

Decision Support Databases Essentials

Antonio Albano

University of Pisa
Department of Computer Science

Copyright © 2015 by Antonio Albano

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that the first page of each copy bears this notice and the full citation including title and authors. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission from the copyright owner.

February 13, 2015

CONTENTS

Preface	iii
1 Decision Support Systems	1
1.1 Information Systems	1
1.2 Types of Information Systems	2
1.3 Data Warehouse: a Decision Support Database	4
1.4 Data Warehousing Architecture	7
1.5 What to Model	9
1.6 Concluding Remarks	13
1.7 Summary	14
2 Data Warehouse Modeling	15
2.1 Conceptual Multidimensional Model	15
2.2 Multidimensional Relational Model	26
2.3 Multidimensional Cube Model	29
2.4 Summary	34
3 Data Warehouse Design	35
3.1 Introduction	35
3.2 Data Warehouse Design Approaches	36
3.3 A Case Study	52
3.4 Project Quality Control	61
3.5 Summary	65
4 Data Analysis	67
4.1 OLAP Systems Solutions	67
4.2 Data Analysis Using SQL	69
4.3 Simple Reports with SQL	70
4.4 Moderately Difficult Reports with SQL	77
4.5 Very Difficult Reports Without Analytic SQL	83
4.6 Summary	97
A Case Studies	99
A.1 Hospital	99
A.2 Airline Companies	101
A.3 Airline Flights	102
A.4 Inventory	104

A.5 Hotels	107
B Case Studies: Solutions	109
B.1 Hospital	109
B.2 Airline Companies	112
B.3 Airline Flights	114
B.4 Inventory	119
B.5 Hotels	122
C Glossary	125
Bibliography	131

PREFACE

In our information-based society one of the most important applications for computers is *information storage and management* with a DBMS to support organizations both in *performing the business* and in *bringing the business* to the Web to allow new routes to market.

It is also well known that information overload is a huge challenge for businesses, but it is also an enormous opportunity in making smarter decisions based mainly on data analysis to improve productivity. Consequently, starting from the 1990s another important application of *information storage and management* to support organizations is *analyzing the business* with data-driven *Decision Support Systems* designed for summarizing large amounts of data into a form that is useful and easily interpretable to help managers to analyze the performance of *key business processes, worthy of improvement*.

Decision support applications involve quite complex analysis which cannot be efficiently executed against operational databases, optimized for online transaction processing. For this reason, organizations maintain a separate database, called a *data warehouse*, which is specifically organized for such complex analysis. The term *data warehouse* is a metaphor: a warehouse is a large structure where things are stored and organized for easy accessibility. However, a *data warehouse* is not only a large repository for historical data extracted from operational systems, but is organized to create the right models for measurable key business processes, to support informed decisions about how to improve them.

Organization

Chapter 1 provides a general overview of the purpose of *decision support systems*, and of the concepts of *data warehouse* and of *data warehousing process*. We also introduce the reason why data warehouses are used to analyze key business processes that are measurable and worthy of improvements, and consequently what is modeled in a data warehouse.

Chapter 2 presents the fundamental concepts about a conceptual model for designing data warehouses, and the logical data model to implement them.

Chapter 3 presents a data warehouse design process and a methodology for the implementation of a logical design schema.

Chapter 4 focuses on the extensions of SQL for online analytic processing, called *Analytic SQL*, the fundamental user-oriented relational languages to analyze data for producing interesting reports to evaluate the performance of the modeled key processes in order to improve them.

Appendix. A set of case studies is presented to apply the concepts presented in the chapters of the book.

A. A.

DECISION SUPPORT SYSTEMS

An overview of *decision support systems* is given below, particularly those data-driven, used to synthesize large amounts of data into a form that is useful to manage the business. The data are first organized in a special database called *data warehouse* and then analyzed with appropriate techniques, called *On Line Analytical Processing (OLAP)* or with semi-automatic and exploratory techniques, called *data mining*. Finally, the characteristics of systems for managing data warehouse are presented, which information should be represented using an appropriate *data model*, and how data can be used for decision-making.

1.1 Information Systems

Organizations have used information systems for centuries and they have used a variety of technologies to deal with information.

■ Definition 1.1

An *information system* is an organized collection of resources, people, and procedures finalized to collect, store, process and communicate the information needed to support the on-going activities.

Nowadays, information is considered to be a critical resource of any organization, as fundamental as capital or machinery, and, in fact, the majority of the labor force in the industrialized countries works in some way with information.

Information can be represented as data, images, text, and voice. Clearly, different types of organizations will have differing needs with respect to the kinds of information they use. However, the attention here will be on information represented as *structured data* shared by a variety of users within an organization, and managed using computers. Reductions in the costs of computer technology, improvements in performances, and new facilities to support the development of applications have created an increasing demand for data processing systems. We use the term *computerized information system* to refer to the hardware and software which is used for storing, retrieving, and processing the information which supports the functions of an organization.¹ In the following, in brief, we use the term *information system* for *computerized information system*.

1. Frequently in the literature, *information system* is used as synonym of *computerized information*

Over time, there are continuous interactions between the two components of the information system and the rest of the organization that will change each other, and this requires attention of management to plan the evolution of both the organizational structure and the employee tasks (Figure 1.1)

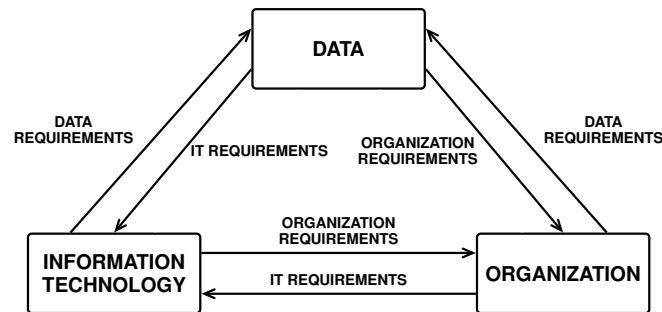


Figure 1.1: System Conception of an Information Systems

For example, new requirements of the organization may include the need of new categories of data being managed by the information system, and to adapt the information technology to provide new services (e.g, think of an organization that decides to offer web services). New categories of data to be managed can result in (a) a review of the organizational structure to review, for example, the tasks and professional employees, (b) an adjustment of the information technology used. The evolution of information technology can enable new opportunities for data management and set new requirements for employees and the organizational structure, and so on. The data, the organization, and the technology all interact and change each other.

1.2 Types of Information Systems

Information systems can be classified in several ways, but for our purposes it is useful to classify them in the following categories on the basis of the business activities that are required to support.

■ Definition 1.2 *A Taxonomy of Information Systems*

Information systems can be classified into the following categories:

- *Operational*, to *perform* the business operational activities.
- *Decision support*, to provide the information that managers need for *analyzing the business*.
- *Web-based for E-commerce*, to *bring the business to the Web* to allow new routes to market.

These information system categories are all ongoing and in a constant state of improvement. They use different technology, have different objectives and require different skills to develop. In the following the attention will be on the first two categories, and in particular on the *decision support* one, a key driver in the business world today (Figure 1.2).

system. Here, we prefer to make a distinction between the two terms to evidence the fact that a *computerized information system* will never completely substitute the *global information system* of an organization.

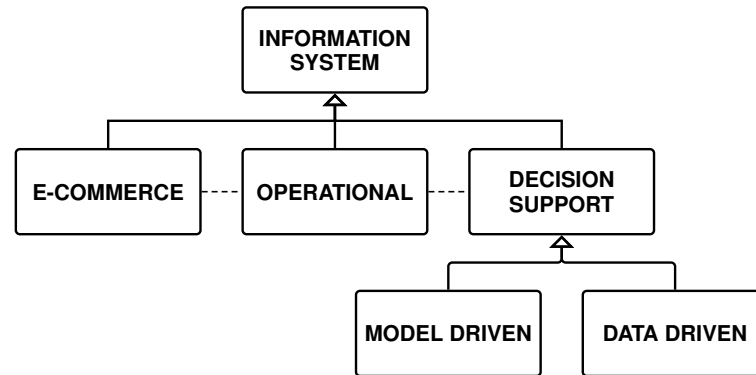


Figure 1.2: Types of Information Systems

1.2.1 Operational Systems

When an information system is implemented using the database technology, it will consist of an *operational database* and a collection of *application programs (transactions)* which are used to access and update the data quickly and efficiently (Figure 1.3). The main goal of such a transaction processing system is to maintain the correspondence between the database and the real-world situation it is modeling, as events occur in the real world.

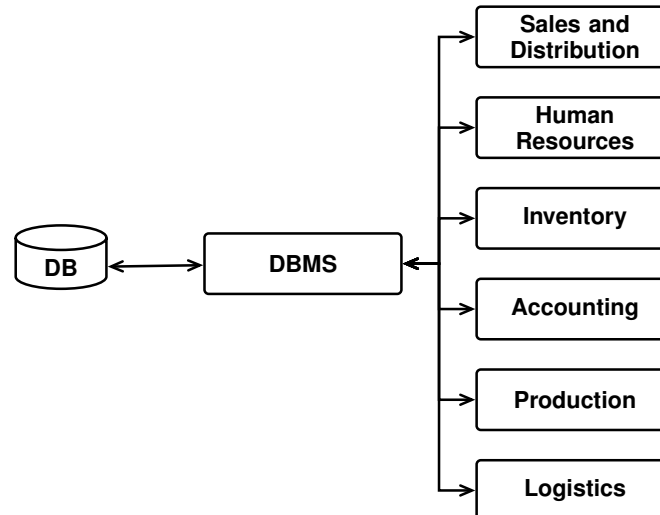


Figure 1.3: Transaction processing system

The data are under the control of a *Data Base Management System (DBMS)*, a centralized or distributed software system, which provides the tools to define the database, to select the data structures needed to store and retrieve the data easily, and to access the data, interactively or by means of a programming language.

1.2.2 Decision Support Systems

Decision support information systems can be classified in *Management Information Systems* and *Decision Support Systems*.

The first one is used by middle tactical or administrative managers in monitoring and controlling their units to correct problems by making decisions based on comparing the actual performance and the planned performance (*variance report*). Decision support systems are used to make strategic decisions about the future directions of the business enterprise, using both historical internal data and external data.

For brevity, in the following we will use the term *Decision Support Systems (DSS)* for both types of decision support information systems.

DSS have been introduced in the organizations since the late '70s to help managers to make decisions of three types:

- *Structured*, when a well-defined decision-making procedure exists.
- *Unstructured*, when a well-defined decision-making procedure does not exist and the experience and creativity of the manager are required.
- *Semistructured*, when the decision-making procedure is partly defined and so it is also required the manager's creative intervention.

There is no strict correspondence between types of decision and levels of decision-making processes, however, at the operational level decisions tend to be more structured, at the tactical level decisions are mainly semistructured and at the strategic level decisions are typically unstructured.

The DSS have very different characteristics, but it is useful to classify them into two main types: *model-driven*, to take structured or semistructured decisions, or *data-driven* to take unstructured decisions.

The model-driven DSS are an evolution of the first proposals made at the end of the 70s for decision support systems and their value depends on the quality of the model used. The simplest solutions utilize spreadsheets for analysis of "what if", while more sophisticated models are used from operations research, simulation and artificial intelligence.

The data-driven DSS are designed to synthesize large amounts of data into a form that is useful and easily interpretable to help managers to assess the performance of business processes and make decisions to address and resolve any critical issues found. Their value depends on the type and quality of data generated using synthetic instruments called *Business Intelligence*. The term *intelligence* is used with the meaning of *investigating to find out something interesting*, like in *Intelligence Service*.

The operational data accumulated over time, integrated with those from external sources, are a potential source of information used by managers regardless of their decision-making level in the organization. The information is derived from the data summarized in an appropriate form and its relevance depends on the recipient. When experience, competence, and attitude are added to information, knowledge is created, and actions can be taken. To become actionable, knowledge should also be closely integrated with an organization's business processes.

In the following, data-driven DSS will be considered to see how they can be designed to support informed decisions.

Decision support applications involve quite complex analysis which cannot be efficiently executed against operational databases, optimized for *On Line Transaction Processing (OLTP)*. For this reason, organizations maintain a separate database, called *data warehouse*, which is specifically organized for such complex analysis.

1.3 Data Warehouse: a Decision Support Database

The first and still now the most widely cited definition of data warehouse was provided by William Inmon in 1990:

■ Definition 1.3

A data warehouse is a subject-oriented, integrated, nonvolatile, and time-varying collection of data in support of management's decisions.

Let us examine each of these distinctive aspects of a data warehouse.

1. **Subject-oriented.** A data warehouse stores data by subject, not by applications, which is what distinguishes a data warehouse from an operational database, that stores information in order to optimize transaction processing. Business subjects differ from organization to organization. They are the critical subjects for an organization. For example, for a manufacturing company, these would include, sales, shipments, returns, and inventory.
A data mart is database that has the same characteristics as a data warehouse, but is usually smaller and is focused on the data for one subject.
2. **Integrated.** Data are gathered into the data warehouse from a variety of sources and merged into a coherent whole. For example, a bank can collect different data on customers for the management of loans, current accounts, or stocks, but they must then be integrated for the purposes of the analysis of the services offered to customers.
3. **Time-variant.** For an operational system, the stored data contains the *current* values. On the other hand, the data in the data warehouse is meant for analysis and decision support, and is thus historical data identified with a particular time period. An operational system contains *current* data, while a data warehouse contains historical data over long time for analysis and decision support, therefore a time dimension is explicitly included in data so that trends and changes over time can be analyzed.
4. **Non-volatile.** The data in a data warehouse is primarily for query and analysis, and it is never changed interactively. This enables management to gain a consistent picture of the business. Periodically, new data may be added or those considered obsolete may be removed.
5. **Decision support.** The primary function of the data warehouse is for decision support, and so it must be specifically designed to answer business questions. Data is the reality that a computer records, stores, and processes. The lowest level in the perception of reality is sometimes referred to as “raw data”. This data is of little benefit unless it can be turned into useful information and knowledge. Data must be *condensed into a more informative format* in such a way that managers (or more in general *knowledge workers* – executives, managers, and analysts) can get the essence of the underlying data.

Three categories of decision support can be provided. Specifically:

- (a) **Reports.** Reporting is considered the lowest level of decision support. A reporting facility capable of generating informative reports for managers in time to be useful is of the utmost importance for the successful operation of any business.
- (b) **Multidimensional data analysis**, sometimes called *On Line Analytic Processing* (OLAP). Data analysis is usually accomplished interactively with some kind of data analysis tool. The goal of data analysis is to get useful information from the data.
- (c) **Exploratory data analysis.** This data analysis technique is very different from reports and multidimensional analysis: it uses what is called a *discovery technique of useful data models* with *data mining* algorithms. That is, the

user does not ask a particular question about data, but rather he uses specific algorithms that analyze the data and report what they have discovered. Unlike reports and multidimensional analysis, where the user has to create and execute queries based on hypotheses, data mining algorithms search for answers. A comparison of the two approaches is shown in Table 1.1 with some example queries. Data mining algorithms are beyond the scope of this book.

Table 1.1: Comparison between OLAP and Exploratory data analysis

OLAP Query	Exploratory data analysis
Which customers spent most with us in the past year?	Which types of customer are likely to spend most with us in the coming year?
How much did the bank lose from loan defaulters in the past two years?	What are the characteristics of the customers most likely to default on their loans before the year is out?
What were the highest selling fashion items in our London stores?	What additional products are most likely to be sold to customers who buy sportswear?

A data warehouse is usually separated from an operational database for the following reasons:

- **Performance.** Special data organization, access and implementation methods are needed to support multidimensional views and data analysis which usually requires complex queries that would degrade the performance of operational transactions. Moreover, concurrency control and recovery DBMS modes are not compatible with data analysis.
- **Function.** Decision support requires (a) historical data, which operational databases do not typically maintain, (b) consolidation (aggregation, summarization) of data from heterogeneous sources, such as operational databases, external sources, and (c) different sources typically use inconsistent data representations, codes and formats which have to be reconciled to enforce data quality.

Table 1.2 summarizes the differences between the traditional applications that use databases (*On Line Transaction Processing, OLTP*), and the decision support applications that use data warehouses (*On Line Analytical Processing, OLAP*).

OLAP is a term that was coined in an unpublished 1993 white paper, “Providing OLAP to User Analysts: An IT Mandate”, by E. F. Codd. By introducing this new term as a play on the then-familiar term on-line transaction processing (OLTP), the paper signaled a shift in the paradigm for business analysis, in parallel with the shift that had already occurred for transaction processing. Instead of reviewing piles of static reports printed on green-bar paper, the OLAP analyst could explore business results interactively, dynamically adjusting the view of the data – asking questions and getting answers almost immediately. This freedom from static answers to fixed questions on a fixed schedule allows business analysts to operate more effectively and to effect improvements in business operations. In the white paper, the authors outlined 12 characteristics of an OLAP system. In a 1995 update to the white paper, six more characteristics were added.

In the 2004, Nigel Pendse, an analyst with Business Intelligence Ltd. who publishes The OLAP Report, provides another valuable point of view. In a Web page entitled “What Is OLAP”, Pendse introduces a simpler model, FASMI (*Fast Analysis*

Table 1.2: Comparison between OLTP and OLAP

	OLTP	OLAP
Function	Operational processing	Decision support
Users	Clerk, IT professional	Knowledge worker
DB design	Application-oriented	Subject-oriented
Usage	90% repetitive	90% ad-hoc
Data	Current, detailed, relational	Historical, summarized, multidimensional, integrated
Access	Read/write	Complex read query
No of users	A lot	Few
DB size	100 MB to GB	100 GB to TB
Orientation	Transactions	Analysis

of *Shared Multidimensional Information*), to characterize OLAP systems. Although no single definition is likely to receive universal support, Pendse's characterization is much simpler than the Codd rules. Briefly, the FASMI characteristics are:

Fast. In keeping with the spirit of the “O” in OLAP, such systems need to provide results very quickly – usually in just a few seconds, and seldom in more than 20 or 30 seconds. This level of performance is key in allowing analysts to work effectively without distraction.

Analytic. Considering the “A” in OLAP, such systems generally must provide rich analytic functions appropriate to a given application, with minimal programming.

Shared. An OLAP system is usually a shared resource. This means that there is a requirement for OLAP systems to provide appropriate security and integrity features. Ultimately, this can mean providing different access controls on each cell of a database.

Multidimensional. Multidimensionality is the primary requirement for an OLAP system, which must present the data in a multidimensional framework. The motivation for this requirement will be discussed later on.

Information. OLAP systems must allow the user to easily condense large amount of data into a form that is useful to business manager and decision makers.

1.4 Data Warehousing Architecture

The term *data warehousing* is used to refer to the process used to organize data in a data warehouse and then allow end users to analyze them with business intelligence applications. In practice three types of solutions are adopted, depending on the number of data layers employed.

One-Layer Architecture. This solution has only one layer of data handled by the operational system, and the data warehouse is virtual, i.e. it is defined as a view of operational data, possibly materialized), and it is used by the business intelligence applications (Figure 1.4). This solution does not require a specific system for managing the data warehouse and it is usually used as the first low-cost solution for small organizations, but it does not meet the requirement for separation between operational and analytical applications.

