

# Dimensionality Reduction

2-D, 3-D :- scatter plot }  
 4-D, 5-D, 6-D :- pair plot }

10-D, 100-D, 1000-D → what about them how to visualize it?

n-D → 2-D or 3-D

PCA & t-SNE

## Row-vector & column-Vector

flower :- [SL, PL, SW, PW]

→ i<sup>th</sup> point:

$x_i \in \mathbb{R}^d \rightarrow$  d-dimensional column vector, containing Real values  
 $x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ \vdots \\ x_{id} \end{bmatrix}_{d \times 1}$  → Real space

for eg:-  $f_1 = \begin{bmatrix} 2.1 \\ 3.2 \\ 1.6 \\ 1.4 \end{bmatrix} \rightarrow \text{P.L}$   
 $\rightarrow \text{P.W}$   
 $\rightarrow \text{S.L}$   
 $\rightarrow \text{S.W}$

Note:- The default vector is always column vector.

$x_i := [2.1, 3.2, 1.6, 1.4]_{1 \times 4}$  : row-vector.

## How to represent dataset

$D = \{x_i, y_i\}_{i=1}^n \rightarrow$  # data points

data-points  $\leftarrow x_i \in \mathbb{R}^d ; x_i \in \mathbb{R}^4$

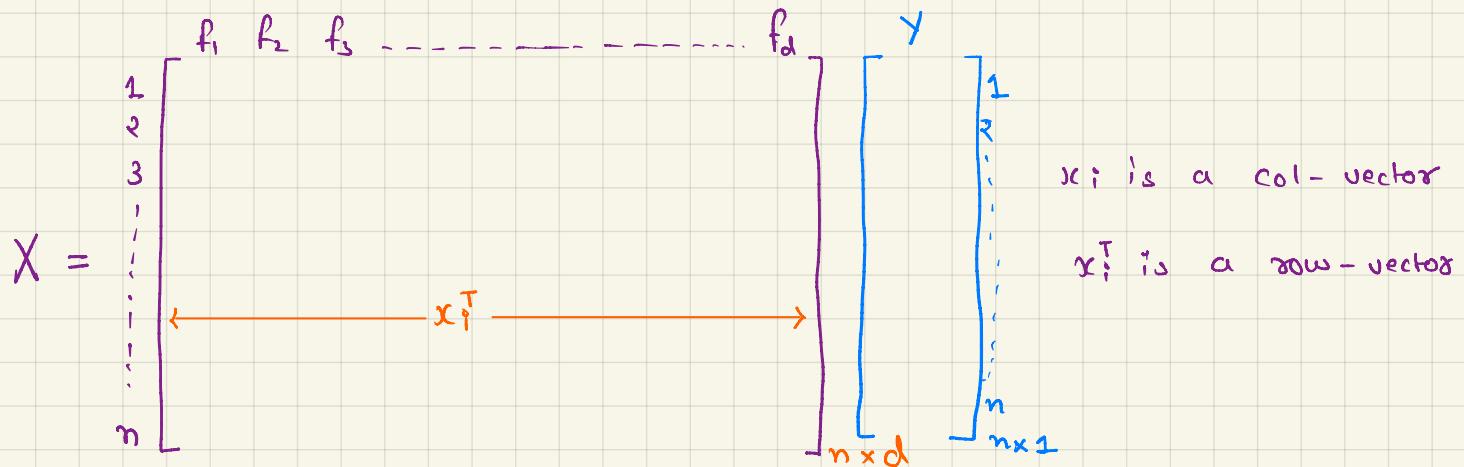
$$x_i = \begin{bmatrix} SL \\ SW \\ PL \\ PW \end{bmatrix}$$

class-labels  $\leftarrow y_i \in \{\text{setosa}, \text{versicolor}, \text{virginica}\}$

## Dataset as a data-matrix

$$D = \{x_i, y_i\}_{i=1}^n, \quad x_i \in \mathbb{R}^d, \quad y_i \in \{S, V_i, V_e\}$$

