

Q2

Data preprocessing is crucial in any data mining process as they directly impact success rate of the project.

This reduces complexity of the data under analysis as data in real world is unclear.

• Method involved in data processing

### (1) Data cleaning

The data can have many irrelevant and missing parts. To handle this part, data cleaning is done.

It involves handling of missing data, noisy data etc.

#### (a) Missing data

The situation arises in case



of missing data. It can be handled in various ways:-

Some of them are:

- (i) Ignore the tuples
- (ii) Fill the missing values.

(b)

### Noisy data

Noisy data is a meaningless data that can't be interpreted by machines. It can be generated due to faulty data collection, data entry error etc.

• It can be handled in following ways:

(i) Binning method:

The method works on sorted data in order to smooth it.

The data is divided into segments of equal size and then various methods are performed.



(ii) Regression:

data can be made smooth by fitting it to a regression function. The regression used may be linear, or multiple.

(iii) Clustering

The approach groups the similar data in a cluster.

The ~~outliers~~ outliers may be marked or it will fall outside the clusters.

(2) Data transformation

This step is taken in order to transform the data in appropriate forms suitable for mining process.

It involves the following ways:

(i) Normalization

It is done in order to scale the data values in a specified range ( $-1.0$  to  $1.0$  or  $0.0$  to  $1.0$ )



## (ii) Attribute ~~set~~ Selection:

In this strategy, new attributes are constructed from the given set of attributes to help the mining process.

## ✓ (iii) Discretization

This is done to replace the raw values of numeric attribute by interval levels or conceptual levels.

## (iv) Concept Hierarchy Generation

Here attributes are converted from lower level to higher level in hierarchy. For ex: The attribute "city" can be converted to "country".



### ③ Data Reduction

Since data mining is a technique that is used to ~~manage~~ handle huge amount of data.

The various steps to data reduction are:

#### (i) Data cube Aggregation:

Aggregation operation is applied to data for the construction of the data cube.

#### (ii) Attribute subset selection

The highly relevant attributes should be used; rest all can be discharged.

#### (iii) Dimensionality Reduction

This enable to store ~~the~~ the model of data instead of whole data, for ex: Regression models.



(iv)

Dimensionality

Reduction

Reduction

This reduce the size of  
data by encoding mechanism.  
It can be lossy or lossless.