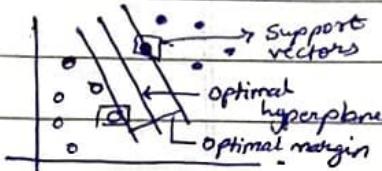


# SVM

- Determining location of hyperplane that produce optimal separation of classes among infinite no. of boundaries
- Optimal hyperplane - constructed by searching for maximal marginal hyperplane
  - ⇒ Sum of distances to the hyperplane from closest points of two classes
  - ⇒ Sides of margin are parallel



## Linear SVM

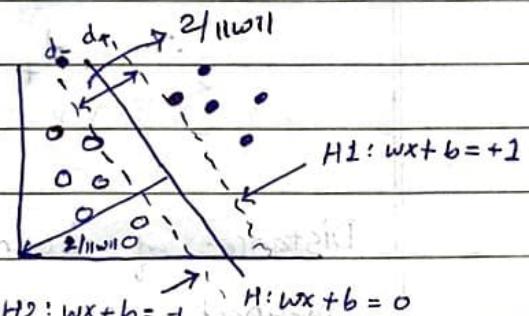
Dataset  $D = \{1-d \rightarrow \{x_i, y_i\}\}_{i=1}^d = \{1, 1\}$

(no. of samples)

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

Eqn of separating hyperplane :

$$\vec{w} \cdot \vec{x} + b = 0$$



1<sup>st</sup> distance of  $(x_0, y_0)$  from hyperplane is :

$$\frac{w_1 x_0 + w_2 y_0 + b}{\sqrt{w_1^2 + w_2^2}} \neq \frac{|b|}{\|w\|}$$

$$\therefore 1^{st} \text{ distance from origin : } \frac{|b|}{\|w\|}$$

For points lying on or above hyperplane,

$$H_1 : \vec{x}_i \cdot \vec{w} + b \geq +1 \quad \text{for } y_i = +1$$

$$H_2 : \vec{x}_i \cdot \vec{w} + b \leq -1 \quad \text{for } y_i = -1$$

Both equations can be combined into one set of inequalities as

$$\boxed{y_i(\vec{x}_i \cdot \vec{w} + b) - 1 \geq 0} \quad \text{for every } i$$

Distance of bounding hyperplane  $H_1$  from origin with normal  $w$  is,

$$\text{dist}(H_1) = \frac{|w_1x_1 + w_2x_2 + b - 1|}{\sqrt{w_1^2 + w_2^2}} = \frac{|b - 1|}{\|w\|}$$

$$\therefore d_+ = \frac{|b - 1|}{\|w\|} - \frac{|b|}{\|w\|} = \frac{|1|}{\|w\|}$$

Distance of bounding hyperplane  $H_2$  from origin with normal  $w$  is,

$$\text{dist}(H_2) = \frac{|b + 1|}{\|w\|}$$

$$\therefore d_- = \frac{|b|}{\|w\|} - \frac{|1 + b|}{\|w\|} = \frac{-1}{\|w\|}$$

Hence, margin bet<sup>n</sup> two bounding hyperplanes  
is  $\frac{2}{\|w\|}$ .

Objective fn  $\rightarrow$  Maximize  $\frac{2}{\|\omega\|}$

$$\sim \text{Minimize } \frac{\|\vec{\omega}\|}{2}$$

$$\sim \text{Minimize } \frac{\|\vec{\omega}\|^2}{2}$$

Thus, SVM solves objective fn  $\Phi(\omega) = \frac{1}{2} \|\vec{\omega}\|^2$

subject to constraint  $y_i(\vec{x}_i \cdot \vec{\omega} + b) \geq 1$

$\rightarrow$  Convex quadratic programming problem which has single global optimum soln. Solved in dual space, the space of Lagrange multipliers

Lagrangian of this problem is given by

$$L(\vec{\omega}, b, \alpha) = \frac{1}{2} \|\omega\|^2 - \sum_{i=1}^l \alpha_i (y_i(\vec{\omega} \cdot \vec{x}_i + b) - 1) \quad \dots \dots \dots \text{(i)}$$

Lagrange dual function of optimization problem is given by

$$\max_{\alpha} w(\alpha) = \max_{\alpha} (\min_{\omega, b} L(\vec{\omega}, b, \alpha))$$

Diff (i) wrt  $\omega$  &  $b$  and setting derivatives equal to 0,

$$\frac{\partial L(\vec{\omega}, b, \alpha)}{\partial \omega} = \vec{\omega} - \sum_{i=1}^l \alpha_i y_i \vec{x}_i = 0$$

$$\therefore \vec{\omega} = \sum_{i=1}^l \alpha_i y_i \vec{x}_i \quad \dots \dots \text{(ii)}$$

$$\frac{\partial L(\vec{w}, b, \alpha)}{\partial b} = \sum_{i=1}^l \alpha_i y_i = 0 \quad \xrightarrow{\text{(iii)}}$$

Substituting (ii) and (iii) in (i),

$$w(\alpha) = \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j - \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j - b \cdot (0) + \sum_{i=1}^l \alpha_i$$

