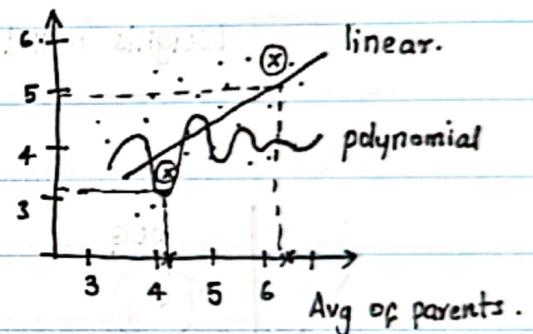


Suppose, we were given data of (average of heights of parents) and (Avg of heights of children)

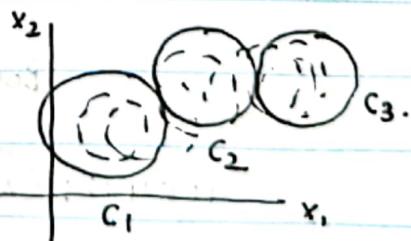
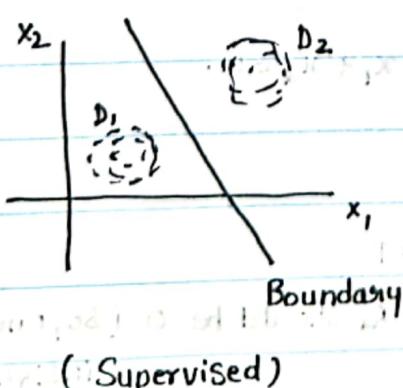
We can establish a relation between Avg of children. a datasets. it can be linear / polynomial.

And in both cases, it can be error-prone.



To avoid errors, we need more data but processing more data can be difficult.

→ In Supervised learning, in order to classify the data we need a boundary that separates the datasets. whereas, in unsupervised learning we form clusters to determine / classify.

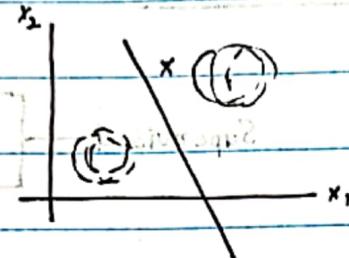


$C_1, C_2, C_3$  : clusters  
(unsupervised).

• Supervised Learning

provided the data with labels, we need a boundary to classify the input data.

• Decision Boundary (A line with points belonging to different classes)



Left side of boundary

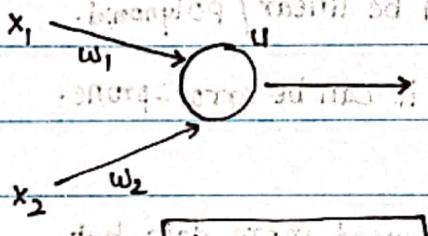
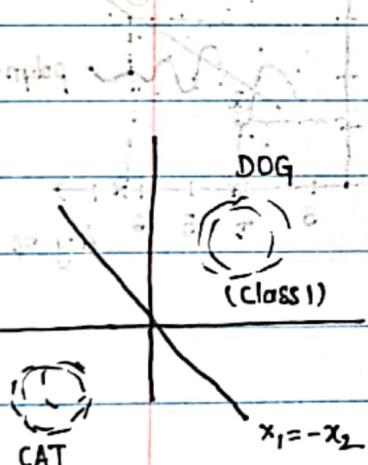
input lies left Side of Boundary

↳ belongs to dataset 1

Boundary Right Side of boundary

↳ belongs to D2.

Weights will be assigned to each dataset ( $w_1, w_2$ )



$$U = x_1 w_1 + x_2 w_2$$

if  $U \geq 0 : y = 1$       When  $U = 0$ ; we get  
else  $y = 0$ .      a boundary line.

(class 0)      (class 1)      (input lies left Side of boundary)

Let's consider the case of classification  $x_1 w_1 + x_2 w_2 = 0$ .

(if  $w_1, w_2 \neq 0$ )  $\Rightarrow x_1 = -x_2$

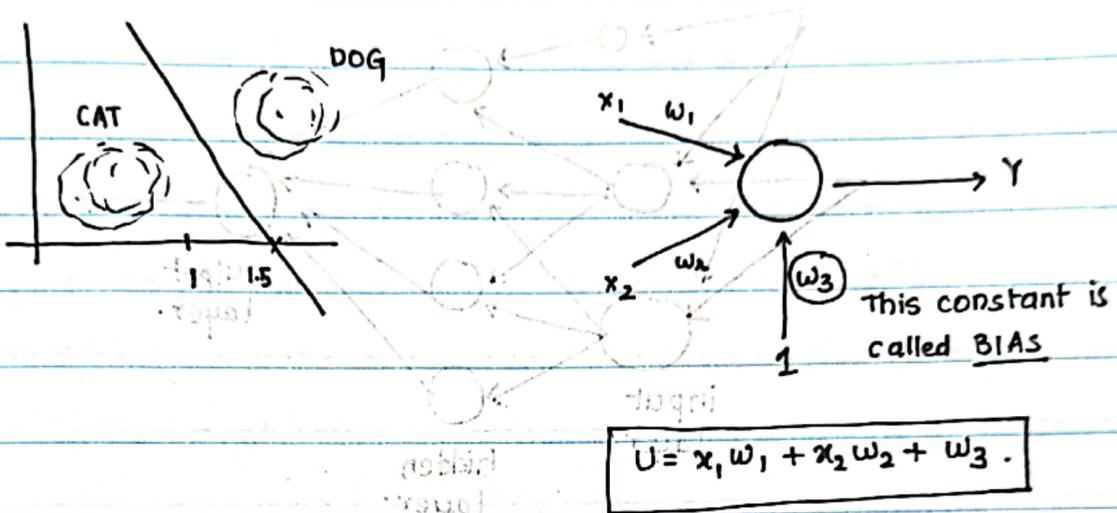
Boundary line:  $x_1 + x_2 = 0$

For  $y=1$ ;  $x_1 + x_2 = 1$

$$x_1 w_1 + x_2 w_2 = 1$$

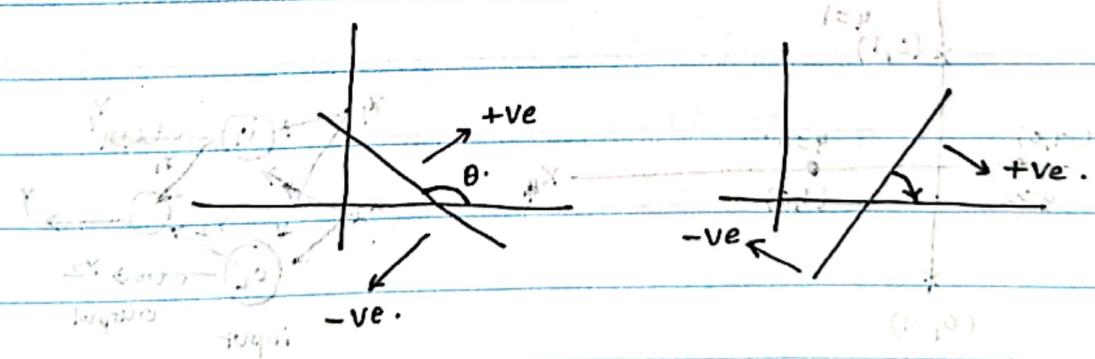
when  $w_2 = 1$ ;  $x_1$  should be 0 (So, one dataset will be insignificant).

