

Data Warehouses contain huge volume of data. OLAP Systems demand that decision support queries be answered in the order of seconds.

Therefore, it is crucial for data warehouse systems to support highly efficient cube computation techniques, access methods & query processing techniques.

#### \* Efficient Data Cube Computation :-

Multidimensional data Analysis is the efficient computation of aggregation across many sets of dimensions.

In SQL, these aggregations are referred as "group-by", where each group-by can be represented by a "Cube", set of group-by's forms a lattice of cubes defining a data cube.

\* Compute Cube Operators :- The "Compute Cube" operators compute aggregate over all subsets of the dimension specified in the operation.

They can require excessive storage space, especially for large number of dimensions.

→ Consider the example of "AllElectronics" sales data cube that contains the following "city", "Item", "Year" & "Sales-In-Dollars" you want to be able to analyze the data, with queries such as foll

- \* "Compute the sum of Sales, grouping by City & Item"
- \* "Compute the sum of Sales, grouping by City"
- \* "Compute the sum of Sales, grouping by Item"

- The total no of Cuboids, or Group-by's that can be computed for this data cube is  $2^3 = 8$ . (City, Item, Year)
  - The possible group-by's are all : { (City, Item, Year), (City, Item, ), (City, Year), (Item, Year), () } . Where "()" mean that the group-by is Empty.

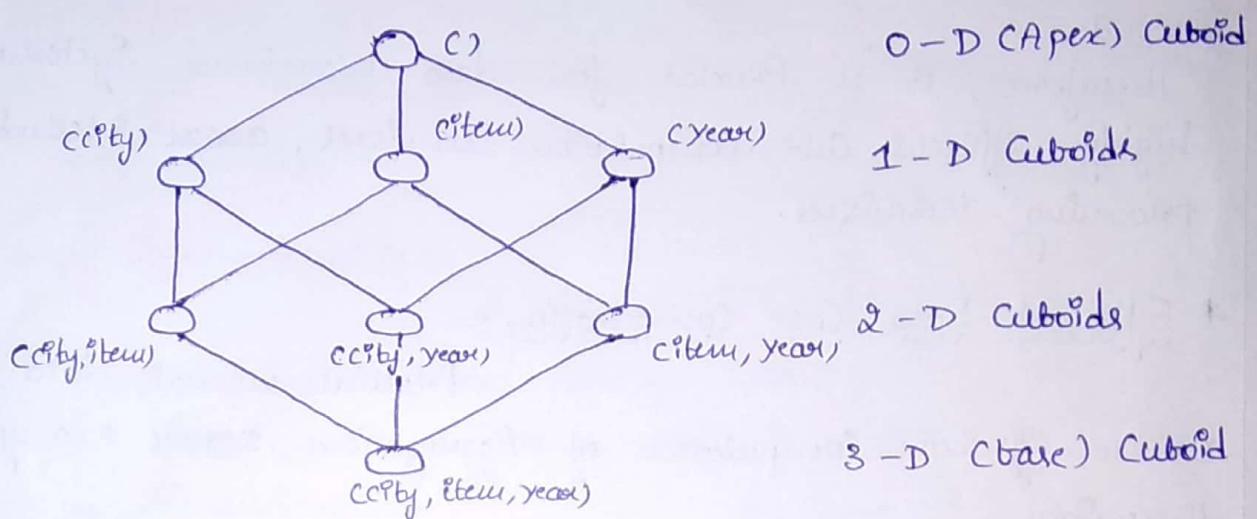


Fig Lattice of Cuboids

- The "Base Cuboid" contains all three dimensions "City", "Item" & "Year". It can return the total Sales for any combination of three dimensions.
- The "Apex Cuboid" or "0-D" refers to the case where the group-by is Empty. It contains the total sum of all Sales.
- The Base Cuboid is the Least generalized (Most Specific) of the cuboids.
- The Apex Cuboid is the Most generalized (Least Specific) of the Cuboids.
- An SQL query containing no group by (e.g. "compute the sum of total sales") is a zero dimensional operation.
- An SQL query containing one group by (e.g. "compute the sum of sales, grouping by City") is a one dimensional operation.
- A cube operator on n dimensional is equivalent to a collection of group-by statements.

→ The SQL Syntax could be defined as

"define Cube Sales-Cube [City, year, Item] : Sum (Sales-In-dollars)"

→ A Statement Such as

"Compute Cube Sales-Cube"

would Explicitly need to Compute the Sales Aggregate Cuboid for all Eight Subsets of the set [City, year, Item], Including the Empty Subset.

### \* Curse of Dimensionality :-

A Major challenge related to precompilation is that the required Storage Space may Explode if all Cuboids in data cube are precomputed. Especially when the cube has Many dimensions.

The Storage requirements are Even More Excessive when Many of dimensions have associated Concept hierarchies. This problem is referred as "Curse of Dimensionality".

→ For Example, time is usually Explored not at Only one Conceptual Level (e.g. Year)

But rather at Multiple Conceptual Levels Such as In the hierarchy "day < Month < Quarter < Year".

→ For an n-dimensional Data cube, the total no of Cuboids that can be generated is

$$\boxed{\text{Total no of Cuboids} = \prod_{i=1}^n (L_i + 1)}$$

where  $L_i$  be the number of levels associated with the Dimension i. One is added to  $L_i$  to include the Virtual top level, all.

## \* partial Materialization & Selected Computation of Cuboids :-

There are

three choices for Data Cube Materialization

→ "No Materialization": Do not precompute any of the "non-base" Cuboids.

