \frac{1}{1 + e^{-y_{ink}}}$$

if sigmoidal

② Error back propagation step

optimizing
adom

$$\delta_1 = (t_1 - y_1) \cdot f'(y_{in1})$$

Actual predicted

$$f(y_{in1}) = f(y_{in1})(1-f(y_{in1}))$$

$$\delta_K = (t_K - y_K) \cdot f'(y_{inK})$$

$$f(y_{inK}) (1-f(y_{inK}))$$

2. BP error to hidden layers I

$$\delta_{inj} = \sum_{k=1}^m \delta_k \cdot w_{jk}$$

calculate error and
send back to last layer

3. Error at hidden layer I

$$\delta_j = \delta_{inj} \cdot f(I_{inj})$$

③ Update the weights learning rate 0.3 to 0.5

$$v_{01}(n) = v_{01}(0) + \alpha \delta_j$$

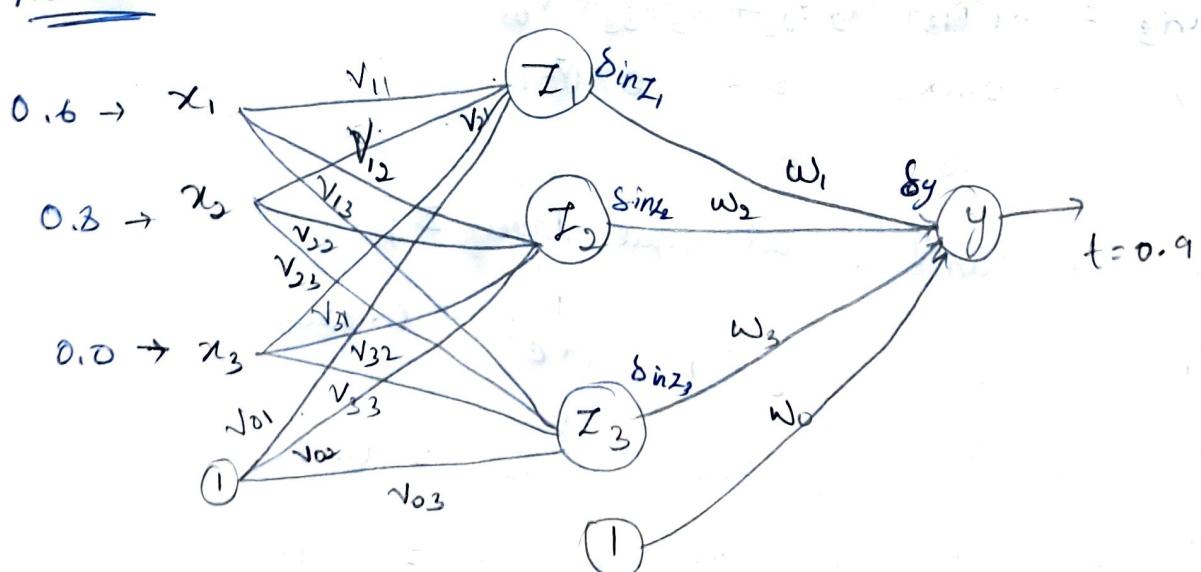
$$v_{0i}(n) = v_{0i}(0) + \alpha \delta_j \cdot x_i^{bias}$$

$$\text{first layer } v_{ij}(n) = v_{ij}(0) + \alpha \delta_j x_i$$

$$\text{next layer } v_{0k}(n) = w_{0k}(0) + \alpha \delta_k$$

$$v_{jk}(n) = w_{jk}(0) + \alpha \delta_k I_j$$

Problem



$$[\omega_1 \ \omega_2 \ \omega_3] = [-1 \ 1 \ 2]$$

$$\omega_0 = -1 \quad \alpha = 0.3$$

Activation = sigmoid

$$\begin{bmatrix} v_{01} & v_{03} & v_{03} \\ v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \end{bmatrix} = \begin{bmatrix} 0 & 0 & -1 \\ 2 & 1 & 0 \\ 1 & 2 & 2 \\ 0 & 3 & 1 \end{bmatrix}$$

① feed forward network \rightarrow net inputs / Activation - output

② Error - Back propagation

③ update weights

④ if N *② net inputs at first layer / hidden layer / input layer
net input at $Z_i \Rightarrow Z_{in_i}$

$$Z_{in_1} = x_1 v_{11} + x_2 v_{21} + x_3 v_{31} + v_{01}$$

$$Z_{in_1} = (0.6 \times 2) + (0.3 \times 1) + (0 \times 0) + 0$$

$$Z_{in_1} = 2$$

$$Z_{in_2} = x_1 v_{12} + x_2 v_{22} + x_3 v_{32} + v_{02}$$

$$Z_{in_2} = (0.6 \times 1) + (0.3 \times 2) + (0 \times 0) + 0$$

$$Z_{in_2} = 2.2$$

$$Z_{in_3} = x_1 v_{13} + x_2 v_{23} + x_3 v_{33} + v_{03}$$

$$Z_{in_3} = (0.6 \times 0) + (0.3 \times 2) + (0 \times 1) + 0$$

$$Z_{in_3} = 0.6$$

activation function at input layer Z ;

$$Z_1 = f(Z_{in_1}) = \frac{1}{1 + e^{-Z_{in_1}}} = \frac{1}{1 + e^{-2}} = 0.8807$$

$$Z_2 = f(Z_{in_2}) = \frac{1}{1 + e^{-Z_{in_2}}} = \frac{1}{1 + e^{-2.2}} = 0.9002$$

$$Z_3 = f(Z_{in_3}) = \frac{1}{1 + e^{-Z_{in_3}}} = \frac{1}{1 + e^{-0.6}} = 0.6456$$

net input at $y =$

$$y_{in} = z_1 w_1 + z_2 w_2 + z_3 w_3 + w_0$$

$$y_{in} = 0.8808 \times -1 + 0.9002 \times 1 + 0.6556 \times 2 - 1$$

$$y_{in} = 0.3106$$

$$y_{out} = f(y_{in}) = \frac{1}{1+e^{-y_{in}}} = \frac{1}{1+e^{-0.3106}} = 0.577$$