When the boundary isn't passing through the centre, we need an additional constant to differentiate with the boundary.



→ (For the line, slope always +ve)

Slope will be calculated w.r.t x-axis; so the values towards x-axis angle is +ve.

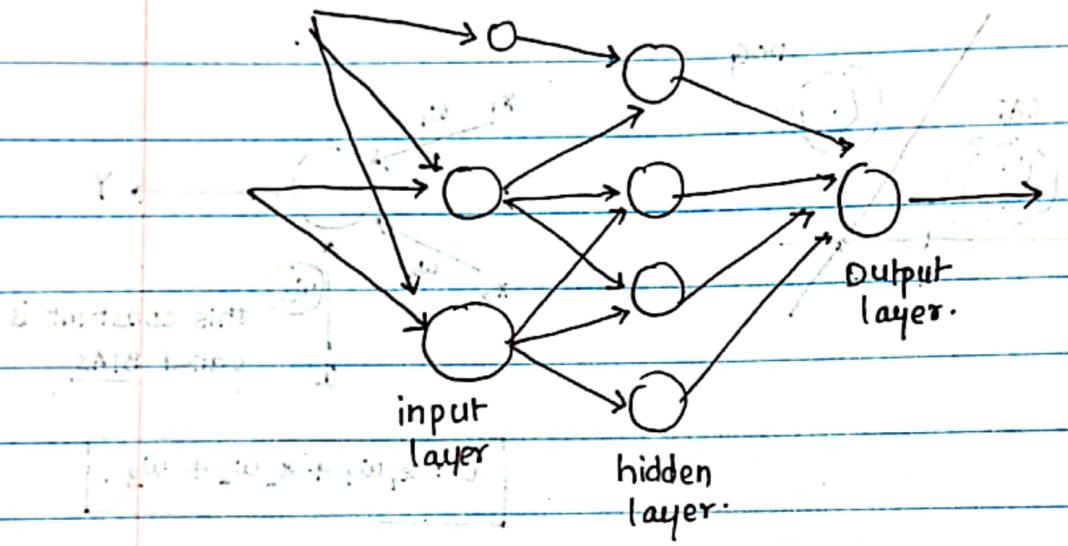


(angle made up)

Input and output variables are?

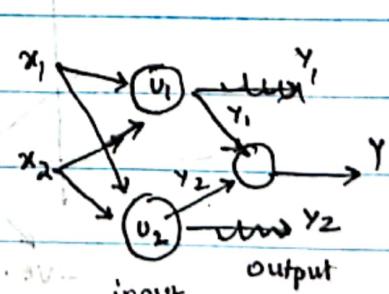
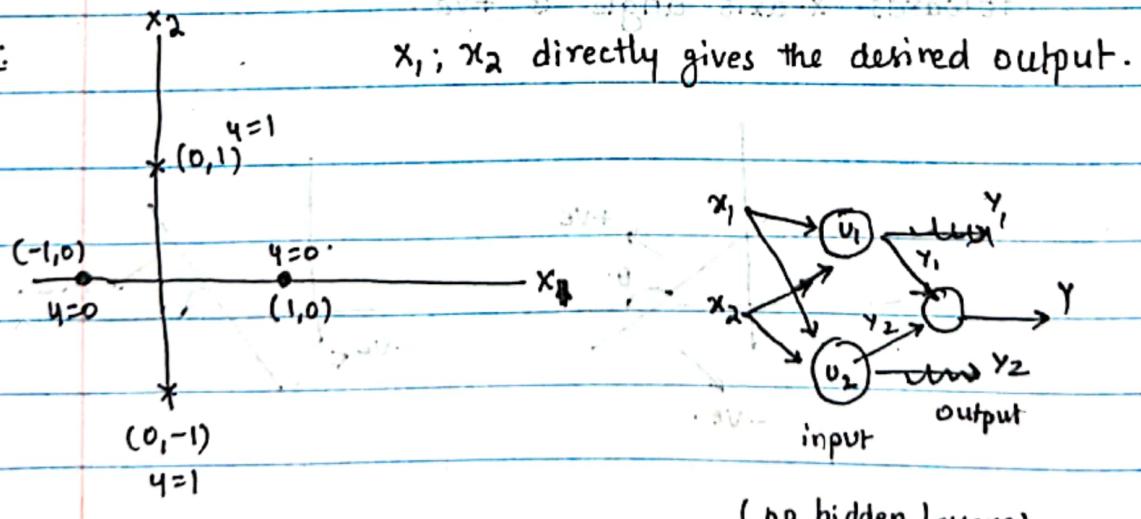
The layers which we don't have access to is called hidden layer.

→ That are often difficult to find without training the model.



When the result can be given without any lines (neurons), i.e. there are no hidden layers we don't need neurons (lines).

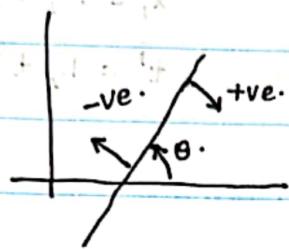
Ex:



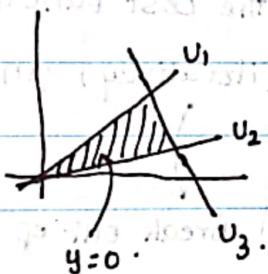
(no hidden layers)

→ So, no neurons in hidden layer

The angular region (region between x-axis and our reference line) will be positive and remaining part will be negative.



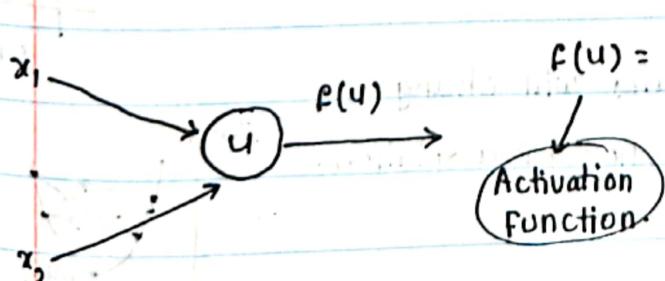
When finding common region to distinguish datasets ( $y=0$  or  $y=1$ ), consider all the neurons/lines making up region into same sign. So that it will be easier to distinguish/get line equations.



towards inside (for region) consider -ve for all neurons/lines.

