

CH - Dimensionality Reduction and visualization

* Dimensionality reduction:- We learned that we can use scatter plot for 2-D, 3-D, for 4-D, 5-D, we can use pair plot, But for n-D, we try to reduce the Dimensionality of the data.

→ Row vector and column vector :-

$x_i \in \mathbb{R}^d \rightarrow d\text{-dim. column vector}$

$\underbrace{y}_{\text{Real Number}}$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}$$

→ Dataset :- $D = \{x_i, y_i\}_{i=1}^n$

$x_i \in \mathbb{R}^d, y_i \in \{\text{setosa, versicol, virginica}\}$

↓ ↓
data points class label

the dataset contains n data points
 x_i is d dim vector of Real values.
e.g. In iris dataset $x_i \in \mathbb{R}^4$

→ Data as a data-matrix

$$D = \{x_i, y_i\}_{i=1}^n$$

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

$$y_i \in \{\text{Set, Verg, Ver}\}$$

$$\Rightarrow \begin{array}{c|cccc|c} & f_1 & f_2 & \dots & f_d & Y \\ \hline 1 & \xleftarrow{x_1^T} & & & \xrightarrow{x_1^T} & | \\ 2 & \xleftarrow{x_2^T} & & & \xrightarrow{x_2^T} & | \\ 3 & & & & & | \\ \vdots & & & & & | \\ n & \xleftarrow{x_n^T} & & & \xrightarrow{x_n^T} & | \\ & & & & & | \\ & & & & & n \times d \\ & & & & & n \times 1 \end{array}$$

→ Data preprocessing : column-normalization

$$\begin{array}{c|c|c} & f_1 & f_2 & \dots & f_d & Y \\ \hline x_1 & \xleftarrow{x_1^T} & & & \xrightarrow{x_1^T} & | \\ x_2 & \xleftarrow{x_2^T} & & & \xrightarrow{x_2^T} & | \\ \vdots & & & & & | \\ x_n & \xleftarrow{x_n^T} & & & \xrightarrow{x_n^T} & | \\ & & & & & | \\ & & & & & n \times d \\ & & & & & n \times 1 \end{array}$$

preprocessing refers some type of mathematical operation performed on data before doing real task.
Column-normalization is one of the preprocessing technique.

Column normalization:- we know each column ~~expected feature~~, take a column, and take all the values of that column as follows.

$$\begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_d \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

take max value of a_i in f_j

$$f_j = \left[\frac{a_1}{\min(a_1)}, \frac{a_2}{\min(a_1)}, \dots, \frac{a_n}{\min(a_1)} \right]$$

$$a_{\min} \quad a_{\max}$$

$$a'_i = \frac{a_i - a_{\min}}{a_{\max} - a_{\min}}$$

by doing this we are ensuring that $a'_i \in [0, 1]$.

Do this for each column.

But why are we doing this? By doing this we are getting no scale, every thing comes in $[0, 1]$ scale.

* Mean of Data Matrix:-

$$x = \begin{bmatrix} f_1 & f_2 & f_3 & \dots & f_d \end{bmatrix} \quad x_i \in \mathbb{R}^d$$

$$x_1 = [2.2, 4.2] \in \mathbb{R}^2$$

$$x_2 = [1.2, 3.2] \in \mathbb{R}^2$$

$$x_1 + x_2 = [3.4, 7.4]$$

$$\bar{x} = \left[\frac{3.4}{2}, \frac{7.4}{2} \right]$$

$$\therefore \bar{x} = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

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

* Data precession: Column Standardisation.

We already seen a precessing technique i.e. column stand. col-stand. is more often used ~~than~~ than normalisation just like Col. Nor. we scale all column one by one

$$y = \left[\dots \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \dots \right] \rightarrow [a_1, a_2, \dots, a_n]$$

\downarrow col-std.

$$a'_1, a'_2, \dots, a'_n \quad (\leftarrow \text{mean of } \{a'_i\}_{i=1}^n)$$

In column stands we convert the column such that the mean of converted data = 0 & $\sigma = 1$