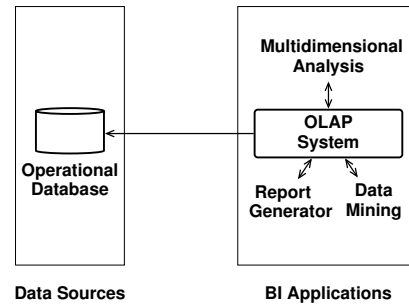


Figure 1.4: One-layer architecture

Two-Layer Architecture. This solution is more general than the previous one, because a data warehouse exists separated from the operational database and managed by a specific system. The data warehouse is loaded with data extracted with *Extract, Transform and Load (ETL)* applications from the operational database, and any other structured data sources, to bring them to a consistent form (Figure 1.5). While the data sources are updated continuously by operational applications, the data warehouse is updated periodically with the ETL applications. This situation typically arises when there are high quality operational databases with schemas sufficiently similar to that of the data warehouse.

This solution separates

- the system for operational database management from the system for data warehouse management and decision support,
- the operational applications from the business intelligence applications, so that business analysis would not interfere with and degrade the performance of operational applications.

Metadata is information about the structure, content and interdependencies of data warehouse components, to support developers, administrators responsible for the data warehouse and the business intelligence applications.

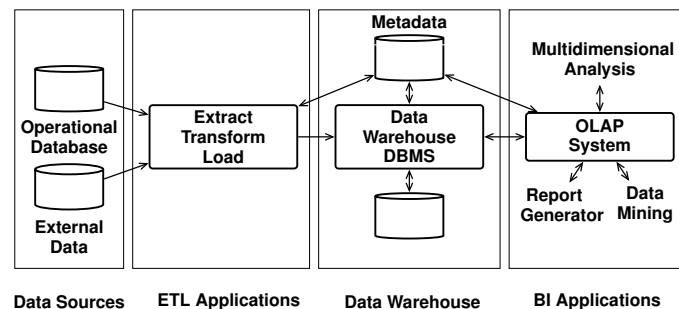


Figure 1.5: Two-layer architecture

Three-Layer Architecture. This solution is the most general with three data layers: the *data sources*, the *data staging* and the *data warehouse*. The data staging contains data obtained from the integration of different data sources and prepared for loading into the data warehouse (Figure 1.6). The data staging may just be a set of files or, at other extreme, a fully developed relational database. The complexity of data staging layer depends on the quality of the data sources.

This solution separates the process of extraction and integration of data sources from the process of data reorganization and loading into the data warehouse.

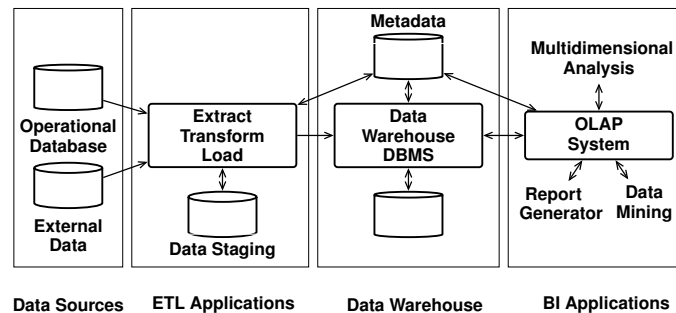


Figure 1.6: Three-layer architecture

1.5 What to Model

According to [Artz, 2005], to support managers in decision-making, data must be organized taking into consideration how they use such data to support their decisions about the performance of key business processes.

■ Definition 1.4

A database is designed to represent some aspects of a reality in terms of the information available about collections of entities with properties and relationship sets between them, while a data warehouse is a *specialized database* designed to represent some aspects of *key business processes* in terms of collections of facts about the interesting process measurements, that represent how the processes are being performed, and a set of dimensions, which provide the context of the facts, to be used for analyzing the process performances.

Let us describe more precisely what to model to help managers in analyzing a business process.

*Managers are interested in analyzing collections of **facts** about the performance of a key business process, measurable and worthy of improvement.*

A fact, in this context, is represented by a set of numerical attributes (hereafter *measure*) by which the process performance is tracked and measured in order to maintain or improve their efficiency. In data warehousing terminology, the interval at which we take measurements is called the *grain*.

Examples of measures for a sale of a product are Quantity, Price, and Revenue. However, without some context, the measures are useless.

*Managers think in terms of **business dimensions**, which give facts their context, and are used to analyze them to evaluate their effects.*

For sales data, the dimensions could include Product, Date, and Store. Dimensions contain the *descriptions* of the subject being measured. Examples of questions managers use to ask for decision-making are: “Show me the total sale revenue by product, year, and store”, “Show me the current and previous year-to-date sales revenue, and percentage change, by product and by store”.

*Managers analyze measurable business process performance using summary data (called **metrics**) obtained by grouping facts by different dimensions and combinations of dimensions, and then aggregating measures into useful forms.*

Common metrics are about economic and financial indicators, but when they are about efficiency and quality of process, are called **Key Performance Indicators, KPI**, because they help understanding how a business process is doing against an objective.

*Managers are interested in analyzing metrics in various levels of details, by exploiting the fact that some dimensions have a set of associated attributes that can be structured as a **hierarchy**.*

A date dimension, for example, with attributes Day, Month, Quarter and Year, could have a hierarchy Day < Month < Quarter < Year, with the meaning that Year is the highest level of generality within the hierarchy, the second level Quarter tells us that more than one quarter is contained in an year, and so on. The combination of a multidimensional and a hierarchical view allows managers to get a good deal of information from data analysis. For example, managers first see the total sales revenue for the entire year by product, then they move down to quarters to look at the sales by quarter and product.

Example 1.1

Let us consider the sales data stored in the relational table Sales(Product, Store, Date, Quantity), where Quantity is the measure and the other attributes are the dimensions that describe a sale fact. A data analysis usually does a *dimensionality reduction (grouping)* to partition a set of rows whose membership is characterized by the fact that all of the rows in a single group agree on the values of the dimensions that are left out. Each group is then aggregated by a function to compute a metric from the measure values. By *aggregation* means to compute a single value from a list of values using an aggregate function such as SUM, COUNT, MIN, MAX, AVG.

Let us look at some examples of an interactive multidimensional data analysis concerning the total quantity of products sold (the metric) to be analyzed by a subset of the dimensions Product, Store, Date. The point is that the user begins with a business question to which wants to answer with the data, gets the results, analyzes the results, uses this new information to formulate another business question, and so on. Later on we will see how to express business questions in SQL to produce the results.

1. The total sales quantity by product, to determine which product is sold best

Product	Total Sales Qty
P1	27 407
P2	5 179
P3	3 446

2. The total sales quantity by product and by store, to determine where it is best to sell certain products

Product	Store	Total Sales Qty
P1	S1	13 945
	S2	9 875
	S3	3 587
P2	S1	1 950
	S2	2 500
	S3	729
P3	S1	1 000
	S2	1 200
	S3	1 246

3. A common type of analysis is a generalization of the former: we want to aggregate the measures on some dimensions and also provide the subtotals for each value of all dimensions. This analysis produces a report such as the one in which shows the total sales quantity by product and by store, extended with subtotals for products, for stores, and with the overall total.

Product	Store	Total Sales Qty
P1	S1	13 945
	S2	9 875
	S3	3 587
P1	Total	27 407
P2	S1	1 950
	S2	2 500
	S3	729
P2	Total	5 179
P3	S1	1 000
	S2	1 200
	S3	1 246
P3	Total	3 446
Total		36 032

4. Starting with the results of a previous analysis, we can proceed to a more detailed one. For example, after a look at the percentage change of annual quantity sales of products we can also do an analysis by store to understand the decrease in sales of the product 'P2'.

Product	Total Sales Qty 2009	Total Sales Qty 2010	Change (%)
P1	12 845	14 562	13
P2	2 753	2 426	-12
P3	1 567	1 879	20
Total	17 165	18 867	10

Product	Store	Total Sales Qty 2009	Total Sales Qty 2010	Change (%)
P1	S1	6 445	7 500	16
	S2	4 225	5 650	34
	S3	2 175	1 412	-35
P1	Total	12 845	14 562	13
P2	S1	900	1 050	17
	S2	1 200	1 300	8
	S3	653	76	-88
P2	Total	2 753	2 426	-12
P3	S1	450	550	22
	S2	580	620	7
	S3	537	709	32
P3	Total	1 567	1 879	20
Total		17 165	18 867	10

The reports for decision support are usually represented in a very different form from the ones shown above. They are much more visually pleasing and intuitive, like the dashboard of a vehicle, using graphics and color-coded alarms to highlight trends, exceptions or values lower than predefined ones (Figure 1.7). Microstrategy, a very active company in the Business Intelligence arena, has some interesting examples on <http://www.microstrategy8.com>.

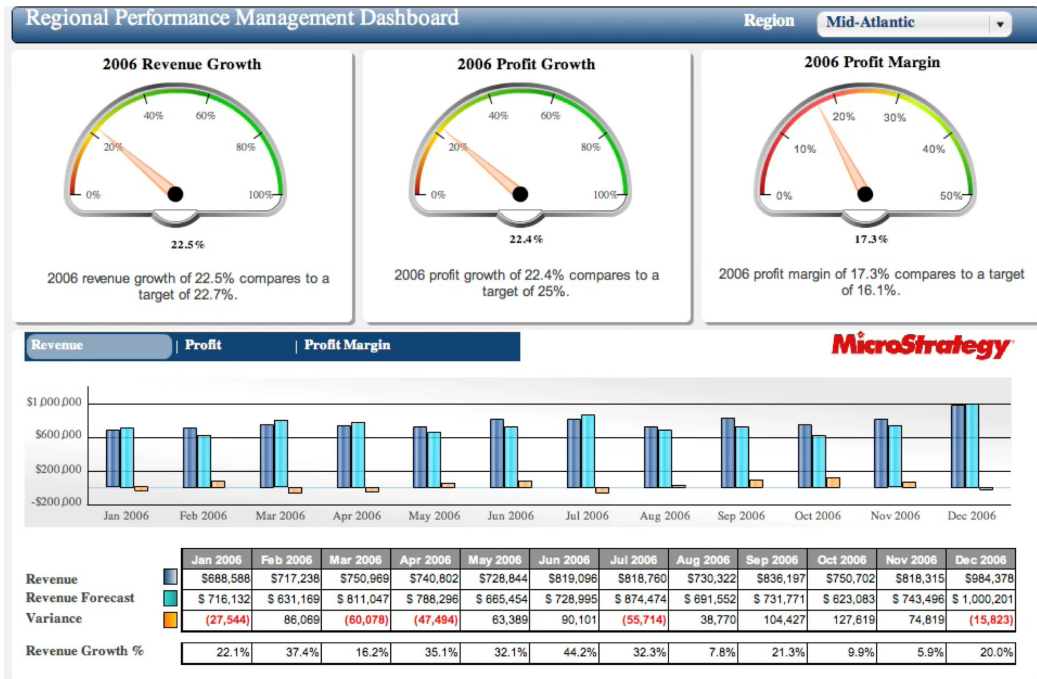


Figure 1.7: Example of Scorecard & Dashboard

1.6 Concluding Remarks

Decision support systems, designed to synthesize, with business intelligence tools, large amounts of data in ways useful to make more rapid and objective decision making, had a growing popularity in recent years for their value strategic and competitive. There has been three very interesting analysis of this trend:

- T. H. Davenport, G. C. Harris, *Competing on Analytics: The New Science of Winning*, Harvard Business School Press, Boston 2007, for the American context.
- Monitoring Business Intelligence, Report 2007-2008, SDA Bocconi, for the Italian context.
- T. Burelli, A. Marzona, M. Pighin, *From intuition to knowledge*, Arachne, Roma, 2007, for the Italian context of small and medium businesses.

It is also interesting to read an article that appeared in print on April 23, 2011 of the *The New York Times* edition with the heading *When There's No Such Thing as Too Much Information*, by Steve Lohr. Here is an excerpt of what he says:

Information overload is a headache for individuals and a huge challenge for businesses. Companies are swimming, if not drowning, in wave after wave of data — from increasingly sophisticated computer tracking of shipments, sales, suppliers and customers, as well as e-mail, Web traffic and social-network comments. These Internet-era technologies, by one estimate, are doubling the quantity of business data every 1.2 years.

Yet the data explosion is also an enormous opportunity. In a modern economy, information should be the prime asset — the raw material of new products and services, smarter decisions, competitive advantage for companies, and greater growth and productivity.

Is there any real evidence of a data payoff across the corporate world? It has taken a while, but new research led by Erik Brynjolfsson, an economist at the Sloan School of Management at the Massachusetts Institute of Technology, suggests that the beginnings are now visible.

Mr. Brynjolfsson and his colleagues, Lorin Hitt, a professor at the Wharton School of the University of Pennsylvania, and Heekyung Kim, a graduate student at M.I.T., studied 179 large companies. Those that adopted data-driven decision making achieved productivity that was 5 to 6 percent higher than could be explained by other factors, including how much the companies invested in technology, the researchers said.

In the study, based on a survey and follow-up interviews, data-driven decision making was defined not only by collecting data, but also by how it is used — or not — in making crucial decisions, like whether to create a new product or service. The central distinction, according to Mr. Brynjolfsson, is between decisions based mainly on data and analysis and on the traditional management arts of experience and intuition.

After having presented *what is modeled* in a data warehouse, in the following the attention will be on:

- *How to model*: which data model is used to model a data warehouse.
- *How data warehouses are designed*: which methodology is used for the design of a data warehouse.

- *How data are analyzed*: which operators are available to analyze data.
- *How to implement a data warehouses system*: which relational DBMS technological extensions are needed to support operations on data warehouses.

1.7 Summary

- An *information system* is a system whose purpose is to store, process, and communicate information.
- The focus of an *operational information system* is the *execution* of business processes, the focus of a *decision support information system* is the *evaluation* of the processes, while the focus of a *web-based information system* is the *use of internet web* to allow new routes to market.
- A data warehouse is a decision support database with historical, nonvolatile data, to facilitate analysis of the performance of key business processes, worthy of improvement.
- Data warehouses and operational databases provide different functions and require different kinds of data, therefore they need to be maintained separately.

DATA WAREHOUSE MODELING

The purpose of a data warehouse is not *just to store data* but rather to *facilitate decision making*. As such, the first step is to model a data warehouse on the basis of the relevant types of business analyses. Data warehouse modeling is a process that produces a well-organized abstract dimensional data model to understand the structure and contents of the data to best support the needs of the business users. In the following sections, three examples of data models are presented that are relevant in dimensional modeling, using the basic concepts of facts, measures, dimensions and hierarchies:

- A *conceptual multidimensional model*, useful to reason about the characteristics of data at a conceptual level, independent of implementation concerns, as it happens with the Entity-Relationship model for databases.
- A *multidimensional relational model*, the traditional logical model to represent data in data warehouse systems.
- A *multidimensional cube model*, useful to show the basic operators for data analysis.

2.1 Conceptual Multidimensional Model

While it is universally recognized that a data warehouse is based on a multidimensional model, there is no agreement on the approach to the conceptual modeling. In what follows, we will present a simplified version of the *Dimensional Fact Model* (DFM), proposed in [Golfarelli et al., 1998], a graphical conceptual model for data warehouses, aimed at

- effectively supporting conceptual design,
- enabling communication between the designer and the final user in order to refine requirements specification,
- supplying a stable platform for logical design, and
- providing an expressive and non-ambiguous design documentation.

The formalism enables the representation of the following basic information.

Facts

The most important abstraction mechanism of the conceptual model is the collections

of facts, i.e., the collection of observations of the performance of a business process. Facts are modeled by a rectangle divided in two parts, which contain the facts name and the set of *measures*. A measure is a numerical property of a fact that describes one of its quantitative aspects of interests for analysis.

Sometimes, facts are without measures, and are usually called *factless facts*, but in accordance with our terminology we call them *measureless facts*. This happens when facts represent events that only need to be counted.

Dimensions

Dimensions give facts their context, and are used to analyze them. Dimensions are represented by lines emanating from the rectangle of facts and ending with a circle (Figure 2.1). In general a dimension is described by a set of attributes used to qualify, categorize, or summarize facts in reports. For example, the dimension Date has the attributes Day, Week, Month, Quarter, and Year, while the dimension Store has the attributes City, State and Country. Dimensional attributes are represented as shown in Figure 2.2a, and the same names should not be used for attributes of different dimensions.

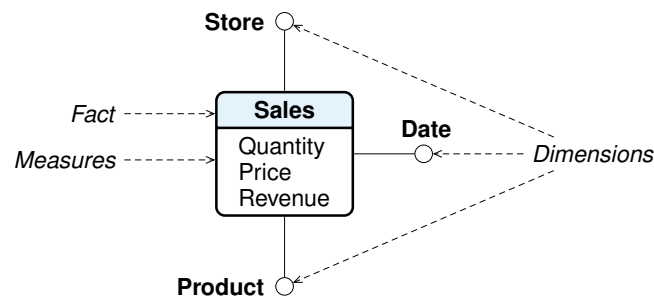


Figure 2.1: A conceptual design without dimensional attributes

Dimensional Hierarchies

In the presence of dimensional attributes, an interesting aspect to model, for the purposes of the data analysis, is a particular hierarchical relationship between their values, i.e., a many-to-one association between pairs of dimensional attributes. For example, the values of Month are in the hierarchy with those of Quarter and Year (Month → Quarter → Year), in the sense that a year is made up of more quarters, and a quarter is made up of more months, and, viceversa, a month corresponds to a single quarter, and a quarter corresponds to a single year. For this reason it is said that Year is more general than Quarter, and Quarter is more general than Month. In the terminology of the relational data model, each arc of the hierarchy models a *functional dependency* between two attributes.

Dimensional hierarchies are represented as shown in Figure 2.2b, with a directed tree, rooted in a dimension, and leaves representing the most general attributes.

Since a week usually crosses the boundary of two consecutive months, it is usually not treated as a lower abstraction of month. Instead, it is treated as a lower abstraction of year, since a year contains approximately 52 weeks.

The presence of a hierarchy between the dimensional attributes increases the possibilities of data analysis from different perspectives (*Multidimensional Analysis*). For example, once the sales of products have been analyzed by year, we can have a deeper analysis at a different level of detail to analyze product sales by quarter.

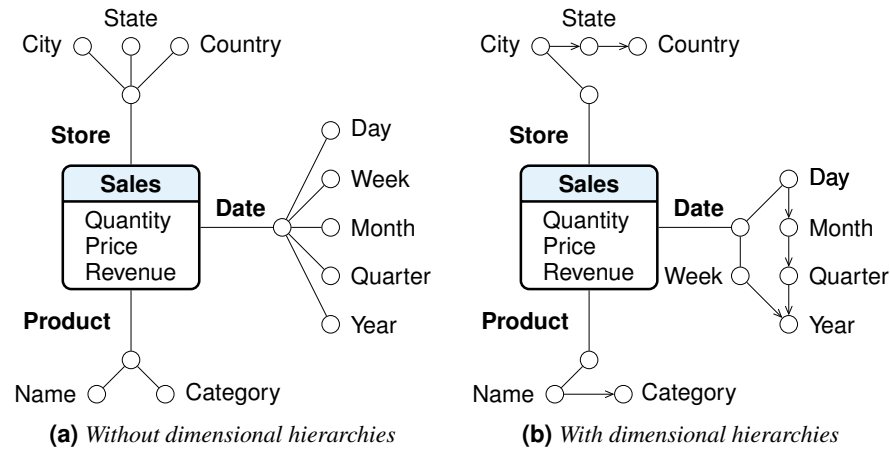


Figure 2.2: A conceptual design with dimensional attributes

The formalism enables the representation of other information. Let us see some examples (Figure 2.3).