$t = 0.9 \neq y_{out} = 0.577$

② Error at output layer y

$$\delta_y = (t - y_{out}) \cdot f'(y_{in})$$

$$\delta_y = (t - y_{out}) f'(y_{in}) (1 - f(y_{in}))$$

$$\delta_y = (0.9 - 0.577) (0.577) (1 - 0.577) = 0.078$$

Error inputs to input layer z (back propagating)

$$\delta_{in z_1} = \delta_y \times w_1 = 0.078 \times -1 = -0.078$$

$$\delta_{in z_2} = \delta_y \times w_2 = 0.078 \times 1 = 0.078$$

$$\delta_{in z_3} = \delta_y \times w_3 = 0.078 \times 2 = 0.156$$

Error at input layer after receiving error inputs

$$\delta_{z_1} = \delta_{in z_1} f'(z_{in_1})$$

$$= \delta_{in z_1} \times f(z_{in_1}) \times (1 - f(z_{in_1}))$$

$$\delta_{z_1} = -0.078 \times 0.8808 \times (1 - 0.8807)$$

$$\delta_{z_1} = -0.00819$$

$$\delta_{z_2} = \delta_{in z_2} \times f'(z_{in_2})$$

$$= 0.078 \times 0.9002 \times (1 - 0.9002) = 0.007$$

$$\delta_{z_3} = \delta_{in z_3} \times f'(z_{in_3}) = 0.0356$$

③ Update the weights \rightarrow from last layer to first

$$w_{0(1)} = w_{0(0)} + \alpha \delta_y = -1 + 0.3 \times 0.078 = -0.9766$$

$$w_1(n) = w_{1(0)} + \alpha \delta y z_1 = -1 + 0.3 \times 0.078 \times 0.8808 = -0.9$$

$$w_2(n) = w_{2(0)} + \alpha \delta y z_2 = 1 + 0.3 \times 0.078 \times 0.9002 = 1.02$$

$$w_3(n) = w_{3(0)} + \alpha \delta y z_3 = 2 + 0.3 \times 0.078 \times 0.6456 = 2.015$$

input weight updation

$$v_{01}(n) = v_{01(0)} + \alpha \delta z_1 = 0 + 0.3 \times (-0.00819) = 0.00245$$

$$v_{02}(n) = v_{02(0)} + \alpha \delta z_2 = 0 + 0.3 \times (0.007) = 0.0021$$

$$v_{03}(n) = v_{03(0)} + \alpha \delta z_3 = -1 + 0.3 \times (0.0356) = -0.9893$$

$$v_{11}(n) = v_{11(0)} + \alpha \delta z_1 x_1 = 2 + (0.3 \times (-0.00819) \times 0.6) = 1.998$$

$$v_{21}(n) = v_{21(0)} + \alpha \delta z_1 x_2 = 1 + (0.3 \times (0.007) \times 0.8) = 0.9980$$

$$v_{31}(n) = v_{31(0)} + \alpha \delta z_1 x_3 = 0 + (0.3 \times (0.0356) \times 0) = 0$$

$$v_{12}(n) = v_{12(0)} + \alpha \delta z_2 x_1 = 1 + (0.3 \times (0.007) \times 0.6) = 0.9980$$

$$v_{22}(n) = v_{22(0)} + \alpha \delta z_2 x_2 = 2 + (0.3 \times (0.007) \times 0.8) = 2.00169$$

$$v_{32}(n) = v_{32(0)} + \alpha \delta z_2 x_3 = 3 + (0.3 \times (0.007) \times 0) = 3$$

$$v_{13}(n) = v_{13(0)} + \alpha \delta z_3 x_1 = 0 + (0.3 \times (0.0356) \times 0.6) = 0.064$$

$$v_{23}(n) = v_{23(0)} + \alpha \delta z_3 x_2 = 2 + (0.3 \times (0.0356) \times 0.8) = 2.008$$

$$v_{33}(n) = v_{33(0)} + \alpha \delta z_3 x_3 = 1 + (0.3 \times (0.0356) \times 0.0) = 1$$

Clustering in machine learning

In real world, not every data we work upon has a target variable. Popular under unsupervised learning is clustering analysis. Goal is to group similar data points into a group.

clustering - organization of unlabeled data into similarity groups called clusters. Data items within a cluster are 'similar'.

↓ unsupervised learning Data items between clusters are 'dissimilar'

clustering algorithms

dimension reduction algorithm → columns deducing

generative algorithms → data points (rows)

Example - Email dataset, when there is no label for spam/ham then clustering.
applications - market segmentation, genetics, medical imaging, video recommendations

popular clustering algorithms

① K-means - fast (Euclidean distance approach)

② Hierarchical

③ Density based spatial clustering (DBSCAN)

④ Gaussian mixture models - soft clustering

Types of clustering - ① Hard clustering → each data point belongs to a cluster completely or not.

one data point = one cluster

eg - if 4 data points, must have 2 clusters

② Soft clustering → instead of putting each data point into separate cluster, probability of likelihood of data to be clustered is assigned to each data point.

Distance measures used in clustering → find dissimilarity between two data points.

1. Minkowski distance family

- Euclidean distance (clustering) = $\sqrt{\sum (q_i - p_i)^2}$

- Manhattan distance = $\sum_{i=1}^d |x_{i1} - x_{j1}|$

2. Cosine distance (text)

3. Mahalanobis distance

→ cosine of the angle between two vectors in multidimensional - high dimensional data

$$\cos \theta = \frac{\text{dot product}}{(|a| \times |b|)} = \frac{a \cdot b}{\|a\| \|b\|}$$

K-Means Clustering

MacQueen, 1967

partitioning clustering algorithm - dividing a dataset into distinct groups of clusters. Each cluster has centre called centroid.

partitioning-based
nonoverlapping groups

K-means algorithm

Step-1 : choose the number of clusters (K) using elbow method

Step-2 : select at random K points, the centroids (not necessarily from your dataset)

Step-3 : Assign each data point to the closest centroid → that forms K clusters

Step-4 : Compute and place the new centroid of each cluster.



Step-5 : Reassign each data points to the new closest centroid if any reassignment took place, go to step 4, otherwise stop.

Example

Determine which medicine belong to cluster 1 and which medicine to other cluster

Objects	Attribute-1 index weight	Attribute-2 PH
Medicine A	1	1
Medicine B	2	1
Medicine C	3	3
Medicine D	4	4

$$\textcircled{1} \quad K=2 \quad C_1(1,1) \quad C_2(2,1)$$

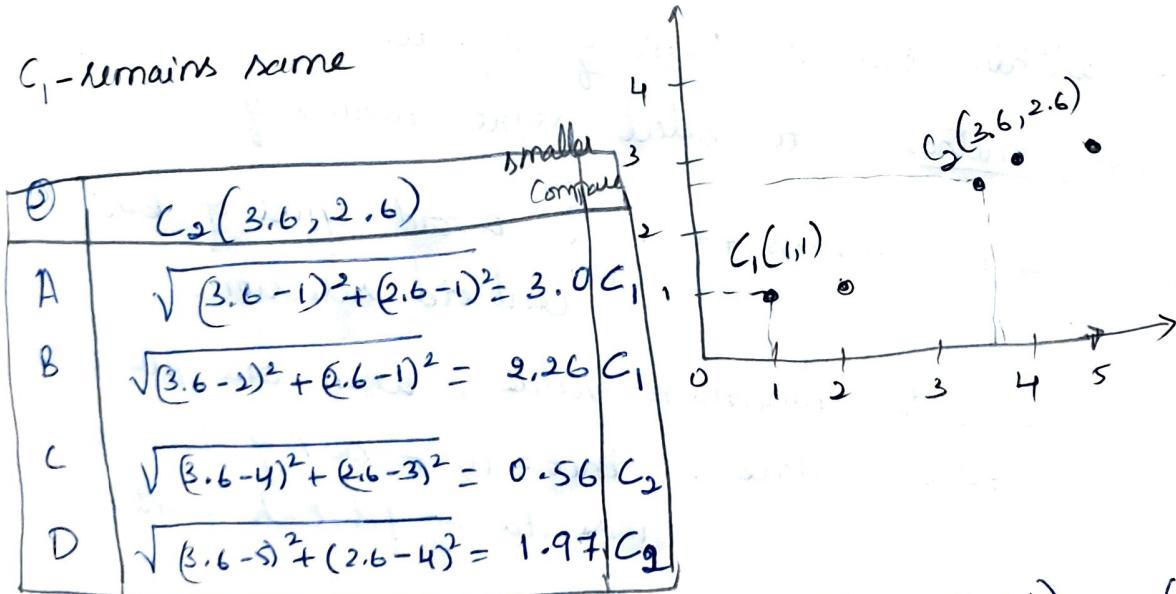
	$C_1(1,1)$	$C_2(2,1)$	smaller (compare C_1 & C_2 values)
A	$(1,1)$	$\sqrt{(1-1)^2 + (1-1)^2} = 0$	C_1
B	$(2,1)$	$\sqrt{(2-1)^2 + (1-1)^2} = 1$	C_2
C	$(4,3)$	$\sqrt{(4-1)^2 + (3-1)^2} = 3.6$	C_2
D	$(5,4)$	$\sqrt{(5-1)^2 + (4-1)^2} = 5$	C_2