This leads to computing Expensive Multidimensional Aggregates on-the-fly, which can be Extremely Slow.

→ "Full Materialization": precompute all the Cuboids. The Hasse-Poset Lattice of Computed Cuboids is referred as "Full Cube".

This choice typically requires huge amounts of Memory Space in order to store all precomputed Cuboids.

→ "partial Materialization": Selectively compute a proper subset of the whole set of possible Cuboids.

We may compute a subset of the cube, which contain only those cells that satisfy some user-specified criterion, such as where the tuple count of each cell is above some threshold.

→ The partial Materialization of Cuboids or Subcubes should consider three factors

- \* Identify the subset of Cuboids or Subcubes to Materialize.
- \* Exploit the Materialized Cuboids or Subcubes during query processing.
- \* Efficiently update the Materialized Cuboids or Subcubes during Load & Refresh.

→ Selection of Subset of Cuboids or Subcubes to Materialize should take into account the queries in workload, their frequency, & their access cost. In addition to storage requirements.

→ A popular Approach is to Materialize the Cuboid Set on which other frequently referenced Cuboids are based.

We can compute an "icing cube" which is a data cube that

Stores only those Cube cells with an aggregate value.

- Another common strategy is to materialize a "Shell Cube", this involves precomputing the cuboids for only a small number of dimensions (e.g., three to five) of a Data Cube.
- Once the selected cuboids have been materialized, it is important to take advantage of them during query processing. Such as
  - \* How to determine the relevant cuboid from the materialized cuboids
  - \* How to use available index structures on materialized cuboids.
  - \* How to transform the OLAP operations onto the selected cuboid.
- Finally, during load & refresh, the materialized cuboid should be updated efficiently.

#### \* Indexing OLAP Data:-

Indexing will be done on OLAP data to facilitate efficient data accessing. Most data warehouse systems support index structures & materialized views (cuboids) for efficient data accessing.

There are two ways of indexing OLAP data

- \* Bitmap Indexing
- \* Join Indexing

\* Bitmap Indexing :- Is a popular indexing method, because it allows quick searching in data cube.

→ In the bitmap index for a given attribute, there is a direct bit vector (BV),

For each value 'v' in the attribute's domain. If a given cell's domain consists of  $n$  values, then ' $n$ ' bits are needed for each entry in the bitmap index.

- If the attribute has the value V for a given row in the data table, then the bit representing that value is set to '1' in corresponding row of bitmap index.  
All other bits for that row are set to '0'.
- Consider an example of "All Electronics" data warehouse, Suppose the dimension "Item" at the top level has four values such as "Home Entertainment", "Computer", "Phone" & "Security".

Base table			Item bitmap table				City bitmap table			
RID	Item	City	RID	H	C	P	S	RID	V	T
R <sub>1</sub>	H	V	R <sub>1</sub>	1	0	0	0	R <sub>1</sub>	1	0
R <sub>2</sub>	C	V	R <sub>2</sub>	0	1	0	0	R <sub>2</sub>	1	0
R <sub>3</sub>	P	V	R <sub>3</sub>	0	0	1	0	R <sub>3</sub>	1	0
R <sub>4</sub>	S	V	R <sub>4</sub>	0	0	0	1	R <sub>4</sub>	1	0
R <sub>5</sub>	H	T	R <sub>5</sub>	1	0	0	0	R <sub>5</sub>	0	1
R <sub>6</sub>	C	T	R <sub>6</sub>	0	1	0	0	R <sub>6</sub>	0	1
R <sub>7</sub>	P	T	R <sub>7</sub>	0	0	1	0	R <sub>7</sub>	0	1
R <sub>8</sub>	S	T	R <sub>8</sub>	0	0	0	1	R <sub>8</sub>	0	1

- Bitmap Indexing is advantageous compared to "hash" or "tree indices"  
Because comparison, join & aggregation operations are reduced to bit arithmetic which substantially reduce the processing time.
- Bitmap indexing provides significant reduction in space & input/output (I/O) since a string of characters can be represented by a single bit.

- \* Join Indexing - Method gained popularity from its use in relational database query processing.
- Traditional indexing maps the value in a given column to a set of rows having the value.
- Join indexing retrieves the joinable rows of two relations from a relational database.  
For ex If two relations R(RID, A) and S(B, SID) join on

the attribute 'A' and 'B', then the Join Index record contains the pair (RID, SID), where RID & SID are record identifiers from the R and S relations.

- Join Index records can identify joinable tuples without performing costly join operations.
- The Star Schema Model of data warehouse makes indexing attractive for cross-table search.  
Because the linkage between a fact table & its corresponding dimension table comprises the fact table's foreign-key & the dimension table's primary key.
- Join indexing maintains relationship between attribute values of a dimension (within dimension table) & the corresponding row in the Fact table.
- Consider an example of "AllElectronics" star schema of a join index relationship between the "Sale" fact table & the "Location" & "Item" dimension tables

