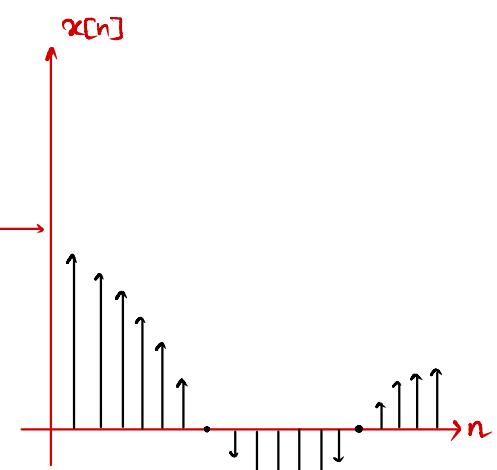
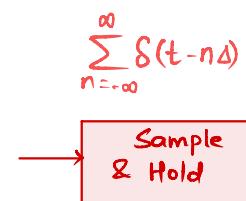


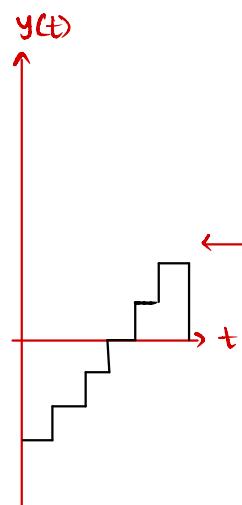
Continuous in time  $\equiv t \in \mathbb{R}$

Continuous in magnitude  $\equiv x \in \mathbb{R}$



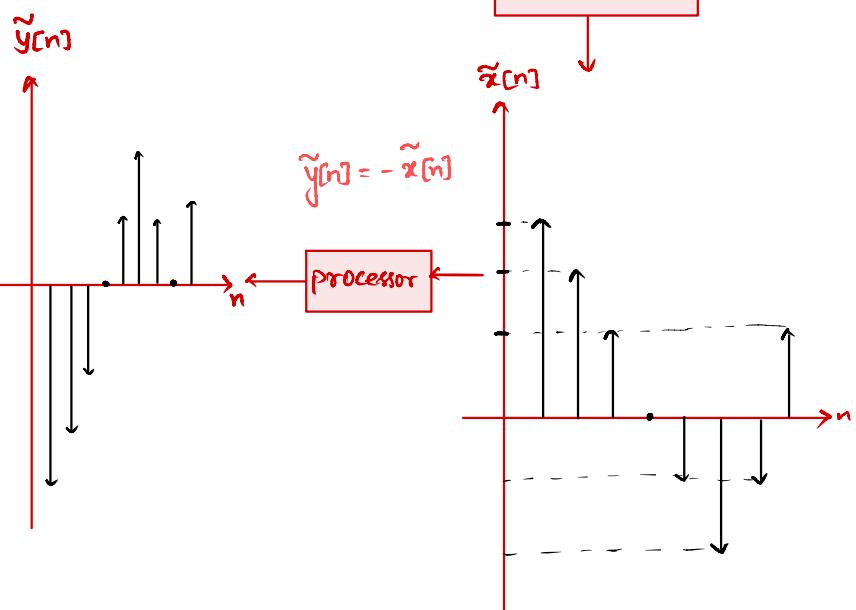
Discrete in time  $\equiv n \in \mathbb{Z}$

Continuous in magnitude  $\equiv x \in \mathbb{R}$



Continuous in time  $\equiv t \in \mathbb{R}$

Discrete in amplitude  $\equiv y \in \mathbb{Z}$



Digital Signal

Discrete in time  $\equiv n \in \mathbb{Z}$

Discrete in magnitude  $\equiv \tilde{x} \in \mathbb{Z}$

## Representation Formats:

6 sensors

Analog

Digital

$$\begin{array}{ccc}
 & \xrightarrow{\hspace{1cm}} & \\
 x^{(0)}(t) & & x^{(0)}[n] \\
 x^{(1)}(t) & & x^{(1)}[n] \\
 x^{(2)}(t) & \longleftrightarrow & x^{(2)}[n] \\
 \vdots & & \vdots \\
 x^{(5)}(t) & & x^{(5)}[n]
 \end{array}$$

$$\text{Sampling Rate} = 10 \text{ sps} = 10 \text{ Hz}$$

$$\text{No. of bits / sample (bps)} = 16 \text{ bits}$$

$$\begin{aligned}
 \text{Data Size in 1s} &= \text{No. of Sensors} \times \text{Sampling Rate} \times \text{Time} \times \text{bps} \\
 &= 6 \times 10 \text{ samples s}^{-1} \times 1s \times 16 \text{ bits sample}^{-1} \\
 &= 960 \text{ bits} \\
 &= 120 \text{ bytes}
 \end{aligned}$$

$$\begin{aligned}
 \text{Data in 1 year} &= 365 \text{ days} \times 24 \text{ hrs day}^{-1} \times 3600 \text{ s hr}^{-1} \\
 &\quad \times 120 \text{ bytes s}^{-1} \\
 &\approx 3.84 \text{ B}
 \end{aligned}$$

Data per year in Nalanda Complex =  $(3.8 \times 96 \times 5)$  GB  
 $\approx 1.8$  TB

### Sensors

- CO<sub>x</sub>
- NO<sub>x</sub>
- SO<sub>x</sub>
- Temperature
- Humidity

$$\text{No of cables} = (6 \times 96) \\ \approx 580$$

If  $K \ll D$ ,  
 $\downarrow$        $\downarrow$   
 $2$        $580$

### Data Compression:

$$\underline{x}[n] = [x^{(0)}[n], x^{(1)}[n], \dots, x^{(d)}[n], \dots, x^{(D-1)}[n]]^T \in \mathbb{R}^D$$

↳ Vector valued sample at  $n^{\text{th}}$  instance

$$\underline{X} = [\underline{x}[0], \underline{x}[1], \underline{x}[2], \dots, \underline{x}[n], \dots, \underline{x}[N-1]] \rightarrow \\ = [x_0, x_1, \dots, x_n, \dots, x_{N-1}]$$

$$= \left[ \begin{array}{c} x^{(0)}[0] \\ x^{(1)}[0] \\ \vdots \\ x^{(d)}[0] \\ \vdots \\ x^{(D-1)}[0] \end{array} \right] \left[ \begin{array}{c} x^{(0)}[1] \\ x^{(1)}[1] \\ \vdots \\ x^{(d)}[1] \\ \vdots \\ x^{(D-1)}[1] \end{array} \right] \cdots \left[ \begin{array}{c} x^{(0)}[N-1] \\ x^{(1)}[N-1] \\ \vdots \\ x^{(d)}[N-1] \\ \vdots \\ x^{(D-1)}[N-1] \end{array} \right]$$

$$= \begin{bmatrix} x_0^{(0)} & x_1^{(0)} & \dots & \dots & x_n^{(0)} & \dots & x_{N-1}^{(0)} \\ x_0^{(1)} & x_1^{(1)} & & & & & \\ x_0^{(2)} & x_1^{(2)} & & & & & \\ \vdots & \vdots & & & & & \\ x_0^{(d)} & x_1^{(d)} & & & x_n^{(d)} & & \\ \vdots & \vdots & & & \vdots & & \\ x_0^{(D-1)} & x_1^{(D-1)} & \dots & \dots & x_n^{(D-1)} & \dots & x_{N-1}^{(D-1)} \end{bmatrix}$$

$E \in \mathbb{R}^{D \times N}$

Number of samples

Dimensions of data samples

$$y = A^T x \rightarrow \mathbb{R}^{D \times 1}$$

$\downarrow$

$E \in \mathbb{R}^{K \times D}$

Q) Prove:-  $y = A^T x$  is linear.

$$y_1 = A^T x_1$$

$$y_2 = A^T x_2$$

$$y = y_1 + y_2$$

$$= A^T x_1 + A^T x_2$$

$$y = A^T x$$

When  $x$  is replaced with  $Kx$  ( $K \in \mathbb{R}$ )

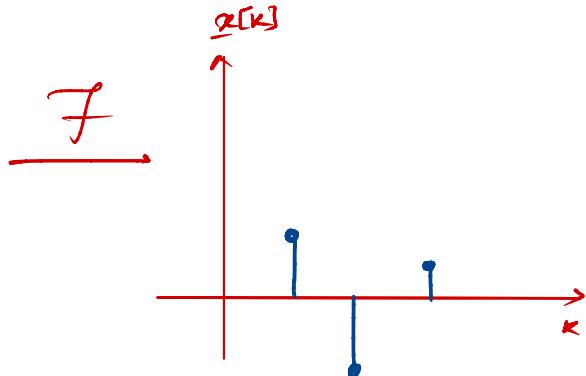
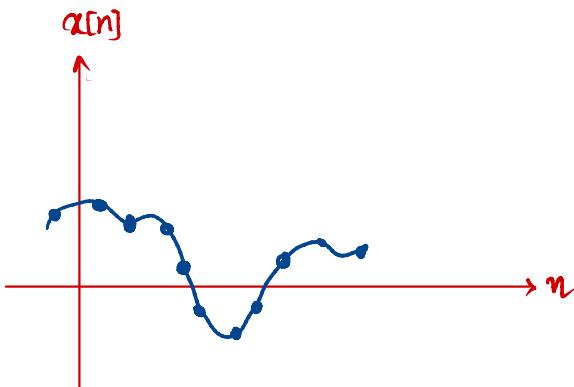
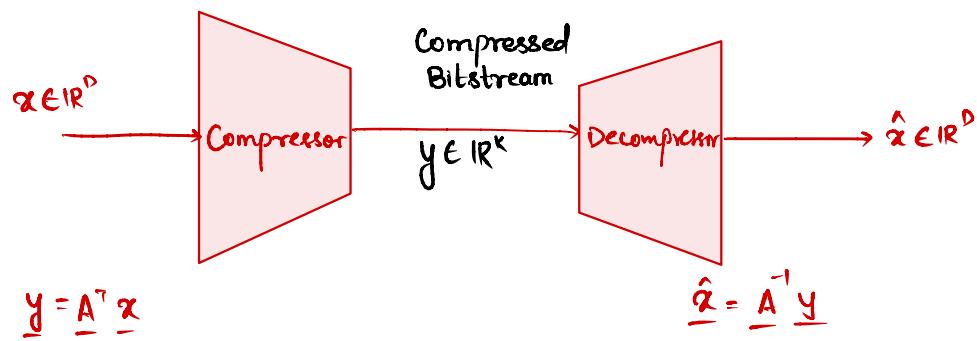
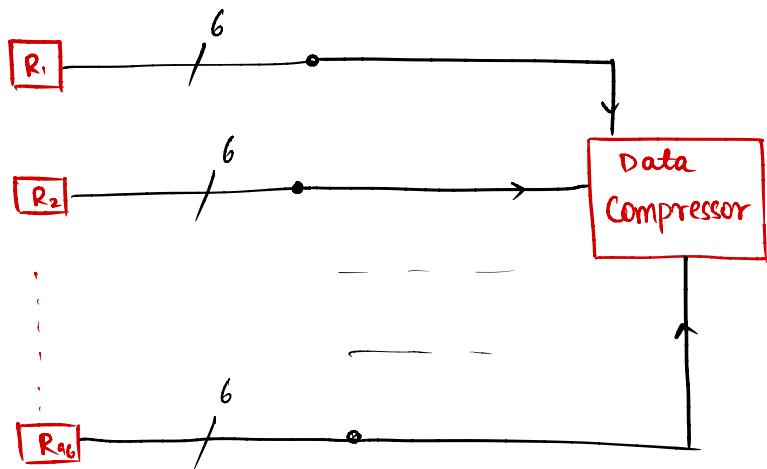
$$y' = A^T (Kx)$$

$$= K A^T x$$

$$= Ky$$

$$= \underline{A}^T (\underline{x}_1 + \underline{x}_2)$$

$$= \underline{A}^T \underline{x}$$



$$\underline{x} = [x[-2], x[-1], x[0], x[1], \dots, x[N-3]] \in \mathbb{R}^N$$

$\downarrow F$

$$\underline{y} = [y[2], y[1], y[0], y[1], \dots, y[N-3]] \in \mathbb{R}^N$$

↓ ↓ ↓ ↓ ↓

$$\underline{y} = F(\underline{x}) \quad (K \ll N)$$

$\downarrow \mathbb{R}^N$

$$\underline{y} = \underline{A}^T \underline{x}$$

↳ DFT coefficient matrix

(Analytically computed in case of FFT, DCT, DWT)

Q] Can we obtain  $\underline{A}^T$  from  $\underline{x}$  s.t.  $\underline{y} = \underline{A}^T \underline{x}$ ,  $\underline{y} \in \mathbb{R}^{K \times N}$ ,  $\underline{x} \in \mathbb{R}^{D \times N}$ ,  $K \ll D$ .

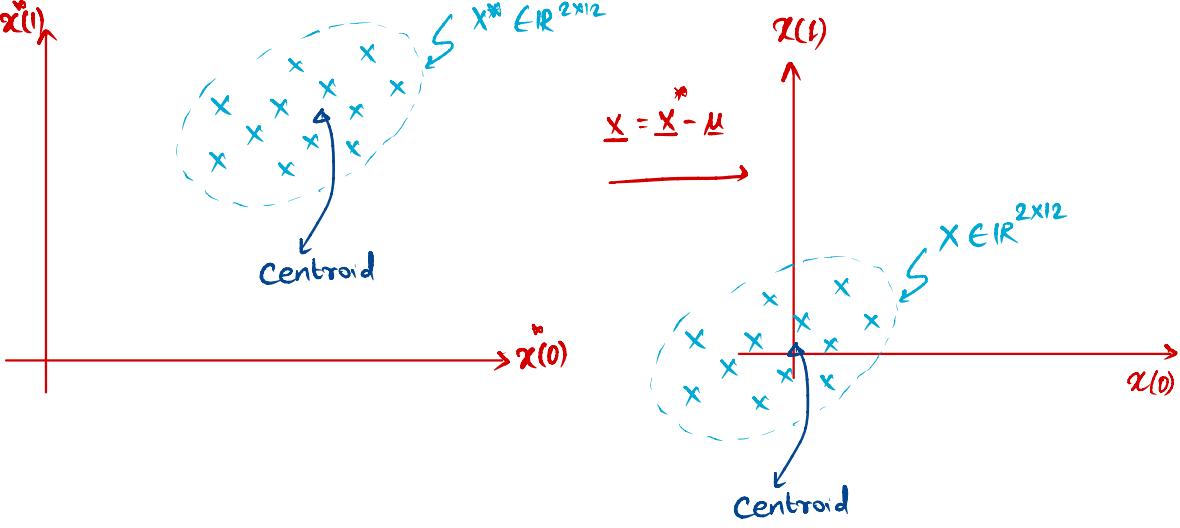
$$D=2$$

$$\underline{x}_n = [x_n^{(0)} \ x_n^{(1)}]^T \in \mathbb{R}^2$$

$$\underline{x} = [\underline{x}_0 \ \underline{x}_1 \ \dots \ \underline{x}_{N-1}] \in \mathbb{R}^{2 \times N}$$

$$\underline{y}_n = [y_n^{(0)}]^T \in \mathbb{R}$$

$$\underline{y} = [\underline{y}_0 \ \underline{y}_1 \ \underline{y}_2 \ \dots \ \underline{y}_{N-1}] \in \mathbb{R}^{1 \times N}$$



$$\underline{\mu} = [\mu^{(0)} \ \mu^{(1)} \ \dots \ \mu^d, \dots, \mu^{(D-1)}]$$

$$\text{bt } \mu^{(d)} = \frac{1}{N} \sum_{n=0}^{N-1} x_n^{(d)}$$

$$\underline{a}^{(0)} = [a_{00}, a_{10}] \in \mathbb{R}^{2 \times 1}$$

$$\underline{a}^{(1)} = [a_{01}, a_{11}] \in \mathbb{R}^{2 \times 1}$$

$$y_n^{(0)} = \underline{a}^{(0)\top} \underline{x}_n$$

$$y_n^{(1)} = \underline{a}^{(1)\top} \underline{x}_n$$

$$\underline{A}^\top = \begin{bmatrix} \underline{a}^{(0)\top} \\ \underline{a}^{(1)\top} \end{bmatrix} = \begin{bmatrix} a_{00} & a_{10} \\ a_{01} & a_{11} \end{bmatrix}$$

$$\underline{y}_n = \underline{A}^T \underline{x}_n$$

$$\underline{A} = \begin{bmatrix} a_{00} & a_{10} \\ a_{01} & a_{11} \end{bmatrix} = \begin{bmatrix} a^{(0)} & a^{(1)} \\ \lambda^{(0)} & \lambda^{(1)} \end{bmatrix} \in \mathbb{R}^{2 \times 2}$$

↓  
Eigenvalues

$$\xrightarrow{\quad \text{xxxx} \quad + \quad \text{xxx} \quad} y^{(0)}$$

Higher  $\lambda$

$$\xrightarrow{\quad \text{xxxx} \quad + \quad \text{xxxx} \quad} y^{(1)}$$

Lower  $\lambda$

$$\underline{x}_n \in \mathbb{R}^D \quad \underline{y}_n \in \mathbb{R}^k ; \quad k \ll D$$

$$\underline{y}_n = \underline{A}^T \underline{x}_n \quad \underline{x}_n \in \mathbb{R}^D$$

$\mathbb{R}^k$   
 $\mathbb{R}^{D \times D}$

$$\left[ \begin{array}{c} y^{(0)} \\ y^{(1)} \\ \vdots \\ y^{(d)} \\ \vdots \\ y^{(D-1)} \end{array} \right] = \left[ \begin{array}{cc|cc|cc} a_{00} & a_{10} & a_{00} & a_{10} & a_{00} & a_{10} \\ a_{01} & a_{11} & a_{01} & a_{11} & a_{01} & a_{11} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{0d} & a_{1d} & a_{0d} & a_{1d} & a_{0d} & a_{1d} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{0,D-1} & a_{1,D-1} & a_{0,D-1} & a_{1,D-1} & a_{0,D-1} & a_{1,D-1} \end{array} \right] \left[ \begin{array}{c} x^{(0)} \\ x^{(1)} \\ \vdots \\ x^{(d)} \\ \vdots \\ x^{(D-1)} \end{array} \right]$$

$$\underline{\tilde{y}_n} = [\underline{y}^{(0)}, \underline{y}^{(1)}, \dots, \underline{y}^{(d)}, \dots, \underline{y}^{(n)}]$$

$$\hat{\underline{x}}_n = (\underline{A}^{nT})^T \underline{\tilde{y}_n}$$

$$||\underline{x}_n - \hat{\underline{x}}_n|| = \epsilon_n$$

$$\epsilon_n^2 = \sum_{j=k}^{n-1} \hat{x}_j^2$$

## Singular Value Decomposition (SVD)

PCA:  $\underline{x}_i \rightarrow \underline{\tilde{y}_i}$

$$\underline{\tilde{y}_i} = \underline{\tilde{A}^T} \underline{x}_i$$

$$\begin{array}{c} \underline{x} \rightarrow \underline{\tilde{Y}} \\ \underline{\tilde{Y}} \rightarrow \underline{A^T} \underline{x} \end{array}$$

$$\underline{x}_i \in \mathbb{R}^{D \times 1}, \underline{\tilde{y}_i} \in \mathbb{R}^{K \times 1}$$

$$\underline{x} = \underline{U} \begin{bmatrix} \Delta^{1/2} & \underline{0}_a \\ \underline{0}_b & \underline{0}_c \end{bmatrix} \underline{V}^H$$

$\underline{U} \in \mathbb{R}^{D \times D}$   
 $\Delta^{1/2} \in \mathbb{R}^{D \times D}$   
 $\underline{V}^H \in \mathbb{R}^{N \times N}$   
 $\underline{0}_a \in \mathbb{R}^{D \times (N-D)}$   
 $\underline{0}_b \in \mathbb{R}^{D \times D}$   
 $\underline{0}_c \in \mathbb{R}^{0 \times (N-D)}$

$$\begin{array}{l} \underline{x} \in \mathbb{R}^{D \times 1} \\ \Delta^{1/2} \in \mathbb{R}^{D \times D} \\ \underline{U} \in \mathbb{R}^{D \times D} \\ \underline{V}^H \in \mathbb{R}^{D \times N} \end{array}$$

$$\Delta^{1/2} = \begin{bmatrix} \sqrt{\lambda_0} & 0 & 0 & 0 & \cdots & 0 \\ 0 & \sqrt{\lambda_1} & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & \cdots & \sqrt{\lambda_{D+1}} \end{bmatrix} \quad \left( |\lambda_0| > |\lambda_1| > |\lambda_2| > \cdots > |\lambda_{D+1}| \right)$$

- $\lambda^d$  is the Eigenvalue of  $\underline{x} \underline{x}^\top$

$$X^T \in \mathbb{R}^{N \times D}$$

$$\underline{X} \in \mathbb{R}^{D \times N}$$

$$\underline{X} \underline{X}^T \in \mathbb{R}^{D \times D}$$

- U is the set of Eigenvectors of  $XX^T$  ( $XX^T \in \mathbb{R}^{N \times N}$ )

$$\underline{u} = [\underline{u}_0 \quad \underline{u}_1 \quad \underline{u}_2 \quad \dots \quad \underline{u}_d \quad \dots \quad u_{D-1}]$$

$EIR^{D \times 1}$

- $V^H$  is the set of Eigenvectors of  $\underline{X}^T \underline{X}$  ( $\underline{X}^T \underline{X} \in \mathbb{R}^{N \times N}$ )

$$\underline{V}^H = [\underline{v}_0 \quad \underline{v}_1 \quad \underline{v}_2 \quad \dots \quad \underline{v}_{n-1} \quad \dots \quad \underline{v}_{N-1}]$$

$\downarrow$   
 $\in \mathbb{R}^{N \times 1}$

p-bytes per number

(K<<D , K<<N)

Double - 8 bytes

$$\text{sizeof } (\underline{x}) = \text{PDN}$$

`sizeof (U) + sizeof (A10) + sizeof (V)`

$$= P(DK + K + KN)$$

$$= \mathcal{P}(OK + K + KN) \times \mathcal{P}(DN)$$

$$\Rightarrow (KD + K + KN) < DN$$

$\underline{U} \in \mathbb{R}^{D \times D}$  $\underline{\Lambda}_k \in \mathbb{R}^{K \times K}$  $\underline{V} \in \mathbb{R}^{N \times N}$ 

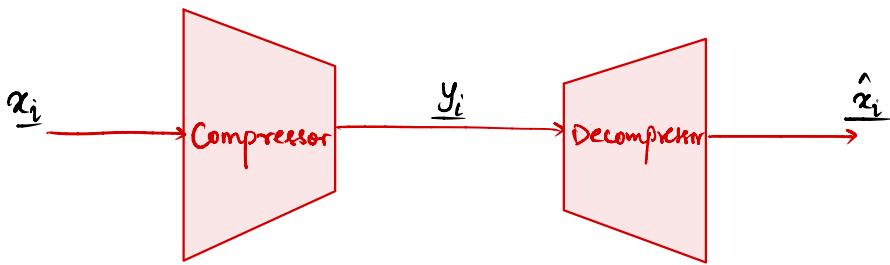
$$\tilde{\underline{X}}_k = \begin{bmatrix} & 1 \\ \underline{U}_k \in \mathbb{R}^{D \times n} & | \\ & 0 \end{bmatrix}^D \times (D-k) \quad \boxed{\begin{bmatrix} \sqrt{\lambda_1} & 0 & 0 & \cdots & 0 \\ 0 & \sqrt{\lambda_2} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \cdots & \sqrt{\lambda_K} \end{bmatrix}} \quad 0^{K \times (N-K)} \quad \underline{V}_k^H \in \mathbb{R}^{K \times N}$$

Find  $k$  such that  $\min ||\underline{X} - \tilde{\underline{X}}||$  and  $\max \left( \frac{\text{sizeof}(\underline{X})}{\text{sizeof}(\underline{U}_k, \underline{\Lambda}_k, \underline{V}_k)} \right)$

↓  
Compression  
Factor (CF)

## Data Transform

- Learnable Transform
- Fixed Transform



PCA:

$$\underline{y}_i = \tilde{A}^T \underline{x}_i$$

$$\underline{y}_i \in \mathbb{R}^{k \times 1}$$

$$\underline{x}_i \in \mathbb{R}^{D \times 1}$$

$$\hat{\underline{x}}_i = \frac{\hat{A} \underline{y}_i}{L}$$

$$\hat{\underline{x}}_i \in \mathbb{R}^{D \times 1}$$

$$\underline{y}_i \in \mathbb{R}^{k \times 1}$$

$(k < D)$

Compression Factor (CF) =  $\frac{\text{Sizeof}(X)}{\text{Sizeof}(Y)}$

SVD:

$$\underline{U}_k \Lambda_k^{1/2} \underline{V}_k^H \approx \underline{X}$$

$$\hat{\underline{X}} = \underline{U}_k \Lambda_k^{1/2} \underline{V}_k^H$$

Compression Factor (CF) =  $\frac{\text{Sizeof}(X)}{\text{Sizeof}(U) + \text{Sizeof}(\Lambda_k^{1/2}) + \text{Sizeof}(\underline{V}_k^H)}$

PDN → Precision of number (bits per number)

double → precision floating point number (64 bits)

single → 32 bits

half → 16 bits

quarter → 8 bits

$$\text{Compression Factor (CF)} = \frac{\text{Size of } (X)}{\text{Size of } (Y)} = \frac{P_{DN}}{q_{DN}} = \frac{P}{q}$$

Dynamic Range of Representing a number

$$X \in \mathbb{R}^{D \times N} \longrightarrow Y \in \mathbb{R}^{D \times N}$$

## Information Theory:-

$X \in [0, 255] \cap \mathbb{R} \longrightarrow \text{Infinite bits}$

$\epsilon [0, 255] \cap \mathbb{Z} \longrightarrow 8 \text{ bits}$

$$2^0 - 1 = 1$$

$$2^1 - 1 = 3$$

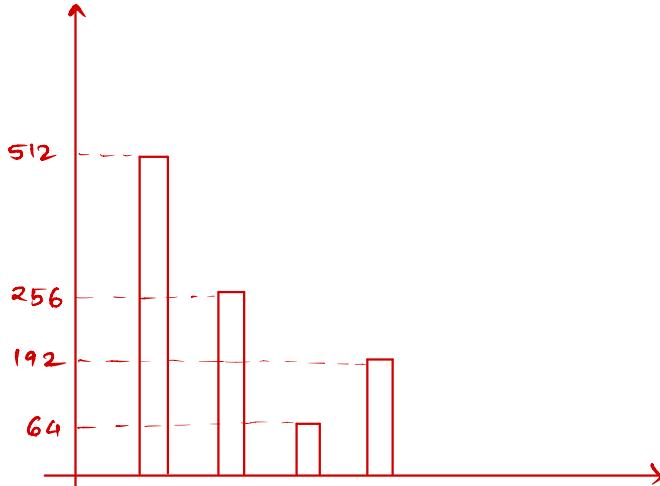
:

:

$$2^8 - 1 = 255$$

$$X = [x_0 \ x_1 \ x_2 \ \dots \ x_{N-1}] \in \mathbb{R}^N$$

Count



for representation of this no. of bits are required is \_\_\_\_\_.

