

# Dimensionality Reduction

2D, 3D  $\Rightarrow$  Scatter plot

4D, 5D, 6D  $\Rightarrow$  Pair plot

10D, 100D, 1000D  $\rightarrow$  dimension reduction techniques

PCA

t-SNE

\* Dataset  $\Rightarrow$

$\nearrow$  data points  
 $\nearrow$  class labels

iris  $\Rightarrow$   $D = \{x_i^a, y_i^a\}^n \rightarrow$  data points

$x_i^a \in \mathbb{R}^d$  ;  $x_i^a \in \mathbb{R}^4$

$y_i^a \in \{\text{setosa, versicolor, virginica}\}$

\* Dataset as a data-matrix  $\Rightarrow$

$$\tilde{X} = \begin{matrix} & f_1 & f_2 & f_3 & f_j & f_d \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ n \end{matrix} & \left[ \begin{array}{cccccc} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{array} \right] \end{matrix}$$

$x^T$  (row vector)  $\rightarrow$   $n \times d$  matrix

each datapoint = row

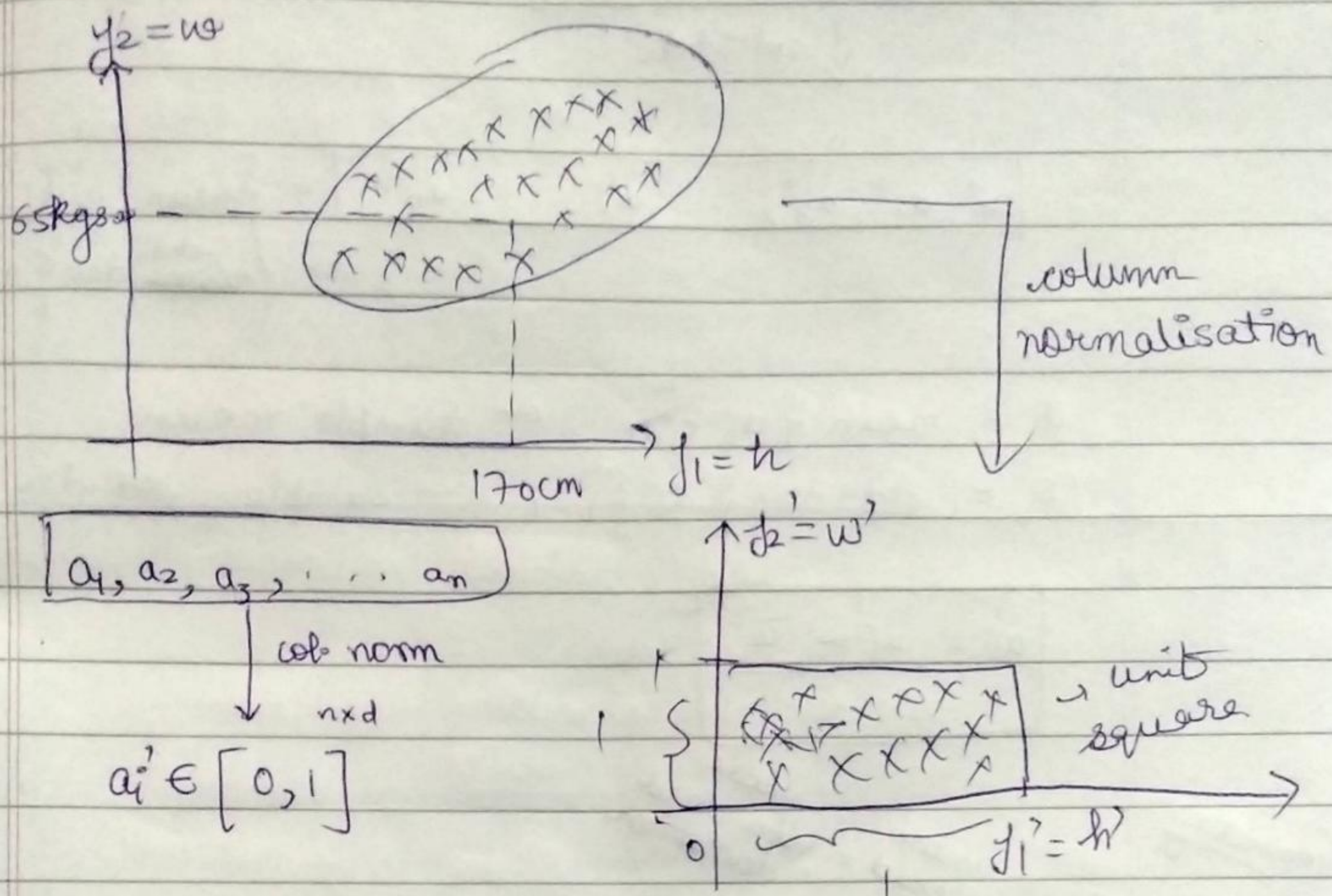
each column = feature.