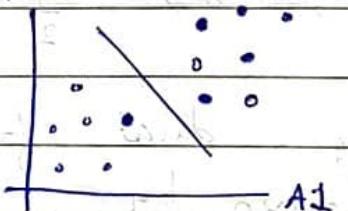
∴ Dual problem can be now stated as :-

$$\text{Maximize: } w(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

Subject to  $\sum_{i=1}^l \alpha_i y_i = 0$  and  $\alpha_i \geq 0$

#

### Soft Margin SVM



- Non separable data - contains noisy data
- Allows mislabelled data points while still maximizing margin
- Slack variables - measure degree of misclassification

Constraints are now:  $\vec{x}_i \cdot \vec{w} + b \geq +1 - \xi_i$  for  $y_i = +1$   
 $\vec{x}_i \cdot \vec{w} + b \leq -1 + \xi_i$  for  $y_i = -1$   
 $\xi_i \geq 0, \forall i$

Combining constraints, it becomes

$$y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0$$

Objective fn now becomes,

$$\text{Minimize : } \Phi(w, b, \xi_i) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i$$

$$\text{Subject to: } y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 - \xi_i \quad \& \quad \xi_i \geq 0$$

$C$  = Regularization parameter  
Controls trade off between  
maximizing margin & minimizing  
error term.

→ Lagrangian of primal form of objective fn is given:

$$L(w, b, \alpha, \beta, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \alpha_i (y_i (\vec{w} \cdot \vec{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^l \beta_i \xi_i$$

→ Solution of above eqn is given as :

$$\max_w L(\alpha, \beta) = \max_{w, b, \xi} (L(w, b, \alpha, \beta, \xi))$$

$$\text{Diff. we get, } \frac{\partial L}{\partial w} = w = \sum_{i=1}^l \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^l \alpha_i y_i = 0$$

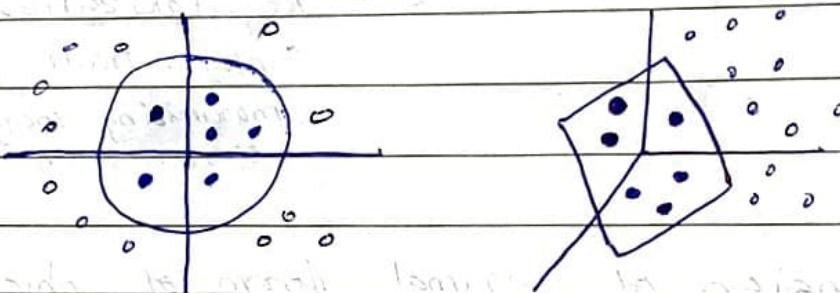
$$\frac{\partial L}{\partial \xi} = \alpha_i + \beta_i = C$$

Substituting this into primal form, we get

$$\text{Maximize : } w(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i y_j y_i \vec{\vec{x}}_i \vec{x}_j$$

Subject to :  $\sum_{i=1}^l \alpha_i y_i = 0$  and  $c \geq \alpha_i \geq 0$

### NON LINEAR SVM



Step I Transform input vectors into high dimensional feature space

Step II Use SVM to find hyperplane of maximal margin in new feature space

"Cover's Theorem" - A complex pattern classification problem cast in high dimensional space non linearly is more likely to be linearly separable than in a low-dimensional space.

Let  $\Phi$  be nonlinear mapping function,

$$\Phi: \mathbb{R}^d \rightarrow \mathcal{H} \quad \text{where } x_i \in \mathbb{R}^d \text{ and } \Phi(x_i) \in \mathcal{H}$$

With this transformation, eq<sup>n</sup> of hyperplane becomes

$$\vec{w} \cdot \Phi(x) - b = 0$$

Using  $\Phi(\cdot)$  function, weight becomes

$$\vec{w} = \sum_{i=1}^{N_s} \alpha_i y_i \Phi(\vec{x}_i)$$

Decision function becomes

$$f(x) = \text{sgn} (\Phi(x) \cdot \vec{w} + b)$$

$$f(x) = \text{sgn} \left( \sum_{i=1}^{N_s} \alpha_i y_i \Phi(\vec{x}_i) \cdot \Phi(\vec{x}) + b \right)$$

Dual problem can be rewritten as

$$w(x) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \Phi(\vec{x}_i) \cdot \Phi(\vec{x}_j)$$

→ Feature mapping function always appears as dot products  $\Phi(\vec{x}_i) \cdot \Phi(\vec{x}_j)$

Computing it in high dimensional can be complex, costly or suffer from curse of dimensionality.

e.g.  $\Phi(\vec{x}_i) = (1, x_{i1}^2, \sqrt{2}x_{i1}x_{i2}, x_{i2}^2, \sqrt{2}x_{i1}, \sqrt{2}x_{i2})$

$$\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^6$$

$$\therefore \Phi(\vec{x}_1) \cdot \Phi(\vec{x}_2) = 1 + 2x_{11}x_{21} + 2x_{12}x_{22} + (x_{11}x_{21} + x_{11}x_{21})^2$$

$$\text{Instead, } (\vec{x}_1 \cdot \vec{x}_2 + 1)^2 = (x_{11}x_{21} + x_{11}x_{21})^2 + 2x_{11}x_{21} + 2x_{12}x_{22} + 1 \\ = \Phi(\vec{x}_1) \cdot \Phi(\vec{x}_2)$$

### Kernel Functions

→ SVM learning req only on dot product  $\Phi(\vec{x}_i) \cdot \Phi(\vec{x}_j)$  bet<sup>n</sup> training examples.