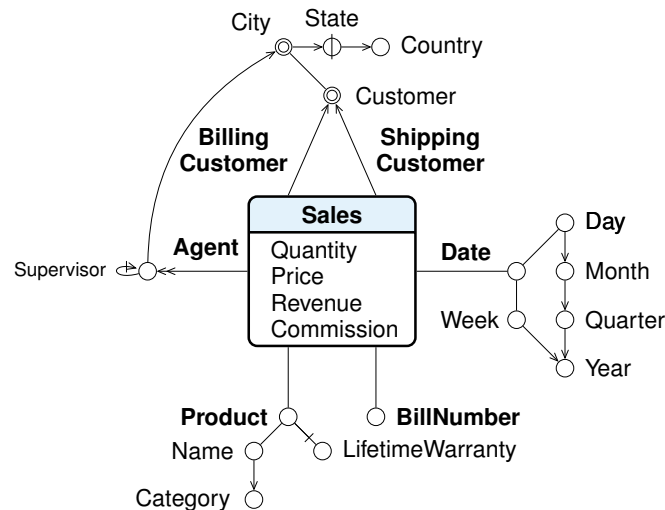


Figure 2.3: A conceptual design with other Dimensional Fact Model features

1. *Descriptive attributes.* Dimensions and dimensional attributes are usually represented with arcs ending with a circle to model that their values may be used in data analysis for selecting or grouping facts data. However there are cases in which dimensions and dimensional attributes are to be considered *descriptive* in the sense that in the data analysis is used only for selecting data, or to show their values in the report result, but not for grouping or aggregating data. A descriptive attribute is represented with an arc without a circle.
2. *Degenerate dimensions.* Dimensions without any attributes are called *degenerate dimensions*. Usually these are transaction-based numbers which describe the fact, but are not measures because it is meaningless to aggregate them. A typical example is a *Bill number* (Figure 2.3).

3. *Optional attributes or dimensions*. When the value of an attribute or a dimension may be undefined, the corresponding arcs are “cut”.
4. *Types of hierarchy*. A hierarchy among dimensional attributes can be of the following types:
 - *Balanced*, when the possible levels are a predefined number and the attribute values are always defined. For example, the Date attributes Month, Quarter and Year belong to a balanced hierarchy with three levels.
 - *Ragged*, when the values of one or more attributes may be undefined. A ragged hierarchy is graphically denoted by marking with a dash the attributes whose values may be undefined. For example, a location dimension with attributes Country, State and City, is balanced in the US, but it is ragged for most European countries where State is non used.
 - *Recursive (unbalanced)*, when the possible levels are a variable number. For example, in the dimension Agent there is the attribute Supervisor representing a recursive hierarchy among agents.

In the conceptual schema, a ragged hierarchy is represented by cutting the circle of the interested attribute, and a recursive hierarchy is represented with a loop.

5. *Shared hierarchy*. The dimensions can share some hierarchy attributes, such as City and Customer. To avoid ambiguity the circle is doubled and the arcs are oriented. Another typical example is the date hierarchy: a fact may have more than one Date dimension, with different semantics, and it may be useful to share among them the hierarchy month-quarter-year.
6. *Multivalued dimension or attribute*. A fact may be associated with more than one value of a dimension. For example, the fact sale is associated with several sales-people who have promoted it. In this case, the outgoing arc from the fact ends with a double arrow. Besides dimensions, dimensional attributes also may be multivalued, and represented in the same way.

2.1.1 Considerations on the Conceptual Modeling of a Data Mart

While in a database project, the focus is on collections of entities, their properties, associations and hierarchies between collections, in a data mart project, the focus is on the collection of facts, their measures, their dimensions, attributes and hierarchies [Kimball and Ross, 2002a], [Adamson and Venerable, 1998].¹

Let us present the key steps in conceptual design of a data mart, assuming that the business process of interest has already been identified together with the key analysis to be performed on the data to get the necessary information to make better business decisions. The example is about the business process of registration of customer orders. The objective is the analysis of the portfolio to adapt marketing strategies, promotion and inventory management.

Step 1: Identify the Granularity of the Fact

When modeling a data mart, the first fundamental decision to be taken is the *meaning of the fact*, because from this choice derive *the measures that characterize the fact and the dimensions for its analysis*. This is a classic problem in the design of data marts: you must carefully choose the right *grain* of the fact, i.e., *the precision with which the measurements are taken*.

In the case of customer orders, it could be said that the interesting thing is the Order, but, thinking about its meaning, we find that there is a problem, because an order is

1. We are grateful to Nicola Ciaramella for his contribution to the preparation of this section.

composed of a header and one or more lines, and we should decide if the fact is the header, about all products ordered, or the line for each product ordered.

As a general rule, it is best to choose a fine grain, even if it increases the number of facts to be treated, and so to choose a order line as a fact, because later in the analysis with aggregation functions, you can always go from the measures about the lines to the measures about the orders. If, instead, we focus on orders, there is no way to do the analysis in reverse order to move from measures about the orders to measures of individual lines.

Another consideration to keep in mind for the choice of the granularity of the fact is its nature, which may be of the following types (Figure 2.4).

■ Definition 2.1

A *Transaction Fact* represents the information on a specific event that occurred at a specific point in time during the execution of a business process.

For example, a fact is a transaction (withdrawal or deposit) on a bank account.

■ Definition 2.2

A *Periodic Snapshot Fact* represents the information on a series of events that have occurred over a period of time.

For example, a fact is the monthly summary of all transactions on a bank account.

■ Definition 2.3

An *Accumulating Snapshot Fact* represents the information on the lifetime of an evolving event that has a duration and change over time.

For example, the fact is about a mortgage application which is processed with the following phases: 1) presentation of the documents by the applicant as requested by the bank, 2) the bank's assessment of the documentation made available by the applicant, 3) approval of the practice and the initiation of investigative procedures; 4) mortgage completion. At the end of each phase, the fact about a mortgage application changes with the specification of the relevant information about its state. A solution for this case is presented in the appendix on case studies.

Cardinality of the facts

The grain of the fact determines the size of the set of facts that can be estimated using the estimates of the number of possible values for each dimension.

For example, let us consider the monthly Sales facts of the last five years, with the following total number of dimension values: Date ($12 \times 5 = 60$), Product (5 000), and Store (200).

If all the products are sold by all stores during a given month, there is a fact for each combination of dimension values, and so the number of sales facts is obtained by multiplying the dimension sizes ($60 \times 5000 \times 200 = 60\,000\,000$). Thus, the number of facts is many times larger than the dimension sizes.

In general, the number of combinations that actually appear in the set of facts is much less than this maximum number, because only some products are likely to be sold by each store during a given month. This property is referred to as *facts sparsity*.

Feature	Transaction	Periodic	Accumulating
Time period represented	Instant of time	Regular interval	Indeterminate period of time, usually of short duration
Grain	One fact per transaction	One fact per time period	One fact for the entire lifetime of an event
Update	No	No	For each state change
Measures	Related to transaction activities.	Related to periodic activities.	Related to activities which have a definite lifetime.
Dimension <i>Date</i>	Event Date	Date at the end-of-period	Multiple date dimensions to show the achievement of different milestones

Figure 2.4: Comparison of fact types

The cardinality of the facts depends on both the number of dimensions and the grain of the fact. Suppose that the marketing department requests that *daily* sales must be considered as facts. With the grain of sales changed to daily, the number of fact sales becomes 1 825 000 000. In this way, a fine granularity could result in a huge cardinality of the facts. Conversely, a too coarse granularity could result in facts that are not detailed enough for users to perform meaningful analysis.

Step 2: Identify the Fact Measures

Once the fact to represent has been chosen, the numerical measurements of interest are defined. A measure describes one of the fact’s quantitative aspects of interest for analysis.² *Facts may be also without measures, when used only to represent the occurrence of an event, such as the attendance of a student in a course.*

In choosing a measure we need to ask whether it makes sense to aggregate them with the function SUM, for analysis of the type “total value of the measure *M*, grouping data by dimension *D*”, which is usually expressed in the abbreviated form “total of *M*, by *D*”. In general, the aggregations with the function SUM are the most used in the analysis, but do not fall into the trap of believing that everything that can add up is an interesting measure, or that the sum is always meaningful. In general, the following measure types are considered.

■ Definition 2.4

An **additive measure** (also called a *flow* or *rate* measure) *can be meaningfully aggregated with the function SUM by any dimension.*

An additive measure is the most common type of measures. It refers to a time period and it is evaluated at the end of the period to record the cumulative effect over the period. For example, the number of products sold in a day or the monthly income.

2. The measurements are referred to as *measures* or *facts*, but we prefer the term *measures* because it is more descriptive.

■ Definition 2.5

A **semi-additive measure** (also called a *stock* or *level* measure) can be meaningfully aggregated with the function *SUM* by certain dimensions, but not all.

A semi-additive measure refers to a particular point in time and it is evaluated to record the state of an event. For example, the monthly account balance or the monthly inventory quantity-on-hand.

■ Definition 2.6

A measure M is *semi-additive with respect to a dimension D_1* when it can not be aggregated with the function *SUM* for groups of data with different values of D_1 .

Therefore, it makes sense to perform an analysis of the type “total of M , by D_1 ” — but not by a different dimension D_2 — or to perform analyses of the type “total of M of data with a certain value of D_1 , by D_2 ”. However, M may be aggregated with other functions such as *AVG*, *MIN*, *MAX*, for groups of data with different values of D_1 .

For example, the bank measure *Account balance* is *semi-additive* with respect to a dimension *Date*, but adding the *Account balance* for a particular day by the dimension *Customer*, or *Branch*, or *Account* can provide a meaningful information for the total amount of money the bank is holding at a given point in time.

Example 2.1

Let us consider the monthly *Quantity-on-hand* measure for different products and store at the end of every month. *Quantity-on-hand* is *semi-additive* with respect to both the dimension *Month*, and the dimension *Product*. In fact, it is not meaningful to total the *Quantity-on-hand* by *Product* because we would total values of different months, but it is also not meaningful to total the *Quantity-on-hand* by *Month* because we would total values of different products. It is correct to total the *Quantity-on-hand* by *Month* and by *Product*, or to total the *Quantity-on-hand* of the product *P1* by *Month*, or to total the *Quantity-on-hand* of the month *M1* by *Product*.

Inventory			
Product	Store	Month	Quantity-on-hand
P1	D1	M1	300
P1	D2	M1	100
P2	D1	M1	500
...
P1	D1	M12	100
P1	D2	M12	0
P2	D1	M12	900

■ Definition 2.7

A **non-additive measures** (also called *value-per-unit* measures) cannot be aggregated with the function *SUM* by any dimension.

Non-additive measures are usually the result of ratios. The only calculation that can be made for such a measure is counting the number of facts with such measures. Examples of non-additive measures include:

- *Per-unit price* cannot be added by any dimensions, while an extended price, such as *Per-unit price* × *Quantity purchased*, it is correctly additive by all dimensions.
- *Percentages* and *ratios*. A ratio, such as $\text{Gross Margin} = \text{Margin} / \text{Revenue}$, is non-additive. Whenever possible, such measures should be replaced with the underlying calculation measures (numerator and denominator) so that the calculation is made in the analysis as a metric. It is also very important to understand that when adding a ratio, it is necessary to take the sum of numerator and denominator separately and these totals should be divided.
- *Measure of intensity* such as the room temperature.
- *Averages* such as average sales price.

■ Definition 2.8

A **calculated measure** is a measure calculated on the basis of other measures.

It is strongly suggested that standard calculated measures are defined to avoid having users perform these calculations, because often they do not agree on their semantics and may perform wrong calculations. Moreover, having users doing standard calculations runs the risk of making the data warehouse seem unfriendly and complex, and, much worse, if the answers are wrong or inconsistent, the data warehouse will be viewed as wrong.

Example 2.2 Let us consider the following interesting measures for the fact OrderLines

OrderLines
Quantity
ExtendedPrice
ExtendedCost
Discount
Revenue
Margin

- The Quantity is the total number of products ordered.
- The ExtendedPrice and the ExtendedCost are calculated as follows

$$\text{ExtendedPrice} = \text{Unit Price} \times \text{Quantity}$$

$$\text{ExtendedCost} = \text{Unit Cost} \times \text{Quantity}$$

- The Discount is the value to be subtracted from the extended price.
- The Revenue and the Margin are calculated as follows

$$\text{Revenue} = \text{ExtendedPrice} - \text{Discount}$$

$$\text{Margin} = \text{Revenue} - \text{ExtendedCost}$$

Step 3: Identify the Fact Dimensions

The dimensions are chosen to provide context for facts. Without context, facts are impossible to analyze. To choose the dimensions, it is useful to consider the classic suggested questions to analyze the facts of everyday life (the **5W-1H rule**: *who, what, when, where, why, how*).

Who is the fact about?

With reference to the orders, they are generated by customers and, for example, it is interesting to analyze the order lines by customers to compute the total revenue. Thus, we select a Customer dimension and then later on we will define its attributes of interest. The question about *who* has another interesting answer: an order involves both the customer and the sales person that promotes the order on behalf of the company, and therefore SalesPerson is another relevant dimension.

The choice of a dimension is not always clear. However, it is useful to ask a question like this to find the right answers in the specific case under consideration.

What is a fact about?

As regards the order lines, a fact is about a product. Therefore, there is a Product dimension, and the choice is justified by the fact that it is meaningful and interesting to analyze order lines by products involved. This dimension is used to analyze the total revenue and cost of order lines by products of a certain category.

We wonder now if there are other interesting answers to the question of *what is an order line about*. No other relevant answers immediately come to mind.

When did a fact take place?

For *when* the answer is that we identify an instant in time or a time period. The two choices are not equivalent. In the case of customer orders, if we consider the order as an instant in time, then we can always perform an analysis by a time period, the converse is not true. For an analyst of business trends the time period is more interesting: to know the hour and minute of an order has its operational importance, but how orders are going is unlikely to be relevant; the preferred analysis will be by day or even better by month.

For our example we decided to choose the Date dimension for order lines, as it usually happens in any multi-dimensional model used in companies.

Where did a fact take place?

The question involves the definition of a Location dimension, another dimension that appears very often in real multidimensional models. For our example, this dimension is not considered because it is assumed that the location information is the customer city.

Similarly we could proceed further by asking questions such as *Why did a fact happen?*, *How did a fact happen?* to discover other dimensions, but the example does not suggest interesting answers.

We have a multidimensional model with the fact OrderLines and the dimensions Customer, SalesPerson, Product, Date, and now it is necessary to establish the attributes of dimensions and any hierarchy among them.

Before proceeding to the definition of the dimension structure it is useful to ask whether it is appropriate to associate to the facts *descriptive attributes* or *degenerate dimensions*. For the fact OrderLines, in addition to the measures discussed above, we consider useful also the degenerate dimension OrderNumber and the descriptive

attribute OrderLineNo (Figure 2.5): analysis by OrderNumber is useful for finding the average revenue by order, and (OrderLineNo, OrderNumber) is useful for identifying each line on an order.

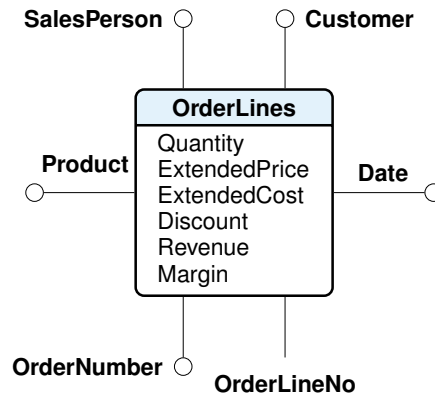


Figure 2.5: The data mart OrderLines conceptual design: the dimensions

Step 4: Identify Dimensional Attributes

Dimensions are the qualifiers that make the measures of the facts meaningful, because they answer the 5W-1H aspects of a question. To perform a more interesting analysis, it is generally necessary to describe each dimension with attributes relevant to the analysis that must be performed, and thus such that for each their value a subset of the facts on which a measures aggregation is somewhat interesting can be identified.

Let us consider the Date dimension. It is easy to imagine that among the requirements there may be both analysis of order lines by the Day attribute to compute the sum of revenues, and analysis by other date attributes, such as Month, Quarter and Year.

Let us consider then the Customer dimension and ask if it makes sense and is relevant to group customers by city of residence. The answer is yes because the information that the customers of a city have issued orders for a total amount higher than those of other cities helps to make a decision of whether to intervene on customers of different cities in a different way with different promotions.

This reasoning draws similar concepts that underlie the *segmentation of customers*. A customer segmentation is useful if

- segments behave differently with respect to their buying behavior;
- segments have a certain homogeneity behavior;
- it is possible to operate on segments with differentiated promotion actions.

Theories of customer segmentation also require other properties of a good segmentation, such that the segments are large enough to warrant different marketing actions, but the three properties listed capture the essence of the idea of segmentation: identifying customer groups that have a common behavior, which is very different from that of other groups, and so different marketing effort must be studied.

This statement is simply the principle underlying the clustering, one of the most important and interesting strategies for data mining. If we apply these principles to the structuring of dimensions, it turns out that the grouping of data may also be done

following other criteria. For example, it is usually not necessary that the Date dimension is structured into periods, but if the sales are about products with strong seasonality, then it will be interesting to divide time into seasons defined according to the logic inherent to the phenomenon to be analyzed. Suppose that some products are sold almost exclusively in the pre-Christmas period, in this case the year can be divided into two periods, one from the beginning of December until Christmas, the other covering the rest of the year. Another example is that it may be useful to make a distinction between sales on the weekend and those on other days.

If we think of a dimensional model of the data of an urban public transport company, we discover that we need to move from the day to time periods, such as entry and exit from offices or schools. The definition of time periods is not standard: it is a decision that must be taken according to the logic inherent to the movements of travelers, but also according to the logic of company operations. For example, early morning (6 a.m. to 8 a.m.), late morning hours (8 a.m. to 11 a.m.), rush hour (11 a.m. to 1 p.m.), lunch hour (1 p.m. to 2 p.m.), and so on.

If the company cannot change its way of operating at night, then breaking the night time in time periods serves only to satisfy curiosity, but ends up complicating the report without any real added value for decision-making.

In general, the structure of a dimension should therefore reflect two logics:

- *The logic of the event to be analyzed*: the values of dimensional attributes at every level of the hierarchy, are used to group fact data so that groups are internally homogeneous and different between themselves with respect to the values of the measures, to help the analyst to understand what the factors that influence the event are.
- *The logic of company operations*: the values of dimensional attributes at every level of the hierarchy are used to group fact data so that groups are internally homogeneous and different between themselves with respect to their reaction to the actions of the company, to help the decision maker to revise their actions in order to influence the event.