Symbol      Bitcodes

10	00
100	01
190	10
250	11

$\underline{x} \rightarrow \underline{y}$  (Encoder)

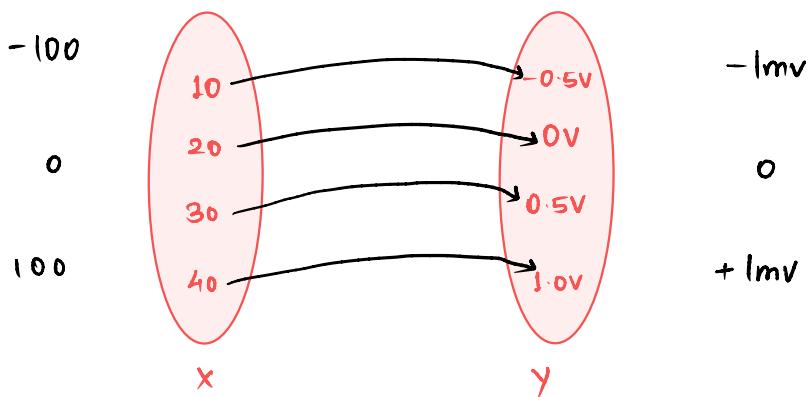
$\underline{y} \rightarrow \tilde{\underline{x}}$  (Decoder)

$$|\underline{x} - \tilde{\underline{x}}| = 0$$

↳ So, lossless system

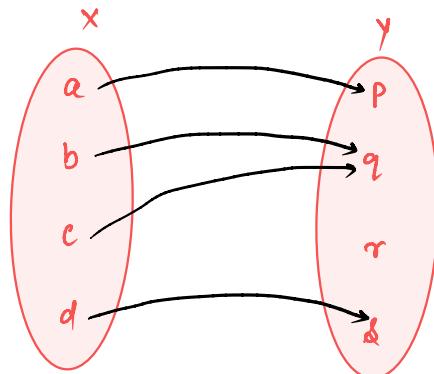
- In worst case, CF is equal to 1.

Claude Elwood Shannon -



1-1 mapping

Thickness - 3m



$$W = H + GS$$

Restrained  
Channel  
Capacity

$$\text{Entropy} = \frac{\text{Heat}}{\text{Temp}}$$

$$\text{of gen} = KT$$

Boltzmann's  
Constant

$$10 - a - 512 \quad P(a) = 1/2$$

$$100 - b - 256 \quad P(b) = 1/4$$

$$190 - c - 64 \quad P(c) = 1/16$$

$$250 - d - 192 \quad P(d) = 3/16$$

$$\sum = 1024$$

- Information carried by a signal,

$$I(a) = \log_r P(a)$$

$$r = 2 \text{ bits}$$

$$r = e \text{-nats}$$

$$H(X) = E(I(s)) = - \sum_{s \in X} p(s) \log_2(p(s)) \geq 0$$

↓  
Entropy

Joint Entropy:  $H(X,Y) = - \sum_{s \in X, r \in Y} p(s,r) \log_2(p(s,r)) \geq 0$

<del>x</del> <del>y</del>	a	b	c	d	
P	1/8	1/16	1/32	1/32	1/4 = P(P)
q	1/16	1/8	1/32	1/32	1/4 = P(q)
r	1/16	1/16	1/16	1/16	1/4 = P(r)
s	1/4	0	0	0	1/4 = P(s)
	1/2	1/4	1/8	1/8	
	"	"	"	"	
	p(a)	p(b)	p(c)	p(d)	

$$H(X) = - \left( \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{4} \log_2 \frac{1}{4} + \frac{1}{8} \log_2 \frac{1}{8} + \frac{1}{8} \log_2 \frac{1}{8} \right)$$

$$= 2.25 \text{ bits}$$

$$H(Y) = - \left( \frac{1}{4} \log_2 \frac{1}{4} + \frac{1}{4} \log_2 \frac{1}{4} + \frac{1}{4} \log_2 \frac{1}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right)$$

$$= 2 \text{ bits}$$

$$H(X,Y) = \frac{27}{8} \text{ bits}$$

$$H(X, Y) = H(X) + H(Y/X) = H(Y) + H(X/Y)$$

Relative Entropy  $\rightarrow$  Kullback - Leiber Divergence

Divergence  $\rightarrow$  Effective distance

$$\begin{aligned} D(P \parallel Q) &= \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} \\ &= E_P \log \frac{P(x)}{Q(x)} \end{aligned}$$

Mutual Information :-

$$E_{P(X,Y)} \log \frac{P(X,Y)}{P(X)P(Y)} = H(X) - H(X/Y)$$

$$I(X;Y) = D(P(X,Y) \parallel P(X)P(Y))$$

$$I(X,X) = H(X)$$

- For symmetric channel  $= H(Y) - H(Y/X)$

Huffman Compression :-

Coconutcurry

$$X = ['c', 'o', 'c', 'o', 'n', 'u', 't', 'c', 'u', 'r', 'r', 'y] \in \Sigma^{8 \times 12}$$

Symbol	No of occurrence	P(Symbol)
o c	3	3/12
o o	2	2/12
o n	1	1/12
o u	2	2/12

0 1 1 0	t	1	1/12
1 1 1	r	2	2/12
0 1 1 1	y	1	1/12
		$\sum = 12$	

Reorder  
Descending

$$c - 3/12$$

$$o - 2/12$$

$$u - 2/12$$

$$r - 2/12$$

$$n - 4/12$$

$$t - 4/12 \\ y - 4/12 \left\{ \begin{array}{l} \{t,y\} - 2/12 \end{array} \right.$$

$$\{u,r\} - 4/12$$

$$c - 3/12$$

$$\{\{t,y\}, n\} - 3/12 \\ o - 2/12 \left\{ \begin{array}{l} \{\{\{t,y\}, n\}, o\} - 8/12 \end{array} \right.$$

$$c - 3/12$$

$$o - 2/12$$

$$u - 2/12$$

$$r - 2/12$$

$$\left. \begin{array}{l} \{t,y\} - 2/12 \\ n - 4/12 \end{array} \right\} \begin{array}{l} n_5 \\ \{ \{t,y\}, n \} - 3/12 \end{array}$$



$$c - 3/12$$

$$\{\{t,y\}, n\} - 3/12 \quad n_4$$

$$o - 2/12$$

$$u - 2/12$$

$$r - 2/12$$

$$\left. \begin{array}{l} \{u,r\} - 4/12 \end{array} \right\}$$

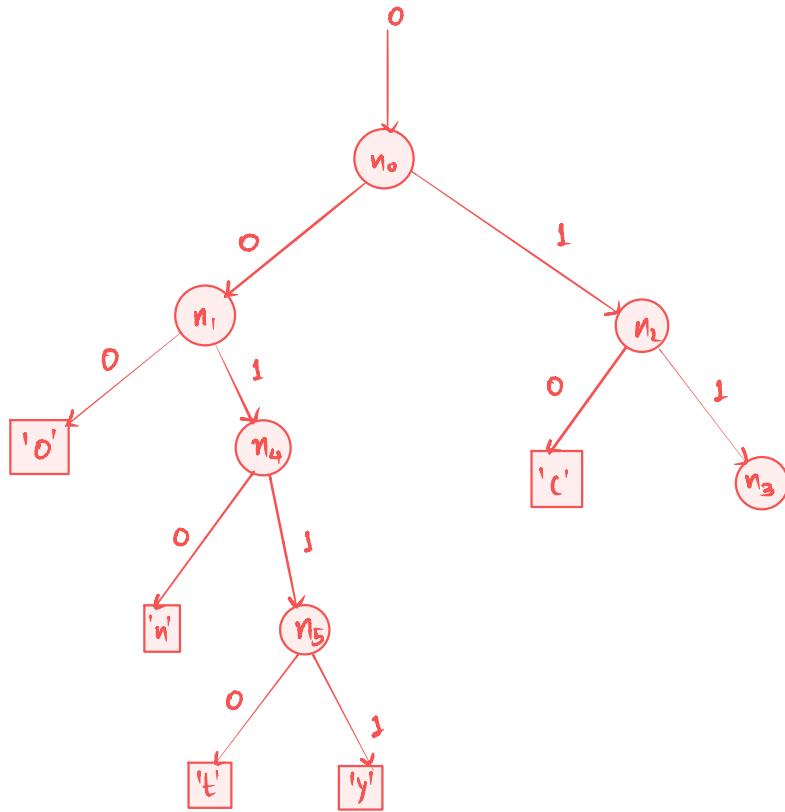
$$\{ \{ \{ t, y \}, n \}, o \} - 5/12$$

$$\left. \begin{array}{l} \{ u, r \} = 4/12 \\ c = 3/12 \end{array} \right\} \begin{array}{l} \{ \{ u, r \}, c \} = 7/12 \\ \{ \{ \{ t, y \}, n \}, o \} = 5/12 \end{array}$$

→

$$\begin{array}{l} \{ \{ u, r \}, c \} = 7/12 \\ n_2 \end{array}$$

$$\begin{array}{l} \{ \{ \{ t, y \}, n \}, o \} = 5/12 \\ n_1 \end{array}$$

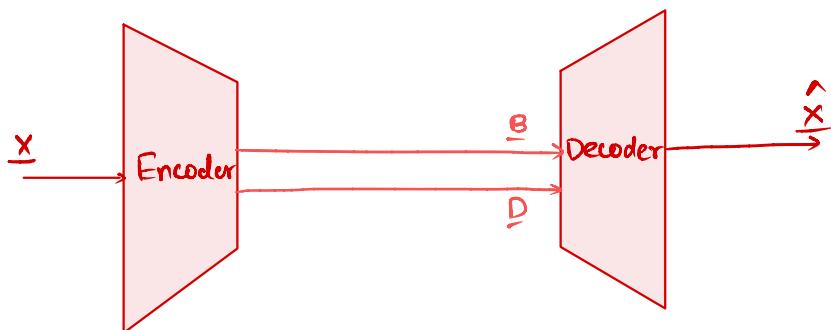


$$\text{No of bits} = 33$$

$$\begin{aligned} \text{Avg no. of bits / symbol} &= \frac{33}{12} \\ &= 2.75 \end{aligned}$$

$$H(x) = - \sum_{\text{symbol}} p(\text{symbol}) \log_2 p(\text{symbol})$$

$$= - \left( \frac{3}{12} \log_2 \frac{3}{12} + \frac{6}{12} \log_2 \frac{6}{12} + \frac{3}{12} \log_2 \frac{1}{12} \right) = 2.69$$



Header

D  
q-bits

Data

B  
33-bits

$MSE(X, \hat{X}) = 0 \rightarrow \text{Lossless}$

Compression Factor (CF) =  $\frac{\text{No of bits in } X}{\text{No of bits in } B}$

$$= \frac{12 \times 8}{33}$$

$$\approx 2.9$$

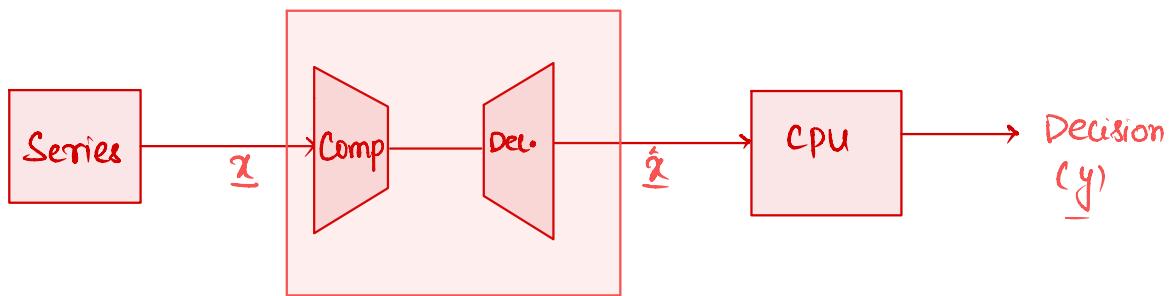
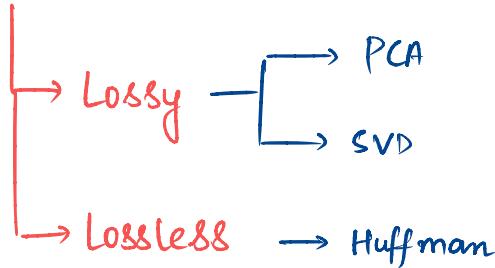
Max CF =  $\frac{\text{Avg. no of bits in } X}{H(X)}$

$$\text{Coding Efficiency} = \frac{H(x)}{\text{Avg. bits per symbol}}$$

$$= \frac{2.69}{2.75}$$

$$\approx 97\%$$

## Data Compression



## Bayes' Decision Theorem:-

$$P(y_j | \underline{x}) = \frac{P(\underline{x} | y_j) \cdot P(y_j)}{P(\underline{x})}$$

↓ Likelihood      ↓ Prior  
 posterior      Evidence

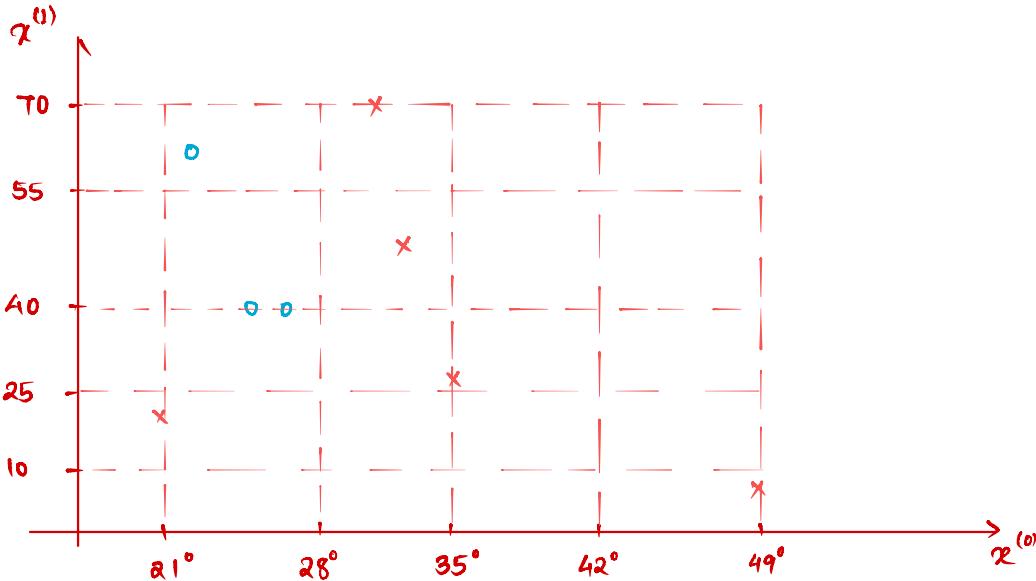
$\underline{x}^{(0)}$	$\underline{x}^{(1)}$		$y^{(0)}$	$y^{(1)}$	
Temp	Humidity		Comfort	N-Comfort	
$\underline{x}_0^T$	31°C	70 % rh	0	1	$\underline{y}_0$
$\underline{x}_1^T$	24°C	40 % rh	1	0	$\underline{y}_1$
$\underline{x}_2^T$	22°C	60 % rh	1	0	$\underline{y}_2$
$\underline{x}_3^T$	49°C	10 % rh	0	1	$\underline{y}_3$
$\underline{x}_4^T$	35°C	30 % rh	0	1	$\underline{y}_4$
$\underline{x}_5^T$	32°C	50 % rh	0	1	$\underline{y}_5$
$\underline{x}_6^T$	27°C	40 % rh	1	0	$\underline{y}_6$
$\underline{x}_7^T$	21°C	20 % rh	0	1	$\underline{y}_7$

$$\underline{x}_i \in \mathbb{R}^{2 \times 1}$$

$$\underline{X} = [\underline{x}_0, \underline{x}_1, \underline{x}_2, \dots, \underline{x}_7] \in \mathbb{R}^{2 \times 8}$$

$$Y_i \in \mathbb{B}^{2 \times 1}$$

$$\underline{Y} = [\underline{y}_0, \underline{y}_1, \dots, \underline{y}_7] \in \mathbb{B}^{2 \times 8}$$



$$P(y = \text{comfort} | \underline{x}) = \frac{P(\underline{x} | y = \text{comfort}) \cdot P(y = \text{comfort})}{P(\underline{x})}$$

$$P(y = \text{N. Comfort} | \underline{x}) = \frac{P(\underline{x} | y = \text{N. Comfort}) \cdot P(\text{N. Comfort})}{P(\underline{x})}$$

$$P(\text{comfort}) = \frac{3}{8} = (P(y = [s, o])) \\ = (P(y = \underline{\omega}_0))$$

$$P(\text{N. Comfort}) = \frac{5}{8} = (P(y = [0, s])) \\ = (P(y = \underline{\omega}_1))$$

Class  
Label

$$h(\underline{x}) = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

$$h(\underline{x} | y=w_0) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$h(\underline{x} | y=w_1) = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

$$\underline{x} = [27, 48]$$

$$P(\underline{x} | y=w_0) = \frac{h(\underline{x} | y=w_0)}{\sum_{\underline{x}} h(\underline{x} | y=w_0)} = \frac{2}{3}$$

$$P(\underline{x} | y=w_1) = \frac{h(\underline{x} | y=w_1)}{\sum_{\underline{x}} h(\underline{x} | y=w_1)} = \frac{0}{5} = 0$$

$$P(y=w_0) = \frac{3}{8}$$

$$P(y=w_1) = \frac{5}{8}$$

$$P(\underline{x}) = P(\underline{x} | y = \omega_0) \cdot P(y = \omega_0) + P(\underline{x} | y = \omega_1) \cdot P(y = \omega_1)$$

$$= \left(\frac{2}{3}\right)\left(\frac{3}{8}\right) + 0 \cdot \frac{5}{8}$$

$$= \frac{1}{4}$$

$$P(y = \omega_0 | \underline{x}) = \frac{P(\underline{x} | y = \omega_0) \cdot P(y = \omega_0)}{P(\underline{x})}$$

$$= \frac{\frac{2}{3} \times \frac{3}{8}}{\frac{1}{4}}$$

$$= 1$$

$$P(y = \omega_1 | \underline{x}) = \frac{P(\underline{x} | y = \omega_1) \cdot P(y = \omega_1)}{P(\underline{x})}$$

$$= \frac{0 \cdot \frac{5}{8}}{\frac{1}{4}}$$

$$= 0$$

$$\underline{x} = [38, 47]$$

$$P(\underline{x} | y = \omega_0) = \frac{h(\underline{x} | y = \omega_0)}{\sum_{\underline{x}} h(\underline{x} | y = \omega_0)}$$

$$= \frac{0}{3}$$

$$= 0$$

$$P(\underline{x}|y=w_i) = \frac{h(\underline{x}|y=w_i)}{\sum_{\underline{x}} h(\underline{x}|y=w_i)}$$

$$= 0$$

$$P(\underline{x}) = \sum_{w_0, w_1} P(\underline{x}|y_j) \cdot p(y=j)$$

$$= P(x|y=w_0) \cdot p(y=w_0) + P(x|y=w_1) \cdot p(y=w_1)$$

$$= 0 \cdot \frac{3}{8} + 0 \cdot \frac{5}{8}$$

$$= 0$$

$$P(y=w_0|\underline{x}) = \frac{P(\underline{x}|y=w_0) \cdot p(y=w_0)}{P(\underline{x})}$$

$$= \frac{0 \cdot \frac{3}{8}}{0}$$

$$= 0\% \text{ (Undefined)}$$

$$P(y=w_1|\underline{x}) = \frac{P(\underline{x}|y=w_1) \cdot p(y=w_1)}{P(\underline{x})}$$

$$= \frac{0 \cdot \frac{5}{8}}{0}$$

$$= 0\% \text{ (Undefined)}$$

## Bayes Decision:-

$$P(y = w_j | \mathbf{x}) = \frac{P(\mathbf{x} | y = w_j) \cdot P(y = w_j)}{P(\mathbf{x})}$$

Likelihood  
Non-zero  
Non-zero

## Gaussian Likelihood :-

$$P(\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\mathbf{x}-\mu)^2}{2\sigma^2}} ; \quad \mathbf{x} \in \mathbb{R}, \quad D \times 2$$

