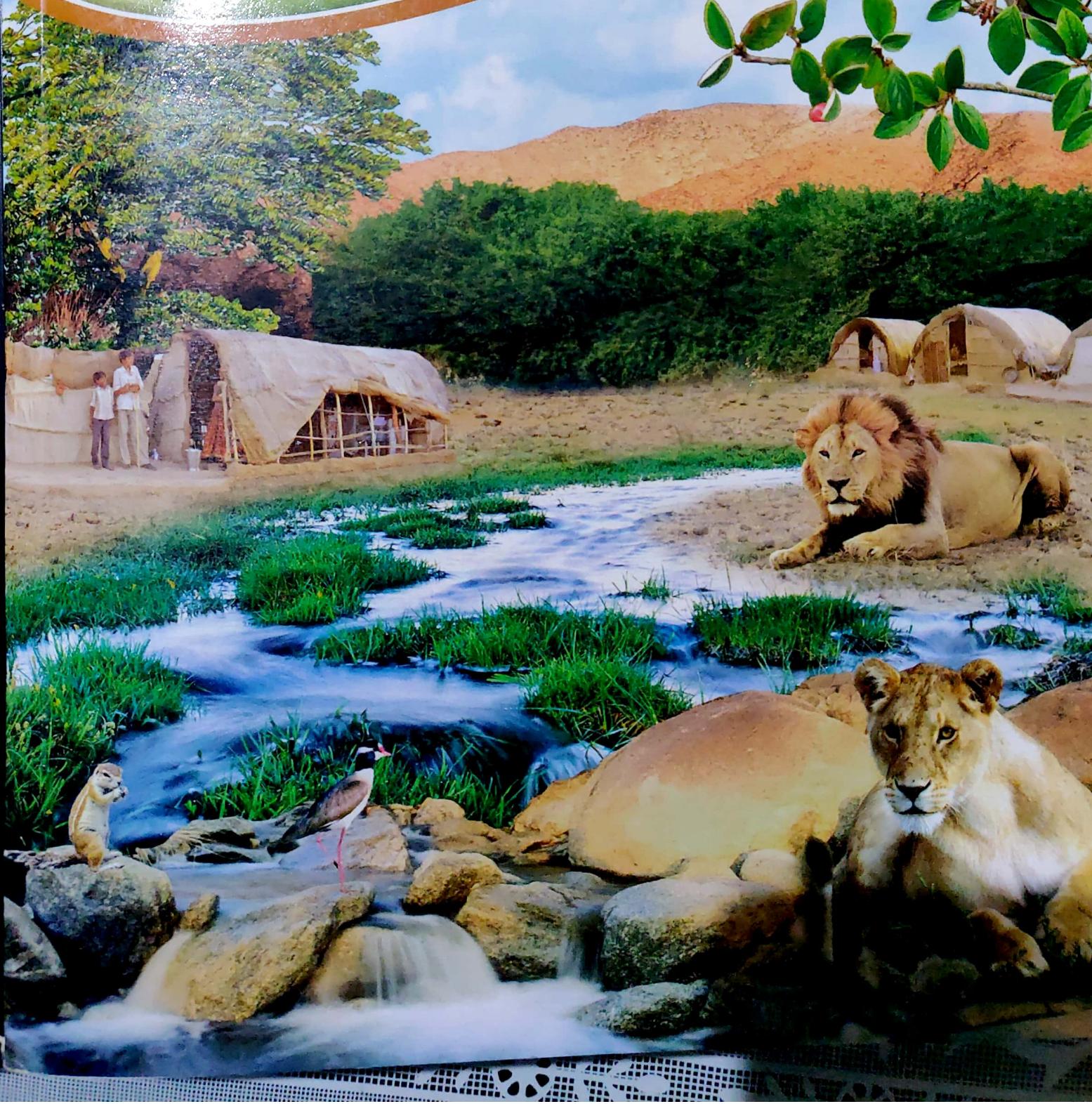




# અનંદિકાર

## નવી રંગા નો કલબ



-Kaushal P. Shah  
- Machine Learning  
(2)

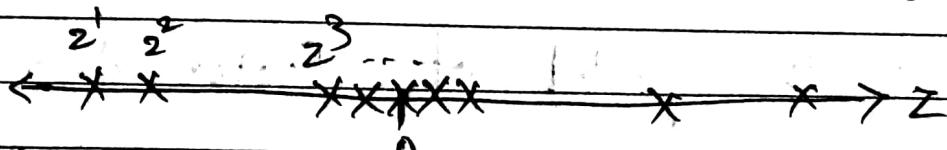
\* Example of factor Analysis

→ let  $z \in \mathbb{R}^3$ ,  $x \in \mathbb{R}^2$ , also

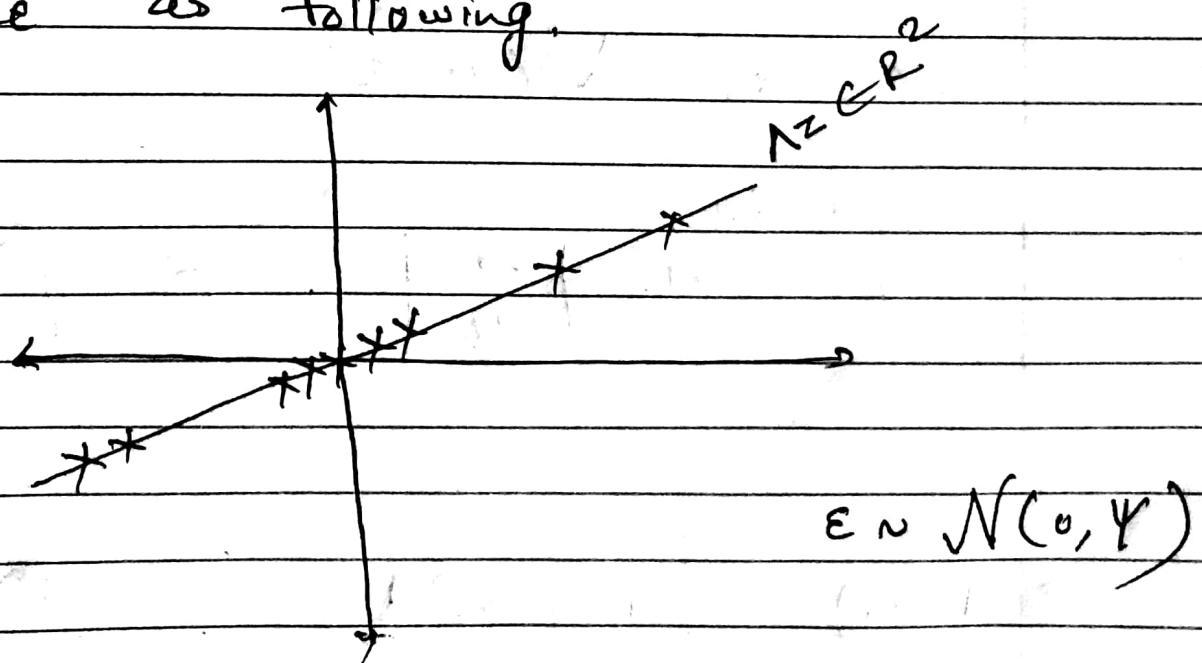
$$z \sim N(0, I), \quad A = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad u = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Psi = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