Step 5: Identify the Dimensional Attribute Hierarchies

Dimensional attributes are useful for generating readable reports, but the most interesting attributes for interactive multidimensional analysis are organized into hierarchies to allow groupings of facts data and aggregations of the measures at different levels of generality, as usually required in practice. For example, in the case of the Date dimension, the hierarchy Day → Month → Quarter → Year is relevant. The hierarchies of interest for the other dimensions are shown in Figure 2.6.

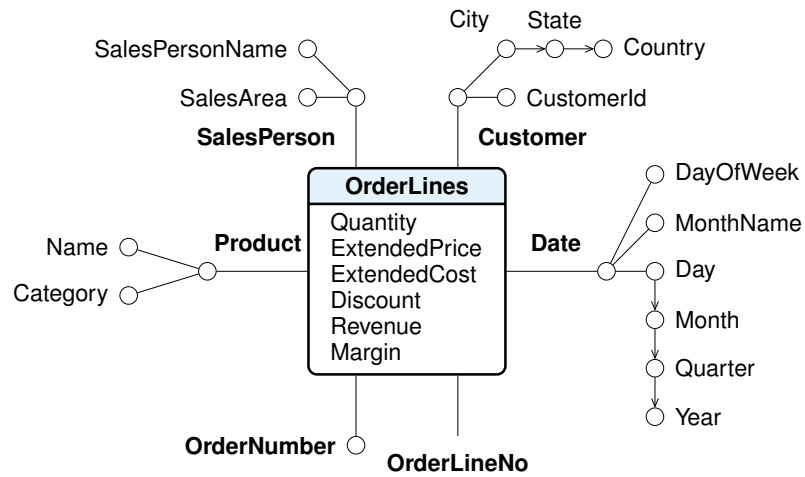


Figure 2.6: The data mart OrderLines conceptual design with dimensional hierarchies

2.2 Multidimensional Relational Model

A conceptual multidimensional schema is transformed into a relational logical schema by applying a set of mapping rules, as will be described in the following chapter. The result depends on the complexity of the conceptual schema, and in this section, we show only the basic idea of the structures of specialized schemas usually used, called *star schema*, *snowflake schema* and *constellation schema*.

■ Definition 2.9

A *star schema* consists of a *fact table*, which contains the data about the facts to be analyzed, and a set of *dimension tables*, one for each dimension. Each of the dimension tables has a single attribute primary key which has a one-to-many relationship with a foreign key in the fact table. Usually a dimensional primary key is a simple integer *surrogate key* that is numbered sequentially from 1 to the number of records in the dimension table. Usually a meaningful integer surrogate key of the form YYYYMMDD is used for a date with the granularity of a day (e.g. 20140926 for September 26, 2014).

The fact table is at the center of the “star”, whose tips are the dimension tables (Figure 2.7). A *star schema* is an intentional simplification of the database design that would be achieved by following the standard rules of normalization.

Note that using the surrogate key YYYYMMDD for the dimension Date, the Day attribute is useless and a value of the Month attribute is an integer of the form YYYYMM to represent correctly the dimensional hierarchy Month → Year.

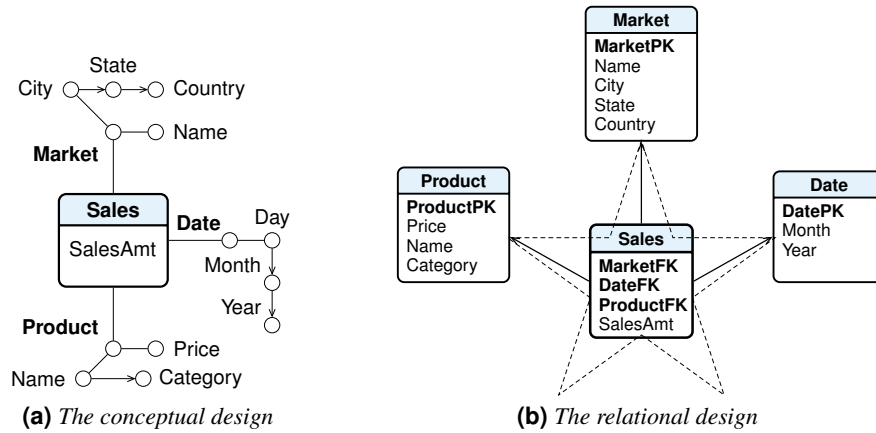


Figure 2.7: Example of a star schema

■ **Definition 2.10**

A *snowflake schema* is a variant of the star schema, where some dimension tables are *normalized*, thereby further splitting the data into additional tables.

The saving of space is usually negligible in comparison to the typical magnitude of the fact table (Figure 2.8). Furthermore, the snowflake structure can increase the time to execute queries that require hierarchies to be traversed, since more joins will be performed to execute them. Hence, although the snowflake schema reduces redundancy, it is not as popular as the star schema in data warehouse design.

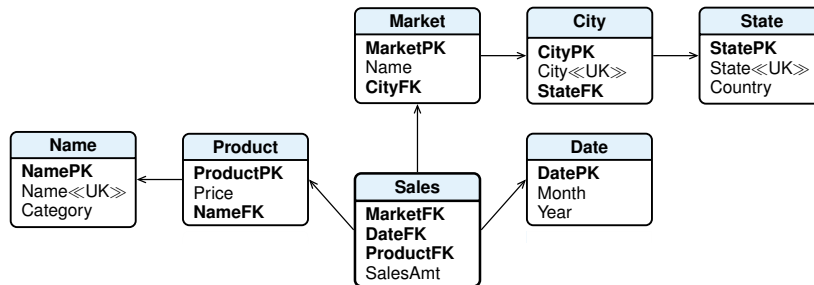


Figure 2.8: Example of a snowflake schema

■ **Definition 2.11**

A *constellation schema* has multiple fact tables that share dimension tables.

The example given in Figure 2.9 has two fact tables Sales and Returns sharing the Date and Product dimensions.

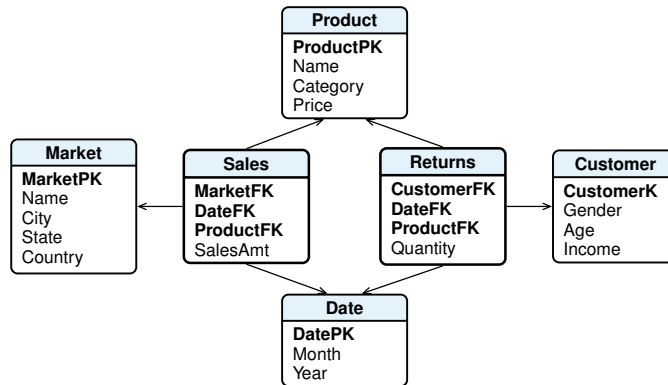


Figure 2.9: Example of a constellation schema

The main relational DBMS vendors provide OLAP servers that map operations on multidimensional data to standard relational operations on specialized relational DBMS to store and manage data warehouses. Such servers are referred to as **ROLAP** (*Relational OLAP*).

2.3 Multidimensional Cube Model

■ Definition 2.12 *Cube Model*

A multidimensional cube model (*data cube*) represents facts with n dimensions by points in an n -dimensional space. A point (a fact) is identified by the values of dimensions and has an associated set of measures.

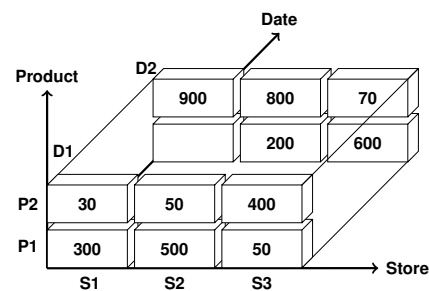
Such a multidimensional view is an intuitive way to think about OLAP queries and their results. For the sake of simplicity, we will consider a cube with at most three dimensions and one measure.

Example 2.3

Let us consider the analysis of the daily sales of different products in different stores over different days. Let us assume that data are stored into a fact table such as that shown in the figure (a). Store identifies a store, Product identifies a product, Date identifies a day, and Qty identifies the quantity sold of that product at that store in that time period.

Store	Product	Date	Qty
S1	P1	D1	300
S2	P1	D1	500
S3	P1	D1	50
S1	P2	D1	30
S2	P2	D1	50
S3	P2	D1	400
S2	P1	D2	200
S3	P1	D2	600
S1	P2	D2	900
S2	P2	D2	800
S3	P2	D2	70

(a) Fact Table



(b) Data Cube

We can view this sales data as *3-dimensional*, because the value of the *measure* Qty is a function of the Store, Product, and Date attributes, which form the so-called *dimensions*. Consequently, we can also think of the data in a fact table as being arranged in a 3-dimensional cube shown in the figure (b). For example, the cell ('S1', 'P1', 'D1') contains the sales for the product P1 on date D1 by the store S1.

The 3-dimensional cube is a generalization of a 2-dimensional cross-tabulation commonly used to give a basic picture of how two attributes inter-relate because it helps to search for patterns of interaction.

Product	Store		
	S1	S2	S3
P1	300	500	50
P2	30	50	400

When dimensions have attributes and hierarchies, the multi-dimensional cube is more complex. We will assume that additional information about the dimensions are stored in tables, which describe dimensions' attributes.

Some vendors provide OLAP servers that implement the fact table as a data cube using a specialized data structure. Such implementations are referred to as **MOLAP** (*Multidimensional OLAP*).

2.3.1 OLAP Operations in the Multidimensional Data Model

Let us show some typical OLAP operations for multidimensional data. Each of the operations described below is illustrated in Figure 2.10.

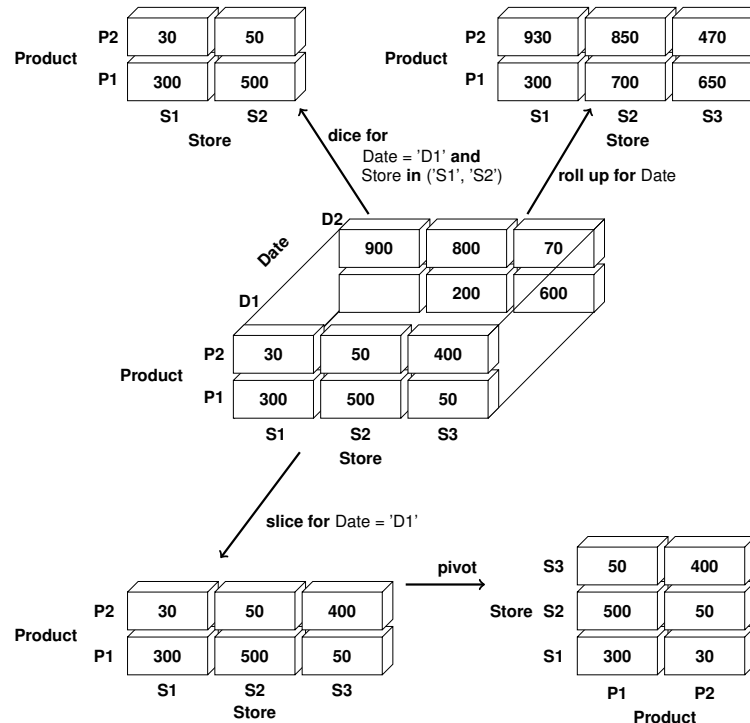


Figure 2.10: Examples of typical OLAP operations on multidimensional data

Slice and dice

The operators *slice* and *dice* generate sub-cubes by selections, but they do not change the measures values, that is they do not make summarizations:

- The *slice* operator selects a cross section that cuts across a cube with a selection on one dimension (Figure 2.10).
- The *dice* operator selects a sub-cube with a selection on two or more dimensions (Figure 2.10).

Roll-up and Drill-down

The *roll-up* operator, also called *drill-up*, performs summarizations at different levels of details either by dimension reduction or by climbing up dimension hierarchy.

Figure 2.10 shows the result of a roll-up operation by removing the date dimension, summarizing the the quantity sold by product and by store.

The *drill-down* operator is the reverse of roll-up. It produces more detailed data from less detailed data. Drill-down can be used by either stepping down a hierarchy for a dimension or introducing additional dimensions.

Pivot

The *pivot* operator (also called *rotate*) performs a rotation of the data axes to provide an alternative presentation of data (Figure 2.10).

2.3.2 The Extended Cube

Let us assume that each dimension is extended with an additional value “*”. This value has the intuitive meaning “all”, and it represents summarization along the dimension in which it appears. A cube can be extended with new “borders” made of cells containing the value of aggregate functions (we consider here only the SUM) as shown in Figure 2.11.

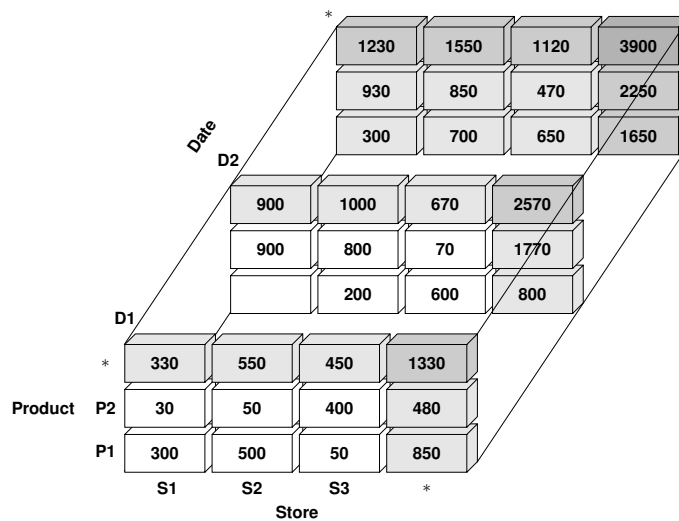


Figure 2.11: Three-dimensional cube extended with cuboids

For example, using the notation Sales(Store, Product, Date, Qty) for a cube with dimensions Store, Product, Date and a measure Qty, we can denote subcubes as follows:

- ('S1', 'P1', 'D1') is the cell that contains 300, the sales for the product P1 on date D1 by the store S1;
- ('S1', *, 'D1') is the cell that contains 330, the sum of sales for all products on date D1 by the store S1;
- ('S1', *, *) is the cell that contains 1 230, the sum of sales for all products over all time by the store S1;
- When a dimension is used as a coordinate instead of one of its values, the notation denotes a so called *cuboid*. For example (Store, Product, *) is the cuboid “roll up for Date” in Figure 2.10, with two dimensions with the cells that contain the sum of sales over all time by the dimensions Store and Product (in SQL terms, the sales data are grouped by Store and Product, and the aggregate function SUM(Qty) is computed).

In Figure 2.11, the border with the lightest shading represents aggregates in one dimension, darker shading for aggregates over two dimensions, and the darkest cuboid in the corner for summarization over all three dimensions. In general the border represent only a small addition to the volume of the data cube (the white cuboid).

Table 2.1: Sales extended cross-tabulation

Product	Store			Total
	S1	S2	S3	
P1	300	500	50	850
P2	30	50	400	480
Total	330	550	450	1330

The 3-dimensional extended cube is a generalization of a 2-dimensional extended cross-tabulation (Table 2.1).

To speed up data analysis, commercial data cube systems precompute all or some of the cuboids and store them as *materialized views* of the data cube. The problem of selecting the cuboids to precompute will be studied in a later chapter.

The total number of cuboids for a data cube with three dimensions is $2^3 = 8$. The possible cuboids can also be denoted without using the “*” as follows: (Store, Product, Date), (Store, Product), (Store, Date), (Product, Date), (Product), (Date), (Store), (). (Store, Product, Date) denotes the data cube, while () denotes the total sum of all sales.

These cuboids can be represented as a lattice, also called the *data warehouse lattice*, as shown in Figure 2.12. We say that the cuboid C_1 is below the cuboid C_2 , written $C_1 \preceq C_2$, if and only if C_1 can be computed from C_2 . The cuboids are named using the abbreviations P for Product, S for Store, D for Date.

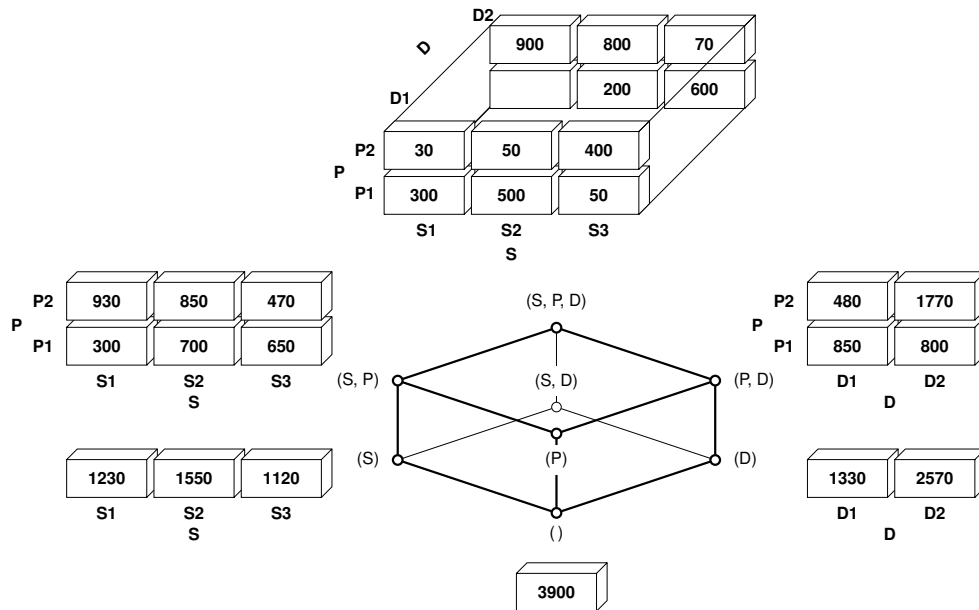


Figure 2.12: Lattice of cuboids

In general the computation of the cuboid C_1 from C_2 depends on the aggregate function used, which can be of one of the following types.

■ **Definition 2.13** *Distributive Aggregate Functions*

An aggregate function f on a multiset of values V is *distributive* if there is a *local* aggregate function f_l and a *global* aggregate function f_g , such that for any k -partition $\{V_1, \dots, V_k\}$ of V we have

$$f(V) = f_g(\{f_l(V_1), \dots, f_l(V_k)\})$$

For example, the SQL functions SUM, MIN, MAX and COUNT are distributive aggregate functions:

- $SUM(V) = SUM(\{SUM(V_1), \dots, SUM(V_k)\})$
- $MIN(V) = MIN(\{MIN(V_1), \dots, MIN(V_k)\})$
- $MAX(V) = MAX(\{MAX(V_1), \dots, MAX(V_k)\})$
- $COUNT(V) = SUM(\{COUNT(V_1), \dots, COUNT(V_k)\})$

■ **Definition 2.14** *Algebraic Aggregate Functions*

An aggregate function is *algebraic* if it can be computed from a finite algebraic expression defined over distributive functions.

For example, the functions average (AVG), variance (VAR), and standard deviation (STDEV) are algebraic aggregate functions, which can be computed on V using the following distributive aggregate functions on the multiset of a k -partition $\{V_1, \dots, V_k\}$

- $n_i = COUNT(V_i)$
- $s_i = SUM(V_i)$
- $s_i^2 = SUM(V_i^2)$, where V_i^2 is the set of the squares of the various elements of V_i .