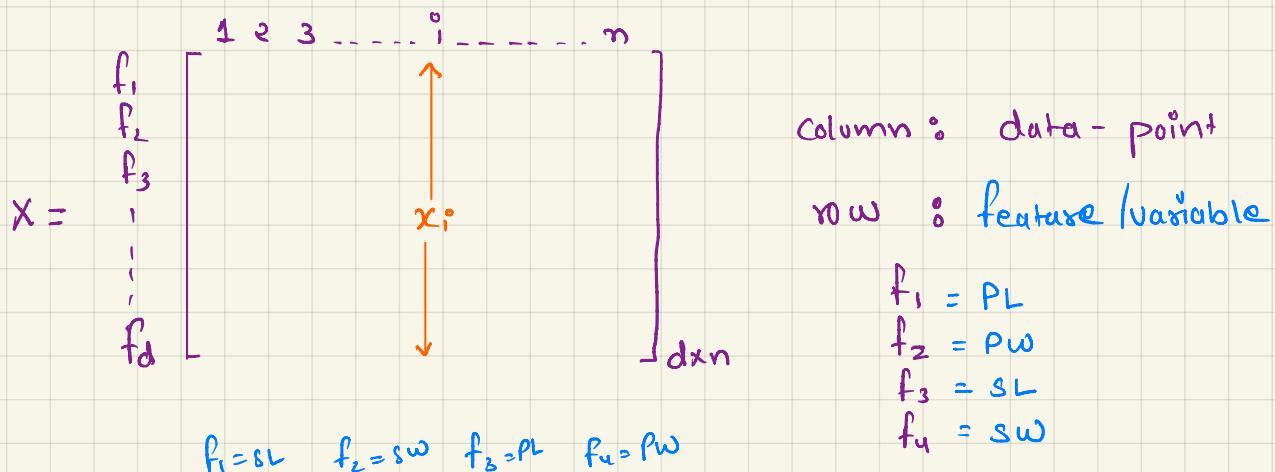
d-features



Each datapoint : row

each column : feature

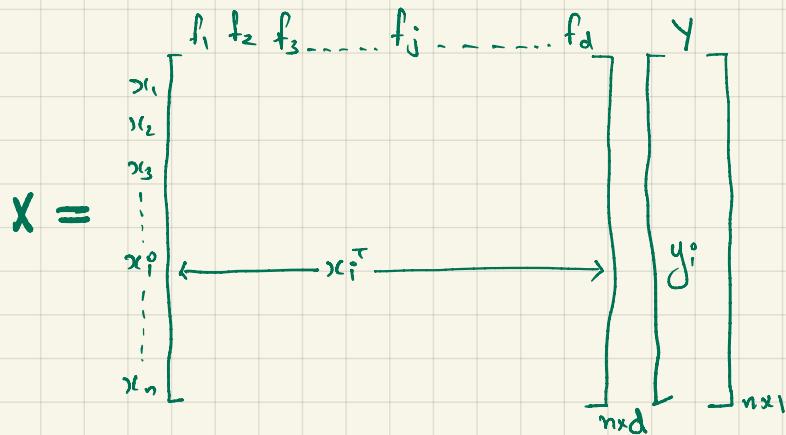
## Another way



Eg:-

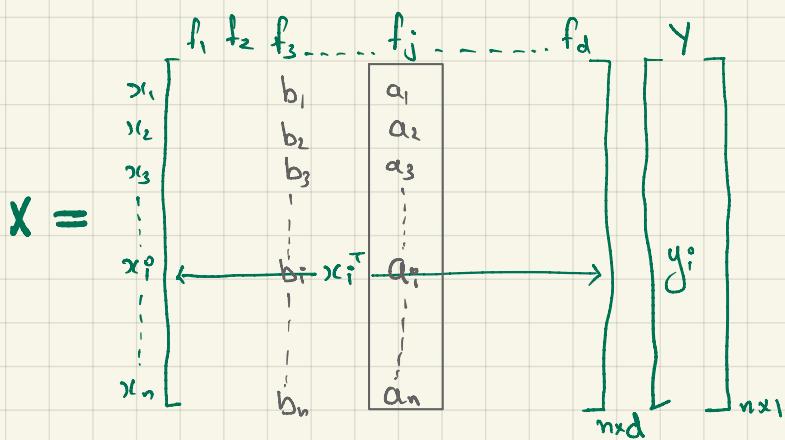
$r_{11}$			
$r_{12}$			

## Data preprocessing: feature/column Normalization



Obtain data  $\rightarrow$  pre-processing  $\rightarrow$  data modeling (dim-reduction)

$\downarrow$  column normalization  
 $\vdots$



$a_1, a_2, \dots, a_n = n\text{-value of } f_j$

$$\max(a_i) = a_{\max} \geq a_i \quad (i: 1 \rightarrow n)$$

$$\min(a_i) = a_{\min} \leq a_i \quad (i: 1 \rightarrow n)$$

$\rightarrow a'_1, a'_2, a'_3, \dots, a'_i, \dots, a'_n$

