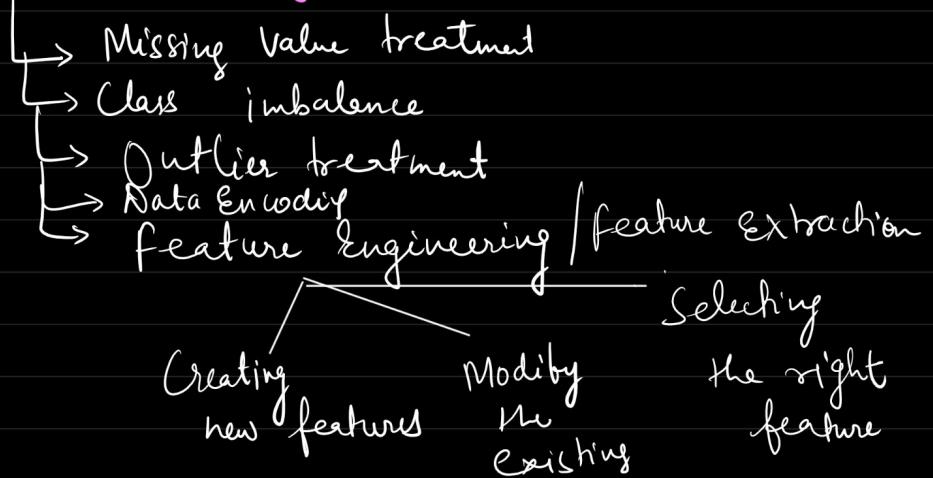


- \* Understanding of problem statement
- \* Data Ingestion.
- \* Data Preparation | Preprocessing & Exploratory data Analysis



\* Feature extraction — process of selecting and extracting the relevant features from the raw data.

1000 features → Most important features  
 ↓  
 Used to train the model.

\* Curse of dimensionality

↳ with increase in no. of feature :-

- ① Model training becomes computationally expensive
- ② Model interpretation becomes complex.

<u>Area of house</u>	<u>Price of a house</u>	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	---	$f_{1000}$
1800	60							
1900	70							
—	—							
—	—							

## ① Creating new feature

	distance	time	Speed = Distance / time
Ex-1	50	2	$50/2$
	100	3.5	$100/3.5$
	150	5.8	-
	110	2.1	-
	-	-	-

	temp	3 month moving avg temp
Ex-2	21	Nan
	22	Nan
	23	$\rightarrow 21 + 22 + 23 = 22^{\circ}$
	18	$\rightarrow \frac{21 + 22 + 18}{3}$
	20	$\rightarrow \frac{22 + 23 + 20}{3}$
	21	
	22	
	18	

new feature  
↓

## Ex: Ratio

$$\text{No. of room} = \frac{\text{eff area}}{\text{Area of 1 room}}$$

	No. of Room	Area of 1 room	effective area
	2	100	200
	3	200	600

## ② Modifying the existing feature

ⓐ Changing the datatype

Age	modified Age
15 yrs	15
20 yrs	20
25 yrs	25
-	-
-	-

install	install-modified
{ 15+	15
{ 20+	20
{ 25+	25

eg :

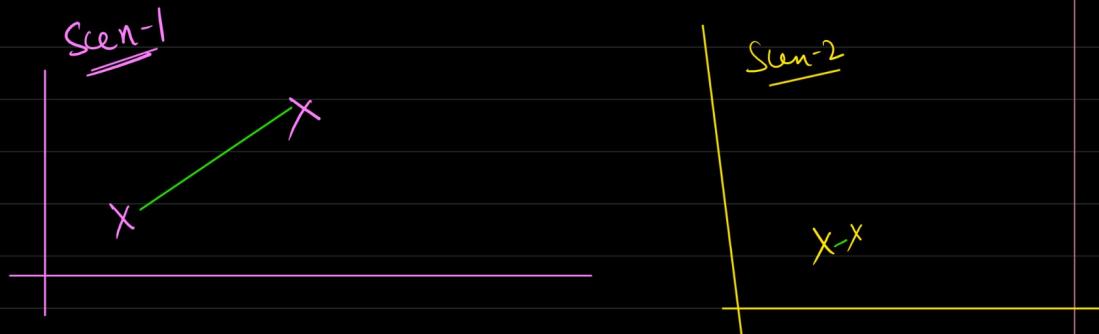
	Date	Day	Month	Year
15-2-2022	15	2	2022	
16-2-2022	-	-	-	
+ -	-	-	-	
- -	-	-	-	
- =	-	-	-	
	-	-	-	