Let $n = COUNT(V) = SUM(\{n_1, \dots, n_k\})$.

The functions $AVG(V)$, $VAR(V)$ and $STDEV(V)$ are computed as follows:

- $AVG(V) = SUM(\{s_1, \dots, s_k\})/n$
- $VAR(V) = \frac{SUM(\{s_1^2, \dots, s_k^2\}) - (SUM(\{s_1, \dots, s_k\}))^2/n}{n - 1}$
- $STDEV(V) = \sqrt{VAR(V)}$

■ **Definition 2.15** *Holistic Aggregate Functions*

An aggregate function is *holistic* if it can not be computed from other aggregate functions.

For example MEDIAN, MODE, RANK.

2.4 Summary

- A data warehouse conceptual model is the best support for discussing, verifying, and refining user specifications since it achieves the optimal trade-off between expressivity and clarity.
- A multidimensional relational model is used to implement data warehouses. This model can adopt a *star schema*, *snowflake schema* or a *constellation schema*. The core of a multidimensional model is a *fact table* and a set of *dimensional tables*.
- The core of a *multidimensional model* is the *data cube*. An extended cube consists of a lattice of cuboids, each corresponding to a different degree of summarization of data. *Full* or *partial materialization* refers to the pre-computation of all or some of the cuboids in the lattice. Commercial systems use different strategies both about which cuboids to materialize, and how to store them.

DATA WAREHOUSE DESIGN

The purpose of a data warehouse (DW) is not just to store data but rather to facilitate decision making. Therefore, a data warehouse must be designed taking into account the different types of analyses that are needed by the business users to make better decisions about key business processes worth of improvements. Since a data warehouse design process is complex, a methodology organized in phases is presented, like the one that is used to design operational databases, to highlight the importance of conceptual design and shows how to transform the conceptual design into the logical one using the relational model.

3.1 Introduction

A data warehouse must be designed to provide the information needed to solve a business problem. If the problem is solved there should be some economic gain in order to allow a cost benefit analysis for the data warehousing project.

As happens for databases, it has become fairly standard to divide the DW design process into the following four phases:

1. *Requirements Analysis*. The goal is to produce a description of the business processes, the typical information analysis activities with which users are involved, and the measures and dimensions of interest. Typically, requirements at this stage are documented rather informally.
2. *Conceptual Design*. The goal is to produce a formal description of the data to be analyzed in high-level-term using a conceptual data model. We will use the *Dimensional Fact Model*, *DFM*, to describe facts, dimensions, dimensional attributes and attribute hierarchies.
3. *Logical Design*. The goal is to transform the conceptual design into the logical structures used for storing the DW in a relational DBMS.
4. *Physical Design*. The goal is to define the data structures needed for storing the database tables created by the logical design. The main issues are what indexes and materialized views to define to optimize the overall performance of the system.

Once the DW has been implemented, data must be extracted from the operational and external systems, transformed into a usable format for the DW, and finally loaded into the DW in order to be usable for query processing and analysis. These *Extract, Load, Transform (ETL)* processes have historically been batch-oriented.

In general the data to load in the DW are processed with two important and complex kinds of operations:

1. *Transform*. When the data come from different sources, their *formats* are revised to align them by eliminating syntactic and semantic differences.
 - (a) *Syntactic transformation*. The same data can have both attributes with different names, and the names are not those to be used in the DW, and different types. For example, a code is defined in some cases of a type string, and in others a type integer; a gender is defined as (M, F), (m, f), (0, 1) or (male, female); a value is defined with different units of measurement, and so on.
 - (b) *Semantic transformation*. The data in source databases may have been used with a different meaning. For example, sales can be daily or weekly.
2. *Cleaning*. The data are analyzed in order to eliminate errors of representation or to complete missing information. For example, in the case of addresses the zip code can be wrong or the name of the town can be written in different ways (Busto Arsizio also written as BustoArsizio or BARSizio)

The information generated during the design and implementation of a DW is organized and stored as metadata using appropriate specialized tools or taking advantage of the capabilities of DW systems that provide, as any DBMS, a *catalog* that contains information about the logical and physical organization of the data managed. In the case of DW metadata are about other aspects of the data and, in short, can be classified into the following main categories:

- *Business metadata*. Concern the meaning of the terms used to define the logical structure of data in corporate terminology. This type of metadata is usually used by users to understand the nature of the data available.
- *Structural metadata*. Concern the logical structure of facts and dimensions, types of attribute values, hierarchies, dimensions and meaningful aggregations. This type of metadata is usually used by users to understand what types of analysis can be performed on the data.
- *Technical metadata*. Are concerned with the physical data property, such as storage structures, data sources, date of loading, transformations applied etc. This type of metadata is usually used by technicians for the maintenance and development of a DW.
- *Operational metadata*. Concerns the types of predefined analysis reports and what parameters should be used.
- *Design metadata*. Concern the results of the DW design phases.

Loading metadata is only partially automated and has a cost: the time that the personnel involved in the design and implementation of the DW must dedicate to the problem. Finally, note that metadata is useful not only to gather information on the data, but also to be exchanged between different OLAP tools. For this reason, proposals have been made to define standards like *Common Warehouse Metamodel (CWM)*.

In the following, the focus will be on how to proceed in the conceptual and logical design phases of a DW.

3.2 Data Warehouse Design Approaches

According to [Artz, 2005] and [Ballard et al., 2006], the approaches to DW design can be of the following types (Figure 3.1):

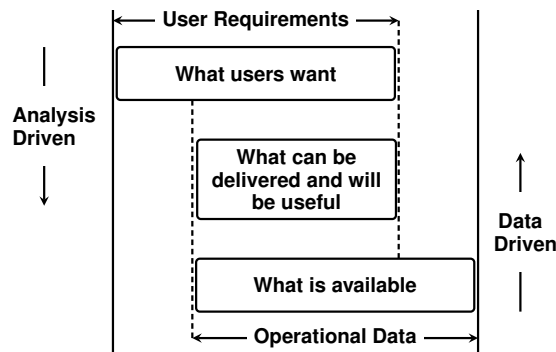


Figure 3.1: Data warehouse design approaches

1. *Data-driven*. This approach was originally proposed by Inmon, one of the first authors on the subject of data warehousing, which he describes as *top-down*, and whose supporters are referred to as “Inmonites”. The goal is to design first an enterprise DW based on the data available in the operational information system, and then the data marts are created from the DW. This is done by analyzing a conceptual model of data, if one is available, or the actual logical record layouts and selecting data elements deemed to be of interest. This approach is the only possible when the demand for information from a DW does not exist until the DW is actually available. An initial DW design on the basis of the data available can help both users to discover new ways in which to use the available data, and the designer to identify areas on which to focus data warehouse development efforts. The disadvantage of this approach is that without user involvement there is the risk of producing a non interesting result.
2. *Analysis-driven*. This approach was originally proposed by Kimball, a well-known author on data warehousing, which he describes as *bottom-up*, and whose supporters are referred to as “Kimballites”. The goal is to design first the data marts based on the data analysis that the users want to perform, and then the data marts are integrated to build the DW. The major advantage to this approach is that the focus is on providing what is really needed, rather than what is available. In general, this approach has a smaller scope than the data-driven approach. Therefore, it generally produces useful data marts in a shorter time span. The disadvantage of this approach is the risk that some of the data that the analysis needs are not available. Moreover, if a user is too tightly focused, it is possible to miss useful data that is available in the operational information system.

Both the approaches can be useful in certain cases. In the following we will use a combination of the two using the following design phases for each data mart of interest:

1. *Requirements analysis*
2. *Initial analysis-driven data mart conceptual design*
3. *Candidate data-driven data mart conceptual design*
4. *Final data mart conceptual design*
5. *Data mart and DW logical design*

3.2.1 Requirements Analysis

The requirements analysis phase is divided into two main sub-phases, characterized by the different language used for the preparation of the documents they produce. The first sub-phase, *Requirements gathering*, produces a natural language specification of requirements. The second sub-phase, *Requirements specification*, produces a description of the requirements for data analysis outlining the salient features to be modeled then with the conceptual design.

1. *Requirements gathering*

- (a) Analysis of the problem domain for which the modeling will be done.
- (b) Analysis of the business processes to select, with end-user interviews, those more interesting to consider for designing the data warehouse.
- (c) Business questions that end-users issue and try to answer in the course of their information analysis activities.

Examples of frequently encountered categories of business questions are:

- *Existence checking analysis*, such as “A given product has been sold to a particular customer.”
- *Item comparison analysis*, such as “Compare the value of purchases of two customers over the last six months,” or “Compare the number of items sold for a given product category, by store, and by week.”
- *Trend analysis*, such as “The growth in item sales for a given set of products, over the last 12 months.”
- *Queries to analyze ratios, rankings, and clusters*, such as “Rank our best customers in terms of dollar sales over the last year.”
- *Statistical analysis*, such as “The average item sales by product category, by sales region.”

Note that the business questions must usually be “interpreted” in order to express them in a form more useful for designing the data warehouse. For example, a business question a manager of a company wish to ask of their data might be: “Why are our sales not meeting our targets”. The business question might be interpreted in a more useful form as “For the current year, the cumulative quantity sold and targets, by product”. That is, an interpreted business question should be expressed with the types of reports to be produced, or phrases that reveal the following information:

- The constraints on data to analyze.
- The requested dimensional attributes and the metrics (*aggregation operation*) to compute.
- The coordinates (dimensions) against which the fact must be analyzed.
- The result sorting criteria and if metrics’ partial value are needed.

When business questions are expressed by means of phrases, the use of the following form is suggested: “For a data *subset* to use, the *metrics* to compute, *by dimension*, . . . , *by dimension*. How the result should be presented”. For example, “For the year 2010 in Italy, the total of sale revenue, by region, by month, by customer name. The result must be sorted by region, month, and customer name, and must include partial totals for all regions.”

It is important also to check that the information analysis requirements need data that are currently available or can be obtained as external data that exist outside the enterprise. If there are multiple data sources, the analysis is complicated by the need to reconcile the likely differences in representation of information. In the following we will not consider this aspect.

2. *Requirements Specification*. The business questions are specified using a set of worksheets with the following structure:

Business Process Requirements

				Process
N	Business questions	Dimensions	Measures	Metrics

Each business question is analyzed to identify the fact measures, the preliminary dimensions, and the metrics to be computed.

Fact Description

			Fact
Description	Preliminary dimensions	Preliminary measures	

The meaning of the fact is described, together with its grain and type (transaction, periodic, or accumulating), and the preliminary measures and dimensions.

Dimensions

			Dimensions
Name	Description	Granularity	

The meaning of each dimension is described, together with its name, a description, and the grain.

Dimensional Attributes

		Dimension
Attribute	Description	

Of all dimensions the attributes and their description are listed. Dimensional attributes must be chosen carefully to express the analysis in a natural way and display results in a comprehensible way. If there are attributes with values that are codes, providing the description of the code is also suggested. It is also advisable not to use the same names for attributes of different dimensions that have different meaning.

Dimensional Hierarchies

			Dimensional Hierarchies
Dimension	Hierarchy description	Hierarchy Type	

It describes the dimensional hierarchies for each dimension, and their type (balanced, ragged, recursive).

Dimensional Attributes Changes

			Changing Dimensions
Name	Changing Attribute	Treatment of changes	

It is important to understand how the business wants to deal with the dimension attributes that can change over time. Consequently, for each of them, besides the name, the type of strategy is specified to deal with them in the logical design phase and data loading.

For example, suppose that the dimension Customer of facts Order contains the attribute Residence, with a value that can change over time, and that there are several sales involving a customer from Lucca, which to a certain date changes residence to Pisa. How can we carry out sales analysis to account for this change? What should the result be of analysis such as “How many sales are made in Pisa?”.

The strategy to be specified depends on the analysis’ objectives, and for this reason it should be documented in the requirements specification. We consider four options, of which the first three are considered for *slowly changing dimensions*:

Type 1 (overwriting the history) If a dimensional attribute changes its value, only the latest value is required to be held in the data mart. This means that there is no need to preserve the previous value.

For example, if a customer changes address, the new one replaces the present value of the dimension Customer. It is the easiest solution, but the history of customer addresses is lost. For example, if a customer at a certain date changes their address from Pisa to Lucca, all sales concerning him before and after this date are considered made in Lucca, and this changes the outcome of analyses such as “How many sales are made in Pisa?”

Type 2 (preserving the history) If a dimensional attribute changes its value, both the old and the new value are required to be held in the data mart, but the structure of the dimension must not be changed. It is the solution commonly used.

Type 3 (preserving one or more versions of history) If a dimensional attribute changes its value, the structure of the dimension is extended with additional attributes to keep the tracking history with both old and new values. Moreover, we also add another attribute EffectiveDate for the date of the change. This solution is rather quirky and it is rarely used. We will not consider it further.

Type 4 (fast changing) The dimensional attributes change frequently, and must not be treated with one of the previous solutions.

Measures

Fact measures			
Measure	Description	Aggregability	Calculated

It describes each measure of the fact identified from the requirements, how it is calculated from other measures, and the aggregate functions that are applicable to the measure when the data are grouped by dimensions.

Descriptive attributes of the facts

Descriptive attributes	
Attribute	Description

It contains each descriptive attribute of the facts, with a description of what they represent.

Summary of Dimensions and Measures

Facts Dimensions			
Dimension	Fact ₁	...	Fact _n

Facts Measures			
Measure	Fact ₁	...	Fact _n

If the requirements concern different facts, the worksheets specify in which facts the selected dimensions and measures are used. The worksheet about dimensions, called the *data warehouse bus architecture*, is useful to identify which dimensions are used by multiple data marts, and therefore they must have a unique interpretation and representation (*conformed dimensions*) to be then shared in the DW. If the dimensions must be kept different, they must be renamed.

3.2.2 Initial Analysis-Driven Data Mart Conceptual Design

An *initial* data mart conceptual design is defined from the analysis that the users perform (*the design from what the users want*), without any claim to completeness but useful as a formal description of requirements. In the conceptual design the dimensional hierarchies are modelled together with their type (balanced, incomplete, recursive), degenerate dimensions and descriptive attributes of the facts.

3.2.3 Candidate Data-Driven Data Mart Conceptual Design

A method is described for developing a *candidate* data mart conceptual design from the operational database relational schema (*the design from what is available*). This approach to data mart design ensures that its schema reflects the underlying structure of the data available. The following steps are based on the proposal presented in [Moody and Kortink, 2000]:

1. **Operational data analysis.** In this step, the relational database schema is analyzed to perform two actions:
 - (a) Standardize terminology and units of numerical quantities that have an identical time reference.
 - (b) Delete tables, and attributes, considered not relevant to the analysis of the data. For example, information such as the tax code and telephone number are usually not relevant for the analysis of the data.
2. **Tables classification.** In this step, the relational database tables are classified in three categories to identify the possible facts, measures, dimensions and hierarchies between dimensional attributes.
 - (a) **Transaction Entities.** These are tables with records that represent events of interest for the business process to be analyzed (orders, insurance claims, salary payments, sales, hotel booking, etc.). Transaction entities have two fundamental characteristics:
 - Describe events that occur at a point in time.

- Contains measurements or quantities that may be summarized (e.g. monetary value, quantity, weight, volume).

It is very important to correctly identify the pertinent transaction entities because they are the natural candidates to be considered later for the definition of facts that decision makers want to understand and analyze. However, it must be kept in mind that not all transaction entities will be of interest for decision support, so they must be analyzed carefully with users to identify which of them are important.

- (b) **Component Entities.** These are tables directly related to a transaction entity via a *one-to-many relationship* (Figure 3.2). These entities define the details or components for each transaction event, and so are useful to answer the *who, what, when, where, why* and *how* of a business event. Component entities are the basis for defining dimensions in the data mart conceptual design.

An important component entity of any transaction entity should be the one that represents the time: the historical analysis, in fact, play a key role in all the DW, but usually in the operational database time is not represented with a table but with an attribute of type *Date* and this must be taken into account when defining the data mart design.

Note that the definition of component entity does not allow us to isolate multi-valued dimensions, resulting instead from tables directly related to a transaction entity via a *many-to-one relationship*. In general, this type of table should be considered in the choice of a component entity, as will be shown in a following example.

- (c) **Classification Entities.** These are tables related to a component entity by a chain of *one-to-many relationships* (Figure 3.2). These entities usually represent hierarchies embedded in the data schema. Their interesting attributes are collapsed into component entities to define then in the data mart design the dimensional attributes and hierarchies.

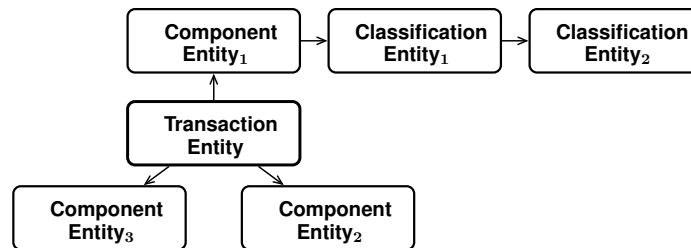


Figure 3.2: Tables classification

In some cases, entities may fit into multiple categories. We therefore define a precedence hierarchy for resolving such ambiguities:

- Transaction entity (highest precedence).
- Classification entity.
- Component entity (lowest precedence).

For example, if an entity can be classified as either a classification entity or a component entity, it should be classified as a classification entity. In practice, some entities will not fit into any of these categories. Such entities do not fit the hierarchical structure of a dimensional model, and cannot be included in the conceptual design.

Example 3.1

Figure 3.3 shows an operational database schema for an orders sales application, assuming that a row of an order may have dealt with more than one shipment.

- Transaction entity: it is interesting to collapse Order into OrderLine, with the OrderLine granularity to define the data mart facts. Other possible transaction entities might be Invoice, Product or Shipment, but they are not considered of interest for decision support.
- Component entity: Customer, Invoice and Product.
- Classification entity: ProductCategory is collapsed into Product.

The table Shipment does not satisfy the condition to be considered as a component entity of OrderLine because of the *many-to-one relationship*, but it might be considered to define a multivalued dimension.

