

Data warehouse & mining

Data mining - Process of extracting information to identify patterns, trends and useful data, used by organisation to solve business prob.

Type of data -

- Relational data warehouse.
- Data repository
- Transactional data.

Application -

- Sales analysis
- Medical analysis
- Credit card fraud.
- Healthcare
- Education
- Finance

Advantage -

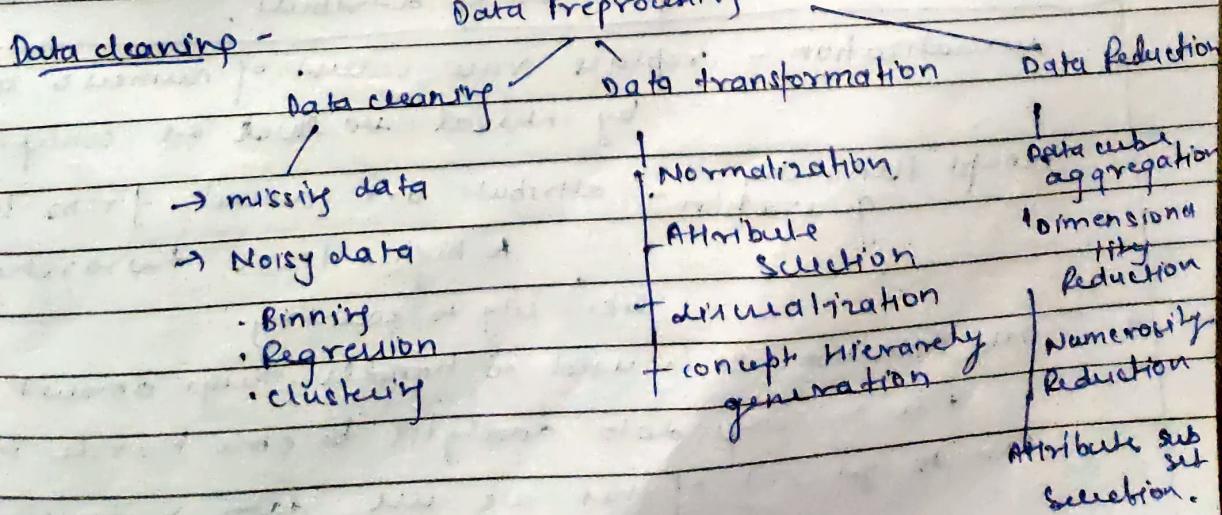
- cost efficient, decision making
- help businesses in modification operation/strategy

disadvantage -

- time consuming
- need training to work on analytics
- Not precise.
- right mining tool is challenging to find.

Data warehouse - refers to a place where data can be stored for useful mining. It is like computer with lot of storage capacity. data from various sources e.g. are copied to warehouse. where it can be further fetched.

+ Data preprocessing - is a data mining technique which is used to transform the raw data in useful and efficient format.



c) Data cleaning - data can have irrelevant and missing parts, to handle this cleaning is done.

a) Missing data → ignore tuple
→ fill missing val. ← mean
median

b) Noisy data - Noisy data is meaningless data cannot be interpreted by machine, generated due to faulty data collection, error.

- binning method - works on sorted data
- divided into segments and various methods are applied.
- value can be replaced by mean-boundary values.

- Regression - data can be made smooth by fitting it to regression function.

- clustering - this approach groups data in cluster, outliers may be undetected or will fall outside cluster.

* - Data transformation - Step is taken to transform data in appropriate forms suitable for mining.

- Normalization - scale data in specified range.

decimal scaling

$$v' = \frac{v_i}{10^j} \quad v' = \frac{v - \min(A)}{\max(A) - \min(A)} \quad \begin{cases} \min & \\ \max & \end{cases} \quad z\text{-score}$$

$$v' = \frac{v - \bar{A}}{\sigma_A} - \text{s.d}$$

Attribute selection - new attribute are constructed from given set of attributes.

discretization - replace raw values of numeric attributes by interval so level or conceptual level.

concept hierarchy generation - attribute are covered from lower level to higher level hierarchy.
like - city to country.

* Data Reduction - It is used to handle huge amount of data due to large data analysis become harder in such case to get rid of this we use data reduction.

- Data cube aggregation - It is applied for construction of data cube , like quarterly, yearly.
- Attribute subset selection - highly relevant attribute should be used rest can be discarded.
- Numerosity redu - This enables to store the model of data instead of whole data ex. Regression model. Actual data is replaced by mathematical model.
- Dimensionality reduction - Reduce size of data by encoding mechanism. Ex- PCA- Principal component analysis. Such reduction can be lossy or lossless.
- ❖ Data integration - It is process of combining data from multiple heterogeneous sources into single unified view.

↳ Two approaches to data mining -

- 1) Tight coupling - Data from various sources is combined by process of ETL - Extract", Transfer", Loading".
- 2) loose coupling - an interface is provided that takes query from user and transform in a way Database can understand.

