

Topics for the day

- ① Numpy and Pandas Intro
- ② Data preprocessing techniques

Numpy → Numerical Python
deal with arrays. (Optimized - memory consumption, CPU speed) utilization

- ① One dimensional array

[1, 2, 3, 4]

- ② Two dimensional array

$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$

length	width	no bedrooms	city	price

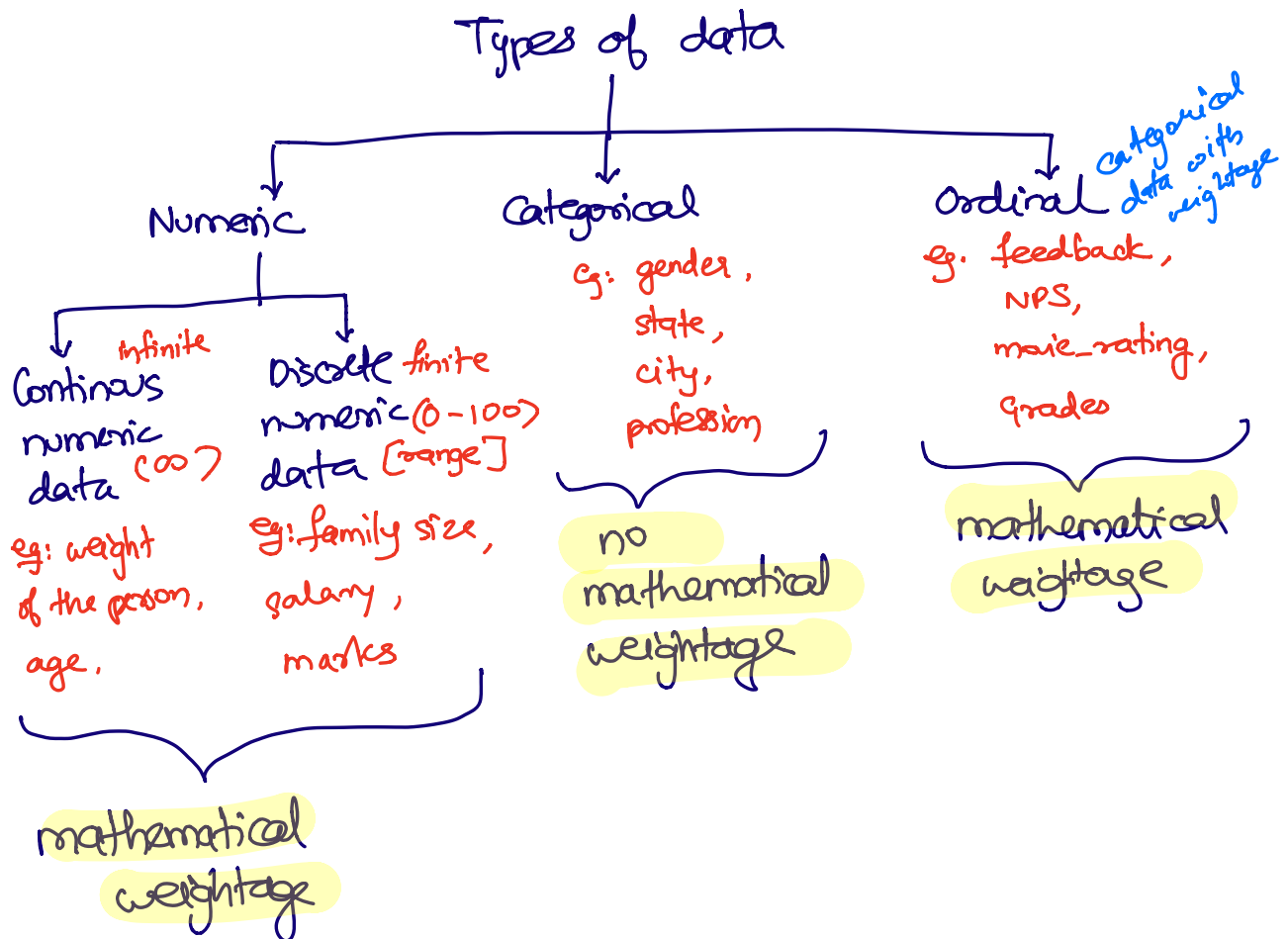
area	no bedrooms	city	price

} pandas

Machine Learning.