$$C_1 = (1,1)$$

$$C_2 = (2,1), (4,3), (5,4)$$

$$C_2 = \left(\frac{2+4+5}{3}, \frac{1+3+4}{3} \right) = (3.6, 2.6)$$

C_1 - remains same



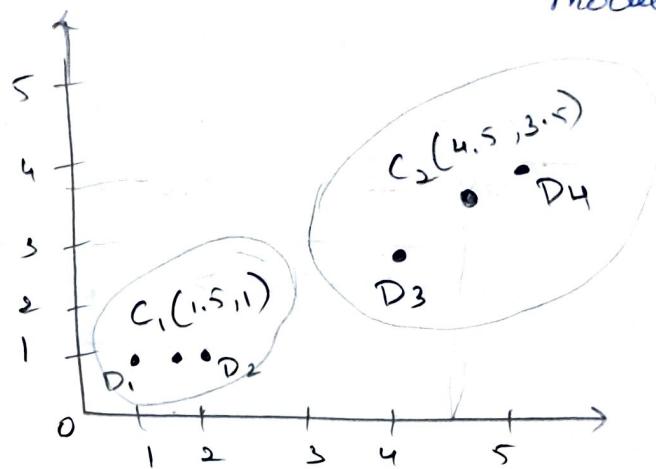
$$C_1 = (1,1) + (2,1) = \left(\frac{1+2}{2}, \frac{1+1}{2} \right) = (1.5, 1)$$

$$C_2 = (4,3) + (5,4) = \left(\frac{4+5}{2}, \frac{3+4}{2} \right) = (4.5, 3.5)$$

(3) (complete)

		$C_1(1.5, 1)$	$C_2(4.5, 3.5)$
A	(1, 1)	$\sqrt{(1.5-1)^2 + (1-1)^2} = 0.5$	$\sqrt{(4.5-1)^2 + (3.5-1)^2} = 4.3$
B	(2, 1)	$\sqrt{(2-1.5)^2 + (1-1)^2} = 0.5$	$\sqrt{(4.5-2)^2 + (3.5-1)^2} = 3.5$
C	(4, 3)	$\sqrt{(4-1.5)^2 + (3-1)^2} = 3.2$	$\sqrt{(4.5-4)^2 + (3.5-3)^2} = 0.7$
D	(5, 4)	$\sqrt{(5-1.5)^2 + (4-1)^2} = 4.6$	$\sqrt{(5-4.5)^2 + (4-3.5)^2} = 0.7$

clusters & centroid do not change so stopping &
model is ready.

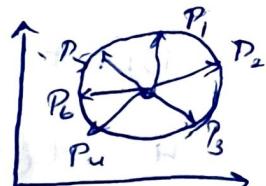


K-means clustering :-

WCSS - within cluster sum of square - Inertia
Using elbow method - to select right number of clusters
K-means \rightarrow K-means ++ . To avoid effect of bad
centroid selection.

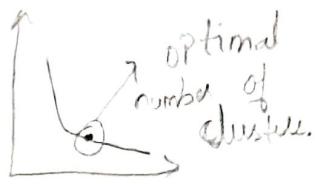
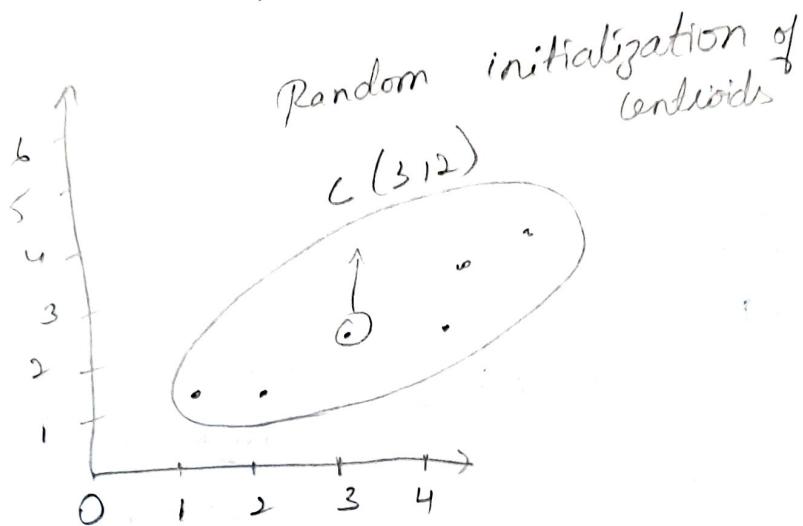
* WCSS \rightarrow one with minimum value, when one placed to
correct place - init = 10 - no. of locations
max-ita - for each point.

$$\sum_{i=1}^{n=s} (P_i, c)^2 = ? \text{ distance}$$



k-means = KMeans (n-clusters = 1, init =
'K-means++', max-iter = 300, n-init = 10, random = 0)

- Steps
- ① finding out right place/location of cluster
 - ② finding right number of clusters.