→ Operations in high dimensional space do not need to be performed explicitly if we find a function  $K(\vec{x}_i, \vec{x}_j)$  such that

$$K(\vec{x}_i, \vec{x}_j) = \Phi(\vec{x}_i) \cdot \Phi(\vec{x}_j)$$



### Kernel function

Allows to calculate dot product of  $(\Phi(\vec{x}_i), \Phi(\vec{x}_j))$  without explicitly applying function to input vectors.

Dual Problem is now defined using Kernel function as:

$$w(d) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j K(\vec{x}_i, \vec{x}_j)$$

subject to  $\sum_i \alpha_i y_i = 0$  and  $c \geq \alpha_i \geq 0$

Classification function becomes

$$f(\vec{x}) = \text{sgn} \left( \sum_{i=1}^{N_s} \alpha_i y_i K(\vec{x}_i, \vec{x}) + b \right)$$

→ kernel function  $K$  is valid iff

$$K(\vec{x}, \vec{y}) = \Phi(\vec{x}) \cdot \Phi(\vec{y})$$

→ If  $K$  is valid kernel, then it is:

\* Symmetric  $[K(\vec{x}, \vec{y}) = K(\vec{y}, \vec{x}) \text{ for any two objects } \vec{x}, \vec{y} \in \mathbb{R}]$

\* Positive (semi-) definite

$$[z^T K z = \sum_{i=1}^l \sum_{j=1}^l z_i z_j K(\vec{x}_i, \vec{x}_j) \geq 0]$$

### Some Kernels

Linear  $K(\vec{x}_i, \vec{x}_j) = (\vec{x}_i \cdot \vec{x}_j)$

Polynomial  $K(\vec{x}_i, \vec{x}_j) = (\vec{x}_i \cdot \vec{x}_j + 1)^d$

Radial Basis Function (RBF)  $K(\vec{x}_i, \vec{x}_j) = e^{-\|\vec{x}_i - \vec{x}_j\|^2 / 2\sigma^2}$

Sigmoid  $K(\vec{x}_i, \vec{x}_j) = \tanh(K(\vec{x}_i \cdot \vec{x}_j) + \beta)$

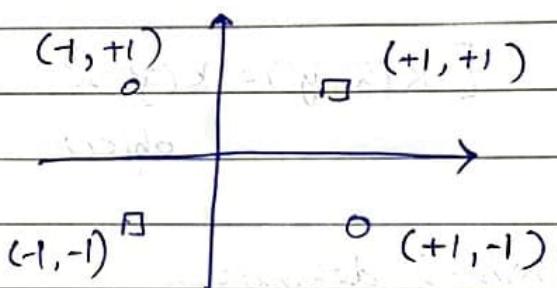
### Advantages of SVM

- Complexity is characterized by no. of support vectors rather than dimension of data. Hence less prone to overfitting
- Good generalization capability by maximizing margin in input space
- No. of SVs found can be used to compute upper bound on expected rate of classifier.

## # Non linear SVM example (XOR problem)

Training set :-	Input vector $X$	Output $y$
$x_1$	(-1, -1)	-1
$x_2$	(-1, +1)	+1
$x_3$	(+1, -1)	+1
$x_4$	(+1, +1)	-1

Plotting training points in 2D space



We can see they are not linearly separable.  
Let's take

$$K(\vec{x}_i, \vec{x}_j) = (1 + \vec{x}_i \cdot \vec{x}_j)^2$$

Image of input vector  $\vec{x}_i$  induced in feature space is therefore deduced to be

$$\Phi(\vec{x}_i) = (1, x_{i1}^2, \sqrt{2}x_{i1}x_{i2}, x_{i2}^2, \sqrt{2}x_{i1}, \sqrt{2}x_{i2})$$

Now objective function of dual form can be written as

$$w(x) = \sum_{i=1}^4 \alpha_i - \frac{1}{2} \sum_{i=1}^4 \sum_{j=1}^4 \alpha_i \alpha_j y_i y_j K(\vec{x}_i, \vec{x}_j)$$