## Naïve Bayes :-

$$P(\mathbf{z}) = P(z^{(0)}) \cdot P(z^{(1)}) \cdot P(z^{(2)}) \cdots P(z^{(D-1)})$$

Gaussian Poisson Rayleigh

(Assumption is  $z^{(d)}$  is independent)

## Multivariate Normal Distribution :-

$$P(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right)$$

$$\begin{bmatrix} \mathbf{x} \end{bmatrix} \rightarrow \begin{array}{l} \mathbf{x} \in \mathbb{R}^{D \times 1} \\ \mu \in \mathbb{R}^{D \times 1} \\ \Sigma \in \mathbb{R}^{D \times D} \end{array}$$

$$\mu = [\mu^{(0)} \mu^{(1)} \cdots \mu^{(n)} \cdots \mu^{(D-1)}]$$

$$\Sigma = \begin{bmatrix} \sigma_{(1)}^2 & \sigma_{(1,2)}^2 \\ \sigma_{(2,1)}^2 & \ddots \end{bmatrix}$$

$$\sigma^{(g,h)} = \frac{1}{N} \sum_{i=0}^{N-1} (\mathbf{x}_i^{(g)} - \mu^{(g)}) (\mathbf{x}_i^{(h)} - \mu^{(h)})$$

	$\underline{x}^{(0)}$	$\underline{x}^{(1)}$	$y$
$\underline{x}_0$	$x_0^{(0)}$	$x_0^{(1)}$	$w_0$
$\underline{x}_1$	$x_1^{(0)}$	$x_1^{(1)}$	$w_1$
$\underline{x}_2$	$x_2^{(0)}$	$x_2^{(1)}$	$w_2$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\underline{x}_{N-1}$	$x_{N-1}^{(0)}$	$x_{N-1}^{(1)}$	$w_{N-1}$

Maximum Likelihood Estimation :-

$$P(\underline{x} | y = w_j) = \frac{1}{(2\pi)^{d/2} |\sum_{y=w_j}|^{1/2}} \exp^{-\frac{1}{2} ((\underline{x} - \underline{\mu}_{y=w_j})^T \sum_{y=w_j}^{-1} (\underline{x} - \underline{\mu}_{y=w_j}))}$$

$$\begin{aligned}\underline{\mu}_{y=w_j} &= [\mu_{y=w_j}^{(0)} \quad \mu_{y=w_j}^{(1)} \quad \cdots \quad \mu_{y=w_j}^{(d)} \quad \cdots \quad \mu_{y=w_j}^{(D-1)}] \\ \underline{\mu}_{y=w_j}^{(d)} &= \frac{1}{N_j} \sum_{y_i \in w_j} x_i^{(d)}\end{aligned}$$

$$N_j = |\underline{y} = w_j|$$

$$\sum_{y=w_j} = \begin{bmatrix} \sigma_{y=w_j}^{(0,0)} & \sigma_{y=w_j}^{(0,1)} & \cdots & \cdots \\ \sigma_{y=w_j}^{(1,0)} & \sigma_{y=w_j}^{(1,1)} & \cdots & \cdots \\ \vdots & \vdots & \ddots & \ddots \end{bmatrix}$$

$$\sigma_{y=w_j}^{(g,h)} = \frac{1}{N_j} \sum (x_i^{(g)} - \mu_{y=w_j}^{(g)}) (x_i^{(h)} - \mu_{y=w_j}^{(h)})$$

$$P(y = w_j | \underline{x}) = \frac{P(\underline{x} | y = w_j) \cdot P(y = w_j)}{P(\underline{x})}$$

$$P(y = w_0 | \underline{x})$$

$$P(y = w_1 | \underline{x})$$

⋮

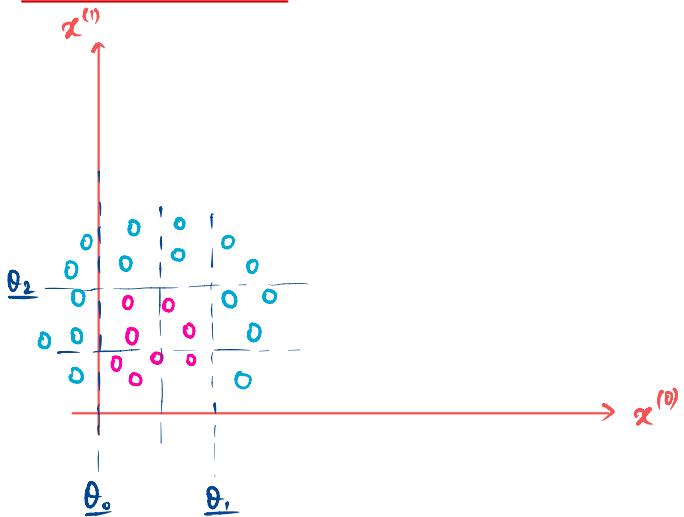
Maximum A posteriori (MAP) :-

$$\hat{y} = \underset{\forall w_j \in \Omega}{\operatorname{argmax}} \quad p(y = w_j | \underline{x})$$

$$\underline{\Omega} = [w_0 \quad w_1 \quad w_2 \quad \dots \quad w_{K-1}]$$

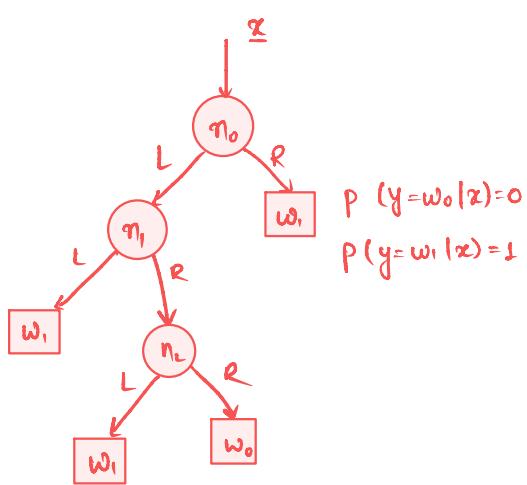
$\underline{x}^{(0)}$ Temp	$\underline{x}^{(1)}$ Humidity	$y^{(0)}$ Comfort	$y^{(1)}$ N. Comfort	
$\underline{x}_0^T$ $35^\circ C$	$T_0 \% \text{ rh}$	0	1	$\underline{y}_0 = w_1$
$\underline{x}_1^T$ $24^\circ C$	$40 \% \text{ rh}$	1	0	$\underline{y}_1 = w_0$
$\underline{x}_2^T$ $22^\circ C$	$60 \% \text{ rh}$	1	0	$\underline{y}_2 = w_0$
$\underline{x}_3^T$ $49^\circ C$	$10 \% \text{ rh}$	0	1	$\underline{y}_3 = w_1$
$\underline{x}_4^T$ $35^\circ C$	$30 \% \text{ rh}$	0	1	$\underline{y}_4 = w_1$
$\underline{x}_5^T$ $32^\circ C$	$50 \% \text{ rh}$	0	1	$\underline{y}_5 = w_1$
$\underline{x}_6^T$ $27^\circ C$	$40 \% \text{ rh}$	1	0	$\underline{y}_6 = w_0$
$\underline{x}_7^T$ $21^\circ C$	$20 \% \text{ rh}$	0	1	$\underline{y}_7 = w_1$

## Decision Trees :-



$o - w_0$

$o - w_1$



$$\begin{aligned} P(y=w_0|x) &= 0 \\ P(y=w_1|x) &= 1 \end{aligned}$$

$n_0$  : if  $x^{(0)} > v(\underline{\theta}_0)$

    Goto L-child

else Goto R-child

$n_1$  : if  $x^{(1)} > v(\underline{\theta}_1)$

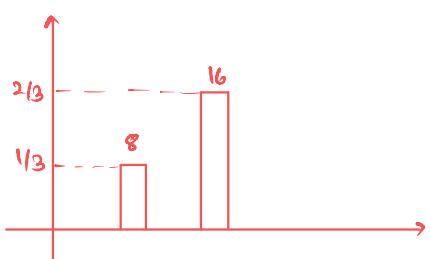
    Goto L-Child

else Goto R-child

$n_2$  : if  $x^{(0)} > v(\underline{\theta}_2)$

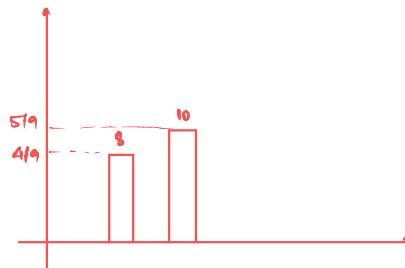
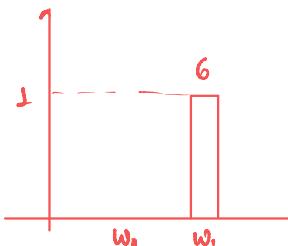
    Goto L-child

else Goto R-child



$$P(y=w_0|x) = -\frac{|y_j=w_0, n_k|}{|n_k|}$$

$$\begin{aligned} H(n_0) &= -\frac{1}{3} \log_2 \left(\frac{2}{3}\right) - \frac{2}{3} \log_2 \left(\frac{2}{3}\right) \\ &= 0.9183 \end{aligned}$$

$n_0 \rightarrow L$  $n_0 \rightarrow R$ 

$$\begin{aligned} H(n_0 \rightarrow L) &= -\frac{4}{9} \log_2 \left(\frac{4}{9}\right) - \frac{5}{9} \log_2 \left(\frac{5}{9}\right) \\ &= 0.9911 \\ &= H(n) \end{aligned}$$

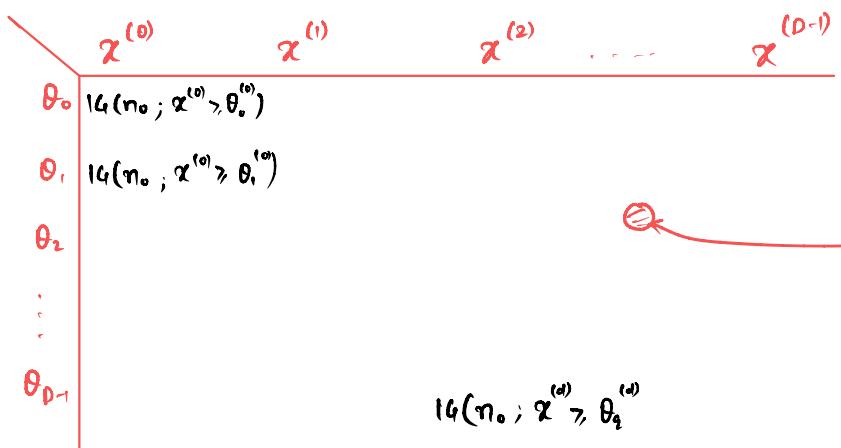
$$H(n_0 \rightarrow R) = -0 \log_2 0 - 1 \log_2 1 = 0$$

$$\begin{aligned} \text{Information Gain (IG)} &= H(n_0) - (P(n_0 \rightarrow L) H(n_0 \rightarrow L) + P(n_0 \rightarrow R) H(n_0 \rightarrow R)) \\ &= 0.9183 - \left( \frac{18}{24} \cdot 0.9911 + \frac{6}{24} \cdot 0 \right) \\ &= 0.2576 \text{ bits} \end{aligned}$$

$$(v, \underline{\theta}_0) = \underset{\sqrt{v(\underline{\theta}_j)} \in S(n)}{\operatorname{argmax}} \left\{ IG(v(\underline{\theta}_j)) \right\}$$

Given a node  $n_k$

$$X(n_k) = \{(x_0, y_0), (x_1, y_1), \dots\}$$



- Minimum number of samples to split
  - ↳ Minimum Leaf Criterion

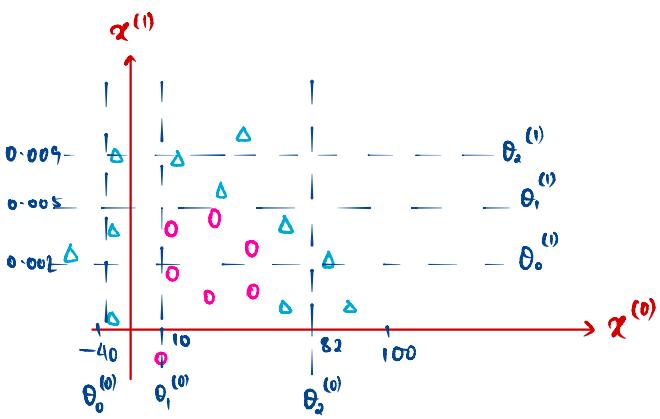
$$p(y = w_0 | \underline{x}) = \frac{|\{y = w_0; n_k\}|}{|\underline{x}_k|}$$

Class Decision Trees :-

$$\underline{x} = [\underline{x}_0, \underline{x}_1, \dots, \underline{x}_n, \dots, \underline{x}_{N-1}] \in \mathbb{R}^{D \times N}$$

$$\underline{y} = [y_0, y_1, \dots, y_n, \dots, y_{N-1}] \in \Sigma^N$$

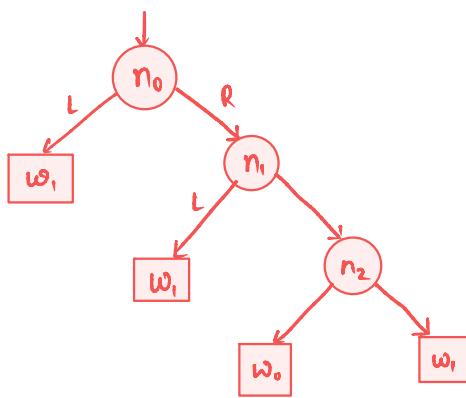
$$\Sigma = \{w_0, w_1\}$$



Axis aligned split -

(a) which axis to split  $x^{(0)}$  or  $x^{(1)}$

(b) what is the split value  $\underline{\theta}_0$



$n_0$ : if  $x^{(0)} > \theta_1^{(0)}$   
 $\theta_1 = 0.0005$   
 $n_0 \rightarrow L$

else  $n_0 \rightarrow R$

$n_1$ : if  $x^{(0)} > \theta_1$   
 $\theta_1 = 68$   
 $n_1 \rightarrow L$

else  $n_1 \rightarrow R$

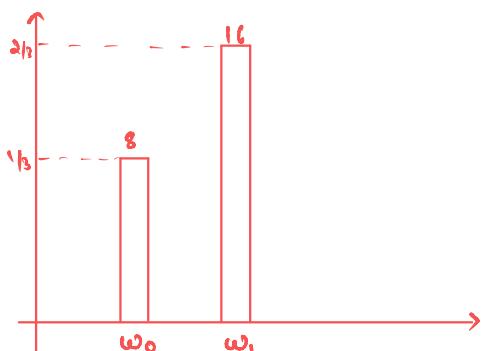
$n_2$ : if  $x^{(0)} > \theta_2$   
 $\theta_2 = 0$

$n_2 \rightarrow L$   
 else  $n_2 \rightarrow R$

	$x^{(0)}$	$x^{(1)}$
$\theta_1$	0.05	0.03
$\theta_2$	0.017	0.25 $\rightarrow \text{max}$
$\theta_3$	0.07	0.05

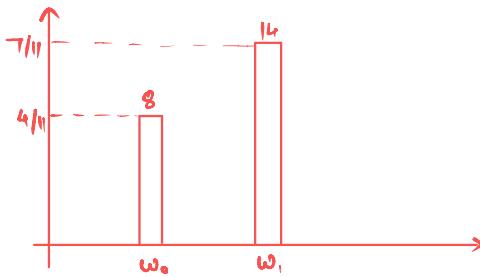
IG ( $x^{(0)}, \theta_i^{(d)}$ )

Splitting with  $\theta_0^{(0)}$  and  $x^0$



$$H(n_0) = -\frac{1}{3} \log_2 (1/3) - \frac{2}{3} \log_2 (2/3)$$

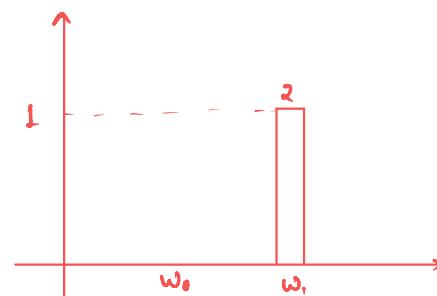
$$= 0.9183 \text{ bits}$$



$n_0 \rightarrow L$   
if  $\chi^{(0)} \geq \theta_0^{(0)}$

$$H(n_0 \rightarrow L) = -\frac{4}{11} \log_2 \left(\frac{4}{11}\right) - \frac{7}{11} \log_2 \left(\frac{7}{11}\right)$$

$$= 0.9456 \text{ bits}$$



$n_0 \rightarrow R$   
if  $\chi^{(0)} < \theta_1^{(0)}$

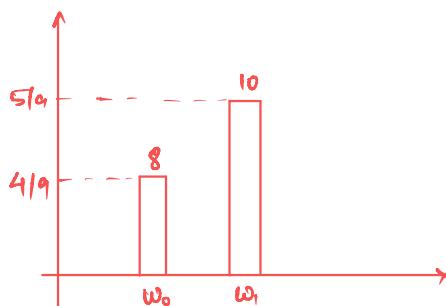
$$H(n_0 \rightarrow R) = H(n_0) - (p(n_0 \rightarrow L) H(n_0 \rightarrow L)$$

$$+ p(n_0 \rightarrow R) H(n_0 \rightarrow R))$$

$$= 0.9183 - \left( \frac{22}{32} (0.9456) + \frac{1}{12} (0) \right)$$

$$= 0.05 \text{ bits}$$

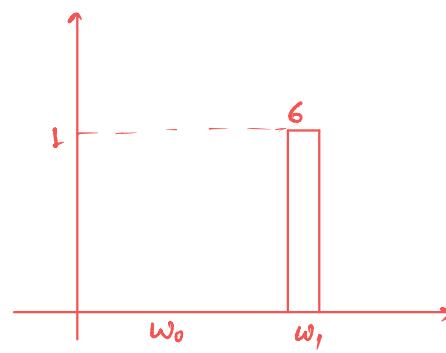
Splitting with  $\theta_1^{(0)}$  on  $\chi^{(0)}$



$n_0 \rightarrow L$   
if  $\chi^{(0)} \geq \theta_1^{(0)}$

$$H(n_0 \rightarrow L) = 0.9911$$

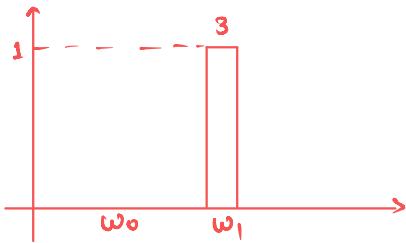
$$\text{IG } (\chi^{(0)}, \theta_1^{(0)}) = 0.17 \text{ bits}$$



$n_0 \rightarrow R$   
if  $\chi^{(0)} < \theta_1^{(0)}$

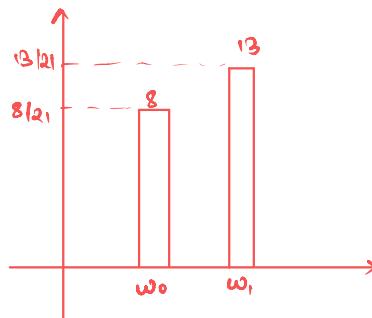
$$H(n_0 \rightarrow R) = 0$$

## Splitting with $\theta_2^{(0)}$ on $x^{(0)}$



$n_0 \rightarrow L$   
if  $x^{(0)} > \theta_2^{(0)}$

$$H(n_0 \rightarrow L) = 0$$

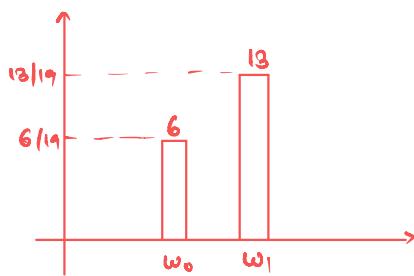


$n_0 \rightarrow R$   
if  $x^{(0)} < \theta_2^{(0)}$

$$\begin{aligned} H(n_0 \rightarrow R) &= -\frac{8}{21} \log_2 \left(\frac{8}{21}\right) - \frac{13}{21} \log_2 \left(\frac{13}{21}\right) \\ &= 0.9587 \text{ bits} \end{aligned}$$

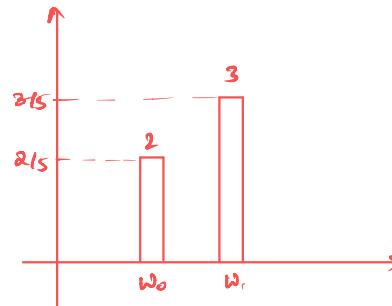
$$\begin{aligned} H(x^{(0)}, \theta_2^{(0)}) &= 0.9183 - \left( \frac{3}{24}(0) + \frac{21}{24}(0.2286) \right) \\ &= 0.074 \end{aligned}$$

## Splitting with $\theta_0^{(1)}$ on $x^{(1)}$



$n_0 \rightarrow L$   
if  $x^{(1)} > \theta_0^{(1)}$

$$H(n_0 \rightarrow L) = -\frac{6}{19} \log_2 \left(\frac{6}{19}\right) - \frac{13}{19} \log_2 \left(\frac{13}{19}\right) = 0.899$$



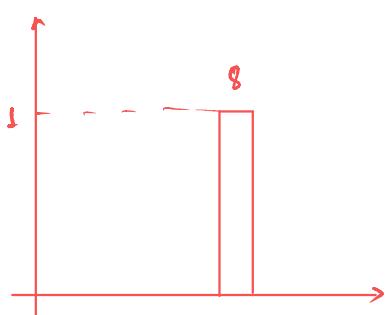
$n_0 \rightarrow R$   
if  $x^{(1)} < \theta_0^{(1)}$

$$H(n_0 \rightarrow R) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.971$$

$$H(x^{(1)}, \theta_0^{(1)}) = 0.9183 - \left[ \frac{19}{24} (0.899) - \frac{5}{24} (0.971) \right]$$

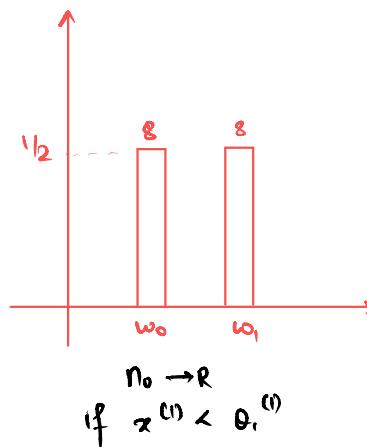
$$= 0.0043$$

Splitting with  $\theta_1^{(1)}$  on  $x^{(1)}$



$$\begin{aligned} n_0 \rightarrow L \\ \text{if } x^{(1)} > \theta_1^{(1)} \end{aligned}$$

$$H(n_0 \rightarrow L) = 0$$

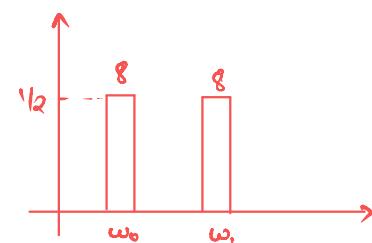


$$\begin{aligned} n_0 \rightarrow R \\ \text{if } x^{(1)} < \theta_1^{(1)} \end{aligned}$$

$$\begin{aligned} H(n_0 \rightarrow R) &= -\frac{1}{2} \log_2 (1/2) - \frac{1}{2} \log_2 (1/2) \\ &= 1 \end{aligned}$$

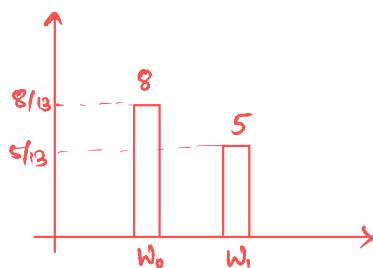
$$\begin{aligned} H(x^{(1)}, \theta_1^{(1)}) &= 0.9183 - \left( \frac{8}{24} (0) + \frac{16}{24} (1) \right) \\ &= 0.252 \end{aligned}$$