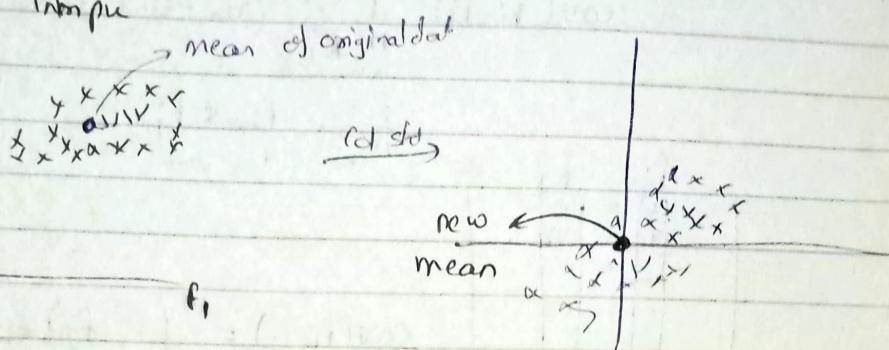
Let's see how to do it

$$\text{Let } \bar{a} = \text{mean } \{a_i\}_{i=1}^n \quad \leftarrow \text{Sample mean}$$

$$s = \text{std } \{a_i\}_{i=1}^n \quad \leftarrow \text{std-dev}$$

$$a'_i = \frac{a_i - \bar{a}}{s} \rightarrow \text{This way we are guaranteeing that mean of } a'_i = 0 \text{ and } \sigma = 1$$

Geometrical Intuition



The mean moves to origin and σ at any axis becomes 1.

* Covariance of Data matrix

What is Data matrix?

$$D = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}_{n \times d} \quad \text{for this Data matrix we have a covariance matrix of size } d \times d$$

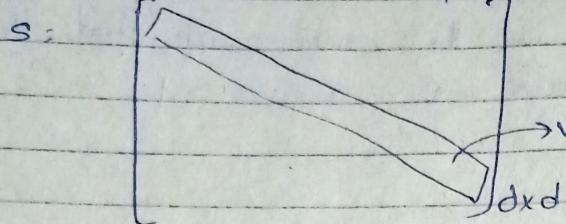
$$S_{i,j} = \text{cov}(f_i, f_j) \quad \left\{ \begin{array}{l} S_{i,j} = \text{covariance of feature } i \times \text{feature } j \\ i: 1 \rightarrow d \\ j: 1 \rightarrow d \end{array} \right.$$

$$\text{We know that } \text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$(\text{cov}(x, x) = \text{var}(x))$$

$$\text{cov}(f_i, f_j) = \text{cov}(f_j, f_i)$$

square matrix
symmetric matrix



variance of features

d x d

let assume our dataset $X =$

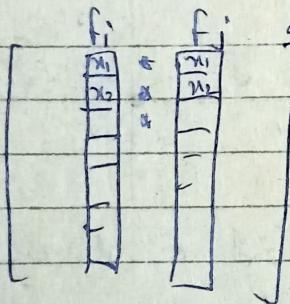
$$\begin{bmatrix} f_1 & f_2 & \dots & f_d \end{bmatrix}^T \quad n \times d$$

if suppose X is column standardized,

that means mean of each col = 0 $\times \sigma = 1$
hence

$$\text{cov}(f_i, f_j) = \frac{1}{n} \sum_{k=1}^n (f_{ik} - 0)(f_{jk} - 0)$$

$$= \frac{1}{n} \sum_{k=1}^n f_{ik} * f_{jk}$$



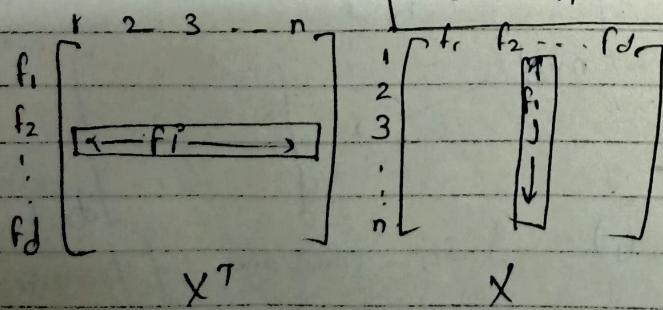
$$\text{cov}(f_i, f_j) = (f_i^T f_j) * \frac{1}{n}$$

