# Data Preprocessing

Data preprocessing is used to transform raw data into a useful and efficient format.

## Data Preprocessing

- Data cleaning
  - missing data
  - noisy data
  - outlier
- Data Transformation
  - standardization
  - normalization
  - discretization
  - encoding categorical data
- Data Reduction
  - dimensional reduction

## ☐ Data Cleaning :-

### Mining data :

Missing values

- remove CCA (Complete Case Analysis)
  - (simple imputer) →
    - numeric
      - mean/median
      - random
      - arbitrary
    - categorical
      - mode
      - mining
- impute (fill)
  - univariate
  - multivariate
    - KNN imputer
    - iterative imputer

## Outlier :

An outlier is a data point that is noticeably different from the rest. They represent errors in measurement, bad data collection, etc.
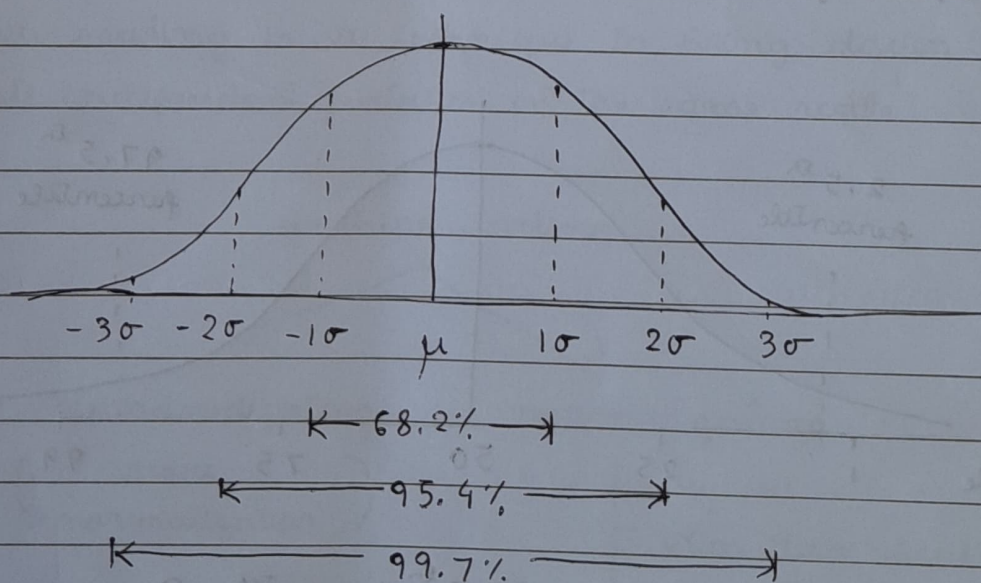
Outliers have major impact on these algorithms —

- Linear Regression
- Logistic  "
- Adaboost
- Deep Learning

Techniques for outlier detection and/ retmoval —

i) Z-score treatment —

(Applied for normally distributed data)



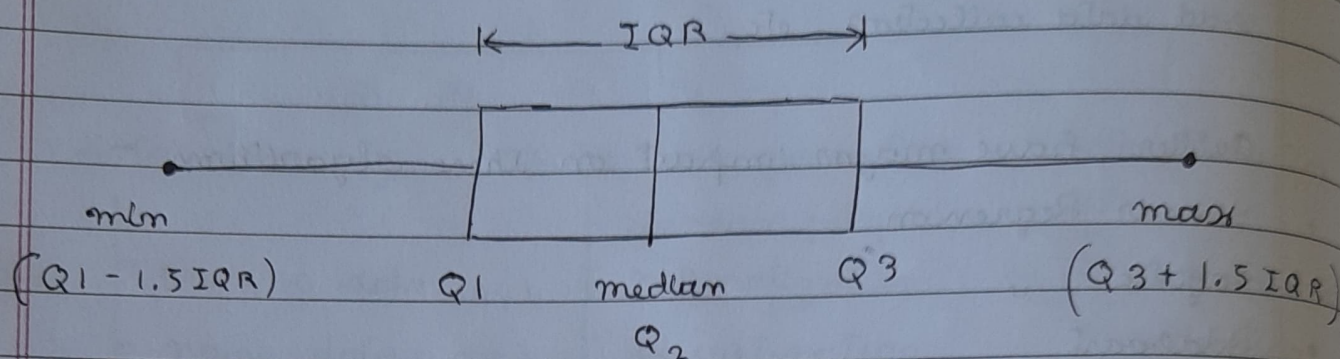$\leftarrow$ 68.2% $\rightarrow$

$\leftarrow$ 95.4% $\rightarrow$

$\leftarrow$ 99.7% $\rightarrow$

data points that are $> (\mu + 3\sigma)$  
$< (\mu - 3\sigma)$ } are considered outliers