Splitting with  $\theta_2^{(1)}$  on  $x^{(1)}$



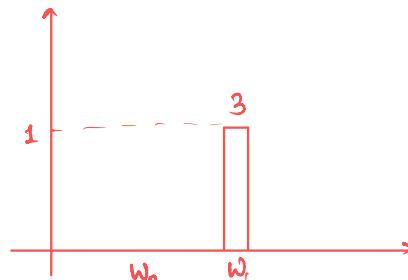
$$\begin{aligned} H(n_0) &= -\frac{1}{2} \log_2 (1/2) - \frac{1}{2} \log_2 (1/2) \\ &= 1 \text{ bit} \end{aligned}$$

## Splitting with $\theta_0^{(0)}$ and $\alpha^{(0)}$



$n_i \rightarrow L$

if  $\alpha^{(0)} > \theta_1^{(0)}$



$n_i \rightarrow R$

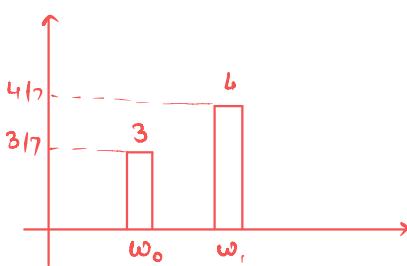
if  $\alpha^{(0)} < \theta_1^{(0)}$

$$H(n_i \rightarrow L) = 0.9612$$

$$H(n_i \rightarrow R) = 0$$

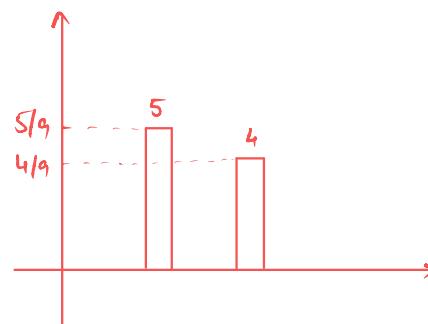
$$IG = 0.21$$

## Splitting with $\theta_0^{(1)}$ on $\alpha^{(0)}$



$$H(n_i \rightarrow L) = -\frac{3}{7} \log_2 (3/7) - \frac{4}{7} \log_2 (4/7)$$

$$= 0.985$$



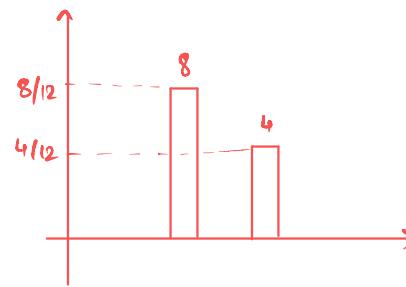
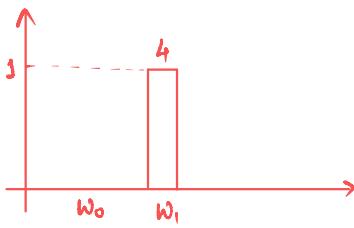
$$H(n_i \rightarrow R) = -\frac{5}{9} \log_2 (5/9) - \frac{4}{9} \log_2 (4/9)$$

$$= 0.991$$

$$IG = 3 - \left( \frac{1}{16} (0.985) + \frac{9}{16} (0.991) \right)$$

$$= 0.012$$

## Splitting with $\theta_0$ on $x^{(0)}$

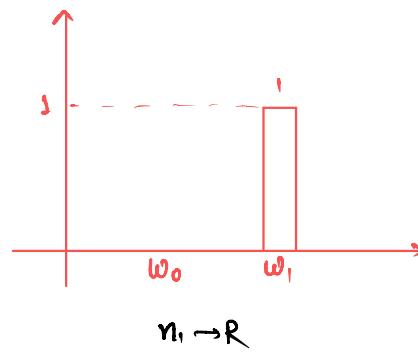
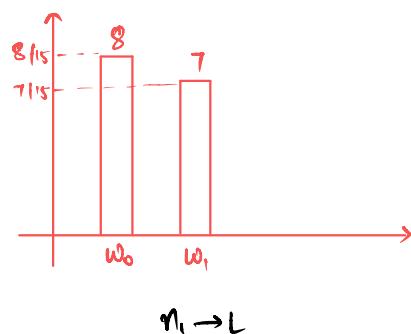


$$H(n_i \rightarrow L) = 0$$

$$\begin{aligned} H(n_i \rightarrow R) &= -\frac{8}{12} \log_2(8/12) - \frac{4}{12} \log_2(4/12) \\ &= 0.9183 \end{aligned}$$

$$\begin{aligned} IG &= 1 - \left( 0 + \frac{12}{16} (0.9183) \right) \\ &= 0.31 \end{aligned}$$

## Splitting with $\theta_0^{(1)}$ on $x^{(1)}$



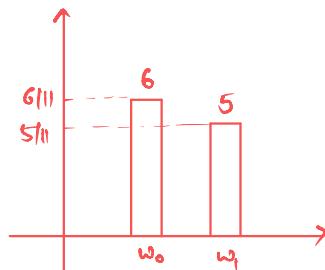
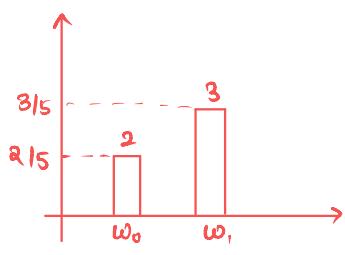
$$\begin{aligned} H(n_i \rightarrow L) &= -\frac{8}{15} \log_2(8/15) - \frac{7}{15} \log_2(7/15) \\ &= 0.997 \end{aligned}$$

$$H(n_i \rightarrow R) = 0$$

$$IG = 1 - \left( \frac{15}{16} (0.997) + 0 \right)$$

$$= 0.065$$

Splitting with  $\theta_1^{(0)}$  on  $x^{(1)}$



$$H(n_r \rightarrow L) = -\frac{2}{5} \log_2 (2/5) - \frac{3}{5} \log_2 (3/5)$$

$$= 0.9709$$

$$H(n_r \rightarrow R) = -\frac{6}{11} \log_2 (6/11) - \frac{5}{11} \log_2 (5/11)$$

$$= 0.9940$$

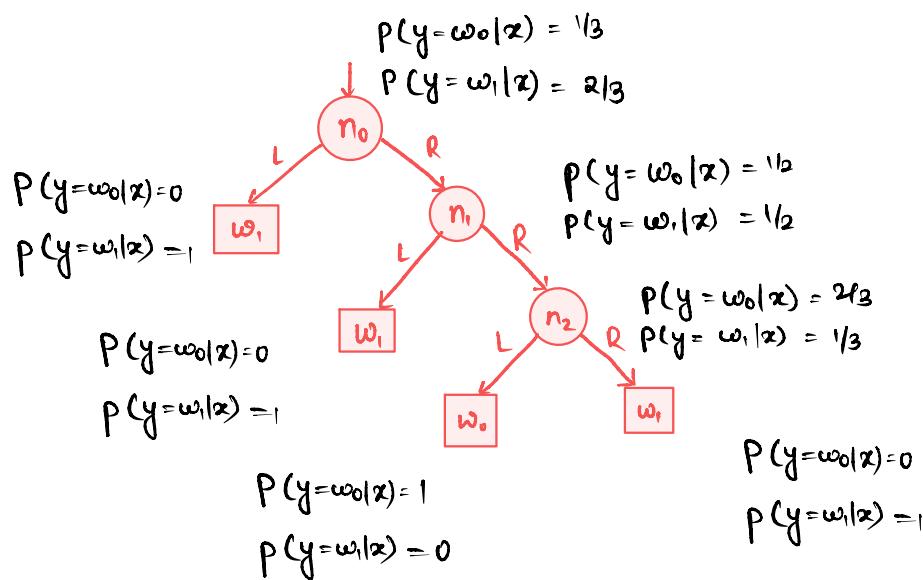
$$IG = 0$$

	$x^{(0)}$	$x^{(1)}$
$\theta_1$	0.21	0.06
$\theta_2$	0.01	0
$\theta_3$	0.31	0.06

Min leaf = 15

Min. no of samples required such that node doesn't become a leaf node.

$$N_0 = 24$$



### Inference With a Decision Tree:

$$\underline{x} = (83, 0.0068)$$

$x^{(0)} \swarrow \quad \searrow x^{(1)}$

$P(y = w_0|x) = 0$   
 $P(y = w_1|x) = 1$

$P(y = w_0|x) = 0$   
 $P(y = w_1|x) = 1$

$$\underline{x} = (5, 0.0023)$$

$x^{(0)} \swarrow \quad \searrow x^{(1)}$

$P(y = w_0|x) = 2/3$   
 $P(y = w_1|x) = 1/3$

$P(y = w_0|x) = 2/3$   
 $P(y = w_1|x) = 1/3$

### Training - grow with the tree (DT):

# Compute  $\alpha$  No. of nodes in a DT

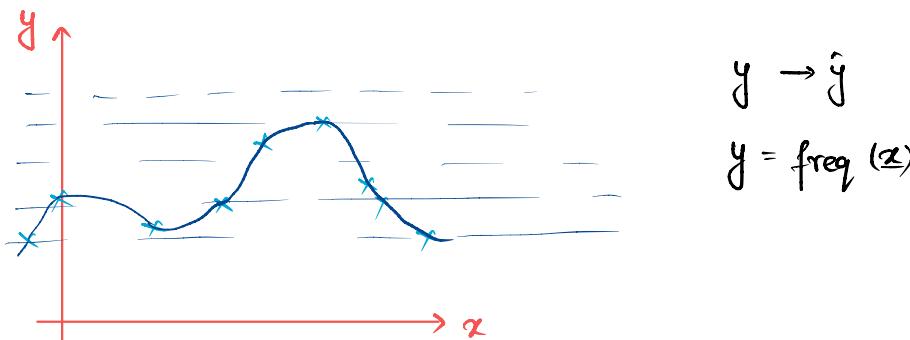
# Space  $\propto$  No. of nodes in a DT.

## Inference - Infer with a DT

# Compute  $\alpha$  Height of DT

# Space  $\alpha$  Nodes in a DT

## Regression problems :-



$$\underline{x} = [\underline{x}_0, \underline{x}_1, \dots, \underline{x}_n, \dots, \underline{x}_{N-1}] \in \mathbb{R}^{D \times N}$$

$$\underline{y} = [\underline{y}_0, \underline{y}_1, \dots, \underline{y}_n, \dots, \underline{y}_{N-1}] \in \mathbb{R}^{B \times N}$$

$$\hat{\underline{y}} = [\hat{\underline{y}}_0, \hat{\underline{y}}_1, \dots, \hat{\underline{y}}_n, \dots, \hat{\underline{y}}_{N-1}] \in \mathcal{D}^N$$

$$\mathcal{L} = \{w_0, w_1, \dots, w_j\}$$

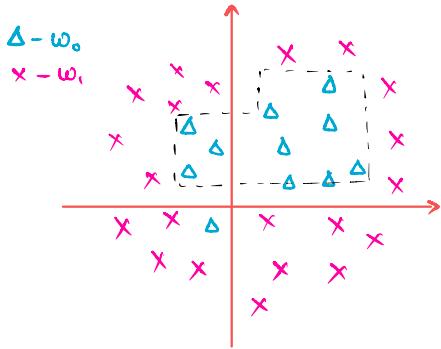
$$w_j = [y_0^{(0)} - y_i^{(0)}, y_1^{(1)} - y_j^{(1)}, \dots]$$

## Random Forest :-

- i) Ensemble ML
- ii) Classifier Combination
- iii) Distribution Consensus

$$\underline{X} = \{\underline{x}_0, \underline{x}_1, \dots, \underline{x}_n, \dots, \underline{x}_{N-1}\} \quad \underline{x}_n \in \mathbb{R}^{Dx1}$$

$$\underline{Y} = \{\underline{y}_0, \underline{y}_1, \dots, \underline{y}_n, \dots, \underline{y}_{N-1}\} \quad \underline{y}_n \in \Sigma = \{\omega_0, \omega_1, \dots, \omega_{k-1}\}$$



$$P(y = \omega_0 | \underline{x})$$

$$\hat{y} = \underset{\forall \omega_k \in \Sigma}{\operatorname{argmax}} \{ P(y = \omega_k | \underline{x}) \}$$

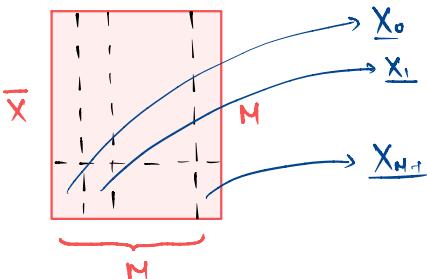
$$\underline{X} = \underline{x}_0 \cup \underline{x}_1 \cup \underline{x}_2 \cup \dots \cup \underline{x}_{m-1}$$

such that  $\underline{x}_m \cap \underline{x}_0 = \emptyset$ ,  $m \neq 0$



$\underline{y}_m$

$$P(y = \omega_k) \Big|_{\forall y \in \underline{y}_m} \approx P(y = \omega_k) \Big|_{\forall y \in \underline{y}_0} \approx P(y = \omega_k) \Big|_{\forall y \in \underline{y}}$$



$$T_m = DT(\underline{x}_m, \underline{y}_m)$$

$$\{\underline{x}_m, y_m\} \xrightarrow[\text{DT}]{\text{Train}} T_m$$

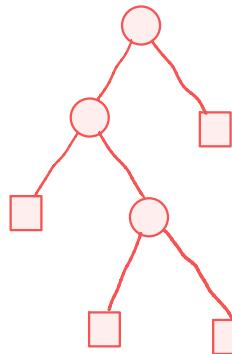
$$\underline{x} \xrightarrow{\text{Inference}} T_m(\underline{x}) = \left[ P(y=w_k | \underline{x})_m \right]_{w_k \in \Sigma}$$

$$P(y=w_0 | \underline{x})_m$$

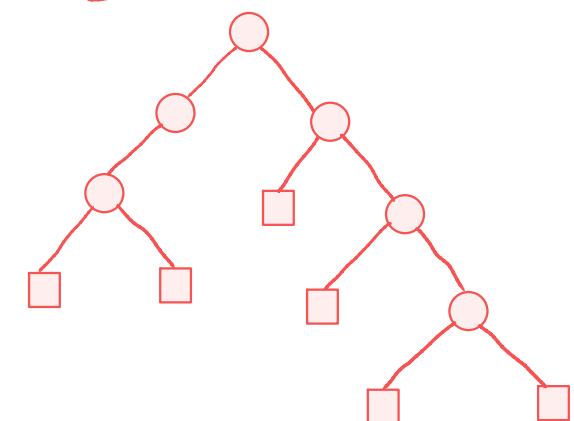
$$P(y=w_1 | \underline{x})_m$$

$$P_{RF}(y=w_k | \underline{x}) = \frac{1}{M} \sum_{m=0}^{M-1} P(y=w_k | \underline{x})_m$$

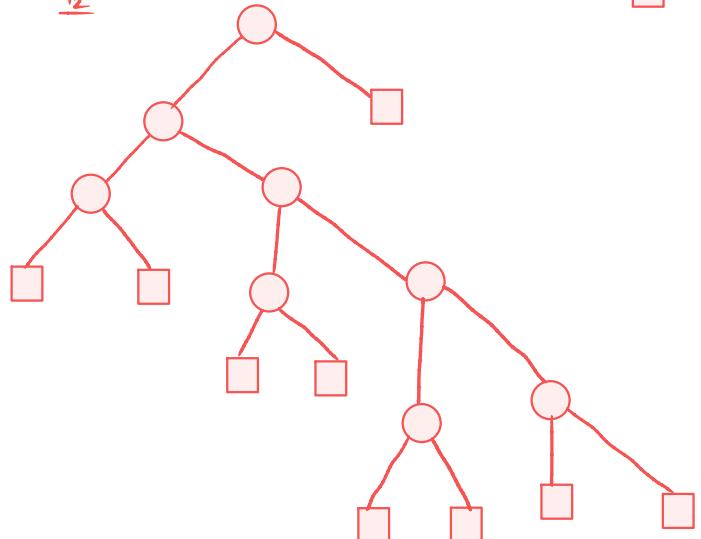
$T_0$



$T_1$



$T_2$



Bootstrap Sampling Aggregation  $\rightarrow$  Bagging

## Mixture Model

$$P(y=w_k | \underline{x}) = \frac{P(\underline{x} | y=w_k) \cdot P(y=w_k)}{P(\underline{x})}$$

likelihood  
 ↓  
 Posterior

Prior  
 ↓  
 Evidence

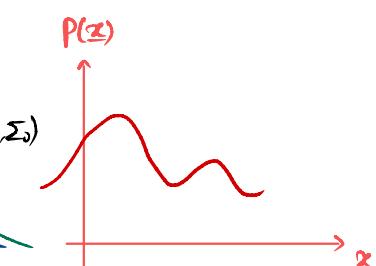
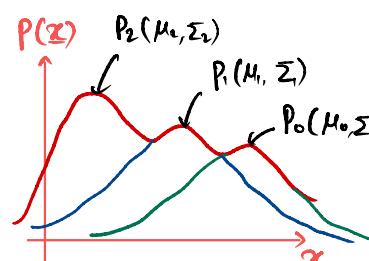
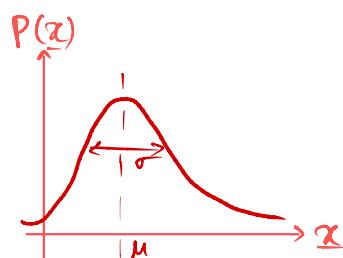
$$P(\underline{x}) = \sum_{w_k \in \mathcal{R}} P(\underline{x} | y=w_k) \cdot p(y=w_k)$$

$$\hat{y} = \underset{\sqrt{w_k \in \mathcal{R}}}{\operatorname{argmax}} \{P(y=w_k | \underline{x})\}$$

$P(\underline{x} | y=w_k) \sim$  Multivariate Normal Distribution

$$\sim \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\underline{x}-\underline{\mu})^\top \Sigma^{-1} (\underline{x}-\underline{\mu})\right)$$

Assumption is that  $p(\underline{x} | y=w_k)$  is unimodal

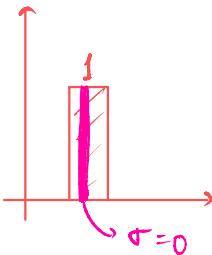


$$P(\underline{x}) = \sum_{j=0}^2 P_j \mathcal{N}(\mu_j, \Sigma_j)$$

$$\int_{\mathbb{R}} P(\underline{x}) d\underline{x} = 1, \quad \sum_{j=0}^2 P_j = 1$$

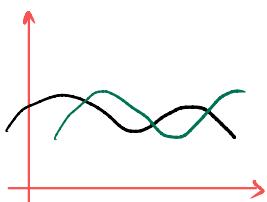
- You have only samples  $\underline{x}, y = (\underline{x}, y)$

You have to model  $P(\underline{x} | y=w_k)$



$$\begin{aligned}\underline{\mu}_0, \underline{\Sigma}_0, p_0 \\ \underline{\mu}_1, \underline{\Sigma}_1, p_1 \\ \underline{\mu}_2, \underline{\Sigma}_2, p_2\end{aligned}$$

$\rightarrow \Theta$  (Parameter Vector of the GMM)



$$B, \underline{x} \in \mathbb{R}^D$$

Each dim has 10 bins if  $D=10$  then I have  $10^{10}$  bins.

$$P(\Theta) \rightarrow Q(\Theta)$$

property denoting the goodness of the GMM.

$$\frac{\partial Q(\Theta)}{\partial \Theta} = 0$$

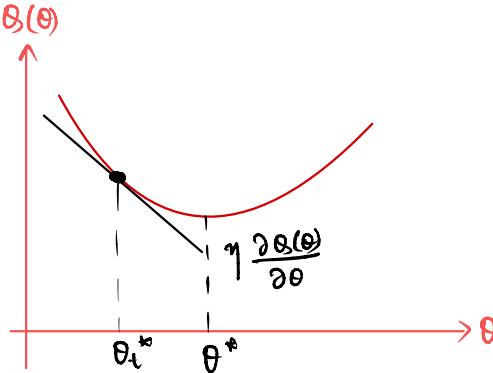
Gaussian Mixture Model (GMM) Solve Expectation-Maximization (EM) algorithm.

E Step: Random Guess of  $\Theta$ .

$$\underline{\theta}_t(\underline{\theta}, \underline{\theta}_t) = E \left[ - \sum_n \ln (P(\underline{x}_n; \underline{\theta} | \underline{\theta}_t)) \right]$$

$$\underline{x} = [\underline{x}_0, \underline{x}_1, \dots, \underline{x}_n, \dots, \underline{x}_{N-1}]$$

$$\underline{M Step: } \frac{\partial \underline{\theta}(\underline{\theta}, \underline{\theta}_t)}{\partial \underline{\theta}} = 0 \quad \left| \begin{array}{l} \underline{\theta}_{t+1} = \underline{\theta}_t - \eta \frac{\partial \underline{\theta}(\underline{\theta}, \underline{\theta}_t)}{\partial \underline{\theta}} \\ \text{Learning Rate / Gradient Control} \end{array} \right.$$



$$X = [x_0, x_1, \dots, x_n, \dots, x_{N-1}] \in \mathbb{R}^{D \times N}$$

$$Y = [y_0, y_1, \dots, y_n, \dots, y_{N-1}] \in \mathbb{R}^N$$

$$\Omega = \{\omega_0, \omega_1, \dots, \omega_{k-1}\}$$

$$S-I: X|_{\omega_0}, X|_{\omega_1}, \dots, X|_{\omega_{k-1}}$$

$$= [x_0 | y_0 = \omega_0, x_1 | y_1 = \omega_1, \dots, x_n | y_n = \omega_n]$$

S-II: Estimate  $p(x|y=\omega_k)$  using GMM.

Select the no. of components:  $j=3$

Set initial  $\theta_{t=0}|_{\omega_k} = \{\underline{\mu}_0, \underline{\Sigma}_0, p_0, \underline{\mu}_1, \underline{\Sigma}_1, p_1, \underline{\mu}_2, \underline{\Sigma}_2, p_2\}$

Solve EM algorithm and get  $\theta^*|_{\omega_k}$

S-III:  $p(x|y=\omega_k) \approx \text{GMM}(\theta^*|_{y=\omega_k})$

$$p(y=\omega_k), p(x) = \sum_{\forall \omega_k \in \Omega} p(x|y=\omega_k) \cdot p(y=\omega_k)$$

$$\hat{y} = \underset{\forall \omega_k \in \Omega}{\operatorname{argmax}} (p(y=\omega_k|x))$$

$$\text{where } p(y=\omega_k|x) = \frac{p(x|y=\omega_k) \cdot p(y=\omega_k)}{p(x)}$$

There is a crucial assumption for j.

↳ Random guess of  $J_s = 0$ .

