## What is Dimensionality Reduction ?

Dimensionality Reduction techniques are built upon the idea of linear algebra.

In EDA, We learn that we can visualize our data in 2D and 3D using Scatter Plots.
For 4D,5D- 6D we can leverage Pair Plots. (nc2)
But For nD (10D) : Pair plots won't work
We reduce the dimensionality to (2d or 3d) to make it understandable so that we can visualize.
Some of techniques are *t-SNE* (almost state of art) and *PCA*(old technique)

## Row Vectors and Column Vector

For our iris flower dataset,
We were given 4 features or 4 variables : [SL,PL,SW,PW]

$\mathbb{R}$ : real space

Row -vector & column -vector

flower: $[SL, PL, SW, PW]$ — real -values

i\th point: $x_i \in \mathbb{R}^{(d)}$ → d-dim. column vector

$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{id} \end{bmatrix}_{d \times 1}$ : column-vector

$f1 = \begin{bmatrix} 2.1 \\ 3.3 \\ 1.6 \\ 4.3 \end{bmatrix}$

column-vector

$x_i \in \mathbb{R}^d$
↳ column-vector

$x_i = [2.1, 3.2, 4.6, 1.2]_{1 \times 4}$ : row-vector

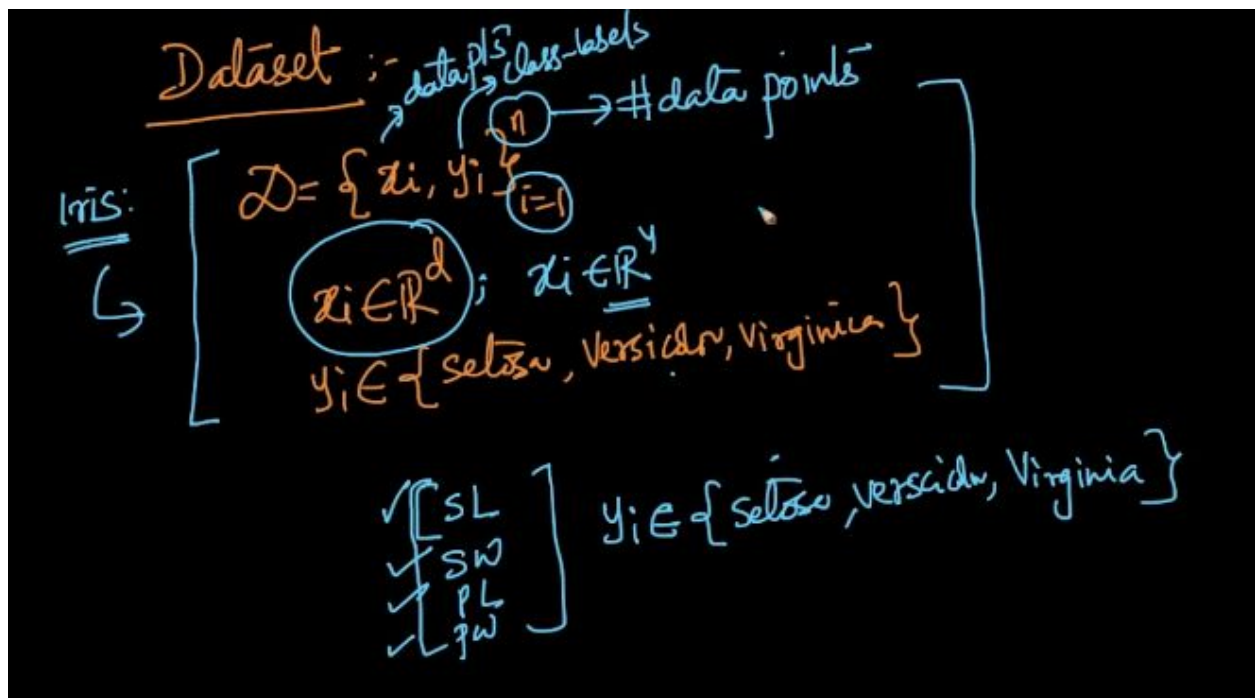## How to represent a dataset?

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

D is a collection of data points $x_i$ and class levels $y_i$ .