Data points	$c(3, 2)$	$c(1, 1)$	$c(5, 3)$
(1, 1)	$(3-1)^2 + (2-1)^2 +$	$(1-1)^2 + (1-1)^2 +$	$(5-1)^2 + (1-1)^2 +$
(2, 1)	$(3-2)^2 + (2-1)^2 +$	$(1-2)^2 + (1-1)^2 +$	$(5-2)^2 + (3-1)^2 +$
(4, 3)	$(3-4)^2 + (2-3)^2 +$	$(1-4)^2 + (1-3)^2 +$	$(5-4)^2 + (3-3)^2 +$
(5, 4)	$(5-3)^2 + (2-4)^2$ = 17	$(1-5)^2 + (1-4)^2$ = 39	$(5-5)^2 + (3-4)^2 =$ = 35

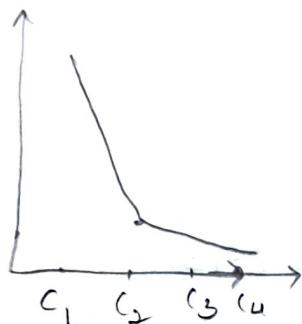
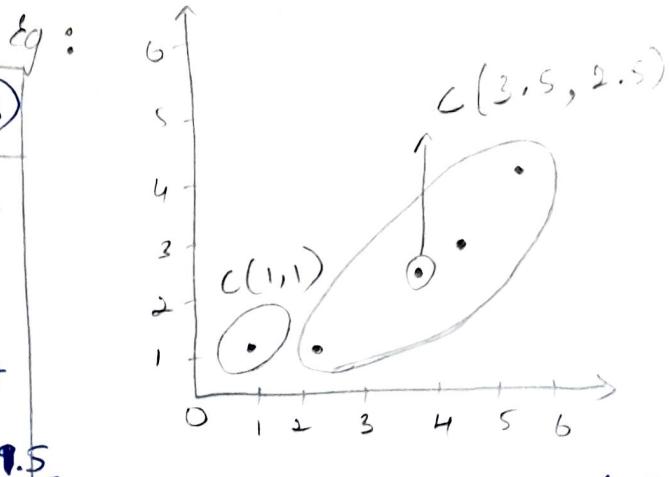
$$dc_1 = (3-1)^2 + (2-1)^2$$

$$dc_2 = 39$$

$$dc_3 = 35.$$

DPS	$c_1(1, 1)$	$c_2 = (3.5, 2.5)$
(1, 1)	$(1-1)^2 + (1-1)^2 = 0$	$WCSS_2 =$ $(3.5-2)^2 + (2.5-1)^2$
(2, 1)		$+ (3.5-4)^2 + (2.5-3)^2$
(4, 3)		$+ (3.5-5)^2 + (2.5-4)^2$
(5, 4)		$= 9.5$

DPS	$C(1.5, 1)$	$C(4.5, 3.5)$
(1, 1)		$(4.5 - 1)^2 +$
(2, 1)		$(3.5 - 1)^2 +$
(4, 3)		$(4.5 - 4)^2 +$
(5, 4)		$(3.5 - 3)^2 +$ $(4.5 - 5)^2 +$ $(3.5 - 4)^2 = 9.5$



$\rightarrow C_2$ has elbow

Strengths of k-means - easy to understand

Time complexity

$$: O(tKn)$$

n - no. of data points

K - no. of clusters

t - no. of iterations

k - means most popular clustering algorithm.

Since both k and t are small. k-means is considered a linear algorithm.

Problem faced

outliers - Handling with outliers

Remove some data points that are much farther away from

import matplotlib.pyplot as plt

from sklearn.cluster import KMeans

Create arrays that resemble 2 variables in a data set.

$$\begin{aligned} X &= \\ Y &= \end{aligned}$$

Turn the data into a set of points

data = list(zip(x, y))

Hierarchical Clustering

Find successive clusters using previously established clusters.

2 types → Agglomerative (bottom-up)
Divisive (top-down)

Agglomerative Algorithm - Begin with each data element as a separate cluster and merge them into successively larger clusters.

implent
Divisive Algorithm - Begin with the whole set and proceed to divide it into successive small clusters

applications same as k-means

Agglomerative algorithm

① Step 1 → make each point a single point cluster → that forms no clusters.

② Step 2 → take two closest data points

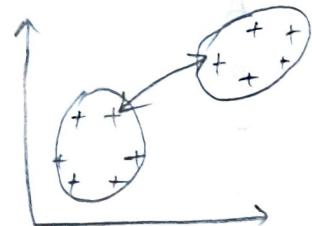
Euclidean or manhattan

③ Step 3 → take closest clusters and make them one cluster

④ Step 4 → Repeat step 3 until there is only one cluster

↓

Finish



Distance between clusters

* option 1 : closest points

Min (single) Linkage - WARD → within cluster variance

* option 2 : furthest points

Max (complete) Linkage

* option 3 : Average Distance
Linkage

* option 4 : Distance b/w centroids
centroïd of
WARD method

Dendograms - How dendograms works

(P₂-P₃)

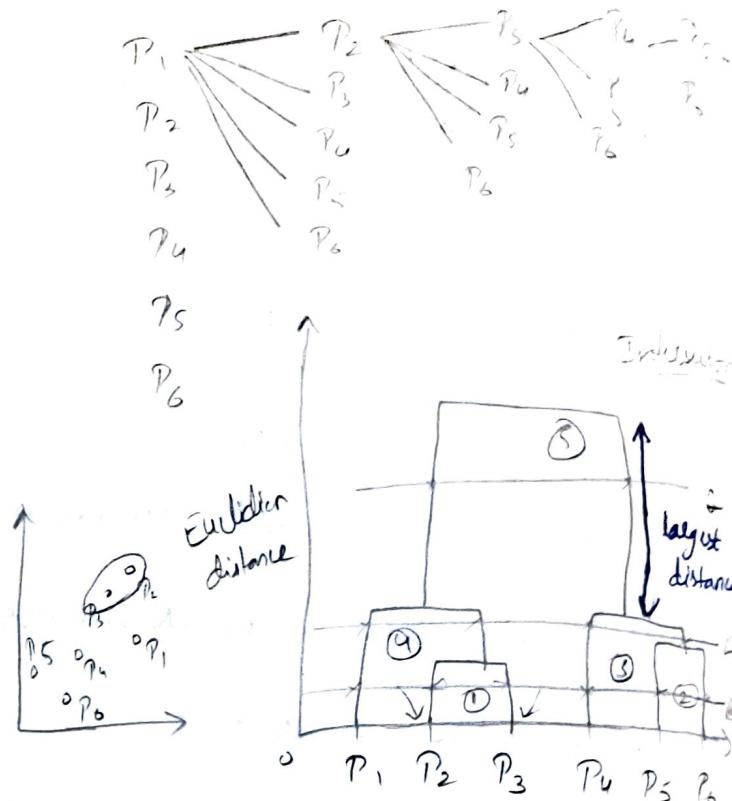
	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0					
P ₂		0				
P ₃	0.5		0			
P ₄				0		
P ₅					0	
P ₆					0.5	0

(P₂-P₃)
①

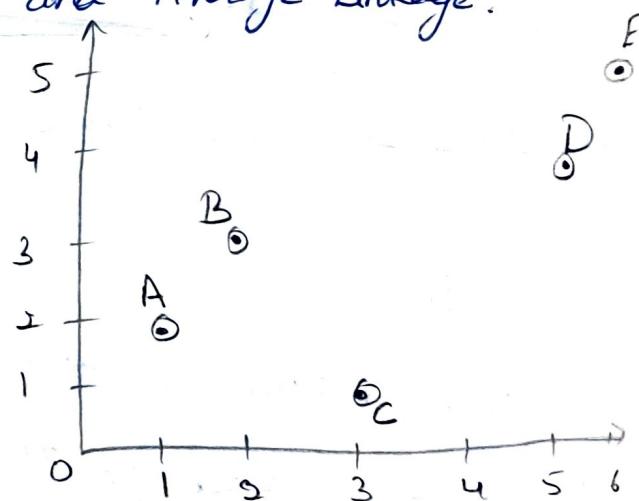
Example : Distance Between
matrices used in hierarchical clustering.