E-Step: for  $J_s$

| E - Step for  $\Theta(t)$

| M - Step for  $\Theta_{t+1}$

M-Step for  $J_{s+1}$

Unsupervised Learning for finite Mixture Models

↳ Figueredo and Jain - IEEE PAMI 24(3), 2002

### Feature Selection :-

$$\underline{x} \in \mathbb{R}^{D \times 1}$$

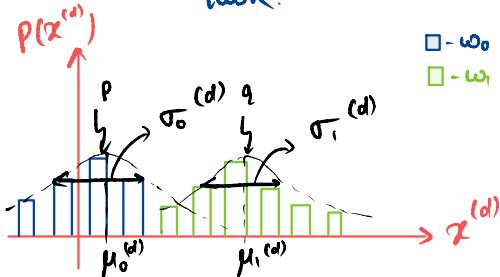
$$y \in \mathbb{Z}^{K \times 1}$$

$$\omega = \{\omega_0, \omega_1, \dots, \omega_{K-1}\}$$

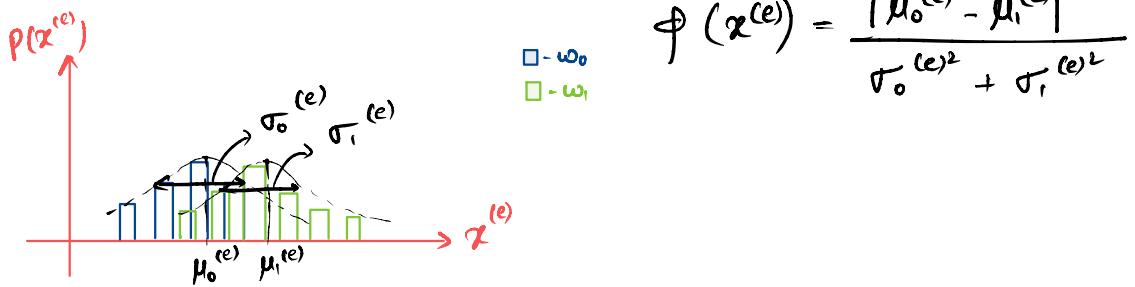
$$\underline{x} = [x^{(0)}, x^{(1)}, \dots, x^{(d)}, \dots, x^{(D-1)}]$$

- Each feature is not equivalently relevant for a classification task.

$q(x^{(d)}) \rightarrow$  Relevance of a feature  $x^{(d)}$  to solve a classification task.



$$\begin{aligned} q(x^{(d)}) &= \frac{|\mu_0^{(d)} - \mu_1^{(d)}|^2}{\sigma_0^{(d)2} + \sigma_1^{(d)2}} \\ &= \frac{\text{Inter Class Variance}}{\text{Intra Class Variance}} \end{aligned}$$



- Relevance of  $x^{(e)}$  is lesser than relevance of  $x^{(d)}$  since  $|\mu_0^{(e)} - \mu_1^{(e)}| < |\mu_0^{(d)} - \mu_1^{(d)}|$ .

### Kulback Leiber Divergence:

$$d_{KL}(P||Q) = \sum_i P_i \log\left(\frac{P_i}{Q_i}\right)$$

$$d_{KL}(P||Q) \neq d_{KL}(Q||P)$$

### Jensen-Shannon Divergence:

$$d_{JS}(P||Q) = \sum_i \left( \frac{P_i}{C_i} \log_2 \left( \frac{P_i}{C_i} \right) + \frac{Q_i}{C_i} \log_2 \left( \frac{Q_i}{C_i} \right) \right); \quad C_i = \frac{P_i + Q_i}{2}$$

### Filter Model of Feature Selection:

- for  $d=0, \dots, (D-1)$

$$\phi(x^{(d)}) = \phi_{JS}(P||Q);$$

$$P \sim h(x^{(d)} | y=\omega_0);$$

$$Q \sim h(x^{(d)} | y=\omega_1);$$

end

histogram

$$\underline{x}^* = [x^{*(0)} \ x^{*(1)} \ x^{*(2)}, \dots, x^{*(d)}, \dots, x^{*(D-1)}]$$

$$\phi(x^{(e)}) = \frac{|\mu_0^{(e)} - \mu_1^{(e)}|^2}{\sigma_0^{(e)2} + \sigma_1^{(e)2}}$$

such that  $\varphi(x^{*(0)}) \geq \varphi(x^{*(1)}) \geq \dots \geq \varphi(x^{*(d)}) \geq \dots \geq \varphi(x^{*(D-1)})$

### Wrapper Model of Feature Selection:

$$\hat{y} = F(\underline{x}), \quad F: \mathbb{R}^{D \times 1} \mapsto \Omega^{k \times 1}$$

$$x^{(0)}, x^{(1)}, \dots, x^{(d)}, \dots, x^{(D-1)}$$

↓ pick any one

$$\tilde{x}^{(j)} \rightarrow F_j(\cdot) \xrightarrow{\text{Accuracy}} a_j$$

$$F_0(x^{(0)}) \rightarrow a_0$$

$$F_1(x^{(1)}) \rightarrow a_1$$

:

$$F_{D-1}(x^{(D-1)}) \rightarrow a_{D-1}$$

$$\left. \begin{array}{l} \\ \\ \end{array} \right\} \xrightarrow{\text{argmax} = x^{(1)}} \tilde{x}_{(0)} = [x^{(1)}]$$

↓ pick any one not present  $\tilde{x}_{(0)}$

$$F_0(x^{(0)}) \rightarrow a_0$$

$$F_1(x^{(1)}) \rightarrow a_1$$

:

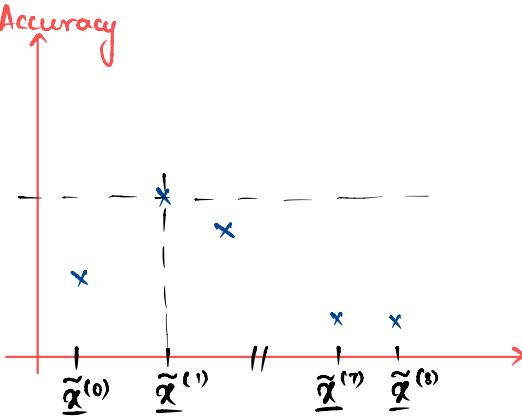
$$F_{D-1}(x^{(D-1)}) \rightarrow a_{D-1}$$

$$\left. \begin{array}{l} \\ \\ \end{array} \right\} \xrightarrow{\text{argmax} = x^{(0)}} \tilde{x}_{(0)} = [x^{(0)}]$$

$$\tilde{x}_{(1)} = [x^{(1)}, x^{(0)}]$$

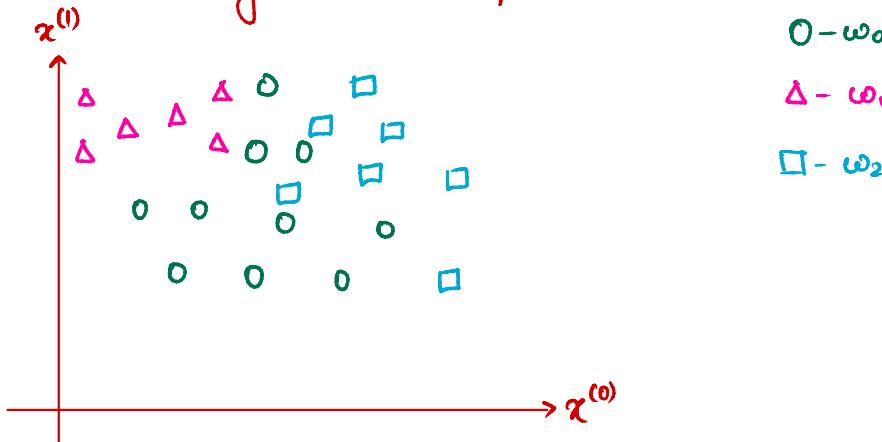
↓  
Accuracy

Optimal set of features relevant to F.



Sequential Forward Search (SFS)

K- Nearest Neighbour Classification (kNN) :-



$$\underline{X} = [\underline{x}_0, \underline{x}_1, \dots, \underline{x}_n, \dots, \underline{x}_{N-1}] \in \mathbb{R}^{D \times N}$$

such that

$$\underline{Y} = [y_0, y_1, \dots, y_n, \dots, y_{N-1}] \in \mathbb{R}^N$$

such that  $\mathcal{L} = \{\omega_0, \omega_1, \omega_2\}$

$$d(\underline{x}_n, \underline{x}_j) = \|\underline{x}_n - \underline{x}_j\|_2$$

• Sort based on  $d(\underline{x}_n, \underline{x}_j)$

$$\tilde{\underline{X}} = [\tilde{\underline{x}}_0, \tilde{\underline{x}}_1, \dots, \tilde{\underline{x}}_n, \dots, \tilde{\underline{x}}_{N-1}] \in \mathbb{R}^{D \times N}$$

such that  $\tilde{\underline{x}}_n \in \mathbb{R}^{D \times 1}$ ,  $d(\tilde{\underline{x}}_0, \underline{x}_j) \leq d(\tilde{\underline{x}}_1, \underline{x}_j) \leq d(\tilde{\underline{x}}_2, \underline{x}_j) \dots \leq d(\tilde{\underline{x}}_{N-1}, \underline{x}_j)$

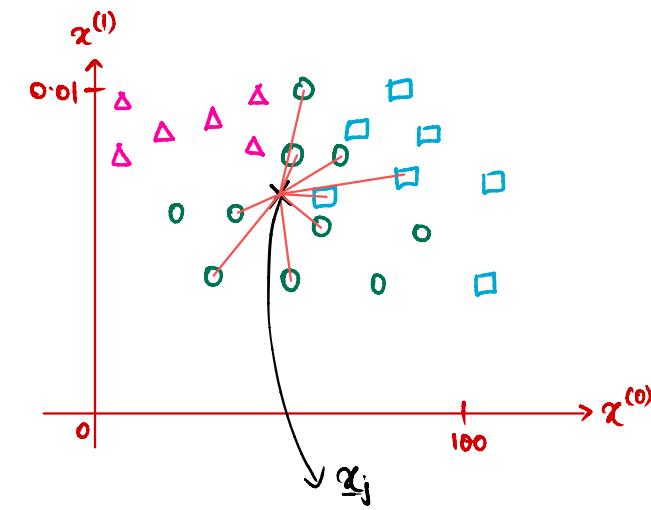
$$\tilde{\underline{Y}} = [\tilde{y}_0, \tilde{y}_1, \dots, \tilde{y}_n, \dots, \tilde{y}_{N-1}] \in \Omega^N$$

such that  $\Omega = \{\omega_0, \omega_1, \omega_2\}$

•  $\tilde{\underline{X}}_k \subset \tilde{\underline{X}}$

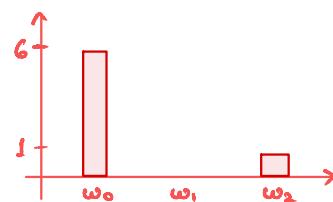
$$\tilde{\underline{X}}_k = [\tilde{\underline{x}}_0, \tilde{\underline{x}}_1, \dots, \tilde{\underline{x}}_{k-1}]$$

max votes  $(\tilde{\underline{Y}}_k = [\tilde{y}_0, \tilde{y}_1, \dots, \tilde{y}_{k-1}]) = \omega_0$



$K=8$

↳ Issue of Selection  
of  $K$



$$\bullet \cos \theta = \frac{|\underline{x}_n \cdot \underline{x}_j|}{\|\underline{x}_n\|_2 \cdot \|\underline{x}_j\|}$$

→ Cosine similarity distance

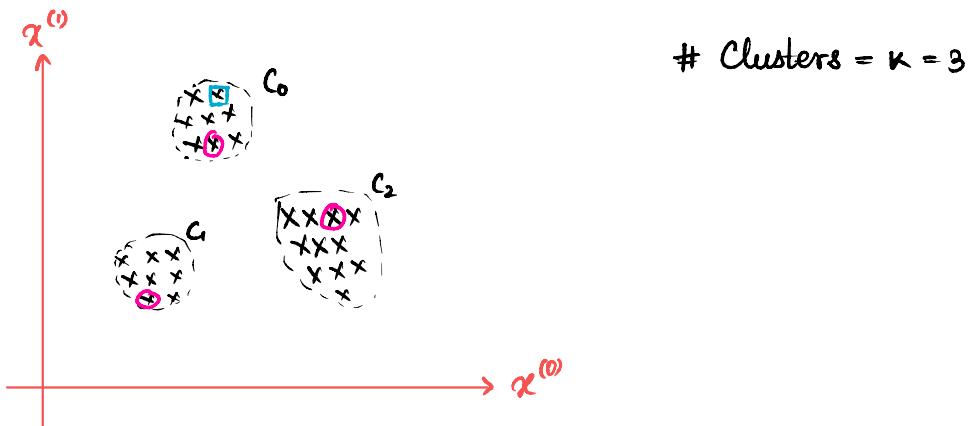
## Clustering :-

- Unsupervised method

$$\underline{X} = [\underline{x}_0, \underline{x}_1, \dots, \underline{x}_n, \dots, \underline{x}_{N-1}] \in \mathbb{R}^{D \times N}$$

such that  $\underline{x}_n \in \mathbb{R}^{D \times 1}$

Y = Not Available



$$\tilde{\underline{Y}} = [y_0, y_1, \dots, y_n, \dots, y_{N-1}] \in \mathbb{C}^N$$

$$C = \{C_0, C_1, \dots, C_n, \dots C_{K-1}\}$$

$\underline{\mu}^{(t)}$  - State of cluster centroid at  $t^{\text{th}}$  iteration.

$$\underline{\mu}^{(t=0)} = [\underline{\mu}_0^{(t=0)}, \underline{\mu}_1^{(t=0)}, \underline{\mu}_2^{(t=0)}] \in \mathbb{R}^{D \times K} \sim \underline{X}$$

such that  $\mu_k^{(t)} \in \mathbb{R}^{D \times 1}$

$$\underline{x} = \{\underline{x}\} - \{\underline{\mu}^{(t=0)}\}$$

for  $t=1 : N-K-1$

$$\underline{x}^{(t)} \xrightarrow[\text{rep}]{w/o} \{\tilde{\underline{x}}\}$$

for  $K=0 : (K-1)$

$$d_k^{(t)} = \|\mu_k^{(t)} - \underline{x}^{(t)}\|$$

end

$$\tilde{y}^{(t)} = \operatorname{argmin} \{d_k^{(t)}\}$$

$$\mu_{\tilde{y}^{(t)}}^{(t+1)} = \frac{1}{|\tilde{y}^{(t)}|} \left( (|\tilde{y}^{(t)}|-1) \mu_{\tilde{y}^{(t)}}^{(t)} + \underline{x}^{(t)} \right)$$

$$|C_0|_{t=0} = 1 \quad |C_0|_{t=1} = 2 \quad |C_1|_{t=1} = 1 \quad |C_2|_{t=1} = 1$$

## • Hyperparameters

- $K$  - No. of clusters
- Initial Seed

## • K-Means - Clustering

## • Iterative Agglomerative Clustering.

$$\theta = 0$$

$$\underline{\mu}^{(\theta=0)}_{\text{k-pls}} \sim \{ \underline{x} \}$$

do

$$\underline{\mu}^{(t=0)} = \underline{\mu}^{(\theta=0)}$$

Call : K-Means clustering

$$\underline{\mu}^{(\theta+1)} = \underline{\mu}^{(t=N-k)}, \theta = \theta + 1$$

$$\text{while } \| \underline{\mu}^{(\theta)} - \underline{\mu}^{(\theta+1)} \| > \epsilon$$

→ Hyperparameter / Tolerance

• E-Step:

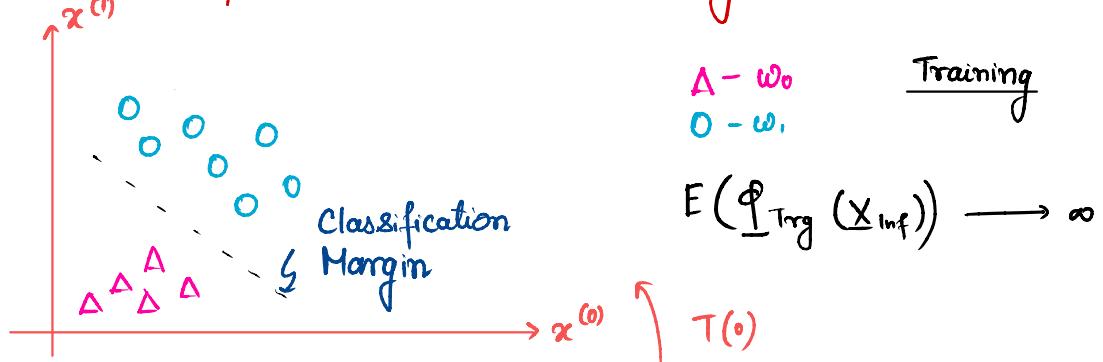
Guess K

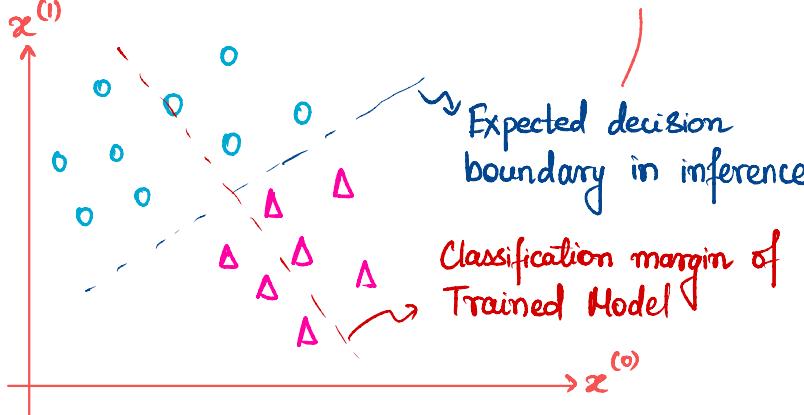
M-Step:

Solve GOF with K

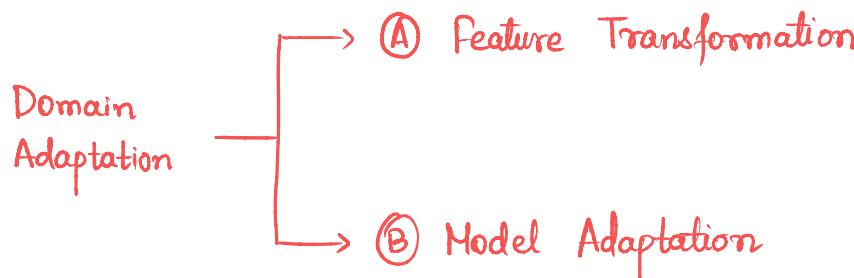
• Cluster Compactness Index (CCI) =  $\frac{\text{Intra Cluster Index}}{\text{Inter Cluster Index}}$

Domain Adaptation and Transfer Learning:-





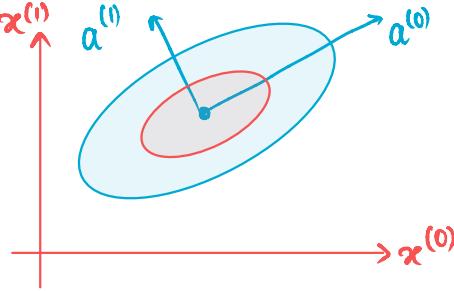
$$\underline{x}_i \in \underline{X}_{\text{Inf}}$$



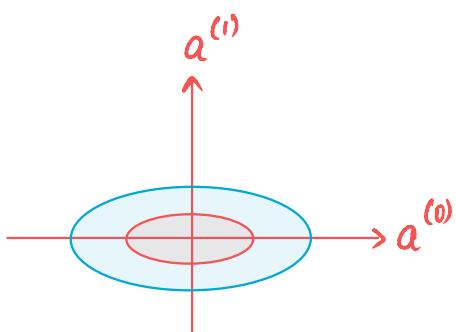
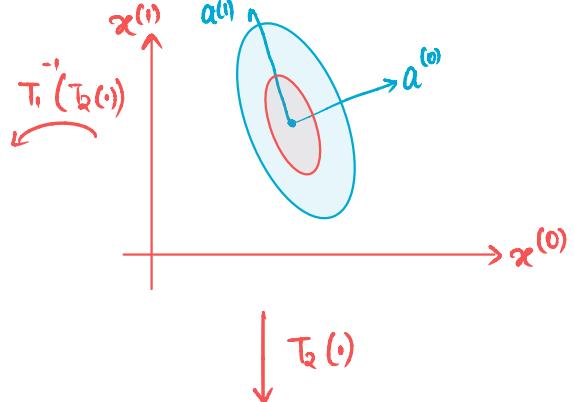
$\Phi$ : Mathematical Representation of the Classifier model.