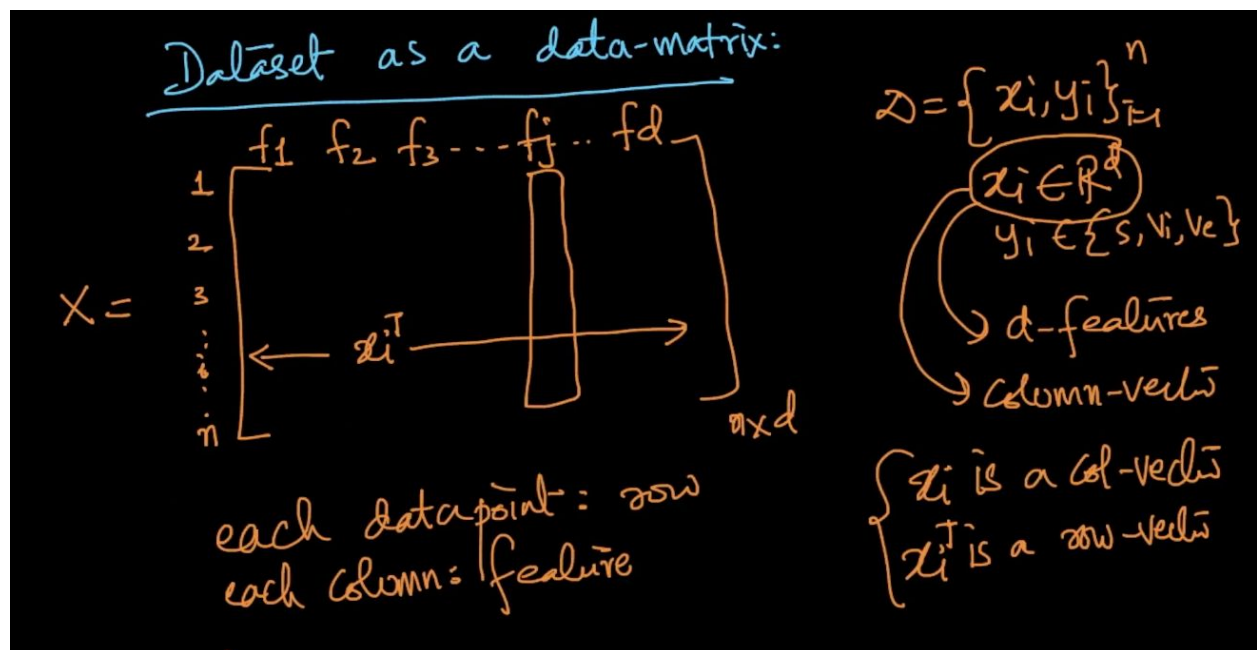Each data points belongs to $|R^d$

In case of iris data set, $X_i \varepsilon |R^4$
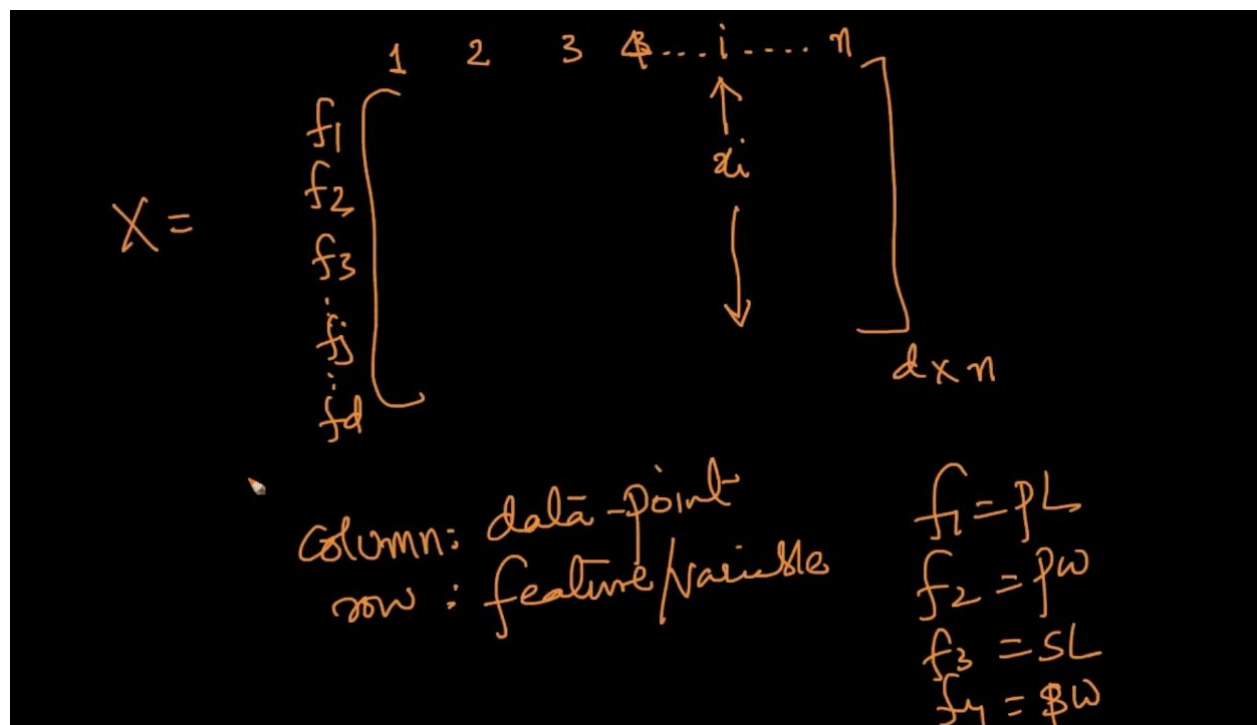
[SL,SW,PL,PW]

