

Handling Missing Values

1) Drop missing Values

2) Impute missing Value

Numerical

- ① Mean
- ② Median (Outliers)
- ③ Mode
- ④ Constant

Categorical

- ① Mode
- ② New Category

Time-series → B-Fill, F-Fill, Simple AVG, WAP

5 10 15 20 1000

mean is 500
median is 15

* Python for Imputing:

① Fillna()

② Simple Imputer

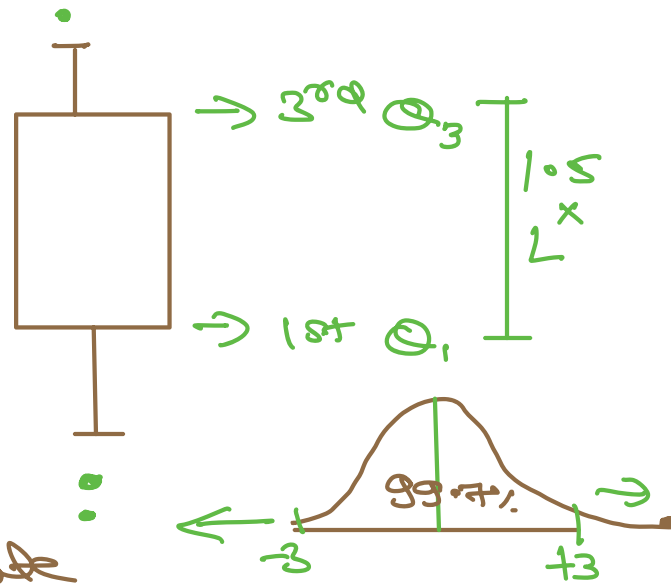
Outlier Treatment

Visualization

Uni-variate

① Box Plot

② Histogram



Bivariate / Multi-Variate

① Scatterplot

② DSSCAN

③ Kmeans

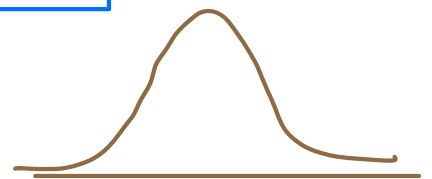
Methods to Find Outliers

① IQR Based

$$IQR = Q_3 - Q_1$$

② Z-score Based

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



Categorical to Numerical

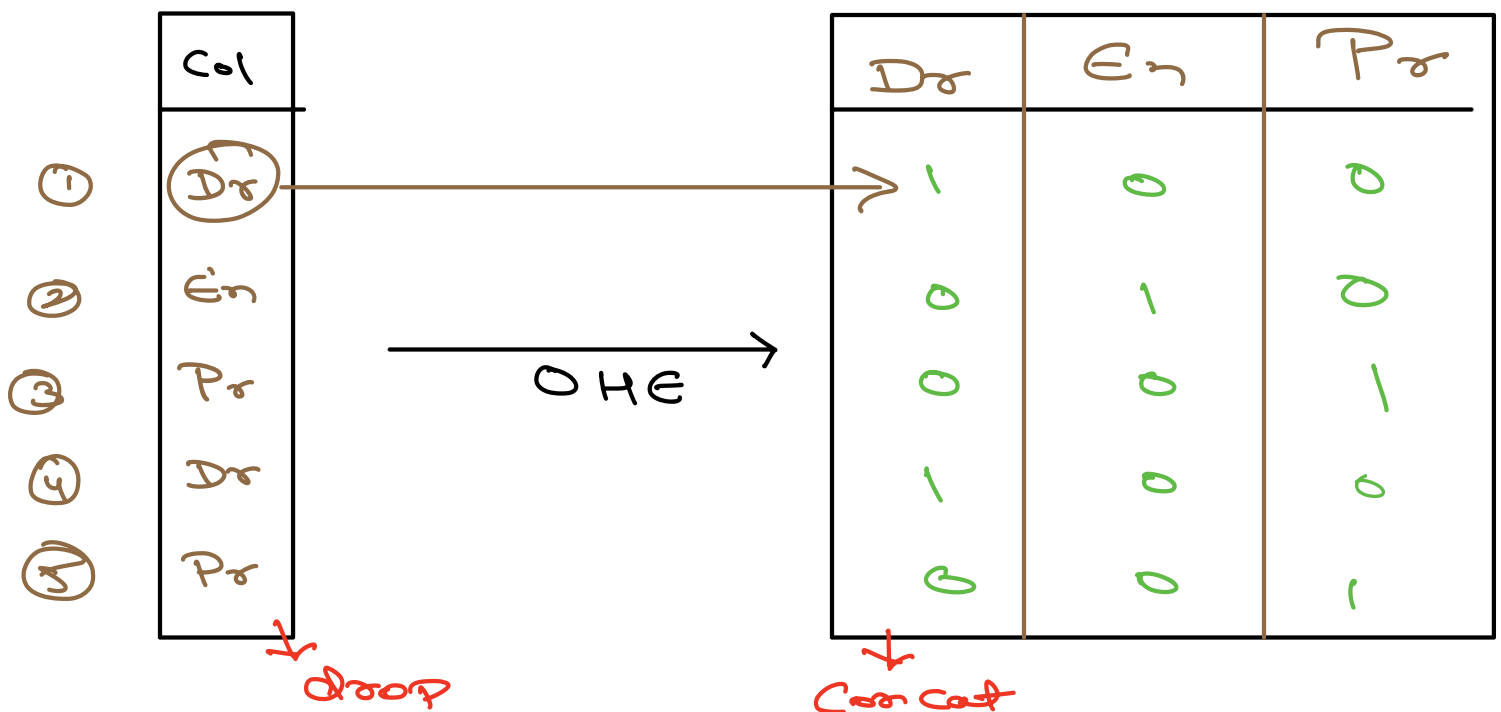
* ML Models process Data only in Numerical format.