→ Then  $z$  can be represented by



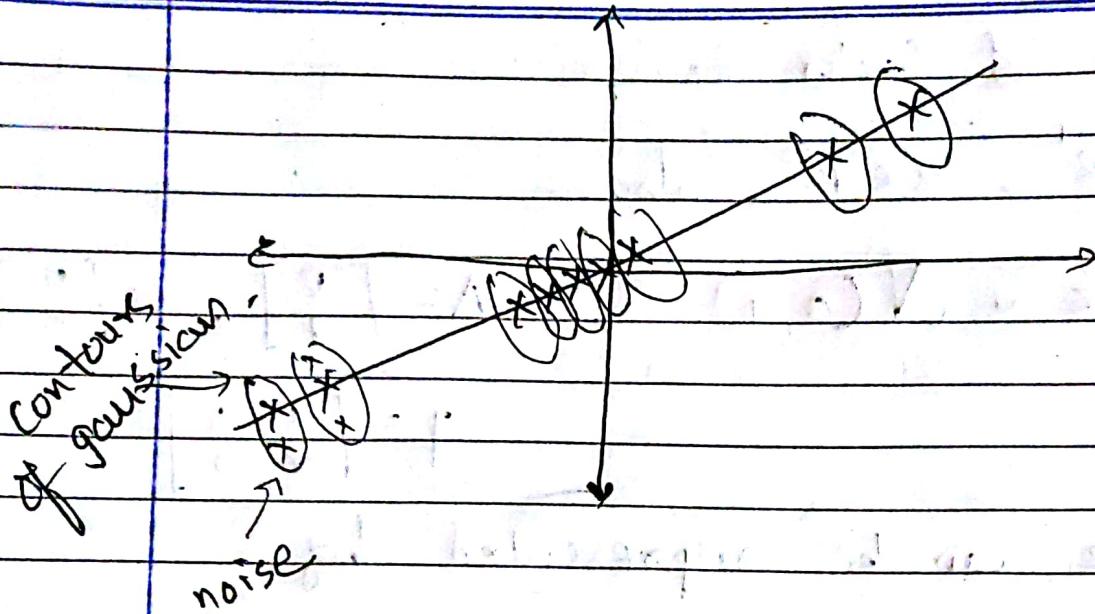
→ Now,  $Az$  will be transformed into 2D space as following.



→  $\epsilon$  is noise term, which is also gaussian. so  $x = u + Az + \epsilon$  can be represented by :-

Note - If  $x \sim N(\mu, \sigma^2)$  and  $y \sim N(0, \sigma^2)$   
 then  $x = u + y$ .

PAGE NO.	
DATE :	



→ This is the model of factor Analysis.

→ Let  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  where  $x_i \in \mathbb{R}$ ,  $x \in \mathbb{R}^{r+s}$

also;  $x \sim N(\mu, \Sigma)$

where,

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \begin{matrix} \leftarrow r \\ \uparrow \downarrow \\ \rightarrow s \end{matrix}$$

$$\text{ex. } \Sigma_{12} \in \mathbb{R}^{r \times s}$$

→  $P(x_1) = \int_{x_2} P(x_1, x_2) dx_2$  ← Marginal distribution  
 Gaussian

$$x_1 \sim N(\mu_1, \Sigma_{11})$$

PAGE NO. \_\_\_\_\_  
DATE \_\_\_\_\_

conditional  
distribution  
of Gaussian

→ also,  $\leftarrow \mathcal{N}(\mu, \Sigma)$

$$P(x_1 | x_2) = \frac{P(x_1, x_2)}{P(x_2)} \leftarrow \mathcal{N}(\mu_2, \Sigma_{22})$$

$$x_1 | x_2 \sim \mathcal{N}(\mu_{1/2}, \Sigma_{1/2})$$

$\star$   
→ where

$$\mu_{1/2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\Sigma_{1/2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

→ Under factor analysis, some random variable  $x$  and  $z$ ,

$$\begin{pmatrix} z \\ x \end{pmatrix} \sim \mathcal{N}(\mu_{zx}, \Sigma) \quad \leftarrow \mathcal{N}(\mu_z, \Psi)$$

$$z \sim \mathcal{N}(0, I)$$

$$x = \mu + \Lambda z + \varepsilon \quad \leftarrow \varepsilon \sim \mathcal{N}(0, \Psi)$$

$$\rightarrow E[z] = 0 \quad E[x] = E[\underbrace{\mu}_{0} + \underbrace{\Lambda z}_{0} + \underbrace{\varepsilon}_{0}]$$

$$E[x] = \textcircled{*} \mu$$

$$\therefore \mu_{zx} = E\begin{bmatrix} z \\ x \end{bmatrix} = \begin{bmatrix} 0 \\ \mu \end{bmatrix} I^n \stackrel{n \times d}{\leftarrow} \in \mathbb{R}$$

$$\Sigma = \begin{bmatrix} E(z - Ez)(z - Ez)^T & E(z - Ez)(x - Ex)^T \\ E(x - Ex)(z - Ez)^T & E(x - Ex)(x - Ex)^T \end{bmatrix}$$

$$\Sigma_{11} = \text{Cov}(z) = I$$

$$\begin{aligned} \cancel{\Sigma_{21}} &= E(x - Ex)(z - Ez)^T \\ &= E[(u + \Lambda z + \varepsilon - u)(z)] \\ &\quad (\because Ez = 0) \end{aligned}$$

$$\begin{aligned} &= E[\Lambda z z^T] - E[\varepsilon z] \\ &= \Lambda E[z z^T] \end{aligned}$$

$$= \Lambda \Lambda^T + \Psi$$

$$\begin{aligned} \Sigma_{22} &= E[(u + \Lambda z + \varepsilon - u)(u + \Lambda z + \varepsilon - u)^T] \\ &= \Lambda \Lambda^T + \Psi \end{aligned}$$

$$\text{So, } \begin{bmatrix} z \\ x \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \vec{0} \\ u \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Psi \end{bmatrix}\right)$$

$\mu_{zx}$

→ So, from this what is distribution of  $x$  ( $P(x) = ?$ )

$$x \sim N(\mu, \Lambda^T + \Psi)$$

→ Likelihood of  $x$  is given by

$$\prod_{i=1}^m P(x^{(i)}; \Lambda, \mu, \Psi)$$

$$= \prod_{i=1}^m \frac{1}{(2\pi)^{1/2} |\Lambda^T + \Psi|} \exp\left(-\frac{1}{2}(x - \mu)^T (\Lambda^T + \Psi)^{-1} (x - \mu)\right)$$

→ Now, if we do max. likelihood estimation, we won't be find the solution.  
[Take log and then take derivatives and set each derivative to zero]

→ So, we need to use EM algo. to estimate the parameters.

\* E-step (EM for continuous  $z^{(i)}$ ).

$$\rightarrow \text{Find } Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \Theta)$$

→ M-step:-

$$\arg \max_{\Theta} \sum_i \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \Theta)}{Q_i(z^{(i)})} dz^{(i)}$$

Expectation = Mean.

PAGE NO.

DATE :

→ Here  $z^{(i)}$  is Gaussian Random variable, so, there is Integration instead of Summation.

## EM for Factor Analysis

→ E-step:

$$z^{(i)} | x^{(i)} \sim N(\mu_{z^{(i)}|x^{(i)}}, \Sigma_{z^{(i)}|x^{(i)}})$$

where,

$$\mu_{z^{(i)}|x^{(i)}} = \vec{\theta} - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} (x^{(i)} - \mu)$$

$$\Sigma_{z^{(i)}|x^{(i)}} = I - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \Lambda$$

→ M-step:

$$\int_{Z^{(i)}} Q_i(z^{(i)}) z^{(i)} dz^{(i)}$$

→ Now, Don't write  $Q_i(z^{(i)})$  in terms of Gaussian, and try to integrate, just observe that above eqn is expectation of  $z^{(i)}$ .

$$E_{z^{(i)} \sim Q_i} [z^{(i)}] = \mu_{z^{(i)}|x^{(i)}}$$

→ So, we've

$\mu, \lambda, \psi$

$$E_{z^{(i)} \sim Q_i} \left[ \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

$$= E_{z^{(i)} \sim Q_i} \left[ \log \left( \frac{P(x^{(i)} | z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) \right]$$

$$E_{z^{(i)} \sim Q_i} \left[ \log \frac{P(z^{(i)})}{Q_i(z^{(i)})} \right]$$

only this term depends on parameters, so only maximize it.

