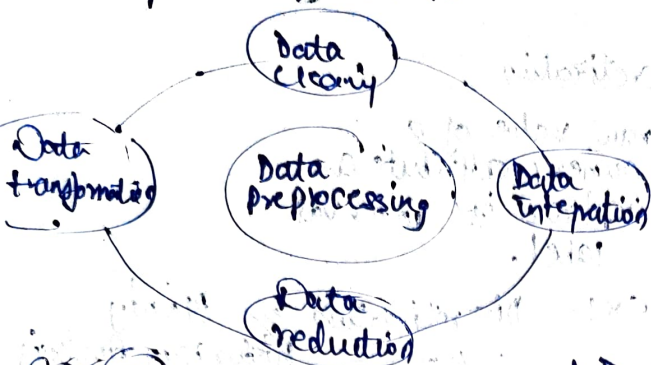


## → Data Pre-processing

↳ is a data mining technique which converts raw data into understandable format. Real-world data is often incomplete, inconsistent, or lacking in certain trends and contains error. Data pre-processing is proven method to resolve such issues.

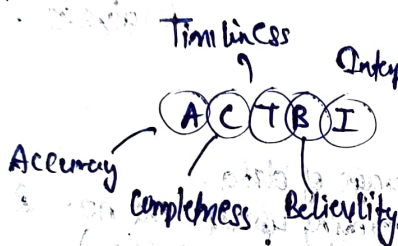


→ Data processing have 4 major task.

- Data cleaning
- Data integration
- Data Reduction
- Data transformation

CRIT  
cleaning integrn

→ Data pre processing is important to achieve accuracy (whether data entered is correct or not).



- Completeness (—//— data is available or not recorded)
- Timeliness (data should be updated correctly)
- Believability (—//— trustworthy)
- Interpretability (understandability of data)

## → Data Cleaning

- ↳ is the process to remove incorrect data, inaccurate data from incomplete data dataset.
- ↳ it also replaces the missing values.
- ↳ it also smooth out the noise while identifying outliers.
- ↳ correct the inconsistency.

## → Handling Missing values:

- Ignore the row (tuple)
- used when class label or primary key is missing
- not very effective unless tuple contain several attribute with missing value.
- Not available or 'na' can be used to replace the missing values.
- Filling manually.
  - time consuming
  - not feasible
- Use global constant
  - replace all missing attribute by same constant like 'Unknown' or '00'.
  - Simple but not foolproof.

- use a measure of central tendency
  - use of mean or median to fill the missing value
  - for normal symmetric data distribution mean is used.
  - for skewed  $\rightarrow$  median is used.
- use the most probable value to fill in the missy value
  - determined with regression, inference base tools using Bayesian formalism or decision tree.

→ Noisy Data, means random error or containing unnecessary data points.

① Binning → This method is to smooth or handle noisy data. First data is sorted and then sorted value are separated and stored in the form of bins.

→ 3 method of smoothing data.

Smoothing by bin mean method

→ value in bin replaced by mean.

Smoothing by bin median method

→ value in bin replaced by median.

Smoothing by bin boundary

→ max & min value are taken at each bin & replaced by closest boundary value.

② Regression

→ used to smooth data & handle data when unnecessary data is present.

→ for analysis purpose, regression helps to decide the variable.

③ Clustering

→ used for finding the outliers & also in grouping the data.

→ used in unsupervised learning.

→ Data Reduction

→ helps to reduction of volume of the data which make analysis easier and produces almost the same result.

→ helps to reduce storage space.

→ 3 techniques.

- Dimensionality red<sup>n</sup>
- Numerosity red<sup>n</sup>
- Data comp<sup>n</sup>.





## Smoothing

- used to remove noise from data
- technique include binning, regression, & clustering.

②

## Attribute Construction

- also known as feature construction
- where new attrib are constructed & added from given set of attribute.

## Normalization

- attribute data are scaled so as to fall with smaller range such as -1.0 - 1.0

## Aggregation

- where summary or aggregation oprn is applied
- used in constructing data cube for analysis of data at multiple level of abstraction.

## Discretization

- raw value of a numeric attribute are replaced by interval label.

ex:- Discretization by Binning  
→ histogram Analysis  
→ by cluster, Decision Tree and Correlation Analysis.

$$\text{Mean} = \frac{X_1 + X_2 + X_3 + \dots + X_n}{n}$$

Average of data.

(Summing up all the no. & then dividing by no. of obs.)

$$= \frac{10 + 20 + 30 + 40 + 50}{5} = 30 \text{ Mean.}$$

ex: 10, 20, 30, 40, 50

## Median

$$\frac{n}{2} \text{ (Even)}$$

$$\frac{n+1}{2} \text{ (Odd)}$$

∴ n - no. of obs.

10, 20, 30, 40, 50

$$\frac{n+1}{2} = \frac{5+1}{2} = 3 \Rightarrow \text{Median} = 30$$

10, 20, 30, 40, 50, 60

$$\frac{n}{2} = 3 \rightarrow 30, 40 \rightarrow \frac{30+40}{2}$$

$$\text{Median} = 35$$

## Mode

1, 3, 4, 6, 7, 3, 3, 5, 10, 3

$$\text{Mode} = 3$$

{ highest frequency }

## IQR $Q_3 - Q_1$

(second half)

(first half)

ex: 48, 52, 57, 61, 64, 72, 76, 77, 81, 85, 88

48, 52, 57, 61, 64, 72

$$Q_1 = \text{median of 1st half} = \frac{57+61}{2}$$

72, 76, 77, 81, 85, 88

$$Q_3 = \text{median of second half} = \frac{77+81}{2}$$

Range = Maximum value - Minimum value  
 82, 41, 28, 54, 35, 26, 23, 33, 38, 40  
 $54 - 23 = \textcircled{31} \text{ Range}$

Standard deviation  
 ↳ deviation from mean value.

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{n}}$$

$n$  :- no. of popul?  
 $\mu = \text{mean} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$

Variance

∴ squared deviation

$$\sigma^2 = \frac{\sum (x - \mu)^2}{n}$$

Popl<sup>n</sup> variance

∴  $\mu = \text{mean}$

$$\sigma^2 = \frac{\sum (x - \bar{x})^2}{n}$$

Sample varian

$\bar{x} = \text{average}$ .

→ Box Plot

- standardized way of displaying the distribution of data based on five no. of summary ("minimum", first quartile ( $Q_1$ ), median, third quartile ( $Q_3$ ), "maximum")
- It can tell you about your outlier & where the value.
- also tell about your data is symmetrical, skewed, grouped.