$$a'_i = \frac{a_i - a_{\min}}{a_{\max} - a_{\min}} \quad a'_i \in [0, 1]$$

$$\text{e.g.: } a'_{\min} = \frac{a_{\min} - a_{\min}}{a_{\max} - a_{\min}} = 0 \quad ; \quad a'_{\max} = \frac{a_{\max} - a_{\min}}{a_{\max} - a_{\min}} = 1$$

$a_1, a_2, a_3, \dots, a_i, \dots, a_d; \quad a_i \in \mathbb{R}$

$\downarrow$  column normalization

$a'_1, a'_2, a'_3, \dots, a'_i, \dots, a'_d; \quad s.t. \quad a'_i \in [0, 1]$

why

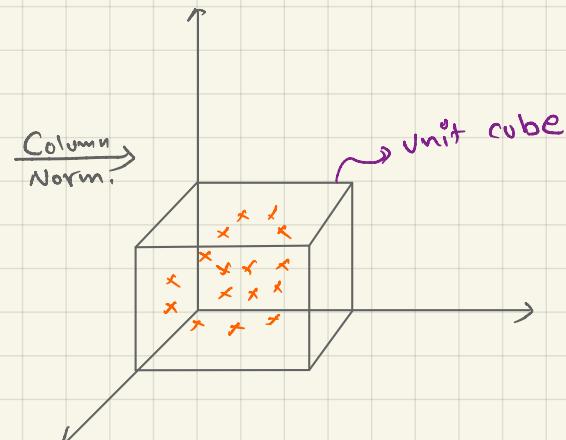
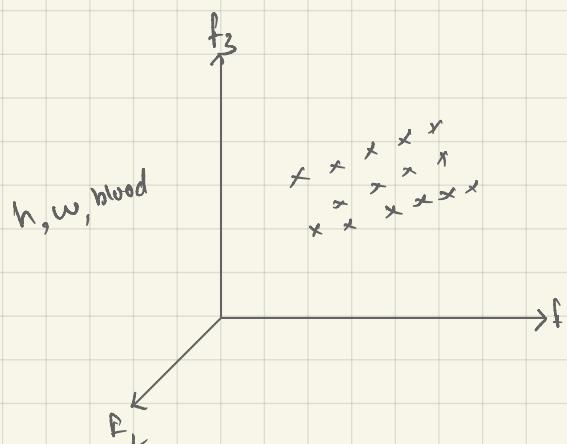
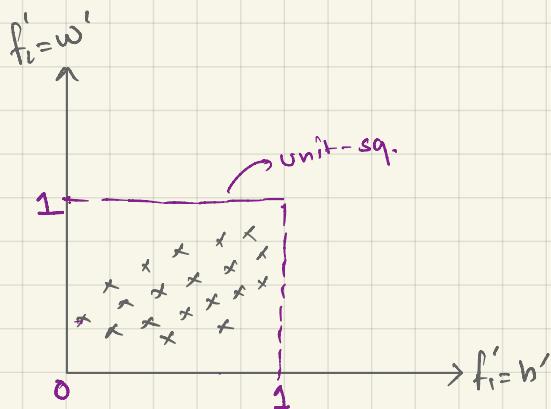
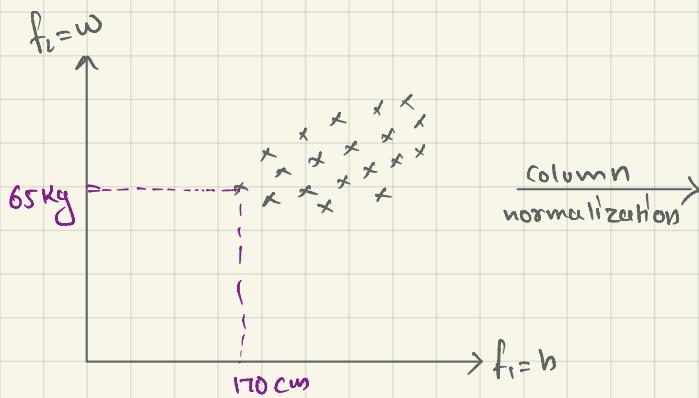
$h$	$w$	$f_1' = h'$	$f_2' = w'$
Student $\rightarrow 162$	56	2	1
$\therefore 2 \rightarrow 172$	72	2	2