→ So, In M-step, maximize

$$\max \sum_{i=1}^m E_{z^{(i)} \sim Q_i} \left[ \log \frac{P(x^{(i)} | z^{(i)}; \mu, \lambda, \psi)}{Q_i(z^{(i)})} \right]$$

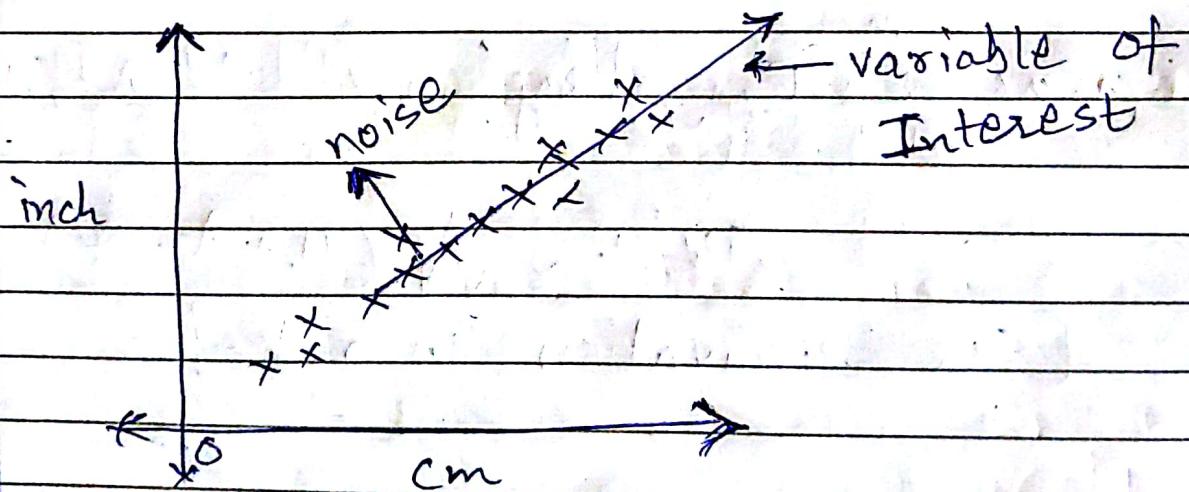
$$x^{(i)} | z^{(i)} \sim N(\mu + \lambda z, \psi)$$

[See Derivation in notes].

## \* Principle Component Analysis (PCA).

→ Given  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ ,  $x^{(i)} \in \mathbb{R}^n$ , in this algo. we reduce the dimension of each dataset  $x^{(i)}$  from  $n$  to  $k$ , where ( $k < n$ ).

→ Why this is required?  
For ex- if we've collected the ~~the~~ data of height of persons in both 'cm' and 'inch'. So this is redundant dimension.



→ The graph will be a straight line (not exact straight line due to rounding of cm to inch or vice-versa).

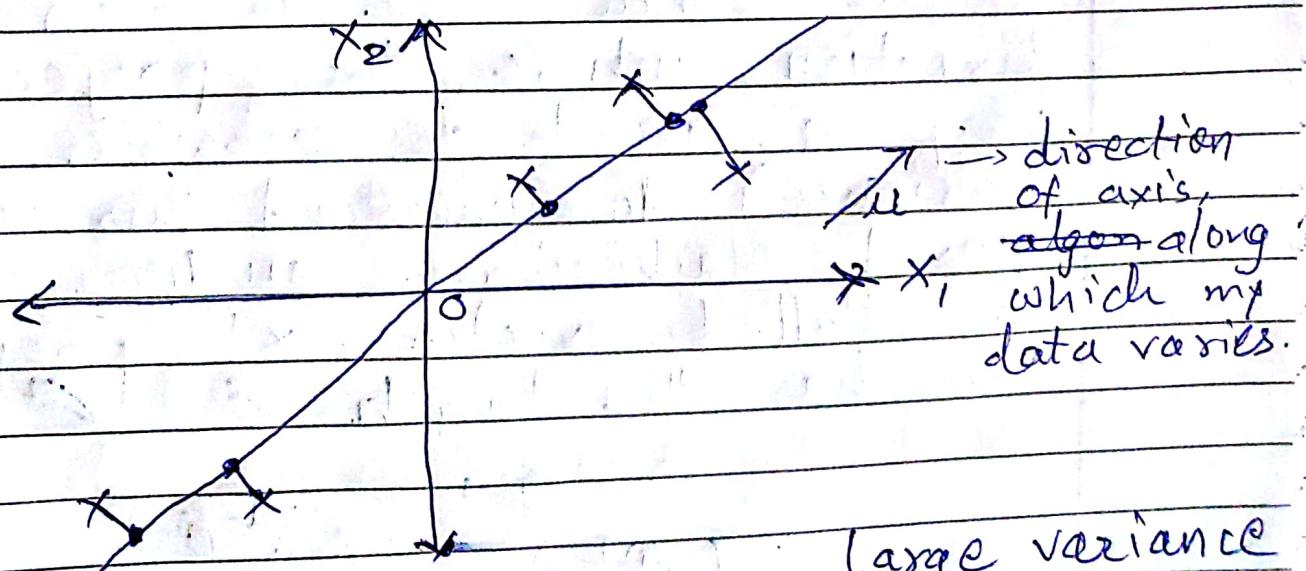
AB

→ preprocessing the data before running the algo:-

- 1.) set  $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$  } zero out the mean
- ~~2.) Replace~~
- 2.) Replace  $x^{(i)}$  with  $x^{(i)} - \mu$  }
- 3.) set  $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)})^2$  } Normalize to unit variance.
- 4.) Replace  $x_j^{(i)}$  with  $\frac{x_j^{(i)}}{\sigma_j}$  }

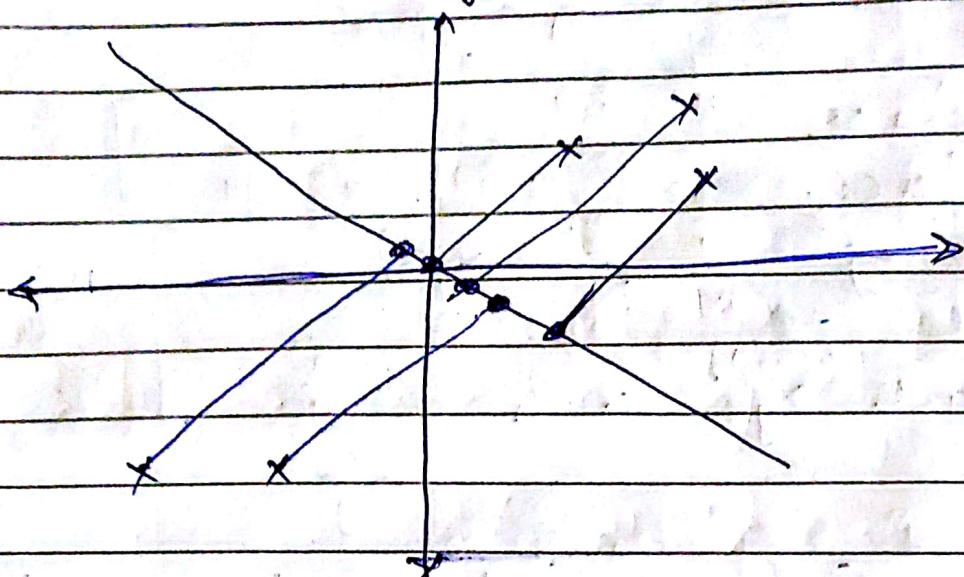
here  $j$  represent feature. So,  $\sigma_j^2$  is variance of feature  $j$ .

→ Now, how to find the mean axis or principle axis, alg. along which the data varies?



→ Note that the projection (dots) have very ↑  
(large variance)

→ If we choose some other direction, like following,



- Then the projections (dots) have very low variance.
- So to formalize this idea of how to find the axis, we can say that the axis should be in such a direction, where the projection (dots) vary as much as possible (high variance). To understand this, read the saved stackexchange answer of PCA.
- If  $\|u\|=1$ , vector  $x^{(i)}$  projected on  $u$  has the length  $= \underline{\underline{x^{(i)} u}}$

$$(x^{(i)T} u)$$

$(x^T u)$  is dot product and hence scalar. see the stack-exchange question for confusion.

