



# Data Warehouse and Data Mining

Concepts and Techniques

Nitu Sharma

***Dedicated to:***

My Nation

&

My Parents

*(Mr. Rakesh Kr. Sharma & Ms. Shashi Sharma)*





## **GLOBAL VISION PUBLISHING HOUSE**

F-4, 1st Floor, 'Hari Sadan' 20, Ansari Road,

Daryaganj, New Delhi-110002 (INDIA)

Tel.: 23261581, 23276291, 43037885 Mob: 9810644769

Email: nsingh\_2004@vsnl.net, info@globalvisionpub.com

Website: www.globalvisionpub.com

*Data Warehouse and Data Mining: Concepts and Techniques*

© Author

First Edition 2013

ISBN: 978-81-8220-592-5

[Responsibility for the facts stated, opinions expressed, conclusions reached and plagiarism, if any, in this title is entirely that of the author. The publisher bears no responsibility for them whatsoever. All rights reserved. No part of this book may be reproduced in any manner without written permission.]

**PRINTED IN INDIA**

---

*Published by* Dr. N.K. Singh for Global Vision Publishing House, New Delhi-2 and *Printed at* G. S. Offset, Naveen Shahdara, Delhi-32.



# **Data Warehouse and Data Mining**

*Concepts and Techniques*



# CONTENTS

---

## PREFACE

(vii)

## 1. HISTORY OF DATA WAREHOUSING AND MINING

1

### *Chapter Overview*

- 1.1. Introduction
- 1.2. History of Data Warehousing
- 1.3. Functionality of Data Warehouses
- 1.4. Databases vs. Data Warehouses
- 1.5. Dimensional vs. Normalized Approach for Storage of Data
- 1.6. Evolution in Organization use
- 1.7. Top-Down Versus Bottom-up Design Methodologies
- 1.8. Data Warehouses vs. Operational Systems
- 1.9. Different Types of Mining
- 1.10. Summary
- 1.11. Short Type Questions with Answers
- 1.12. Exercise

## 2. DECISION SUPPORT SYSTEM AND DATA WAREHOUSING

13

### *Chapter Overview*

- 2.1. Introduction
- 2.2. Decision Support System (DSS)
- 2.3. Operational Data Store (ODS)
- 2.4. Operational vs. Informational System
- 2.5. Executive Information System (EIS)



- 2.6. Data Warehousing
- 2.7. Comparing Operational Database to Data Warehouse
- 2.8. Characteristics of Data Warehouse
- 2.9. Need for Data Warehousing
- 2.10. Purpose of Data Warehouse
- 2.11. Major Elements and External Entities of Data Warehouse
- 2.12. Information Extracting from a Data Warehouse
- 2.13. Summary
- 2.14. Short Type Questions with Answers
- 2.15. Exercise

### **3. DATA WAREHOUSE ARCHITECTURE**

27

#### *Chapter Overview*

- 3.1. Introduction
- 3.2. Architecture of Data Warehouse
- 3.3. Benefits of Data Warehousing
- 3.4. Data Marts
- 3.5. Reasons for Creating a Data Mart
- 3.6. Data Mart vs. Data Warehouse
- 3.7. Client Server Computing Model and Data Warehousing
- 3.8. Processing Approaches in Client/Server Computing
- 3.9. Generation of Client/Server Computing
- 3.10. Difference between Datawarehouse and OLTP
- 3.11. Summary
- 3.12. Short Type Questions with Answers
- 3.13. Exercise

### **4. MULTIPROCESSOR SYSTEM AND DISTRIBUTED MEMORY ARCHITECTURE**

43

#### *Chapter Overview*

- 4.1. Introduction
- 4.2. Rise Versus Cisc
- 4.3. Multiprocessor System
- 4.4. Distributed Memory Architecture
- 4.5. Cluster System
- 4.6. RDBMS Architecture for Scalability