If f_i & f_j have been standardised then:

$$\text{cov}(f_i, f_j) = \frac{f_i^T f_j}{n}$$

Hence If data set is column standardized then

covariance matrix, $S_{d \times d} = \frac{1}{n} (X^T) X$



$$X^T \cdot X = \begin{bmatrix} f_1 + f_1 & f_1 + f_2 & f_1 + f_d \\ f_2 + f_1 & f_2 + f_2 + \dots & f_2 + f_d \\ \vdots & \vdots & \vdots \end{bmatrix}$$

This is not S
But S

Note: $X = \begin{bmatrix} f_1 & f_2 & f_3 & \dots & f_d \\ 1 & 2 & 3 & \vdots & n \end{bmatrix}$

$S = \text{cov_matrix} = \begin{bmatrix} f_1 & f_2 & \dots & f_d \\ f_2 & f_1 & \dots & f_d \\ \vdots & \vdots & \ddots & \vdots \\ f_d & f_d & \dots & f_1 \end{bmatrix}$

→ Square sym matrix
→ Variance of feature
 $d \times d$

If X is not column star the

$$S_{ij} = \frac{1}{n} \sum_{k=1}^n (f_{ik} - \text{mean}(f_i)) (f_{jk} - \text{mean}(f_j))$$

If X is column standardize

$$S = \frac{1}{n} (X^T)(X)$$

* MNIST dataset

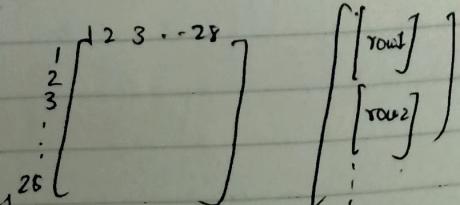
- * Dataset of hand written digits.
- * 28×28 pixel images.
- * 60K training data points
- * 10K test data points

$$D = \{(x_i, y_i)\}_{i=1}^{60K} \quad x_i = \begin{bmatrix} 1 & 2 & \dots & 28 \\ \vdots & \vdots & \ddots & \vdots \\ 28 & 28 & \dots & 1 \end{bmatrix} \quad y_i = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

Task: given a new 28×28 pixel image classify that in one of the (0-9) character.

We learned the x_i should be vector, but here x_i is 28×28 pixel. So how do we convert x_i to a 1D col vector.

We do flattening



Ex:

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{bmatrix}$$

We can do either row-wise or column wise.

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \xrightarrow{\text{row wise}} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{bmatrix} \xrightarrow{\text{column wise}} \begin{bmatrix} 1 \\ 4 \\ 7 \\ 2 \\ 5 \\ 8 \\ 3 \end{bmatrix}$$

In our MNIST data now our $\begin{bmatrix} \quad \end{bmatrix}_{28 \times 28}$
row flatten we get $\begin{bmatrix} \quad \end{bmatrix}_{784 \times 1}$ size column vector.

Now over all data set looks like.

$$X = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \\ 60K \end{bmatrix} \xrightarrow{\text{dim}} 784 \xrightarrow{\text{60K}} 60K \times 784$$
$$Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ 60K \end{bmatrix} \xrightarrow{\text{60Kx1}}$$

A PCA (Principal Component analysis (PCA))

PCA is basically used for many other thing but one of its application is Dimensionality Reduction.

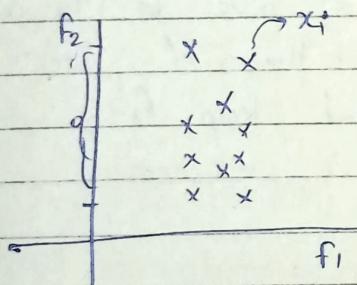
Data given d -dim \rightarrow d' -dim
s.t. $d' < d$

for Ex:- MNIST \rightarrow 2D
 $784 \rightarrow 2D$

We will also see other applications of dimensionality reduction, for now we see only visualisation point of you.

→ Geometric interpretation of PCA:-

(d) for example $2D \rightarrow 1D$



Let f_1 = thickness of hair $\in \mathbb{R}$

f_2 = height of people

Indian data

One thing we can note is height have larger spread. while f_1 have small variance.

