

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

1000 features → All the features you cannot use model training.

* Curse of dimensionality

With increase in no of features:-

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

$f_1 \ f_2 \ f_3 \dots \ f_{1000}$

① Creating new feature

ex.1	distance	time	$\left\{ \begin{array}{l} \text{Speed} = \text{Distance} / \text{Time} \\ 50/2 \\ 100/3.5 \\ \vdots \end{array} \right. \right\}$
	50	2	
	100	3.5	
	-	-	

ex.2	temp	3 months moving avg temp.
	21°	-NaN
	22°	NaN
	23°	
	24°	$\frac{21+22+23}{3} = 22^{\circ}$
	25°	
	-	
	-	

No of rooms	Area of room	effective carpet area
1	100	200
2	200	600
3	100	-
4	-	-
5	-	-
-	-	-
-	-	-
-	-	-

② Modifying existing feature

ⓐ Change the data type

Age	Modified age
15 year	15
20 yrs	20
25 yrs	25
-	-
-	-
-	-

eg install	Modified-install
15+	15
20+	20
25+	25
-	-

eg Date	Day	Month	Year
15-2-2022	15	2	2022
16-2-2022			
-			
-			

* feature Scaling (optional) / transformation

→ To bring all the features on same scale.

Why?



Scaling / not do scaling
↓ no affect

<u>Area of room</u>	<u># of rooms</u>	<u>Parking area</u>	<u>Y (price of house)</u>	<u>Prediction</u>
1800	2	100	66	
2800	3	50	65	
3000	4	25	100	
6000	5	1	-	
((((
))))	

ML → Mathematical relationship

↓
Computation

$$\begin{cases} 2 \times 3 = 6 \\ 5 \times 2 = 10 \\ 3 \times 5 = 15 \end{cases}$$

$250 \times 328 = ?$

Why Scaling?

- Computation is less expensive.
- Gradient descent ↴ Optimisation becomes faster
- ⇒ Interpolation becomes easy. faster



* Type of scaling

① Standardisation (ML algorithms)

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

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

Age Age-Scaled

$$25 \rightarrow \frac{25 - 24.6}{1.03}$$

$$26 \rightarrow \frac{26 - 24.6}{1.03}$$

$$23 \rightarrow \frac{23 - 24.6}{1.03}$$

$$24 \rightarrow \frac{24 - 24.6}{1.03}$$

$$\bar{x} = 24.6$$

$$\sigma = 1.03$$

② Normalisation (min-max Scaler) \rightarrow DL Algorithms
 ↓
 $[0, 1]$

Age Age Scaled

25	$\rightarrow \frac{25-23}{26-23} = \frac{2}{3}$	
26	$\rightarrow \frac{26-23}{26-23} = 1$	
23	$\rightarrow \frac{23-23}{26-23} = 0$	
24	$\rightarrow \frac{24-23}{26-23} = \frac{1}{3}$	

$\min = 23$ $0.$

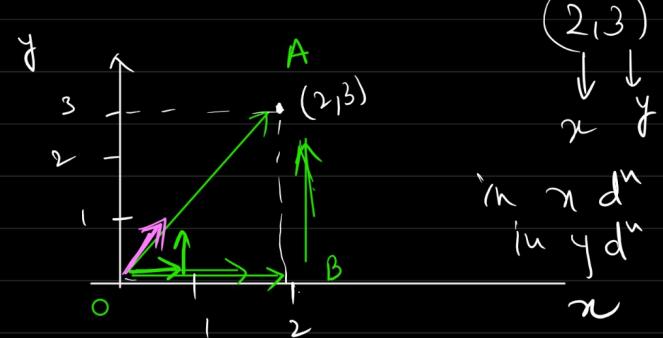
$\max = 26$ 1

$$y_{\text{scale}} = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

③ Unit Vector

$$\begin{aligned} OA^2 &= OB^2 + AB^2 \\ \overrightarrow{OA} &= \sqrt{OB^2 + AB^2} \\ \| \overrightarrow{OA} \| &= \sqrt{4+9} = \underline{\sqrt{13}} \end{aligned}$$

$$\hat{u} = \left(\frac{2}{\| \overrightarrow{OA} \|}, \frac{3}{\| \overrightarrow{OA} \|} \right)$$



$$(2, 3) \rightarrow \underline{\underline{2\hat{i} + 3\hat{j}}}$$

in x direction = 2 units
 in y direction = 3 units

$$\begin{matrix} \sqrt{13} \\ 2 \end{matrix} \rightarrow \begin{pmatrix} \hat{i} \\ \hat{j} \end{pmatrix}$$

$$\sqrt{\left(\frac{2}{\sqrt{13}}\right)^2 + \left(\frac{3}{\sqrt{13}}\right)^2} = \sqrt{\frac{4}{13} + \frac{9}{13}} = \sqrt{\frac{13}{13}} = 1$$

Selecting the right feature

1000 - feature \rightarrow Top 10 feature

- ① Filter method
- ② Embedded method
- ③ Wrapper method.