ii) IQR - based filtering -

(applied for skewed data)

$$\leftarrow \text{------} IQR \text{------} \rightarrow$$



min
$(Q1 - 1.5 IQR)$        Q1        median        Q3        $(Q3 + 1.5 IQR)$        max
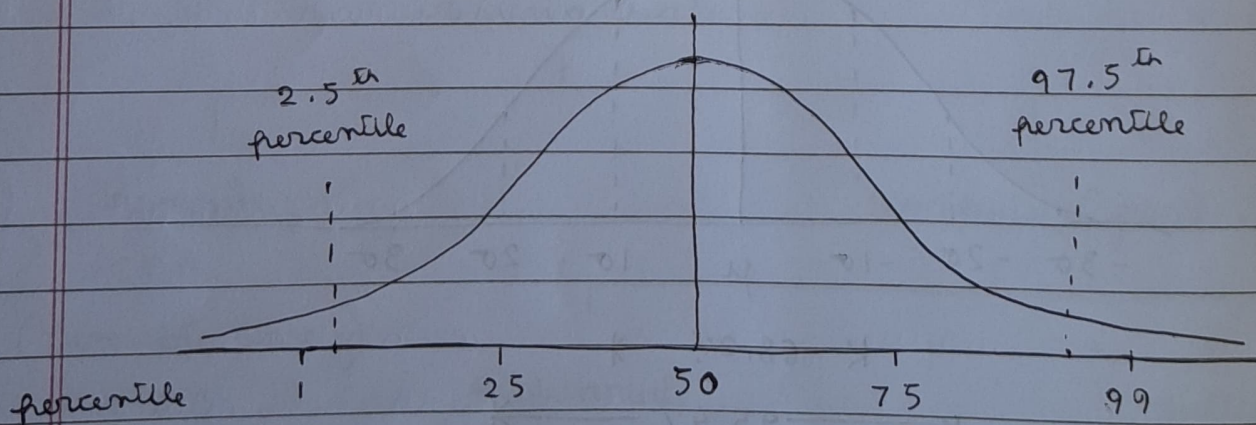                                              $Q_2$

data points that are $> (Q3 + 1.5 IQR)$ } are considered outliers
                      $< (Q1 - 1.5 IQR)$

iii) Percentile method -

(applied for other distributions)



$2.5^{th}$ percentile                          $97.5^{th}$ percentile

percentile        25        50        75        99

data points that lie above $97.5^{th}$ percentile } are considered outliers
                     below $2.5^{th}$ percentile

In this method, we can decide any percentile threshold value.

Techniques for outlier treatment -

i) Trimming - Here the outlier rows are removed from the dataset

ii) Capping - The outliers below the lower limit are replaced with the lower limit value and the outliers above the upper limit are replaced with upper limit value.

# Data Transformation :-

## Feature scaling :

Feature scaling is a technique to bring down the values of all independent features on the same scale.

```
                    feature scaling
                  /                \
        standardization          normalization
        ( z-score      )            ├ Min Max Scaling
        ( normalization )           ├ Mean normalization
                                    ├ Max Abs scaling
                                    ├ Robust scaling
```

i) Standardization -

Standardization makes the value of features have mean = 0 and standard deviation = 1

$$x_i' = \frac{x_i - \bar{x}}{\sigma}$$

$x_i'$ → transformed value

$\bar{x}$ → mean

$\sigma$ → standard deviation

Standardization is done for the following algorithms -

- K-means
- KNN
- Principal Component Analysis (PCA)
- Logistic Regression
- Gradient descent
- Artificial Neural Network

ii) Normalization -

In Normalization, the value of features are rescaled between the range 0 to 1.