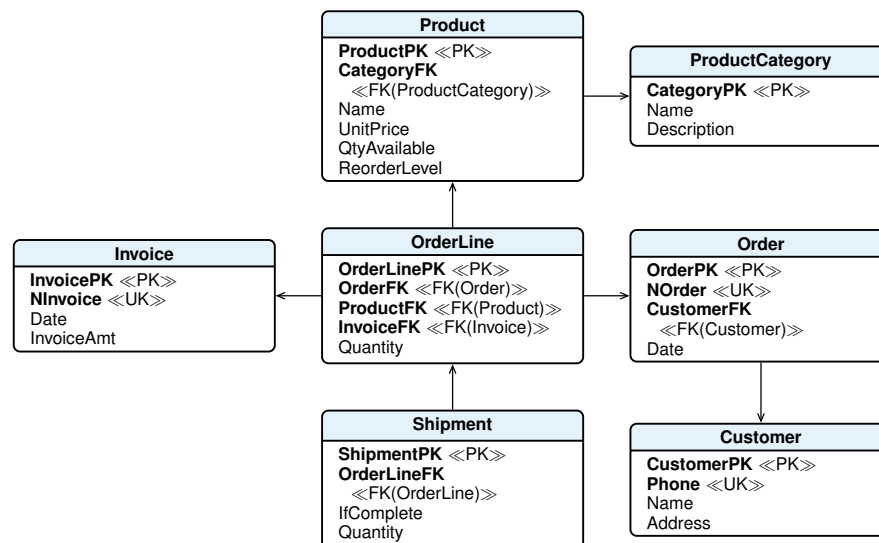


Figure 3.3: A database for order management

3. Candidate data mart conceptual design.

The *Candidate* data mart conceptual designs are defined in the following way:

- For each transaction entity, a data mart fact is defined.
- A dimension is formed for each component entity, by collapsing hierarchically related classification entities into it. The dimensional attributes are analyzed to decide possible hierarchies. For example, an attribute Date of a transaction entity is usually substituted with attributes Day, Month, Year, and a hierarchy is defined among them; an attribute Address may be replaced by City, and Region, and a hierarchy is defined among them.
- Where hierarchical relationships exist between transaction entities, the child entity inherits all dimensions (and key attributes) from the parent entity. This provides the ability to “drill down” between transaction levels.

3.2.4 Final Data Mart Conceptual Design

From a comparison of the initial and candidate conceptual designs the *final* data marts are defined (*the design of what can be delivered and will be useful*), which in general will be an extension of the common parts.

3.2.5 Data Mart and Data Warehouse Logical Design

Assuming that the multi-dimensional model is implemented with a ROLAP system, firstly each final conceptual data mart design is translated in a relational schema, deciding whether to make a star or snowflake schema, and then integrating the various data marts schemas in a single DW schema, considering the following possibilities:

- Standardize and share the fact tables with the same dimensions.
- Standardize and share common dimension tables.

In the definition of the relational tables of the data mart logical schema, the following problems will be considered, among others that may arise [Kimball and Ross, 2002b].

Primary keys of dimension tables

The primary key of each dimension table should be an attribute with numerical values automatically generated (*artificial or surrogate key*) in addition to any primary key used in the original data, if it is considered relevant to also keep this information in order to determine from which original data it comes from, but which is not necessarily a key for the dimension table.

For the Date dimension table with the granularity of the day, usually present in every data mart, it is useful not to use a surrogate key for the primary key, but an integer representing a day in the form YYYYMMDD. With a similar format it is useful to represent attribute values in the dimensional hierarchy Month → Quarter. Usually there are also other attributes useful to show in reports, such as DayName, MonthName, Week Number, etc.

Foreign keys in the fact table

When modeling the facts, foreign keys for dimension tables have the values of surrogate primary keys, and it is useful to assume that foreign keys are always defined, or that their values are not equal to Null. To deal with cases in which for a fact record the dimension value may be unknown, a common solution is to add into the dimension table a special record with an attribute, different from the primary key, with a default value such as “Not Found”, and then the foreign key of fact record without the dimension data points to the row “Not Found”.

As in the case of dimension tables, the fact table too may have descriptive attributes, such as the primary key used in the operational database to know the source of the fact.

Changing dimensions

For a slowly changing dimension, we adopt the following solutions on the basis of the strategy specified in the requirements:

Type 1 (overwriting the history) The new attribute value replaces the old value in the record of the dimension table.

Type 2 (preserving the history) A new record is inserted in the dimension table, with a different surrogate key. For example, if a customer changes residence, a new record is inserted in the dimension Customer, as if there were two customers with different surrogate keys. The orders of the past relate to the customer with the old residence, the orders of the future will refer to the customer with the new residence.

This solution is an example that motivates the use of surrogate primary keys, but creates a problem for the analysis that requires counting the number of different customers who have made a certain order: if the count in the analysis is done with a `COUNT(DISTINCT CustomerFK)`, customers who have changed address would be counted several times and, therefore, the result would not be correct.

The problem is solved by adding an attribute to the dimension table with a value appropriate to establish that records with different surrogate keys relate to the same customer who changed residence. Possible solutions are (Figure 3.4): (a) use a customer “natural” key different from the surrogate, like the *Social Security* number (SSN), (b) use the first surrogate key that was assigned to the customer and, (c) to avoid having to do some frequent analysis with junctions, this information is stored in the fact table as a degenerate dimension.

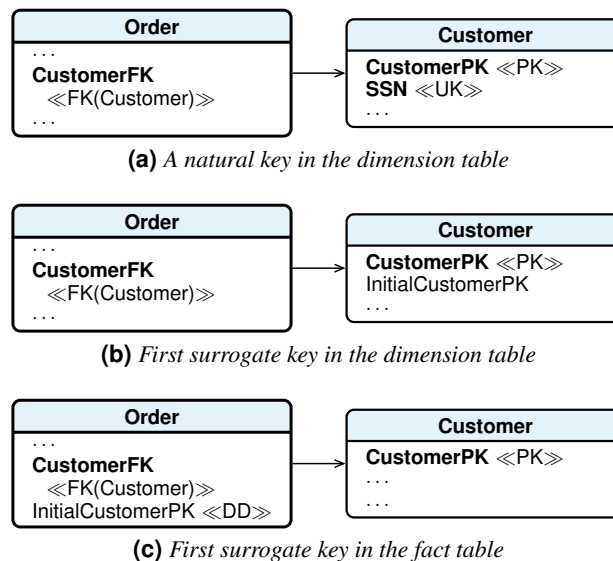


Figure 3.4: Slowly changing dimension

Type 3 (preserving one or more versions of history) In the dimension two attributes are added, one for the new value and the other for the modification date. For example, if a customer changes residence, the Customer dimension structure three attributes are used for the residence: Residence, NewResidence, ChangeDate. If the residence changes again, a new record may be inserted as the solution of the Type 2. Other solutions are possible on the basis of expected analysis, but they all make the solution and the queries for the analysis more complex, and their use should be considered carefully.¹

If a dimension changes frequently due to numerical attributes, an alternative to consider to the previous ones is the following:

1. For examples see http://en.wikipedia.org/wiki/Slowly_changing_dimension

Type 4 (fast changing dimensions) A dimension is considered to be a fast changing dimension if one or more of its attributes changes frequently and in many rows, such as age or income. A fast changing dimension can grow very large if we use the Type-2 approach to track numerous changes. Fast changing dimensions are also called *rapidly changing dimensions*.

An appropriate approach for handling very fast changing dimensions is to break off the fast changing attributes into one or more separate dimensions, called *mini-dimensions*. For example, the dimension is stored in two tables, one with the attributes that do not change (or change slowly) and the other with only those attributes that change frequently, and defined by range of values (e.g. with strings like “From-To”), agreed with users based on the type of analysis to be done. The fact table would then have two foreign keys: one for the primary dimension table and another for the fast changing attributes.

Shared Hierarchies

If there is a shared hierarchy, its attributes are stored in a separate table.

Recursive Hierarchies

If a dimension contains a recursive hierarchy, there is a problem because of limitations in some versions of SQL that does not allow us to define recursive queries. For example, let us assume that in the dimension Agent of Order there is an attribute Supervisor to represent a recursive relationship among agents (Figure 3.5).

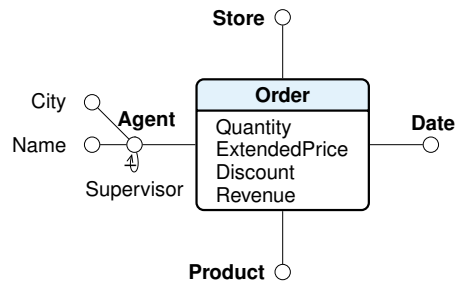


Figure 3.5: A dimension with a recursive hierarchy

If interested in analyzing the total number of orders placed by each agent, including subordinates for which it is responsible at every level, this can not be expressed with SQL versions that do not allow recursion or with BI systems for data analysis and generating reports that are not able to generate it using the relational representation of Figure 3.6a.

To solve the problem a solution is used that involves the following tables (Figure 3.6b):

- The table Agent, also without the attribute SupervisorFK, has one row for each agent.
- The fact table Order with the foreign key AgentFK for Agent.
- The auxiliary table ForTheHierarchy, called *bridge table*, contains one row for each pair of (Supervisor, Subordinate), as well as a row for each agent with itself, and has the following structure:
 - SupervisorFK, a foreign key for the table Agent that represents the supervisor agent.

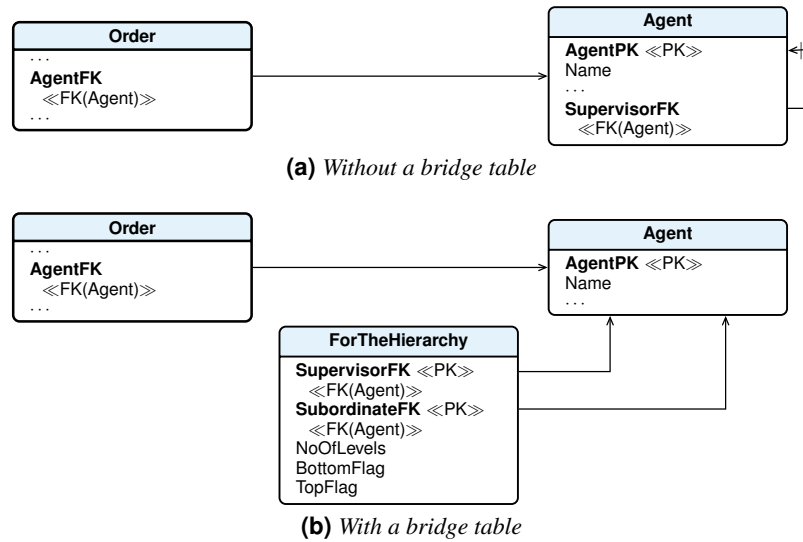


Figure 3.6: A bridge table to represent a recursive hierarchy

- SubordinateFK, a foreign key for the table Agent that represents the subordinate agent or itself (SubordinateFK = SupervisorFK). (SupervisorFK, SubordinateFK) is the primary key.
- NoOfLevels with a value of the number of nodes, minus one, of the path from the supervisor to the subordinate.
- BottomFlag, with value T if the subordinate is not the supervisor of others, otherwise F.
- TopFlag, with value T if the agent (SubordinateFK = SupervisorFK) is at the top of the hierarchy, and so does not have a supervisor, otherwise F.

For example, the records of the table ForTheHierarchy, for the agents hierarchy of Figure 3.7, are the following:

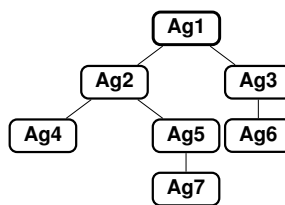


Figure 3.7: A hierarchy for an organization chart

ForTheHierarchy				
SupervisorFK	SubordinateFK	NoOfLevels	BottomFlag	TopFlag
1	1	0	F	T
1	2	1	F	F
1	3	1	F	F
1	4	2	T	F
1	5	2	F	F
1	6	2	T	F
1	7	3	T	F
2	2	0	F	F
2	4	1	T	F
2	5	1	F	F
2	7	2	T	F
3	3	0	F	F
3	6	0	T	F
4	4	0	T	F
5	5	0	F	F
5	7	1	T	F
6	6	0	T	F
7	7	0	T	F

The curious thing is how to use the tables in queries, in particular the bridge table ForTheHierarchy, to move up or down in the hierarchy: the joins operations are between the ForTheHierarchy table and the table Agent, using primary and foreign keys, and between the ForTheHierarchy table and the fact table Order, using the foreign key AgentFK for agent and a foreign key of the ForTheHierarchy table. Note that between the fact table Order and the bridge table ForTheHierarchy there is a *many-to-many relationship*.

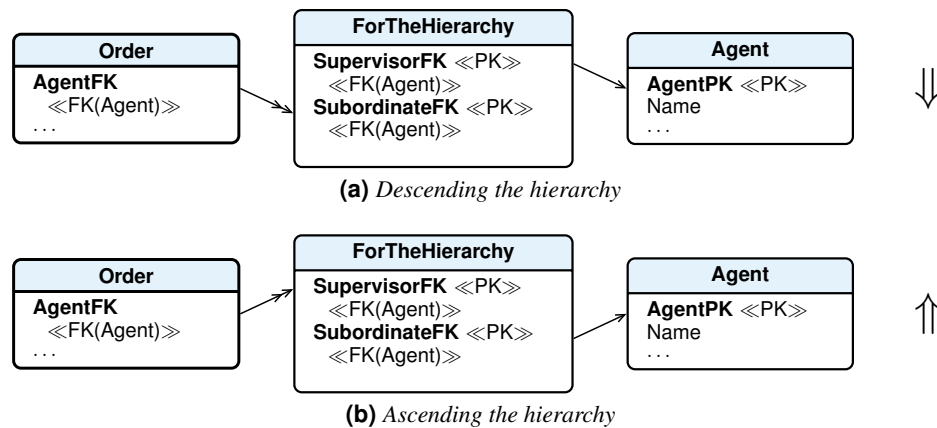


Figure 3.8: Ways of joining a bridge table to fact and dimension tables

- For ascending the hierarchy, the joins are made as shown in Figure 3.8b. For example, to generate the report to find out the total order revenue for agent 6 and all its supervisors (3 and 1), the necessary SQL query is:

```

SELECT  A.Name, SUM(Revenue)
FROM    Order O, ForTheHierarchy H, Agent A
WHERE   O.AgentFK = H.SupervisorFK AND H.SubordinateFK = A.AgentPK
        AND A.Name = 'Ag6'
GROUP BY A.Name;

```

- For descending the hierarchy, the joins are made as shown in Figure 3.8a. For example, to generate the report to find out the total order revenue for agent 2 and all its subordinates (4, 5, and 7), the necessary SQL query is:

```

SELECT   A.Name, SUM(Revenue)
FROM     Order O, ForTheHierarchy H, Agent A
WHERE    O.AgentFK = H.SubordinateFK AND H.SupervisorFK = A.AgentPK
           AND A.Name = 'Ag2'
GROUP BY A.Name;

```

The data hierarchy may be restricted to those of a certain level ($NoOfLevels < 2$), to the bottom of the hierarchy ($BottomFlag = 'T'$), to the top of the hierarchy ($TopFlag = 'T'$) etc.

The disadvantages of this approach are (a) the bridge table data is complex to maintain, and (b) the bridge table design is too complex to be used directly by the users.

Multivalued Dimensions

If there is a multivalued dimension, for example, an order item has been promoted by several agents (Figure 3.9), one of the following relational representations might be used (other solutions are presented in [Song et al., 2001]):

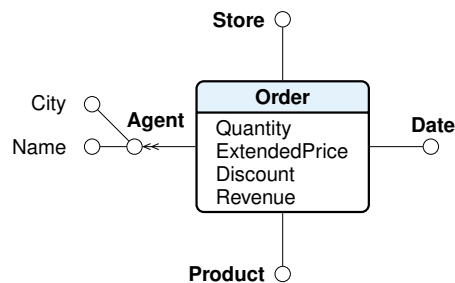


Figure 3.9: A multivalued dimension

1. The *many-to-many relationship* between the fact table and the dimension table is represented with a traditional auxiliary table (Figure 3.10a), called even in this case *bridge table*. This solution, however, violates the properties of a scheme to be a star and may not be accepted by some BI systems (e.g. SQL Server Analysis Services).
2. The *many-to-many relationship* between the fact table and the dimension table is represented with another type of an auxiliary table GroupMembers (Figure 3.10b), with attributes
 - Group, the code of a group of agents,
 - AgentFK, the foreign key for the table Agent, and
 - Allocation with a value between 0 and 1, which represents the contribution to the order promotion credited to each group member, such that the sum of all the allocation factors belonging to a single group is exactly 1.

The table GroupMembers has, for each group, as many elements as are the agents of the group. An agent may be present in several groups. This solution may not be accepted by some BI systems (e.g. SQL Server Analysis Services).

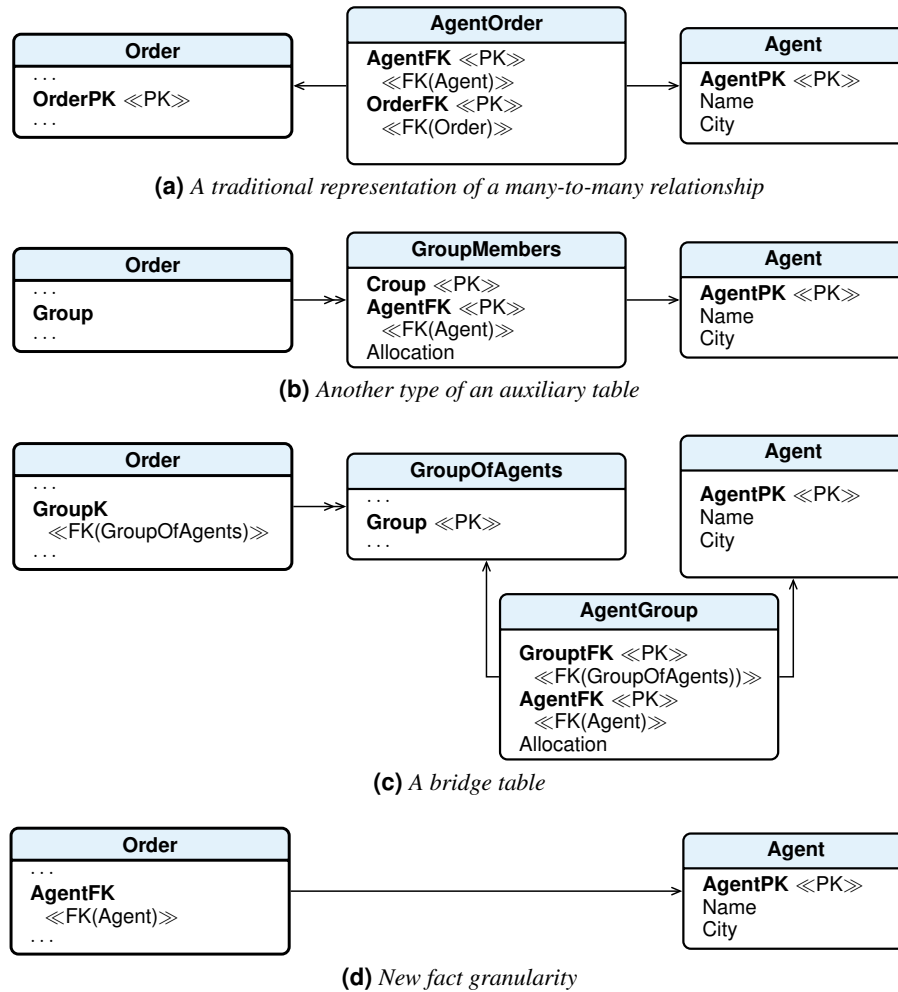


Figure 3.10: An auxiliary table to associate any number of agents with an order

In the fact table *Order* there is the attribute *Group* which indicates the group of agents of the table *GroupMembers* involved in a particular order. The relationship between the fact table and the *GroupMembers* table in Figure 3.10b is *many-to-many*. This is not a mistake: if the same group of agents collaborate again for another order, the same group number will be used.