→ So, our problem will be

choose  $u$  such that:

$$\max_{\|u\|=1} \frac{1}{m} \sum_{i=1}^m (x^{(i)T} u)^2$$

$$= \max_{\|u\|=1} \frac{1}{m} \sum_{i=1}^m (u^T x^{(i)}) (x^{(i)T} u)$$

$$= \max_{\|u\|=1} u^T \left[ \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \right] u.$$

Now  $\frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}$  is co-variance matrix.  $\Sigma$ .

$$\text{so, } \Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}$$

$$\text{now, } \max_{\|u\|=1} u^T \Sigma u.$$

→ We can apply lagrange's multipliers to maximize constrained optimization problem.

$$\text{so, constraint is } \|u\| = 1 = u^T u$$

→ Lagrangian will be ~~Lagrangian~~ <sup>Lagrange multiplier</sup>

$$L(u, \lambda) = u^T \Sigma u - \lambda(u^T u - 1)$$

→ differentiating  $L$ ,

$$\cancel{\frac{\partial L}{\partial u}} \cdot \Sigma u - \lambda u = 0$$

→ which is eigen-vector equation.  
Substitute this in objective function,

$$\begin{aligned} & u^T \Sigma u - \lambda(u^T u - 1) \\ &= u^T \Sigma u - \lambda u^T u + \lambda \\ &= u^T \Sigma u \quad (\because u^T u = 1) \\ &= \lambda u^T u \\ &= \lambda \end{aligned}$$

→ So, to maximize  $f'$ ,  $\lambda$  must be the largest eigen value.

→ For any matrix  $A$ , and vector  $u$ ,  
if

$$Au = \lambda u, \text{ then}$$

$\lambda$  is eigen value and  
 $u$  is eigen vector of  $A$ .

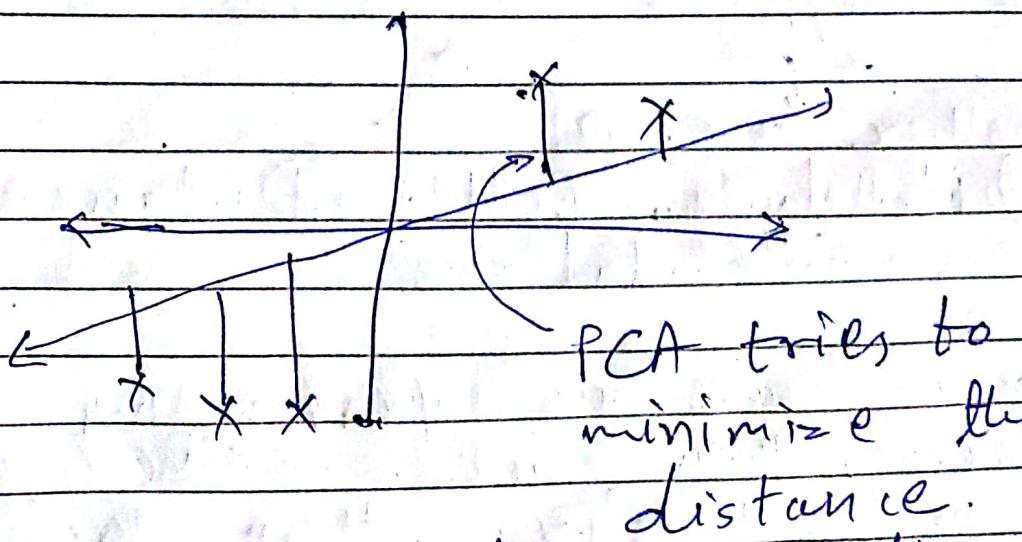
→ More generally, if we want ~~to~~  
 k-dimensional subspace, then  
 choose  $u_1, u_2, \dots, u_k$  to be k top  
 eigenvectors of  $\Sigma$ . That means  
 eigenvectors corresponding to top  
 k eigen values.

→ So, If we've  $x^{(i)} \in \mathbb{R}^n$ , now new  
 representation of data in  
 $\{u_1, u_2, \dots, u_k\}$  basis will be

$$y^{(i)} = (u_1^T x^{(i)}, u_2^T x^{(i)}, \dots, u_k^T x^{(i)}),$$

$$y^{(i)} \in \mathbb{R}^k, 1 \leq k \leq n$$

→ Another way to think of PCA is that  
 it tries to minimize the distance  
 between data points and the axis.



→ Which is also called reconstruction error

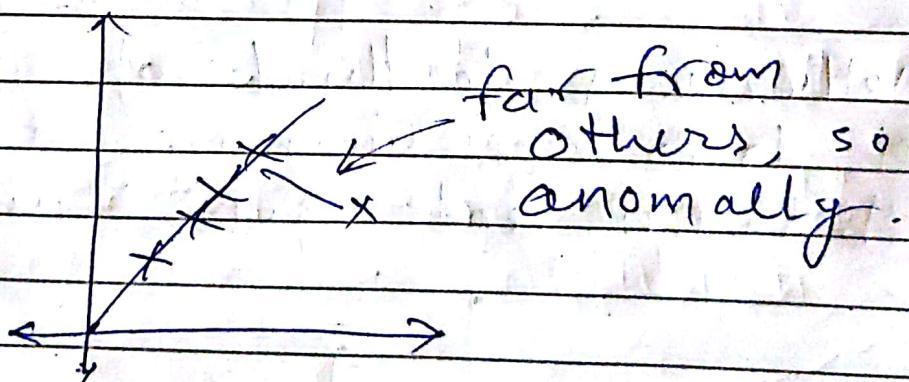
## \* Applications of PCA:

1) Visualization: Very often we've very high-dimensional data sets. It's not possible to plot it. So PCA is used to reduce the dimension & plot & visualize the datasets.

2) Compression.

3) Learning algo. runs much faster. Also, with PCA, we've less features, so less prone to overfitting.

4) Anomaly Detection.



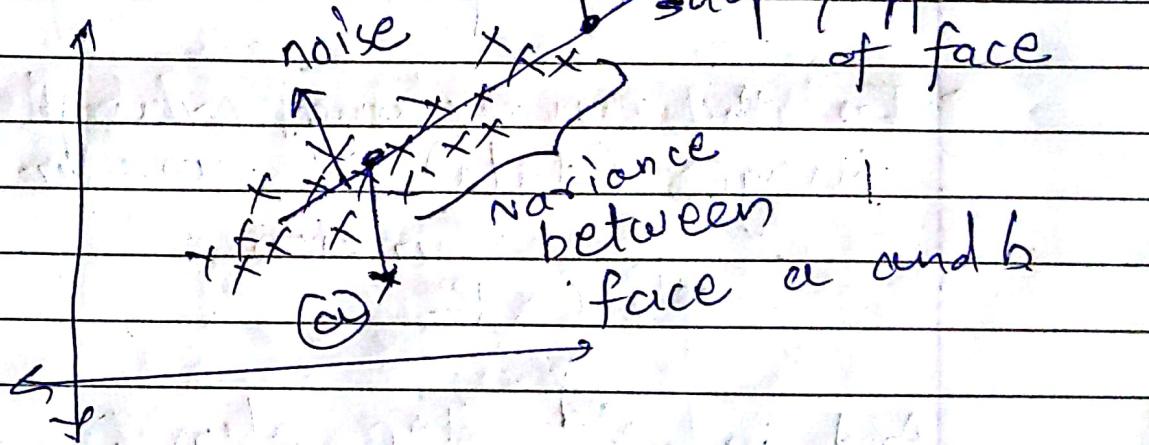
5) Matching, Better Distance Calculations.

→ For ex. we are performing face recognition and we've  $100 \times 100$  images of faces, then each data point  $x^{(i)}$  will be of

$100 \times 100 = 10,000$  dimensional vector

$$x^{(i)} \in \mathbb{R}^{10,000}$$

→ If we reduce it to ~~100~~ approx. 50 dimensions, then we can use it efficiently to match different faces.