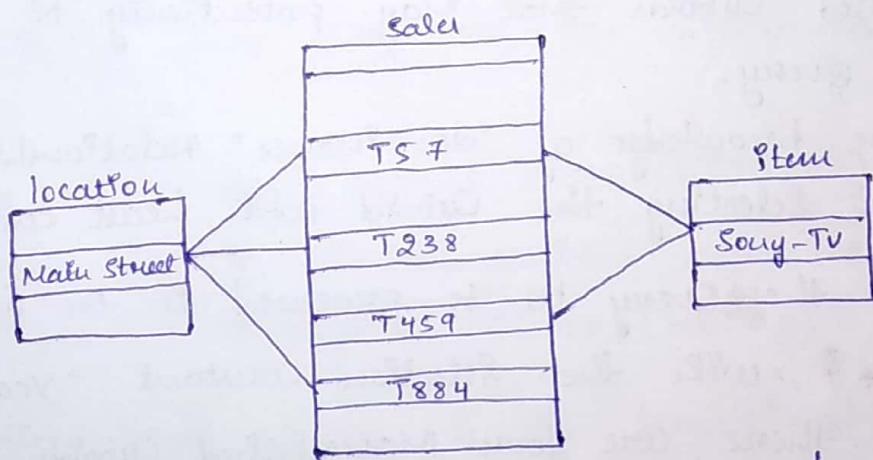


Fig: Linkage between Sale fact table & Location & Item Dimension table

location	sale-key
...	....
Main Street	T57
Main Street	T238
Main Street	T884
....	....

item	sale-key
...	....
Sony-TV	T57
Sony-TV	T459
....	....

location	item	sale-key
...	...	...
MainStreet	SonyTV	T57
....	....	....
....	....	....

→ To Speed-up query processing, the JOIN Indexing & bitmap indexing Methods can be Integrated to form "bitmapped JOIN Indexer".

\* Efficient processing of OLAP queries -

The purpose of Materialization Cuboids & constructing OLAP Index Structures is to Speed up query processing in data cube.

The query processing in Data cube should proceed as follows.

\* Determine which operation should be performed on available Cuboids - The Involve transformation of any selection, projection, roll-up or drill-down operations Specified in the query into Corresponding SQL & OLAP operation.

\* Determine to which Materialized Cuboid(s) the relevant operation should be performed - The Involve Identifying all of the Materialized Cuboids that may potentially be used to answer the query.

By using the knowledge of "dominance" relationships among the Cuboids & Selecting the Cuboid with least cost.

→ Suppose that the query to be processed is on {Brand, province-one-state}, with the Selection Constant "year = 2010" Also Suppose that there are four Materialized Cuboids available, as follows:

\* Cuboid 1: {Year, Item-name, City}

\* Cuboid 2: {Year, Brand, country}

\* Cuboid 3: {Year, brand, province-one-state}

\* Cuboid 4: {Item-name, province-one-state} where year = 2010

"Which of these four Cuboids should be selected?"

→ Finer-granularity data cannot be generated from Coarser-granularity data.

Therefore, "Cuboid 2" cannot be used because "Country" is a more general concept than "Province-or-State".

"Cuboids 1, 3 & 4" can be used to process the query.

→ "How would the cost of each cuboid compare if used to process the query?"

\* Cuboid 1 would cost the most because both "Item-name" & "Cpty" are at lower level than "Province-or-State" & "Brand".

\* Cuboid 3 be smaller than Cuboid 4 & thus Cuboid 3 should be chosen to process the query.

If efficient indices are available for Cuboid 4, then Cuboid 4 may be a better choice.

\* OLAP Server Architectures & ROLAP require MOLAP values while HOLAP cubes in a data warehouse are stored in three different modes.

\* Relational OLAP (ROLAP)

\* Multidimensional OLAP (MOLAP)

\* Hybrid OLAP (HOLAP)

\* Multidimensional OLAP is the traditional mode in OLAP analysis. In MOLAP data is stored in form of multidimensional cubes & not in relational database.

→ The advantage of this mode is that it provides excellent query performance & the cubes are built for fast data retrieval.

All calculations are pre-generated when the cube is created & can be easily applied while querying data.

→ The disadvantages of this Model are that it can handle only a limited amount of data.  
Since all calculations have been pre-built when the cube was created, the cube cannot be derived from a large volume of data.

\* Relational OLAP- Data in the model is stored in relational database. Since the data is stored in relational database, this model give the appearance of Traditional OLAP's slicing & dicing functionality.

→ The advantage of this model is that it can handle a large amount of data & can leverage all the functionalities of the relational database.

→ The disadvantages are that the performance is slow & each ROLAP report is also limited by SQL functionalities.

\* Hybrid OLAP- Approach combine ROLAP & MOLAP technology. Benefiting from the greater Scalability of ROLAP & faster computation of MOLAP.

→ HOLAP leverage cube technology & for drilling down into details it uses the ROLAP model.

The Microsoft SQL Server 2000 supports a hybrid OLAP server.

\* Specialized SQL Servers- To meet the growing demand of OLAP processing in relational database, some database system vendors implement specialized SQL servers to implement:

- \* Advanced query language
- \* Advanced query processing

Over "Star" & "Snowflake" Schema in ready-only environment.

### \* Data mining :-

Data Mining is a technology that blends traditional Data analysis Method with Sophisticated Algorithms for processing large volume of data.

→ It has also opened up exciting opportunities for Exploring & Analyzing new types of data & for analyzing diff types of data in new ways.

"Data Mining is a process of automatically discovering useful information in large Data Repository".

\* Data Mining as knowledge Discovery :- Data Mining is an integral part of knowledge Discovery in Database (KDD), which is a overall process of converting raw data into useful info.

→ This process consists of a series of Transition Steps from "Data Preprocessing" to "post processing of Data Mining Results".