## (b) feature scaling (optional)

→ feature Scaling  
tries to bring  
all the features  
on the same  
scale

<u>Area of house</u>	# of rooms	Parking area	<u><math>y</math> (price of a house)</u>
1800	2	100	60
2500	3	50	65
3000	4	25	100
6000	1	30	-
-	5	30.5	-
-	-	-	-
-	-	-	-

Why?



\* Many of Algo's are distance based, that's why if high magnitude it becomes computationally expensive.

\* optimisation becomes faster.

\* Interpretation becomes easier.

## Types of feature scaling

### ① Standardisation (ML algo)

$$\text{Z score} = \frac{x_i - \bar{x}}{\sigma}$$

$$SND \Rightarrow \mu = 0, \sigma = 1$$

<u>Age</u>	<u>Age_Scaled</u>	
25	$\frac{25 - 24.6}{1.03}$	/
26	$\frac{26 - 24.6}{1.03}$	-
23	$\frac{23 - 24.6}{1.03}$	=
24	-	-
25	-	-
$\bar{x} = 24.6$		
$\sigma = 1.03$		

② Normalization  $[\min \max \text{ scalar}] \Rightarrow [0, 1] \rightarrow \text{DL Algorithm} \rightarrow \underline{\text{PCA}}$

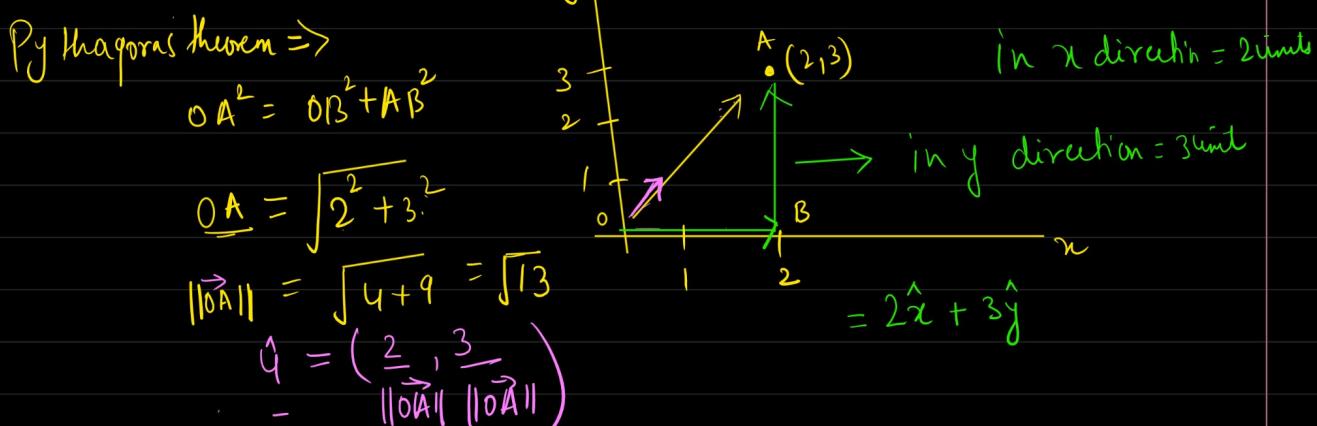
$$\begin{array}{ll} \text{Age} & \text{Age_Scaled} \\ 25 & \frac{25 - 23}{26 - 23} = \frac{2}{3} \\ 26 & \frac{26 - 23}{26 - 23} = 1 \\ 23 & \frac{23 - 23}{26 - 23} = 0 \\ 24 & \\ 25 & \end{array} \quad \chi_{\text{scaled}} = \frac{\chi_i - \chi_{\text{min}}}{\chi_{\text{max}} - \chi_{\text{min}}} = [0, 1]$$

