

Lectures in Data Warehouse

(4th I-Systems)(Work 1)

Lecture (1)

A data warehouse is constructed by integrating data from multiple heterogeneous sources. It supports analytical reporting, structured and/or ad hoc queries and decision making. This tutorial adopts a step-by-step approach to explain all the necessary concepts of data warehousing.

Audience

This tutorial will help computer science graduates to understand the basic-to-advanced concepts related to data warehousing.

Prerequisites

Before proceeding with this tutorial, you should have an understanding of basic database concepts such as schema, ER model, structured query language, etc.

Copyright & Disclaimer

□ Copyright 2014 by Tutorials Point (I) Pvt. Ltd.

المحاضرات أدناه موزعة على خمسة عشر أسبوعا (كل فصل لكل أسبوع)،
والفصل (١٦) تضمن أسئلة وأجوبة لمراجعة للمادة المعطاة..

مدرس مادة مستودع البيانات: أ.د. مرتضى محمد حمد
الصف الرابع/ قسم نظم المعلومات
كلية علوم الحاسوب وتكنولوجيا المعلومات/جامعة الأنبار

Table of Contents

About the Tutorial	i
Audience.....	i
Prerequisites.....	i
Copyright & Disclaimer	i
 Table of Contents.....	 ii
1. OVERVIEW.....	1
Understanding a Data Warehouse	1
Why a Data Warehouse is Separated from Operational Databases.....	2
Data Warehouse Features.....	2
Data Warehouse Applications.....	3
Types of Data Warehouse	3
2. CONCEPTS	5
What is Data Warehousing?	5
Using Data Warehouse Information.....	5
Integrating Heterogeneous Databases	5
Functions of Data Warehouse Tools and Utilities.....	6
3. TERMINOLOGIES.....	8
Metadata	8
Metadata Repository	8
Data Cube	9
Data Mart	11
Virtual Warehouse	12

Table of Contents

4. DELIVERY PROCESS.....	13
Delivery Method	13
IT Strategy.....	14
Business Case	14
Education and Prototyping	14
Business Requirements	14
Technical Blueprint	15
Building the Version.....	15
History Load	15
Ad hoc Query	16
Automation	16
Extending Scope.....	16
Requirements Evolution.....	16
5. SYSTEM PROCESSES	18
Process Flow in Data Warehouse	18
Extract and Load Process	18
Clean and Transform Process	19
Backup and Archive the Data	20
Query Management Process	20
6. ARCHITECTURE	21
Business Analysis Framework	21
Three-Tier Data Warehouse Architecture.....	21
Data Warehouse Models.....	22
Load Manager	24
Warehouse Manager	25

Table of Contents

Query Manager	26
Detailed Information	27
Summary Information.....	28
7. OLAP	29
Types of OLAP Servers.....	29
Relational OLAP	29
Multidimensional OLAP	29
Hybrid OLAP.....	29
Specialized SQL Servers.....	30
OLAP Operations.....	30
OLAP vs OLTP	35
8. RELATIONAL OLAP.....	37
Relational OLAP Architecture.....	37
9. MULTIDIMENSIONAL OLAP	39
MOLAP Architecture	39
MOLAP vs ROLAP	40
10. SCHEMAS	41
Star Schema	41
Snowflake Schema	42
Fact Constellation Schema	42
Schema Definition.....	43
11. PARTITIONING STRATEGY	46
Why is it Necessary to Partition?	46
Horizontal Partitioning.....	47
Vertical Partition.....	49
Identify Key to Partition.....	51

Table of Contents

12. METADATA CONCEPTS	52
What is Metadata?	52
Categories of Metadata	52
Role of Metadata	53
Metadata Respiratory.....	54
Challenges for Metadata Management.....	55
13. DATA MARTING	56
Why Do We Need a Data Mart?	56
Cost-effective Data Marting.....	56
Designing Data Marts.....	58
Cost of Data Marting.....	59
14. SYSTEM MANAGERS.....	61
System Configuration Manager.....	61
System Scheduling Manager	61
System Event Manager	62
System and Database Manager.....	63
System Backup Recovery Manager	64
15. PROCESS MANAGERS.....	65
Data Warehouse Load Manager.....	65
Warehouse Manager	66
Query Manager.....	67
16. INTERVIEW QUESTIONS.....	69

1. OVERVIEW

The term "Data Warehouse" was first coined by Bill Inmon in 1990. According to Inmon, a data warehouse is a subject-oriented, integrated, time-variant, and non-volatile collection of data. This data helps analysts to take informed decisions in an organization.

An operational database undergoes frequent changes on a daily basis on account of the transactions that take place. Suppose a business executive wants to analyze previous feedback on any data such as a product, a supplier, or any consumer data, then the executive will have no data available to analyze because the previous data has been updated due to transactions.