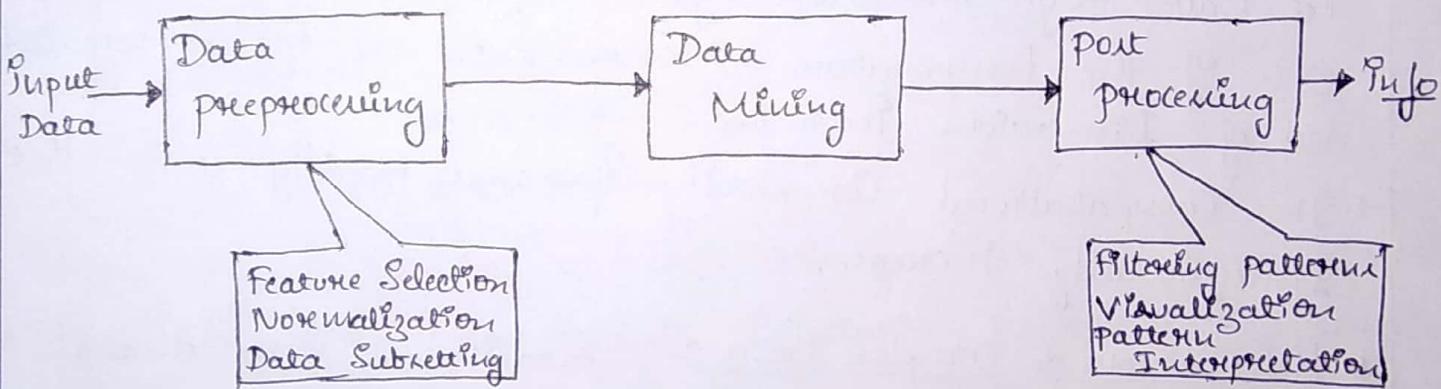


Fig :-

The Input data can be stored in a variety of file formats such as c files, Spread Sheets or Relational tables.

→ The purpose of preprocessing is to transform the raw input data into an appropriate format for subsequent analysis.

The steps involved in Data preprocessing are

- \* Fusion Data from Multiple Sources.

- \* Cleaning Data to remove noise.
  - \* Selecting records & features to data Mining, this is the most time consuming step in knowledge discovery process.
- "post processing" steps that ensure that only valid & useful results are incorporated into the Decision Support System.

#### \* Challenges -

The motivating challenges of Data Mining are.

- \* Scalability- Advances in Data generation & collection of data sets with sizes of Exabytes, Petabytes or even Petabytes are becoming common.  
→ If the data mining algorithms are to handle these massive data sets, then they must be Scalable.
- \* High Dimensionality- It is now common to encounter data sets with hundreds or thousands of attributes.  
→ Data sets with temporal or spatial components also tend to have high dimensionality.  
Ex:- If the temperature measurements are taken repeatedly no of dimensions increase.
- \* Computational Complexity increases rapidly with the dimensionality increase.
- \* Heterogeneous & Complex Data- Traditional Data Analysis Methods often deal with data sets containing attributes of the same type, either continuous or categorical.  
→ The role of Data Mining in Business, Science, Medicine & other fields have grown, so has the need for techniques that handle heterogeneous attributes.
- \* Data Ownership & Distribution- Sometimes the Data Model for analysis is not stored in one location or owned by one org instead data is geographically distributed.

→ The key challenges faced by distributed data Mining Algorithm include

- \* How to reduce the amount of communication required to perform distributed computation.
- \* How to efficiently consolidate the Data Mining results obtained from multiple sources.
- \* How to address Data Security issues.
- \* Non-Traditional Analysis - Traditional Statistical Approach is based on a "hypothesize-and-Test" paradigm.

→ Current Data Analysis tasks often requires the generation & evolution of thousands of hypotheses & consequently the development of some data mining techniques has been motivated. It is a non-traditional approach.

#### \* Data Mining Tasks :-

The Data Mining tasks can be classified generally into two types based on what a specific task tries to achieve.

- \* predictive Task
- \* Descriptive Task.

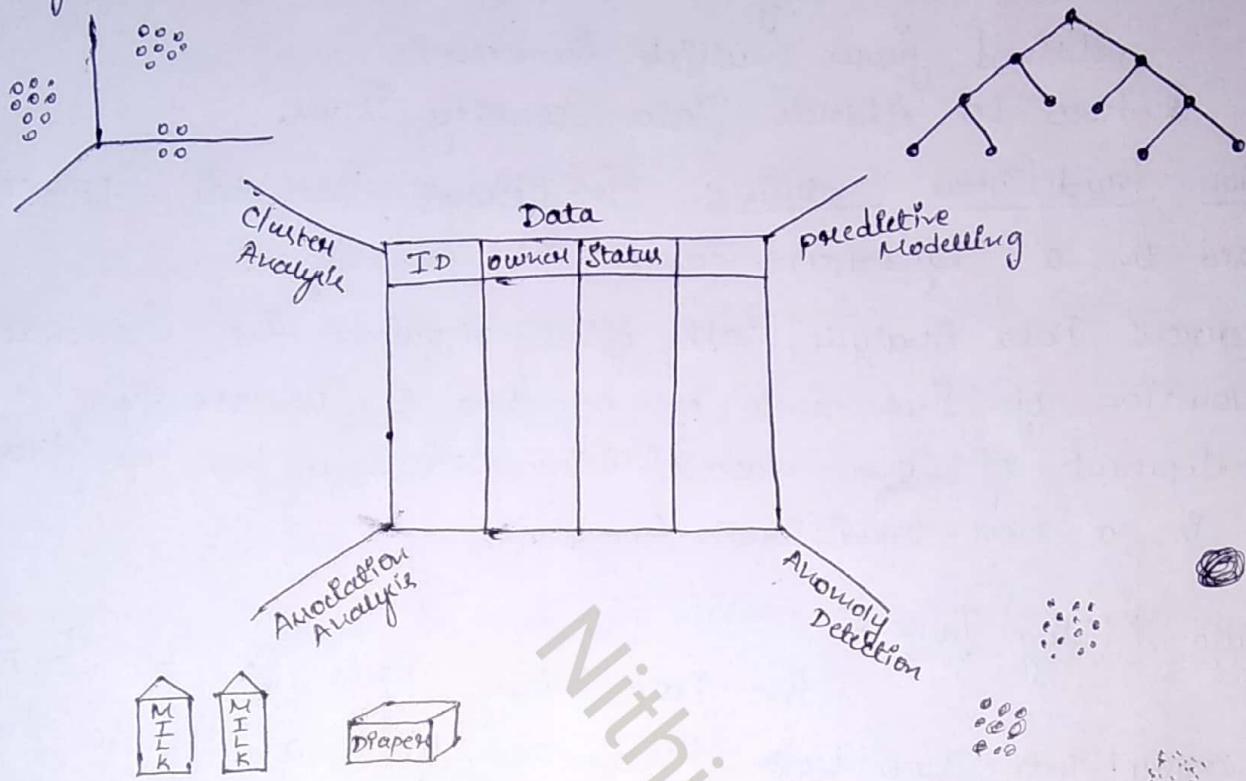
There are a number of Data Mining tasks such as classification, prediction, Time-Series Analysis, Association, Clustering, Summarization etc. All these tasks are either predictive tasks or Descriptive tasks.