$$K(\vec{x}_1, \vec{x}_2) = (1+0)^2 = 1$$

$$K(\vec{x}_1, \vec{x}_1) = (1+2)^2 = 9$$

$$K(\vec{x}_1, \vec{x}_3) = (1+0)^2 = 1$$

$$K(\vec{x}_1, \vec{x}_4) = (1+2)^2 = 9$$

$$K(\vec{x}_2, \vec{x}_3) = (1+1-2)^2 = 1$$

$$K(\vec{x}_2, \vec{x}_4) = (1+0)^2 = 1$$

$$K(\vec{x}_3, \vec{x}_3) = (1+2)^2 = 9$$

$$K(\vec{x}_3, \vec{x}_4) = (1+0)^2 = 1$$

$$K(\vec{x}_4, \vec{x}_4) = (1+2)^2 = 9$$

$$\therefore w(\alpha) = \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 - \frac{1}{2} \left( 9\alpha_1^2 - \alpha_1\alpha_2 - \alpha_1\alpha_3 + \alpha_1\alpha_4 - \alpha_1\alpha_2 + 9\alpha_2^2 + \alpha_2\alpha_3 - \alpha_2\alpha_4 - \alpha_1\alpha_3 + \alpha_2\alpha_3 + 9\alpha_3^2 - \alpha_3\alpha_4 + 9\alpha_4^2 + \alpha_1\alpha_4 - \alpha_2\alpha_4 - \alpha_3\alpha_4 \right)$$

$$= \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 - \frac{1}{2} \left( 9\alpha_1^2 + 9\alpha_2^2 + 9\alpha_3^2 + 9\alpha_4^2 - 2\alpha_1\alpha_2 + 2\alpha_1\alpha_3 - 2\alpha_1\alpha_4 - 2\alpha_2\alpha_3 + 2\alpha_2\alpha_4 - 2\alpha_3\alpha_4 \right) \quad (i)$$

Optimising  $w(\alpha)$  w.r.t. Lagrange multipliers gives following eq's:

$$9\alpha_1 - \alpha_2 - \alpha_3 + \alpha_4 = 1$$

$$-\alpha_1 + 9\alpha_2 + \alpha_3 - \alpha_4 = 1$$

$$-\alpha_1 + \alpha_2 + 9\alpha_3 - \alpha_4 = 1$$

$$\alpha_1 - \alpha_2 - \alpha_3 + 9\alpha_4 = 1$$

Hence optimal values of Lagrange multipliers are

$$\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = \frac{1}{8}$$

Substituting values of  $\alpha$  in (i), we get

$$w(x) = \frac{1}{8} \frac{1}{4}$$

$$\therefore \frac{1}{2} \|\vec{w}\|^2 = \frac{1}{4(4+1)}$$

$$\therefore \|\vec{w}\|^2 = \frac{1}{2} \Rightarrow \|w\| = \frac{1}{\sqrt{2}}$$

Now  $w = \sum_{i=1}^4 \alpha_i y_i \Phi(x_i)$

$$= \frac{1}{8} [y_1 \Phi(\vec{x}_1) + y_2 \Phi(\vec{x}_2) + y_3 \Phi(\vec{x}_3) + y_4 \Phi(\vec{x}_4)]$$

$$= \frac{1}{8} [-\Phi(\vec{x}_1) + \Phi(\vec{x}_2) + \Phi(\vec{x}_3) - \Phi(\vec{x}_4)]$$

$$= \frac{1}{8} \left[ -\begin{bmatrix} 1 \\ \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} + \begin{bmatrix} 1 \\ -\frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} + \begin{bmatrix} 1 \\ -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} - \begin{bmatrix} 1 \\ \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix} \right]$$

$$= \begin{bmatrix} 0 \\ 0 \\ -\frac{1}{\sqrt{2}} \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

bias b is 0 because first element of w is 0.

∴ Optimal hyperplane becomes :

$$\vec{w} \cdot \Phi(\vec{x}) = [0 \ 0 \ -\frac{1}{\sqrt{2}} \ 0 \ 0 \ 0] \begin{bmatrix} 1 \\ x_1^2 \\ x_2^2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_1 x_2 \end{bmatrix} = 0$$

$\therefore \boxed{-x_1 x_2 = 0}$  — Ans

## KNN

- Essential classification algs in ML
- Supervised learning domain
- Non parametric meaning doesn't make any underlying assumptions about distribution of data
- Instance based or lazy learning where  $f^n$  is only approx. locally & all computation is deferred until classification.

Process :

- Let  $m$  be no. of training data samples. Let  $p$  be unknown point.

1. Store training samples in an array of datapoints arr[ ]. This means each element of this array represents tuple  $(x_1, x_2, \dots, x_n)$ .
2. for  $i=0$  to  $m$  : calculate cosine similarity  $\cos\phi = \frac{(a \cdot b)}{\|a\| \|b\|}$  where  $a$  is input tuple &  $b$  is targeted tuple.
3. Make set  $S$  of  $K$  smallest distances obtained. Each of these distances corresponds to an already classified datapoint.
4. Return majority label / distance weighted label among  $S$ .

## Cosine similarity

$$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|} = \frac{\sum_1^n a_i b_i}{\sqrt{\sum_1^n a_i^2} \sqrt{\sum_1^n b_i^2}}$$

where  $\vec{a} \cdot \vec{b} = \sum_1^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$  is  
dot product of two vectors

### # Solved example

History Science Research offers Students Hall Class

An	0.537	0.477	0.673	0.177	A
Ar			0.981	0.195	0.196 B
Bi	0.317	0.924		0.111	0.112 A
Ch	0.195			0.155	0.158 A
Comm			0.78	0.626	0 B
Comp	0.989			0.130	0.067 A
Cr				1	0 B
Ec		1		0	0 A
En			0.98	0	0.199 B
Ge	0.849			0.529	0 A
Hi	0.991		0.135	0	0 B
Ma		0.616	0.549	0.198	0.201 A
Mo	0			0.17	0.373 B
Mu	0.97			0	0.172 B
Phi	0.741		0.458	0.315	0.136 B
Phy		0.894		0.062	0.318 A
Po	0.933	0.348	0.387	0.313	0.053 A
Ps		0.852	0.57	0.0459	0.162 A
So		0.639		0.967	0.237 A
Th	0.537	0.477	0.967	0.253	0.177 A
An			0.673		0.177 A

Find  $\text{class}(T_h) = ?$

Calculate cosine similarity with each training sample & return results in descending order.