$\Phi_{\text{Trg}}$ :  $\Phi$  trained on  $\underline{X}_{\text{Trg}}$

$\Phi_{\text{Inf}}$ :  $\Phi$  trained on  $\underline{X}_{\text{Inf}}$



$\downarrow T_1(\cdot)$  PCA projection



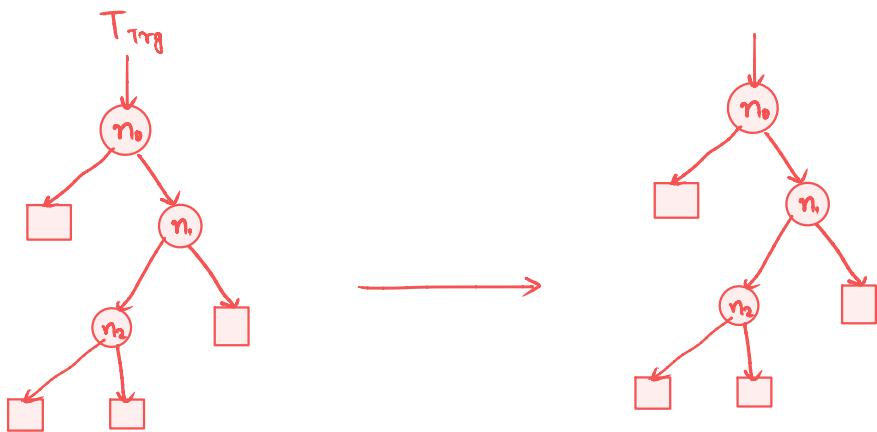
$$\underline{X} \neq \underline{X}_{\text{Inf}}$$

$$\underline{\Phi}_{\text{Trg}} \neq \underline{\Phi}_{\text{Inf}}$$

$\downarrow \text{TL}$

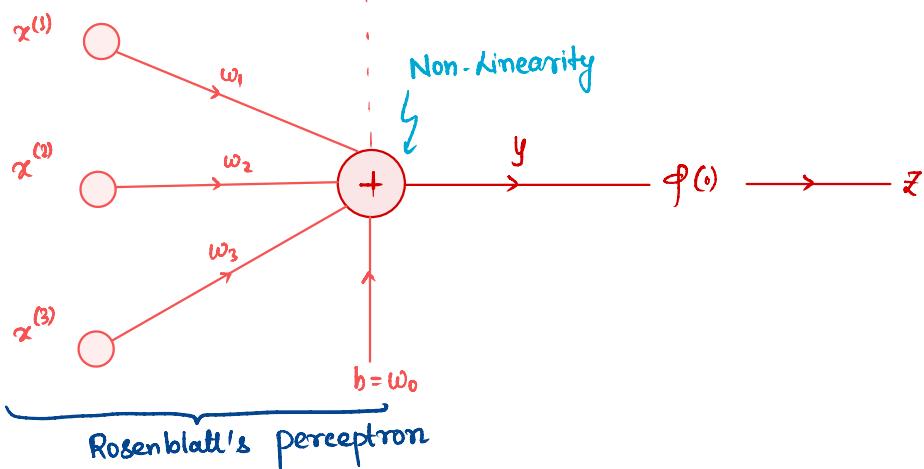
$$\underline{\Phi}_{\text{Trg} \rightarrow \text{Inf}} \simeq \underline{\Phi}_{\text{Inf}} = E(\underline{\Phi}_{\text{Trg} \rightarrow \text{Inf}}(\underline{X}_{\text{Inf}})) = E(\underline{\Phi}_{\text{Inf}}(\underline{X}_{\text{Inf}})) + \epsilon$$

such that  $\epsilon \rightarrow 0$ .



$$n_0 := f_{n_0}(x) = \begin{cases} 1 & \text{if } x^{(d)} > \theta_{n_0}^{(d)} \\ 0 & \text{otherwise} \end{cases}$$

## Neural Networks:-



$$\underline{x} = [x^{(0)} \quad x^{(1)} \quad \dots \quad x^{(d)} \quad \dots \quad x^{(0-1)}] \in \mathbb{R}^D \quad (\text{Not for NN})$$

$$\underline{x} = [x^{(1)} \quad x^{(2)} \quad \dots \quad x^{(d)} \quad \dots \quad x^{(0)}] \in \mathbb{R}^D \quad (\text{For NN})$$

$$y = \omega_1 x^{(1)} + \omega_2 x^{(2)} + \omega_3 x^{(3)} + b$$

*Bias*

$$= [\omega_0 \quad \omega_1 \quad \omega_2 \quad \omega_3] \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix} + b$$

↓  
 Weights  
 ↓  
 Tunable / Variable factors

$$\therefore b = \omega_0,$$

$$y = [\omega_0 \quad \omega_1 \quad \omega_2 \quad \omega_3] \begin{bmatrix} 1 \\ x^{(1)} \\ x^{(2)} \\ x^{(3)} \end{bmatrix}$$

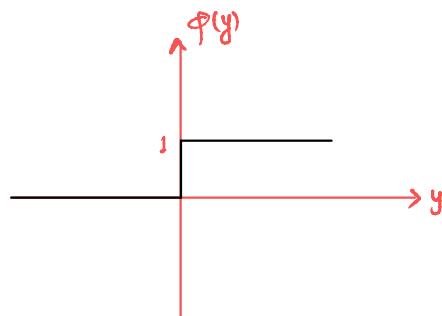
$$= \underline{W} \underline{X} + b$$

$$= \underline{W^*} \underline{X^*}$$

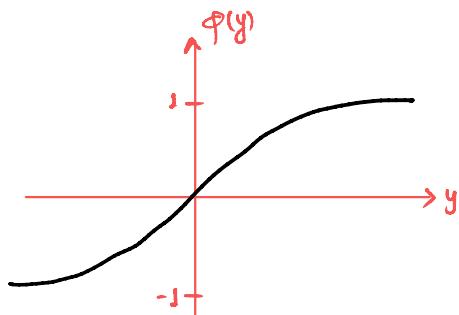
such that  $\underline{W^*} = [\omega_0 \quad \omega_1 \quad \omega_2 \quad \omega_3]$  and  $\underline{X^*} = [1 \quad \underline{X}]^\top$

$$Z = \varphi(y)$$

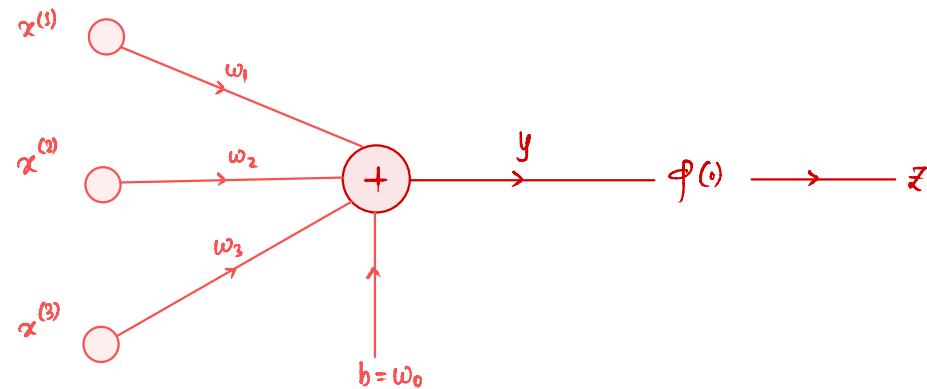
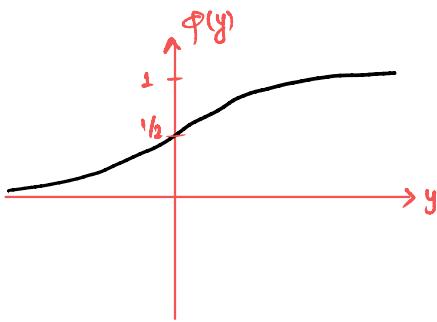
$$\varphi(y) = \begin{cases} 0 & y \leq 0 \\ 1 & \text{otherwise} \end{cases}$$



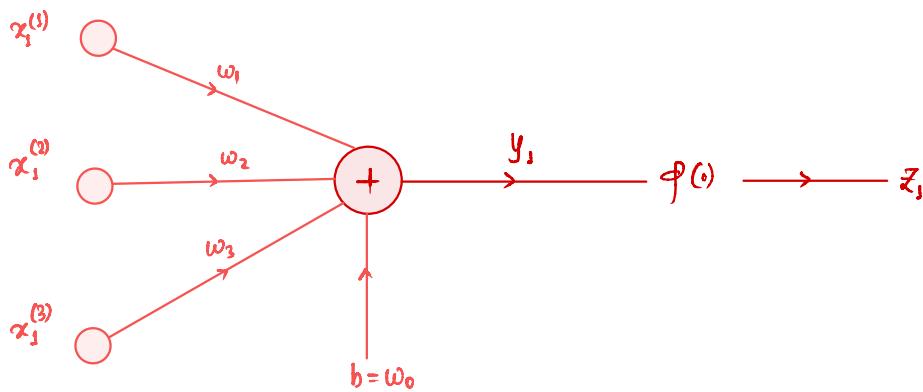
$$\varphi(y) = \tanh y = \frac{e^y - e^{-y}}{e^y + e^{-y}}$$

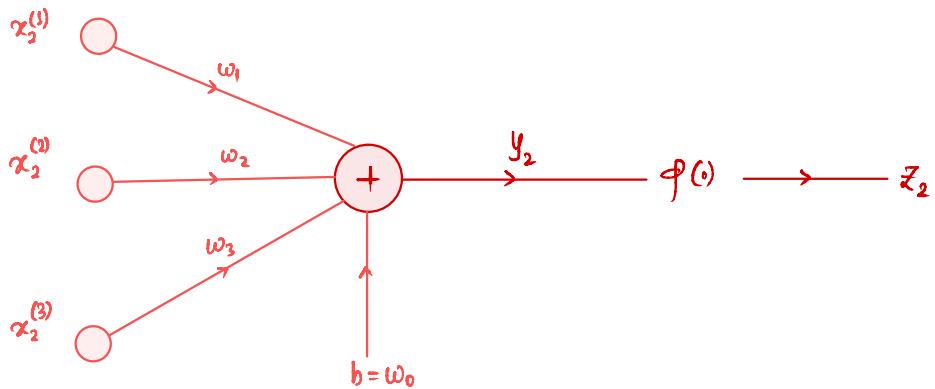


$$\phi(y) = \text{Sigmoid}(y) = \frac{1}{1 + e^{-y}}$$



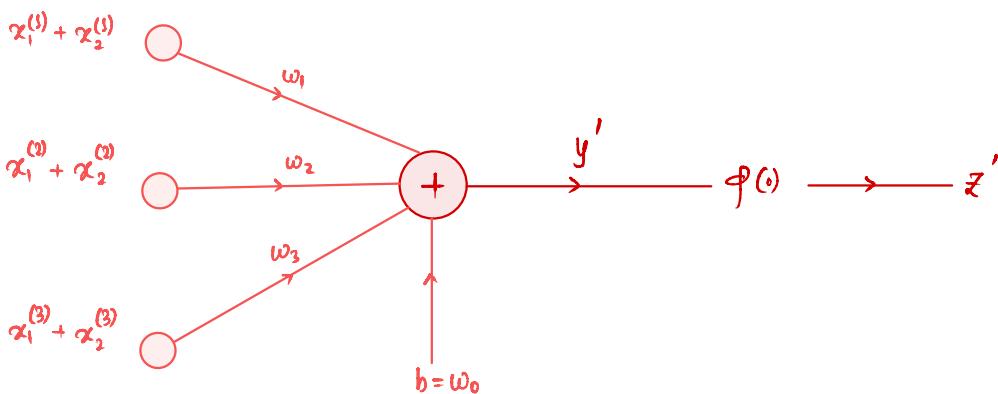
Is this system linear? prove it.



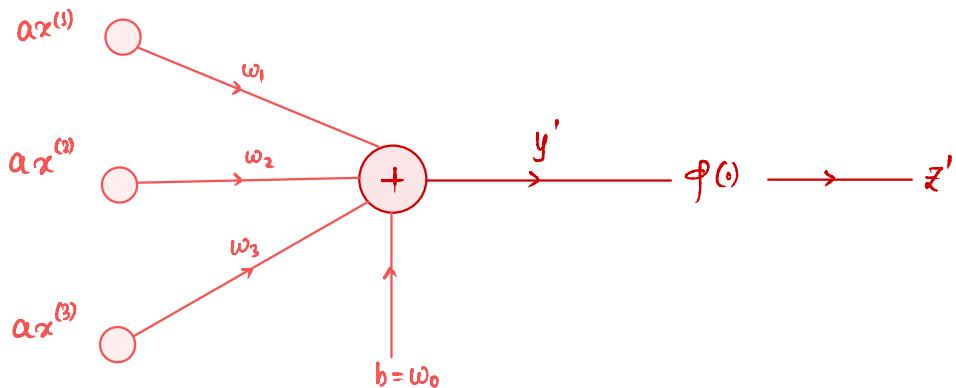


$$y_1 = w_0 + w_1 x_1^{(1)} + w_2 x_1^{(2)} + w_3 x_1^{(3)}$$

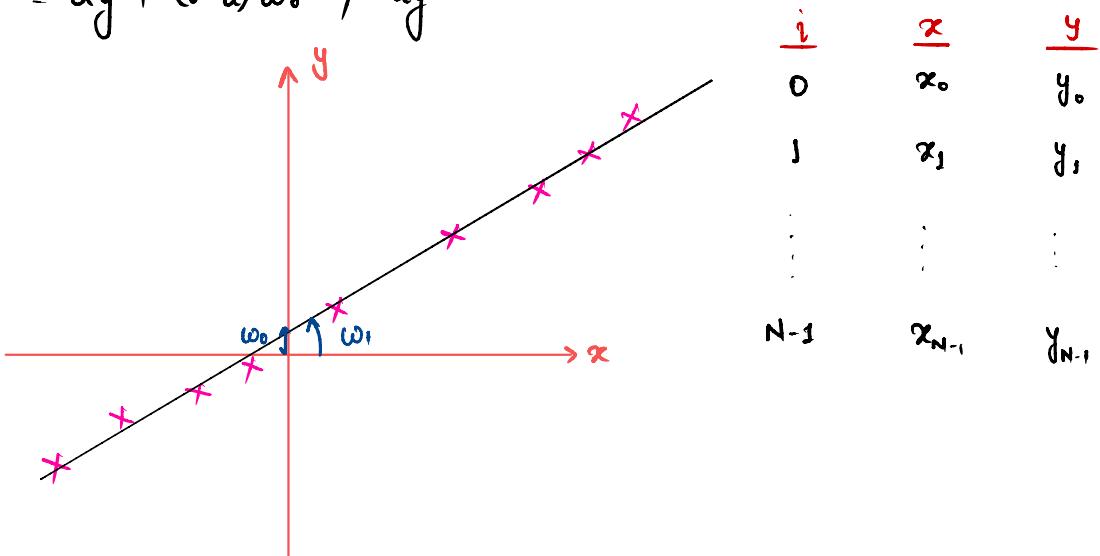
$$y_2 = w_0 + w_1 x_2^{(1)} + w_2 x_2^{(2)} + w_3 x_2^{(3)}$$



$$\begin{aligned}
 y' &= w_1 (x_1^{(1)} + x_2^{(1)}) + w_2 (x_1^{(2)} + x_2^{(2)}) + w_3 (x_1^{(3)} + x_2^{(3)}) + w_0 \\
 &= w_1 x_1^{(1)} + w_2 x_1^{(2)} + w_3 x_1^{(3)} + w_1 x_2^{(1)} + w_2 x_2^{(2)} + w_3 x_2^{(3)} \\
 &\quad + w_0 + w_0 - w_0 \\
 &= y_1 + y_2 - w_0 \neq y_1 + y_2
 \end{aligned}$$



$$\begin{aligned}
 y' &= w_1 \cdot \alpha x^{(1)} + w_2 \cdot \alpha x^{(2)} + w_3 \cdot \alpha x^{(3)} + w_0 + \alpha w_0 - \alpha w_0 \\
 &= \alpha (w_1 x^{(1)} + w_2 x^{(2)} + w_3 x^{(3)} + w_0) + (1-\alpha) w_0 \\
 &= \alpha y + (1-\alpha) w_0 \neq \alpha y
 \end{aligned}$$



$$y = f(x)$$

$\downarrow$

Some function  
implemented with NN

$$f(x) = w_0 + w_1 x = y$$

$$\begin{array}{rcl}
 y &=& mx + c \\
 \parallel && \parallel \\
 w_1 &=& w_0 \\
 \downarrow && \downarrow \\
 \text{Slope} &=& \text{y-axis intercept}
 \end{array}$$

Find  $w_0, w_1$  such that  $\text{MSE}(y, f(x))$  is minimum

MSE: Loss

$$L(\cdot) = J(\cdot) = \text{MSE}(y, \hat{y}) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - f(x_i))^2$$

$$\{w_0, w_1\}^* = \underset{\{w_0, w_1\}}{\operatorname{argmin}} J(\cdot)$$

$$(w_0, w_1)^{(t=0)} \sim \text{Random}, t=1, e=\infty$$

while  $e > \epsilon$

$$e = 0$$

for  $i=0, 1, \dots, N-1$

$$\hat{y}_i = \phi(x_i)$$

$$e = e + (\hat{y}_i - y_i)^2$$

end

$$J(\cdot) = e$$

$$(w_0, w_1)^{(t)} = (w_0, w_1)^{(t-1)} - \eta \frac{\partial}{\partial (w_0, w_1)} J^{(t)}(\cdot)$$

$$t = t+1$$

end

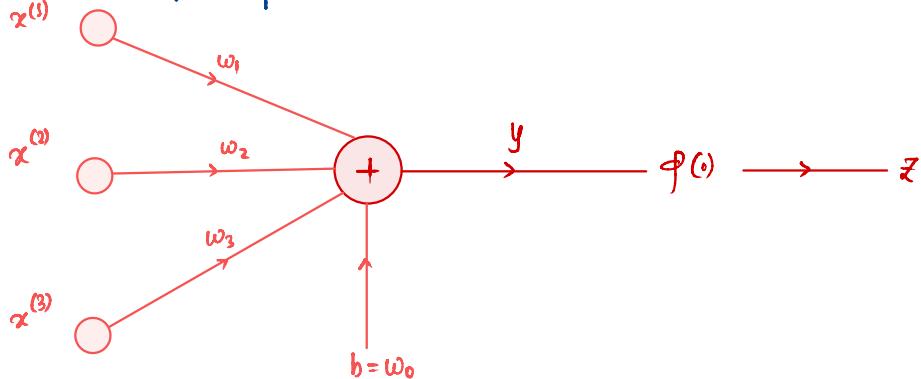
$\downarrow$   
Learning  
Rate

$$\hat{y} = w_0 + w_1 x$$

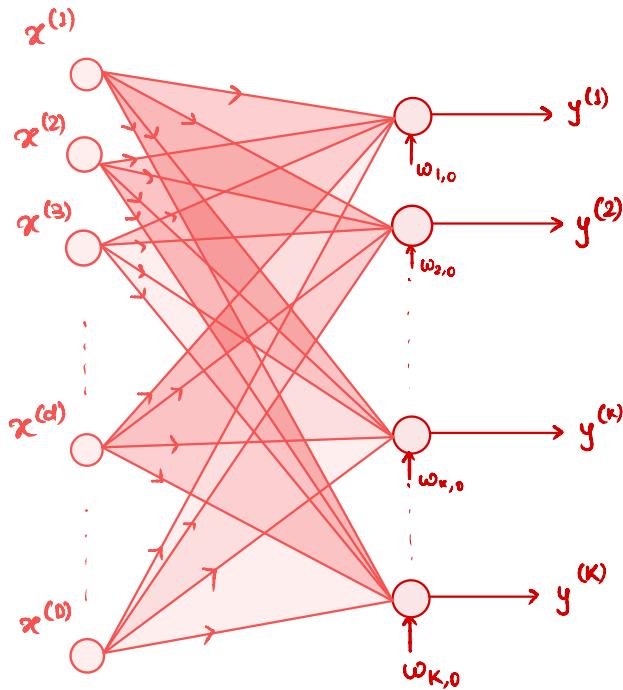
$$\frac{\partial}{\partial w_0} J(\cdot) = \frac{\partial}{\partial \hat{y}} J(\cdot) \cdot \frac{\partial}{\partial w_0} \hat{y} = -\frac{2}{N} \sum_{i=0}^{N-1} (y_i - f(x_i))$$

$$\frac{\partial}{\partial w_1} J(\cdot) = \frac{\partial}{\partial \hat{y}} J(\cdot) \frac{\partial}{\partial w_1} \hat{y} = -\frac{2}{N} \sum_{i=0}^{N-1} (y_i - f(x_i)) x_i$$

## Rosenblatt's perceptron



## Fully Connected Network



net ( $\cdot$ )

$$y^{(1)} = w_{1,0} + w_{1,1} x^{(1)} + w_{1,2} x^{(2)} + \dots + w_{1,d} x^{(d)} + \dots + w_{1,0} x^{(D)}$$

$$y^{(2)} = w_{2,0} + w_{2,1} x^{(1)} + w_{2,2} x^{(2)} + \dots + w_{2,d} x^{(d)} + \dots + w_{2,0} x^{(D)}$$

$$y^{(K)} = \omega_{k,0} + \omega_{k,1} x^{(1)} + \omega_{k,2} x^{(2)} + \dots + \omega_{k,d} x^{(d)} + \dots + \omega_{k,D} x^{(D)}$$

↓ Output Neuron Index      ↓ Output Neuron Index      ↓ Input Neuron Index  
Input Neuron Index

$$y^{(K)} = \omega_{k,0} + \omega_{k,1} x^{(1)} + \omega_{k,2} x^{(2)} + \dots + \omega_{k,d} x^{(d)} + \dots + \omega_{k,D} x^{(D)}$$

$$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(K)} \\ \vdots \\ y^{(K)} \end{bmatrix} = \begin{bmatrix} w_{1,0} & w_{1,1} & \cdots & w_{1,d} & \cdots & w_{1,D} \\ w_{2,0} & w_{2,1} & \cdots & w_{2,d} & \cdots & w_{2,D} \\ \vdots & \vdots & & \vdots & & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,d} & \cdots & w_{k,D} \\ \vdots & \vdots & & \vdots & & \vdots \\ w_{K,0} & w_{K,1} & \cdots & w_{K,d} & \cdots & w_{K,D} \end{bmatrix} \begin{bmatrix} 1 \\ x^{(1)} \\ \vdots \\ x^{(K)} \\ \vdots \\ x^{(K)} \end{bmatrix}$$

$$\underline{Y} = \underline{W} \underline{X} + \underline{b}$$

$\underline{Y} \in \mathbb{R}^{K \times J}$  : Output Activation

$\underline{W} \in \mathbb{R}^{K \times D}$  : Weights

$\underline{X} \in \mathbb{R}^{D \times J}$  : Input Activation

$\underline{b} \in \mathbb{R}^{K \times J}$  : Bias

$$\underline{W}^{(t+1)} = \underline{W}^{(t)} - \eta \frac{\partial}{\partial \underline{W}^{(t)}} J^{(t)}(\cdot)$$

$$\Rightarrow \underline{W}^{(t+1)} = \underline{W}^{(t)} - \eta \nabla \underline{W}^{(t)}$$

$$\frac{\partial}{\partial \underline{W}^{(t)}} J(\cdot) = \underbrace{\frac{\partial}{\partial \underline{y}^{(t)}} J^{(t)}(\cdot)}_{\nabla \text{Loss}(\cdot)} \cdot \underbrace{\frac{\partial}{\partial \underline{W}^{(t)}} \underline{y}^{(t)}}_{\nabla \text{net}(\cdot)}$$

$\nabla \text{Loss}(\cdot)$   
Gradient of loss  
function

$\nabla \text{net}(\cdot)$   
Gradient of  
the network

$$\underline{y} = \underline{W} \underline{x} + \underline{b} \quad \begin{matrix} \text{IR}^{k \times 1} \\ \text{IR}^{k \times D} \end{matrix} \quad \begin{matrix} \text{IR}^{1 \times D} \\ \text{IR}^{1 \times 1} \end{matrix}$$

$$\nabla \underline{W} = \frac{\partial}{\partial \underline{W}} J(\cdot) = \nabla \hat{\underline{y}} \underline{x}^T$$

$$\nabla \underline{b} = \frac{\partial}{\partial \underline{b}} J(\cdot) = \nabla \hat{\underline{y}} \quad \begin{matrix} \text{IR}^{k \times 1} \\ \text{IR}^{1 \times 1} \end{matrix}$$

$$\begin{matrix} \underline{y} & \xleftarrow{\text{R}} & \underline{y} & \xleftarrow{\text{S}} & \underline{x} \\ \underline{Q} & = & \underline{R} \underline{S} & + & \underline{T} \end{matrix}$$

$$\underline{R} \in \text{IR}^{a \times b}$$

$$\underline{b} \times c$$

$$\underline{S} \in \text{IR}^{a \times c}$$

$$\underline{T} \in \text{IR}^{a \times c}$$

$$\underline{Q} \in \text{IR}^{a \times c}$$

$$\frac{\partial}{\partial \underline{R}} J(\cdot) = \nabla \underline{R} = \nabla \underline{Q} \underline{S}^T$$

$$\frac{\partial}{\partial \underline{T}} J(\cdot) = \nabla \underline{T} = \nabla \underline{Q}$$

$$\frac{\partial}{\partial \underline{S}} J(\cdot) = \nabla \underline{S} = \underline{R}^T = \nabla \underline{Q}$$

$$J(\cdot) = \frac{1}{N} \sum_{i=0}^{N-1} (\underline{y}_i - \hat{\underline{y}}_i)^2$$

$$\begin{aligned} \frac{\partial}{\partial \underline{\hat{y}}} J(\cdot) &= \left[ \frac{\partial}{\partial \hat{y}_1} J(\cdot) \quad \frac{\partial}{\partial \hat{y}_2} J(\cdot) \quad \dots \quad \frac{\partial}{\partial \hat{y}_k} J(\cdot) \quad \dots \quad \frac{\partial}{\partial \hat{y}_K} J(\cdot) \right]^T \\ &= -\frac{2}{N} \sum_{i=0}^{N-1} (\underline{y}_i^{(i)} - \hat{\underline{y}}_i^{(i)}) \end{aligned}$$

$\in \mathbb{R}^{K \times 1}$

Update Rule:

$$\underline{W}^{(t+1)} = \underline{W}^{(t)} - \eta \nabla \underline{W}^{(t)}$$

$\mathbb{R}^{K \times D}$        $\mathbb{R}^{K \times D}$        $\mathbb{R}^{K \times D}$

$$\underline{b}^{(t+1)} = \underline{b}^{(t)} - \eta \nabla \underline{b}^{(t)}$$

$\mathbb{R}^{K \times 1}$        $\mathbb{R}^{K \times 1}$        $\mathbb{R}^{K \times 1}$

$\underline{x} \in \mathbb{R}^{D \times 1}$  : Inputs

$\underline{y} \in \mathbb{R}^{D \times 1}$  : Outputs

$\underline{W} \in \mathbb{R}^{K \times D}$  : Weights       $\underline{b} \in \mathbb{R}^{K \times 1}$  : Biases      } Parameters / Learnable parameters.