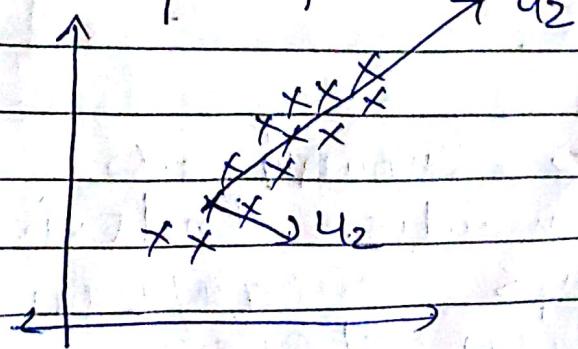
$$x^{(i)} \in \mathbb{R}^{10,000}$$

$$u^{(i)} \in \mathbb{R}^{10,000}$$

eigen faces.

Lecture - 15      15/7/18.

Summary of PCA



### Algorithm

1) Normalize the data to zero mean, unit variance.

2) find covariance Matrix as

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)T}$$

3) Find top k eigenvectors of  $\Sigma$ .

→ In the matching example, if our images (data) is  $100 \times 100$ , then

$$\mathbf{x}^{(i)} \in \mathbb{R}^{100 \times 100}$$

And in 2<sup>nd</sup> step above, our Co-variance matrix will be  $10000 \times 10000$ .

$\Sigma \in \mathbb{R}^{10000 \times 10000}$ , which is of very high dimension.

→ Another use of PCA is to measure similarity between two text documents.

→ This algo. is known as LSI (Latent Semantic Indexing), which usually skips mean/variance normalization steps.

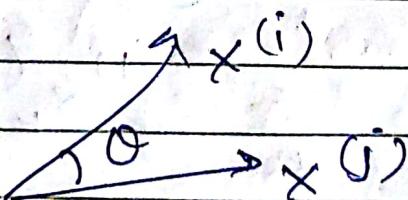
- To measure the similarity between two documents  $x^{(i)}$  and  $x^{(j)}$ , which are represented as a very high-dimensional feature vector like following:

1	a
0	aardwark
0	aardwold
1	learn

where  $x \in \mathbb{R}^{50,000}$  (size of dictionary).

- Then to measure the similarity,

$\text{sim}(x^{(i)}, x^{(j)})$ , we can view these two vectors as following



- One measure of similarity b/w two documents is the angle b/w the above vectors.

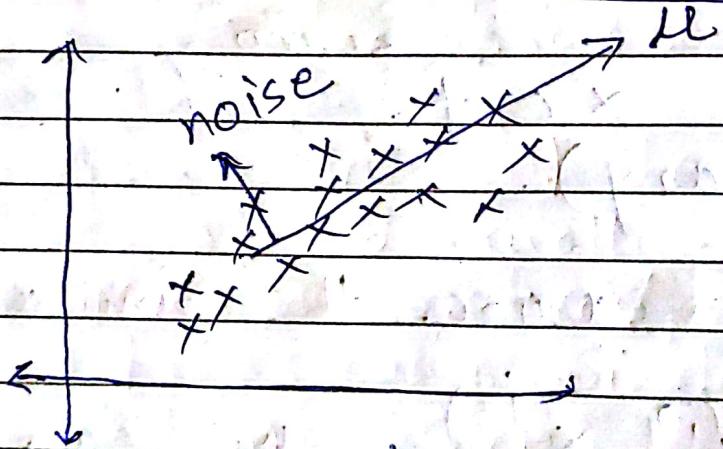
- If  $\theta$  is small  $\rightarrow$  Documents are similar. & vice versa.

→ So, similarity is  $\cos\theta \cdot$  [Angle increases  $\rightarrow$  Similarity decreases, and vice-versa].

$$\cos\theta = \frac{\mathbf{x}^{(i)\top} \mathbf{x}^{(j)}}{\|\mathbf{x}^{(i)}\| \|\mathbf{x}^{(j)}\|}$$

[from Definition of Dot Product]

→ Intuition behind LSI.



→ Numerator of  $\text{sim}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ ; is

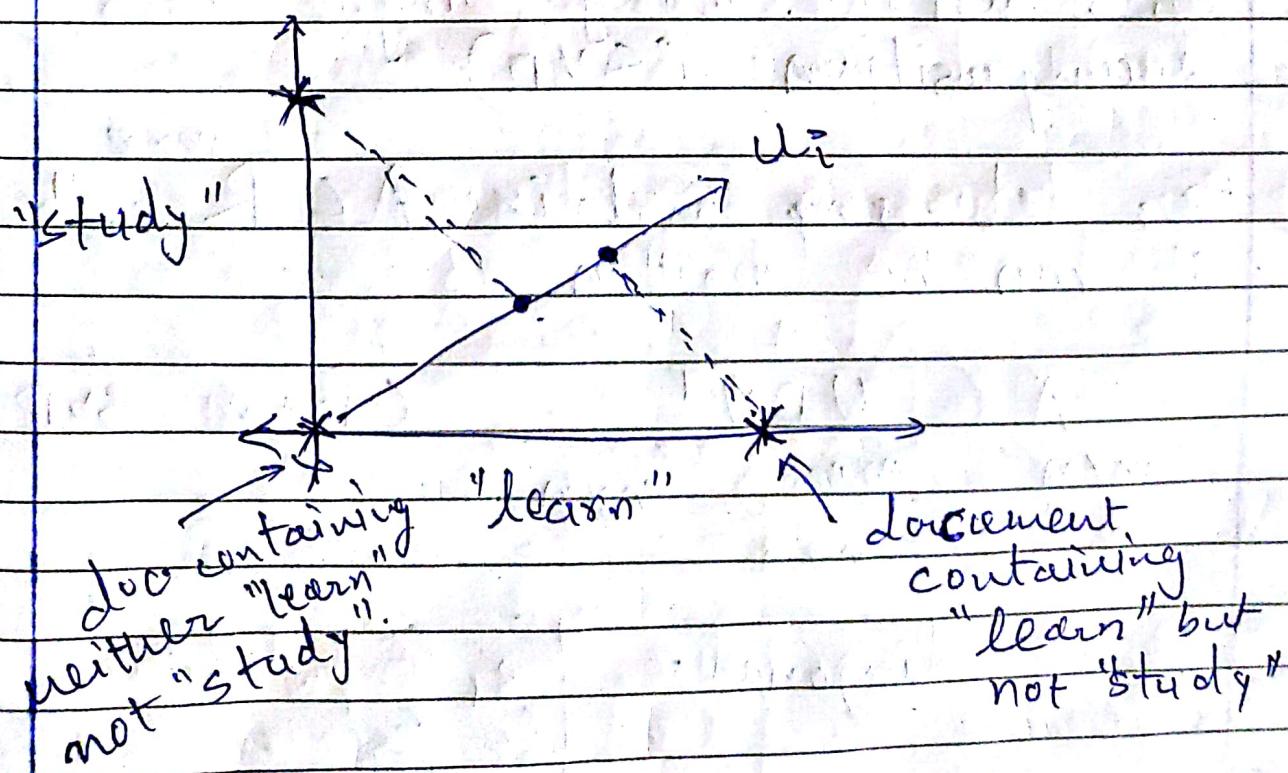
$$\mathbf{x}^{(i)\top} \mathbf{x}^{(j)} = \sum_k x_k^{(i)} x_k^{(j)}$$

which will be zero if documents don't have any word in common. [Dot product], so it can be written as

$$x^{(i)T} x^{(j)} = \sum_k x_k^{(i)} x_k^{(j)}$$

Document i and j both contain word k?

- and in above formula,  $x_k^{(i)}$  means whether document i contains word k, and same for  $x_k^{(j)}$ .
- E.g. if  $x^{(i)} = "study"$  and  $x^{(j)} = "learn"$ , then both documents are considered completely dissimilar documents.



→ So, from LSI, we can say that there is +ve similarity between these two documents.