Document	Class	Similarity to Theatre ( $T_h$ )
Cr	B	0.967075
An	A	0.695979
Co	B	0.605667
Gr	A	0.570389
So	A	0.5
Ph	A	0.38
Ps	A	0.34
Ma	A	0.21
Ar	B	0.23
Mn	B	0.2
Ch	A	0.18
Comp	A	0.14
Bi	A	0.13
Mo	B	0.09
Po	A	0.07
En	O	0.05
Ph	B	0.03
Hi	B	0
Ec	A	0

Hence for 1-NN : B

for 3-NN : B

for 5-NN : A

Distance weighted 3-NN : B

19-NN : A

**Q** How to choose  $k$  ?

→ No straightforward method to calculate  $k$ .

Thumb rule is  $K = \sqrt{n}$

↳ no. of training examples

Imp  $\Rightarrow K$  should be odd

$\Rightarrow K$  should not be multiple of no. of classes

$\Rightarrow$  shouldn't be too small or too large

greater influence  
of noise

biased towards  
highly probable class

## CLUSTERING

→ Organization of unlabeled data into similarity groups called clusters.

→ Cluster is collection of data items which are similar between them & dissimilar to data items in other clusters.

### Evaluation

#### Intra cluster cohesion

→ How near data points in cluster are to cluster centroid

→ Sum of squared error is common measure

#### Inter cluster separation

→ Diff. cluster centroids should be far away from each other

## # K Means Clustering

→ Partitional clustering algorithm

→ Let set of datapoints  $D$  be  $\{x_1, x_2, \dots, x_n\}$

Here  $x_i = (x_{i1}, x_{i2}, \dots, x_{iK})$

( $K \rightarrow$  no. of dimensions)

→ K-means algo partitions given data into  $K$  clusters



Each cluster has  
centre called  
centroid.

### Process

Given  $K$ , algo works as follows:

1. Choose  $K$  (random) datapoints as ini centroids
2. Assign each datapoint to closest centroid
3. Recompute centroids using current cluster membership
4. If convergence criterion is not met,  
repeats steps 2 and 3.

### Stopping criteria

- No reassessments of datapoints to diff cluster
- No change of centroids
- min. decrease in sum of squared error (SSE)

$$SSE = \sum_{j=1}^k \sum_{x \in C_j} d(x, m_j)^2$$

Euclidean distance bet<sup>n</sup>  $x$  & centroid  $m_j$

• Solved example

Use k means algo & euclidean distance to cluster the following 8 examples into 3 clusters:

$$A_1 = (2, 10) \quad A_2 = (2, 5) \quad A_3 = (8, 4) \quad A_4 = (5, 8)$$

$$A_5 = (7, 5) \quad A_6 = (6, 4) \quad A_7 = (1, 2) \quad A_8 = (4, 9)$$

	1	2	3	4	5	6	7	8
1	0	5	6	$\sqrt{13}$	$\sqrt{50}$	$\sqrt{52}$	$\sqrt{65}$	$\sqrt{5}$
2	0	$\sqrt{37}$	$\sqrt{18}$	5	$\sqrt{17}$	$\sqrt{10}$	$\sqrt{20}$	
3			0	5	$\sqrt{2}$	$\sqrt{2}$	$\sqrt{53}$	$\sqrt{41}$
4				0	$\sqrt{13}$	$\sqrt{17}$	$\sqrt{52}$	$\sqrt{2}$
5					0	$\sqrt{2}$	$\sqrt{45}$	$\sqrt{5}$
6						0	$\sqrt{29}$	$\sqrt{29}$
7							0	$\sqrt{58}$
8								0

Suppose initial seeds are  $A_1$ ,  $A_4$  and  $A_7$ . Run k-means algo for 1 epoch. At end of this epoch show:

a) New clusters

b) Centre of new clusters

c) Draw  $10 \times 10$  space with all 8 points & show the clusters after first epoch & new centroids.

d) How many more iterations are needed to converge?

$$\text{Seed 1} = A_1 = (2, 10)$$

$$\text{Seed 2} = A_4 = (5, 8)$$

$$\text{Seed 3} = A_7 = (1, 2)$$

Epoch 1.

