

# Data Preprocessing - its a technique that Involves transforming raw data into a useful and efficient format.

↳ Need -

Real world data are generally -

→ Incomplete - missing certain attributes of importance

→ Noisy - containing errors or outliers

→ Inconsistent - containing discrepancies in codes or names

which may lead to false analysis, so preprocessing is required.

Data preprocessing provides - accuracy - completeness

- consistency - timeliness & - Interpretability to data.

# Major tasks in Data Preprocessing -

1. Data Cleaning - filling in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies.
2. Data Integration - integration of multiple databases, data cubes & files.
3. Data Reduction - dimensionality red<sup>n</sup>, numerosity red<sup>n</sup>, & data compression
4. Data Transformation & Discretization - normalization & concept hierarchy generation.

# DATA CLEANING -

Data in real world is dirty, there is potential<sup>cor.</sup> of incorrect data, its

→ incomplete : lacking of attribute or lacking of value in attribute.

◦ It may be due to -

- equipment malfunction
- inconsistent with other recorded data
- certain data may not be considered imp at time of entry.
- data not entered due to misunderstanding.

◦ how to handle such data -

- replace with global const
- fill in it with attribute mean.
- fill with attribute mean for all samples belonging to same class.
- fill it with most probable value using decision tree.

→ Noisy - random error or variance in a measured variable

• it may be due to -

- data entry problems
- faulty data collection instruments
- data transmission problems
- technology limitation

• how to handle noisy data -

- binning: sort the data and partition into bins & smoothen it
- regression: smooth by fitting data into regression fn.
- clustering: detect and remove outliers
- combined computer & human inspection - detect suspicious values and check by human.

#

## DATA INTEGRATION -

Combines data from multiple sources into a coherent store. It includes detecting and resolving data value conflicts that may arise due to different scales or representations.

↳ Problem - Data Redundancy i.e.

- same object or attribute may have diff names in diff db.
- one attribute may be a derived attribute in another table

↳ Solution - Correlation & Covariance analysis

• Correlation -

1) Numerical Data -

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)\sigma_x\sigma_y} = \frac{\sum_{i=1}^n (x_i y_i) - n\bar{x}\bar{y}}{(n-1)\sigma_x\sigma_y}$$

where  $n$  = no. of tuples

$\bar{x}, \bar{y}$  = respective means

$\sigma_x, \sigma_y$  = respective standard deviations

NOTE - if  $r > 0 \Rightarrow (+)$ ve relation

$r < 0 \Rightarrow (-)$ ve relation

$r = 0 \Rightarrow$  no relation (independent)

2) Nominal Data -  $\chi^2$  test.

larger the value  $\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$   
more likely they are related.





⇒ PCA (Principal Component Analysis) - find a projection that captures the largest amount of variation in data. the original data is projected onto a much smaller space, resulting in dimensionality red<sup>n</sup>.

2) Numerosity Reduction - reduce data volume by choosing alternative, smaller forms of data representation.

⇒ Parametric methods - [Regression]

- Linear - data modeled to fit a straight line
- multiple - allows a response variable  $Y$  to be modeled as a linear f<sup>n</sup> of multi-dimensional feature vector
- log-linear - approximates discrete multidimensional probability distributions.

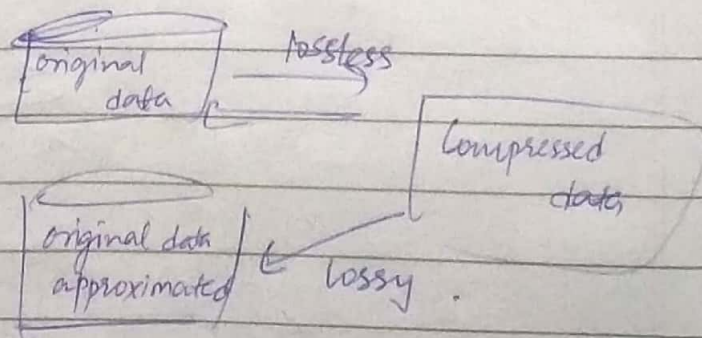
⇒ Non-Parametric -

- Histogram Analysis - divide data into buckets & store the sum/average for each bucket.
- Clustering - partition data set into clusters based on similarity, and store cluster representation
- Sampling - obtaining a small sample  $S$  to represent the whole data set  $N$ .

↳ Types -

1. Simple random sampling
2. Sampling without replacement
3. Sampling with replacement
4. Stratified sampling - partition dataset, and draw samples from each data proportionally.

3) Data Compression -





# DATA TRANSFORMATION - A function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new ones.

⇒ Methods -

1. Smoothing - [Binning] [partitioning the data]

- equi-depth/frequency: each bin has same no. of samples
- equi-width -  $W = (\text{highest value} - \text{lowest value}) / \text{no. of intervals}$
- bin means - fill value with means of value

2. Aggregation - [cube lattice]

3. Normalization -

- min max -

$$V' = \frac{V - \min_i}{\max_i - \min_i} (\text{new max} - \text{new min}) + \text{new min}$$

- Z score -

$$V' = \frac{V - \mu}{\sigma}$$

$\mu \rightarrow \text{mean}$   
 $\sigma \rightarrow \text{SD}$

- Decimal scaling -

$$V' = \frac{V}{10^j}$$

where  $j$  is smallest integer such that  $\max(|V|) < 1$

4. Discretization -

- Correlation

- Covariance

- Concept Hierarchy Generation: organizes concepts

(ie. attribute values) hierarchically and is usually associated with each dimension in a data warehouse. It facilitates drilling & rolling in data warehouses to view data in multiple granularity.