If it's positive, multiply the line equation with -1 to change sign.

Gradient Descent Algorithm: In order to get the minimum errors, we calculate cost function  $J(w)$  [when there's only one variable] and randomly assume values for  $w$  and calculate cost value. And, for value which cost is minimum, that ' $w$ ' value will be considered.



$$e(u) = \begin{cases} 1 & \text{if } u \geq 0 \\ 0 & \text{if } u < 0 \end{cases}$$

$m \Rightarrow \# \text{ of samples}$

$x_i^i \rightarrow i^{\text{th}} \text{ sample}$

$t^i \rightarrow \text{target/output}$

$x_i \rightarrow i^{\text{th}} \text{ attribute}$

$n \rightarrow \text{Dimension}$ .

Example:

$x_1$	$x_2$	target( $t$ )	$m = 2$
0	0	1	$x^1 = (1, 0)$
0	1	0	$x^2 = (0, 1)$

$x_1$   $\rightarrow$   $x_2$

Error / Cost Function / Loss Function =  $J(w) = \frac{1}{2m} \sum_{i=1}^m (f(u^i) - t^i)^2$

Gradient Descent: Based on random selection of  $w_1, w_2$   
we calculate the cost function. And differentiate  
the cost function (iteratively) until you get the minimum  
cost / error value.

$$w_1 = w_1 - \rho \Delta w_1$$

$$w_2 = w_2 - \rho \Delta w_2$$

where  $\rho \rightarrow$  learning rule / learning rate.  
learning rate.

Break out of the loop when

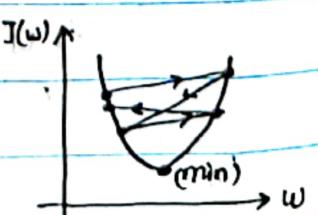
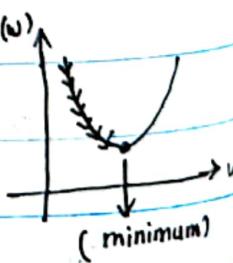
i) Error is very minimal. (like  $\epsilon < 0.0001$ )

(and)

ii) There's no change in  $\Delta w$ .  
(change in  $w$ ).

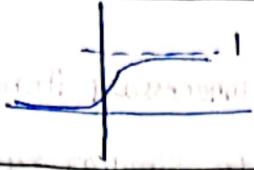
→ if  $\rho$  is small, the steps we take to reach the solution  
will be very small. (Takes more time).

→ if  $\rho$  is high,  $w_1/w_2$  values will change  
drastically and we might not get a solution.



Sigmoid  
function

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



deciding  
function we use to  
determine y value.

→ in case of non-linear classification, we take the data to the higher dimension in order to classify the data.

Example:  $x_3 = x_1^2 + x_2^2$

(con)  $ax^2 + by^2 + c^2 = 0 \Rightarrow$  change  
this to  $ax_1 + by_1 + c^2 = 0$

where  $x_1 = x^2$ ;  $y_1 = y^2$

Using Chain rule,  $\frac{\partial}{\partial w_j} (J(w)) = \frac{\partial J}{\partial f} \cdot \frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial w_j}$

$$\frac{\partial J}{\partial f} = \frac{\partial}{\partial f} \left( \frac{1}{m} \sum_{i=1}^m (f(u^i) - t^i)^2 \right) = 2(f(u) - t)$$

$$\frac{\partial f}{\partial u} = \frac{\partial}{\partial u} ((1+e^{-u})^{-1}) = -1(e^{-u}) \cdot (1+e^{-u})^{-2} = \frac{e^{-u}}{1+e^{-u}} \cdot \frac{1}{1+e^{-u}}$$

$$= 1 - \frac{1}{1+e^{-u}} \cdot \frac{1}{1+e^{-u}}$$

$$\frac{\partial f}{\partial u} = (1-f(u)) \cdot f(u)$$

$$\frac{\partial u}{\partial w_j} = \frac{\partial}{\partial w_j} \left( \sum_{i=1}^n (w_0 x_{0i} + w_1 x_{1i} + \dots + w_n x_{ni}) \right) = x_j$$

$$\Rightarrow \Delta w = \frac{\partial (J(w))}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m [(f(u^i) - t^i) \cdot f(u^i) \cdot (1-f(u^i)) \cdot x_j]$$

$$w_j = w_j - \alpha \Delta w$$

$\alpha$ : Learning rate.

Scaling: In order to avoid the unnecessary iterations / to speed up the procedure we scale our values to similar ranges.

↳ we scale our values to a standard unit

Example:

$x_i$	Deviation	Squared	standard unit
1	-2.5	6.25	$-2.5/1.7 = -1.47$
2	-1.5	2.25	$-1.5/1.7 = -0.88$
3	-0.5	0.25	$-0.5/1.7 = -0.29$
4	0.5	0.25	$0.5/1.7 = 0.29$
5	1.5	2.25	$1.5/1.7 = 0.88$
6	2.5	6.25	$2.5/1.7 = 1.47$
<hr/>		Avg: $\frac{17.5}{6} = 2.9$	Variance
Average: $\frac{21}{6} = 3.5$		Standard deviation = $\sqrt{2.9} \approx 1.7$	

→ we can't consider deviation as standard unit, as in most cases summation will be zero which doesn't help us.

↳ we use,  $\frac{x_i - \mu_i}{sd_i}$  =  $x_i - \text{Average}$   
standard deviation as our standard unit.  
to scale our data.

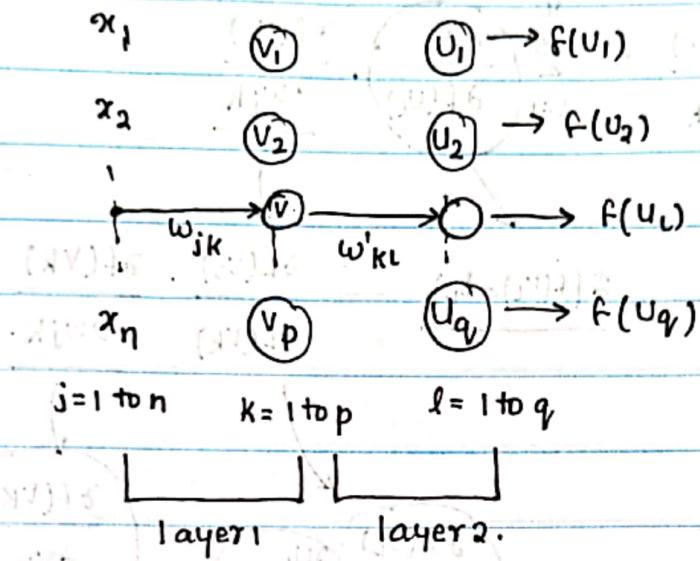
$x_i$ : input values for feature(i)

$\mu_i$ : Average of feature (i)

$sd$  = standard deviation of feature (i)

↳ This is called MEAN NORMALIZATION method.