① Filter method

\rightarrow Correlation

$$\text{Correlation } (f_1, y) \quad (f_2, y)$$

\rightarrow Multicollinearity (VIF)

f_1	f_2	Price of house (y)
# of houses	# of parking	

$$\left\{ \begin{array}{l} x_1 = -x_2 \\ 2x_1 + 3x_2 \\ 3x_1 + 7x_2 \\ 10x_1 \\ 10x_{x_2} \end{array} \right\}$$

$\Rightarrow x_1 \approx (x_2 x_3 x_4)$ Correlation among feature themselves

x_1	x_2	x_3	x_4	y
-------	-------	-------	-------	-----

Say $x_1 = x_2$

$$0x_1 + 3x_2 + 5x_3 + 4x_4 = y$$

② Wrapper method \rightarrow Recursive feature elimination.

Forward Select Backward Selection

③ Embedded method

\rightarrow Lasso Regression, Elastic Net.

(L1)

Decision tree \rightarrow MyBoost, LightGBM

Date Encoding

ML Algorithms

Numerical values

- ① Nominal / OHE
- ② Label and ordinal
- ③ Target guided ordinal encoding.

- ① Nominal / OHE | Dummy Variable creation.

→ Categorical data to numerical data

→ No order in the data.

Status	dummy style		den. man	numerical	Single		married		in a reln.	
	1	0			1	0	1	0	1	0
Single	1	0	0	0	Single	1	0	1	0	0
Married	0	1	1	1	Married	0	1	0	1	0
In relationship	0	0	0	0	In a reln	0	0	0	0	1
Single	1	0	0	0						
Married	0	1	1	1						
Single in a reln.	0	0	1	1						
=	0	0								

DS DE BA DA

OS | 0 0 D

DE | 0 0 0

BA | 0 0 0 0

DA | 0 0 0 0 0 1

No of Dummy Variables Create = No of Unique Categories in the columns

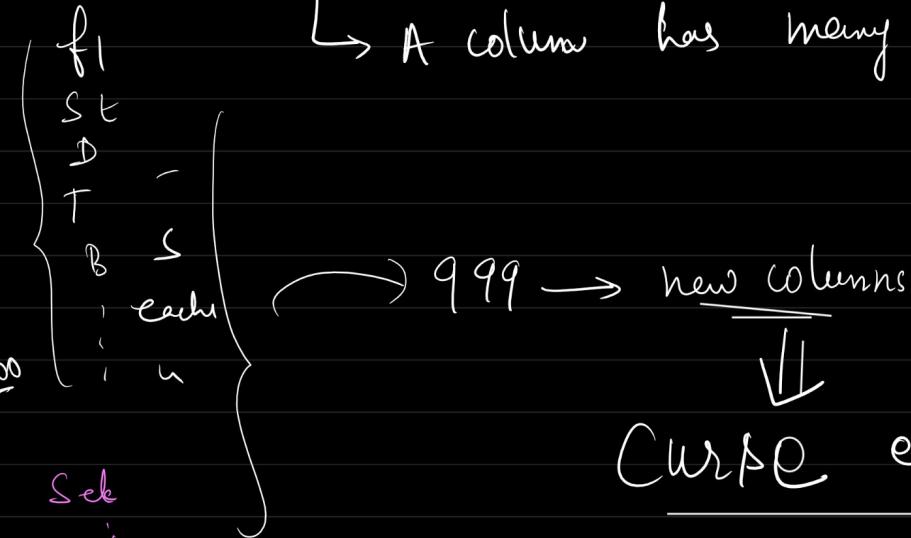
S M R

→ S	1	0	0
→ M	0	1	0
✓ R	0	0	1

For n unique groups = $n-1$ variable
 will be enough to explain the original feature

* disadvantage

→ A column has many categories



* Sol'n ⇒ seek

of it top 5 or top 10 unique categories

in the feature and make that no

→ dummy variables.

Value-counts()

Dum DS	Dum DE	Dum DA	DS - 102
1	0	1	DE - 90
0	1	0	DA - 10
0	0	0	
0	0	0	
0	0	0	
0	0	0	

② Label and Ordinal encoding

* Label encoding — assign numerical label to each category.

Rcd - 1
Exo - 2
Velo - 3

DS - 1
DE - 2
DA - 3

1
2
3
1
2
3
1

Adv → * No problem of curse of dimensional fp.

disadvantage :- ML can assume that there is some inherent order due to numerical ranking.

* Ordinal encoding

High school	- 1
Colleg	- 2
Post Graduate	- 3
Pnd	- 4

Grade	A	- 4
	B	- 3
	C	- 2
	D	- 1

(3) Target guided ordinal encoding (also useful)

for Nominal data with higher no of categories

→ based on their relationship with target Variable

→ When we have large no of Unique Categories in Categorical value.

→ Categorical grouping with mean/median of corresponding target variable

<u>Ex.</u>	<u>tips</u>	<u>time</u>	<u>total_bill</u>
	105	lunch	150
	105	lunch	120
	105	lnd	100
	150	dinn	-
	150	brr	-
	90	de	-
	90	b	-

groupby('time') ['total_bill').mean()

lunch - 105

dinner = 150

brr = 90

