

Development and Design Methodologies in DWM

James Yao

Montclair State University, USA

John Wang

Montclair State University, USA

Qiyang Chen

Montclair State University, USA

June Lu

University of Houston – Victoria, USA

INTRODUCTION

Information systems were developed in early 1960s to process orders, billings, inventory controls, payrolls, and accounts payables. Soon information systems research began. Harry Stern started the “Information Systems in Management Science” column in *Management Science* journal to provide a forum for discussion beyond just research papers (Banker & Kauffman, 2004). Ackoff (1967) led the earliest research on management information systems for decision-making purposes and published it in *Management Science*. Gorry and Scott Morton (1971) first used the term *decision support systems* (DSS) in a paper and constructed a framework for improving management information systems. The topics on information systems and DSS research diversifies. One of the major topics has been on how to get systems design right.

As an active component of DSS, data warehousing became one of the most important developments in the information systems field during the mid-to-late 1990s. It has been estimated that about 95% of the *Fortune 1000* companies either have a data warehouse in place or are planning to develop one (Wixon & Watson, 2001). Data warehousing is a product of business need and technological advances. Since business environment has become more global, competitive, complex, and volatile customer relationship management (CRM) and e-commerce initiatives are creating requirements for large, integrated data repositories and advanced analyti-

cal capabilities. By using a data warehouse, companies can make decisions about customer-specific strategies such as customer profiling, customer segmentation, and cross-selling analysis (Cunningham, Song, & Chen, 2006). To analyze these large quantities of data, data mining has been widely used to find hidden patterns in the data and even discover knowledge from the collected data. Thus how to design and develop a data warehouse and how to use data mining in the data warehouse development have become important issues for information systems designers and developers.

This article presents some of the currently discussed development and design methodologies in data warehousing and data mining, such as the multidimensional model vs. relational entity-relationship (ER) model, corporate information factory (CIF) vs. multidimensional methodologies, data-driven vs. metric-driven approaches, top-down vs. bottom-up design approaches, data partitioning and parallel processing, materialized view, data mining, and knowledge discovery in database (KDD).

BACKGROUND

Data warehouse design is a lengthy, time-consuming, and costly process. Any wrongly calculated step can lead to a failure. Therefore, researchers have placed important efforts to the study of design and development related issues and methodologies.

Data modeling for a data warehouse is different from operational database, for example, online transaction processing (OLTP), data modeling. An operational system is a system that is used to run a business in real time, based on current data. An OLTP system usually adopts ER modeling and application-oriented database design (Han & Kamber, 2006). An information system, like a data warehouse, is designed to support decision making based on historical point-in-time and prediction data for complex queries or data mining applications (Hoffer, Prescott, & McFadden, 2007). A data warehouse schema is viewed as a dimensional model (Ahmad, Azhar, & Lukauskis, 2004; Han & Kamber, 2006; Levene & Loizou, 2003). It typically adopts either a star or snowflake schema and a subject-oriented database design (Han & Kamber, 2006). The schema design is the most critical to the design of a data warehouse.

Many approaches and methodologies have been proposed in the design and development of data warehouses. Two major data warehouse design methodologies have been paid more attention. Inmon, Terdeman, and Imhoff (2000) proposed the CIF architecture. This architecture, in the design of the atomic-level data marts, uses denormalized entity-relationship diagram (ERD) schema. Kimball (1996, 1997) proposed multidimensional (MD) architecture. This architecture uses star schema at atomic-level data marts. Which architecture should an enterprise follow? Is one better than the other? Currently, the most popular data model for data warehouse design is the dimensional model (Bellatreche & Mohania, 2006; Han & Kamber, 2006). Some researchers call this model the data-driven design model. Artz (2006) advocates the metric-driven view, which, as another view of data warehouse design, begins by identifying key business processes that need to be measured and tracked over time in order for the organization to function more efficiently. There has always been the issue of top-down vs. bottom-up approaches in the design of information systems. The same is with a data warehouse design. These have been puzzling questions for business intelligent architects and data warehouse designers and developers. The next section will extend the discussion on issues related to data warehouse and mining design and development methodologies.

DESIGN AND DEVELOPMENT METHODOLOGIES

D

Data Warehouse Data Modeling

Database design is typically divided into a four-stage process (Raisinghani, 2000). After requirements are collected, conceptual design, logical design, and physical design follow. Of the four stages, logical design is the key focal point of the database design process and most critical to the design of a database. In terms of an OLTP system design, it usually adopts an ER data model and an application-oriented database design (Han & Kamber, 2006). The majority of modern enterprise information systems are built using the ER model (Raisinghani, 2000). The ER data model is commonly used in relational database design, where a database schema consists of a set of entities and the relationship between them. The ER model is used to demonstrate detailed relationships between the data elements. It focuses on removing redundancy of data elements in the database. The schema is a database design containing the logic and showing relationships between the data organized in different relations (Ahmad et al., 2004). Conversely, a data warehouse requires a concise, subject-oriented schema that facilitates online data analysis. A data warehouse schema is viewed as a dimensional model which is composed of a central fact table and a set of surrounding dimension tables, each corresponding to one of the components or dimensions of the fact table (Levene & Loizou, 2003). Dimensional models are oriented toward a specific business process or subject. This approach keeps the data elements associated with the business process only one join away. The most popular data model for a data warehouse is the multidimensional model. Such a model can exist in the form of a star schema, a snowflake schema, or a starflake schema.

The star schema (see Figure 1) is the simplest data base structure containing a fact table in the center, no redundancy, which is surrounded by a set of smaller dimension tables (Ahmad et al., 2004; Han & Kamber, 2006). The fact table is connected with the dimension tables using many-to-one relationships to ensure their hierarchy. The star schema can provide fast response

7 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/development-design-methodologies-dwm/11261

Related Content

A Method for Ranking Non-Linear Qualitative Decision Preferences using Copulas

Biljana Mileva-Boshkoska and Marko Bohanec (2012). *International Journal of Decision Support System Technology* (pp. 42-58).

www.irma-international.org/article/method-ranking-non-linear-qualitative/69516

A Comparative Study of Two Models for Handling Transportation Cost in Combinatorial Auctions

Fu-Shiung Hsieh (2020). *International Journal of Decision Support System Technology* (pp. 62-84).

www.irma-international.org/article/a-comparative-study-of-two-models-for-handling-transportation-cost-in-combinatorial-auctions/258563

Money Supply: Predictive Analytics in India

Rituparna Das (2014). *Emerging Methods in Predictive Analytics: Risk Management and Decision-Making* (pp. 334-348).

www.irma-international.org/chapter/money-supply/107912

Knowledge Mobilization for Agri-Food Supply Chain Decisions: Identification of Knowledge Boundaries and Categorization of Boundary-Spanning Mechanisms

Guoqing Zhao, Shaofeng Liu, Sebastian Elgueta, Juan Pablo Manzur, Carmen Lopez and Huilan Chen (2023). *International Journal of Decision Support System Technology* (pp. 1-25).

www.irma-international.org/article/knowledge-mobilization-for-agri-food-supply-chain-decisions/315640

In-Store Stimuli and Impulsive Buying Behaviour: Modeling Through Regression Equation

Chandan Parsad, Sanjeev Prashar, T. Sai Vijay and Mukesh Kumar (2018). *International Journal of Strategic Decision Sciences* (pp. 95-112).

www.irma-international.org/article/in-store-stimuli-and-impulsive-buying-behaviour/208681