→ But in such cases,  $\Sigma$  will be very very high dimensional matrix. For ex -

$x^{(i)} \in \mathbb{R}^m$ , then co-variance matrix will be of dimension  $50,000 \times 50,000$

$$\Sigma \in \mathbb{R}^{m \times m}$$

So, applying PCA will be difficult for such cases. So another way to apply PCA is to decompose the data matrix using Single value decomposition (SVD)

→ Ex: let any matrix  $A \in \mathbb{R}^{m \times n}$ , then it can be written as

$$A = U D V^T \quad (\text{from SVD})$$

$m \times n \quad m \times n \quad \downarrow \quad n \times n$

$$D \text{ is diagonal} = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & \sigma_n \end{bmatrix}$$

$\sigma_i$  in D matrix is called singular values of matrix A.

Eg  $\begin{bmatrix} A \\ mxn \end{bmatrix} = \begin{bmatrix} U \\ mxn \end{bmatrix} \begin{bmatrix} D \\ n \times n \end{bmatrix} \begin{bmatrix} V^T \\ n \times n \end{bmatrix}$

U's columns - Eigenvectors of  $A A^T$   
 V's columns - Eigenvectors of  $A^T A$

→ SVD can be computed by "svd" in matlab or octave.

→ We written  $\Sigma$  as

$$\Sigma = \sum_{i=1}^n \sigma^{(i)} x^{(i)T} \quad (\text{Ignore } L \text{ for } n-1 \text{ row})$$

→ Now, let's writer our data points  $x^{(i)}$  as rows in design matrix X.

$$X = \begin{bmatrix} \cdots & x^{(1)} & \cdots \\ \cdots & x^{(2)} & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & x^{(m)} & \cdots \end{bmatrix}$$

then co-variance matrix can be written as

$$\Sigma = X^T X \quad (\gamma_{n-1} \text{ ignored for now})$$

$$= \begin{bmatrix} & & & \\ & & & \\ & x^{(1)} & x^{(2)} & \dots \\ & & & \end{bmatrix} \begin{bmatrix} -\bar{x}^{(1)} \\ -\bar{x}^{(2)} \\ \vdots \end{bmatrix}$$

→ Now, to get top k eigen vectors,

$$X = UDV^T$$

→ Top k columns of  $\mathbf{U}$  are the top k eigen vectors of  $X^T X = \Sigma$ .  
So this process is same as PCA.

→ So, if we've m training example, and 50,000 features (dictionary size), then dimension of  $X$  would be

$$m \times 50,000.$$

$$X \in \mathbb{R}^{m \times 50,000}$$

→ So, this require less computation power compared to  $50000 \times 50,000$  matrix  $\Sigma$ .

→ So, if we've very high dimensional data then we should apply SVD to perform PCA.

\* PCA, Factor Analysis, and differences and when to use each one?

	Model $P(x)$	Not Probabilistic
"similarity"	"Subspace" Factor Analysis	PCA
"Group clumpings"	Mixture of Gaussian	K-means differences

- When we want to reduce dimensions, use PCA.
- When we want anomaly detection, use factor analysis.
- FA is for modeling the data / model
- Both assumes, that the data lies in subspace.
- When Data lies in different Groups / classes, and we want density estimation, use Mixture of Gaussian and same for K-means

## \* Independent Component Analysis (ICA):

→ In PCA, we're finding the main (principle) axes of variation of data, whereas in ICA, we are finding independent components of variation of data.

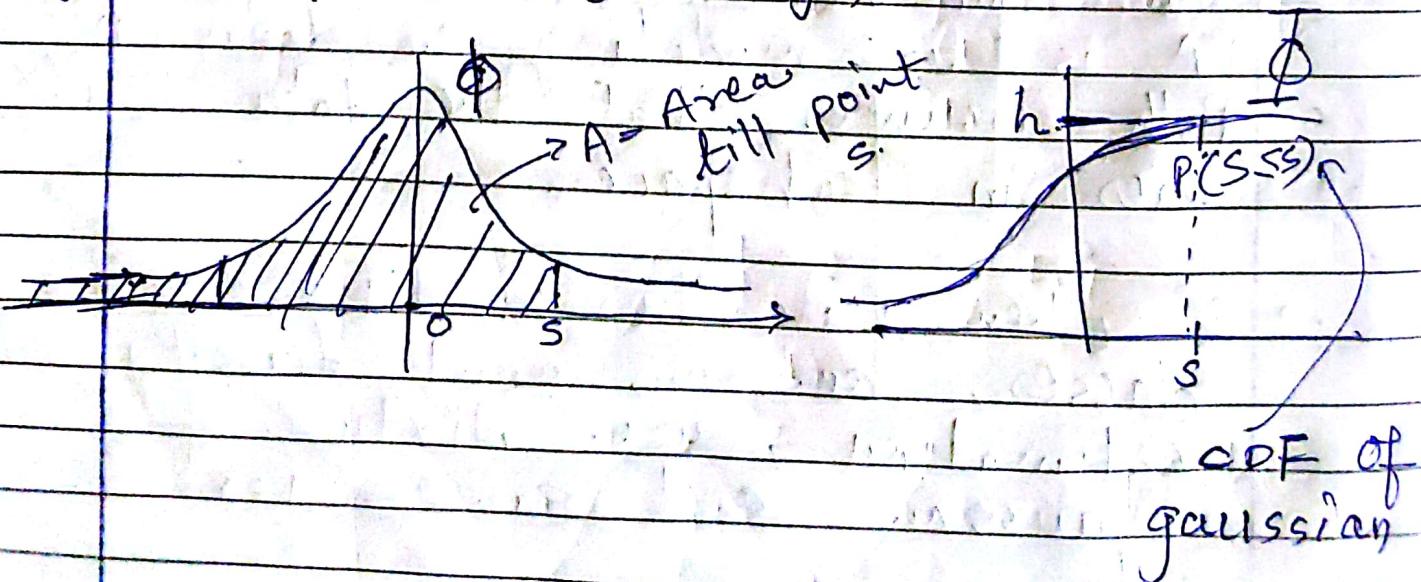
→ One example where we can use ICA is 'cocktail party' problem.

\* CDF (Cumulative distribution function):

→ Let's say we've one-dimensional random variable  $s$  and probability density function  $P_s(s)$ , then

$$\text{CDF } F(s) = P(S \leq s)$$

E.g. For Gaussian density;



→ Height of  $P(S \leq s)$  in  $\Phi(h)$  is equal to the area covered by  $\phi$  (small  $\phi$  means gaussian bell shaped curve). till the point  $s$ .

$$\therefore A = h.$$

→ In CDF, as height approaches to 1 ( $x \rightarrow \infty$ ), the area also approaches to 1.

→ From above intuition, Definition of CDF can be given by:-

→ For any real-valued random variable  $S$ , or distribution function of  $S$ , evaluated at  $s$ , CDF ~~area~~ is the probability that  $S$  will take a value less than or equal to  $s$ .

→ So, ~~area~~ for continuous distribution functions, CDF is given by

$$CDF F(s) = \int_{-\infty}^s P_S(t) dt.$$

→ suppose we've a random variable  $S$  and we want to model the distribution of Random Variable  $S$ , then,

$s \rightarrow$  specify  $P_s(s)$  or  
 $\rightarrow$  specify  $F(s)$

$\Rightarrow$  Because both are related.