4.7. Parallel Relational DBMS Processing	
4.8. Advanced RDBMS Features	
4.9. RDBMS Reliability and Availability	
4.10. RDMS Administration	
4.11. Data Clustering	
4.12. Data Base Architectures of Parallel Processing	
4.13. Parallel DBMS Features	
4.14. Parallel DBMS Vendors	
4.15. Summary	
4.16. Short Type Questions with Answers	
4.17. Exercise	
<b>5. METADATA AND DBMS SCHEMA</b>	<b>64</b>
<i>Chapter Overview</i>	
5.1. Introduction	
5.2. Metadata	
5.3. DBMS Schemas for Decision Support	
5.4. Data Cardinality in Data Warehouse	
5.5. Column Local Storage	
5.6. Data Extraction, Cleanup and Transformation Tools	
5.7. Multirelational Database	
5.8. The Multidimensional Data Model	
5.7. Summary	
5.10. Short Type Questions with Answers	
5.11. Exercise	
<b>6. VARIOUS CONSIDERATION IN DATA WAREHOUSE</b>	<b>80</b>
<i>Chapter Overview</i>	
6.1. Introduction	
6.2. Components of Data Warehousing	
6.3. Why Organization Consider Data Warehouse?	
6.4. Mapping the Data Warehouse Architecture to Multiprocessor Architecture	
6.5. Summary	
6.6. Short Type Questions with Answers	
6.7. Exercise	

**7. OLAP, REPORTING AND QUERY TOOLS AND ALLICATION****91***Chapter Overview*

- 7.1. Introduction
- 7.2. OLAP
- 7.3. Reporting and Query Tools
- 7.4. Categories of OLAP Tools
- 7.5. OLAP Guidelines
- 7.6. OLTP vs. OLAP
- 7.7. OLAP Tools and The Internet
- 7.8. Summary
- 7.9. Short Type Questions with Answers
- 7.10. Exercise

**8. PATTERN, MODEL, STATISTICS, HISTOGRAM AND HYPOTHESIS TESTING****105***Chapter Overview*

- 8.1. Introduction
- 8.2. Pattern
- 8.3. Model
- 8.4. Difference Between Pattern and Model
- 8.5. Missing Data
- 8.6. Statistics
- 8.7. Histograms
- 8.8. Types of Categorical Predictors
- 8.9. Baye's Theorem
- 8.10. Bayesian Classification
- 8.11. Hypothesis Testing
- 8.12. Artificial Intelligence
- 8.13. Expert System
- 8.14. Fuzzy Logic
- 8.15. Summary
- 8.16. Short Type Questions with Answers
- 8.17. Exercise

<b>9. DATA-MINING AND ITS TECHNOLOGIES</b>	<b>122</b>
<i>Chapter Overview</i>	
9.1. Introduction	
9.2. Data Mining	
9.3. A Brife History of Data Mining	
9.4. Data Mining Prqcess	
9.5. History of Data Mining	
9.6. Data Mining Models	
9.7. Data Mining Users and Activities	
9.8. Data Mining Functions	
9.9. Data Mining and Business Intelligence	
9.10. Summary	
9.11. Short Type Questions with Answers	
9.12. Exercise	
<b>10. DATA MINING TECHNIQUES, APPLICATIONS AND ITS PROBLEMS</b>	<b>137</b>
<i>Chapter Overview</i>	
10.1. Introduction	
10.2. Data Mining Techniques	
10.3. Decision Tree	
10.4. Cart	
10.5. Data Mining Problems	
10.6. Application of Data Mining	
10.7. Uses of Data Warehousing and Mining	
10.8. Research Challenges in Knowledge Discovery and Data Mining	
10.9. Summary	
10.10. Short Type Questions with Answers	
10.11. Exercise	
<b>11. DATA MINING TECHNIQUES, CLASSIFICATION AND PREDICTION</b>	<b>158</b>
<i>Chapter Overview</i>	
11.1. Introduction	
11.2. Genetic Algorithm	
11.3. Knowledge Discovery	
11.4. Classification and Prediction	