To generate the report to find out the *total order revenue by agent name*, to avoid a wrong SQL query, we must distinguish whether we are looking for the *total order revenue contribution of a group member (weighted analysis)*, or if we are looking for the *total order revenue of the groups to which an agent belongs (impact analysis)*. In the first case the value of the attribute *Allocation* must be used as follows:

```

SELECT  A.AgentPK, A.Name, SUM(Revenue * GM.Allocation)
FROM    Order O, GroupMembers GM, Agent A
WHERE   O.Group = GM.Group AND GM.AgentFK = A.AgentPK
GROUP BY A.AgentPK, A.Name;

```

while in the second case the attribute *Allocation* is not used, and in general a differ-

ent result is found.

3. Two auxiliary tables are used (Figure 3.10c): GroupOfAgents, which contains one row for each group of agents, with the primary key Group, and AgentGroup to model the *many-to-many* relationship between the tables Agent and GroupsOfAgents.

In the fact table Order there is the foreign key GroupFK for the table GroupOfAgents which indicates the group of agents involved in a particular order. This solution is usually accepted by BI systems (e.g. SQL Server Analysis Services).

4. Another solution is to change the fact granularity: instead of using a record for each order item, a record is used for each agent who has promoted the order item, with the weighted values of the measures, and the attribute Group as a degenerate dimension to recognize groups of records that relate to the same order item (Figure 3.10d). This solution increases the memory occupied by the fact table by a factor equal to the average number of agents that promote an order, while the one with a *bridge table* in general uses less memory.

Multivalued Dimensional Attributes

If a dimension has multivalued attributes (Figure 3.11), the problem is solved as in the case of the multivalued dimension, by treating the dimensional table as the fact table in the previous case, and preserving the relationship between the fact table and dimension table. Figure 3.12 shows only the first solution.

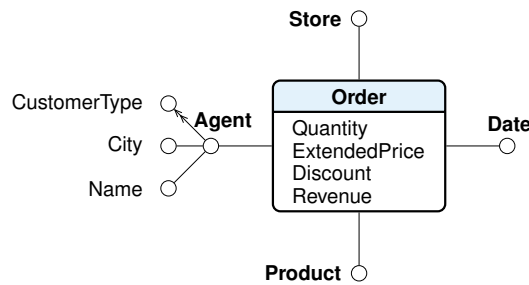


Figure 3.11: A dimension with a multivalued attribute

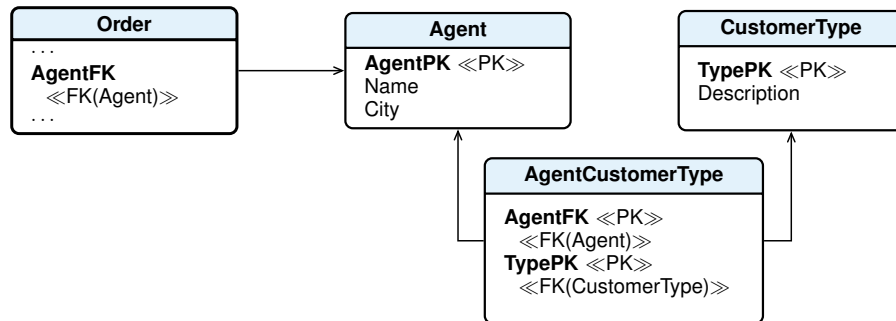


Figure 3.12: An auxiliary table to deal with a multivalued dimensional attribute

3.3 A Case Study

A case study is presented to show how to apply the methodology to design a DW. Do not be misled by the simplicity of the example. In practice the procedure is complicated by the difficulties of the requirements analysis phase for the quantities of the details and the many exceptions that usually occur. Aspects that are neglected in the example.

3.3.1 Requirements Analysis

1. Requirements gathering

- (a) *Analysis of the nature and purpose of the company.* CelPhone is a company that deals with the production and sale of cellular phones with its own sales outlets.

To meet growing market demand the company has expanded by opening new plants and sales outlets. The company growth has started to level off, and management is refocusing on the performance of the organization using a DW to facilitate the analysis of the inventory and revenue from product sales. It has created a team consisting of one data analyst, one process analyst, one manufacturing plant manager, and a sales manager for the project.

- (b) *Business processes analysis.* The products are available in different models and are constructed from a set of common components. Each model may be eligible for discounting, and in this case the salesperson may discount the price if the customer buys a large quantity of the model or a combination of models. The discount must be approved by the manager of the sales outlet.

The plant keeps an inventory of the product models. When the quantity on hand for a model falls below a predetermined level, a work order is created to cause more of the model to be manufactured.

A customer places orders from a sales outlet. Unless a discount is negotiated, the suggested retail price is charged. Each sales outlet, on average, creates 500 orders per day, seven days per week. Each order consists of an average of 10 product models.

- (c) *Collection of business requirements for data analysis (business questions) and verification that the data needed are available.* Let us assume that the expert in DW, after an analysis of the life cycle of a product, inventory and sales processes, organization structure, the meaning of cost, discount and revenue, has interviewed executives interested in the data analysis, and has collected the following set of typical online data analysis of users interest:

Inventory process	
1	Average quantity on hand and reorder level for each model by month, by model identifier and description, by manufacturing plant, name and region.
2	Models that have reached the reorder level at least once in all manufacturing plants of a certain region.

Sales process	
3	The total cost and revenue by model sold, by month, by manufacturing plant, name and region.
4	Percentage of models eligible for discounting, and of those, what percentage are actually discounted when sold, by sales outlet, for all sales this week (or this month, or this quarter).
5	The top five models sold last month by total revenue, or by quantity sold, or by total cost.
6	Total cost and revenue by model Id and description, by month of the last year, by sales outlet and region.
7	Number of customers who last month bought the 5 models that have produced the highest margin, by the region of the sales outlet.

The operational database, which contains all the necessary information for the analysis, is shown in Figure 3.13.

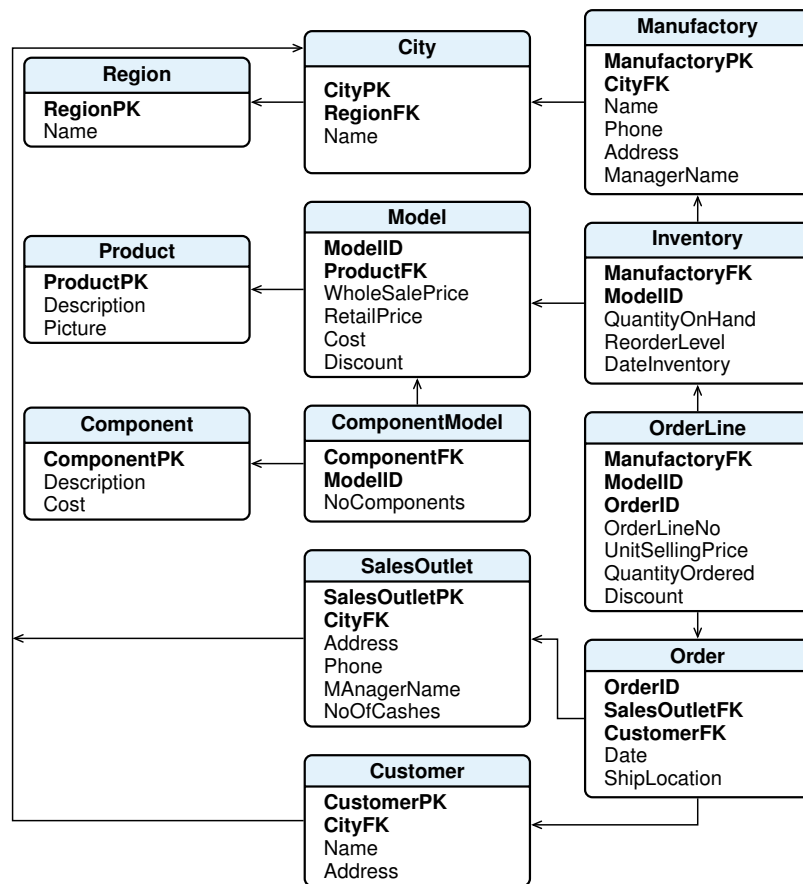


Figure 3.13: The operational database

2. Requirements specification

Each business question is analyzed to identify the preliminary dimensions (in parentheses, the attributes) and measures of interest (for brevity, the metrics are not described), and then the grain of the fact.

			Inventory process	
N	Business questions	Dimensions	Measures	
1	Average quantity on hand and reorder level for each model by month, by model identifier and description, by manufacturing plant, name and region.	Model (ModelID Description), Manufactory (Name, Region), Date(Month)	QuantityOnHand, ReorderLevel	
2	Models that have reached the reorder level at least once in all manufacturing plants of a certain region.	Model, Date(Week), Manufacturing(Region)	ReorderLevel	

With regard to the granularity of the facts, for the Inventory the data of interest are those about each product at the end of the month.

			Inventory fact	
Description	Preliminary Dimensions	Preliminary measures		
A fact is about each product state at the end of the month.	Model, Manufactory, Date	QuantityOnHand, ReorderLevel		

The description of dimensions, attributes, and fact *Inventory* measures follow.

			Dimensions		Date	
Name	Description	Granularity	Attribute	Description		
Date	...	A month	Month	...		
Model	...	A model	Year	...		
Manufactory	...	A manufacturing plant				

		Model		Manufactory	
Attribute	Description	Attribute	Description	Attribute	Description
ModelID	...	Name	...		
Description	...	Region	...		

				Measures	
Measure	Description	Aggregability	Calculated		
QuantityOnHand	...	Semi additive across Date	No		
ReorderLevel	...	Non additive	No		

Sales process			
N	Business questions	Dimensions	Measures
3	The total cost and revenue by model sold, by month, by manufacturing plant, name and region	Model (ModelID, Description), Date(Month), Manufactory (Name, Region)	ExtendedCost, Revenue
4	Percentage of models eligible for discounting, and of those, what percentage are actually discounted when sold, by sales outlet, for all sales this week (or this month, or this quarter)	Model(Discount), SalesOutlet, Date (Week, Month, Quarter)	ExtendedPrice, Discount
5	The top five models sold last month by total revenue, or by quantity sold, or by total cost.	Model, Date(Month)	ExtendedCost, QuantityOrdered, Revenue
6	Total cost and revenue by model id and description, by month of the last year, by sales outlet and region	Model (ModelID, Description), SalesOutlet(Region), Date(Month, Year)	ExtendedCost, Revenue
7	Number of customers who last month bought the 5 models that have produced the highest margin, by the region of the sales outlet.	Customer, Model, SalesOutlet(Region), Date(Month)	Margin

With regard to the granularity of the facts, for the Sales the data of interest are those about each single *line item* of an order.

Sales fact		
Description	Preliminary Dimensions	Preliminary Measures
A fact is about a product sold	Model, Manufactory, Customer, SalesOutlet, Date	QuantityOrdered, ExtendedPrice, Revenue, ExtendedCost, Discount

The description of dimensions, attributes, and fact *Sales* measures follow.

Dimensions			Model	
Name	Description	Grain	Attribute	Description
Model	...	A model	ModelID	...
Date	...	A day	Description	...
Manufactory	...	A plant	Discount	...
Customer	...	A customer		
SalesOutlet	...	A sales outlet		

Date		Manufactory	
Attribute	Description	Attribute	Description
Day	...	Name	...
Month	...	Region	...
Quarter	...		
Year	...		
Week	...		

Customer		SalesOutlet	
Attribute	Description	Attribute	Description
		Region	...

Dimensional Hierarchies		
Dimension	Description	Hierarchy type
Date	Day → Month → Quarter → Year	Balanced

			Measures	
Measure	Description	Aggregability	Calculated	
QuantityOrdered (Q)	...	Additive	No	
ExtendedPrice (P)	$\text{UnitPrice} \times Q$	Additive	Yes	
ExtendedCost (C)	$\text{UnitCost} \times Q$	Additive	Yes	
Discount (D)	ExtendedPrice reduction	Additive	No	
Revenue (R)	$P - D$	Additive	Yes	
Margin	$R - C$	Additive	Yes	

Before moving on to other phases of design, dimensions and measures of the facts are represented in the following tabular form highlighting what measures and dimensions are common to different facts and therefore need to be conformed or renamed. The dimensions Date and Model have different attributes in the two facts, and it is assumed that those of the fact Sales are used.

Fact dimensions		
Dimension	Inventory	Sales
SalesOutlet		X
Model	X	X
Manufactory	X	X
Customer		X
Date	X	X

Fact measures		
Measure	Inventory	Sales
QuantityOnHand	X	
ReorderLevel	X	
ExtendedPrice		X
ExtendedCost		X
Revenue		X
Margin		X
QuantityOrderd		X
Discount		X

3.3.2 Initial Analysis-driven Data Mart Conceptual Design

The data analysis requirements show that the facts are about Inventory and Sales. The attributes used in the data analysis suggest that the two possible initial conceptual data marts are those shown in Figure 3.14.

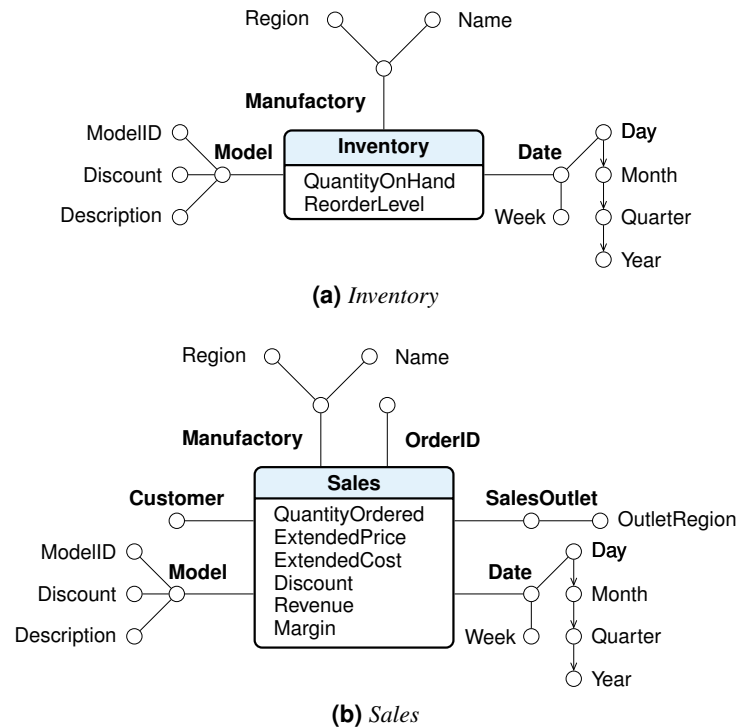


Figure 3.14: Initial data mart conceptual designs

3.3.3 Candidate Source-driven Data Mart Conceptual Design

1. **Operational data analysis.** Based on the data analysis requirements, the relational schema is examined to decide which tables and attributes are interesting, other than primary and foreign keys.

As for the tables, we observe that ComponentModel and Component contain information that is not relevant for the purposes of data analysis.

As for the attributes of other tables, the following considerations apply:

- Region: the attribute Name is of interest.
- City: although not explicitly required, it is good to retain the information about the city, because, as we shall see later, the table City is related to the table Region via a *many-to-one relationship*, and because, in principle, it is always better to consider some more information, potentially relevant, than what is strictly necessary.
- Manufacture: the attributes Phone, Address and ManagerName are not relevant for data analysis; the pertinent attribute is the geographic location (City and Region).
- Inventory: all attributes are of interest for the data analysis of process *Inventory*.

- Model: the relevant attributes are ModelID, Discount and Description.
- Product: the attributes of this table are not of interest for data analysis.
- OrderLine: the attributes of this table are important for data analysis; Discount is a percentage.
- Order: ShipLocation is not useful for our purposes, while Date is of interest.
- SalesOutlet: the relevant attribute is the geographic location (City and Region).
- Customer: we are interested in Name and the geographic location (City and Region).

In this first analysis the keys to the tables were deliberately neglected, because they will be considered later. Once the relevant information has been chosen, we proceed to the next phase of design, namely the classification of entities.

2. **Entity classifications.** We classify the tables in the relational schema based on their content and the relationships between them:

- **Transaction entity:** Recalling the definition of transaction entity, it is quite easy to fit into this category the tables Inventory and the merging of OrderLine and Order, with the granularity of OrderLine. They, in fact, (a) describe events that occur frequently at certain dates and (b) contain numerical attributes that represent possible measures of interest for the analysis of data. Note that the table Model contains other numerical attributes, but they are not relevant for data analysis.
- **Component entity:** They are the tables related to a transaction entity via a *one-to-many relationship*. Analyzing the relational schema it is discovered that
 - for the transaction entity Inventory, the component entities are Manufactory and Model,
 - for the transaction entity OrderLine, merged into Order, the entities component are Customer, SalesOutlet and Inventory.

Finally, as mentioned earlier, among the entities of each entity component, we add the time entity (present in the relational schema with attributes of type Date).

- **Classification entities:** These are the tables related to a component entity by a chain of *one-to-many relationships*. Their interesting attributes are added to those of minimal component entity.

For the entity Manufactory, component of Inventory, the classification entities of interest are City and Region.

For the entity SalesOutlet and Customer, component of OrderLine, the relevant classification entities are City and Region.

For the entity Inventory, component of OrderLine, the relevant classification entities are Models and Manufactory, with City and Region. Since in the requirements analysis of *Sales* process there is no interest in Inventory attributes, the classification entities Models and Manufactory are treated as components of the entity event OrderLine.

3. **Definition of the candidate data mart conceptual designs.** Having identified two interesting event entities, we proceed with the definition of two conceptual designs for the data marts with the relative dimensions (Figure 3.15).

As far as the dimensional hierarchies are concerned, the case contains only one which is explicit, that between City and Region.

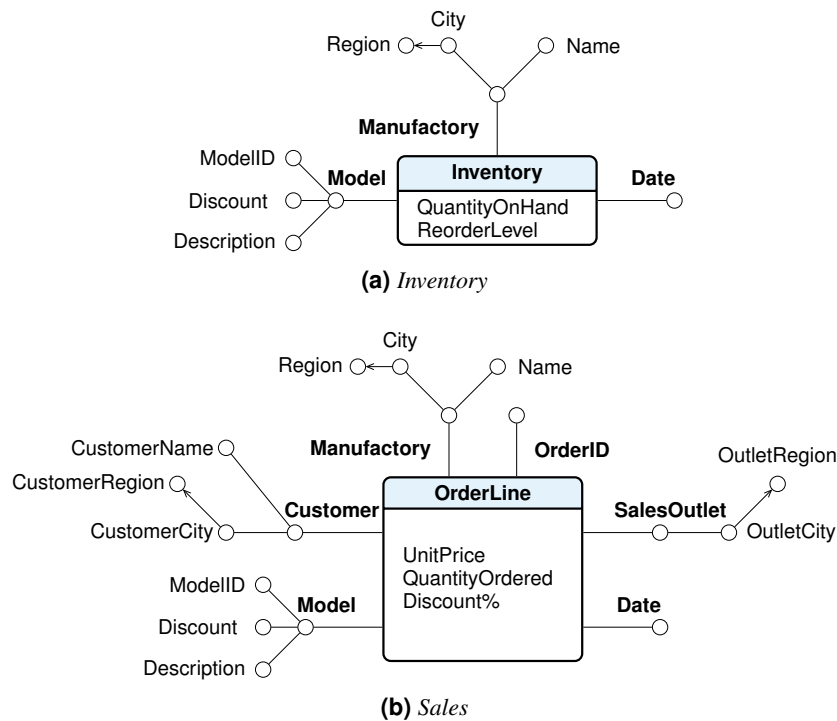


Figure 3.15: The candidate data mart conceptual designs

4. **Analysis of data marts fact granularity and measures additivity.** This step produces no information other than that already known with regard to the fact *Inventory*. For the fact *OrderLines*, however, we note that the measures *UnitPrice* and *Discount%* are not additive, and so in the final step of the conceptual design, the solution used to the fact *Sales* is preferred.