$$P_s(s) = F(s)$$

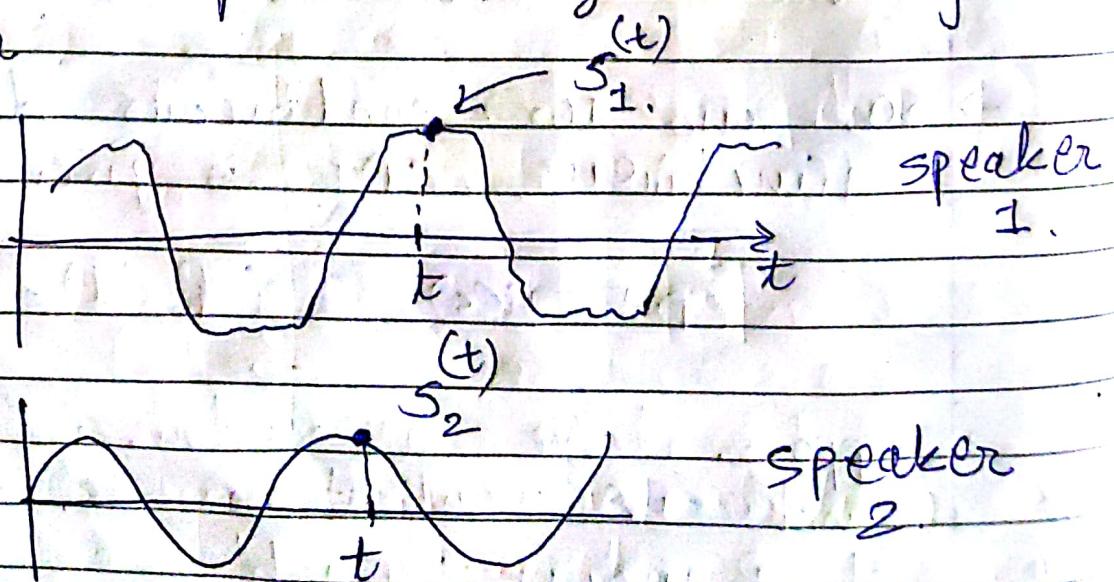
\* ICA :-

$\rightarrow$  Suppose there are  $n$  original data sources. (e.g.  $n$  different speakers in cocktail party).

$$s \in \mathbb{R}^n$$

$s_j^{(i)}$  means signal from speaker  $j$  at time  $i$ .

$\rightarrow$  Sound is represented by following in graph



→ We observe,

$$\mathbf{x}^{(i)} = \mathbf{A}\mathbf{s}^{(i)}, \text{ where } \mathbf{x}^{(i)} \in \mathbb{R}^n$$

assume we've  $n$  microphones. ( $\mathbf{x}^{(i)}$ )

→ So, ~~above~~ below equation means that each microphones records some overlapping linear combination of speakers.

→ For ex.

$\mathbf{x}_j^{(i)}$  (microphone  $j$  records at time  $i$ ).

$$\mathbf{x}_j^{(i)} = \sum_k A_{jk} s_k^{(i)}$$

→ Goal :- Find a matrix  $\mathbf{W} = \mathbf{A}^{-1}$ , so that,

$$\mathbf{s}^{(i)} = \mathbf{W}\mathbf{x}^{(i)}$$

→  $\mathbf{W}\mathbf{x}^{(i)}$  is original recordings of speaker  $s^{(i)}$ .

→ notation of  $\omega$ ,

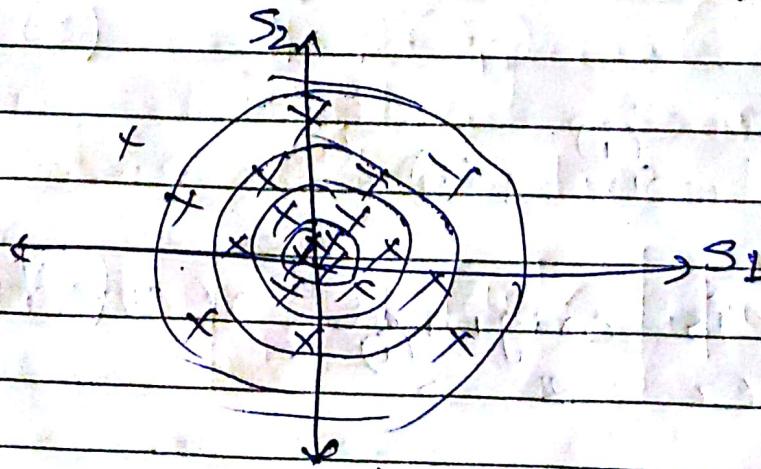
$$\omega = \begin{bmatrix} \omega_1^T \\ \omega_2^T \\ \vdots \\ \omega_n^T \end{bmatrix}$$

### \* Ambiguities of ICA

- 1) Order or permutation of speakers recovered not guaranteed.
- 2) sign of the speakers.

→ The reason for only above 2 ambiguities is that  $s_j^{(i)}$  is non-gaussian.

→ But suppose that  $s_j^{(i)} \sim N(0, I)$ , then the data contours would be :-



→ Now, as gaussian distribution is rotationally symmetric, we can't

recover the original speaker sounds. Because, there is no way we can tell the two axes  $s_1$  and  $s_2$  located.

→ Let  $s \in \mathbb{R}$

Density of  $s$ :  $P_s(s)$

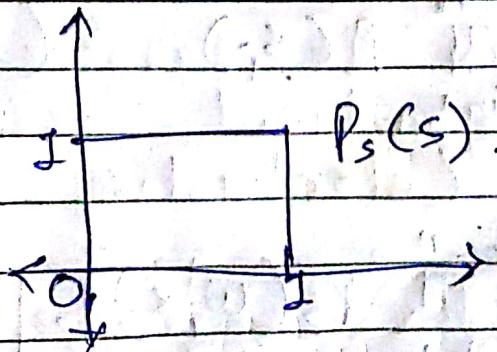
$$x = As = w^{-1}s, \quad s = w x$$

→  $A$  is mixing matrix and  $w$  is unmixing matrix.

$$P_x(x) = P_s(wx) |w|$$

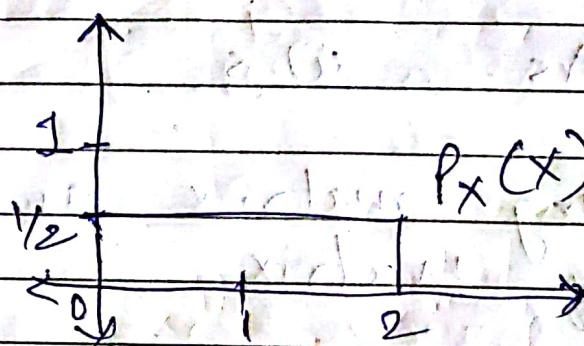
E.g.  $P_s(s) = I\{0 \leq s \leq 1\}$

( $s \sim \text{Uniform}[0, 1]$ )



Let  $x = 2s$ , then  $A = 2$ ,  $\omega = \frac{1}{2}$

- now,  $s$  is uniform distribution over  $[0, 1]$ , so  $x = 2s$  will be uniform distribution over  $[0, 2]$ .
- So, density for  $x$  would be



$$p_x(x) = 1 \cdot \{0 \leq x \leq 2\} / 2$$

- ~~ICA~~ ICA model:  $n$  speakers.
  - \* ICA model:  $n$  speakers.
- $$P(s) = \prod_{i=1}^n P_s(s_i) \quad (\text{$s_i$ are independent})$$

$$P(x) = \left[ \prod_{i=1}^n P_s(w_i^T x) \right] / |w|$$

$$\therefore \omega = A^{-1} = \begin{bmatrix} -\omega_1^T \\ \vdots \end{bmatrix}, \text{ so, } s_i = \omega_i^T x,$$

→ Now, choose density (probability),  $\omega_i$  of what each speaker is speaking.

$$P_s(s_i) = ?$$

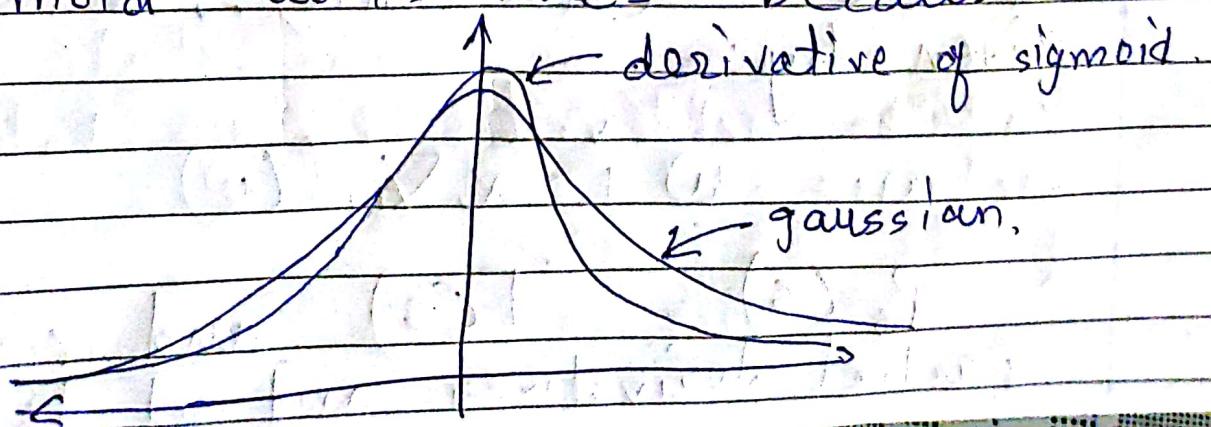
→ As discussed in CDF, we can either choose CDF or density.