$y_i \varepsilon$ {Setosa, versicolor,virginica}



## How to represent a dataset as a Matrix ?

Data Set can be represented as a data matrix:

# Dataset as a data-matrix:

$$X = \begin{array}{c} \\ 1 \\ 2 \\ 3 \\ \vdots \\ n \end{array}
\begin{array}{cccccc} f_1 & f_2 & f_3 & \cdots & f_j & \cdots & f_d \end{array}
\left[ \begin{array}{c} \xleftarrow{\quad x_i^T \quad} \\ \\ \end{array} \right]_{n \times d}$$

each datapoint : row
each column : feature

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

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

$y_i \in \{S, V_i, V_e\}$

→ d-features

→ column-vector

$\begin{cases} x_i \text{ is a col-vector} \\ x_i^T \text{ is a row-vector} \end{cases}$

In Research papers we mostly see:

$$X = \begin{array}{c} f_1 \\ f_2 \\ f_3 \\ \vdots \\ f_j \\ \vdots \\ f_d \end{array}
\begin{array}{cccccc} 1 & 2 & 3 & 4 & \cdots & i & \cdots & n \end{array}
\left[ \begin{array}{c} \\ \uparrow \\ x_i \\ \downarrow \\ \end{array} \right]_{d \times n}$$

column : data-point
row : feature/variable

$f_1 = PL$
$f_2 = PW$
$f_3 = SL$
$f_4 = SW$

Let's Look into big picture:





We will be using this format on our hands-on assignment and course work

## Data Preprocessing: Feature Normalisation/ Column Normalisation
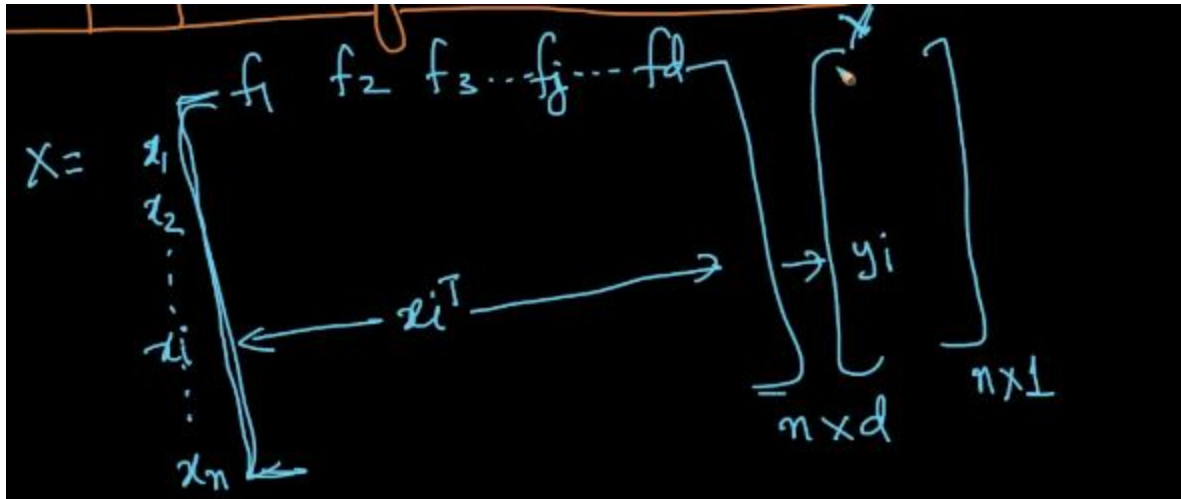
Technique to squash most of the data into unit cube or cuboids with the aim of getting rid of scales such as kg,cm, pounds to make data modelling easier.



Preprocessing means some types of mathematical operations / transformation done on data itself after obtaining data and before doing data modelling(dimension reduction).

