

31-08-2025

### Agenda:

→ ML algo high level overview

→ FE - III

### ML algo high level Overview:

System that can learn & improve from data without explicit  
programming

$f_1$ weight (kg)	$f_2$ Height (metre)	T Value - x	$(\omega / h^2) \rightarrow$
70	1.75	22.88	$\downarrow$
85	1.7	29.41	<u>2 variables</u>
50	1.8	15.43	

A	B	C	D	E	F	value - x
-	-	-	-	-	-	100
-	-	-	-	-	-	102.3

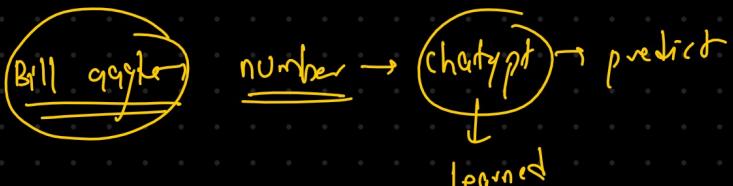
ML →

A	D	value - x
B	E	
C	F	using ML

Learned that formula → pattern / relation  
which will provide us  
value - x given  $[A \dots F]$



Data  
 wiki, stackflow, x, y, z



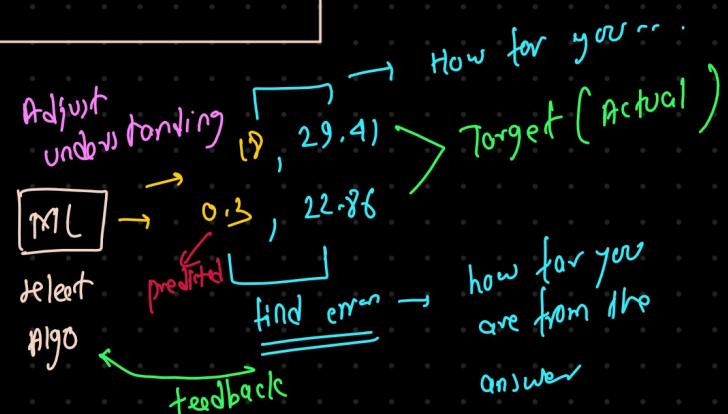
very very complex formula

Features	Target
○	Red
○	Green
○	Blue
○	Yellow

ML

- $I_1 \rightarrow 70, 1.75$
- $I_2 \rightarrow 85, 1.7$
- $I_3 \rightarrow 50, 1.7$
- $I_4 \rightarrow 80, 1.8$

Input values



random  $\rightarrow$  red  
 red

"Dog" & "Cat" → "Animal"

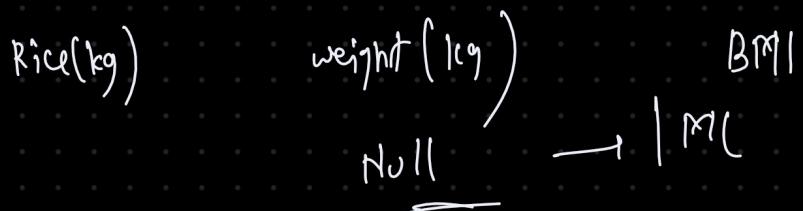
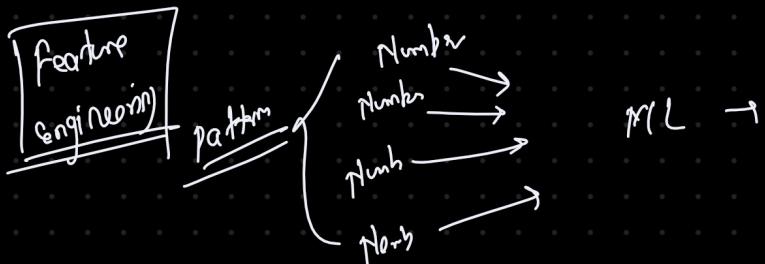
Human & AI & Engineer → AI Engineer

Let  $(v_1, v_2, v_3)$  :



$v_4$

calc & pointer understand → number



## Feature Scaling:-

KNN → distance-based algo

Height (metres)      weights (kg)