- 11.6. OLAP Operation or Features
- 11.7. OLAP vs. Statistical Data Bases
- 11.8. Cluster Techniques and its Types
- 11.9. K-Means
- 11.10. Summary
- 11.11. Short Type Questions with Answers
- 11.12. Exercise

**12. MULTIMEDIA DATABASE, MINING, AND WEB MINING**

**172**

*Chapter Overview*

- 12.1. Introduction
- 12.2. Multimedia Database
- 12.3. Multimedia Mining
- 12.4. Mining the World Wide Web
- 12.5. Web Data Mining
- 12.6. Data Mining Applications
- 12.7. ID3 ALGO
- 12.8. Summary
- 12.9. Short Type Questions with Answers
- 12.10. Expecise

**BIBLIOGRAPHY**

**181**

**INDEX**

**185**

# 1

## HISTORY OF DATA WAREHOUSING AND MINING

---

---

### CHAPTER OVERVIEW

- 1.1. Introduction
- 1.2. History of Data Warehousing
- 1.3. Functionality of Data Warehouses
- 1.4. Databases vs. Data Warehouses
- 1.5. Dimensional vs. Normalized Approach for Storage of Data
- 1.6. Evolution in Organization use
- 1.7. Top-Down Versus Bottom-up Design Methodologies
- 1.8. Data Warehouses vs. Operational Systems
- 1.9. Different Types of Mining
- 1.10. Summary
- 1.11. Short Type Questions with Answers
- 1.12. Exercise

### 1.1. INTRODUCTION

This chapter deals with the history and function of data warehouse, and discuss about the absence of a data warehousing architecture, and various approach for storage of data like normalized and dimensional approach. Various design methodology like top-down, bottom-up, and hybrid design has also been explained. This chapter also includes the various type of field where we can use different mining like sensor, pattern, in science and engineering etc.

## 1.2. HISTORY OF DATA WAREHOUSING

The concept of data warehousing dates back to the late 1980s when IBM researchers Barry Devlin and Paul Murphy developed the “business data warehouse”. In essence, the data warehousing concept was intended to provide an architectural model for the flow of data from operational systems to decision support environments. The concept attempted to address the various problems associated with this flow—mainly, the high costs associated with it. In the absence of a data warehousing architecture, an enormous amount of redundancy was required to support multiple decision support environments. In larger corporations it was typical for multiple decision support environments to operate independently. Each environment served different users but often required much of the same data. The process of gathering, cleaning and integrating data from various sources, usually long existing operational systems (usually referred to as legacy systems), was typically in part replicated for each environment. Moreover, the operational systems were frequently re-examined as new decision support requirements emerged. Often new requirements necessitated gathering, cleaning and integrating new data from “data marts” that were tailored for ready access by users.

Many large businesses found themselves with data scattered across multiple platforms and variations of technology, making it almost impossible for any one individual to use data from multiple sources. A key idea within data warehousing is to take data from multiple platforms/technologies (As varied as spreadsheets, DB2 databases, IDMS records, and VSAM files) and place them in a common location that uses a common querying tool. In this way operational databases could be held on whatever system was most efficient for the operational business, while the reporting/ strategic information could be held in a common location using a common language.

Key developments in early years of data warehousing were:

- 1960— General Mills and Dartmouth College, in a joint research project, develop the terms dimensions and facts.
- 1970— ACNielsen and IRI provide dimensional data marts for retail sales.
- 1970— Bill Inmon begins to define and discuss the term: Data Warehouse
- 1975— Sperry Univac Introduce MAPPER (Maintain, Prepare, and Produce Executive Reports) is a database management and reporting system that includes the world’s first 4GL. It was the first platform specifically designed for building Information Centers (a forerunner of contemporary Enterprise Data Warehousing platforms)
- 1983— Teradata introduces a database management system specifically designed for decision support.
- 1983— Sperry Corporation Martyn Richard Jones defines the Sperry Information Center approach, which while not being a true DW in the Inmon sense, did contain

many of the characteristics of DW structures and process as defined previously by Inmon, and later by Devlin. First used at the TSB England and Wales