• Single Linkage, complete Linkage and Average Linkage.

Point	x	y
A	1	2
B	2	3
C	3	1
D	5	4
E	6	5



optimal no of clusters - 2
after intersection largest distance
vertical line is considered



① Single / Minimum Linkage

Euclidean distance

$$d(AB) = \sqrt{(2-1)^2 + (3-2)^2} = \sqrt{2} = 1.41$$

$$d(AC) = \sqrt{(3-1)^2 + (1-2)^2} = \sqrt{5} = 2.23$$

$$d(AD) = \sqrt{(5-1)^2 + (4-2)^2} = \sqrt{20} = 4.47$$

$$d(AE) = \sqrt{(6-1)^2 + (5-2)^2} = \sqrt{34} = 5.83$$

$$d(BC) = \sqrt{(3-2)^2 + (3-1)^2} = \sqrt{5} = 2.23$$

$$d(BD) = \sqrt{(2-5)^2 + (3-4)^2} = \sqrt{10} = 3.16$$

$$d(BE) = \sqrt{(2-6)^2 + (3-5)^2} = \sqrt{20} = 4.47$$

$$d(CD) = \sqrt{(3-5)^2 + (4-1)^2} = \sqrt{13} = 3.61$$

$$d(CE) = \sqrt{(3-6)^2 + (5-4)^2} = \sqrt{25} = 5$$

$$d(DE) = \sqrt{(5-6)^2 + (4-5)^2} = \sqrt{2} = 1.41$$

	A	B	C	D	E
A	0				
B	1.41	0			
C	2.23	2.23	0		
D	4.47	3.16	3.61	0	
E	5.83	4.47	5	1.41	0

① Initial clusters = {A, B, C, D, E}

② min A and B (distance = 1.41)

③ merge A & B

{(AB), C, D, E} = new clusters

$$d\{(AB), C\} = \min\{d(AC), d(BC)\}$$

$$= \min\{2.23, 2.23\}$$

$$d\{(AB), C\} = 2.23$$

$$d\{(AB), D\} = \min\{d(AD), d(BD)\}$$

$$= \min\{4.47, 3.16\}$$

$$d\{(AB), D\} = 3.16$$

$$d\{(AB), E\} = \min\{d(AE), d(BE)\}$$

$$= \min\{5.83, 4.47\}$$

$$d\{(AB), E\} = 4.47$$

	AB	C	D	E
AB	0			
C	2.23	0		
D	3.16	3.61	0	
E	4.47	5	1.41	0

④ Minimum distance is between D and E

d is 1.41

Merge {D and E}

new clusters = {(AB), C, (DE)}

$d\{AB\} \{DE\}$

$$= \min \{ d(AD), d(AE), d(BD), d(DE) \}$$

$$= \min \{ 4.47, 5.83, 3.16, 4.47 \}$$

$$= 3.16$$

	AB	C	DE
AB	0		
C	2.23	0	
DE	3.16	3.61	0

$d\{(C)\} \{DE\}$

$$= \min \{ d(CD), d(CE) \}$$

$$= \min \{ 3.61, 5 \}$$

$$= 3.61$$

$d\{(ABC)\} \{DE\}$

$$= \min \{ d(AD), d(AE), d(BD), d(BE), d(CD), d(CE) \}$$

$$= \min \{ 4.47, 5.83, 3.16, 4.47, 3.16, 5 \}$$

$$= 3.16$$

	ABC	DE
ABC	0	
DE	3.16	0

clusters

① Initial clusters = {A, B, C, D, E}

② A and B dist = 1.41

③ D & E dist = 1.41

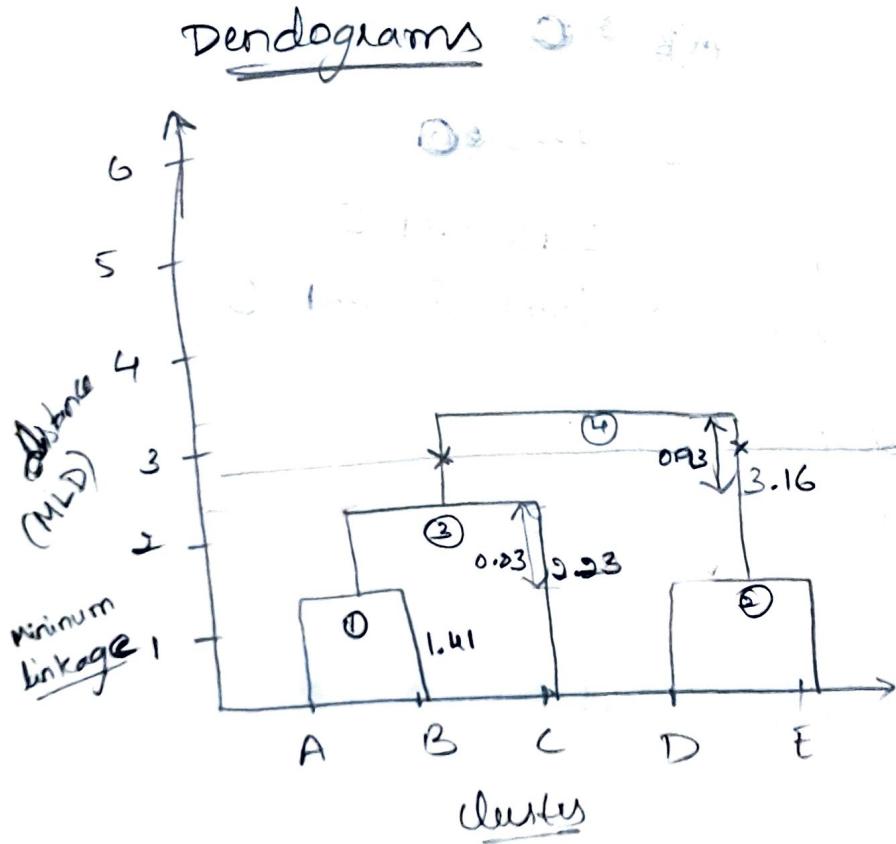
④ (AB) & (C) dist = 2.23

⑤ (ABC) & (DE) dist = 3.16

$$3.16 - 2.23 = 0.93 \checkmark$$

$$2.23 - 1.41 = 0.82$$

= 2 clusters



⑤ Complete linkage (maximum linkage)

