# Ch -2 Data

→ Data mining is the process of extracting extra knowledge 2 hidden patterns from large amount of data

→ An attribute is a property or characteristic of an object that may vary from one object to another or one time to another For eg. eye colour

→ A measurement scale is a rule (function) that associate a numerical or symbolic value with attribute of an object

• Different types of attributes :

① Nominal : That provide info to distinguish one object from the others (=, ≠) Eg. & zip code, employee ID, eye color, gender

② ordinal : They provide info to help order the data (<, >) Eg. {good, better, best}, grades, street no.

③ Interval : The difference b/w values are meaningful i.e a unit of measurement exists (+, -). Eg. calender dates, temp.

④ Ratio : Both difference & ratios are meaningful (*, /) Eg. temp. in Kelvin, counts, age, length

★ Nominal 2 ordinal are called categorical or qualitative Interval & ratio are called numerical or quantitative

→ Discrete attribute has a finite or countably infinit set of values. Eg. zip codes, ID. Binary attributes has only 2 values. Eg. male, female or true/false

→ Continuous values are real no's Eg. temp, height etc.

→ For asymmetric att, only the presence of non-zero value is imp.

• General characteristics of Data set :

① Dimensionality

② Sparsity

③ Resolution : For eg. variation in atmospheric pressure on a scale of hours reflect the movement of storms. On a scale of months it is not detectable

- Types of Dataset

| TId | Refund | Marital status | Taxable Inc. | DB |
|-----|--------|----------------|--------------|-----|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100k | No |
| 3 | No | Single | 90k | Yes |

• Record Data

| TID | ITEMS |
|-----|-------|
| 1 | Bread, Soda, Milk |
| 2 | Beer, Bread |
| 3 | Soda, Diaper, Milk |

• Transaction Data

| Xload | Yload | Dist. | Load | Thin |
|-------|-------|-------|------|------|
| 10.23 | 5.27 | 15.22 | 27 | 1.2 |
| 12.65 | 6.25 | 16.22 | 22 | 1.1 |
| 14.27 | 7.23 | 07.34 | 23 | 1.2 |

• Data Matrix

|       | team | coach | play | won | lost |
|-------|------|-------|------|-----|------|
| Doc 1 | 3 | 0 | 5 | 0 | 2 |
| Doc 2 | 0 | 7 | 0 | 0 | 3 |
| Doc 3 | 0 | 1 | 0 | 2 | 0 |

• Document term matrix

→ It is a type of sparse matrix which tells the word count of a doc

- Types of Ordered Data

① Sequential / Temporal Data is an extension of record data, where each record has a time associated with it.

② Sequence Data consists of data set like a sequence of indivisual entities (words / letter) Eg. genetic code

③ Time Series Data is a series of measurement taken over time. Eg. avg. monthly temp. of Delhi from 1982 to 2020

- temporal autorelation is when two measurements are close in time, their values are similar

④ ~~Data~~ spatial Data is collection of data over various geographical locations eg. temp. all over india

- spatial autocorrelation: close to location, similar value

→ Precision is the closeness of repeated measurements ~~from~~ of thea same quantity to one another. measured by std. deviation

→ A ~~systematic~~ systematic variation of measurements from the quantity being measured. measured by finding the difference b/w mean of measurements & actual quantity

→ Accuracy is the closeness of measurements to the true value of quantity being measured.

→ Outliers are either data objects that have different characters from the rest of data set OR values of attribute that are unusual w.r.t to typical values of that attribute.

→ ways to treat missing values:
① Eliminate data object / Attribute
② Estimate missing values
③ Ignore the missing values during analysis

→ Inconsistent values: Eg. given zip code & city does not match
→ Deduplication: process of dealing with & duplicates

## DATA PREPROCESSING

① Aggregation: Combining of two or more objects into single objects
→ Resulting data sets are smaller, which take less time & memory

→ It can act as change of scope or scale by providing high level view

→ The behaviour of group of objects is much more stable than indivisual objects

★ Disadvantage: Loss of interesting detail

② Sampling: selecting a subset of data to be analyzed.
→ It is too expensive or time consuming to ~~most~~ process all the data
→ A sample is representative if it has appx. the same properties as original data set.

• Simple random sampling: There is an equal probability of selecting any particular item. It can be with or without replacement

• If the data has different types of objects, then it is added to prespecified groups & equal no. of records from each group are selected - Stratified Sampling

• It is difficult to determine sample size, so adaptive/progressive sampling is used in which the size of sample is increased until sufficient : Stop when accuracy levels off

③ Dimensionality Reduction: lowering no. of attributes
→ Many DM algo work better with less dimensionality
→ Can eliminate irrelevant features & reduce noise
→ creates more undestandable model, easily visualized
→ Less time & memory reqd.

★ Curse of Dimensionality : Phenomenon that ~~DM algo~~ many types of Data analysis become significantly harder with increase in dimensionality.

④ Feature Subset Selection: Subset of features/att. are used
→ Redundant & Irrelevant features can reduce accuracy & quality

- Embedded Approaches: Feature selection occurs naturally as a part of DM algo. The algo itself decides which att to use or ignore
- Filter Approaches: Features are selected before DM algo is run
- Wrapper Approaches: The DM algo is used as black box to find best subsets

⑤ Feature Creation: To create a new set of att from original dataset that captures imp. info. This is also known as Feature Extraction

⑥ Discretization & Binarization: To transform a continuous att. into a categorical att is called discretization.
Transformation of conti. or discrete att. into one or more binary att. is binarization

- Binarization technique: If there are m categorical values, then uniquely assign each value to an integer $[0, m-1]$. If the att. is ordinal, order must be preserved. Convert into binary, $n = \lceil \log_2(m) \rceil$

- Problem of discretization is to determine how many split points to choose & where to place them.
→ Equal width divide into user specified no. of intervals with same width
→ Equal freq. tries to put same no. of objects

⑦ Variable Transformation: Tray. applied to values of variable
For eg. changing the magnitude of variable
→ Simple Functions: $x^k$, $\log x$, $e^x$, $\sqrt{x}$, $\frac{1}{x}$, $\sin x$ or $|x|$
→ Normalisation / Standardisation: Att values are normalized by scaling their values so that they fall in specified range

Eg. If our data is scattered, we can use z tray. so that our values lie b/w 0 & 1

$$Zvalue = \left| \frac{u - \bar{u}}{\sigma} \right| \quad \bar{u} = mean$$
$$\sigma = s.d$$

If we have outlier, use median instead of ~~mode~~ mean in calc. of s.d

- Dissimilarities b/w data objects
→ Euclidean distance: $\quad d(x,y) = \sqrt{\sum_{k=1}^{n}(x_k - y_k)^2}$

→ Minkowski distance! $\quad d(x,y) = \left( \sum_{k=1}^{n} |x_k - y_k| \right)^{1/r}$

| | |
|---|---|
| r = 1 | Hamming distance |
| r = 2 | Euclidean " |
| r = ∞ | Supremum " |

- Properties of Euclidean Distance —
① Positivity
   (a) $d(x,y) \geq 0$ for all x & y
   (b) $d(x,y) = 0$ if only x = y
② Symmetry
   $d(x,y) = d(y,x)$ for all x & y
③ Triangle Inequality
   $d(x,z) \leq d(x,y) + d(y,z) \quad \forall x y z$

→ Measures that satisfy all three properties are known as ~~metric~~ metric

- Similarities b/w data objects
1. $s(x,y) = 1$ only if x = 1 $\quad (0 \leq s \leq 1)$
2. $s(x,y) = s(y,x) \quad \forall x y$