

## PCA (Principal Component Analysis)

used for dimensionality reduction.

②  $d$  dimens  $\xrightarrow{\text{convert}}$   $d'$  dimens  
where  $d' < d$ .

### Applications of PCA

① MNIST  $\rightarrow$  784 dim  $\rightarrow$  2 dim (visualise)

\* ~~Geometric~~ Geometric Interpretation of PCA  $\rightarrow$

2 dim - data  $\rightarrow$  1 - dim data

$$X = \begin{matrix} & f_1 & f_2 \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ n \end{matrix} & \begin{bmatrix} & \\ & \\ & \\ & \\ & \end{bmatrix} \end{matrix} \rightarrow X' = \begin{matrix} & f_2 \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ \vdots \\ n \end{matrix} & \begin{bmatrix} & \\ & \\ & \\ & \\ & \end{bmatrix} \end{matrix}$$

preserving the direction with maximal spread / variance / information in it which is  $f_2$  here.



Dropping the features which have low spread & keeping the features which have high spread.

⇒ We want to find a direction  $f_1'$  such that the variance of  $x_i$ 's projected into  $f_1'$  is maximal.

rotating my axis to find  $f_1'$  with max variance and drop  $f_2'$ .

⇒ We have to find a unit vector  $u_1$  such that variance  $\left\{ \text{proj}_{u_1} x_i \right\}_{i=1}^n$  is maximal.

$$\Rightarrow \text{var} \{x_i'\}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u_1^T x_i)^2$$

objective of an optimisation problem

$$\max_{u_1} \frac{1}{n} \sum_{i=1}^n (u_1^T (x_i)) ^2$$

→ var  $\{x_i'\}$   
→ data matrix  
↓  
optimisation problem.

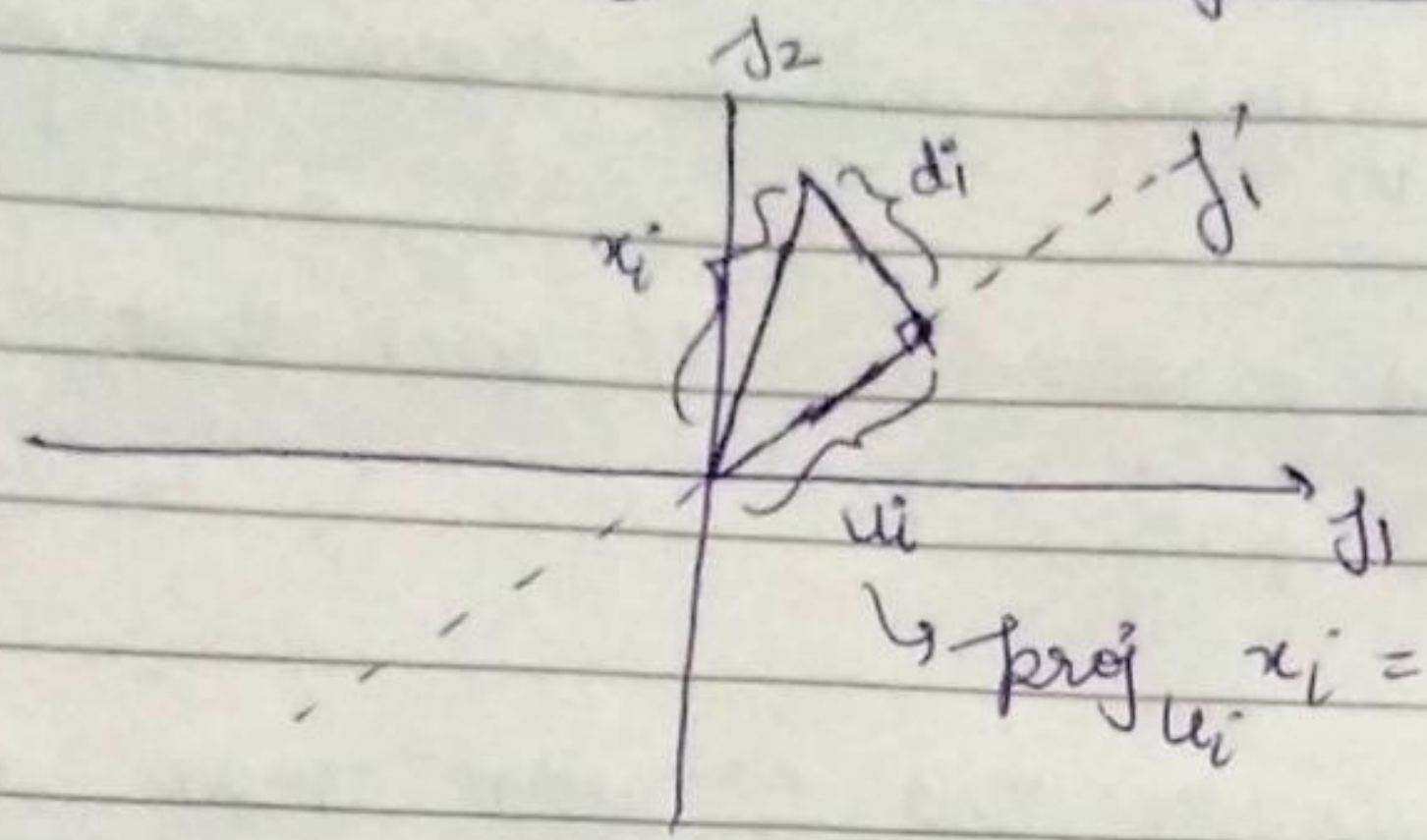
so that,  $u_1^T u_1 = 1 = |u|^2$   
↳ constraint      ↳  $u_1$  is a unit vector