(i) first take minimum $\min(AB) = 1.41$

initial clusters = { (AB), C, D, E }

$$\begin{aligned} \text{(ii)} d(AB, C) &= \max\{d(AC), d(BC)\} \\ &= \max\{2.23, 2.23\} \\ &= 2.23 \end{aligned}$$

$$\begin{aligned} d(AB, D) &= \max\{d(AD), d(BD)\} \\ &= \max\{4.47, 3.16\} = 4.47 \end{aligned}$$

$$\begin{aligned} d(AB, E) &= \max\{d(AE), d(BE)\} \\ &= \max\{5.83, 4.47\} = 5.83 \end{aligned}$$

$$\begin{aligned} d(AB, DE) &= \max\{d(AD), d(AE), d(BD), \\ &\quad d(BE)\} \\ &= 5.83 \end{aligned}$$

$$d(DE, C) = \max\{d(DC), d(EC)\} = 5$$

clusters

$$A \in B = 1.41$$

$$D \in E = 1.41$$

$$AB \in C = 2.23$$

$$ABC \in DE = 5.83$$

	AB	C	D	E
AB	0			
C	2.23	0		
D	4.47	3.61	0	
E	5.83	5	1.41	0

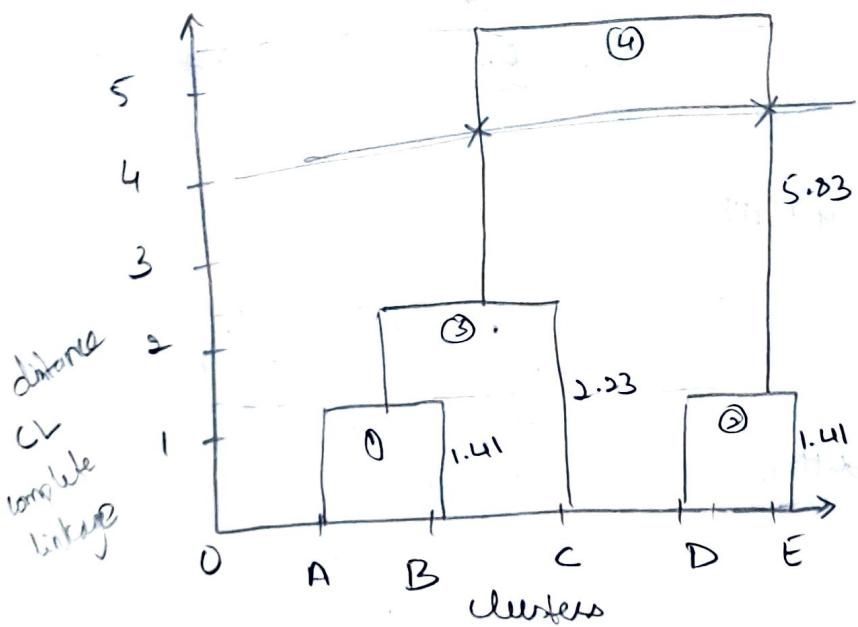
	AB	C	DE
AB	0		
C	2.23	0	
DE	5.83	5	0

	ABC	DE
ABC	0	
DE	5.83	0

$$5.83 - 2.23 = 3.59 \checkmark$$

$$2.23 - 1.41 = 0.82$$

= 2 clusters



Hierarchical Clustering

Use case

targeting audience will form cluster.

it can belong to multiple clusters [soft clustering] and send advertising.

Expectation-Maximization Algorithm

Semi-supervised learning

- if features are observable then learning
if features are unobservable then predict

↓ finds local maximum likelihood of parameters, where latent variables are observed (no labels)

Algorithm

① Expectation step (E-step): algorithm computes the latent variable expectation of the log-likelihood using the current parameter estimates.

② Maximization step (M-step):

With label

θ_1 be probability

Example: C_1 & C_2 of getting head with C_1 , θ_2 be head from C_2

$$\theta_1 = \frac{\text{no of heads with } C_1}{\text{Total no of flips using } C_1}$$

$$\theta_2 = \frac{\text{no of heads with } C_2}{\text{Total no of flips using } C_2}$$

B	H	θ_1	-	-	Coin A	Coin B
A	H	-	-	-	9H, 1T	5H, 5T
A	H	-	-	-	8H, 2T	
B	H	-	-	-		4H, 6T
A	T	-	-	-	7H, 5T	
					24H, 6T	9H, 11T

$$\theta_1 = \frac{24}{24+6}$$

$$\underline{\theta_1 = 0.8}$$

$$\theta_2 = \frac{9}{9+11}$$

$$\underline{\theta_2 = 0.45}$$

Without label



⇒ EM clustering algorithm

H	T	T	H	H	T	H	T
H	H	H	T	H	H	H	H
H	T	H	H	H	H	J	H
H	T	H	H	H	H	H	J
T	H	H	H	H	H	H	H

① Initial values of $\theta_A^{(0)} = 0.6$ $\theta_B^{(0)} = 0.5$

② Expectation - step $P(A) = \frac{L(A)}{L(A)+L(B)}$

$$L(B) = \theta_B^k (1-\theta_B)^{n-k}$$

$$L(B) = (0.5)^6 (1-0.5)^{10-5}$$

$$L(B) = \underline{0.0009765}$$

$$L(A) = \theta_A^k (1-\theta_A)^{n-k}$$

$$L(A) = 0.6^5 (1-0.6)^{10-5}$$

$$L(A) = \underline{\underline{0.0007962}}$$

$k = \text{no of heads}$

$n = \text{total no of clauses}$

$$P(A) = \frac{0.0007962}{0.0007962 + 0.0009765}$$

$$P(A) = \underline{\underline{0.4491}} = 0.445$$

$$P(B) = \frac{0.0009762}{(0.0007962) + 0.0009765}$$

$$P(B) = \underline{\underline{0.55}}$$

Coin A - H	T
$P(A) \times 5$	$P(A) \times 5$
$0.45 \times 5 = 2.25$	$0.45 \times 5 = 2.25$
L(H)	L(T)

Coin B - H	T
$P(B) \times 5$	$P(B) \times 5$
$0.55 \times 5 = 2.75$	$0.55 \times 5 = 2.75$
L(H)	L(T)

Usage of EM algorithm - It can be used to fill the missing data in a sample

- It can be used for the purpose of estimating the parameters of Hidden Markov Model (HMM).

- It can be used for discovering the values of latent variables.

Advantages - likelihood will increase with each iteration.

E-step and M-step are often easy to implementation

Disadvantages - slow convergence.

Clustering Metrics in Machine Learning

Linear regression - mean square error \rightarrow performance measure
 Root mean square error
 Squared error

$$\text{MSE} = \min_{\text{parameters}}$$

} supervised learning

Linear classification - confusion matrix

Roc - Receiver operating characteristics

Clustering (quality of the clusters)

↳ quantitative metrics

- ① Silhouette score
- ② Davis - Bouldin Index
- ③ Calinski - Harabasz Index

! Silhouette analysis - study the separation distance between the resulting clusters.
↳ how close each point in one cluster

+1 = sample is far away from the neighbouring clusters

0 = not properly grouped (overlapping), can be further divided.