Col-normalization

$[0,1]$        $[0,1]$

cm      kg

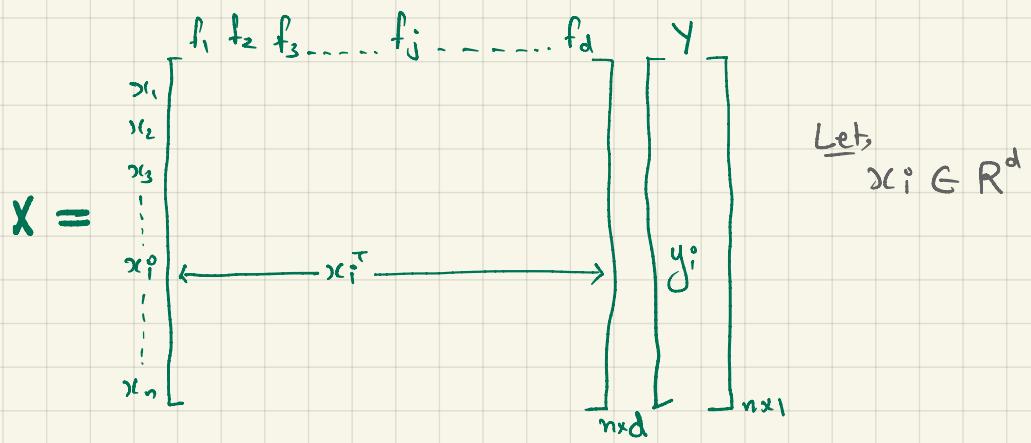
you are getting rid of scale ( $\text{cm} \rightarrow \text{m}$  or  $\text{kg} \rightarrow \text{ton}$ ) by doing normalization. we are putting everything in b/w of 1.



$\Rightarrow$  anywhere in n-dim. space  $\xrightarrow{\text{Col. norm.}}$  Unit - hyper cube in the same n-dim. space.

by doing that we are getting rid of scale.