- Min - Max scaling

$$X_i' = \frac{X_i - X_{min}}{X_{max} - X_{min}}$$

- Mean normalization
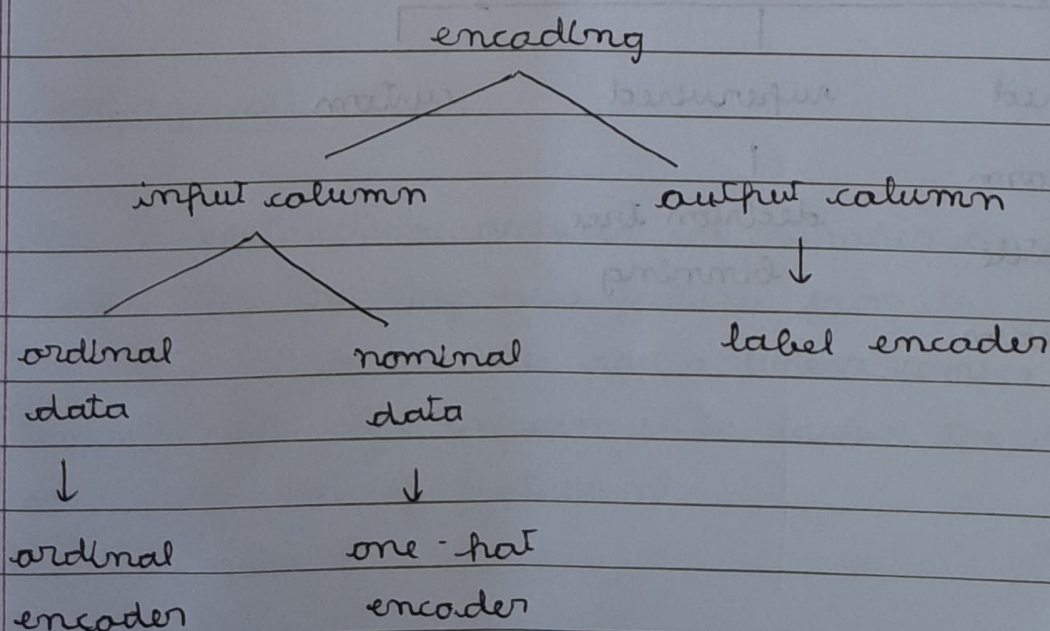
$$X_i' = \frac{X_i - \bar{X}}{X_{max} - X_{min}}$$

- Max - Absolute scaling (used in case of sparse data)

$$X_i' = \frac{X_i}{|X_{max}|}$$

- Robust scaling (used when data has outliers)

$$X_i' = \frac{X_i - X_{median}}{IQR}$$

Encoding categorical data:



encoding

input column                    output column

ordinal        nominal                  label encoder
data           data

↓              ↓

ordinal        one-hot
encoder        encoder

Encoding categorical data is the process of converting categorical data into integers so that it can be fed to the ML model.
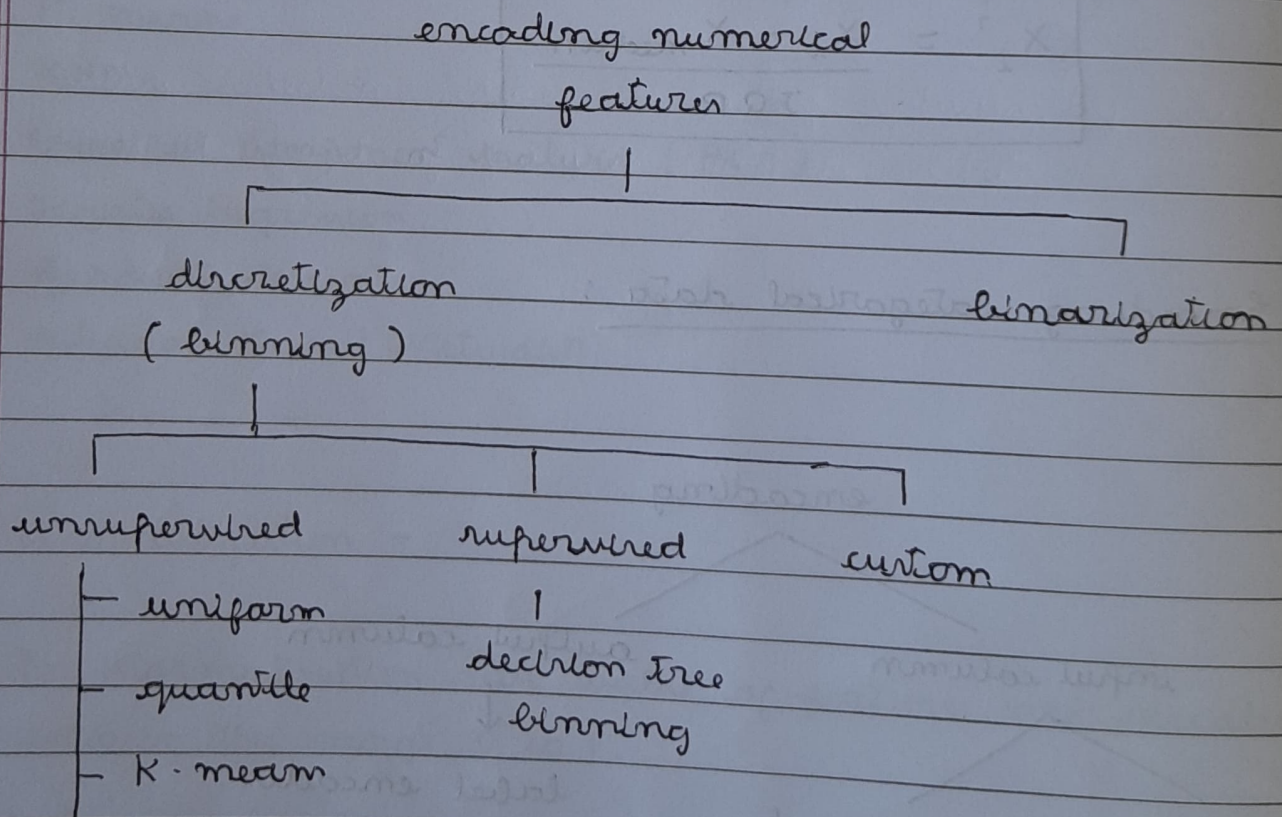
Categorical data is of 2 types -
ordinal - There is a specific order
nominal - " " no " "