## Training Two Layer Network:



m : number of samples.

$$\text{cost function } J(w) = \sum_{i=1}^m \sum_{l=1}^q (f(u_l) - t_l)^2$$

$$\Delta w'_{kl} = \frac{\partial J(w)}{\partial w'_{kl}} = \frac{\partial J(w)}{\partial F(u_l)} \cdot \frac{\partial F(u_l)}{\partial u_l} \cdot \frac{\partial u_l}{\partial w'_{kl}}$$

$$u_l = \sum_{k=1}^p f(v_k) \cdot w'_{kl}$$

$$2(f(u_l) - t_l) \cdot (f(u_l) \cdot (1 - f(u_l)) \cdot f(v_k))$$

$$\Delta w'_{kl} = 2(f(u_l) - t_l) \cdot (f(u_l) \cdot (1 - f(u_l)) \cdot f(v_k))$$

Layer 02

$$\boxed{\Delta w'_{kl} = \sum_{i=1}^m \sum_{l=1}^q (f(u_l^i) - t_l^i) (f(u_l^i) \cdot (1 - f(u_l^i)) \cdot f(v_k))}$$

For Layer 01:

$$\Delta w_{jk} = \frac{\partial J(w)}{\partial w_{jk}} = \sum_{l=1}^q \left( \frac{\partial J}{\partial f(u_l)} \cdot \frac{\partial f(u_l)}{\partial w_{jk}} \right)$$

$\frac{\partial (f(u_l) - t_l)}{\partial w_{jk}} \cdot \frac{\partial f(u_l)}{\partial w_{jk}} \cdot \frac{\partial f(v_k)}{\partial v_k}$

$\frac{\partial f(u_l)}{\partial u_l} \cdot \frac{\partial u_l}{\partial f(v_k)} \cdot \frac{\partial f(v_k)}{\partial v_k} \cdot \frac{\partial v_k}{\partial w_{jk}}$

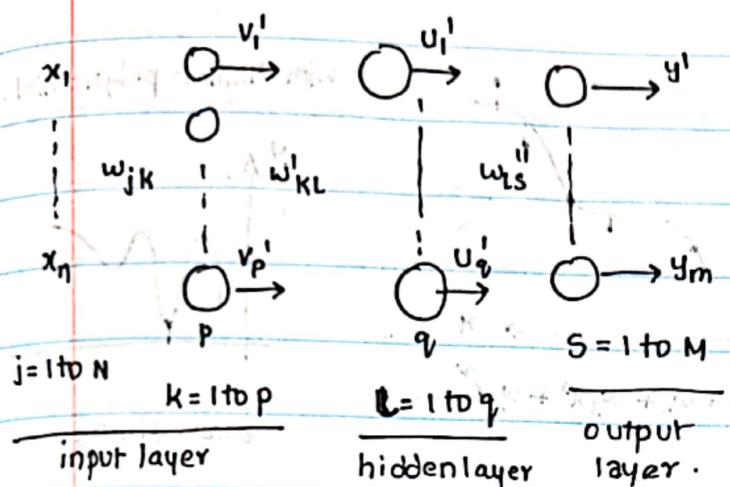
$\frac{(f(u_l) \cdot (1-f(u_l)))}{w_{kl}} \cdot (f(v_k) \cdot (1-f(v_k)))$

$$\boxed{\Delta w_{jk} = \alpha x_j (f(v_k) \cdot (1-f(v_k))) \sum_{l=1}^q (f(u_l) - t_l) \cdot f(u_l) (1-f(u_l)) \cdot w_{kl}}$$

→ So, we have to calculate the  $\Delta w$  for layer 2 and then we have to calculate for the layers until convergence.

→ we call this **Back propagation Training Algorithm.**

Training 3 layers:  $\rightarrow$  Back propagation.



$$\Delta w_{ls}''' = \rho u_l' y_s (1-y_s) (y_s - t_s)$$

$$\text{let } y_s (1-y_s) (y_s - t_s) = \delta_s$$

$$\Delta w_{kl}' = \rho v_k' u_l' (1-u_l') \sum_{s=1}^M w_{ls}''' y_s (1-y_s) (y_s - t_s)$$

$$\Delta w_{kl}' = \rho v_k' u_l' (1-u_l') \sum_{s=1}^M w_{ls}''' \delta_s$$

$$\text{let } u_l' (1-u_l') \sum_{s=1}^M w_{ls}''' y_s (1-y_s) (y_s - t_s) = \delta_l'$$

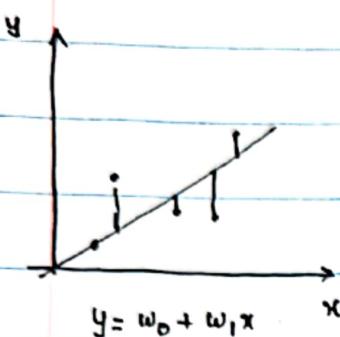
$$\Delta w_{jk} = \rho x_j v_k' (1-v_k') \sum_{l=1}^q w_{kl}' \delta_l'$$

$$\text{let } v_k' (1-v_k') \sum_{l=1}^q w_{kl}' \delta_l' = \delta_k''$$

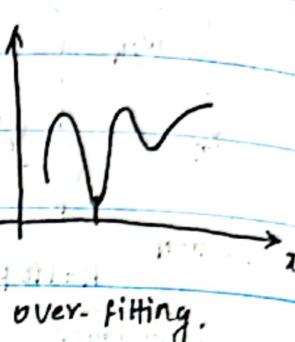
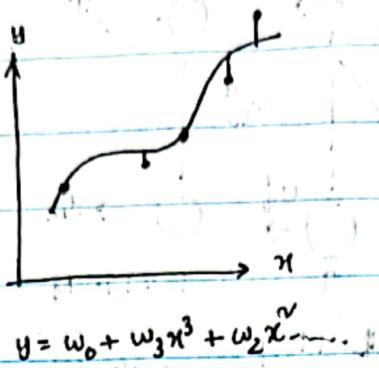
$$\Delta w_{ij} = \rho x_i \delta_k''$$

→ Overfitting may minimise the error data and may give error response for the new data.