**What**



$$X = \frac{1}{2}, 3 \quad \cdots \quad 150 = n$$

columns: $f_1 \, f_2 \, f_3 \, \cdots \, f_j \, (IPL) \, \cdots \, f_d$

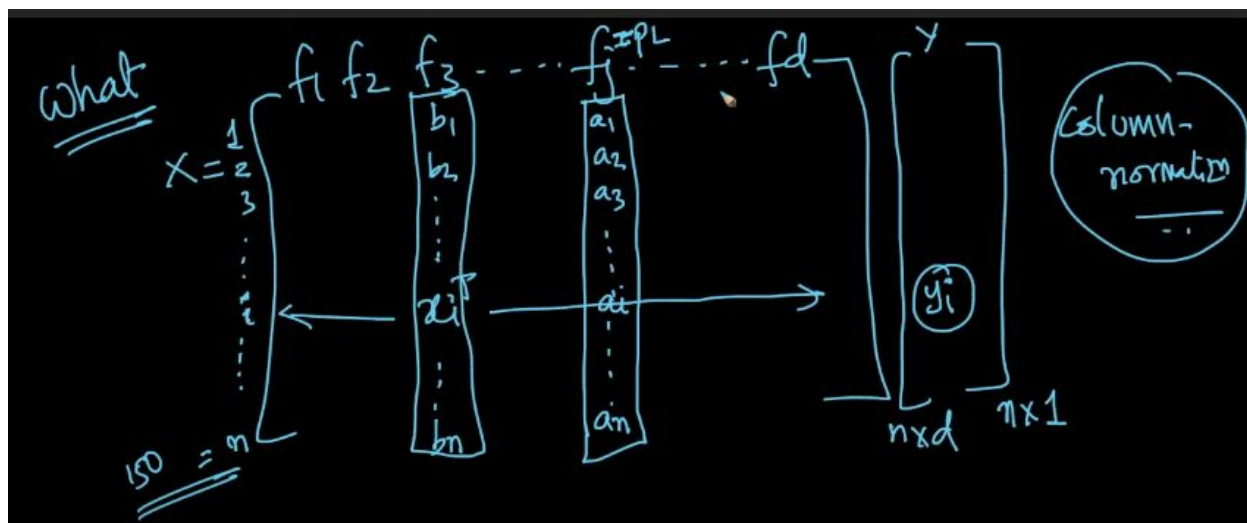$$\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ a_i^T \\ \vdots \\ b_n \end{bmatrix} \quad \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_i \\ \vdots \\ a_n \end{bmatrix}$$

$y$: $\begin{bmatrix} \vdots \\ y_i \\ \vdots \end{bmatrix}$

$n \times d$    $n \times 1$

Column-normalism

---

column: 1.2   1.3   1.4   1.9   1.7

$(a_1, a_2, \cdots, a_i, \cdots, a_n)$ → $n$-values of $f_j$

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

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

$(a_1', a_2', a_3', a_4' \cdots, a_i', \cdots a_n')$

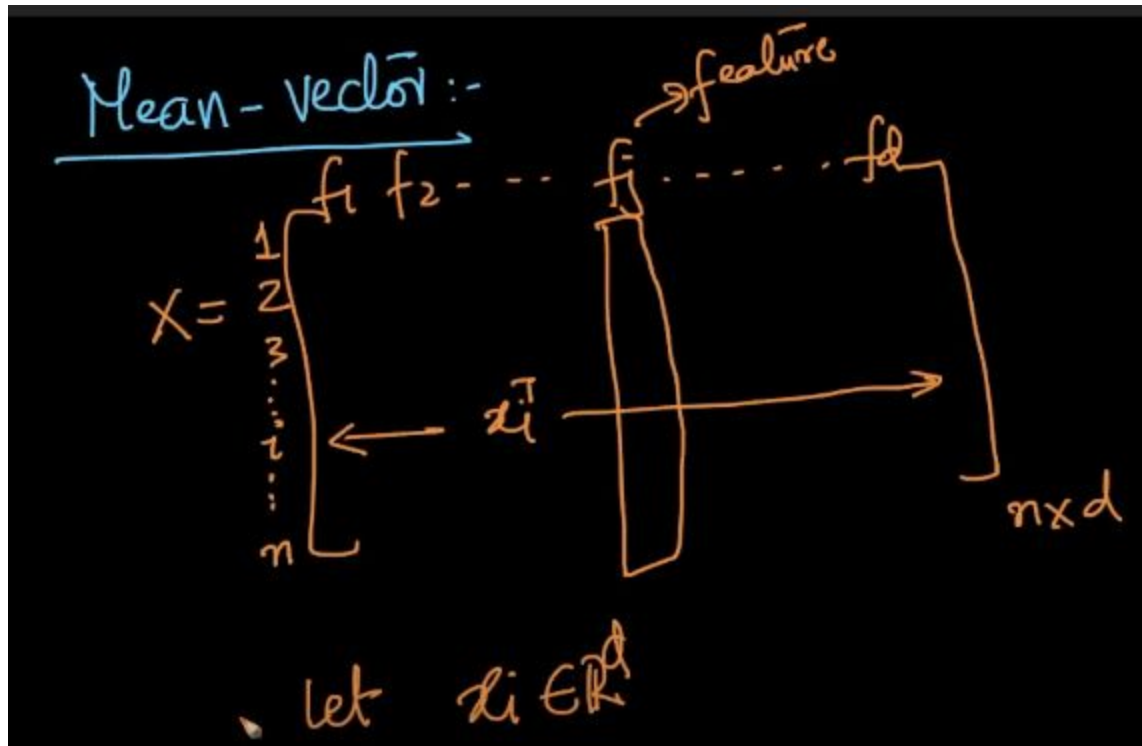$$a_i' = \frac{(a_i) - a_{min}}{a_{max} - a_{min}}$$