why  $u_1$  should be a unit vector?

⇒ Otherwise if we will take  $u_1$  as  $[0, \infty]$  then it will be always maximum and we want to make it maximum when it is.



⇒ Alternative formulation of PCA ⇒



$u_i = \text{unit vector}$   
 $u_i^T u_i = 1 = |u_i|^2$

$$d_i^2 = |x_i|^2 - (u_i^T x_i)^2$$

dist min PCA

$$\min_{u_i} \sum_{i=1}^n \underbrace{\left( (x_i^T x_i) - (u_i^T x_i)^2 \right)}_{d_i^2}$$

such that  $u^T u = 1$

\* Solution to our optimisation problems ⇒

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

Steps: ① Col standardisation of  $X$  is done.

②  $S_{d \times d} = X^T X$

③ eigen values & vectors of  $S$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

$$v_1, v_2, v_3, \dots, v_d$$

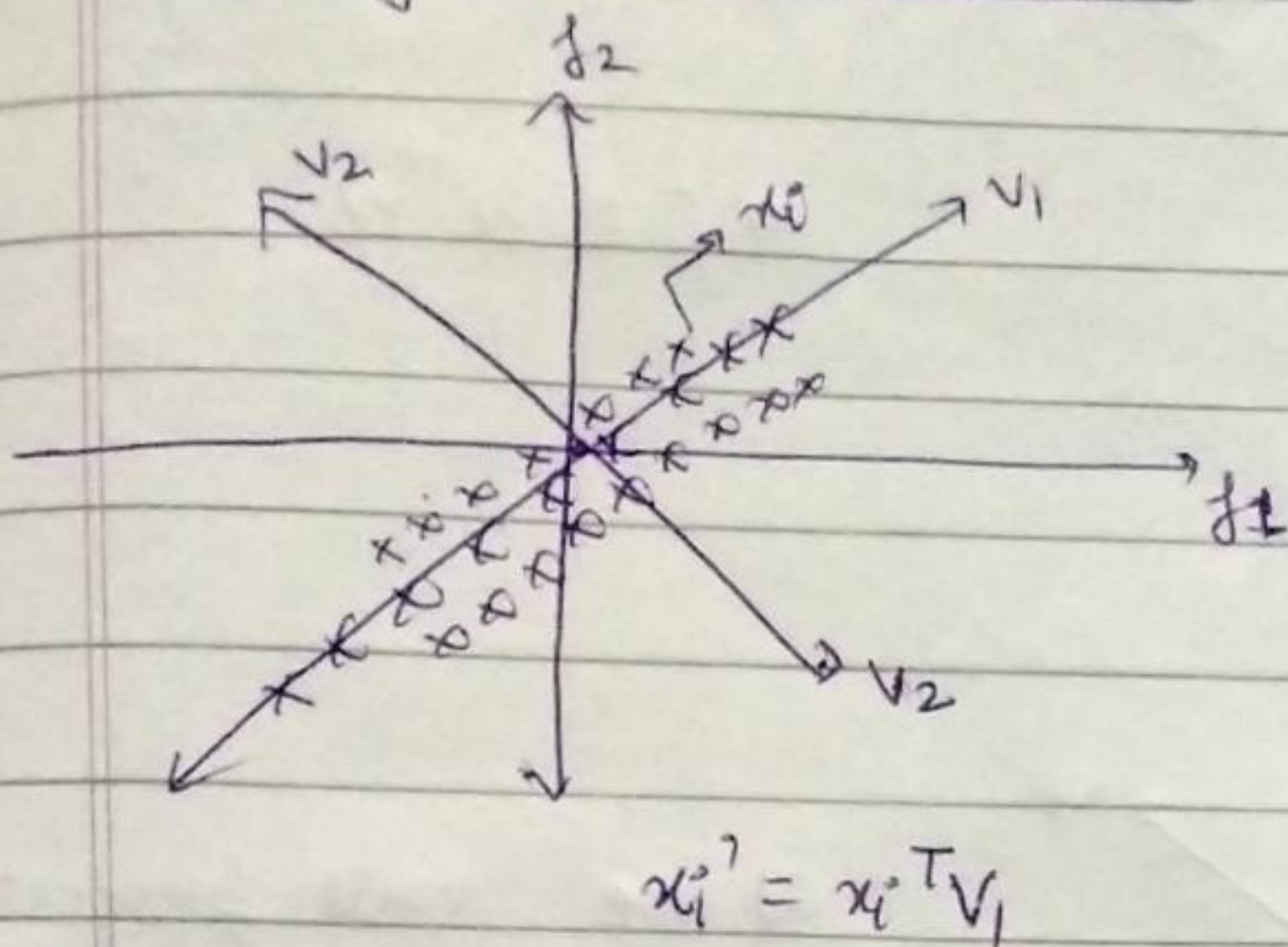
④  $u_1 = v_1$

$v_1, v_2 \rightarrow$  tells about direction of max variance  
 $\lambda_1, \lambda_2 \rightarrow$  tells about the data spread that is  
 it at only one axis or more axis



$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = 1 = \% \text{ age of variance explained.}$$