Now, given this 2D dataset, we want to convert it to 1D. As spread of f_1 is less, we can skip the f_1 and make another 1D dataset of only f_2 .

$$x = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix} \rightarrow x' = \begin{bmatrix} f_2 \\ \vdots \\ f_n \end{bmatrix}$$

Now suppose another example:

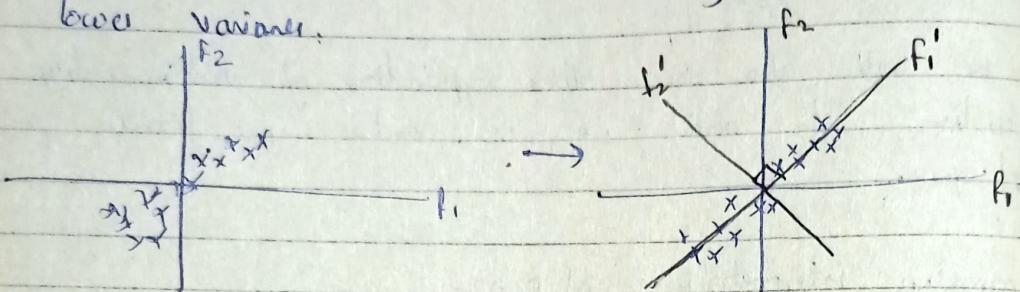
$X = 2\text{-dim column standardize dataset}$

$$X_2 = \begin{bmatrix} f_1 & f_2 \end{bmatrix}$$

$$\text{mean}\{f_1\} = \text{mean}\{f_2\} = 0$$

$$\sigma\{f_1\} = \sigma\{f_2\} = 0$$

So, in this case we can't simply drop ~~one~~ one with lower variance.



Let's draw another set of axes like above.
Now $f_2' \ll$ variance than f_1' so we can simply drop f_2' and keep f_1' .

step 1: Find $f_1' \perp f_2'$

② drop f_2'

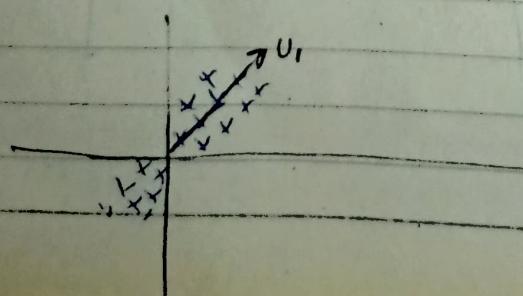
③ project x 's onto f_1' then $2D \rightarrow 1D$

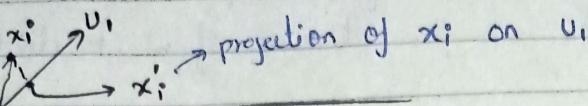
→ we want to find a direction f_1' s.t. variance of x 's projected onto ~~f1~~ f_1' is maximal.

Now, let's try to write it in mathematical form.

Let $f_1' = \underline{u}_1$,

second thing is we really don't need the whole f_1' line we just need direction or we can say a unit vector.





$$x_i^{\circ} = \text{proj}_{u_i} x_i = \frac{u_i \cdot x_i}{\|u_i\|^2}$$

unit vector
 $\therefore \|u_i\| = 1$

$$= u_i^T x_i$$

so each of $x_i^{\circ} = u_i^T x_i$

$$\bar{x}' = u_i^T \bar{x}$$

$\underbrace{\bar{x}'}_{\text{mean of } x_i^{\circ}}$ $\underbrace{\bar{x}}_{\text{mean of } x_i}$

so, find u_i s.t. the variance of points x_i projected on u_i is maximal.

$$\text{Var}\{\text{proj}_{u_i} x_i\}_{i=1}^n \Rightarrow \text{Var}\{u_i^T x_i\}_{i=1}^n$$

\Downarrow

$$\frac{1}{n} \sum_{i=1}^n (u_i^T x_i - u_i^T \bar{x}_i)^2$$

≥ 0 {as X is rd stn
less mean = 0}

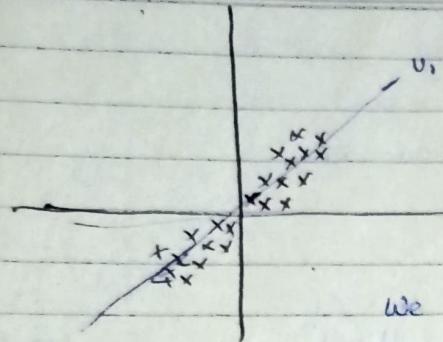
$$\text{we want to max} \rightarrow \frac{1}{n} \sum_{i=1}^n (u_i^T x_i)^2 \quad \text{s.t. } u_i \text{ is unit vector}$$

so, we want to find u_i s.t. $\frac{1}{n} \sum_{i=1}^n (u_i^T x_i)^2$ is max and u_i is unit vector.
 This is an optimization problem.