$$A_1: d(A_1, \text{seed 1}) = 0 \quad A_2: d(A_2, \text{seed 1}) = 5 \quad A_3: d(A_3, \text{seed 1}) = 6$$

$$d(A_1, \text{seed 2}) = \sqrt{13}$$

$$d(A_2, \text{seed 2}) = \sqrt{18}$$

$$d(A_3, \text{seed 2}) = 5$$

$$d(A_1, \text{seed 3}) = \sqrt{65}$$

$$d(A_2, \text{seed 3}) = \sqrt{50}$$

$$d(A_3, \text{seed 3}) = \sqrt{53}$$

$$\rightarrow A_1 \in \text{cluster 1} \quad \rightarrow A_2 \in \text{cluster 3} \quad \rightarrow A_3 \in \text{cluster 2}$$

$$A_4: d(A_4, \text{seed 1}) = \sqrt{13} \quad A_5: d(A_5, \text{seed 1}) = \sqrt{50} \quad A_6: d(A_6, \text{seed 1}) = \sqrt{52}$$

$$d(A_4, \text{seed 2}) = 0$$

$$d(A_5, \text{seed 2}) = \sqrt{13}$$

$$d(A_6, \text{seed 2}) = \sqrt{17}$$

$$d(A_4, \text{seed 3}) = \sqrt{52}$$

$$d(A_5, \text{seed 3}) = \sqrt{45}$$

$$d(A_6, \text{seed 3}) = \sqrt{29}$$

$$\rightarrow A_4 \in \text{cluster 2} \quad \rightarrow A_5 \in \text{cluster 2} \quad \rightarrow A_6 \in \text{cluster 2}$$

$$A_7: d(A_7, \text{seed 1}) = \sqrt{65} \quad A_8: d(A_8, \text{seed 1}) = \sqrt{5}$$

$$d(A_7, \text{seed 2}) = \sqrt{52}$$

$$d(A_8, \text{seed 2}) = \sqrt{2}$$

$$d(A_7, \text{seed 3}) = 0$$

$$d(A_8, \text{seed 3}) = \sqrt{58}$$

$$\rightarrow A_7 \in \text{cluster 3}$$

$$\rightarrow A_8 \in \text{cluster 2}$$

End of epoch 1. new clusters :

$$1 : \{A_1, A_3\}$$

$$2 : \{A_4, A_5, A_6, A_8\}$$

$$3 : \{A_2, A_7\}$$

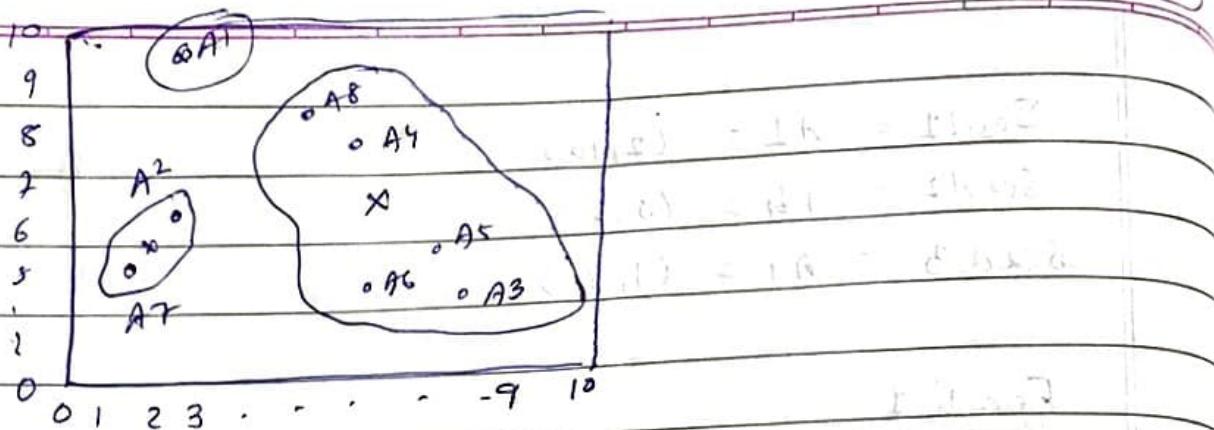
Centres of new clusters :

$$C_1 = (2, 10)$$

$$C_2 = \left( \frac{8+5+7+6+4}{5}, \frac{4+8+5+4+9}{5} \right)$$

$$= (6, 6)$$

$$C_3 = \left( \frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$



We would need 2 more epochs. After 2<sup>nd</sup> epoch results would be:

1: {A1, A8}, 2: {A3, A4, A5, A6}, 3: {A2, A7}  
with centres  $C_1 = (3, 9.5)$ ,  $C_2 = (6.5, 5.25)$  &  $C_3 = (1.5, 3.5)$

After 3<sup>rd</sup> epoch, results:

1: {A1, A4, A8}, 2: {A3, A5, A6}, 3: {A2, A7}  
with centres  $C_1 = (3.66, 9)$ ,  $C_2 = (7, 4.33)$  and  $C_3 = (1.5, 3.5)$

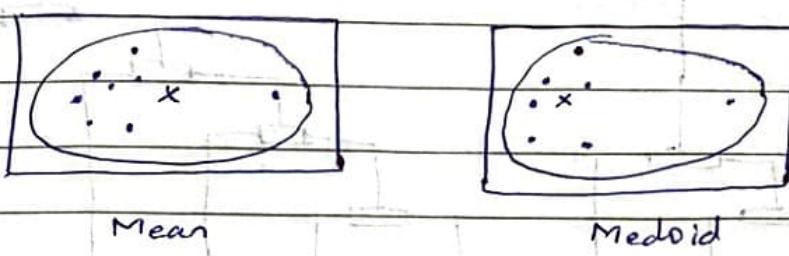
Make two plots one for each 2<sup>nd</sup> & 3<sup>rd</sup> epoch as above.