Cat_str \rightarrow Cat_num

Methods to Convert

- ① One - Hot Encoding
- ② Label Encoding
- ③ Target Encoding

One Hot Encoding



Issue: Many Sparse columns are added leading to 'Curse of Dimensionality'

Label Encoder

Col
Dr
En
Pr
Dr
Pr

Label Encoder

map $\{$ Dr:0
En:1
Pr:2 $\}$

Col
0
1
2
0
2

Issue: Nominal Categorical Data \rightarrow Ordinal Categorical Data

Target Encoding mean Encoding

Col
Dr
En
Pr
Dr
Pr

Target Encoder

New-Col
Avg(Dr-T)
Avg(En-T)
Avg(Pr-T)
Avg(Dr-T)
Avg(Pr-T)

Scaling the Data

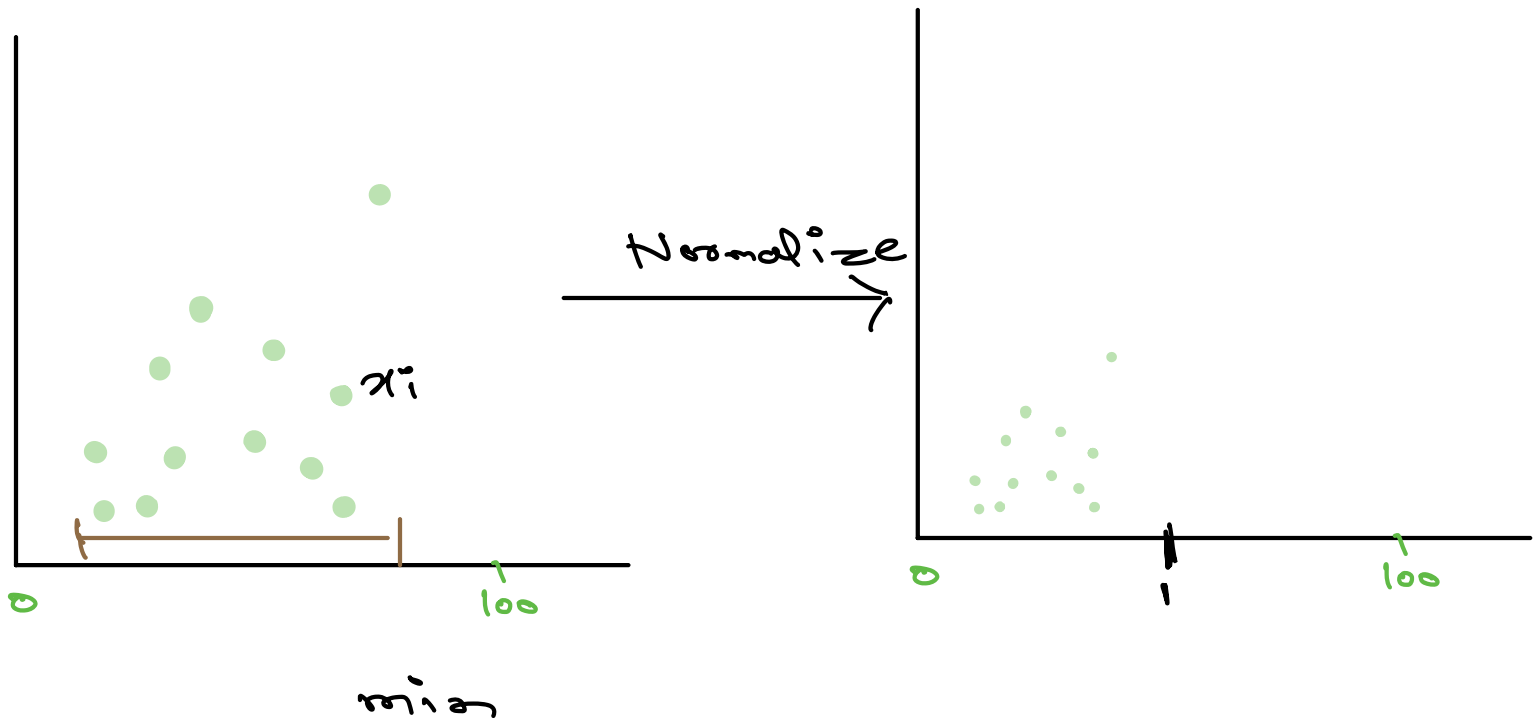
① Feature 1 \rightarrow (10, 100)

