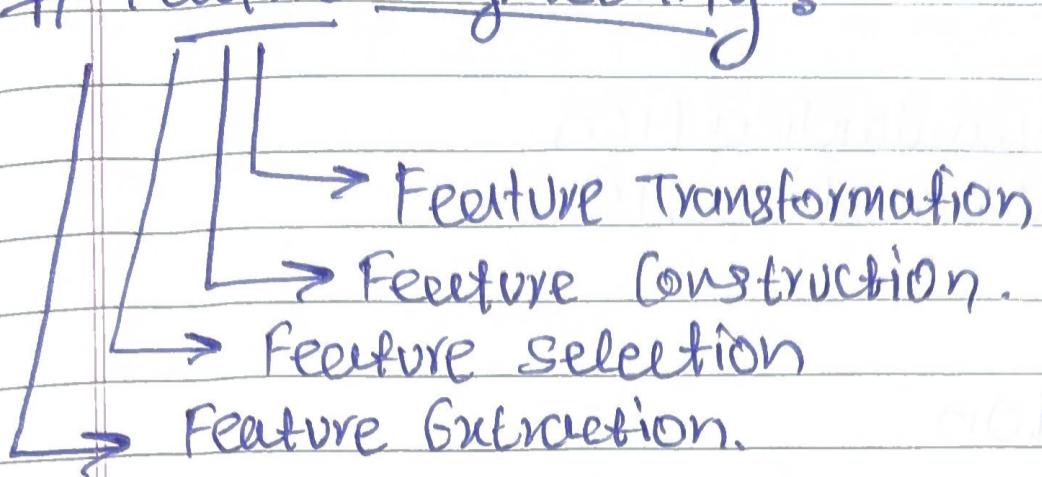


# DAT-5

## # Feature Engineering :

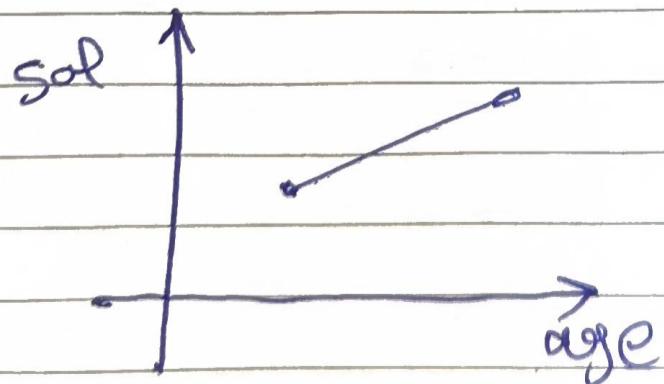


### ⇒ Feature Transformation :

- ↳ Missing value imputation.
- ↳ Handelling categorical features.
- ↳ Outlier Detection.
- ↳ Feature Scaling

### i) Feature Scaling :

| age | Salary | purchase |
|-----|--------|----------|
| 50  | 83000  | 1        |
| 27  | 48000  | 0        |



- for e.g: KNN  
when we take the euclidean distance the distance may vary too much for salary. so we scale in the same range, so that it works as a same unit.

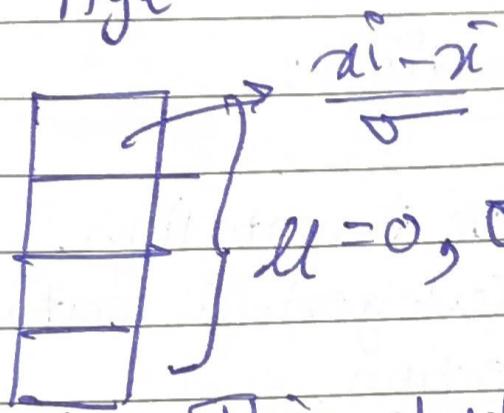
## Feature Scaling

- Scaling : 1) Standardization  
2) Normalisation.

### 1) Standardization.

Age    sal    Age'

|    |   |
|----|---|
| 27 | 3 |
| 15 | 3 |
| 33 | 2 |
| 63 | 0 |



→ This delta will satisfy

$$\text{Standardization} := \frac{x_i - \bar{x}}{\sigma} \quad \boxed{\mu = 0, \sigma = 1}$$