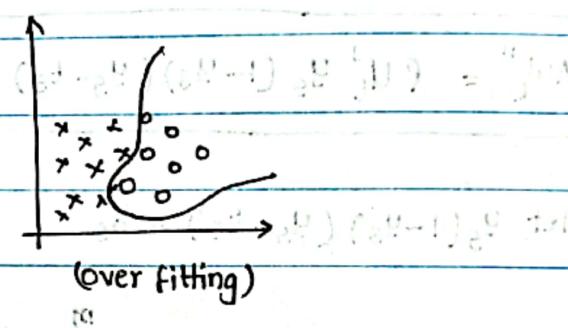
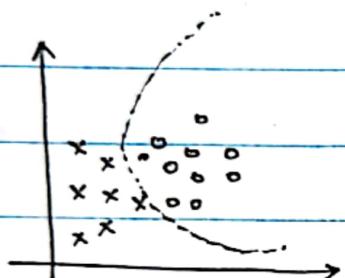
Linear Regression.



with higher polynomial.



Logistic regression



Regression Line:

$$\text{Standard unit} = \frac{x - \text{Avg of } x}{\text{standard deviation}} = \left[ \frac{x - \bar{x}}{s} \right]$$

$$Y_{\text{est}} = r \cdot X_{\text{given}}$$

in Original scale:

$$Y_{\text{estimate}} = r \cdot \frac{SD_y}{SD_x} \cdot X_{\text{given}} + Y_{\text{average}} - r \cdot \frac{SD_y}{SD_x} \cdot X_{\text{average}}$$

Slope of  
the regression  
line.

r: Correlation

Intercept of  
the regression line

## Convolution Network:

Consider a  $5 \times 5$  image, with a pattern (let's say  $X$ )

-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1

consider 2 Filters.

1	-1
-1	1

-1	1
1	-1

let dark pixels

value = 1

and other pixels = -1

use each Filter, traverse through the image (by each column)  
and do the elementary multiplication and update the values.

→ do this for both filters.

With Filter 01,

1	-0.5	0.5	0
-0.5	1	0	0.5
0.5	0	1	0
0	0.5	0	1

use max-pooling

(Take max value  
in each  $2 \times 2$   
block)

1	0.5
0.5	1

FLAT

1	0.5
0.5	1
1	0.5
0.5	1
1	0.5

With Filter 02,

-1	0.5	-0.5	1
0.5	-1	1	-0.5
-0.5	1	-1	0.5
1	-0.5	0.5	-1

max-pooling.

0.5	1
1	0.5

(max pooled pixels to Flat  
feature vectors).

→ In training perceptron models, with images, pixels acts as features and filters acts as weights.

④ To adjust the error response caused by over fitting of trained data, we would make higher polynomial weights to zero, so that few attributes are minimized removed.

Example:

$$w_0x_0 + w_1x_1 + w_2x_2 + w_3x_3^2 + w_4x_4^4$$

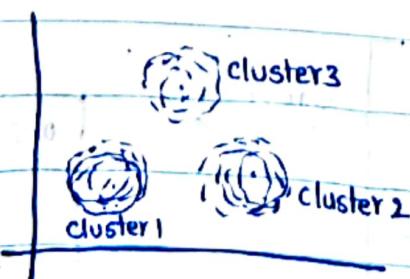
(minimize the weights)

Text Classification: To identify / train a model to identify a statement of string to be True / False, we convert the words to vector.

Example: I like reading book in night.  
↓                    ↓  
high                centre                high

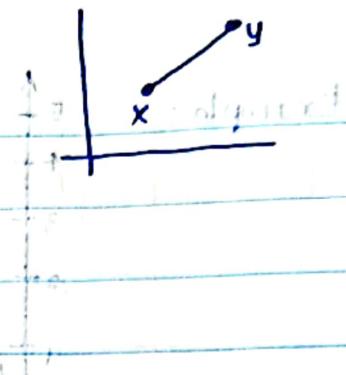
→ We focus mainly on middle word and its neighbouring words to be converted to feature vector.

Hierarchical clustering: To group / cluster any data we measure the distance between data points to identify the cluster they belong to.



## Distance Measurement:

Euclidean Distance:  $\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$

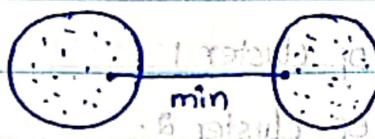


Squared Euclidean Distance:  $\sum_{i=1}^n (x_i - y_i)^2$

Manhattan Distance:  $\sum_{i=1}^n |x_i - y_i|$

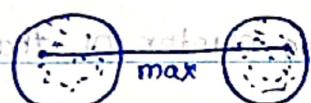
## Measure Similarity between clusters:

### Single Linkage:



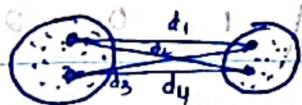
minimum / closest distance between points from 2 clusters.

### Complete Linkage:



maximum distance between points from 2 clusters.

### Group Linkage:



= average of the points

$$\frac{d_1 + d_2 + d_3 + d_4}{4}$$

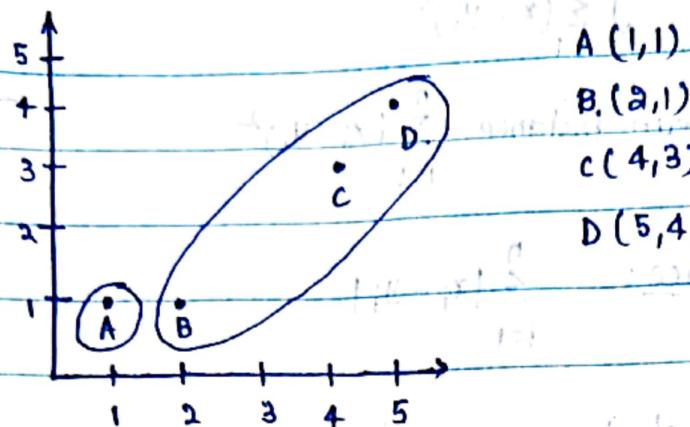
### Average Group Linkage:



Distance between centre data points of the clusters.

$$(c_1, c_2) = (E(S), E(B))$$

Example:



A(1,1)

B(2,1)

C(4,3)

D(5,4)

clustering of given data points A, B, C, D.

lets consider number of clusters = 2

$c_1 = (1,1)$  centre of cluster 1

$c_2 = (2,1)$  centre of cluster 2.

if distance of any point is minimum to any centre then that data point belongs to the corresponding cluster of that centre.