-1 = assigned data points to wrong clusters

Calculating Silhouette coefficient

P_1 and P_2 are in cluster 1

P_3 and P_4 are in cluster 2

Point	cluster label
P_1	1
P_2	1
P_3	2
P_4	2

if no distance given - calculate using Euclidean distance

Distance matrix

Point	P_1	P_2	P_3	P_4
P_1	0	0.1	0.65	0.55
P_2	0.1	0	0.7	0.6
P_3	0.65	0.7	0	0.3
P_4	0.55	0.6	0.3	0

$$\textcircled{1} \quad S_C = 1 - \frac{a}{b}$$

point P_1 - within cluster
 $a = \text{Avg}[(P_1 \rightarrow P_2)]$



$a = 0.1$
from one point to all other clusters

$$b = \text{Avg}[(P_1 \rightarrow P_3), (P_1 \rightarrow P_4)]$$

$$b = \frac{0.65 + 0.55}{2} = 0.6$$

$$\text{Silhouette score :}$$

$$\textcircled{2} \quad S_C \text{ for point } P_1 = 1 - \frac{a}{b}$$

$$= 1 - \frac{0.1}{0.6}$$

$$= 0.833$$

point P_2

$$a = \text{Avg}[(P_2 \rightarrow P_1)]$$

$$a = 0.1$$

$$b = \text{Avg}[(P_2 \rightarrow P_3), (P_2 \rightarrow P_4)]$$

$$b = \underline{0.65} \quad (0.7 + 0.6)$$

$$S_C \text{ for point } P_2 = 0.8461$$

① point P₃

$$\alpha = \text{Avg}(P_3 \rightarrow P_4) = 0.3$$

$$b = \text{Avg}[(P_3 \rightarrow P_1)(P_3 \rightarrow P_2)] = \frac{0.65 + 0.7}{2} = 0.675$$

$$SC = 1 - (0.3 / 0.675) = \underline{\underline{0.55}}$$

highly coupled (0.83 & 0.84)

$$\text{Avg SC for } C_1 = \frac{P_1 + P_2}{2} = \underline{\underline{0.8395}}$$

loosely coupled

$$\text{Avg SC for } C_2 = \frac{P_3 + P_4}{2} = \underline{\underline{0.5171}}$$

$$\text{Average of overall cluster} = \frac{0.8395 + 0.5171}{2} = \underline{\underline{0.6723}}$$

② Point P₄

$$\alpha = \text{Avg}(P_4 \rightarrow P_3) = 0.3$$

$$b = \text{Avg}[(P_4 \rightarrow P_1)(P_4 \rightarrow P_2)] = \frac{0.55 + 0.8}{2} = 0.575$$

$$SC = 1 - (0.3 / 0.575) = \underline{\underline{0.4722}}$$

SVM (Support Vector Machine)

- * Maximum margin linear separators
- * finding maximum margin separators
- * kernels methods for non-linear functions
- * varying length pattern classification using SVM.

SVM is a supervised learning model.

Approach 1 - create a label using keywords, for instance "ugent", "cooper", "help". → drawback: mix, hard to implement rules.

Approach 2 - use a supervised machine learning algorithm

label classification. Logistic Regression

SVM is used for classification (DT, NBC, XGBoost, Random Forest, ANN)

SVR (Support Vector Regression) is used for regression.

Vapnik and Chervonenkis → "Maximal margin classifier" → SVM was born
also developed "kernel trick", to classify non-linear data.

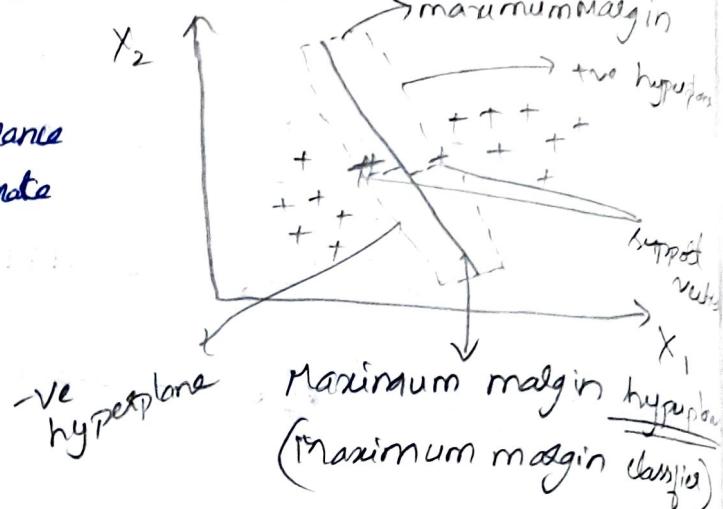
Classification → "soft margin classifier" ⇒ which allows some kind of misclassification.

SVM \rightarrow It is a supervised machine learning uses concept of maximum margin classification.

↓
linear separation

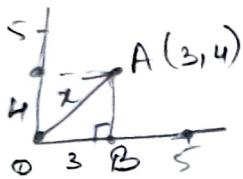
Finding maximum margin between the hyperplanes that means maximum distances between the two classes.

Data points having close boundaries/appearance to other cluster of data points then we note them as vector support vectors.



mathematical calculations

①



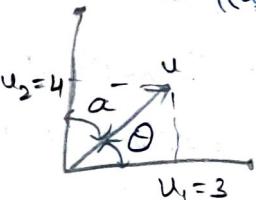
vector = magnitude & directed
A = data point

$$\begin{aligned} OA^2 &= OB^2 + AB^2 \\ OA^2 &= 3^2 + 4^2 \\ OA &= 5 \end{aligned}$$

② The director of a vector $u(u_1, u_2)$ is the vector $w\left(\frac{u_1}{\|u\|}, \frac{u_2}{\|u\|}\right)$

$$w\left(\frac{3}{5}, \frac{4}{5}\right)$$

directed of vector $= w(0.6, 0.8)$



③ sum of two vectors

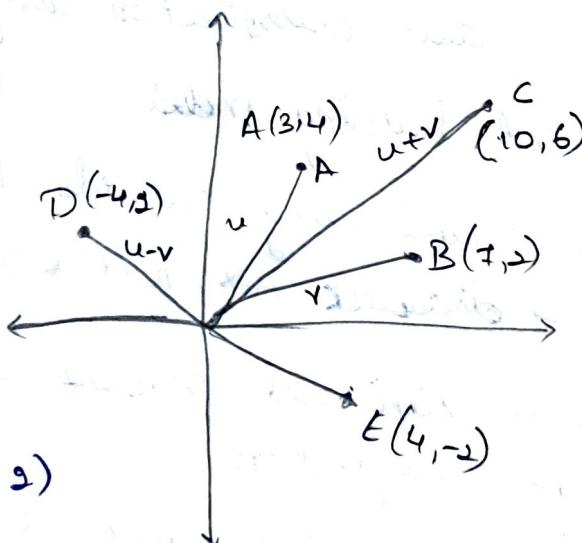
$$u(u_1, u_2) \text{ and } v(v_1, v_2)$$

$$u+v = (u_1+v_1, u_2+v_2)$$

$$\begin{aligned} u_1 &= 3 & u_2 &= 4 \\ v_1 &= 7 & v_2 &= 2 \end{aligned} \Rightarrow (10, 6)$$

$$u-v = (u_1-v_1, u_2-v_2)$$

$$= (3-7, 4-2) = (-4, 2)$$

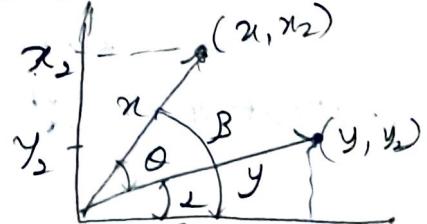


$$v-u = (4, -2)$$

How close the two points are to each other we use dot product.