- 1984— Metaphor Computer Systems, founded by David Liddle and Don Massaro, releases Data Interpretation System (DIS). DIS was a hardware/software package and GUI for business users to create a database management and analytic system.
- 1988— Barry Devlin and Paul Murphy publish the article An architecture for a business and information system in IBM Systems Journal where they introduce the term “business data warehouse”.
- 1990— Red Brick Systems, founded by Ralph Kimball, introduces Red Brick Warehouse, a database management system specifically for data warehousing.
- 1991— Prism Solutions, founded by Bill Inmon, introduces Prism Warehouse Manager, software for developing a data warehouse.
- 1992— Bill Inmon publishes the book Building the Data Warehouse.
- 1995— The Data Warehousing Institute, a for-profit organization that promotes data warehousing, is founded.
- 1996— Ralph Kimball publishes the book The Data Warehouse Toolkit.
- 2000— Daniel Linstedt releases the Data Vault, enabling real time auditable Data Warehouses warehouse.

### 1.3. FUNCTIONALITY OF DATA WAREHOUSES

Data warehouses exist to facilitate complex, data-intensive and frequent adhoc queries. Data warehouses must provide far greater and more efficient query support than is demanded of transactional databases. The data warehouse access component supports enhanced spreadsheet functionality, efficient query processing, structured queries, adhoc queries, data mining and materialized views. Particularly enhanced spreadsheet functionality includes support for state-of-the art spreadsheet applications as well as for Online Analytical Processing applications programs. These provide preprogrammed functionalities such as the following:

1. *Roll-up*: Data is summarized with increasing generalization
2. *Drill-down*: Increasing levels of detail are revealed
3. *Pivot*: Cross tabulation that is, rotation, performed
4. *Slice and dice*: Performing projection operations on the dimensions
5. *Sorting*: Data is sorted by ordinal value
6. *Selection*: Data is available by value or range
7. *Derived or computer attributes*: Attributes are computed by operations on stored and derived values.



## 1.4. DATABASES VS. DATA WAREHOUSES

A database is a collection of related data and a database system is a database and database software together. It is also a collection of information as well as a supporting system. Databases are transactional such as relational, object-oriented, network or hierarchical.

Traditional databases support on-line transaction processing (OLTP), which includes insertions, updates, and deletions, while also supporting information query requirements. Traditional databases are optimized to process queries that may touch a small part of the database and transactions that deal with insertions or updates of a few tuples per relation to process.

Thus databases must strike a balance between efficiency in transaction processing and supporting query requirements (ad hoc user requests). That is, they can't further optimized for the applications such as OLAP, Decision Support System and data mining.

But a data warehouse is typically optimized for access from a decision maker's needs. Data warehouses are designed specifically to support efficient extraction, processing and presentation for analytic and decision-making purposes.

In contrast to databases, data warehouses generally contain very large amounts of data from multiple sources that may include databases from different data models and sometimes files acquired from independent systems and platforms.

Multidatabases provide access to disjoint and usually heterogeneous databases and are volatile. Whereas a data warehouse is frequently a store of integrated data from multiple sources, processed for storage in a multidimensional model and nonvolatile. Data warehouses also support time-series and trend analysis, both of which require more historical data.

In transactional systems, transactions are the unit and are the agent of change to the database, but data warehouse information is much more coarse-grained and is refreshed according to a careful choice of incremental refresh policy. Warehouse updates are handled by the warehouse's acquisition component that provides all required processing. As data warehouses encompass large volumes of data, they are more or less double the size of source databases.

The sheer volume of data likely to be in terabytes is an issue that has been dealt with through enterprise-wide data warehouses, virtual data warehouses and data marts. Enterprise-wide data warehouses are huge projects in need of massive investment of time and resources. Virtual data warehouses are bound to provide views of operational databases that are materialized for efficient access.

A data mart is an easy-to-access repository of a subset of highly focused data for a single function or department (*i.e.*, finance, sales, and marketing) and is considerably smaller