- data remains in actual source database.

Assume - data redundancy

- schema integration

- Data value conflicts.

- Data compression - technique reduces the size of file using different encoding mechanism. (Huffman - encoding).

↳ Lossless compression

↳ Lossy compression

* Data generalization - Process of summarizing data by replacing relatively low level values with higher level concepts.

It is form of descriptive mining.

1) Data cube approach - OLAP approach

efficient

computation and result are stored in data cube.

It involve count(), sum(), avg(), max() used for knowledge discovery

2) Attribute oriented - online data analysis, query oriented
- generalization based on attributes.

tuples are merged and their counts are accumulated.

- It includes attribute removal, attribute generalization,

* Analytical characterisation -

It is statistical approach for preprocessing data to filter out irrelevant attribute or rank relevant attribute. This preprocessing is called analytical characterisation.

Reasons - decide which dimension must be included

- produce high level generalization

- reduce attribute that support pattern easily.

Mining class comparison

class discrimination or comparison that distinguish a target class from its contrasting class. They should be comparable. Ex - person, address and item are not comparable.

comparison method and implementation

- Data collection -

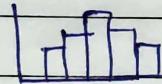
- Dimension relevance analysis -

- Synchronous generalization -

- Presentation of derived comparison - can be visualise in chart (count).

- Ex - we want to perform comparison between graduate and undergraduate students using discriminant rule.
- Measuring central tendency
- Mean, median, mode.
- Measuring dispersion of data.
- Range. - mean duration.
- Standard deviation.
- Graph display of statistical class description

- Histogram



- Line graph



- Bar graph



- Pie chart



Association Rule mining - It finds interesting association and relationship among large sets of data item.

Ex - market basket analysis.

Support count = frequency occurrence of a itemset

$$\text{support}(\text{Milk, Bread, diaper}) = 2.$$

Association rule - An implication of form

$x \rightarrow y$, x, y are two

- Rule Evaluation metrics

ID	data
1	Bread, milk
2	Bread, Diaper
3	Beer, egg
4	Milk, diaper, Beer
5	Bread, milk, diaper, beer

Support - Number of transaction that include x, y as part of rule or total percentage of all transaction. It measures how frequently item together.

- confidence - (It is ratio of no. of transaction that include all item in ΔB) as well as no. of transaction that include all item in A) to the no. of transaction that include all item in ΔA .

$$\text{confidence } (x \rightarrow y) = \text{support}(x \cup y) \div \text{support}(x).$$

- It measure how often each item in y appear in transaction

$$\text{lift} = (x \rightarrow y) = \text{supp}(x \rightarrow y) \div \text{supp}(x)$$

$\{ \text{Milk, diapers} \} \rightarrow \{ \text{Beer} \}$

$$\text{support} = \sigma(\text{Milk, diapers, Beer}) \div (T) = 2/5 = 0.4$$

$$\text{confidence} = \sigma(\text{Milk, diapers, Beer}) \div \sigma(\text{Milk, diapers})$$

$$= 2/3$$

$$\text{lift} = \frac{\text{supp}(\text{Milk, diapers, Beer})}{\text{supp}(\text{Milk, diapers}) * \text{supp}(\text{Beer})}$$

$$= 0.4 / 0.6 + 0.6 = 1.11$$

lift > 1 - appear together more than expected

if ≤ 1

- appear together less than expected

$$\text{lift} = \frac{\text{supp}(x \cup y)}{\text{supp}(x) * \text{supp}(y)}$$

Apriori algorithm

This algorithm is used to calculate association rules b/w object also called frequent pattern mapping. Name is apriori as it uses prior knowledge of frequent itemset property. We apply iterative approach.

Limitation - slow algorithm

- not efficient for large data

go 10^{14} frequent 1-itemset
it generates 10^{17} candidate

into two length which
in turn will be tested

- take memory

ID	item
T1	91, 92, 95
T2	92, 94
T3	92, 93
T4	91, 92, 94
T5	91, 93
T6	92, 93
T7	91, 93
T8	91, 92, 93, 95
T9	91, 92, 93

Step-I create table containing support count

$$91 = 6$$

$$92 = 7$$

$$93 = 6$$

$$94 = 2$$

$$95 = 2.$$

$$\text{min-supp} = \underline{2}$$

Step-II remove sets having support count less than minimum support count.

Step-II - generate candidate set C_2

91 92	4
91 93	4
91 94	1
91 95	2
92 93	4
92 94	2
92 95	2
93 94	0
93 95	1
94 95	0

remove who have less than
two min. support

Step-3 generate candidate set C_3

$$91 92 93 = 2$$

$$91, 92, 95 = 2.$$

Step-4 generate candidate set C_4

we have discovered all frequent item sets now
strong association rule come into picture for that
we need to calculate confidence.