$\underline{t} \in \mathbb{R}^{K \times 1}$  : Target

Forward: (Input  $\rightarrow$  Output Port)

$$\underline{y} = \underline{W} \underline{x} + \underline{b}$$

Backward: (Output Port  $\rightarrow$  Input)

$$\nabla_{\underline{W}} = \nabla_{\underline{y}} \underline{x}^T$$

$$\nabla_{\underline{b}} = \nabla_{\underline{y}}$$

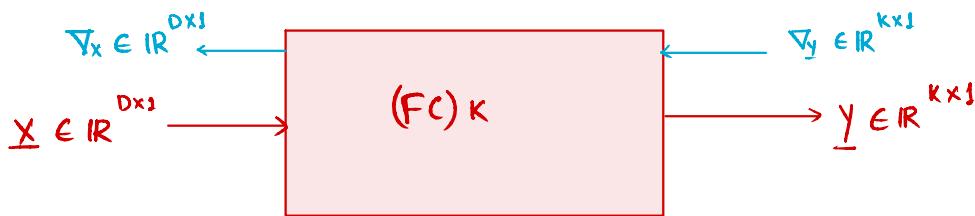
$$\nabla_{\underline{x}} = \underline{W}^T \nabla_{\underline{y}}$$

$$\underline{W} \in \mathbb{R}^{K \times D}$$

$$\underline{b} \in \mathbb{R}^{K \times 1}$$

$$\nabla_{\underline{W}} \in \mathbb{R}^{K \times D}$$

$$\nabla_{\underline{b}} \in \mathbb{R}^{K \times 1}$$



(Fc) := Symbolic Notation

$$\text{Forward} \left[ \begin{array}{l} \underline{y}_i = \underline{W} \underline{x}_i + \underline{b} \\ J(\cdot) = \text{fn}(\underline{t}_i, \underline{y}_i) \end{array} \right] : \text{Forward : net}(\cdot)$$

$$: \text{Forward : Loss}(\cdot)$$

$$\left[ \begin{array}{l} \nabla_{\underline{y}} = \frac{\partial}{\partial \underline{y}} J(\cdot) \end{array} \right] : \text{Backward : Loss}(\cdot)$$

$$\nabla_{\underline{W}} = \frac{\partial}{\partial \underline{W}} J(\cdot) = \frac{\partial}{\partial \underline{W}} \underline{y} \cdot \frac{\partial}{\partial \underline{W}} J(\cdot) = \nabla_{\underline{y}} \underline{x}^T : \text{Backward : Weight}$$

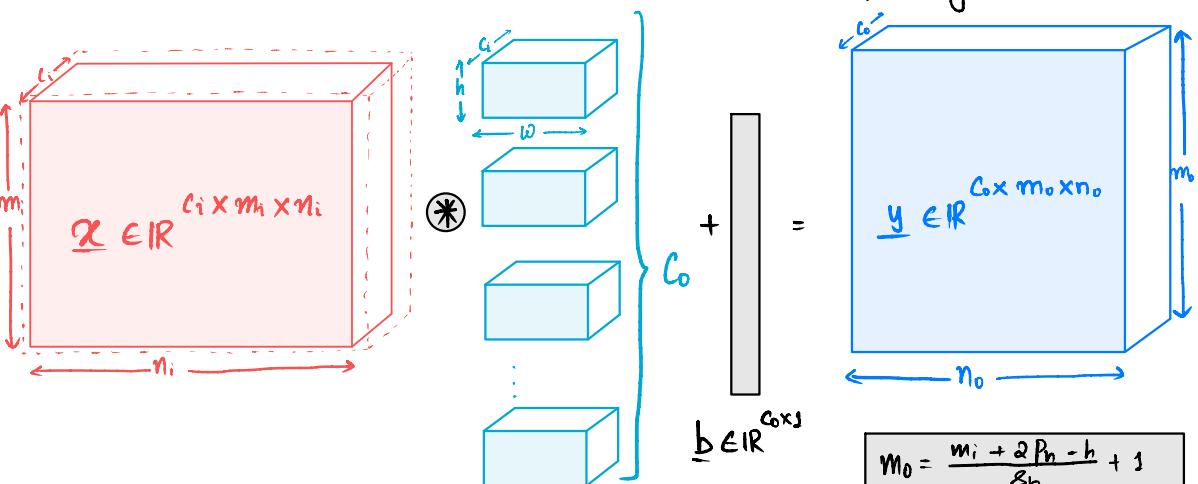
$$\underline{W}^{T+1} = \underline{W}^T - \eta \nabla_{\underline{W}}^{(T)} : \text{Update : weights}$$

Backward

$$\nabla_{\underline{b}} = \frac{\partial}{\partial \underline{b}} J(\cdot) = \frac{\partial}{\partial \underline{b}} \underline{y} \cdot \frac{\partial}{\partial \underline{b}} J(\cdot) : \text{Backward : Biases}$$

$$\underline{b}^{T+1} = \underline{b}^T - \eta \nabla_{\underline{b}}^{(T)} : \text{Update : Biases.}$$

## Convolution Neuron :-



$$m_o = \frac{m_i + 2P_h - h}{\delta_h} + 1$$

$$n_o = \frac{n_i + 2P_w - w}{\delta_w} + 1$$

$$\begin{matrix} x_{11} & x_{12} & x_{13} & \dots \\ x_{21} & x_{22} & x_{23} & \dots \\ x_{31} & x_{32} & x_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{matrix}$$



$$\begin{matrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{matrix}$$

180° rotated  
version of  
Convolution Kernel

$$\begin{matrix} y_{11} & y_{12} & y_{13} & \dots \\ y_{21} & y_{22} & y_{23} & \dots \\ y_{31} & y_{32} & y_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{matrix}$$

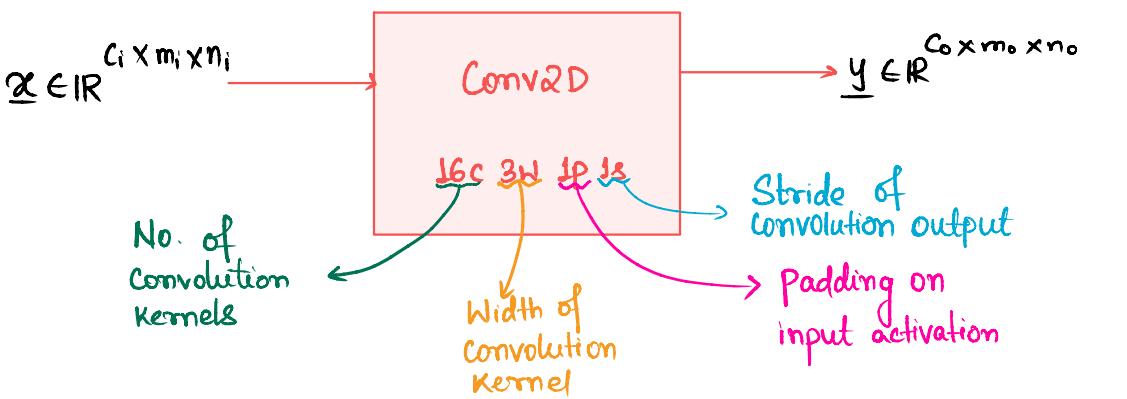
$$y_{11} = x_{11} w_{11} + x_{12} w_{12} + x_{13} w_{13} + x_{21} w_{21} + x_{22} w_{22} + x_{23} w_{23} \\ + x_{31} w_{31} + x_{32} w_{32} + x_{33} w_{33}$$

$$y_{12} = x_{12} w_{12} + x_{13} w_{13} + x_{14} w_{14} + x_{22} w_{22} + x_{23} w_{23} + x_{24} w_{24} \\ + x_{32} w_{32} + x_{33} w_{33} + x_{34} w_{34}$$

Size of Convolution Kernel :  $w, h = 3, 3$

Padding :  $P_w, P_h = 1$

Stride :  $S_w, S_h = 1$



$$P_w, P_h = 3, 3$$

$$S_w, S_h = 3, 3$$

$$(\text{Conv2D}) : \mathbb{R}^{1 \times 32 \times 32} \mapsto \mathbb{R}^{16 \times 32 \times 32}$$

$$\underline{x} \rightarrow (\text{Conv2D}) 16c 5w 1s Op \longrightarrow \underline{y}$$

$$\underline{x} \in \mathbb{R}^{1 \times 32 \times 32}$$

$$\underline{y} \in \mathbb{R}^{16 \times 32 \times 32}$$

If  $\underline{y} \in \mathbb{R}^{16 \times 32 \times 32}$ , then find when  $w \times h = 5 \times 5$ ,

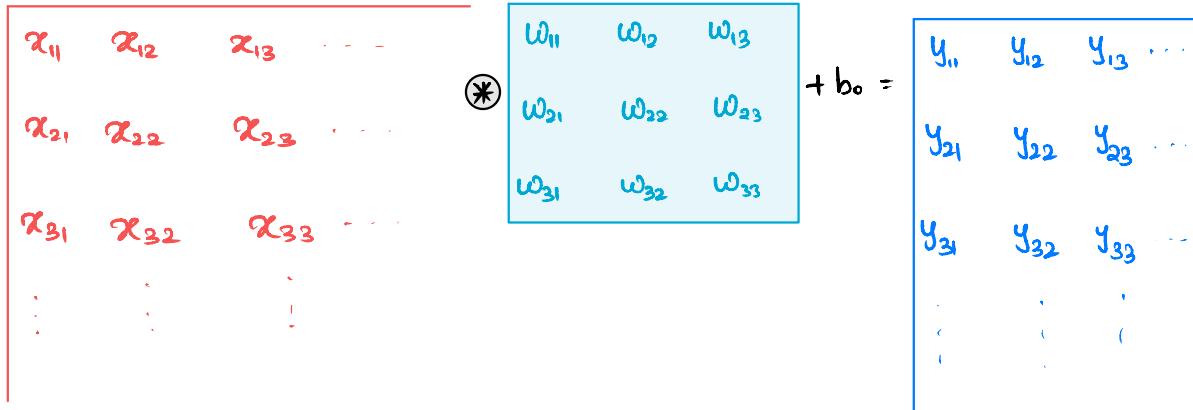
$$C_o = 16$$

$$(\text{Conv2D}): \mathbb{R}^{1 \times 32 \times 32} \mapsto \mathbb{R}^{16 \times 32 \times 32}$$

$$P_w, P_h = 2, 2$$

$$\delta_w, \delta_h = 1, 1$$

$P_w, P_h$	2, 2	33, 33
$\delta_w, \delta_h$	1, 1	3, 3
Proportion (%)	87%	33%



$$Y_{11} = X_{11}w_{11} + X_{12}w_{12} + X_{13}w_{13} + X_{21}w_{21} + X_{22}w_{22} + X_{23}w_{23} \\ + X_{31}w_{31} + X_{32}w_{32} + X_{33}w_{33} + b_0$$



$$\underline{X} \in \mathbb{R}^{m_0 n_0 \times C_i W_i}$$

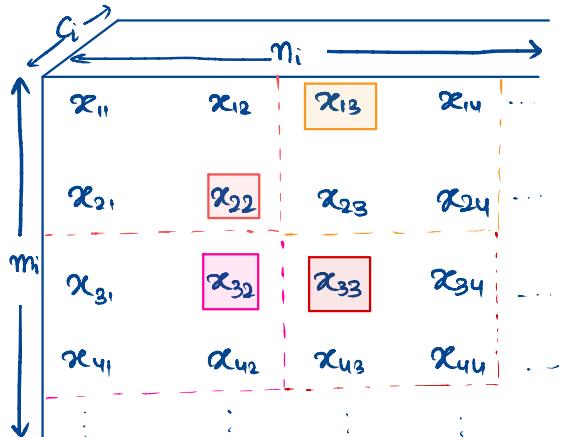
$$\underline{W} \in \mathbb{R}^{C_{out} \times C_i} \quad \underline{Y} \in \mathbb{R}^{m_0 n_0 \times C_{out}}$$

$$\nabla_w = ?$$

$$\nabla_b = ?$$

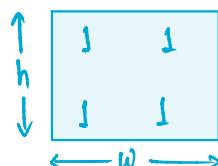
Pooling :-

Max-pooling -

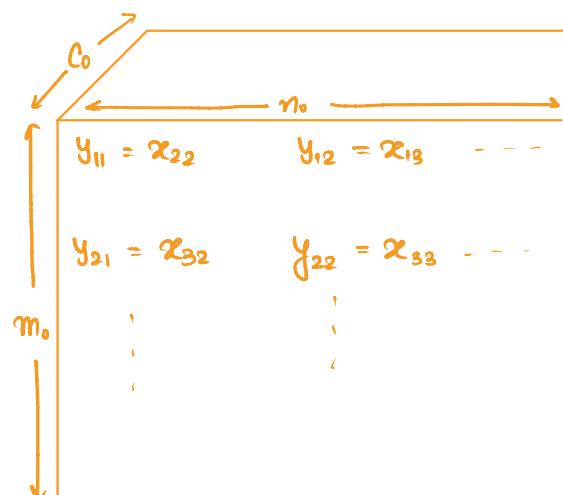


Pooling Kernel

$$\text{Kernel Size : } w, h = 2, 2$$



$$\text{Stride : } \Delta_w, \Delta_h = 2, 2$$

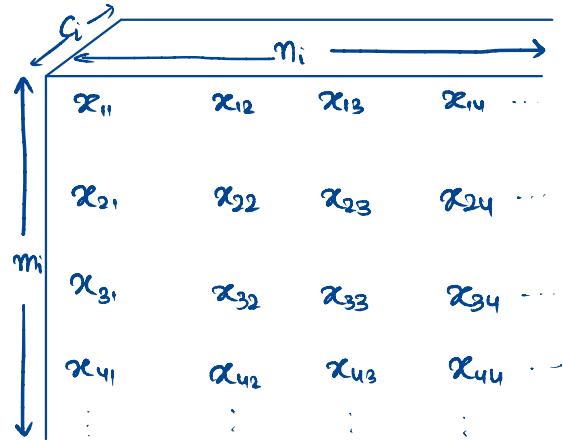


$$m_0 = \frac{m_i - h}{\Delta_h} + 1$$

$$n_0 = \frac{n_i - w}{\Delta_w} + 1$$

$$C_0 = C_i$$

## Avg - pool 2D -



$$s_w, s_h = 2, 2$$

$$y_{11} = \frac{1}{4} (x_{11} + x_{12} + x_{21} + x_{22})$$

$$y_{12} = \frac{1}{4} (x_{13} + x_{14} + x_{23} + x_{24})$$

$$y_{21} = \frac{1}{4} (x_{31} + x_{32} + x_{41} + x_{42})$$

⋮

$$\underline{x} \in \mathbb{R}^{C_i \times m_i \times n_i}$$

$$\underline{y} \in \mathbb{R}^{C_o \times m_o \times n_o}$$

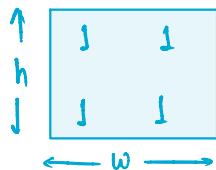
$$C_o = C_i$$

$$m_o = \frac{m_i - h}{s_h} + 1$$

$$n_o = \frac{n_i - w}{s_w} + 1$$

## pooling Kernel

Kernel Size :  $w, h = 2, 2$



$$s_w, s_h = 1, 1$$

$$y_{11} = \frac{1}{4} (x_{11} + x_{12} + x_{21} + x_{22})$$

$$y_{12} = \frac{1}{4} (x_{13} + x_{14} + x_{23} + x_{24})$$

$$y_{21} = \frac{1}{4} (x_{31} + x_{32} + x_{41} + x_{42})$$

$$y_{22} = \frac{1}{4} (x_{33} + x_{34} + x_{43} + x_{44})$$

$$y_{22} = \frac{1}{4} (x_{33} + x_{34} + x_{43} + x_{44})$$

$$\nabla_x \in \mathbb{R}^{C_i \times m_i \times n_i}$$

$$\nabla_y \in \mathbb{R}^{C_o \times m_o \times n_o}$$

 Yann Le Cun → Le → Net → NN → -5 → 5 layers with trainable parameters

net(0)  $\mapsto$  (1: Conv2D) 6c 5w 5h 0p  $\rightarrow$  (2: maxpool2D) 2w 2h  
 $\rightarrow$  (3: Conv2D) 16c 3w 3h 0p  $\rightarrow$  (4: maxpool2D) 2w 2h  
 $\rightarrow$  (5: Fc) 120  $\rightarrow$  (6: Fc) 10  $\rightarrow$  (7: Sigmoid)

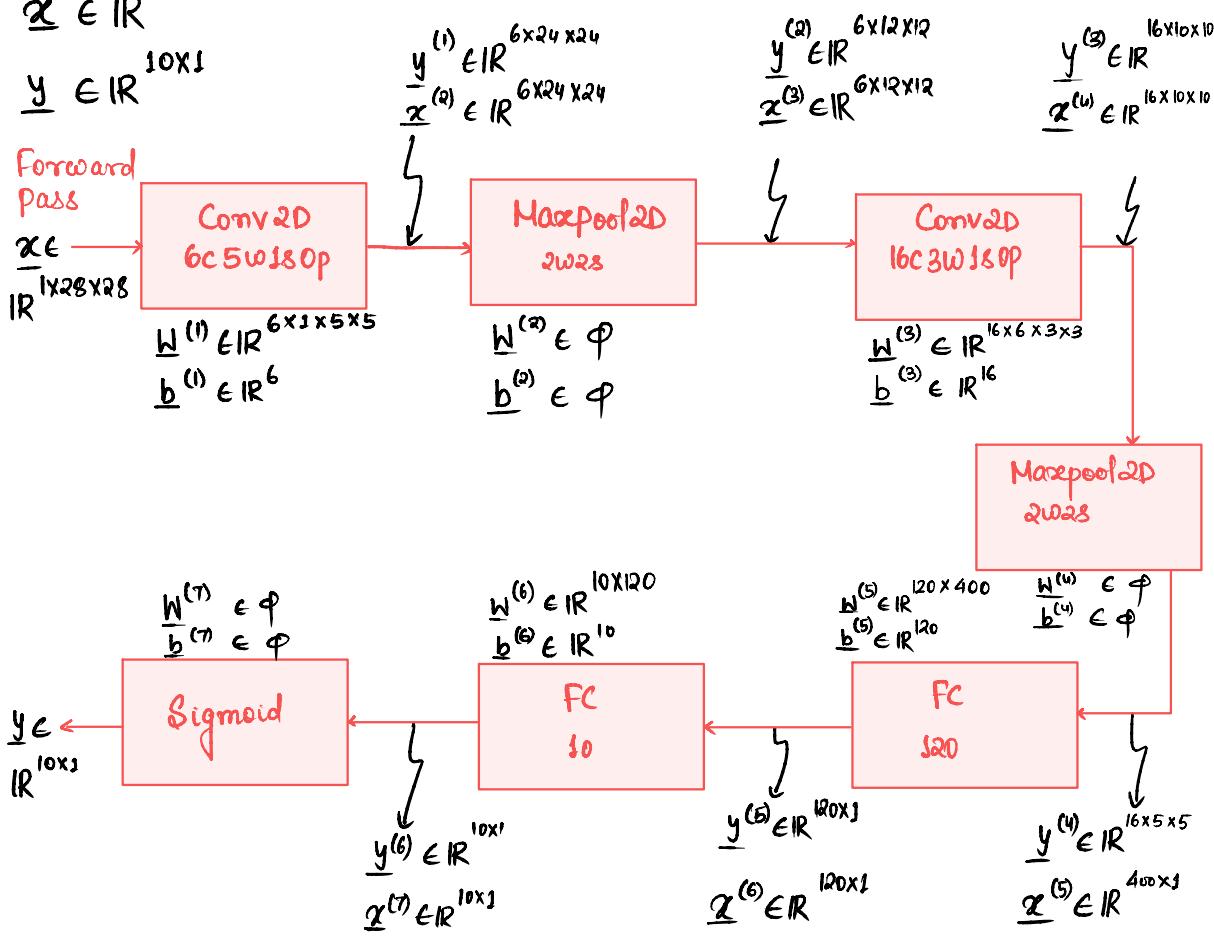
$$y = \text{net } (\underline{x})$$

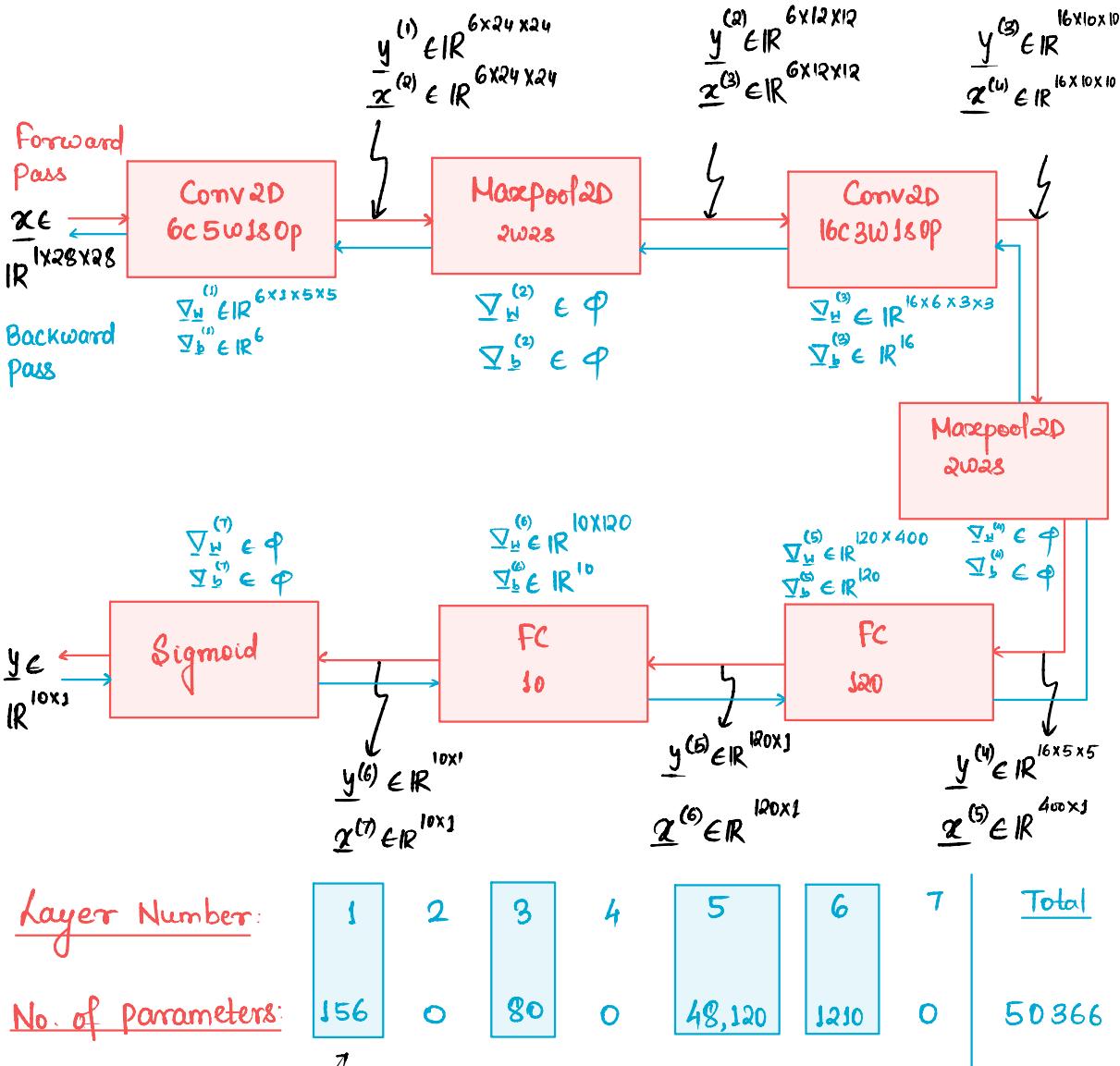
1x28x28

$$\underline{x} \in \mathbb{R}$$

$$\|c\| \in \mathbb{R}^{10 \times 1}$$

## Forward Pass





$| \cdot | \rightarrow$  Set cardinality operator

$\text{sizeof}(\text{net}(\cdot)) = \text{Total number of parameters} \cdot \text{precision}$

$$\text{Float 32} = 50366 \times 4 \text{ Bytes}$$

$$\approx 200 \text{ KBytes}$$

$$\text{Double / Float 64} = 50366 \times 8 \text{ Bytes}$$

$$\approx 400 \text{ KBytes}$$

$$\text{No. of multiplications} = n_o m_o C_o C_i w h$$

$$\begin{aligned}\text{No. of additions} &= C_o m_o n_o ((C_i w h - 1) + 1) \\ &= n_o m_o C_o C_i w h\end{aligned}$$

$$\text{No. of operations} = 2 n_o m_o C_o C_i w h$$

<u>Operations</u>	172800	20736	172,800	25600	9600	2400	90	490426
								$\approx 0.5 \text{ MegaFlops}$