\* PCA for dim reduction  $\Rightarrow$



$$X = \begin{bmatrix} 1 & j_1 & j_2 \\ 2 & & \\ 3 & & \\ \vdots & & \\ n & & \end{bmatrix}$$

$\leftarrow x_i^T \rightarrow$

$$S = X^T X$$

max-  
variance  
method (PCA)

$v_1 \rightarrow$  as  $v_1$  has  
max.  
variance.

$$X^o1 = \begin{bmatrix} 1 \\ 2 \\ 3 \\ \vdots \\ n \end{bmatrix} \begin{bmatrix} x_i^{o1} \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & j_1 & j_2 & \dots & j_{10} \\ 2 & & & & \\ 3 & & & & \\ \vdots & & & & \\ n & & & & \end{bmatrix} \quad n \times 10$$

$\leftarrow x_i^T \rightarrow$

dim reduction  
(PCA)

\* visualize

$$X^o1 = \begin{bmatrix} 1 & v_1 & v_2 \\ 2 & & \\ 3 & & \\ 4 & & \\ 5 & & \\ \vdots & & \\ n & & \end{bmatrix} \quad n \times 2$$

$\leftarrow x_i^T \rightarrow$

$$S = X^T X$$

$$\text{eigen}(S) = \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_{10}$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \dots \quad \downarrow$$

$$\textcircled{v_1} \quad \textcircled{v_2} \quad v_3 \quad \dots \quad v_{10}$$

take top two  
eigen vectors.