\* predictive Tasks → Usually from the available data set that is helpful in predicting unknown or future values of another data set of interest.

A medical practitioner trying to diagnose a disease based on the medical test results of a patient can be considered as a predictive Data Mining task.

\* Descriptive Task → Usually finds data describing patterns & comes up with new, significant information from the available Data Set.

A retailer trying to identify products that are purchased together can be considered as a descriptive data mining task.



\* Predictive Modeling → Is a process that uses Data Mining & probability to forecast outcomes.

→ There are two types of predictive modeling

- \* Classification → used for Discrete Variables

- \* Regression → used for Continuous Variables.

The goal of both tasks is to learn a model, in which error between the predicted & true value is low.

→ For example, a model can predict the income of an employee based on Education, Experience & other demographic factors like place of stay, gender etc.

\* Anomaly Detection :- Is a task of identifying observations whose characteristics are significantly different from the rest of Data.

It is also called as "outlier Detection".

→ A good Anomaly detector must have a high detection rate & a low false alarm rate.

The good anomaly detection algo is to discover real anomalies & avoid false labeling.

→ Typically the anomalous item will translate to some kind of problem such as bank fraud, structural defect, medical problems etc.

\* Association Analysis :- Association discovers the association or connection among a set of items. Association identifies the relationships among objects.

→ Association Analysis is used for Commodity Management, Advertising, Direct Marketing etc.

→ For Example, A retailer can identify the products that normally customers purchase together or even find the customers who respond to the promotion of same kind of products.

\* Clustering Analysis :- Is used to identify data objects that are similar to one another.

→ The similarity can be decided based on a no of factors like purchase behavior, responsiveness, geographic location & so on.

→ For Example, An Insurance Company can cluster its customers based on age, residence, income etc.

\* Types of Data :-

Data Set can be often be viewed as a collection of data objects. Other names for a data object are "record", "point", "vector", "pattern", "Event", "case", "sample" or "observation".

Data can be classified into two broad types such as

- \* Qualitative
- \* Quantitative (numeric)

\* Qualitative Data :- arise when the observations fall into separate distinct categories.

→ For Ex:- colors of Eye - blue, green brown  
Exam result - fail, pass

Such data are inherently "discrete", in that there are finite number of possible categories into which each observation may fall.

\* Quantitative Data :- are numerical data arise when the observations are "counts" or "measurements".

→ The data are said to be "discrete" if the measurements are integers (Ex no of people in a house) & "continuous" if the measurements can take on any value, usually within some range (Ex weight).

\* Attribute of Measurement of Data :- An "Attribute" is a property or characteristic of an object that may vary, either from one object to another or from one time to another.

→ For Ex Eye color varies from person to person, while the temperature of an object varies over time.

\* Measurement Scale is a rule (function) that associates a numerical or symbolic value with an attribute of an object.

For Ex:- we count the no of chairs in a room to see if there will be enough to seat all the people coming to a meeting.

\* Different Types of Attributes- A useful way to specify the type of an attribute is to identify the properties of numbers that correspond to underlying properties of the attribute.

→ Distinguish = ≠ (Nominal) :- the values of a nominal attribute are just different names.

Nominal values provide only enough info to distinguish one object from another (=, ≠).

For ex:- Zip codes, Employee ID.

→ Order < ≤ > & ≥ (ordinal) :- The values of an ordinal attribute provide enough info to order objects.

For ex:- Readers, Street Number.

→ Interval + - (Addition) :- For Interval attributes, the differences between values are meaningful, i.e. a unit of measurement exists (+, -).

For ex:- Calendar dates.

→ Multiplication \* / (Ratio) :- For ratio variables both differences & ratios are meaningful (\*, /).

For ex:- Count, Age, Mass etc.

\* General Characteristics of Data Sets :- There are 3 characteristics that apply to many data sets & have a significant impact on the data mining techniques.

\* Dimensionality

\* Sparsity

\* Resolution

→ Dimensionality :- refers to how many attributes a dataset has.

For ex:- Healthcare data is notorious for having vast amounts of variables (Blood pressure, weight, cholesterol level).