$a_i' \in [0, 1]$

$a_{min}' = \dfrac{(a_{min}) - a_{min}}{a_{max} - a_{min}} = 0 \, ;$

$a_{max}' = \dfrac{(a_{max}) - a_{min}}{a_{max} - a_{min}} = 1$

---

$a_1, a_2, \cdots, a_i \cdots a_d ; \, a_i \in \mathbb{R}$

$\downarrow$ column-normalization

$a_1', a_2' \cdots a_i' \cdots a_d' \; ; \; s.t \; a_i' \in [0,1]$

$f_1 = h$    $W = f_2$

| student → | 162 | 56 | 1 |
| | 172 | 72 | 2 |
| | 182 | 84 | |
| | 150 | 58 | |
| | ⋮ | ⋮ | |
| | | | n |

↑ cm / m / ft    ↑ kg / lbs

col-normaliz$^n$ → (getting rid of scale) (same scale)

$f_1' = h'$    $f_2 = W'$

1, 2, 3 ... → [0, 1]

{0, 1} ✓    ✓

---

**Geom:**

$f_2 = W$

65 kgs - - - [scatter cluster]

170cm → $f_1 = h$

$f_2$

h, w, bloodsf

[scatter points] $f_1$

$f_3$

column normaliz$^n$ →

$f_2' = W'$

unit-sq.

1 [scatter cluster in unit square]

0   1   $f_1' = h'$

$f_2'$

[cube with scatter points] $f_1'$

$f_3'$

---

anywhere in n-dim space   $\xrightarrow[\text{norm}]{\text{col.}}$   unit hyper cube in n-dim-space

# Mean of data matrix

Mean - vector :-

→ feature

$$X = \begin{matrix} & f_1 \; f_2 \cdots \; f_i \cdots \cdots f_d \\ 1 \\ 2 \\ 3 \\ \vdots \\ i \\ \vdots \\ n \end{matrix} \qquad n \times d$$

$\leftarrow x_i^T$

• let $x_i \in R^d$

$$x_1 = [\overset{f_1}{2.2}, \overset{f_2}{4.2}] \in R^3$$

$$x_2 = [1.2, \; 3-2] \in R^2$$

$$\frac{(x_1 + x_2) = [3.4, \; 7.4 \;]}{\bar{x} \in R^d}$$

$$\boxed{\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i} \qquad x_i \in R^d$$