Iteration 0 :

$$D^0 = c_1 \begin{pmatrix} 0 & 1 & 3.61 & 5 \\ 1 & 0 & 2.83 & 4.24 \end{pmatrix}$$

clustering

$$\begin{matrix} c_1 & \begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix} \\ c_2 & \begin{bmatrix} 0 & 1 & 1 & 1 \end{bmatrix} \end{matrix}$$

Find new centres after clustering;

New values of  $c_1$  &  $c_2$ ;

$$c_1 = (1,1)$$

$$c_2 = \left( \frac{2+4+5}{3}, \frac{1+3+4}{3} \right) = \left( \frac{11}{3}, \frac{8}{3} \right)$$

Iteration 01:

$$D_1 = \begin{pmatrix} A & B & C & D \\ 0 & 1 & 3.61 & 5 \\ C_1 & & & \\ C_2 & 3.14 & 2.36 & 0.41 & 1.89 \end{pmatrix}$$

Clustering:

$$\begin{matrix} C_1 & \left[ \begin{array}{cccc} A & B & C & D \\ 1 & 1 & 0 & 0 \end{array} \right] \\ C_2 & \left[ \begin{array}{cccc} 0 & 0 & 1 & 1 \end{array} \right] \end{matrix}$$

New values for  $C_1$  &  $C_2$ :

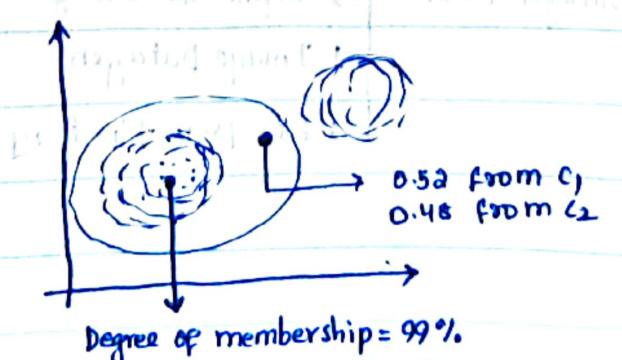
$$C_1 = \left( \frac{1+2}{2}, \frac{1+1}{2} \right) = (1.5, 1)$$

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

Repeat the iterations until  $C_1$  and  $C_2$ , become consistent / doesn't change.

k-Means clustering: Here any given data point belongs to particular cluster (100%) as distance would be min to centre of that cluster than others.

Fuzzy clustering: if any given data point is at equal distances from centres of different clusters then any one of the cluster would be pull the data point based on degree of membership.



## CATS & DOGS

Define Dataset  $\leftrightarrow$  images, classes (keras Image DataGenerator)

model sequential

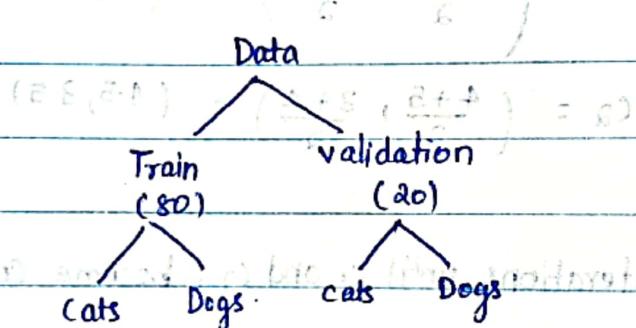
model compile

model fit (Training data, Testing Data)

history

model predict.

Divide the dataset (into folders) [80-20]



→ image dataset can be diverse i.e. differ in background / angle of capture / transformation (Rotation, Zoom etc...) / Effects (B/w, etc...) / ...

split the data into different folders. [path/train/cats

path/train/Dogs]

Keras Image Data generator → Define the changes

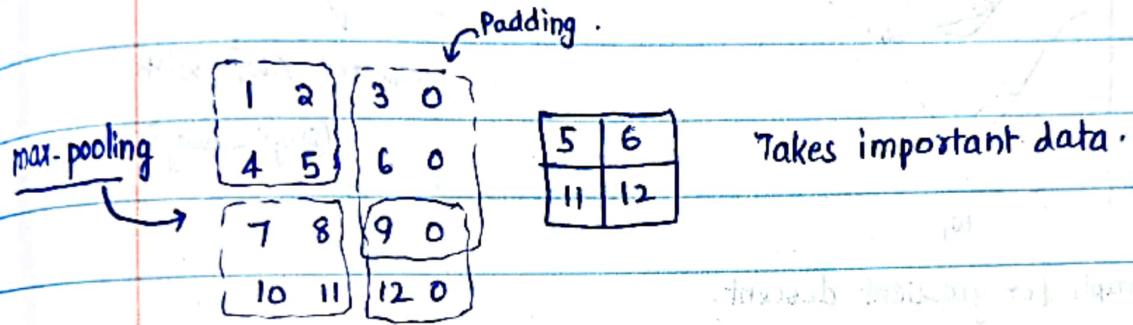
→ Image DataGen

→ Flow from Directory

Training / model: Sequential API for building a model.

Sequential:

Conv2D layer : Conv2D (num of filters,  
filter size,  
activation function, input size)



Flatten: Converts input to a 1D-array / Vector.

Dense: Neuron Layer.

→ Then compile the model with relevant optimizer.

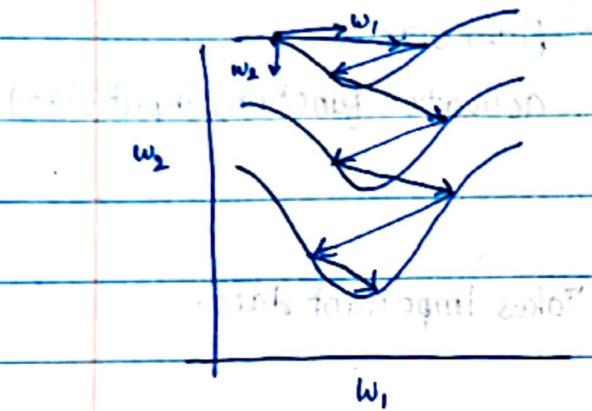
and then train the data [model.fit()]

→ if we have multiple convolution layers, maxpooling layer output of previous layer will be the input of next convolution layer.

## K-Nearest Neighbourhood (Algorithm)

Returns the closest result to the sample.

## Improve Gradient Descent: (using optimization Techniques).



$$w_j^{t+1} = w_j^t - \frac{\Delta w_j^t}{\text{gradient}}$$

replace  $\Delta w_j^t$  with

$$(w_j^t - \Delta w_j^{t-1})$$

contour graph for gradient descent.

Momentum:

$$w^{t+1} = v^t + w^t$$

$\beta$ : momentum co-efficient

$$v^t = \beta v^{t-1} - \alpha \left( \frac{\partial E}{\partial w} \right)_t G^t$$

$$[ \text{for } \beta = 0.9, \text{ then } v^3 = 3 [ 3 (-2G') - \alpha G^2 ] - 2G^3 ]$$

thus we repeat previous step until convergence of all the weights.

repeat until both of the weights are not moving.

Root Mean Square propagation (RMSprop):

$$w^{t+1} = w^t - \frac{\beta}{\sqrt{v^t + \epsilon}} \cdot G^t.$$

$\epsilon$  is very small number which we include to avoid denominator being 0.