0.0 → 7.0      0 → 170

formula:  $a, b$

$\begin{bmatrix} a & 7.0 \\ b & 170 \end{bmatrix}$       ML

7.0 → 70

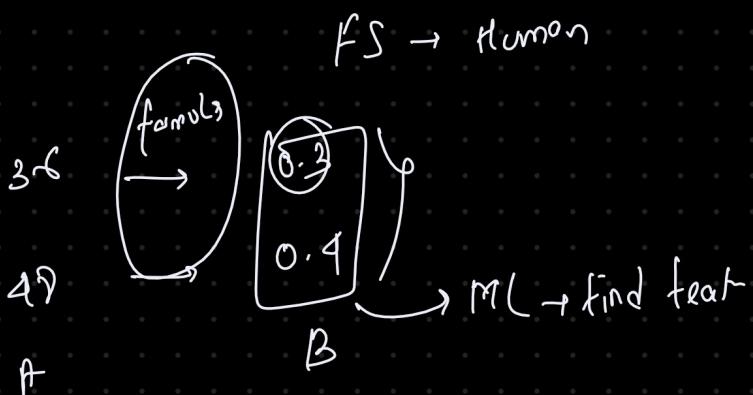
6.5 → 65

5.5 → 55

↓  
10\*

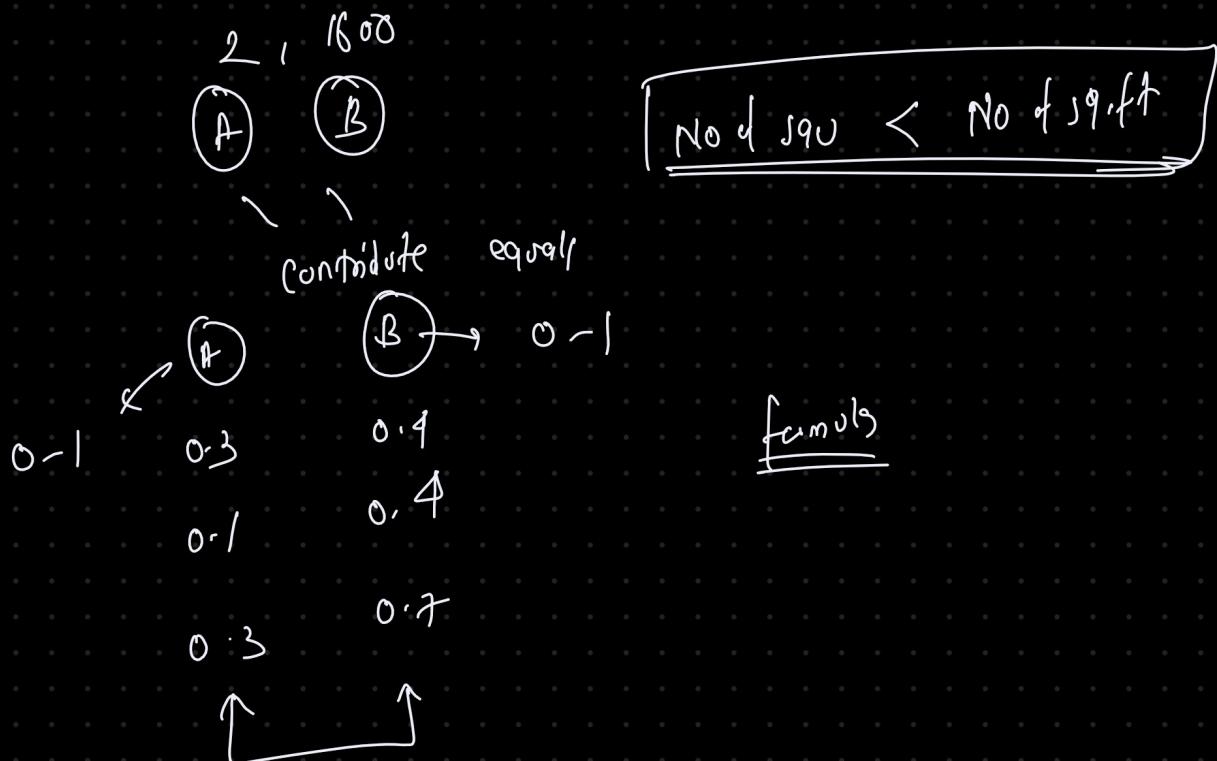
$a \rightarrow 0.7$

$b \rightarrow 0.3$



House (sq.ft)	No. of Bedrooms	Price
900	2	300k \$
900	1	250k \$
1600	2	500k \$

$$2, \overbrace{1, 900}^{\text{formula}} \rightarrow \overbrace{\text{MC}}^{\text{formula}} \rightarrow \overbrace{2, 900}^{\text{formula}}$$



	Age	Income
A	20	100,000
B	25	100,000
C	20	150,000

$$\text{distance} = \sqrt{(age_1 - age_2)^2 + (income_1 - income_2)^2}$$

$$\text{distance}(A, B) = 5 \rightarrow \text{close} \quad | \quad 5 \quad 50,000$$

$$\text{distance}(A, C) = 50,000 \rightarrow \text{far far}$$

	Age	Income
A	0.037	0.040
B	0.135	0.040
C	0.037	0.066

Distance  $(A, B) \approx 0.097$  ] )  $\xrightarrow{\text{min-max scaling}}$   
 Distance  $(A, C) \approx 0.026$  ] )

prob:

Math (out of 100)	Science (out of <u>200</u> )
(80) <u>0.75</u>	(150) <u>0.71</u>
60	120
90	180
0-1	0-1

Scaling #1 Min-Max Scaling (Normalization)  
 → scale values into  $[0, 1]$   
 → sensitive outliers

- A
- 70  $\rightarrow$  0
- 80  $\rightarrow$  0.2
- 90
- 120  $\rightarrow$  1

$$\text{formula} = x_{\text{new}} = \frac{x_{\text{old}} - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

$$x_{\text{new}} = \frac{70 - 70}{120 - 70} = \frac{0}{50} = \text{something}$$

$$1 = \frac{120 - 70}{120 - 70} = 1$$

$$= \frac{80 - 70}{120 - 70} = \frac{10}{50} = 0.2$$

## (2) Standardization (Z-score scaling)

\* handles outlier

better than Min-Max

$$x_{\text{new}} = \frac{x_{\text{old}} - \mu}{\sigma}$$

$$\text{Mean} = 0, \quad \text{Std Dev} = 1$$

## (3) Robust scaler (Median & IQR based)

$$x_{\text{new}} = \frac{x_{\text{old}} - \text{Median}(x)}{\text{IQR}(x)} \quad \text{IQR} = Q_3 - Q_1$$

## (4) log transformation

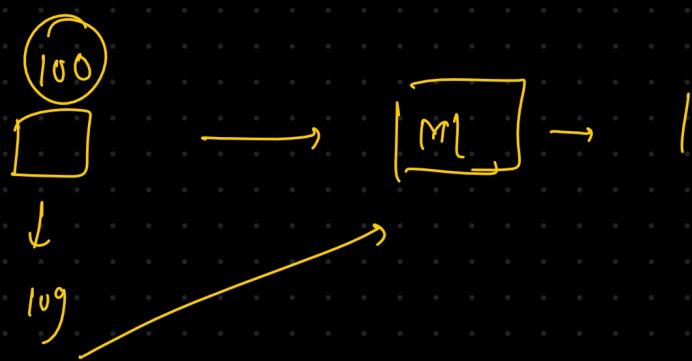
$$x_{\text{new}} = \log(1 + x)$$

\* right skewness

## (5) Others

* square root	$\rightarrow x_{\text{new}} = \sqrt{x_{\text{old}}}$	* Box-Cox ( $\lambda$ )
* cube-root	$\rightarrow x_{\text{new}} = x^{1/3}$	$x_{\text{new}} = (x_{\text{old}}^\lambda - 1)/\lambda, \forall \lambda \neq 0$
* Yeo-Johnson		$\log(x)$ if $\lambda = 0$

	<u>No. of apples</u>		
A	100	2.5	Y
B	20	1.9	ML
C	40	1.9	NN
D	600	5.6	CV
E	584	4.2	NLP



Raw file count per plate      Chapatji

Conclusion : —

- doesn't change the underlying representation
- Algo can learn patterns effectively
- Make feature more consistent across data.

## Data Encoding : —

→ "Blue", "green", "red"



### #1 Label Encoding (Nominal)

"Blue"		"green"		"red"
↓		↓		↓
0		1		2

- Data-order dependant

### #2 Label Encoding (Ordinal) / Ordinal Encoder

"Red"	,	"Green"	,	"Yellow"
↓		↓		↓
0		1		2

Danger sort them in Desc : —      "Red", "Yellow", "Green"  
     ↓                                    ↓                                    ↓  
     0                                    1                                    2

Student :

Marks	Habit	Weight	fav-colour	behavior
1	1	1	blue	good
2	2	2	pink	bad
3	3	3	white	average



$\text{Red} > \text{Green}$   
 $4 > 3$

# 3 OHE (One-hot encoding)

	A	B	C	D
	0	1	2	3
A	0	0	0	1
B	0	0	1	0
C	0	1	0	0
D	1	0	0	0

category into  
binary dummy variable

$\text{id} \rightarrow \underline{\text{color}} \rightarrow \begin{matrix} \text{Blue} \\ \text{Green} \\ \text{Yellow} \end{matrix}$

<u><math>\text{id}</math></u>	<u><math>\text{color}</math></u>	<u><math>\text{color\_red}</math></u>	<u><math>\text{color\_green}</math></u>	<u><math>\text{color\_yellow}</math></u>
1	Red	0	0	1
2	Green	0	1	0
3	Yellow	1	0	0

- disadvantage

## Curse of dimensionality

dimension of the data

Column → 1000 unique string

$$5 \rightarrow \begin{array}{ccccccccc} | & | & | & | & | & | & | & | & | \end{array}$$

$$+ \boxed{1000}$$