This data could be represented in spreadsheet, with one column representing each dimension.

- Sparsity :- In one  $\Sigma$  in which a relatively high percentage of the Variable cells do not contain actual data.  
For ex:- a new variable dimensioned by MONTH for which you do not have data for part Month.
- Resolution :- It is frequently possible to obtain data at different levels of resolution, & often the properties of the data are different at different resolutions.  
For ex:- the Surface of Earth seems very uneven at a resolution of few meters, but is relatively smooth at a resolution of tens of kilometers.

\* Type of Data Sets :- There are Many types of data Sets & as the field of Data Mining Develops & Matures, a greater Variety of Data Sets become available for Analysis.

- Record Data :- is a collection of records (data objects), Each of which consists of a fixed set of data fields (Attributes).  
there is no explicit relationship among records or data fields. Every record has the same set of attributes.
- Record data is usually stored either in flat files or in relational database.
- Different types of Record data are
  - \* Transaction Data → where Each record (Transaction) involves a set of Items.
  - \* Data Matrix → collection of Data with same fixed set of Numeric attributes.
- \* Document-Term Matrix → where Each term is a component (Attribute) of the Vector.
- Graph-Based Data :- A graph can sometimes be a convenient & powerful representation of Data.

- There are two diff types of Graph-Based Data
  - \* Data with relationship among objects → Relationship among objects frequently convey important information.  
Ex: Web pages in World Wide Web
  - \* Data with objects that are Graphs → If objects contain Subobjects that have relationships, then such objects frequently represented as Graphs.

→ Ordered Data is one in which the attributes have relationship that involve order in time or space.

→ Different types of Ordered Data are

- \* Sequential Data → an Extension of Record data, where each record has a time associated with it.  
Ex: retail transaction data with time of transaction.
- \* Time Series Data → is a special type of Sequential Data in which each record is a time series.  
a series of measurements taken over time.
- \* Spatial Data → Some objects have spatial attributes, such as positions or areas, as well as other types of attributes.  
Ex: Weather Data.

\* Data Quality :-

Todays real-world databases are highly susceptible to noisy, missing & incomplete Data due to their typically huge size & their likely origin from multiple heterogeneous source.

→ The Data Mining Applications focus on

- \* the detection & correction of Data quality problems
  - \* use of algorithms that can tolerate poor Data Quality
- "Detection" & "correction" is often called "Data cleaning".

\* Measurement & Data Collection Errors - refer to any problem resulting from the Measurement process.

A common problem is that the value recorded differs from the true value to some extent.

→ The term "data collection errors" refers to errors such as omitting data objects or attribute values or inappropriate including data object.

\* Noise and Artifacts - Noise is the random component of a measurement error. It may involve the distortion of a value or addition of spurious objects.



fig: Time Series

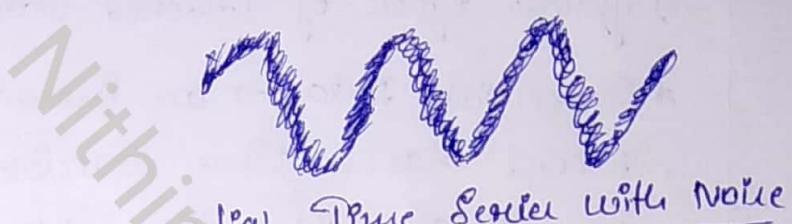


fig: Time Series with Noise

→ Data Errors may be the result of a more deterministic phenomenon. Such deterministic distortions of the data are often referred to as "Artifacts".

Ex: A streak in same place on a set of photographs.

\* Precision, Bias & Accuracy - "precision" is a closeness of repeated measurements to one another.

"Bias" is a systematic variation of measurements from the quantity being measured.

→ precision is often measured by the standard deviation of a set of values, while bias is measured by taking the difference between the mean of the set of values & quantity being measured.

→ It is common to use the general term, "Accuracy" to refer to the degree of measurement error in data.

Accuracy depends on precision & bias.

\* Outlier :- An outlier is an observation point that is distant from other observations (Data).

An outlier may be due to variability in the measurement or it may indicate experimental error.

→ The latter are sometimes, outliers will be excluded from the data set.

An outlier can cause serious problems in analysis.

\* Missing Value :- Imagine you need to analyze an enterprise data, you note that many tuples have no recorded value for several attributes. In such case, fill the missing values by following methods.

1. Ignore the tuple → usually done when class label is missing, this method is very ineffective.
2. Fill the missing value manually → It's very time consuming & not feasible.
3. Use a global constant to fill missing value → Replace all missing values by some constant such as "Unknown" or "00" this method is simple, but not fool proof.
4. Use a measure of central tendency for the attribute (Mean or Median) to fill missing values.
5. Use the most probable value to fill the missing value → this may be determined with regression, interface based tool using Bayesian or decision tree.

\* Inconsistent Values :- Data can contain inconsistent values.

Some types of inconsistent values are easy to detect.

For ex :- A person's height should not be negative.

Once an inconsistency has been detected, it is sometimes possible

to Connect the data. The Connection of an Inconsistency Requires additional or redundant Information.

\* Duplicate Data :- A data Set May Include data objects that are duplicates, or almost duplicates of one another.

→ To Detect & Eliminate Such duplicates, two Main Issues Must be addressed

- \* If there are two objects that actually represent a Single object with different values, & these Inconsistent values Must be resolved
- \* Care to be taken to avoid accidentally Combining data - objects that are similar, but not duplicate

The "Deduplication" Is often used to deal with duplicate data.