Ex:  $\epsilon \approx 0.00001$

$$v^t = \rho v^{t-1} + (1-\rho)(G^t)^2.$$

For  $v^0 = \phi$  (lets say).

$$v^1 = (1-\rho)(G^1)^2$$

$$v^2 = \rho(1-\rho)(G^1)^2 + (1-\rho)(G^2)^2$$

$$v^3 = \rho [\rho(1-\rho)(G^1)^2 + (1-\rho)(G^2)^2] + (1-\rho)(G^3)^2$$

Depending on each slope (during descent), for each weights change will be dynamic.

### Adaptive Moment Optimization [ADAM]

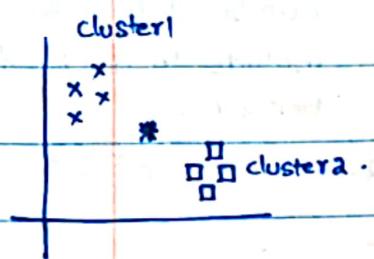
ADAM combines both momentum and RMSprop.

$$w^{t+1} = w^t - \beta \cdot \frac{s^t}{\sqrt{v^t + \epsilon}} \cdot G^t.$$

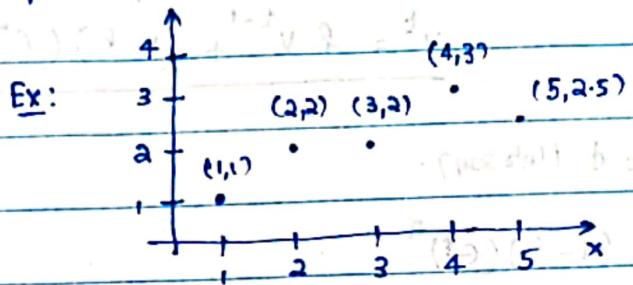
where  $v^t = \rho v^{t-1} + (1-\rho)(G^t)^2.$

$$s^t = \rho s^{t-1} - (1-\rho)G^t.$$

Principal Component Analysis (PCA): Technique used to project higher dimensional data to lower dimension.

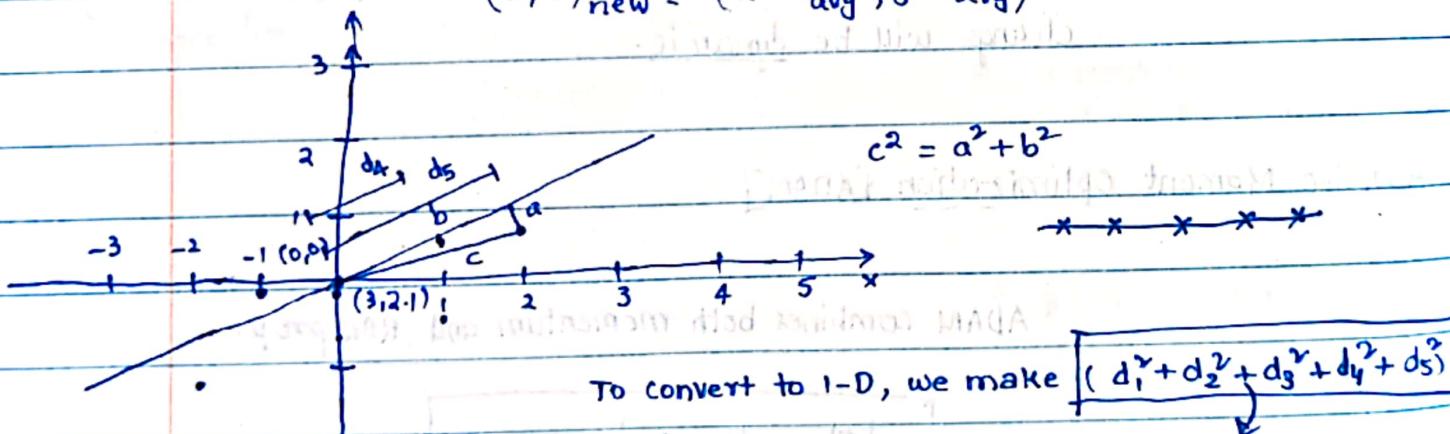


Consider the average point, and shift origin to that point.



shift origin to  $(x_{avg}, y_{avg}) = (3, 2.1)$

$$(x, y)_{new} = (x - x_{avg}, y - y_{avg})$$

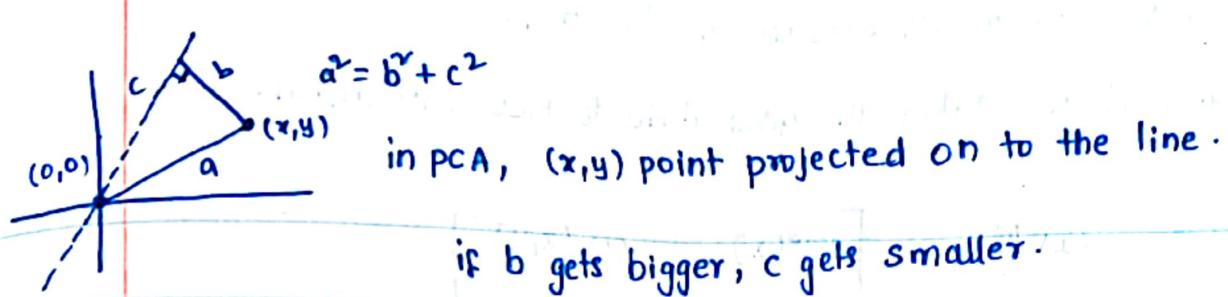


To convert to 1-D, we make  $(d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2)$  maximum.

(SS)

We draw random line, through Origin and find if the line fits the data or not.

$\{d_1, d_2, \dots, d_5\}$  is the distance between origin to the points projected on to the random line.



in PCA,  $(x, y)$  point projected on to the line.

if  $b$  gets bigger,  $c$  gets smaller.

(or) if  $b$  gets smaller,  $c$  gets bigger.

So, PCA can minimize the distance to the line (or) maximize the distance of the projected point to the origin.

Consider the line, which has maximum value of ss.

we call this line principal component 1 (PC1)

and PC2 will be  $1^{\text{st}}$  to PC1.

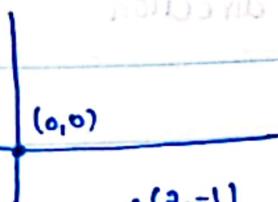
calculate Co-variance Matrix based on PC1 & PC2.