$\min = 23$   
 $\max = 26$

### ③ Unit Vector

$$(2, 3)$$

$\downarrow$      $\downarrow$   
 $x$      $y$



$$\hat{u} = \left( \frac{2}{\sqrt{13}}, \frac{3}{\sqrt{13}} \right) \Rightarrow \sqrt{\left(\frac{2}{\sqrt{13}}\right)^2 + \left(\frac{3}{\sqrt{13}}\right)^2}$$

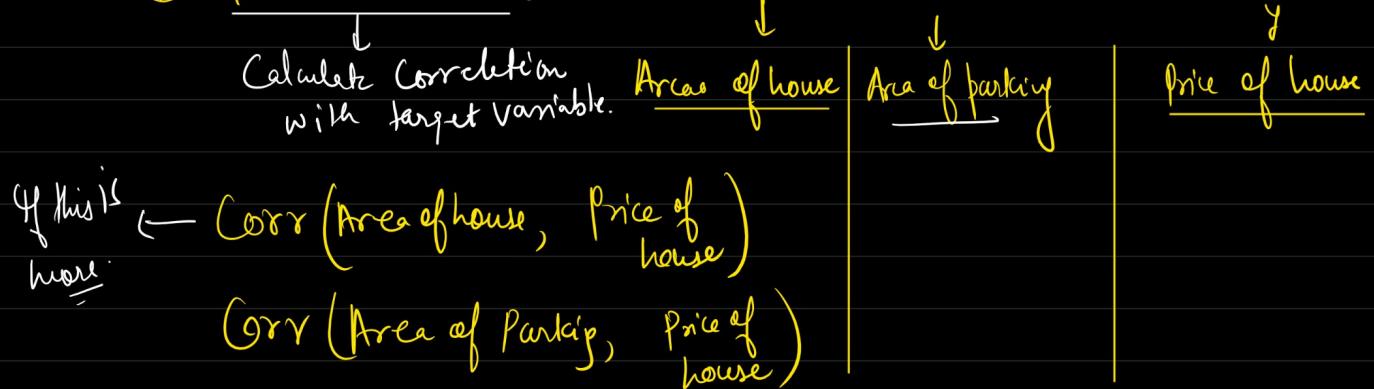
$$= \sqrt{\frac{4}{13} + \frac{9}{13}} = \sqrt{1} = 1$$

### ③ Selecting the right feature (feature selection)

1000 features → Top 10 features

↓  
training of model

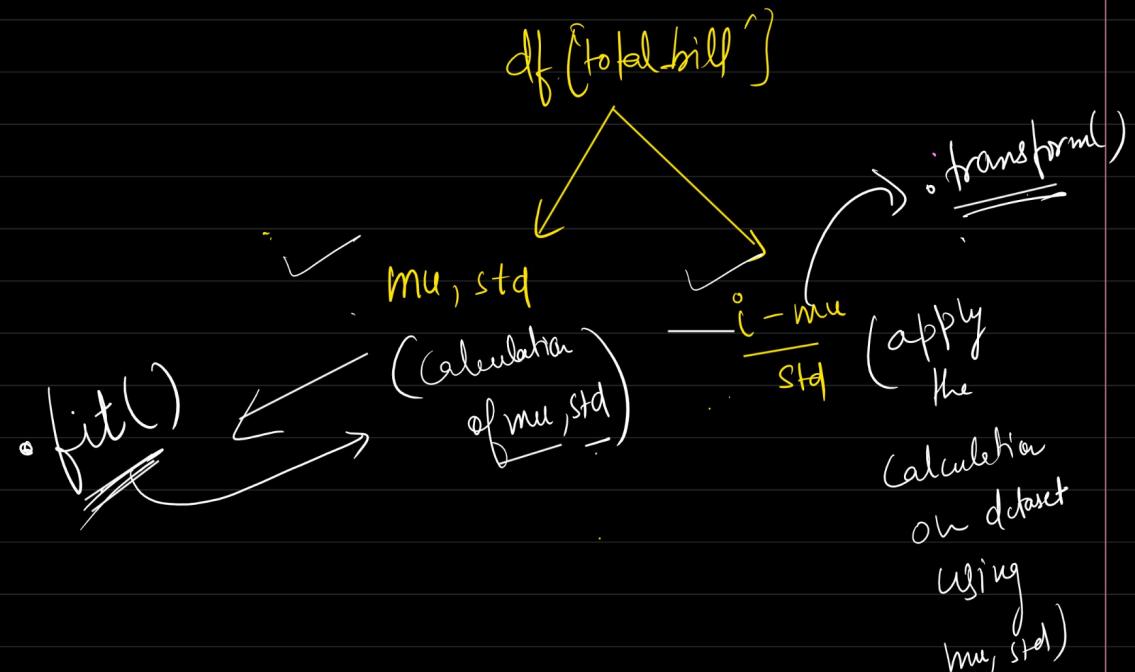
#### ① Filter method :



