

Data Handling

⇒ Missing Values → > 50% → ignore (drop)
column (numerical column)
 ↙ ↘
not skewed skewed → median
 ↓
 mean

categorical column → mode

⇒ categorical column

exg

Gender

Male

Female

① if it is nominal category → dummy variables

a) → `get_dummies()` → pandas

b) → Onehot encoding → `sklearn.preprocessing`

② if it is an ordinal category → label encoding

a) for input \rightarrow Ordinal Encoder

b) for o/p \rightarrow label encoder.

\Rightarrow Feature Scaling

1	2	3
0.1	500	10000
0.2	600	30000
0.3	700	40000

to deal
with this

Scaling

Standardization
(Z-score)
($\mu=0, \sigma=1$)

Normalization
(0-1)

Sklearn preprocessing

StandardScaler
(Standardization)

MinMaxScaler
(Normalization)

Bollywood-dataset \rightarrow Assignment (i) Perform feature scaling
Numerical columns

\Rightarrow Read
docs & blogs

(ii) Make all columns numerical
in nature.

\Rightarrow Class imbalance

\rightarrow resampling \rightarrow Sklearn utils

CGPA	IQ	Placed
10	125	1
9	120	1
8	115	1
7	125	1
6	110	0
6	110	0
6	110	0
6	110	0