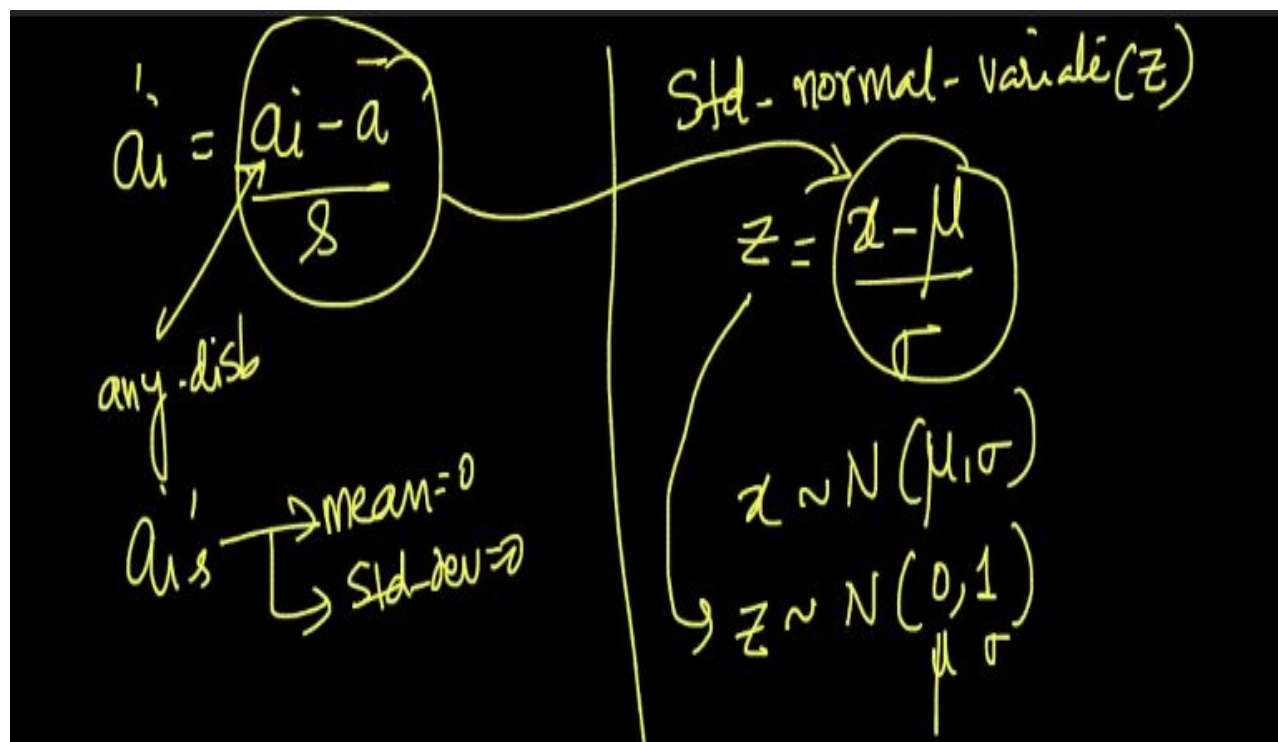$$= \frac{1}{n}\left(x_1 + x_2 + \cdots + x_n\right)$$

Mean-
vector

## Data Preprocessing: Column Standardization

Like Column normalization = [0,1] = get rid of scales of each feature

Column standardization = Most often used in practice.

$\sqrt{a_1, a_2, \ldots a_i, \ldots a_n}$ $\longleftarrow$ any-disb

$\downarrow$ col-std

$a_1', a_2' \ldots a_i' \ldots a_n'$ $\longrightarrow$ $\begin{cases} \text{mean} = 0 \\ \text{std-dev} = 1 \end{cases}$

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

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

$a_i' = \boxed{\dfrac{a_i - \bar{a}}{s}}$

(ex:) $\text{mean}\{a_i'\}_{i=1}^{n} = 0$

$\text{std-dev}\{a_i'\}_{i=1}^{n} = 1$

---

$a_i' = \boxed{\dfrac{a_i - \bar{a}}{s}}$

any-disb

$a_i' \longrightarrow \begin{cases} \text{mean} = 0 \\ \text{std-dev} = 0 \end{cases}$

Std-normal-variable $(z)$

$z = \boxed{\dfrac{x - \mu}{\sigma}}$

$x \sim N(\mu, \sigma)$

$z \sim N(0, 1)$
  $\mu \quad \sigma$

Geom: $f_2=w$

mean-wght

$\sigma=5$   $\sigma=0.5$

$f_1=h$

mean-hght

col·std $\longrightarrow$

$f_2'=w'$   std-dev=1

Col. std $\longrightarrow$

std-dev=1

$f_1'=h'$

mean-vect [0,0] ↑ origin

① moving the mean-vector to origin
② squishing/expanding s.t std-dev fw any feature is 1



Col. Standardizatⁿ :- mean - centering → origin
+ scaling → std-dev = 1 fw all features

## Covariance Matrix of Data Matrix



Co-variance matrix

$f_1\ f_2\ \cdots\ f_j\ \cdots\ f_d$

$X = \frac{1}{2}$

$x_i - i$

$n$

$x_{ij}$

$n \times d$

def:

$S=$ Cov-mat of $X$

$S_{ij}$

$d \times d$

Square-matrix

$f_j$ = col-vector jth feature

$S_{ij}$ = ith row & jth col. element in S

$x_{ij}$ = jth feature for ith data point
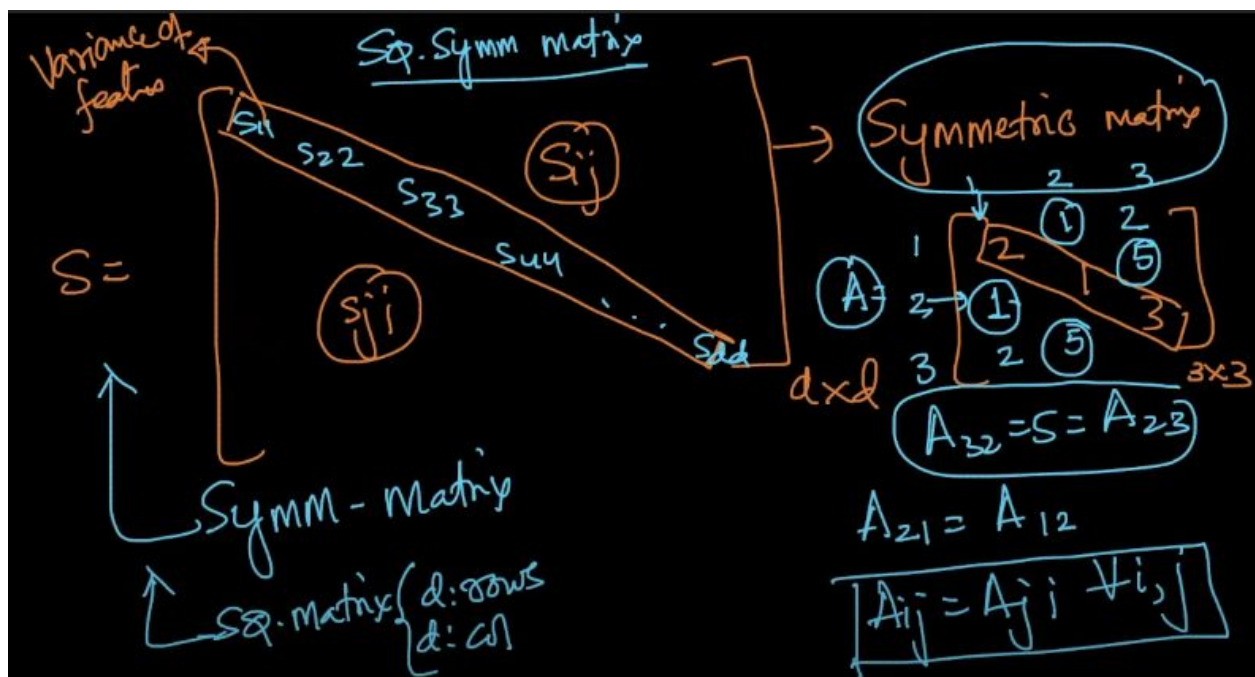
$$S_{ij} = cov(f_i, f_j)$$

$i: 1 \to d$
$j: 1 \to d$

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

$$cov(f_i, f_i) = \cancel{to} Var(f_i)$$

$$\begin{cases} \checkmark \; cov(x, x) = Var(x) \qquad ① \\ \checkmark \; cov(f_i, f_j) = cov(f_j, f_i) \qquad ② \end{cases}$$

Variance of features

Sq. Symm matrix

$$S = \begin{bmatrix} S_{11} & & & & \\ & S_{22} & & & \\ & & S_{33} & & \\ & & & S_{44} & \\ S_{ij} & & & & \ddots \\ & & & & S_{dd} \end{bmatrix}$$

$S_{ij}$

$S_{ji}$

Symm - Matrix

Sq. matrix $\begin{cases} d: rows \\ d: col \end{cases}$

Symmetric matrix

$$A = \begin{bmatrix} & 2 & 3 \\ 1 & 2 & ① & 2 \\ 2 \to & ① & 1 & ⑤ & 3 \\ 3 & 2 & ⑤ \end{bmatrix}$$

$d \times d$   $3 \times 3$

$$A_{32} = 5 = A_{23}$$

$$A_{21} = A_{12}$$

$$\boxed{A_{ij} = A_{ji} \; \forall i, j}$$

$$X = \begin{bmatrix} & f_1 \ f_2 & f_i \text{--} & f_d \\ 1 & & \\ 2 & x_{1i} & \to x_{1j} \\ \vdots & & \to x_{2j} \\ & x_{2i} & \\ n & & \end{bmatrix} \quad n \times d \qquad X \begin{Bmatrix} \to \begin{bmatrix} f_1 & f_2 \\ \Box & \Box \end{bmatrix} \end{Bmatrix}$$

$$\overset{\mu L \; \mu w}{\overset{x \; x}{}}$$

Let $\widehat{X}$ col. standardized $\Rightarrow$ mean$\{f_i\} = 0$
std-dev$\{f_i\} = 1$

$$Cov(f_1, f_2) = \underset{avg}{\frac{1}{n}} \overset{n}{\underset{i=1}{\sum}} \overset{\mu 1}{(x_{i1} - \mu_1)} \overset{\mu w \quad \to mean(f_2)}{(x_{i2} - \mu_2)}$$

