Chapter Name

Data Warehousing and Business Intelligence Solutions

2024-03-18

Traditional Databases and Data Analysis

The importance of data analysis has been steadily increasing from the early 1990s

• Organizations in all sectors are being required to improve their decision-making processes in order to maintain their competitive advantage.

Traditional database systems do not satisfy the requirements of data analysis.

- They are designed and tuned to support the daily operations of an organization, and their primary concern is to ensure fast, concurrent access to data.
- This requires transaction processing and concurrency control capabilities, as well as recovery techniques that guarantee data consistency.
- These systems are known as operational databases or online transaction processing (OLTP) systems.
- The OLTP paradigm is focused on transactions.

OLAP – A New Paradigm Oriented to Analyze Data

- Since OLTP systems must support heavy transaction loads, their design should prevent update anomalies, and thus, OLTP databases are highly normalized.
 - Thus, they perform poorly when executing complex queries that need to join many relational tables together or to aggregate large volumes of data.
- Besides, typical operational databases contain detailed data and do not include historical data.
- The above needs called for a new paradigm specifically oriented to analyze the data in organizational databases to support decision making.
 - This paradigm is called online analytical processing (OLAP).
 - OLAP is focused on queries, in particular, analytical queries.
 - OLAP-oriented databases should support a heavy query load.

2024-03-18

Typical OLAP queries would ask, for example, for the total sales amount by product and by customer or for the most ordered products by customer. These kinds of queries involve aggregation, and thus, processing them will require, most of the time, traversing all the records in a database table. Indexing techniques aimed at OLTP are not efficient in this case: new indexing and query optimization techniques are required for OLAP. It is easy to see that normalization is not good for these queries, since it partitions the database into many tables. Reconstructing the data would require a high number of joins. Therefore, the need for a different database model to support OLAP led to the notion of data warehouses, which are (usually) large repositories that consolidate data from different sources (internal and external to the organization), and follow the multidimensional data model.

The Concept of Data Warehousing The concept of data warehousing has evolved out of the need for easy Data warehouse access to a structured store of quality data that can be used for decision making. Organizations have vast amounts of data • Demographic Client data but have found it increasingly difficult to Product data data access it and make use of it. In-house data **External data** A data warehouse is a database for the specific purpose of data analysis. Data warehouses vs. operational databases: Data warehouses are constructed with use 1. Contain more data of database technology. Data could be obtained from several 2. Data is only for reading sources. 3. Time is particularly respected.

A Database Model to Support OLAP

- Being dedicated analysis databases, data warehouses can be designed and optimized to efficiently support OLAP queries.
 - In addition, data warehouses are also used to support other kinds of analysis tasks, like reporting, data mining, and statistical analysis.
- OLAP is a method of organizing the large data stores, where data are well suited for analysis.
 - Data stored in relational databases are turned into meaningful, easy to navigate business information by creating a data cube.
- Data warehouses and OLAP systems are based on the multidimensional model, which
 - views data in an *n*-dimensional space, usually called a data cube or a hypercube.
- A data cube is defined by dimensions and facts.









Hierarchies (1/3)
The granularity of a data cube is determined by the combination of the levels corresponding to each axis of the cube.
 In our example, the dimension levels are indicated between parentheses: Category for the Product dimension, Quarter for the Time dimension, and City for the Customer dimension.
In order to extract strategic knowledge from a cube, it is necessary to view its data at several levels of detail.
 In the current example, an analyst may want to see the sales figures at a finer granularity, such as at the month level, or at a coarser granularity, such as at the customer's country level.
 Hierarchies allow this possibility by defining a sequence of mappings relating lower-level, detailed concepts to higher-level, more general concepts.
Given two related levels in a hierarchy, the lower level is called the child and the higher level is called the parent.
The hierarchical structure of a dimension is called the dimension schema, while a dimension instance comprises the members at all levels in a dimension.
2024-03-18