$$x \cdot y = \|x\| \|y\| \cos(\theta)$$

$$\cos(\theta) = \frac{\text{adjacent}}{\text{hypotenuse}}$$



$$\theta = \beta - \alpha$$

$$\cos(\beta - \alpha) = \cos(\beta)\cos(\alpha) + \sin(\beta)\sin(\alpha)$$

$$\cos(\beta - \alpha) = \frac{x_1}{\|x\|} \cdot \frac{y_1}{\|y\|} + \frac{x_2}{\|x\|} \cdot \frac{y_2}{\|y\|}$$

$$\cos \theta = \frac{x_1 y_1 + x_2 y_2}{\|x\| \|y\|}$$

$$\boxed{\|x\| \|y\| \cos \theta = x_1 y_1 + x_2 y_2}$$

$$x \cdot y = x_1 y_1 + x_2 y_2 + \dots + x_n y_n$$

$$\cos \beta = \frac{x_1}{\|x\|} \quad \cos \alpha = \frac{y_1}{\|y\|}$$

$$\sin(\beta) = \frac{\text{opposite}}{\text{hypotenuse}} \quad \sin(\beta) = \frac{x_2}{\|x\|}$$

$$\sin(\alpha) = \frac{y_2}{\|y\|}$$

Example

Linear Equation Example

Suppose we are given the following positively labeled data points in \mathbb{R}^2 : $\{(1, 0), (0, 1), (-1, -1), (0, -1)\}$ and the following negatively labeled data points in \mathbb{R}^2 : $\{(-1, 1), (1, 1), (1, -1), (-1, -1)\}$

Support vectors = $S_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, S_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, S_3 = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$

$$\text{Support vectors} = S_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, S_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, S_3 = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

① Add bias

$$\alpha_1 S_1 + \alpha_2 S_2 + \alpha_3 S_3 = -1 = \alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

$$\alpha_1 S_1 + \alpha_2 S_2 + \alpha_3 S_3 = +1$$

$$\alpha_1 S_1 + \alpha_2 S_2 + \alpha_3 S_3 = +1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} -1 \\ -1 \end{pmatrix} = 1$$

$$4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1 \quad \text{--- (C1)}$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 1$$

$$4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1 \quad \text{--- (C2)}$$

$$4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1 \quad \text{--- (C3)}$$

$$4\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$

$$-(4\alpha_1 + 11\alpha_2 + 9\alpha_3) = 1$$

$$-2\alpha_1 - 7\alpha_2 - 5\alpha_3 = -2$$

$$-2\alpha_2 + 2\alpha_3 = 0$$

$$-2\alpha_1 + 9\alpha_2 - 3\alpha_2 = -2$$

$$-2\alpha_1 - 6\alpha_2 = -2$$

$$\alpha_1 + 3\alpha_2 = 1$$

$$4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1$$

$$-(4\alpha_1 + 9\alpha_2 + 11\alpha_3) = 1$$

$$+\alpha_1 + 2\alpha_2 - 2\alpha_3 = 0$$

$$\cancel{\alpha_3} =$$

$$\alpha_2 = \alpha_3$$

$$3\alpha_2 = \alpha_2$$

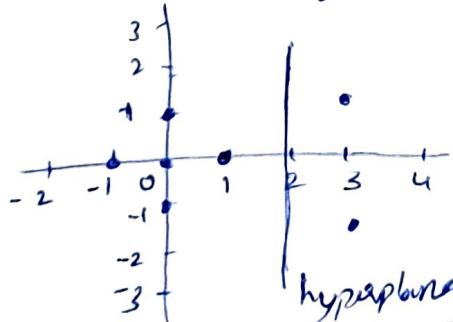
$$\cancel{\alpha_2} = 1$$

$$\alpha_3 = 0.75 \quad \alpha_2 = \frac{3}{4} = 0.75$$

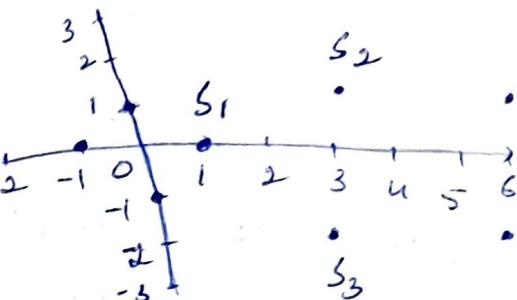
② parallel to y-axis

③ parallel to x-axis

④ origin



hyperplane



$$4\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$

$$4\alpha_1 + 4(0.75) + 4(0.75) = -1$$

$$2\alpha_1 + 3 + 3 = -1$$

$$2\alpha_1 = -1 - 3 - 3 = -7$$

$$\alpha_1 = -\frac{7}{2} = -3.5$$

⑤ weight matrix

$$\bar{w} = \sum_{i=1}^3 \alpha_i \bar{s}_i$$

$$\bar{w} = \alpha_1 \bar{s}_1 + \alpha_2 \bar{s}_2 + \alpha_3 \bar{s}_3$$

$$\bar{w} = -3.5 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 0.75 \begin{pmatrix} 0 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

$$\bar{w} = \begin{pmatrix} -3.5 & 0 & -3.5 \end{pmatrix} \begin{pmatrix} 3.75 & 2.25 \\ 0.75 & 0.75 \end{pmatrix} = \begin{pmatrix} +1 \\ 0 \\ -2 \end{pmatrix}$$

Non linear classifier

mapping to a higher dimension

1D to 2D to 3D
mapping function

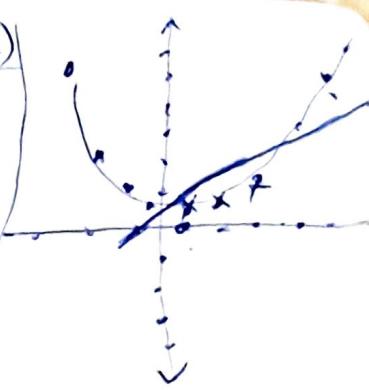


Kernel Trick

Gaussian RBF Kernel

$$K(\vec{x}, \vec{l}^i) = e^{-\frac{1}{2\sigma^2} \|\vec{x} - \vec{l}^i\|^2}$$

X	X-5	$(X-5)^2 / (\sigma^2)$
2	-3	9
3	-2	4
4	-1	1
6	1	1
7	2	4
8	3	9
9	4	16
10	5	25
11		



$$e^{-\alpha} = 0$$

$$e^{\alpha} = 1$$