Eg:

$\begin{array}{c} \times \\ (-1,0) \quad (0,0) \quad (1,0) \end{array}$

$$\text{variance: } \frac{(-1)^2 + (0)^2 + (1)^2}{3} = 2/3.$$

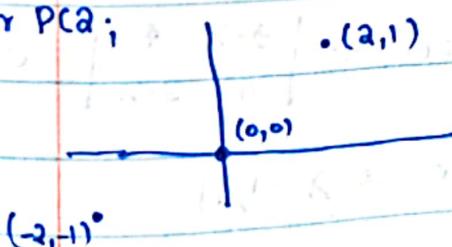
①  $(-2, 1)$



$$x_{\text{var}} = \frac{(-2)^2 + (0)^2 + (1)^2}{3} = 8/3$$

$$y_{\text{var}} = \frac{(1)^2 + (0)^2 + (-1)^2}{3} = 2/3.$$

For PC2:



$$x_{\text{var}} = 8/3$$

$$y_{\text{var}} = 2/3$$

variances are same for both  
PC1 & PC2.

Co-variance: Take average of product of  $x$  &  $y$  values.

$$\textcircled{1} \quad \frac{(-2 \times 1) + (0 \times 0) + (1 \times -1)}{3} = -4/3$$

$$\textcircled{2} \quad \frac{(2 \times 1) + (0 \times 0) + (-2 \times -1)}{3} = 4/3$$

Co-variance matrix gives us the direction.

Matrix transformation: Takes the given data to new dimension.

$$\text{Cov Mat} = \begin{bmatrix} \text{Var}(x) & \text{cov}(x,y) \\ \text{cov}(x,y) & \text{Var}(y) \end{bmatrix}$$

ex: CovMat =  $\begin{bmatrix} 9 & 4 \\ 4 & 3 \end{bmatrix}$

Transformation :  $(x \ y) \begin{bmatrix} 9 & 4 \\ 4 & 3 \end{bmatrix} = (9x+4y, 4x+3y)$

$$(1,0) \Rightarrow (9,4)$$

$$(0,1) \Rightarrow (4,3)$$

$$(2,1) \Rightarrow (22,11) : (2,1) \text{ multiplied by } 11 \rightarrow \text{Eigen value.}$$

(magnified by 11-times)

Eigen vectors: Vectors that doesn't change direction.

$$\text{Eigen vector } (\mathbf{v}) : A\mathbf{v} = \lambda \mathbf{v}$$

$\lambda$ : Eigen value.

$$(A - \lambda I)\mathbf{v} = 0$$

$$\begin{pmatrix} 9-\lambda & 4 \\ 4 & 3-\lambda \end{pmatrix} \mathbf{v} = 0 \Rightarrow \begin{vmatrix} 9-\lambda & 4 \\ 4 & 3-\lambda \end{vmatrix} = 0$$

$$\Rightarrow (9-\lambda)(3-\lambda) - 16 = 0 \Rightarrow \lambda = 1, 11$$

Solve the equations for  $\mathbf{v}$ .

$$\begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = 11 \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

$$\begin{pmatrix} 9v_1 + 4v_2 \\ 4v_1 + 3v_2 \end{pmatrix} = \begin{pmatrix} 11v_1 \\ 11v_2 \end{pmatrix}$$

$$\begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = 1 \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

$$\begin{pmatrix} 9v_1 + 4v_2 \\ 4v_1 + 3v_2 \end{pmatrix} = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

For  $\lambda = 11, \mathbf{v} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$

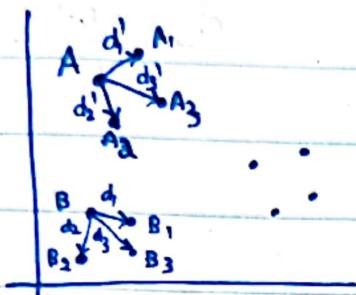
For  $\lambda = 1, \mathbf{v} = \begin{pmatrix} 1 \\ -2 \end{pmatrix}$

## T-SNE (T-Distribution stochastic Neighbour Embedding)

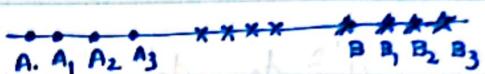
→ Technique to project higher dimension data to lower dimension.

Ex: projecting 2D data into 1D.

Actual Data.



Target projection



Stochastic distribution of the given data.

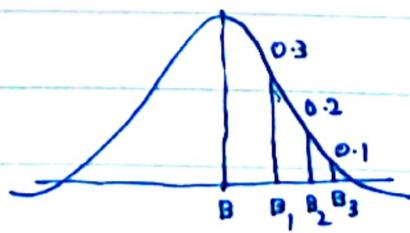
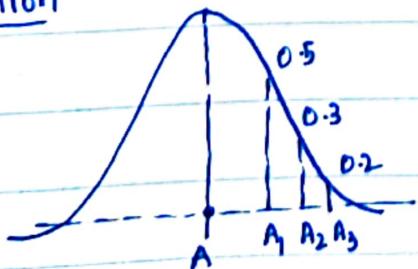
(random distribution)



1D data of stochastic distribution tries to change as points from same group moves closer and points from different group move away from each other. This process happens millions of times to attain the 'Target projection'.

To move the points around within the cluster in specific order we consider similarity among points to be same between groups i.e., similarity from A to  $A_1$  should be same as B to  $B_1$ .

## t-Distribution



Simillarity

$$d_1' = \frac{0.5}{0.5+0.3+0.2} = 0.5$$

$$d_1 = \frac{0.3}{0.3+0.2+0.1} = 0.5$$

$$d_2' = \frac{0.3}{1} = 0.3$$

$$d_2 = \frac{0.2}{0.6} = 0.3$$

$$d_3' = \frac{0.1}{1} = 0.1$$

$$d_3 = \frac{0.1}{0.6} = 0.16$$

$$d_1' \approx d_1, d_2' \approx d_2, d_3' \approx d_3.$$

$\therefore$  points within different groups move around with same force because of having equal simillarities.

Regularization: To overcome the problem of overfitting, we use Regularization.

- i. More Data
- ii. Data Augmentation
- iii. Early stop.

m : # of samples

n: # of features

$$h_w(x) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$\text{Error/Cost } J(w) = \frac{1}{2m} \sum_{i=1}^m (h_w(x^i) - t^i)^2$$

in order to minimize the weights of few features (to zero), we introduce regularization factor ( $\lambda$ ).

updated cost function

$$J(w) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_w(x^i) - t^i)^2 + \lambda \sum_{j=1}^n w_j^2 \right]$$

$$w_j = w_j - \frac{\alpha}{m} \left[ \sum_{i=1}^m (h_w(x^i) - t^i) x_j^i + 2\lambda w_j \right]$$

$$w_j = \underbrace{\left(1 - \frac{2\alpha\lambda}{m}\right)}_{\text{also called as decay-rate.}} w_j - \frac{\alpha}{m} \left[ \sum_{i=1}^m (h_w(x^i) - t^i) x_j^i \right]$$

# in this approach, instead of taking  $w_j$ , we take fraction of  $w_j$  which helps us to minimize  $w_j$  faster.