

Unit-5

- Aggregation in data mining is the process of finding, collecting and presenting the data in summarised format to perform statistical analysis of business schemes.
 - Aggregated data helps in solving relational problems which
 - reduce time in solving queries from dataset.
 - When data is collected from various sources it is necessary to get crucial information.
 - Data aggregator - are system that collect data from numerous sources then process and repackages them to useful data package.
- Help in query and delivery process where customer requests.

working of data aggregator -

- collecting - collect from various sources.
- Processing - Apply various ML & AI algo
- Presentation - show summarised data to user.

Collection of data

Processing of Data

Presentation of Data

- Can be done manually or automatically.

Types of data aggregation

- 1) Time - provide data point for single resource for defined time.
- 2) Spatial - provide data point for group of resources for defined time.

- Time interval for data aggregation - Reporting period - one day

Spatial { Granularity period - month.

Polling period - Every 7 min.

- Application - Ecommerce
travel
business analysis.

- Historical information - it is used to build model using this. new response can be generated based on previous data. ex - credit card fraud.

Query facility - It is a tool that allow client to retrieve data from warehouse by selecting field and some cases like query.

Data mining query language can be used to specify mining task.

+ OLAP function and tool

OLAP stands for on-line Analytical processing is a tool which help analyst, manager to gain insight into information in fast, consistent way.

- Used to analyze multidimensional data

Types-

1) • Relational OLAP - ROLAP server are placed between relational backend server and client frontend tool.

To store and manage data ROLAP use extended RDBMS

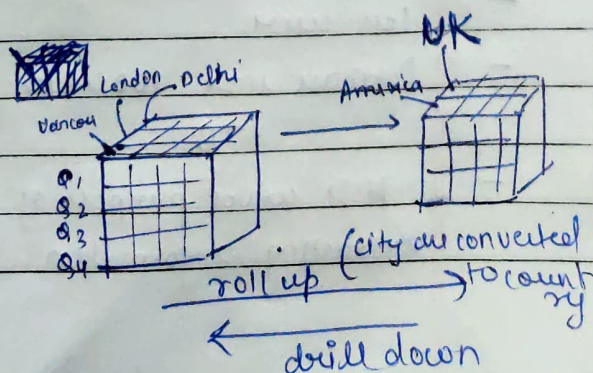
2) Multidimensional OLAP - MOLAP use array based multidimensional engine for multidimensional view of data. Storage utilization is low if data is sparse, therefore MOLAP server used ~~use~~ two level of data storage.

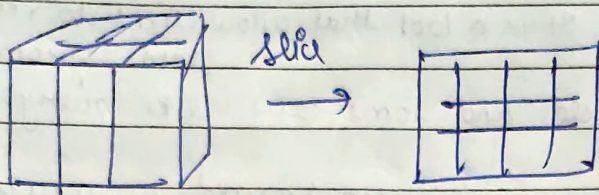
3) Hybrid OLAP - combination of both ROLAP and MOLAP offer high scalability ^{of ROLAP} and fast computation of MOLAP, can store large data

4) Specialised query server - It provide advanced query language and query processing support SQL queries over star and snowflake schema.

OLAP operations - Roll up

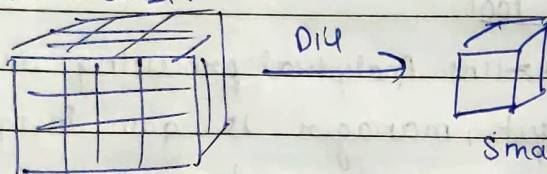
- drill-down
- slice
- dice
- Pivot





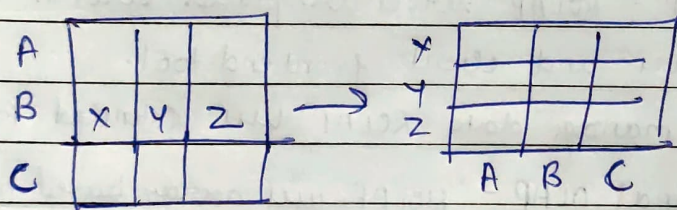
selects some part of data from data cubes.

→ Dice slices two or more dimension and provide a new sub-cube.



Small cube means some part of information.

→ Pivot - It is also known as rotation, it rotates data axes.



rotated

OLAP	OLTP
<ul style="list-style-type: none"> Involve historical processing of data Used by analysts, managers Useful in analyzing Based on star, snowflake and fact constellation schema Less user Database size 100GB to 1TB It serves purpose of extracting information 	<ul style="list-style-type: none"> Involve data processing (day to day) like bank servers. Used by clerk and professional Useful in running business. Based on ER-model more user Database size 100MB to 1GB It serve purpose of insert, update, delete.

- data mining interface -

data mining interface provide medium that allows user to communicate with data mining process.

It basically user interface where we see chart, graphs, or visualise the data to see result of mining process.

- security - The objective of data mining is to make large amount of data easily accessible to user without compromising security and data leakage.

so it is required to ^{have} security features to effect performance of the data warehouse.

security features are difficult to add when data warehouse has gone live.

security affects following features -

User access

Data loading

Data movement

Query management

- Backup - A backup is copy of data from database.

Backup can be physical and logical backup.

Physical backup include physical datafiles, control file.

logical backup contain logical data (ex-table or procedures).

- Recovery - Recovery is retrieving data which could be lost due to any hardware failure or any other issue.

- Tuning data warehouse - A data warehouse keeps evolving and it is unpredictable what query user is going to post in future.

- difficulties in data warehouse tuning
 - data warehouse is dynamic
 - difficult to predict query.
 - Business requirement keep changing.
 - data upload on warehouse also changes with time.
- Data load tuning - If there is delay in transferring data then system get affected so data load is important
- Integrity check - Integrity system avoid performance degrade of data load.

Testing - Testing is important for data warehouse system for data validation and to make them work correctly and efficiently.

- Unit testing - Each component module is tested
- Integration testing - In this testing all module are combined and then tested. done to detect fault in integrated modules.
- System testing - form of testing that validates and test whole data warehouse application.