\Rightarrow Alternative formulation of PCA : dist-minimization.

We have already seen a formulation of PCA i.e. finding a unit vector u_i s.t. $\frac{1}{n} \sum_{i=1}^n (u_i^T x_i)^2$ is max.

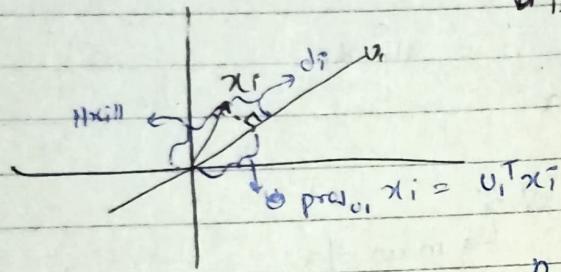
There is another formulation.



for each point x_i calculate d_i i.e. distance of x_i to direction v_1 .

We want to find v_1 which minimize the squared distance of each point x_i to v_1 .

$$\min v_1 \sum_{i=1}^n d_i^2$$



$$\begin{aligned} d_i^2 &= \|x_i\|^2 - (v_1^T x_i)^2 \\ &= (x_i^T x_i) - (v_1^T x_i)^2 \end{aligned}$$

∴ We want $\min v_1 \sum_{i=1}^n (x_i^T x_i - (v_1^T x_i)^2)$

So we get two type of formulation

$$\max_{v_1} \frac{1}{n} \sum_{i=1}^n (v_1^T x_i)^2$$

Variance maximization
PCA

$$\min_{v_1} \sum_{i=1}^n (x_i^T x_i - (v_1^T x_i)^2)$$

dist min PCA

i. Solution to our optimization problem:

Here we only see solution, we don't see how we reach

to $S^{1/n}$. f_1, f_2, \dots, f_d

$$X = \begin{bmatrix} 1 & 2 & \dots & n \end{bmatrix}_{n \times d}$$

Column standardize
($\mu=0, \sigma=1$)

$$S = \frac{1}{n} X^T X \quad \left\{ \because X \text{ is column standardized} \right.$$

Covariance matrix

If we compute the eigen values of $S = \lambda_1, \lambda_2, \lambda_3, \lambda_4, \dots, \lambda_d$
Eigen vectors of S

$$S = \begin{bmatrix} \downarrow \lambda_1 & & & \\ \downarrow \lambda_2 & \ddots & & \\ \downarrow \lambda_3 & & \ddots & \\ \vdots & & & \ddots & \lambda_d \end{bmatrix}$$

Let $\lambda_1 > \lambda_2 > \lambda_3 > \lambda_4 > \dots > \lambda_d$
Eigen value.

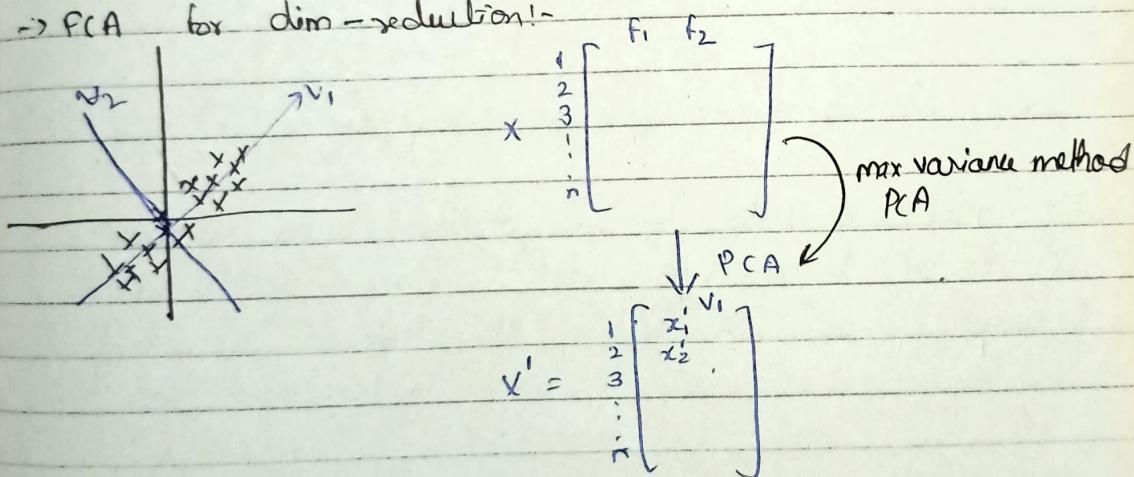
$u_1 = v_1$ = eigen vector of S corresponding to largest eigen value λ_1