\* Data preprocessing :-

Is a Data Mining technique that involves transforming raw data into an understandable format.

Real-world data Is often Incomplete, Inconsistent, and/or Lacking In certain behaviors or trends, & Is likely to contain Many Errors.

Data preprocessing Is a proven Method of resolving Such Issues.

→ Some of the Most Important Approaches for Data preprocessing are

- \* Aggregation
- \* Sampling
- \* Dimensionality Reduction
- \* Feature Subset Selection
- \* Feature Creation
- \* Discretization & Binarization
- \* Variable Transformation

These Approaches fall Into two categories, Selecting data objects & Attributes for Analysis or Creating/Changing the attributes.

\* Aggregation- Aggregation refers to combining two or more attributes (or objects) into a single attribute (or object).  
For ex- Merging daily Sales figures to obtain Monthly Sales figures.

→ Data Aggregation helps in Data Reduction, Allows use of More Expensive Algorithms.

If done properly, Aggregation can act as Scope or Scale, providing a High Level View of data instead of Low Level View.

→ The Aggregate Quantities have "less Variability" than the Individual Objects.

\* Sampling- Is the process of understanding characteristics of Data or Models Based on a Subset of the Original Data.  
It is used Extensively in all aspects of Data Exploration & Data Mining.

→ Obtaining the "Data of Interest" on Entire Set is too Expensive & time consuming

So we can use Sampling, where Entire Set of data may not be necessary.

→ "Representative Sample" is a Representative for a particular operation, if it results in approximately the same outcome as if the Entire data set was used.

→ There are two probability of Selecting any particular Item

\* Sampling without replacement → Once an item is selected it is removed from the population for obtaining Future Samples.

\* Sampling with replacement → Selected item is not removed from the population for obtaining the Future Samples.

→ Even if proper Sampling technique is known, it is important to choose proper Sample Size.

- \* Larger Sample Size increases the probability that a Sample will be "representative".

- \* Smaller Sample Size may result in Mixed or Erroneous patterns Detected.

\* Dimensionality Reduction- "Curse of Dimensionality", the Data Analysis becomes significantly harder as the Dimensionality of the Data Increases.

→ Dimensionality Reduction is a process of Determining the Dimensions that are Important for Modeling.

Many of Data Mining Algorithms work better if the Dimensionality of Data is lower.

→ If Dimensionality Reduction Eliminates Irrelevant features or Reduces Noise, then Quality of Results May Improve. Can lead to More Understandable Model.

→ The Dimensionality Reduction can be done in two ways

- \* Redundant features → Duplicate much or all of the Info contained in one or more attributes.

Ex- purchase price of product & Sales tax paid contain the same info.

- \* Irrelevant features → Contain no Information that is useful for Data Mining task.

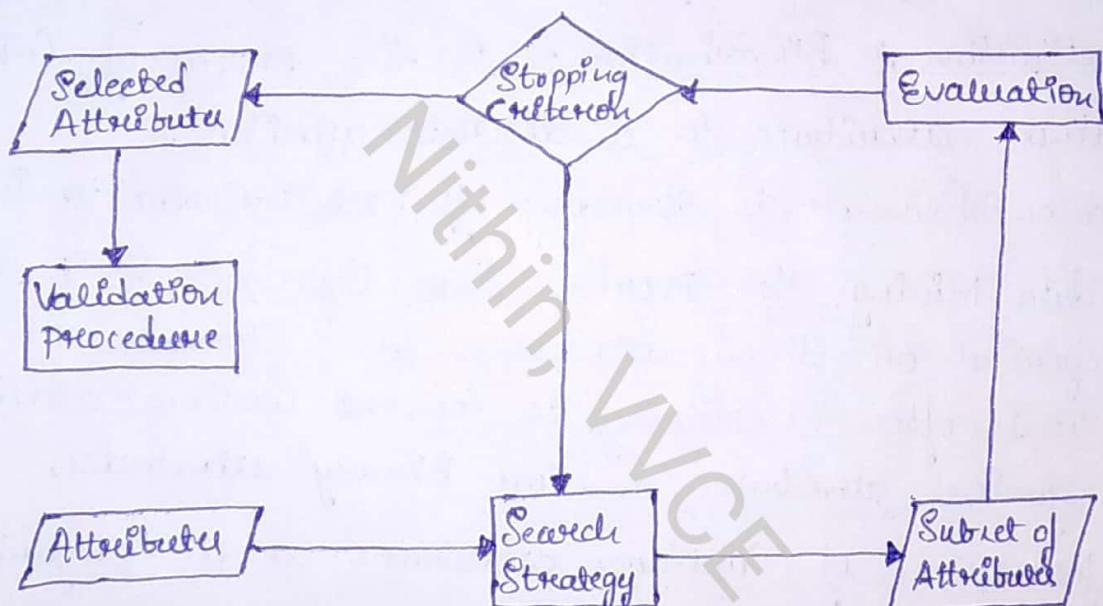
Ex- Student ID is irrelevant to the task of predicting their Marks.

→ Dimensionality Reduction can be done using "Principal Component Analysis (PCA)" is a Linear Algebra technique for continuous attributes that finds new Attributes.

\* Feature Subset Selection:- There are three Standard Approaches to Feature Selection.

- \* Embedded Approach → Feature Selection occurs naturally as part of the Data Mining Algorithm.
- \* Filter Approach → Features are Selected before the Data Mining Algorithm Is Run.
- \* Wrapper Approach → Use the target data Mining Algo as a black box to find the best subset of attributes.