### Limitations :

- Sensitive to outliers as mean is easily influenced by extreme values.
- Sensitive to initial seed

## F K Medoids (Partitioning around medoids (PAM) ) algo

- K-medoids uses an actual point in the cluster to represent it.
- Medoid is most centrally located object of cluster, with minimum sum of distances to other points.



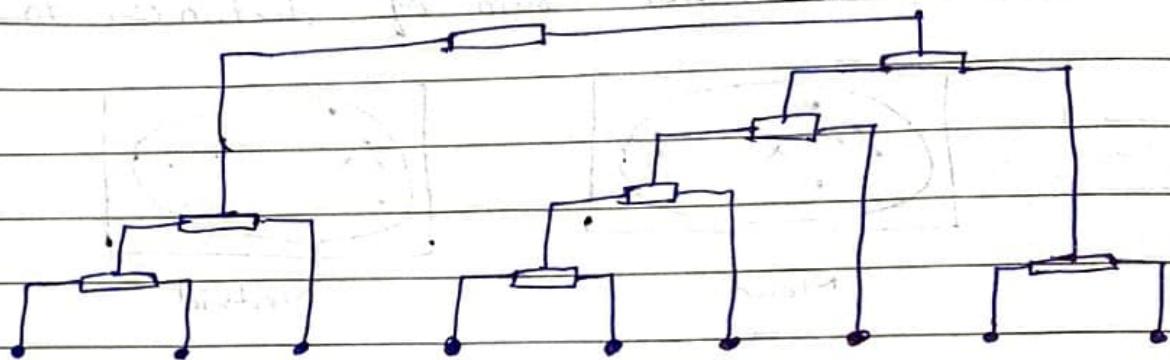
- Mean is greatly influenced by outlier & thus cannot represent correct cluster centre. While medoid is robust to outlier & correctly represents the cluster centre.

## H Hierarchical Agglomerative Clustering

- Works by grouping data objects into hierarchy or "tree" of clusters.
- "Hierarchical Agglomerative approach" uses bottom up strategy.
- Starts by letting each object form its own cluster & iteratively merges clusters into larger and larger clusters, until all the objects are in single cluster or certain termination conditions are satisfied.

## Dendrogram

- Tree that shows how clusters are merged / split hierarchically.
- Each node on tree is cluster; each leaf node is singleton cluster.



## Agglomerative Clustering

### Process :

1. Compute distance matrix
2. Let each data point be cluster
3. Repeat
  1. Merge two closest clusters
  2. Update distance matrix
4. Until only a single cluster remains

### Distance measures

Min distance :  $\text{dist}_{\min}(C_i, C_j) = \min_{p \in C_i, p' \in C_j} \{ \|p - p'\| \}$

Max distance :  $\text{dist}_{\max}(C_i, C_j) = \max_{p \in C_i, p' \in C_j} \{ \|p - p'\| \}$

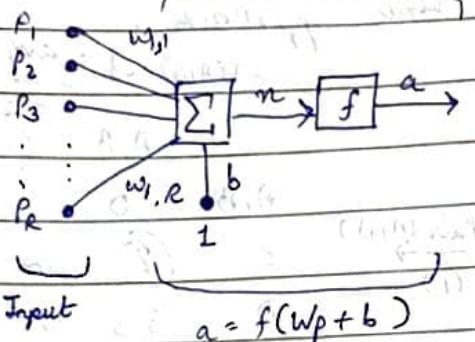
Mean distance :  $\text{dist}_{\text{mean}}(C_i, C_j) = \|m_i - m_j\|$

Avg distance :  $\text{dist}_{\text{avg}}(C_i, C_j) = \frac{1}{n_i n_j} \sum_{p \in C_i, p' \in C_j} \|p - p'\|$

## Neural Networks

### Neuron with vector input

Neuron w Vector Input



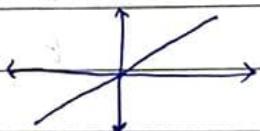
where

R = no. of elements in  
input vector.