Or we can also represent it into  $d \times n$  representation.



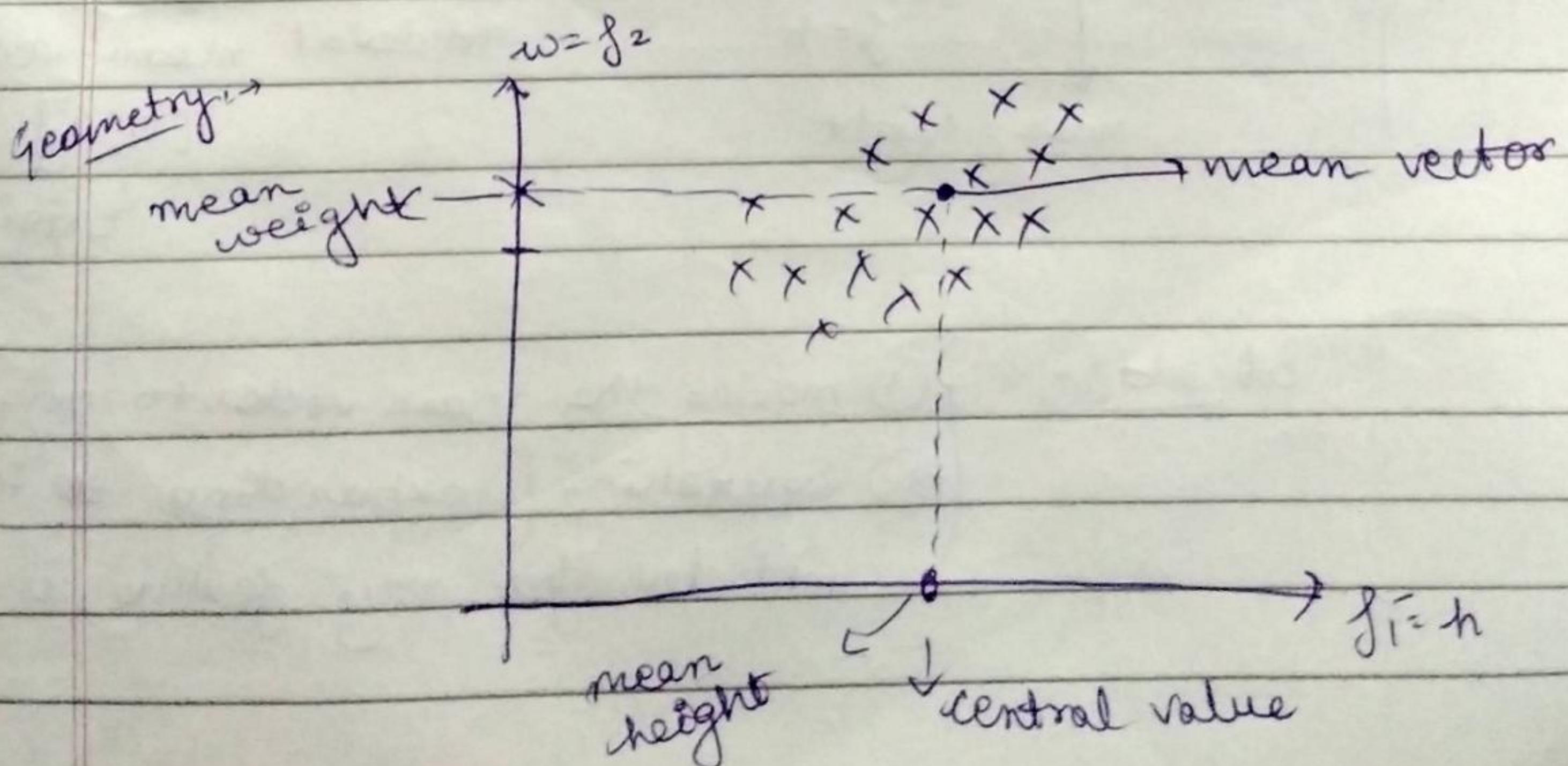
## \* Data Pre-processing : column normalisation

Squashing the data into unit cube or unit cuboid.