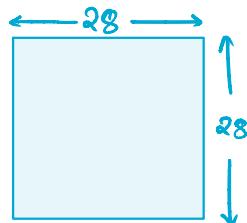
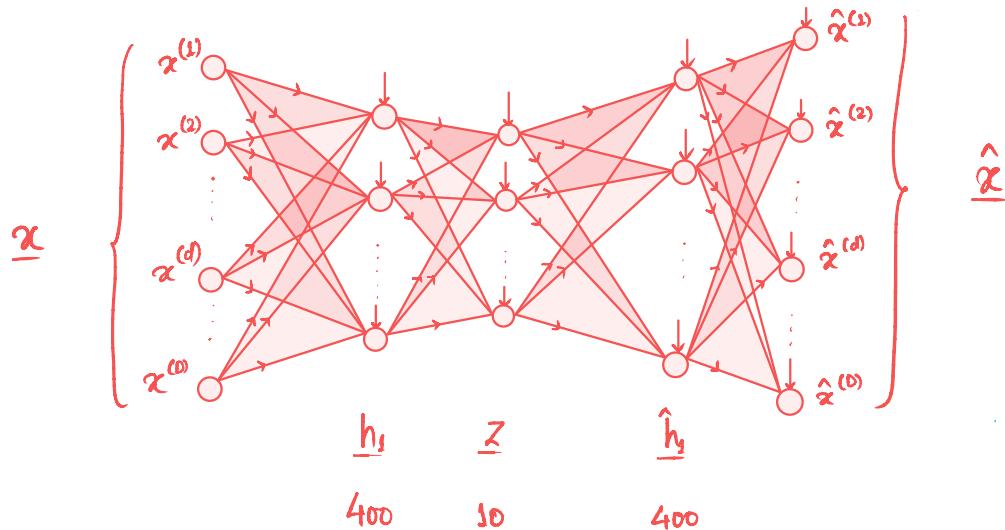
$$\begin{aligned}\text{No of operations per second} &= \frac{10 \times 10^3 \text{ MegaFlops } s^{-1}}{0.5 \text{ MegaFlops}} \\ &= 20,000 s^{-1}\end{aligned}$$

### Generative Models :-

$p(\underline{x} = \omega | \underline{x}) := \text{Discriminative Model}$

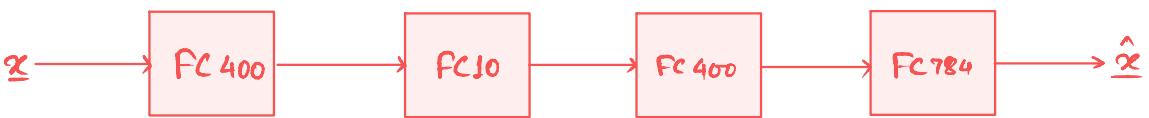
$p(y, \underline{x}) := \text{Generative Model}$

## AutoEncoder-



$$\underline{I} \in \mathbb{R}^{1 \times 28 \times 28}$$

$$\underline{x} \in \mathbb{R}^{784 \times 1}$$



$$\text{net}_{AE}(\cdot) = \text{net}_{\text{Dec}}(\text{net}_{\text{Enc}}(\cdot))$$

$$\hat{\underline{x}}_{AE} = \text{net}_{AE}(\underline{x})$$

$$\underline{z} = \text{net}_{\text{Enc}}(\underline{x})$$

$$\hat{\underline{x}} = \text{net}_{\text{Dec}}(\underline{z})$$

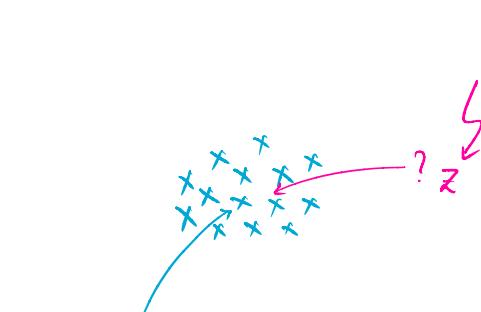
$$J(\cdot) = \frac{1}{N} \sum_{i=0}^{N-1} \|\underline{x}_i - \hat{\underline{x}}_i\|_2^2$$

$$\{\underline{w}^*, \underline{b}^*\} = \underset{\{\underline{w}, \underline{b}\}}{\operatorname{argmin}} J(\cdot)$$

$$\text{net}_{\text{Enc}}(\cdot) \mapsto (1: \text{FC}) 400 \rightarrow (2: \text{FC}) 10 \Rightarrow \mathbb{R}^{784 \times 1} \mapsto \mathbb{R}^{10 \times 1}$$

$$\text{net}_{\text{Dec}}(\cdot) \mapsto (1: \text{FC}) 400 \rightarrow (2: \text{FC}) 784 \Rightarrow \mathbb{R}^{10 \times 1} \mapsto \mathbb{R}^{784 \times 1}$$

$\underline{z}^{(2)}$



Randomly picked  
point

$$\text{net}_{\text{Enc}}(\cdot) : \mathbb{R}^{784 \times 1} \mapsto \mathbb{R}^{2 \times 1}$$

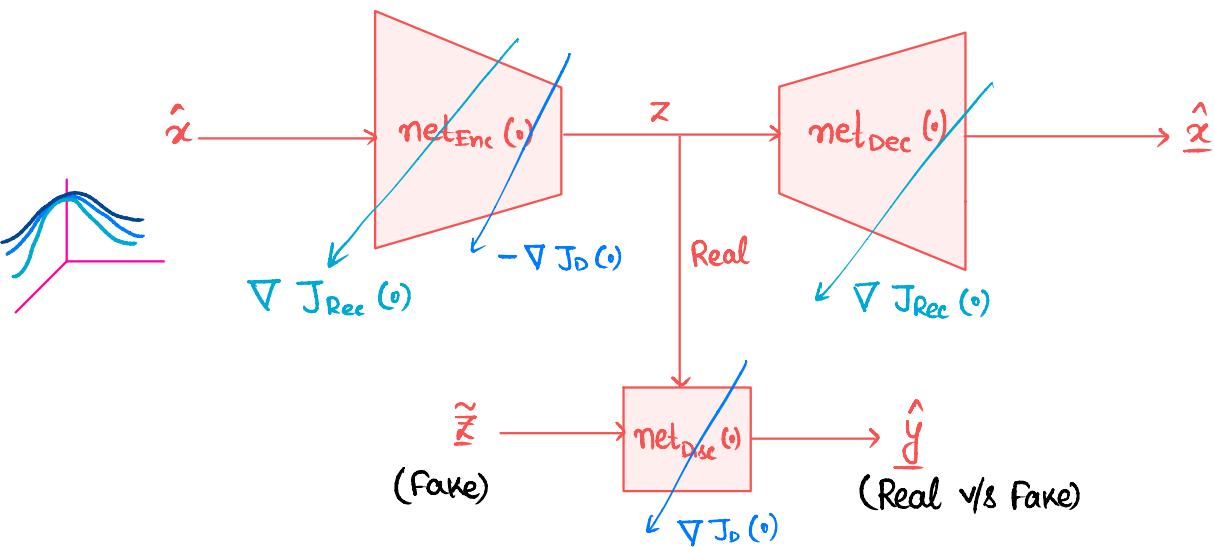
$$\text{net}_{\text{Dec}}(\cdot) : \mathbb{R}^{2 \times 1} \mapsto \mathbb{R}^{784 \times 1}$$

$$\text{net}_{\text{Dec}}(\cdot) : \mathbb{R}^{2 \times 1} \mapsto \mathbb{R}^{784 \times 1}$$

$$\tilde{\underline{z}} \mapsto \tilde{\underline{x}}$$

$$\begin{aligned} \tilde{\underline{x}} &\in \mathbb{R}^{784 \times 1} \\ \hat{\underline{t}} &\in \mathbb{R}^{1 \times 28 \times 28} \end{aligned}$$

$$\tilde{\underline{x}} \neq \{\underline{x}_i\}$$



$$J_{Rec}(\cdot) = \text{MSE}(x_i, \hat{x}_i)$$

$$\hat{y} \in \{\text{Real}, \text{Fake}\}$$

$$J_D(\cdot) = \text{BCE}(y, \hat{y})$$

- $net_{Enc}(\cdot)$ ,  $net_{Dec}(\cdot)$ ,  $net_{Disc}(\cdot)$  → Randomly initialized

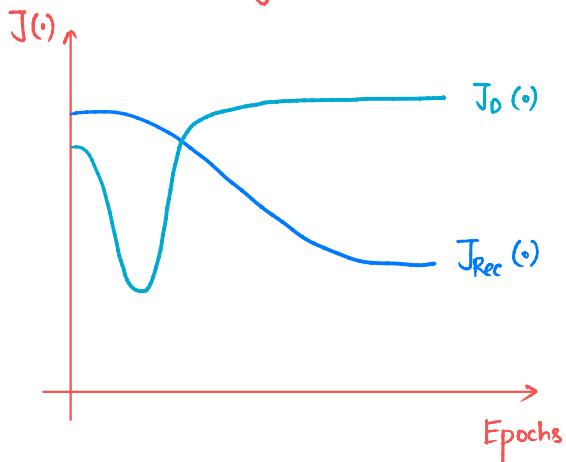
for  $i = 0, 1, \dots, N-1$

$$\underline{z}_i = net_{Enc}(\underline{x}_i)$$

$$\hat{\underline{x}}_i = net_{Dec}(\underline{z}_i)$$

$$e_i = \|\hat{\underline{x}}_i - \underline{x}_i\|_2^2$$

$$s_i = \text{Random}\{\text{Fake}, \text{Real}\}$$



$$\underline{h}_i = \underline{z}_i \cdot \underline{s}_i + \tilde{\underline{z}}_i \cdot (1 - \underline{s}_i)$$

$\{\underline{w}_{\text{enc}}, \underline{b}_{\text{enc}}\} \longleftrightarrow \text{net}_{\text{enc}}(\cdot)$

$$\hat{\underline{y}}_i = \text{net}_{\text{disc}}(\underline{h}_i)$$

$\{\underline{w}_{\text{dec}}, \underline{b}_{\text{dec}}\} \longleftrightarrow \text{net}_{\text{dec}}(\cdot)$

$$d_i = \text{BCE}(\hat{\underline{y}}_i, \underline{s}_i)$$

$\{\underline{w}_{\text{disc}}, \underline{b}_{\text{disc}}\} \longleftrightarrow \text{net}_{\text{disc}}(\cdot)$

$$J_{\text{Rec}}(\cdot) = J_{\text{Rec}}(\cdot) + e_i$$

$$J_D(\cdot) = J_D(\cdot) + d_i$$

end

$$\underline{w}_{\text{dec}}^{(t+1)} = \underline{w}_{\text{dec}}^{(t)} - \eta \frac{\partial}{\partial \underline{w}_{\text{dec}}^{(t)}} J_{\text{Rec}}^{(t)}(\cdot)$$

$\nabla J_{\text{Rec}}(\cdot)$

$$\underline{b}_{\text{dec}}^{(t+1)} = \underline{b}_{\text{dec}}^{(t)} - \eta \frac{\partial}{\partial \underline{b}_{\text{dec}}^{(t)}} J_{\text{Rec}}^{(t)}(\cdot)$$

$$\underline{w}_{\text{enc}}^{(t+1)} = \underline{w}_{\text{enc}}^{(t)} - \eta \frac{\partial}{\partial \underline{w}_{\text{enc}}^{(t)}} J_{\text{Rec}}^{(t)}(\cdot)$$

$$+ \eta_2 \frac{\partial}{\partial \underline{w}_{\text{enc}}^{(t)}} J_D^{(t)}(\cdot)$$

$$\underline{b}_{\text{enc}}^{(t+1)} = \underline{b}_{\text{enc}}^{(t)} - \eta \frac{\partial}{\partial \underline{b}_{\text{enc}}^{(t)}} J_{\text{Rec}}^{(t)}(\cdot)$$

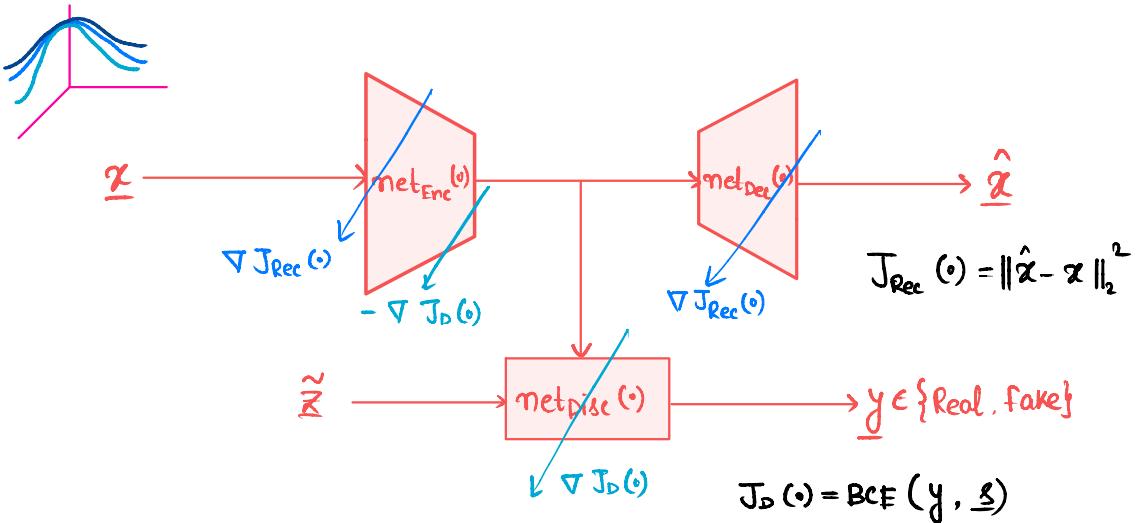
$$+ \eta_2 \frac{\partial}{\partial \underline{b}_{\text{enc}}^{(t)}} J_D^{(t)}(\cdot)$$

$$\underline{w}_{\text{disc}}^{(t+1)} = \underline{w}_{\text{disc}}^{(t)} - \eta \frac{\partial}{\partial \underline{w}_{\text{disc}}^{(t)}} J_D^{(t)}(\cdot)$$

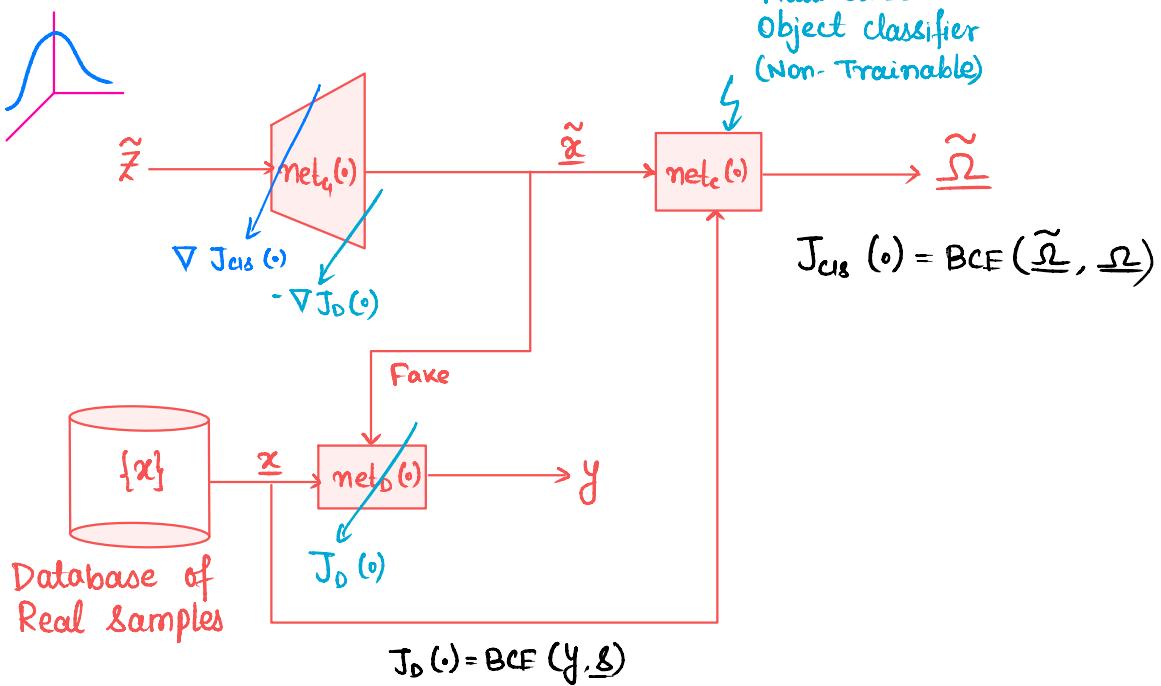
$\nabla J_D(\cdot)$

$$\underline{b}_{\text{disc}}^{(t+1)} = \underline{b}_{\text{disc}}^{(t)} - \eta \frac{\partial}{\partial \underline{b}_{\text{disc}}^{(t)}} J_D^{(t)}(\cdot)$$

## Adversarial Autoencoder (AAE)-



## Generative Adversarial Network-



AAF

$$\text{net}_{\text{Enc}}(\cdot) : \mathbb{R}^{C \times N \times N} \mapsto \mathbb{R}^q$$

$$\text{net}_{\text{Dec}}(\cdot) : \mathbb{R}^q \mapsto \mathbb{R}^{C \times N \times N}$$

$$\text{net}_{\text{Disc}}(\cdot) : \mathbb{R}^q \mapsto \mathbb{R}^2$$

GAN

$$\text{net}_A(\cdot) : \mathbb{R}^2 \mapsto \mathbb{R}^{C \times N \times N}$$

$$\text{net}_B(\cdot) : \mathbb{R}^{C \times N \times N} \mapsto \mathbb{R}^2$$

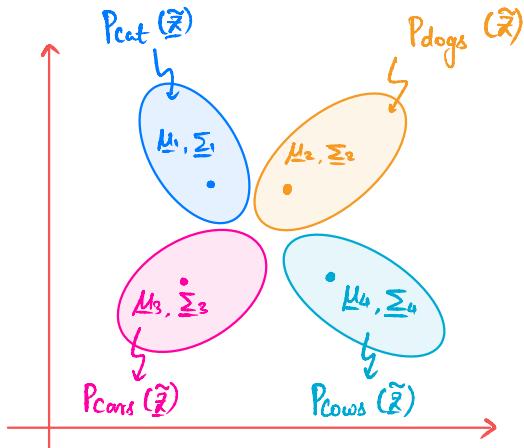
$$\text{net}_C(\cdot) : \mathbb{R}^{C \times N \times N} \mapsto \mathbb{R}^k$$

Conditional Generation:-

$$P_{\text{cats}}(\tilde{x}) \neq P_{\text{dogs}}(\tilde{x}) \neq P_{\text{cows}}(\tilde{x}) \neq P_{\text{cars}}(\tilde{x})$$

$$q=2$$

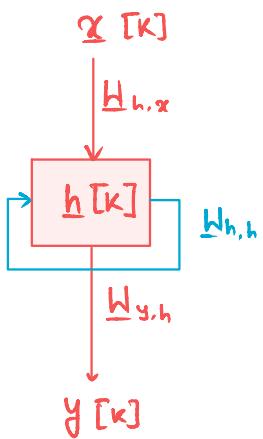
$$P(\tilde{x}) \sim \mathcal{N}(\mu, \Sigma)$$



$$\mu_1, \mu_2, \mu_3, \mu_4 \in \mathbb{R}^2$$

$$\Sigma_1, \Sigma_2, \Sigma_3, \Sigma_4 \in \mathbb{R}^{2 \times 2}$$

## Recurrent Neural Networks (RNN) :-



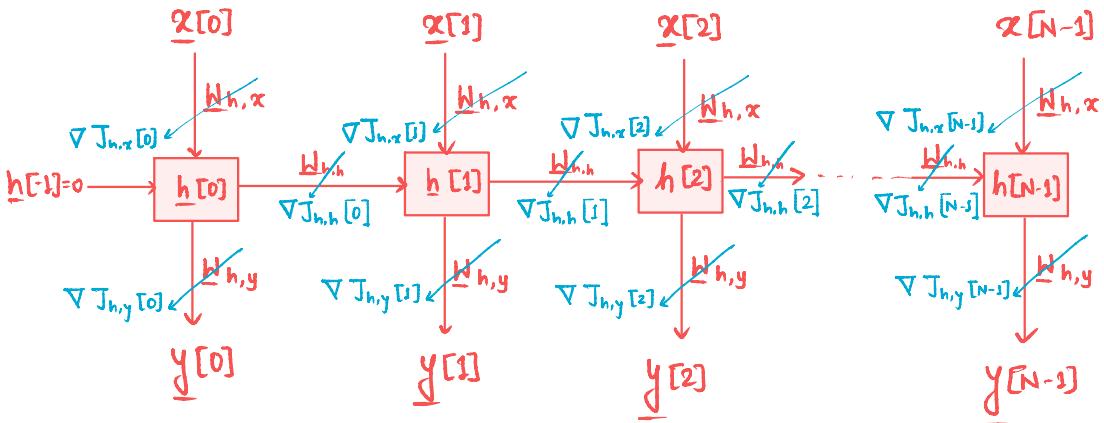
$$\underline{y}[k] = \underline{W}_{y,h} \underline{h}[k] + \underline{b}_y$$

$$\underline{h}[k] = \underline{W}_{h,x} \underline{x}[k] + \underline{b}_{h,x} + \underline{W}_{h,h} \underline{h}[k-1]$$

Many to Many  $\Rightarrow \mathbb{R}^{D \times N} \mapsto \mathbb{R}^{K \times N}$

Many to One  $\Rightarrow \mathbb{R}^{D \times N} \mapsto \mathbb{R}^{K \times 1} \rightarrow \underline{y}[N-1]$

One to Many  $\Rightarrow \mathbb{R}^{D \times 1} \mapsto \mathbb{R}^{K \times N} \rightarrow \underline{x}[0]$



$$\underline{X} = [\underline{x}[0], \underline{x}[1], \underline{x}[2], \dots, \underline{x}[n], \dots, \underline{x}[N-1]]$$

Say,  $\underline{x}[n] \in \mathbb{R}^D$ , then  $\underline{X} \in \mathbb{R}^{D \times N}$

$$\underline{Y} = [y[0], y[1], y[2], \dots, y[n], \dots, y[N-1]]$$

Say,  $y[n] \in \mathbb{R}^k$ , then  $\underline{Y} \in \mathbb{R}^{k \times N}$