② Feature 2 \rightarrow (0, 10000)



Bring all Feature to similar Scale
(Scaling)

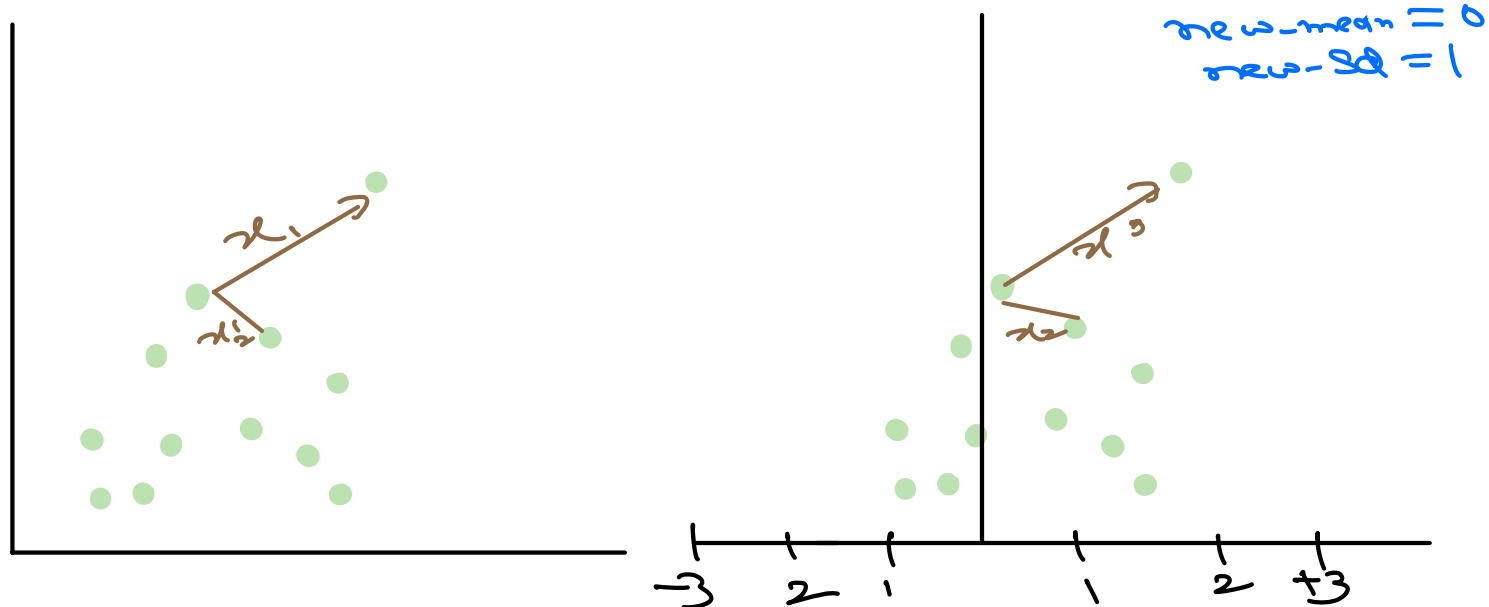
① Normalization \rightarrow (0, 1)



$$x_{i_new} \Rightarrow \frac{x_i - \min}{\max - \min}$$

* Normalize is preferred
Non-Gaussian Variable

② Standardization



$$x_{i_new} = \frac{x_i - \bar{x}}{\text{std}(x)}$$