## \* Mean Vector :-

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$





# \* Data preprocessing - Column standardization

more often used

$$[a_1, a_2, a_3, \dots, a_n] \rightarrow n \text{ values of } f_j$$

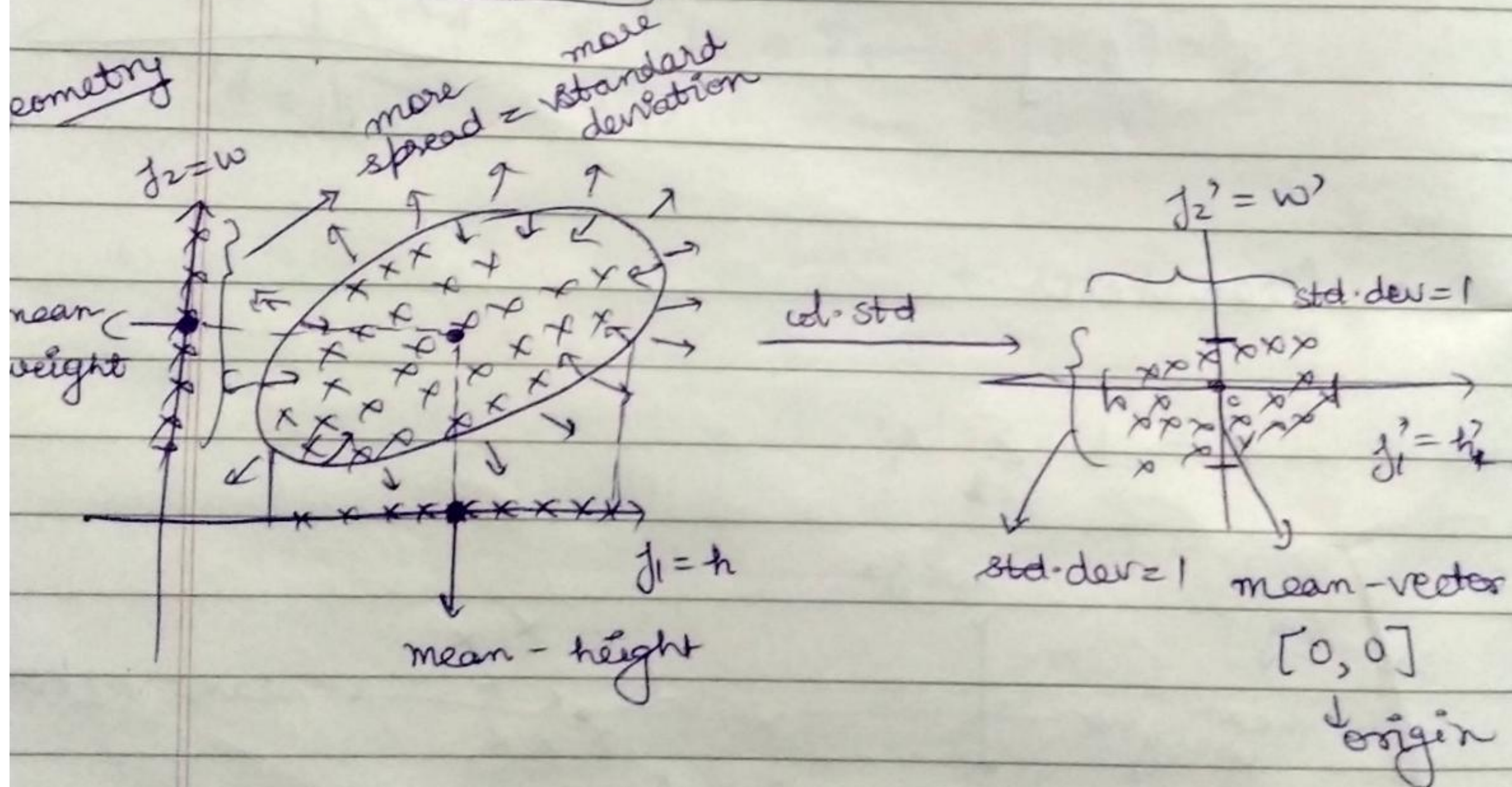
↓ col. std

$$[a_1', a_2', a_3', \dots, a_i', \dots, a_n'] \leftarrow \begin{cases} \text{mean} \{a_i'\}_{i=1}^n = 0 \\ \text{std. dev} \{a_i'\}_{i=1}^n = 1 \end{cases}$$