→ So, let's choose CDF to be sigmoid function. We can't choose gaussian as mentioned before.

$$\text{So, } F(s) = \frac{1}{1 + e^{-s}}$$

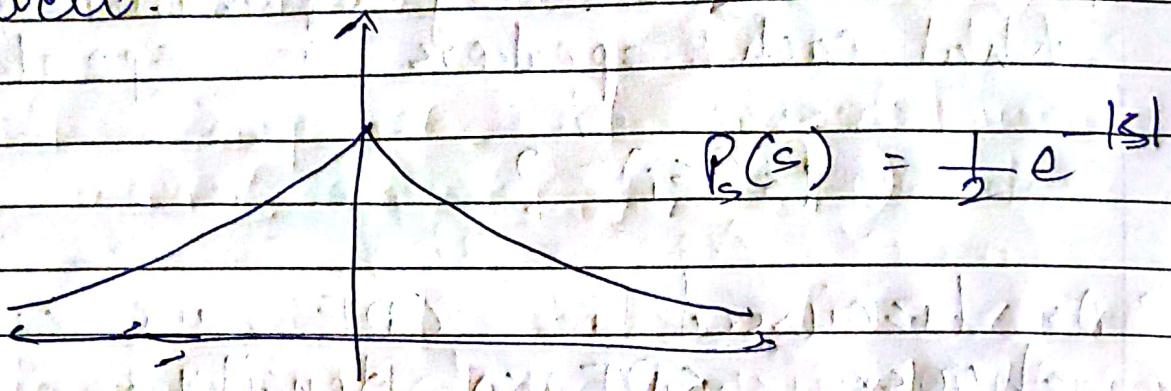
$$P_s(s) = F'(s_i)$$

→ sigmoid works fine because



→ Gaussian decreases by  $e^{-s}$ , while sigmoid's derivation by  $e^{-s^2}$ .

→ Also, Poisson distribution works well.



→ Given  $\{x^{(1)}, \dots, x^{(m)}\}$  training set, log likelihood of examples  $\omega$ 's given by

$$l(\omega) = \sum_i \log \left( \prod_j P_s(\omega_j x^{(i)}) \right)$$

→ And to learn the parameters  $(\omega)$ , use the gradient descent ascent. (stochastic), which is given by:-

$$\omega := \omega + \alpha \nabla_{\omega} l(\omega)$$

→  $P_s(s) = F(s)$  and compute partial derivatives and set to zero,

we find,

$$\nabla_{\omega} ELL(\omega) = \begin{bmatrix} 1 - 2g(\omega^T x) \\ 1 - 2g(\omega_0^T x) \end{bmatrix} + (\omega^T)^{-1}$$

→ when we find  $\omega$ , we can recover each source by

$$s^{(i)} = \omega x^{(i)}$$

## Lecture-16

### \* Reinforcement Learning

→ RL problems model the world using MDP. (Markov decision process)

→ MDP comprises of five elements  
 $MDP(S, A, \{P_{sa}\}, \gamma, R)$

$S \rightarrow$  set of states [In case of autonomous helicopter, set of positions and set of orientations of helicopter]

$A \rightarrow$  set of Action [set of positions in which we can control the helicopter]

$P_{sa} \rightarrow$  State transition distributions

$$\sum_{s'} P_{sa}(s') = 1 \text{ and } P_{sa}(s') \geq 0$$

~~$P_{sa}$  means in what state ( $s'$ ) will be achieved if we take action 'a'.~~

- $P_{sa}$  gives the probability distribution of what state we would transition to, if we take action 'a' in state 's'.
- So,  $P_{sa}(s')$  means the probability distribution of state  $s'$  achieved by us, ~~then~~ when we take action 'a' from state 's'.

→ Discount factor;  $0 \leq \gamma < 1$

R → Reward function.

$$R: S \rightarrow \mathbb{R}$$

$\mathbb{R}$  set of real numbers.

- Example problem:

3				+1
2				-1
1				

1    2    3    +    -1

position (3,1)

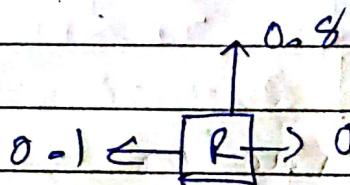
→ MDP for this problem.

$S \rightarrow 11$  states.

$A = \{ \text{North, South, East, West} \}$

→ Robot dynamics for moving in this grid.

If we command the robot to go north,



→ It's 80% chance in which robot goes north, 10% in east & west.

$$\rightarrow P((3,1) \xrightarrow{\text{N}} (3,2)) = 0.8$$

$$P((3,1) \xrightarrow{\text{N}} ((4,1))) = 0.1$$

$$P((3,1) \xrightarrow{\text{N}} ((2,1))) = 0.1$$

$$P((3,1) \xrightarrow{\text{N}} (3,3)) = 0$$

∴ etc.

$$\rightarrow R((4,3)) = +1$$

$$R((4,2)) = -1$$

$$R(s) = -0.02$$

for all other states

Now let's see how MDP works.

- At time zero ( $t_0$ ), and initial state  $s_0$ , Robot starts. It's task is to choose an action  $a_0$ .
- Depending on choice, we get some state,  $s_1 \sim P_{s_0, a_0}$ .
- Now choose action  $a_1$  from new state  $s_1$ .
- So, it get to new state  $s_2 \sim P_{s_1, a_1}$ .
- After doing this for a while, Robot has visited a sequence of states  $s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow \dots$ .
- And to see how well the robot did, sum up the rewards it got.

$$s_0, R(s_0) + R(s_1) + R(s_2) + \dots$$

→ Let's change this total reward as following:

$$\text{Total Payoff} = \gamma^0 R(s_0) + \gamma^1 R(s_1) + \gamma^2 R(s_2) + \dots$$

$\gamma \in [0, 1)$

→ So, what this does is control the reward with time

→ So, Reward we got at time zero (initially on step  $s_0$ ) is more than step  $s_1$  and so on.

→ So the goal of RL algo. is to choose the actions over time ( $a_0, a_1, a_2, \dots$ ), to maximize the expectation of total payoff.

$$E[R(s_0) + \gamma^1 R(s_1) + \gamma^2 R(s_2) + \dots]$$

→ RL algo computes the policy

$$\pi : S \rightarrow A$$

→ Policy is just the mapping from state to action so it tells the algo. that when you are at this state, perform that action

→ Optimal policy for this example is given by

3	→	→	→	+1
2	↑		↑	-1
1	↑	←	↖	↖

→ So, we see when we execute the optimal policy, it will give max. total pay off.

→ MDP is good at making certain trade offs. consider the position (3,1), here we can go North instead of following the longer path:

- But if we go north, there are certain chances (0.1) that Robot end up going right side (Bcz of Robot dynamics we defined previously)
- So, discount factor discourage this and rather takes longer route.

### \* Definitions -

→  $V^{\pi}$  - for any given policy  $\pi$  define value function  $V^{\pi}: S \rightarrow \mathbb{R}$  such that

$V^{\pi}(s)$  is expected total payoff, starting in state  $s$ , and execute policy  $\pi$ .

$$V^{\pi}(s) = E[\text{Total payoff} | \pi, s_0 = s]$$

$\pi$  is not a random variable but the executed policy.