# Assignment - 1

**Q.1)** What are the different types of data that can be mined?

**Sol^n)** The different types are:-

1) Flat files
2) Relational Databases
3) Data Warehouses
4) Transactional Databases
5) Multimedia Databases
6) Spatial Databases
7) Time series Databases
8) World Wide Web (WWW)

**Q.2)** Describe five number Summary in detail.

**Sol^n)** The five number is a Set of descriptive statistics that provides information about a dataset. It consists of the five most imp. sample percentiles.

i) The Sample minimum (smallest observation)
ii) The lower quartile or first quartile
iii) The median (the middle value)
iv) The upper quartile or third quartile
v) The Sample maximum (largest observation)
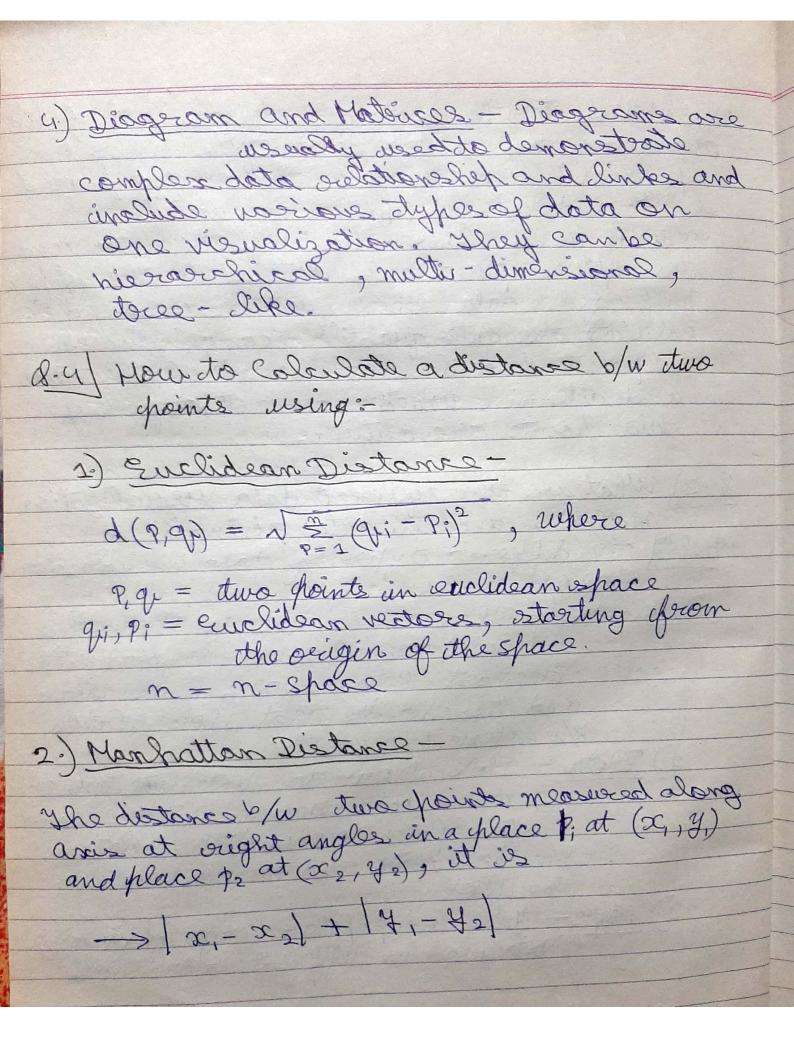
The addition to the median of a single set of data there are two related statistics called the upper and lower quartile. If data are placed in to order, then the lower quartile is central to the lower half of the data and the upper quartile is central to the upper half of the data. These qualities are used to calculate the interquartile range, which helps to describe the spread of the data and determine whether or not ~~any~~ any data points are outliers.

8.3) What are the different data visualization methods used in data mining?

Sol.) The different types of data visualization methods are :-

1) Charts - The easiest way to show the development of one or several sets in chart. Charts vary from bar and line charts that show relation b/w element over time.

2) Plots - Plots allow to distribute two or more data set over 2D or even 3D space to show the relation b/w these sets and parameters on the plot.

3) Maps - Maps are popular way to visualize data used in different industries.

4) Diagram and Matrices — Diagrams are usually used to demonstrate complex data relationship and links and include various types of data on one visualization. They can be hierarchical, multi-dimensional, tree-like.

Q.4] How to Calculate a distance b/w two points using :-

1) Euclidean Distance -

$$d(P, q) = \sqrt{\sum_{P=1}^{m} (q_i - P_i)^2}$$ , where

$P, q$ = two points in euclidean space
$q_i, P_i$ = euclidean vectors, starting from the origin of the space.
$n$ = $n$-space

2) Manhattan Distance —

The distance b/w two points measured along axis at right angles in a place $P_1$ at $(x_1, y_1)$ and place $P_2$ at $(x_2, y_2)$, it is

$$\rightarrow |x_1 - x_2| + |y_1 - y_2|$$

3) **Supremum Distance -**

Distance b/w two points as the maximum diff over any of their axis value. In a 2-D grid for instance, if we have two points $(x_1, y_1)$ and $(x_2, y_2)$

then, $Max(|y_2 - y_1|, |x_2 - x_1|)$

**Q.5** Write a Short note on Data Reduction and Transformation.

**sol:** It is a technique used to obtain a reduced representation of the data set that is much smaller in ~~the~~ volume than the original data but still contains all the critical information. The following methods are used for data reduction.
i) Data Cube Aggregation
ii) Dimensionality Reduction
iii) Data Compression
iv) Numerosity Reduction
v) Discretization and Concept hierarchy generation.

Data Transformation is a process in which data is transformed from one format to another format which is more suitable for data mining. Some data transformation strategies are :-

1) Smoothing
2) Aggregation
3) Generalisation
4) Normalisation
5) Attribute Construction