## Encoding numerical features :

Encoding numerical features is the process of converting continuous numerical column to categorical columns.

encoding numerical
features

discretization                                    binarization
( binning )

unsupervised          supervised          custom
 ├ uniform              |
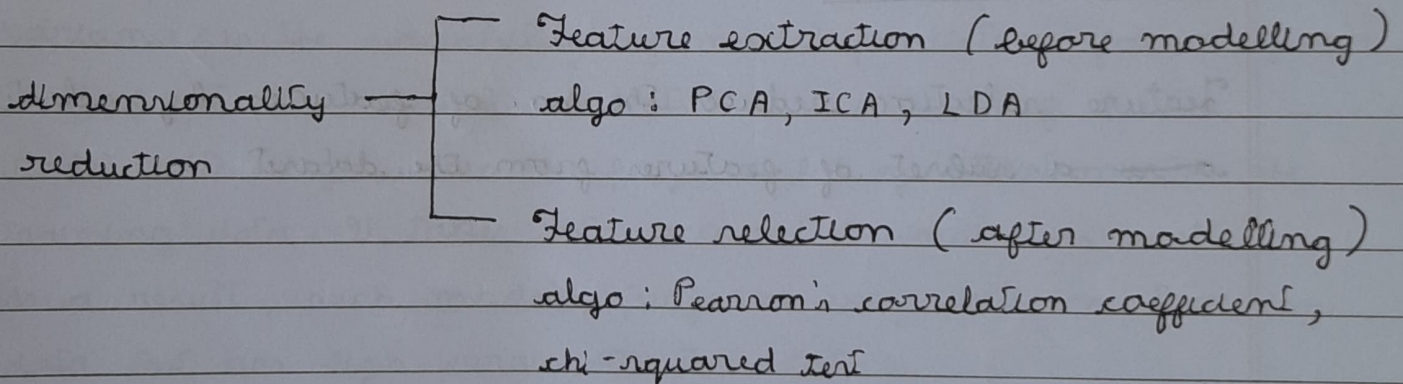 ├ quantile          decision tree
 └ K-mean               binning

i) **discretization** - In the process of converting continuous variables into categorical variables by creating a set of contiguous intervals that span the range of the variable values.

ii) **Binarization** - In the process of converting continuous variables into binary numbers

## Data (dimensionality) reduction :-

```
                         ┌─  Feature extraction (before modelling)
                         │    algo : PCA, ICA, LDA
dimensionality    ───────┤
reduction                │
                         └─  Feature selection (after modelling)
                              algo : Pearson's correlation coefficient,
                              chi-squared test
```

## Feature extraction :

Feature extraction reduces the number of features in a dataset by creating new features from the existing ones. The new set of features are a linear combination of the original features. The aim is to capture the data pattern with fewer no. of features.

| LDA (Linear discriminant analysis) | PCA (Principal Component Analysis) |
|---|---|
| • used for supervised models. | • used for unsupervised models |
| • describes direction of maximum separability. | • describes direction of maximum variance |
| • requires class label info to fit () | • doesn't require class label info to fit (). |

## Feature selection :

Feature selection reduces the no. of features by selecting a subset of features from the dataset