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

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

$$a_i' = \frac{a_i - \bar{a}}{s}$$



col. std  $\rightarrow$

- ① moving the mean vector to origin
- ② Squashing / expanding so that std. dev for any feature is 1.



col-standardization  $\rightarrow$  mean-centering  $\rightarrow$  origin  
 $\downarrow$   
 $+ \text{scaling}$

ensuring that std-dev for  
all the features is 1.

\* Co-Variance Matrix  $\Rightarrow$

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

$$S_{ij} \rightarrow \text{cov matrix} = \text{cov}(f_i, f_j)$$

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

$\rightarrow$  It is a symmetric matrix as  $S_{ij} = S_{ji}$   
and diagonal elements are variance of features.

$\rightarrow$  It is a square matrix as it has  $d$  rows and  $d$  columns.

$$\text{cov}(f_1, f_2) = \frac{1}{n} \sum_{i=1}^n x_{i1} * x_{i2}$$

$$\text{cov}(f_1, f_2) = \underbrace{(f_1^T f_2)}_{f_1 \cdot f_2 \rightarrow \text{dot product}} * \frac{1}{n}$$

$f_1 \cdot f_2 \rightarrow \text{dot product}$

So, if  $f_1$  &  $f_2$  have been column standardised,  

$$\text{cov}(f_1, f_2) = \frac{f_1^T f_2}{n}$$



Date : / /

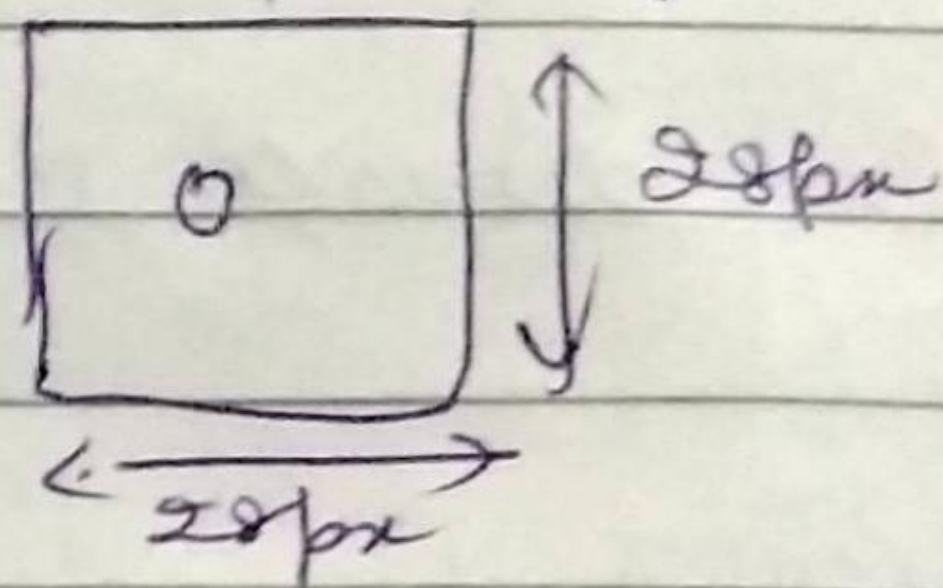
Page No.

$$S_{d \times d} = \begin{matrix} & X^T & \star & X \\ \begin{matrix} d \times n \\ \end{matrix} & & & \begin{matrix} n \times d \\ \end{matrix} \end{matrix} \rightarrow \text{if } x \text{ has been col-stored}$$

\* MNIST dataset  $\rightarrow$  Dataset for visualising digits

$$D = \{x_i, y_i\}_1^{260K}$$

$x_i :$



$$y_i = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$