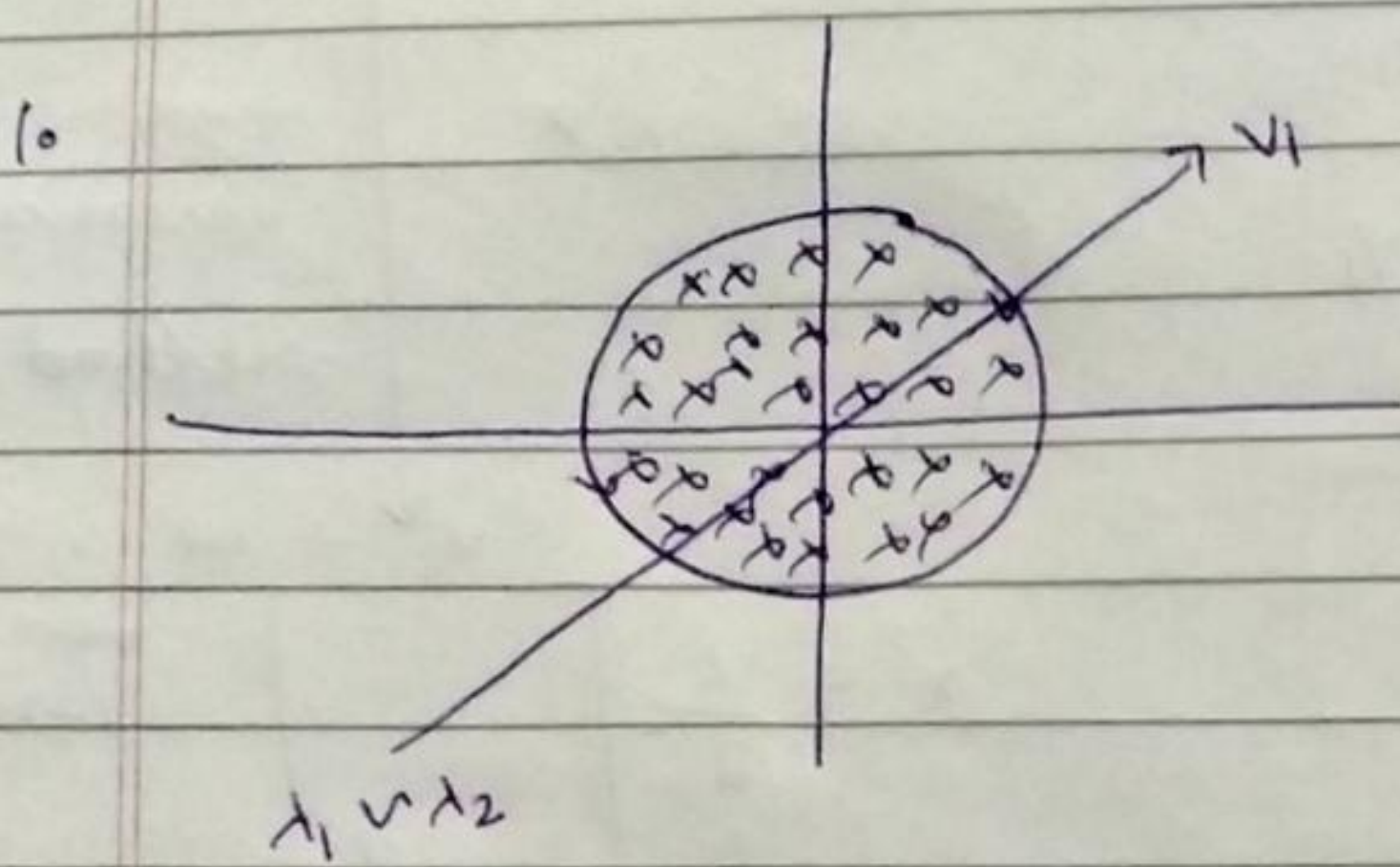
$$x_i^{o1} = [x_i^T v_1, x_i^T v_2]$$

$\textcircled{v_1} \quad \textcircled{v_2}$

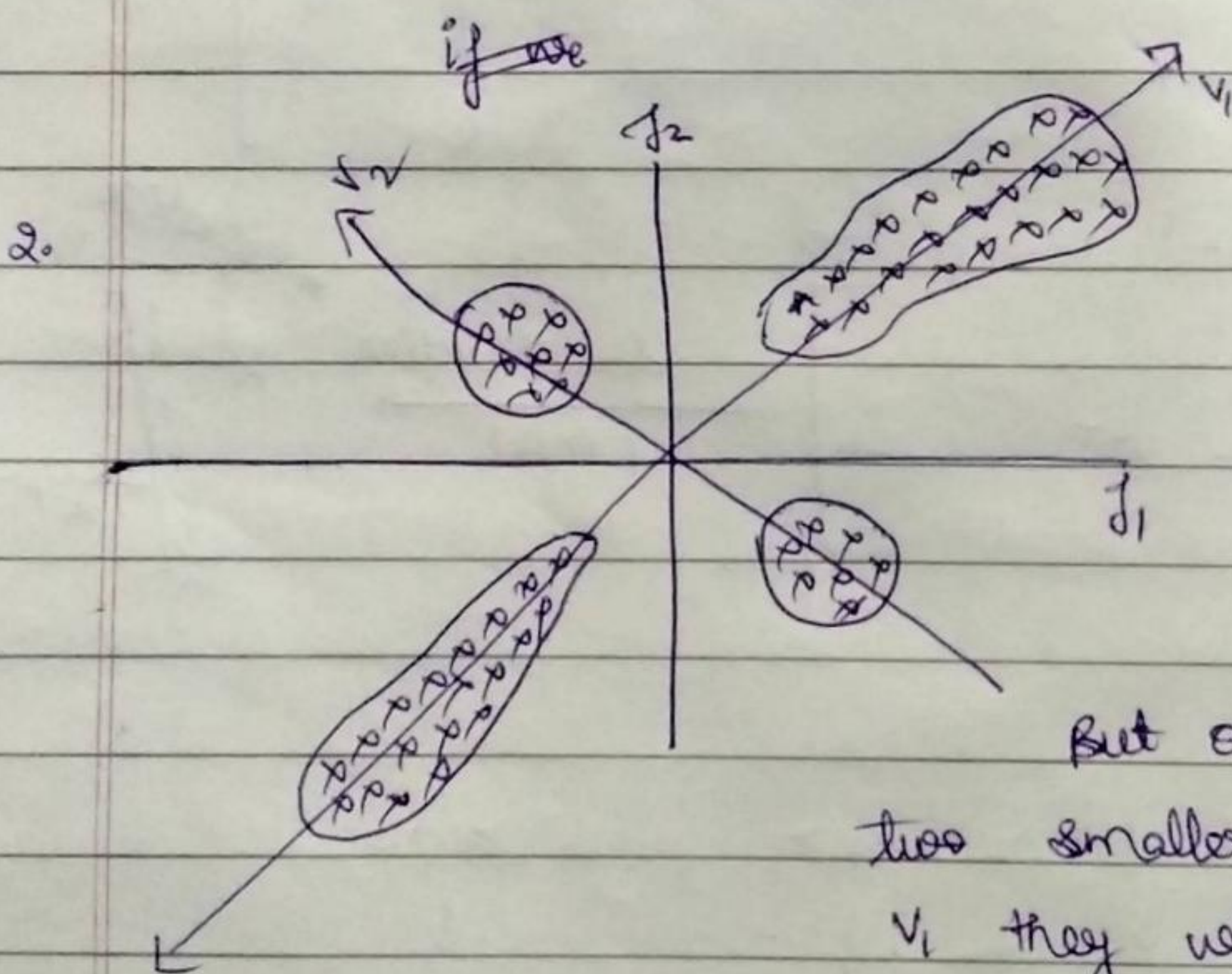


$X = \begin{bmatrix} 1 & d_1 & d_2 & \dots & d_{100} \\ 2 & & & & \\ 3 & & & & \\ \vdots & & & & \\ n & & & & \end{bmatrix}$ 
 $\xrightarrow{\text{PCA}}$ 
 $X' = \begin{bmatrix} x_1' & x_2' & \dots & x_{50}' \\ \vdots & \vdots & & \vdots \\ x_n' & x_n' & \dots & x_n' \end{bmatrix}$ 
 $x_{ij}' = x_i^T v_j$

### \* Limitations of PCA



if we will project all data points on  $v_1$  we will lose almost 50% of useful info using PCA



for projecting two bigger clusters on  $v_1$ , they are well separated

But on projecting two smaller clusters on  $v_1$  they will be cluttered together and info will not be clear and distinguishable.



⇒ Standard Scalar → used to centre and scale each feature independently.

⇒ to make mean = 0

and to scale variance to 1.

Basically to perform column standardisation.

\* Other advantages of PCA is to

convert  $\frac{784}{d} \rightarrow 10 \rightarrow \text{ML models.}$   
 $d \leq d'$

for ex: →

784

→

200 dim

$\times 15000 \times 784$

$\times$

$V_{784 \times 200}$

→

$S_{15K \times 200}$

Top 200 eigen vectors.

Q → How to know how much % of total variance/info will be explained by a lower dimension which will get using PCA?

⇒

for ex: →

784

→

10 dim

% of variance explained in 10 dim =  $\frac{\lambda_1 + \lambda_2 + \dots + \lambda_{10}}{\sum_{i=1}^{784} \lambda_i}$

let's assume it comes to be 0.2

So, 20% of the total variance in 784 dim is explained in 10 dimensions.