#### ② Embedded method

#### ③ Wrapper method.

\* PCA





You don't know this date  $\left\{ \begin{array}{l} \text{representation} \\ \text{of unseen test} \end{array} \right.$   $\left[ \begin{array}{c} \text{transform} \\ = \end{array} \right]$

## Data Encoding

1. Nominal | OHE
  2. Label and Ordinal
  3. Target Guided ordinal encoding.
- ML Algorithms  
↓  
numerical data

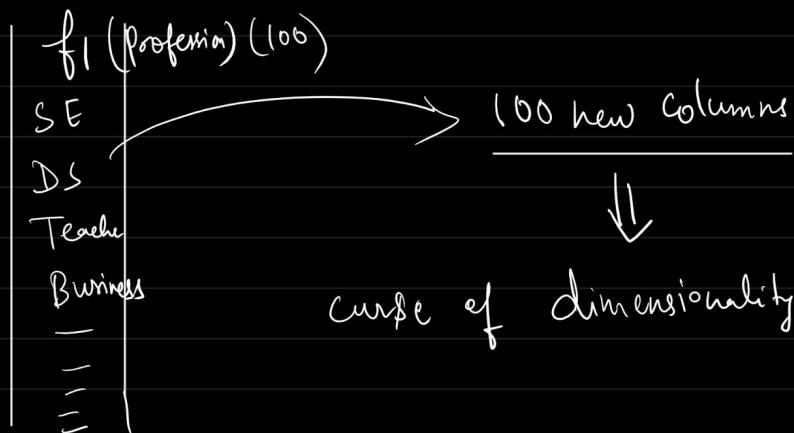
1. Nominal | OHE  
 $\rightarrow$  Categorical data to Numerical data.  
 $\rightarrow$  No order in the data

Status

Single	ML Algorithms
married	
Separate	
Single	
Single	
married	
Separate	
	Single      married      Separated
	1            0            0
	0            1            0
	0            0            1

	red	green	yellow
red	1	0	0
green	0	1	0
yellow	0	0	1

\* disadvantage → A column has many categories



## ② Label and ordinal encoding techniques.

\* Label encoding → assign numerical label to each category.

Red - 1

Green - 2

Yellow - 3

DS - 1

DE - 2

DA - 3

BA - 4

1
2
3
4
1

\* No problem of curse of dimensionality  
disadvantage - It will learn the pattern

## \* ordinal encoding

High School - 1

College - 2

Post Graduate - 3

### ③ target guided ordinal Encoding

- based on their relationship with the target Variable
- Useful when we have large no of Unique Categories in Categorical Value.
- Categorical group with mean |  
median  
of corresponding target variable.

<u>time</u>	<u>total bill</u>
150 ← lunch -	groupby('time').
180 ← dinner	mean
120 ← breakfast	
150 ← lunch	
180 ← dinner	