① File must hold logically structured data.

② All your data must be strictly **numeric**.



Machine learning program implementation steps

① Loading the data

② Preprocessing activities & EDA

③ General statistical summary.

④ Separate data as features and label. (ensure data is in the form of numpy array)

⑤ Deal with missing data.

→ Statistics as a base

① Replace NaN with mean (continuous numeric data)

② Replace NaN with median (discrete numeric data)

③ Replace NaN with Mode (categorical data)

→ Domain knowledge as base

① Statistics Technique (mean, median, mode)

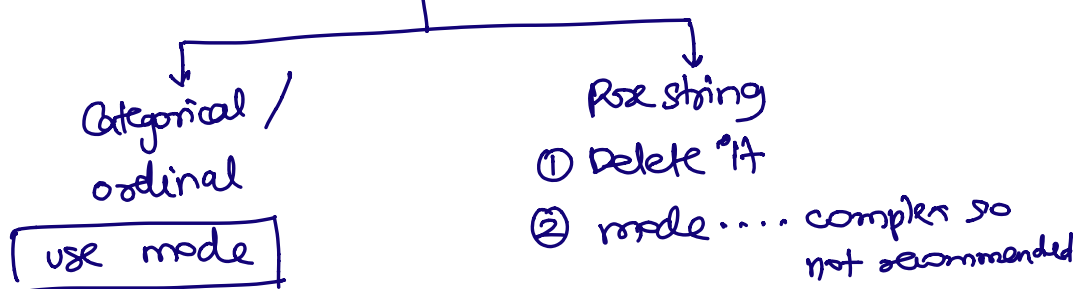
② Decide the default value for replacement.

if column data is string, ideally it is recommended to delete that record.

if label has missing data, delete that specific record

GUIDELINES

String Data (Handling missing data)



② Dealing with Categorical Data (creation of dummy variables)

eg:

eid	esal	ecity
1	1000	mumbai ✓
2	2000	chennai ✓
3	3000	mumbai

↑ mathematical ↑ m ↑ m

① get the unique values of categorical column.

['mumbai', 'chennai']

⇒ ② Sort the list in ascending order

['chennai', 'mumbai']

③ Replace the values of the column with the index values of the list created in step 2.

eid	esal	ecity
1	1000	0
2	2000	1
3	3000	0

↑ m ↑ m ↑ m

LABEL ENCODING

dummy variables

0	1	eid	esal
0	1	1	1000
1	0	2	2000
0	1	3	3000

④ Convert the column values into column itself such that the columns are losing mathematical weightage.

ONE-HOT ENCODING

$[m, c, h, b] \rightarrow \begin{matrix} b, c, h, m \\ 0, 1, 2, 3 \end{matrix}$

0	1	2	3	eid	esal
0	0	1	0	10	1000

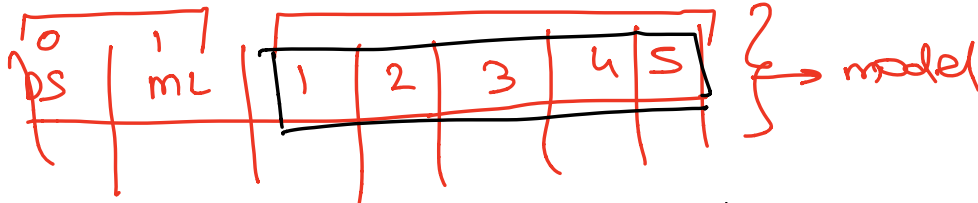
^X Trainer Name	(m-DS) Subject	(1-5) ^{discrete} Rating	Next class
Prashant	ML	4	y
Anu	DS	2	n
.	.	.	.

pre string categorical ordinal
 LE OHE
 one discrete

features → TN, S, R

Label → Next class

follow statistics - ✓
domain - x



Subject	^{ordinal} Rating
DS	Excellent
ML	Good
DS	Bad
ML	OK
.	Better
.	Excellent

Bad - 1
 OK - 2
 Good - 3
 Better - 4
 Excellent - 5

replace string with

Domain

||
 V
OHE