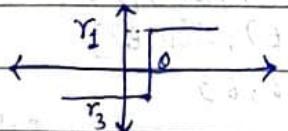
$$n = w_{1,1} p_1 + w_{1,2} p_2 + \dots + w_{1,R} p_R + b$$

### Activation Functions

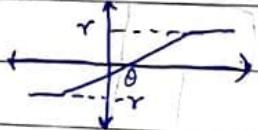
Linear



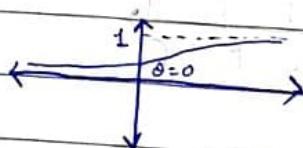
Step



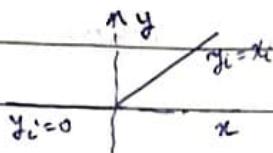
Ramp



Sigmoid



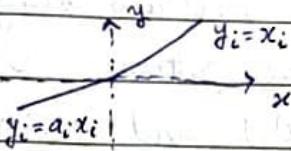
ReLU



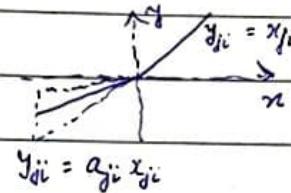
$$f(\text{net}) = \max(0, \text{net})$$

ReLU

Leaky ReLU



$$0.01 \times x_i / \alpha_i \times x_i$$

Randomized  
Leaky ReLU

$$y_{ji} = \begin{cases} x_{ji} & \text{if } x_{ji} \geq 0 \\ \alpha_{ji} x_{ji} & \text{if } x_{ji} < 0 \end{cases}$$

## Linear Separability

\* 1D case :

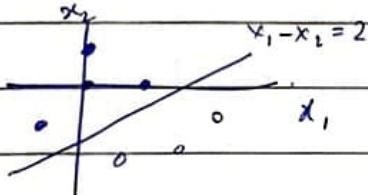


Separable by perceptron



Linearly inseparable

\* 2D case :

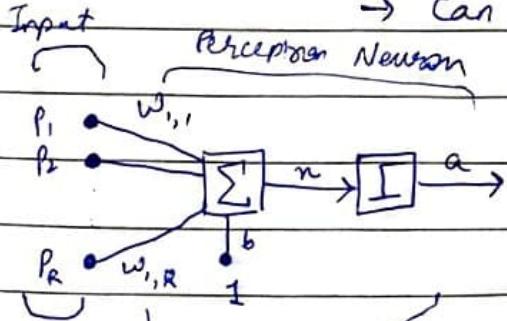


Separable by perceptron

Perception model

→ Machine which can learn to assign input vectors to different classes.

→ Can do 2 class linear classification problem



$$a = \text{hardlim}(w_p + b) \quad \text{hardlim}(n) = 1, \text{ if } n \geq 0; 0 \text{ otherwise}$$



$$W_{\text{new}} = W_{\text{old}} + \eta e p$$

$$b_{\text{new}} = b_{\text{old}} + \eta e, \text{ where } e = \text{target} - \text{actual}$$

- Q. 7 one dimensional input patterns (0.0, 0.17, 0.33, 0.5, 0.67, 0.83, 1.03). Assume that first four patterns belong to class 0 & remaining patterns belong to class 1. Design a perceptron to classify these patterns. Use perceptron learning rule. Assume  $\eta = 0.1$  and initial weight & bias to be (-0.36) & (-0.1) respectively. Show computation for two epochs.

\* Solve yr problem

$y$	$x$	$\hat{y}$	$y - \hat{y}$
15.5	630	10.75	-4.75
7.5	370	4.25	-3.25
13.7	616	10.4	-3.5
18.7	700	12.5	-6.2
8.2	430	5.75	-2.45
13.2	568	9.2	-4
23	1200	25	2
87.3	2936	69.04	<u>-18.29</u>

Initially,  $\beta_0 = -5$ ,  $\beta_1 = 0.025$

$$\alpha = 0.1$$

$$-(65-11 + 1-1)$$

$$\therefore \beta_0 = \beta_0 - \alpha \sum_{m=1}^8 (y - \hat{y})$$

$$\beta_0 = -5 - \frac{(0.1)(-18.29)}{8}$$

$$= \boxed{-5.500}$$

	$x$	$y$	class
I	-1	1	-
II	0	1	+
III	0	2	-
IV	1	-1	-
V	1	0	+
VI	1	2	+
VII	2	2	-
VIII	2	3	+

To be classified  $\rightarrow (1, 1)$

$$\begin{aligned}
 \text{dist}(i, (1, 1)) &\rightarrow \sqrt{(1 - (-1))^2 + (1 - 1)^2} & 2 \\
 \text{dist}(ii, " ) &\rightarrow \sqrt{(1 - 0)^2 + (1 - 1)^2} & 1 \\
 " &\rightarrow \sqrt{(1 - 0)^2 + (1 - 2)^2} & \sqrt{2} \\
 " &\rightarrow \sqrt{(1 - 1)^2 + (1 - 0)^2} & 2 \\
 " &\rightarrow \sqrt{(1 - 1)^2 + (1 - 0)^2} & 1 \\
 " &\rightarrow \sqrt{(1 - 1)^2 + (1 - 2)^2} & 1 \\
 " &\rightarrow \sqrt{(1 - 2)^2 + (1 - 2)^2} & \sqrt{2} \\
 " &\rightarrow \sqrt{(1 - 2)^2 + (1 - 3)^2} & \sqrt{5}
 \end{aligned}$$

Arranging in ascending order

Dist      class

$\sqrt{5}$       +

~~2~~      -

~~2~~      -

$\sqrt{2}$       -

$\sqrt{2}$       -

1      +

$\frac{1}{2}$       +

1      +

$\rightarrow$  By 3-NN, ans is  $(-)$

$\rightarrow$  By 5-NN, ans is  $(-)$