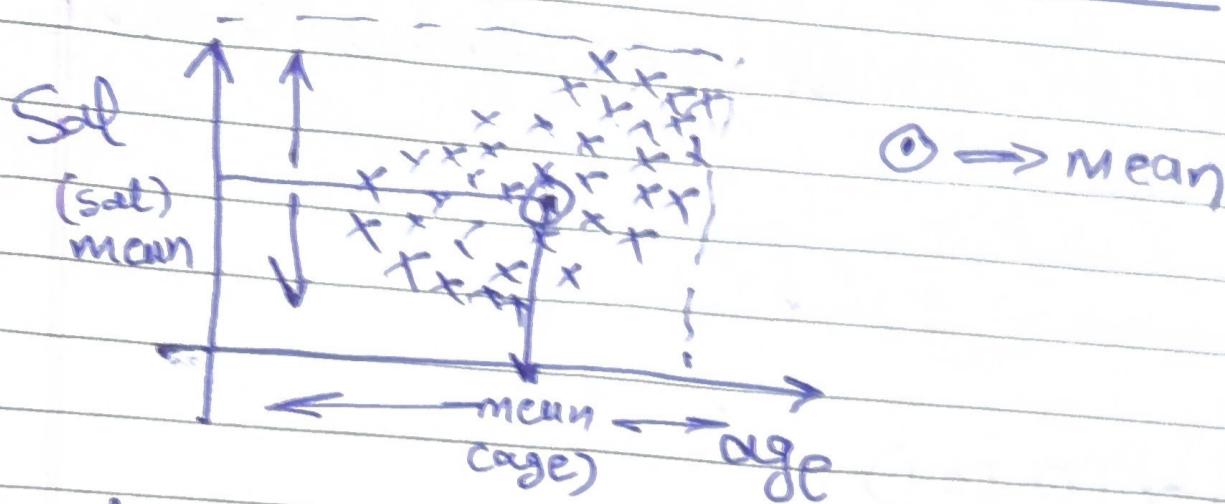
$$= \frac{27 - 32}{10} = -0.5$$

for example

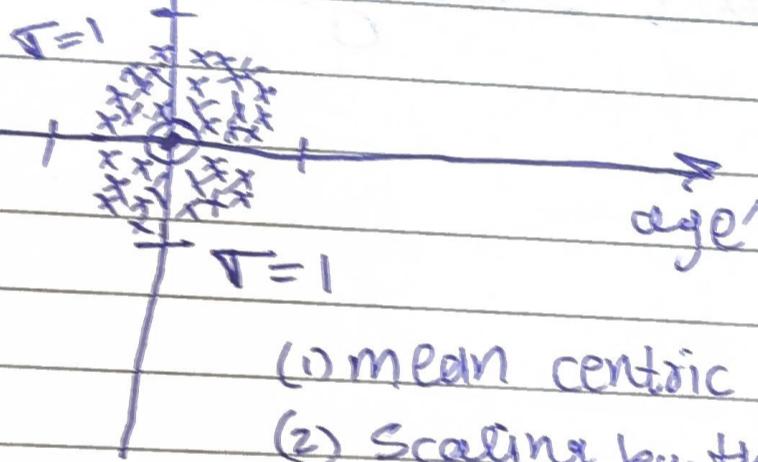
#

# Geometric Intuition of Standardization:

Update

Page No.  
Date

$\Downarrow$  Standardization



## # Normalisation :

$\Rightarrow$  we have to remove or eliminate the unit factors.

Key. N, ---  
 (- to +) or (0 to 1)

Techniques: MinMaxScaler  
 Mean Normalisation  
 Max absolute Scaling.  
 Robust Scaling.

### MinMaxScaler:

Weight

130

67

78

50

48

30

80

:

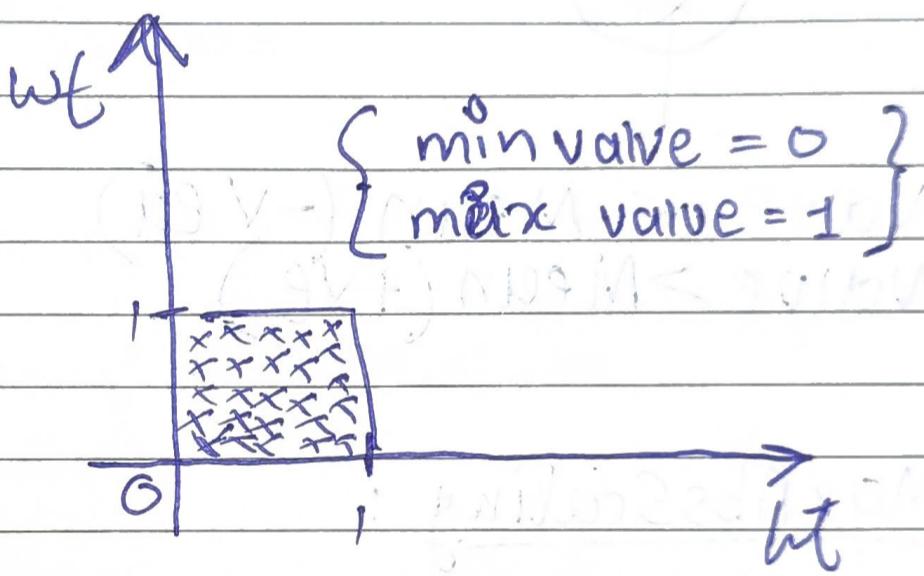
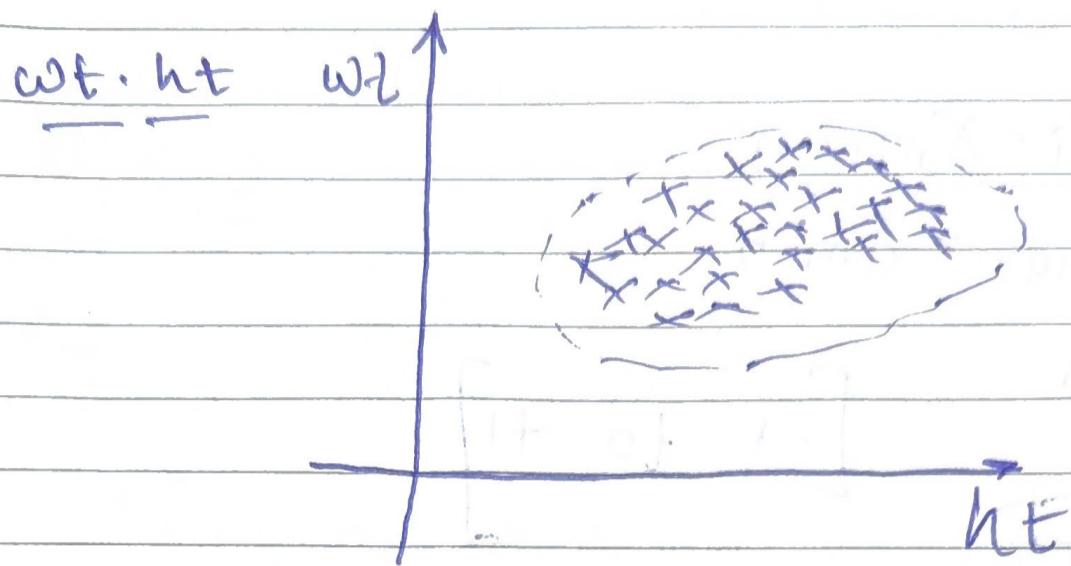
Normalising

$$\rightarrow x_i' = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \rightarrow (0 - 1)$$

$$(130)x_i' = \frac{130 - 30}{130 - 30} = 1$$

$$(50)x_i' = \frac{50 - 30}{130 - 30} = 0.2$$

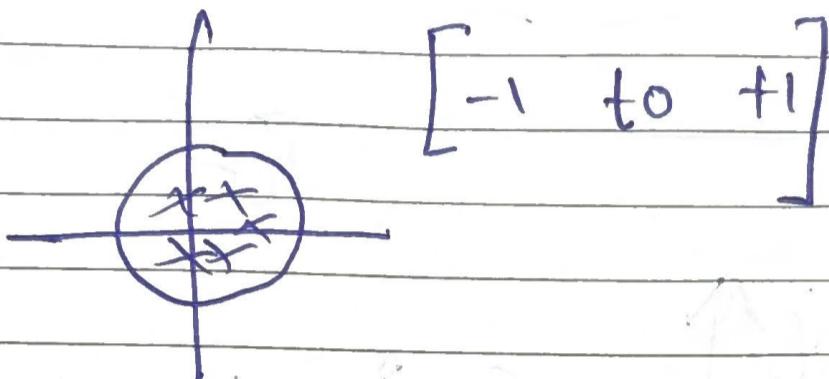
## # Geometrical Intuition:



2]

Mean Normalisation :

$$x_i' = \frac{x_i - \bar{x}_{\text{mean}}}{x_{\text{max}} - x_{\text{min}}}$$



Value &lt; Mean (-ve)

Value ≥ Mean (+ve)

3]

MaxAbsScaling :

$$x_i' = \frac{x_i}{|x_{\text{max}}|}$$

→ MaxAbsScaler  
↳ sklearn

Use : Sparse data.

→ zyadah zeros  
jis data me ho.

## 4) Robust Scaling

$$X'_i = \frac{X_i - \text{Median}}{\text{IQR}}$$

↳ Robust Scaler  
↳ Sklearn

- Robust to outlier.
  - ↳ for lot of outliers data.

## # Normalisation V/s Standardization.

Some questions to ask

- ↳ 1) Scaling required?
- 2) ~~Max~~ Maximum time standardization
- 3) Normalisation
  - ↳ CNN  $\rightarrow$  (0-255)

92