$$\downarrow \; mean(f_1)$$

---

$$Cov(f_1, f_2) = \frac{1}{n} \overset{n}{\underset{i=1}{\sum}} x_{i1} * x_{i2}$$

$$X = \begin{array}{c} 1 \\ 2 \\ 3 \\ 4 \\ \vdots \\ n \end{array} \begin{bmatrix} f_1 & f_2 & \cdots & f_d \\ \boxed{\begin{array}{c} \bullet \\ \checkmark \\ x \\ \\ \boxed{x_{i1}} \end{array}} & \boxed{\begin{array}{c} 1 \\ \checkmark \\ x \\ \\ \boxed{x_{i2}} \end{array}} & & \end{bmatrix}$$

$$\overset{f_1 \cdot f_2}{\overbrace{\phantom{xxxx}}}$$

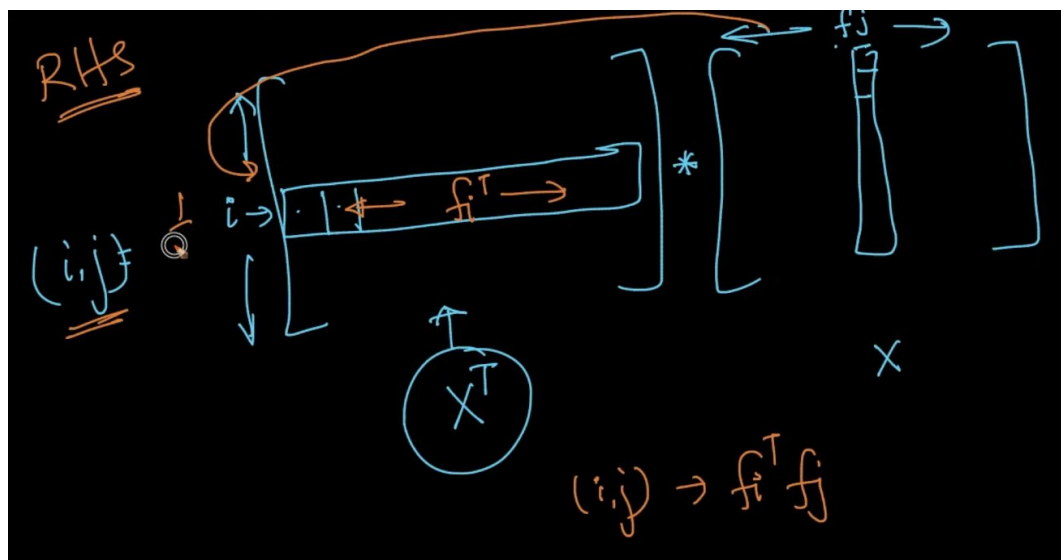$$Cov(f_1, f_2) = \left( f_1^{T} f_2 \right) * \frac{1}{n}$$

if $f_1$ & $f_2$ have been std,

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

$$S_{d \times d} = \frac{1}{n} (X^T)_{d \times n} (X)_{n \times d} = \boxed{d \times d} \checkmark$$

$\llcorner$ data-matrix

(*) assuming X has been col. std

$$S_{ij} = Cov(f_i, f_j) = \frac{f_i^T f_j}{n}$$

RHS

$(i,j)$ 



$$(i,j) \to f_i^T f_j$$

**MNIST Data Sets**

**Additional References:**

http://colah.github.io/posts/2014-10-Visualizing-MNIST/

(MNIST) dataset

Iris :- 4 dim. dataset

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^{60K}$$

$x_i :$ [image of 0] $28$ , $28$

Obj: Classify the written char into one of the 10 numeric char.

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

$x_i =$ [square] $28$ , $\leftarrow 28 \rightarrow$

$\rightsquigarrow \mathcal{D}$

$$x_i = \begin{bmatrix} \\ \\ \end{bmatrix}$$

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

$x_i =$ image $\Rightarrow \begin{bmatrix} & \\ & \end{bmatrix}_{28 \times 28}$

NOT data-Matrix X

Matrix representation of image

numerical/real matrix

$x_i = $ [image] ← 28 →, 28 ↕  →  $x_i = $ [vector]   $x_i \in \mathbb{R}^d$

784 × 1

row-flatting

$x_i = $ image ⇒ $\begin{bmatrix} 0 & 10 & 20 \\ 0 & 1 & 0 \\ 0 & 1 \end{bmatrix}$ 28×28

NOT data-Matrix X
Matrix representation of image

numerical/real matrix

1 2 3 4 5 (row-flatting) → 25 × 1

5×5