A data warehouse provides us generalized and consolidated data in a multidimensional view. Along with a generalized and consolidated view of data, a data warehouse also provides us Online Analytical Processing (OLAP) tools. These tools help us in interactive and effective analysis of data in a multidimensional space. This analysis results in data generalization and data mining.

Data mining functions such as association, clustering, classification, prediction can be integrated with OLAP operations to enhance the interactive mining of knowledge at multiple levels of abstraction. That's why a data warehouse has now become an important platform for data analysis and online analytical processing.

Understanding a Data Warehouse

- A data warehouse is a database, which is kept separate from the organization's operational database.
- There is no frequent updating done in a data warehouse.
- It possesses consolidated historical data, which helps the organization to analyze its business.
- A data warehouse helps executives to organize, understand, and use their data to take strategic decisions.
- Data warehouse systems help in the integration of diverse application systems.
- A data warehouse system helps in consolidated historical data analysis.

Why a Data Warehouse is Separated from Operational Databases

A data warehouse is kept separate from operational databases due to the following reasons:

- ☐ An operational database is constructed for well-known tasks and workloads such as searching particular records, indexing, etc. In contrast, data warehouse queries are often complex and they present a general form of data.
- ☐ Operational databases support concurrent processing of multiple transactions. Concurrency control and recovery mechanisms are required for operational databases to ensure robustness and consistency of the database.
- ☐ An operational database query allows to read and modify operations, while an OLAP query needs only **read only** access of stored data.
- ☐ An operational database maintains current data. On the other hand, a data warehouse maintains historical data.

Data Warehouse Features

The key features of a data warehouse are discussed below:

- ☐ **Subject Oriented** - A data warehouse is subject oriented because it provides information around a subject rather than the organization's ongoing operations. These subjects can be product, customers, suppliers, sales, revenue, etc. A data warehouse does not focus on the ongoing operations, rather it focuses on modelling and analysis of data for decision making.
- **Integrated** - A data warehouse is constructed by integrating data from heterogeneous sources such as relational databases, flat files, etc. This integration enhances the effective analysis of data.
- ☐ **Time Variant** - The data collected in a data warehouse is identified with a particular time period. The data in a data warehouse provides information from the historical point of view.
- ☐ **Non-volatile** - Non-volatile means the previous data is not erased when new data is added to it. A data warehouse is kept separate from the operational database and therefore frequent changes in operational database is not reflected in the data warehouse.

Note: A data warehouse does not require transaction processing, recovery, and concurrency controls, because it is physically stored and separate from the operational database.

Data Warehouse Applications

As discussed before, a data warehouse helps business executives to organize, analyze, and use their data for decision making. A data warehouse serves as a sole part of a plan-execute-assess "closed-loop" feedback system for the enterprise management. Data warehouses are widely used in the following fields:

- ☐ Financial services
- ☐ Banking services
- ☐ Consumer goods
- ☐ Retail sectors
- ☐ Controlled manufacturing

Types of Data Warehouse

Information processing, analytical processing, and data mining are the three types of data warehouse applications that are discussed below:

- **Information Processing** – A data warehouse allows to process the data stored in it. The data can be processed by means of querying, basic statistical analysis, reporting using crosstabs, tables, charts, or graphs.
- **Analytical Processing** – A data warehouse supports analytical processing of the information stored in it. The data can be analyzed by means of basic OLAP operations, including slice-and-dice, drill down, drill up, and pivoting.
- ☐ **Data Mining** - Data mining supports knowledge discovery by finding hidden patterns and associations, constructing analytical models, performing classification and prediction. These mining results can be presented using visualization tools.

Data Warehouse (OLAP)	Operational Database(OLTP)
It involves historical processing of information.	It involves day-to-day processing.
OLAP systems are used by knowledge workers such as executives, managers, and analysts.	OLTP systems are used by clerks, DBAs, or database professionals.

It is used to analyze the business.	It is used to run the business.
It focuses on Information out.	It focuses on Data in.
It is based on Star Schema, Snowflake Schema, and Fact Constellation Schema.	It is based on Entity Relationship Model.
It focuses on Information out.	It is application oriented.
It contains historical data.	It contains current data.
It provides summarized and consolidated data.	It provides primitive and highly detailed data.
It provides summarized and multidimensional view of data.	It provides detailed and flat relational view of data.
The number of users is in hundreds.	The number of users is in thousands.
The number of records accessed is in millions.	The number of records accessed is in tens.
The database size is from 100GB to 100 TB.	The database size is from 100 MB to 100 GB.
These are highly flexible.	It provides high performance.