- According to the way in which they can be aggregated, measures can be classified as follows.
- Additive measures can be meaningfully summarized along all the dimensions, using addition.
 - These are the most common type of measures.
 - For example, the measure Quantity in our cube is additive: it can be summarized when the hierarchies in the Product, Time, and Customer dimensions are traversed.
- Semiadditive measures can be meaningfully summarized using addition along some, but not all, dimensions.
 - A typical example is that of inventory quantities, which cannot be meaningfully aggregated in the time dimension, for instance, by adding the inventory quantities for two different quarters.
- Nonadditive measures cannot be meaningfully summarized using addition across any dimension.
 - Typical examples are item price, cost per unit, and exchange rate.

























omparison Between Operational Databases and Data Warehous					
	Aspect	Operational databases	Data Warehouses		
1	User type	Operators, office employees	Managers, executives		
2	Usage	Predictable, repetitive	Ad hoc, nonstructured		
3	Data content	Current, detailed data	Historical, summarized data		
4	Data organization	According to operational needs	According to analysis needs		
5	Data structures	Optimized for small transactions	Optimized for complex queries		
6	Access frequency	High	From medium to low		
7	Access type	Read, insert, update, delete	Read, append only		
8	Number of records per access	Few	Many		
9	Response time	Short	Can be long		
10	Concurrency level	High	Low		
11	Lock utilization	Needed	Not Needed		
12	Update frequency	High	None		
13	Data redundancy	Low (normalized tables)	High (denormalized tables)		
14	Data modeling	UML, ER model	Multidimensional model		

Data Warehouse Architecture

@ A general data warehouse architecture consists of several tiers:

- The back-end tier is composed of extraction, transformation, and loading (ETL) tools, used to feed data into the data warehouse from operational databases and other data sources, which can be internal or external from the organization, and a data staging area, which is an intermediate database where all the data integration and transformation processes are run prior to the loading of the data into the data warehouse.
- The data warehouse tier is composed of an enterprise data warehouse and/or several data marts and a metadata repository storing information about the data warehouse and its contents.
- The OLAP tier is composed of an OLAP server, which provides a multidimensional view of the data, regardless of the actual way in which data are stored in the underlying system.
- The front-end tier is used for data analysis and visualization. It contains client tools such as OLAP tools, reporting tools, statistical tools, and data mining tools.











Physical Modeling of Data Warehouses

- The physical design of data warehouses is crucial to ensure adequate query response time.
- There are typically three common techniques for improving performance in data warehouse systems: materialized views, indexing, and partitioning.
- A materialized view is a view that is physically stored in a database, which enhances query performance by precalculating costly operations such as joins and aggregations.
- With respect to indexing, traditional techniques used in OLTP systems are not appropriate for multidimensional data. Thus, alternative indexing mechanisms are used in data warehouses, typically bitmap and join indexes.
 - A bitmap index is a special kind of database index that uses bitmaps. They are particularly useful for columns with a low number of distinct values.
 - Join indexes materialize a relational join between two tables by keeping pairs of row identifiers that participate in the join.
- Finally, partitioning or fragmentation divides the contents of a relation into several files, typically based on a range of values of an attribute.

Evolution of DSS into Business Intelligence
In the early days of Decision Support Systems, manager let their staff do some supportive analysis by using DSS tools.
As PC technology advanced, a new generation of managers evolved – one that was comfortable with computing and knew that technology can directly help make intelligent business decision faster.
New tools such as OLAP, data warehousing, data mining and intelligent systems, delivered via Web technology added promised capabilities and easy access to tools, models and data for computer-aided decision making.
These tools started to appear under the names Business Intelligence (BI) and business analytics in the mid-1990s.
Provide the second s
2024-03-18

Business Intelligence, Decision Support Systems and Data Warehousing

- Business intelligence comprises a collection of methodologies, processes, architectures, and technologies that transform raw data into meaningful and useful information for decision making.
- Business intelligence and decision-support systems provide assistance to managers at various organizational levels for analyzing strategic information.
 - These systems collect vast amounts of data and reduce them to a form that can be used to analyze organizational behavior.
 - This data transformation comprises a set of tasks that take the data from the sources and, through extraction, transformation, integration, and cleansing processes, store the data in a common repository called a data warehouse.
- Data warehouses have been developed and deployed as an integral part of decision support systems to provide an infrastructure that enables users to obtain efficient and accurate responses to complex queries.









21