# Mean - Vector



Suppose

$$\vec{x}_1 = [2.2, 4.2] \in \mathbb{R}^2$$

$$\vec{x}_2 = [1.2, 3.2] \in \mathbb{R}^2$$

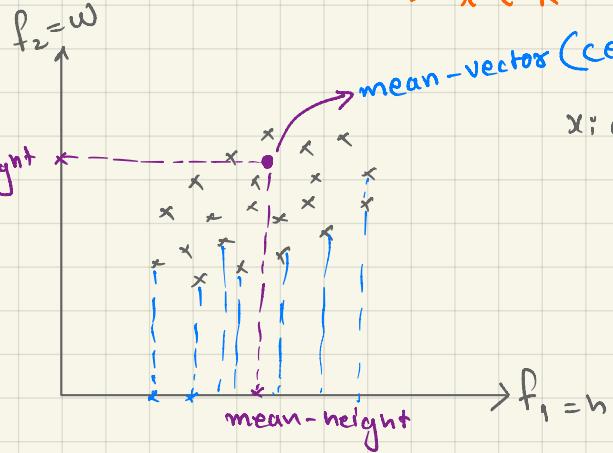
$$\vec{x}_1 + \vec{x}_2 = [3.4, 7.4]$$

→ mean vector

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n \vec{x}_i = \frac{1}{n} (\vec{x}_1 + \vec{x}_2 + \vec{x}_3 + \dots + \vec{x}_n)$$

$$\bar{x} \in \mathbb{R}^d$$

Geo.



mean-vector (central-vector/value)

$$x_i \in \mathbb{R}^2$$

$$[h_i, w_i]$$

$$f_1, f_2$$

$$\bar{x} = [h_{\bar{x}}, w_{\bar{x}}]$$

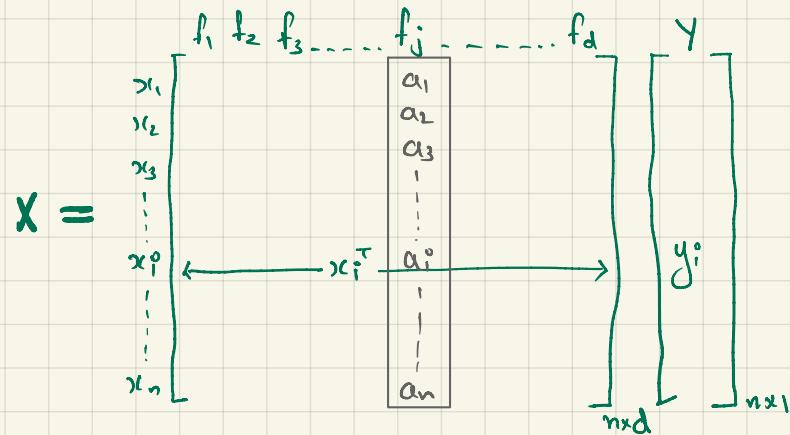
$$h_{\bar{x}} = \text{mean}(h_i)_{i=1}^n$$

$$w_{\bar{x}} = \text{mean}(w_i)_{i=1}^n$$

## Data preprocessing : column - standardization

Column - normalization :-  $[0, 1] \leftarrow$  get rid of scales of each feature

column - standardization is more often used in practice.



$a_1, a_2, a_3, \dots, a_i, \dots, a_n$   $\leftarrow$  n-value of  $f_i$   
 $\leftarrow$  any distribution

col. std.

$$a'_1, a'_2, a'_3, \dots, a'_i, \dots, a'_n \leftarrow \begin{cases} \text{mean } \{a'_i\}_{i=1}^n = 0 \\ \text{std-dev } \{a'_i\}_{i=1}^n = 1 \end{cases} \text{ standardization}$$

$$\bar{a} = \text{mean } \{a'_i\}_{i=1}^n \leftarrow \text{sample mean}$$

$$s = \text{std-dev } \{a'_i\}_{i=1}^n \leftarrow \text{sample std-dev}$$

$$a'_i = \frac{a_i - \bar{a}}{s} \rightarrow \begin{aligned} &\text{mean } \{a'_i\}_{i=1}^n = 0 \\ &\text{std-dev } \{a'_i\}_{i=1}^n = 1 \end{aligned}$$

$\rightarrow$  similar like Std-normal variate (z)

$$z = \frac{x - \mu}{\sigma}$$

$x \sim N(\mu, \sigma^2)$   $\rightarrow$  only diff. it comes from normal distn & above formula come from any distn.  
 $z \sim N(0, 1)$