→ The Architecture for Feature Subset Selection Is



The feature Selection process Is Viewed as consisting of four parts Such as, a Measure for Evaluating a Subset, a Search Strategy that Controls the generation of new Subset of features, a Stopping Criterion & a Validation procedure.

\* Feature Creation:- Sometimes, a Small Number of New attributes can Capture the Important Information In a Data Set Much More Efficiently than the Original attributes.

Such Attributes can be generated by three Methodologies

- \* Feature Extraction → Is a creation of a new, Smaller Set of features from the original Set of features.  
Ex:- In a set of photographs, where each photograph

be to be classified whether it's a human face or not.

- \* Mapping the Data to New Space → Sometimes a totally different view of the data can reveal important & interesting features.

Ex:- Applying Fourier Transform to data to detect time series patterns.

- \* Feature Construction → one or more new features constructed out of the original features.

Ex:- there are two attributes that record volume & mass of a set of objects, i.e. density = mass/volume.

- \* Discretization & Binarization - It is the process of converting a continuous attribute to a discrete attribute.

A common example is rounding off real numbers to integers.

- Some data mining algo requires that the data be in the form of categorical or binary attributes.

Thus, it is often necessary to convert continuous attributes to categorical attributes and/or binary attributes.

- Transformation of continuous attributes to a categorical attributes / binary attributes involves

- \* Deciding how many categories to have

- \* How to map the values of the continuous attributes to categorical attribute.

- \* Attribute Transformation - Refers to a transformation that is applied to all values of an attribute.

i.e. for each object, the transformation is applied to the value of the attribute for that object.

- There are two important types of Attribute Transformation

- \* Simple function Transformation →  $x^k$ ,  $\log x$ ,  $e^x$ ,  $\sqrt{x}$  etc

- \* Standardization or Normalization

## \* Measures of Similarity & Dissimilarity (Similarity Measure)

Similarity & Dissimilarity are

Important Because they are used by a number of Data Mining Techniques, Such as Clustering, Nearest Neighbors Classification & anomaly Detection.

→ The term "proximity" is used to refer to Either Similarity or Dissimilarity

\* Similarity between two objects Is a Numerical Measure of the degree to which the two objects are alike. Conversely, similarities are higher for pairs of objects that are more alike.

Similarities are usually non-negative & are often between "0" (no similarity) & "1" (complete similarity).

\* Dissimilarity between two objects Is a Numerical Measure of the degree to which the two objects are different.

Dissimilarity is lower for more similar pairs of objects.

→ The term "distance" is used as a Synonym for Dissimilarity. Dissimilarities sometimes fall in the interval  $[0, 1]$ , but it is also common for them to range from "0" to " $\infty$ ".

\* proximity Measures e. Especially similarities are defined to have values in the interval  $[0, 1]$ . If the similarity between objects can range from "1" (not similar) to "10" (completely similar).

we can make them fall into the range  $[0, 1]$  by formula.

$$S' = (S - \bar{S}) / q$$

where  $S$  &  $S'$  are original & the new similarity value respectively.

→ The More general case  $S'$  is calculated as

$$S' = (S - \min_S) / (\max_S - \min_S)$$

Where  $\text{min}_d$  &  $\text{Max}_d$  are the Minimum & Maximum  $d$ -  
Similarity Values respectively.

→ Likewise, Dissimilarity Measure with a finite range can be  
Mapped to the Interval  $[0,1]$  using the formula

$$d' = \frac{(d - \text{min}_d)}{(\text{max}_d - \text{min}_d)}$$

\* Similarity & Dissimilarity between Simple attributes - The  
proximity of objects with a number of attributes is defined  
by combining the proximities of individual attributes.  
Consider "p" & "q" are the attribute values for two data  
objects.

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & p=q \\ 1 & p \neq q \end{cases}$	$s = \begin{cases} 1 & p=q \\ 0 & p \neq q \end{cases}$
ordinal	$d = \frac{ p-q }{n-1}$ (Value mapped to Range from 0 to $n-1$ )	$s = 1 - d$
interval or Ratio	$d =  p-q $	$s = -d, s = \frac{1}{1+d}$ or $s = e^{-d}$

\* Distance - Distances are dissimilarities with certain properties  
the "Euclidean Distance"  $d$ , between two points  $x$  &  $y$  in one,  
two or high dimensional Space is given by

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

where  $n$  is the number of Dimensions &,  $x_k$  &  $y_k$   
are respectively, the  $k^{th}$  attribute of  $x$  and  $y$ .

- \* Similarities between Data objects e.g.  $s(x, y)$  is the Similarity between points  $x$  &  $y$ , then typically we will have
  - \*  $s(x, y) = 1$  only if  $x = y$  ( $0 \leq s \leq 1$ )
  - \*  $s(x, y) = s(y, x)$  for all  $x$  &  $y$  (Symmetry)

- \* Similarity Measures for Binary Data e.g. Similarity Measures between objects that contain only Binary attributes are called "Similarity Co-efficients".

→ "Simple Matching Co-efficient (SMC)" one commonly used Similarity Co-efficient is defined as

$$\boxed{SMC = \frac{\text{No of Matching attribute Value}}{\text{No of attributes}}} = \frac{f_{11} + f_{00}}{f_{11} + f_{00} + f_{01} + f_{10}}$$

This Measure counts both presence & absence Equally.

→ "Jaccard Co-efficient" Measures only the presence of an item is relevant

$$\boxed{J = \frac{\text{No of Matching presence}}{\text{No of attributes not involved in no Match}}} = \frac{f_{11}}{f_{11} + f_{10} + f_{01}}$$

Nithin, WCE