so, find steps are.

- I Column-standardize X
- II Compute covariance matrix S
- III Compute eigen value & eigen vector
- IV $u_1 = v_1$ eigen vector of corresponding to largest eigen value.

so, eigen vector is the vector in direction where we get maximal variance. But what does eigen value signifies.
It tells us the percentage of variance on that eigen vector.

→ PCA for dim-reduction:-



Let's take a bigger example:

$$X = \begin{bmatrix} f_1 & f_2 & \dots & f_{10} \\ x_1 & x_2 & \dots & x_n \end{bmatrix}_{n \times 10} \rightarrow S = X^T X \rightarrow \lambda_1, \lambda_2, \lambda_3, \dots, \lambda_{10} \rightarrow X' = \begin{bmatrix} x'_1 & x'_2 & \dots & x'_{10} \end{bmatrix}_{n \times 1}$$

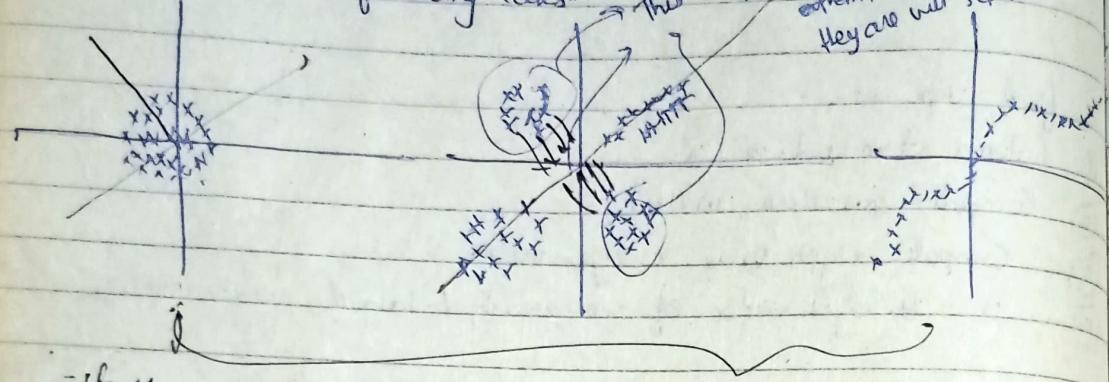
$$x'_i = [x_{i1}, x_{i2}, \dots, x_{i10}] \rightarrow x'_i = [x_i^T v_1, x_i^T v_2]$$

So, someone ask - reduce dimension of data from d to d' - dim. s.t. we preserve 99% of data.

In that case as we know λ tells us the % of variance present by corresponding eigen vector. We will check upto which λ we get 99%. Include that much of eigen vectors.

→ Limitations of PCA:-

Suppose the following cases.



- If we apply PCA to change $2D \rightarrow 1D$ we can see that each we can lost various information.

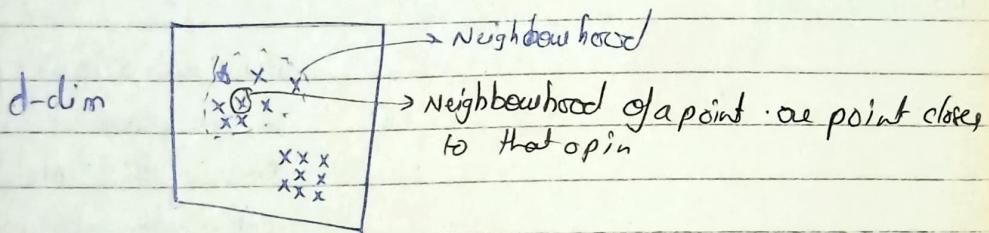
A t-SNE (t-distributed, stochastic Neighbourhood Embedding)

This is one of the best ~~di~~, state of art dimensionality reduction technique specially for visualization.

There are also other techniques of Dimensionality Reduction like, multi-dime mapping, common mapping etc.

→ Neighbourhood, Embedding:

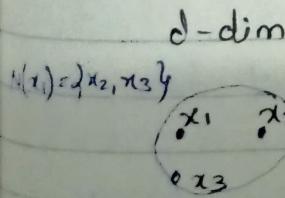
Imagine we have a d-dim space.



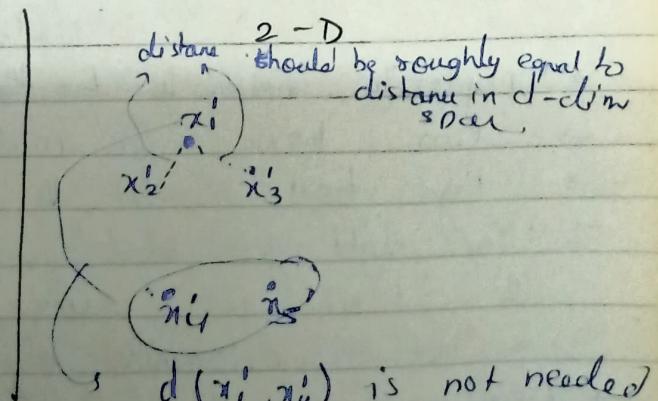
$$N(x_i) = \{x_j \mid x_i \text{ and } x_j \text{ are geometrically close}\}$$

Suppose we want to change $d\text{-D} \rightarrow 2\text{-D}$
Embedding means take each point in $d\text{-D}$
find its corresponding value in 2-D space
or lower-dim space.

→ Geometric Intuition



$$N(x_4) = \{x_5\}$$



$d(x'_1, x'_4)$ is not needed to be same as original distance of nodes.

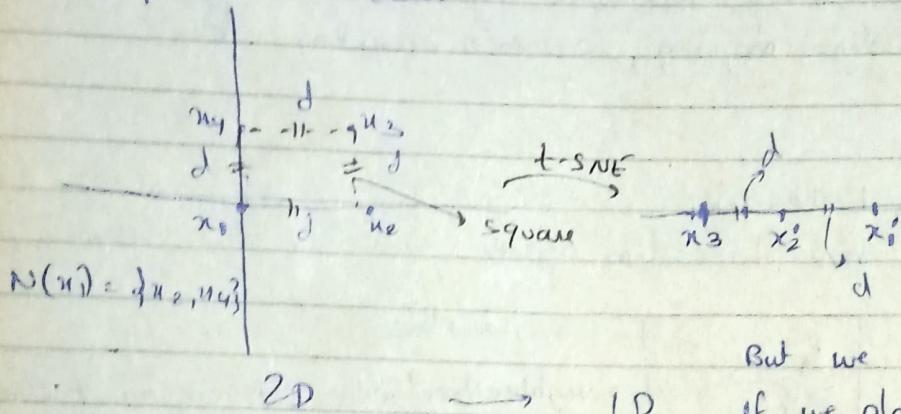
t-SNE only try to preserve same neighbourhoods.

If two nodes are not in neighbor then t-SNE don't care about their distance.

* crowding problem.

We said: t-SNE: preserv dist in a Neighbor

Suppose a case:



But we cannot place x_2
if we place it anywhere.
one of the Neighbors dist
will not preserved.

Some time it is impossible to preserve distance in all
the neighbourhood. Such a problem is called crowding
problem.

* t-SNE:-

2 important parameters

(1) steps (No. of iteration)

(2) perplexity / Neighborhood size (Not exactly but similar idea).

- perplexity should be less than no of data points.
- Run t-SNE for multiple values of perplexity.
- Run t-SNE until the shape does not change.
- You might get different plot each time you run the t-SNE with same parameter. Hence t-SNE is not a deterministic algorithm.
- t-SNE expands the denser Neighbors and shrinks less dense Neighbors. Hence, cluster sizes change after t-SNE.