3.3.4 Final Data Mart Conceptual Design

From a comparison of candidate designs and the initial ones, the terminology is unified and the final designs are those of Figure 3.16.

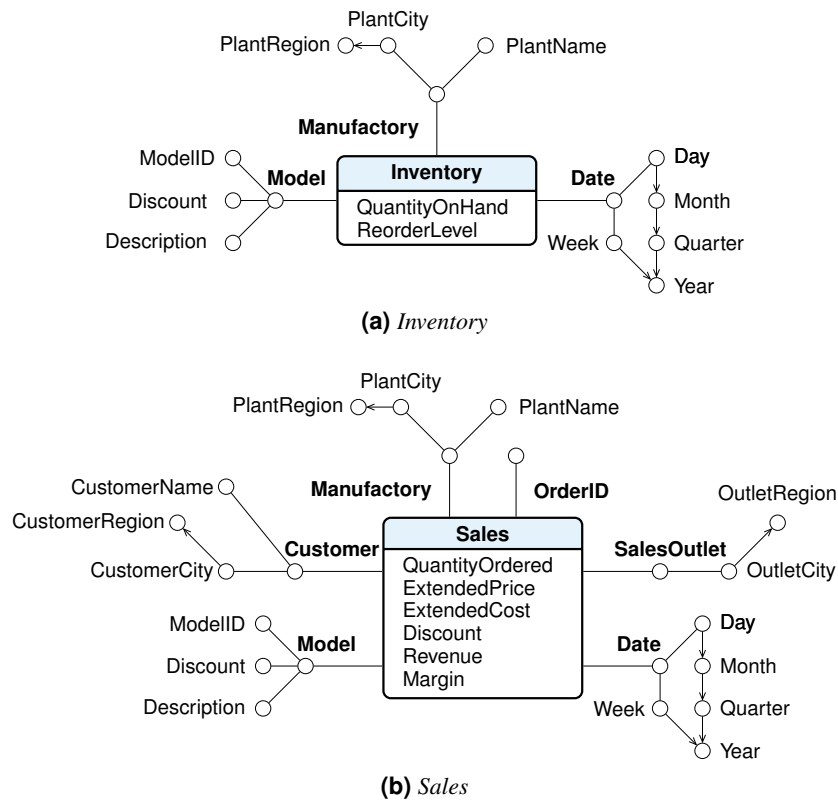
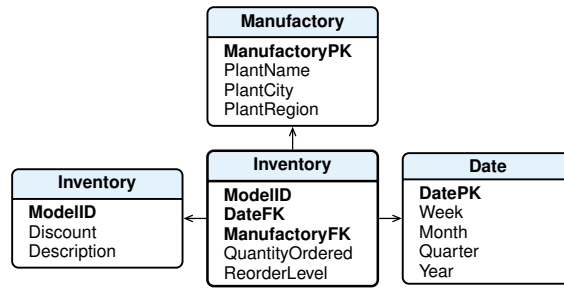


Figure 3.16: The final data mart conceptual designs

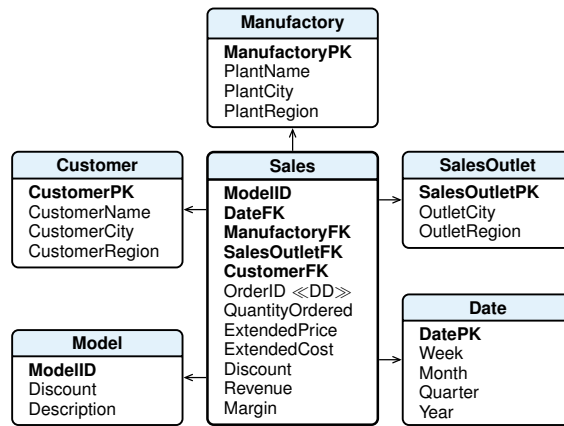
3.3.5 Data Mart and Data Warehouse Logical Design

Two data mart star schemas are defined with a different fact table for each conceptual design, while for each dimension a table is defined in association with the fact table, by defining appropriate surrogate primary keys and foreign keys (Figure 3.17a,b).

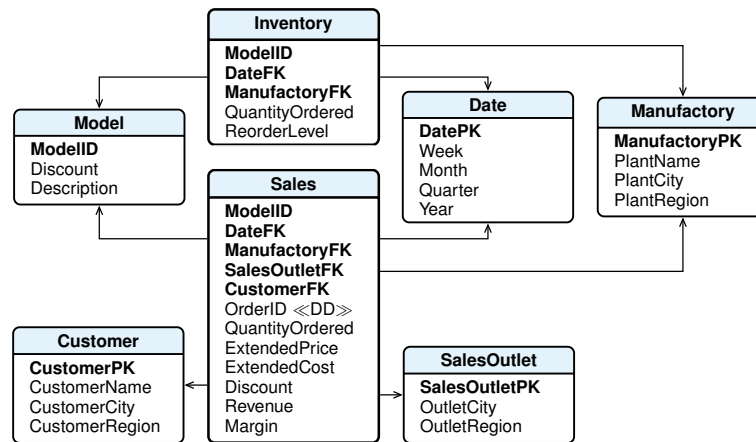
The data marts relational schemas are then integrated to define the DW schema. Note that the two star schemas share the dimensions Date, Model and Manufacture. The structure of the DW is shown in Figure 3.17c.



(a) Inventory data mart star schema



(b) Sales data mart star schema



(c) Data warehouse constellation schema

Figure 3.17: The data mart star schemas and the data warehouse constellation schema

3.4 Project Quality Control

Let us consider some checks for the final review of a project to improve its quality.

Conceptual Design

1. Granularity of the facts

The determination of the grain of the facts is the first key step in the design of a data mart. Choosing the grain means deciding the meaning of a fact, and so the pertinent measures and dimensions.

For example, for orders with multiple lines, if the granularity is the order, it makes sense to consider a measure that concerns the total order value, but not the amount of product ordered on each order line, and so the Product dimension can not be used, but if the granularity is the order line, then it is meaningful to consider measures about the quantity of product ordered, the price and the revenue (but not the total order value), and the Product dimension.

In the data mart conceptual design, the measures have numerical values that can be added across dimensions. Descriptive attributes must be modeled as *degenerate dimensions*.

The grain decision for the facts also determines the grain of each dimensions. For example, if the grain for the PropertySales is an individual property sale, the grain of the Client dimension is the detail of the client who bought a particular property.

2. Measures

Measures are numerical quantities useful for evaluating the performance of the processes to be analyzed. It is also useful to define measures that can be calculated from others, at the time of loading the data. This is particularly true for values fundamental in the analysis, such as revenue and margin, to avoid their being calculated by users incorrectly, or in different ways, within reports at the time of the analysis. If the answers are wrong or inconsistent, the data warehouse will be viewed as wrong.

It is also important to document, as part of the conceptual design, whether they are *additive*, *semi-additive* or *non-additive*, to avoid very common mistakes when they are summed up.

Another common error to avoid is modeling unit amounts (e.g. unit price) as measures rather than extended amounts (e.g. extended price). This does not mean that unit amounts must be excluded from conceptual design, because they may be valuable information for analysis. If there is not a dimension where they can be stored, they may be placed in the data mart conceptual design as *degenerate dimensions*.

The most useful measures are those additives that can be aggregated with any type of function and by combining facts with various dimensions to answer common business questions.

The most critical measures, often to be avoided, are the *non-additive* because they can not be aggregated with the sum. Typical examples are measures defined as rates or percentages. These measures must be broken down into underlying components that are additive, to calculate the *ratio of the sums*, not the sum of the ratio. For example, the *margin rate* is the ratio of the margin to revenue. These components are fully additive, and they are usually defined as measures that can be safely aggregated to any level of detail. The non additive margin rate is computed in a query, or by additional processing logic in the reporting tool, as the ratio of the sums of margin and of revenue.

3. Date and Time

They should always be modeled separately as dimensions. They are modeled as facts descriptive attributes only when are not used for analysis.

4. Dimensions quality

The dimensions should be chosen considering the user's need for examining facts, and the future development of the DW. If the same dimension appears in multiple data marts they must be defined in the same way to be shared (*conformed dimensions*). Examples of these dimensions are time, date, customer and product.

- (a) The dimensional attributes must be useful (a) to analyze facts (restrictions, groupings, and aggregations) and (b) to produce summary reports with the headers using the users vocabulary to facilitate understanding.
- (b) The names of the attributes of different dimensions should be different: a way to disambiguate them is to prefix the attribute with the dimension name. Attribute names should not be those used in the operational databases, but those that are used in the analysis and that appear in reports.
- (c) Dimensional attributes already represented as numerical measures of the facts should not be repeated in dimensions.
- (d) The values of dimensional attributes, usually strings of characters, should facilitate the interpretation of the reports: avoiding codes (0/1, F/M, etc.) or adding attributes that describe them.
- (e) Represent as a string data type attributes such as Date and Address only if there is no interest in exploiting in the analysis the implicit hierarchies among their attribute values.
- (f) If a fact is associated with more elements of a dimension, it must be modeled as multi-valued.
- (g) The dimensions can be *degenerate*, that is they are without attributes because their values are numbers, such as the order number, invoice number etc., or strings of characters.
- (h) The dimensional hierarchies should always be present to make the analysis more useful at several levels of detail.

Logical Design

1. Surrogate keys

Surrogate keys must be used for dimension tables, which may also include the primary key of the source data.

2. No attributes with null values

Default attribute values must be set in the database schema to avoid nulls in the database. The default value for all the fact measures must be zero in the schema. In SQL, NULL plus a number equals a NULL, and the aggregate functions perform a NULL-elimination step, so that NULL values are not included in the final result of the calculation. The only aggregate function that does not implicitly eliminate NULL is the COUNT(*) function. However, an aggregate function AGG(A), with A a set of NULL, returns NULL, while COUNT(A) returns zero.

To deal with cases in which for a fact record the dimension value may be unknown, the dimension table must have a special record with an attribute value "Not Found", and when a fact record is missing a dimension data, the foreign key value is the surrogate key of the record "Not Found".

3. Degenerate dimensions

In the logical schema degenerate dimensions become attributes of the fact table as the foreign keys to other dimensions, if the attributes take up little storage space, otherwise they are stored in separate dimension tables.

Another type of degenerate dimension arises when there are a few attributes that

take different values (status indicators), e.g. order line with status (closed, open, canceled), customer satisfaction, type of delivery, payment terms etc. These attributes could be added (a) in the fact table, increasing its memory size, (b) in the relevant dimension tables, duplicating the records (if a customer pays in cash, by bank transfer, credit card, three records are needed), (c) in different dimension tables of small cardinality, by increasing the number of foreign keys in the fact table, (d) collect them all into a separate dimension table with as many attributes as there are fields with discrete values, and as many elements as are the possible combination of values (*junk dimension*). The latter solution is preferable in the presence of several attributes of state indicators used in the fact table or in different dimension tables.

4. Shared data

If the shared data are dimensional hierarchies, such as geographic hierarchies, in the logical design a way to treat them is to deconstruct the dimension table into a tree structure. So a snowflake dimension is defined, and its advantage must be evaluated considering the savings in space, the greater complexity of the scheme, the execution time analysis, and any ambiguity in analysis, such as “Analysis of the total revenue for the city”: What city does it refer to? Customer, agent, or the warehouse city?

To facilitate understanding of the data mart schema, and to avoid ambiguous analysis, if the tables are small, it is usually preferable to duplicate shared hierarchies in the dimension tables using different attribute names.

If the shared data are dimensions, or there are more dimensions with different attributes that have the same values (e.g. two dates with the same attributes day, month, year, or with different attributes to highlight the role of different dates), such as OrderDate and ReceivedDate, another solution may be used. Instead of using two separate date tables, with the same granularity, two views are created, with different attribute names, from a single Date table.

5. Dimensional hierarchies.

Check that the hierarchies type (balanced, incomplete or recursive) is correctly represented. Moreover, verify that functional dependencies hold over a loaded dimension table with dimensional hierarchies. For example, let the dimension be Date(PkDate, Month, Quarter, Year). If the dimensional hierarchy Month → Quarter is valid, then following query returns an empty result set.

```
SELECT    Month
FROM      Date
GROUP BY Month
HAVING   COUNT(DISTINCT Quarter) > 1;
```

6. Snowflake dimensions

Do not normalize (snowflake) dimension tables, since it will be harder for the users to analyze data. Moreover, in general, there is not a very significant memory saving because of the relatively small cardinality of the dimension tables. Snowflakes are meaningful only when it is necessary to define interesting dimension tables shared among several data marts.

7. Changing dimensions

Recognize the dimensions with attributes that change over time and treat them appropriately.

3.5 Summary

- The design of a data warehouse to support business decisions is a complex task that requires a methodology organized into phases, like that used to design operational databases, but the phases' objectives must be revised properly to adapt them to multidimensional modeling.
- A possible design methodology has been presented, with the documentation to be produced at the end of each phase, to proceed by considering both the requirements analysis and the operational database available.
- The logical design phase has been presented to highlight some critical aspects of the transition from the conceptual design to the relational one, especially for treating dimensions that change over time, multivalued dimensions and multivalued dimensional attributes.
- Finally, some controls have been listed for the final review of a project to improve its quality.

DATA ANALYSIS

Once the data warehouse has been implemented, the final step is *data analysis*, that is to identify and to develop a suite of reports showing how the information provided in these reports can be used by decision makers to improve the business process modeled. Data analysis is usually done interactively with tools that provide graphical interfaces to make the requests, which are then translated automatically into SQL queries on the data warehouse. To facilitate the implementation of complex analysis, the SQL language has been extended with new operators to group and aggregate data using several analytic functions. Some of them will be presented with examples to show how to express in SQL basic OLAP operations on multidimensional data.

4.1 OLAP Systems Solutions

When talking about systems for data analysis, terms are used such as OLAP, ROLAP, MOLAP, HOLAP, DOLAP, OLAP Server, OLAP Services, etc., which can create problems of interpretation for the way they are used by vendors of these types of products.

The term OLAP is used to refer to the activity of multidimensional analysis of large amounts of data, with interactive and intuitive ways of changing the perspectives of analysis and moving to different levels of synthesis of the detailed data.

An *OLAP client* provides graphical environments where the business users can click on actions and perform drag-and-drop operations to provide input to summarize data. More experienced users can also create complex queries with languages such SQL or MDX. An *OLAP client* interacts with the data manager using one of the following solutions (Figure 4.1):

- (a) The data warehouse is stored in a relational database system RDBMS (*Data server*) and the interactions with the *OLAP client* occur in SQL. The benefit of this solution is that it uses a standard technology usually already available. In the past the approach was not considered satisfactory for the performance of RDBMS as systems for data warehouses and limitations of the SQL as an OLAP language. But now the main producers of RDBMS systems have made them more and more specialized for OLAP applications (*OLAP-Aware RDBMS*), aware of managing data warehouses with special relational schemas (star, snowflake or constellation), dimensional hierarchies, specific storage structures, and materialized views (also called *aggregates*, *MQT (materialized query tables)*, or *summary data*).

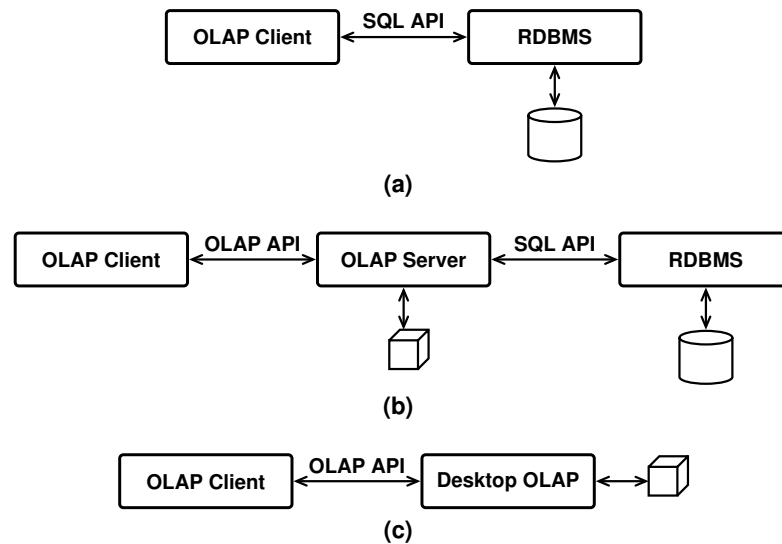


Figure 4.1: OLAP Systems Solutions

(b) An *OLAP client* interacts with an *OLAP server*, a system that provides a multi-dimensional cube vision of a data mart, which can be analyzed with the typical operations slice, dice, drill down, roll up, pivot, etc. An *OLAP server* can be one of the following types:

- **MOLAP** (Multidimensional-OLAP), which stores in the local memory both the data cube, taken from a *Data server*, and the aggregates of the extended cube (materialized views), using a specialized multidimensional arrays structure. A MOLAP server does not support SQL, but proprietary languages not for business users, the most popular being MDX from Microsoft. The solution provides excellent performance, but is not suitable for large amounts of data. Examples of products are Hyperion Essbase, Microsoft Analysis Services, Cognos PowerPlay and DB2 OLAP, using Hyperion Essbase technology.
- **ROLAP** (Relational-OLAP), which stores both the data and the materialized views in the relational *Data server*. ROLAP servers may also need to implement functionality not supported in the SQL of the *Data server*, for example, analytic functions. Examples of products are Informatica, MicroStrategy, Microsoft Analysis Services and SAP BW.
- **HOLAP** (Hybrid-OLAP), which combines ROLAP and MOLAP by splitting storage data in a MOLAP and a relational *Data server*. Splitting the data can be done in different ways. One method is to store the detailed data in the *Data server*, and precomputing aggregated data in MOLAP. Another method is to store more recent aggregate data in MOLAP to provide faster access, and older aggregates in the *Data server*. Microsoft Analysis Services is an example of product that can operate as MOLAP, ROLAP or HOLAP.

Data update is usually done by batch, at predetermined time intervals. There are also systems capable of doing *proactive caching*, updating MOLAP data incrementally at time intervals or after each transaction on an operational database. This permits the use of OLAP in real-time, or near real-time, useful in certain contexts such as, for example, the stock market.

The requests of the *OLAP client* to the *OLAP server* are formulated in SQL or

proprietary languages such as MDX of SQL Server Analysis Services and OLAP DML of Oracle. The results are communicated to the client in proprietary formats or in XML, for example using to the standard XMLA.

- (c) The *OLAP client* interacts with a local **DOLAP** system (*Desktop OLAP*), which manages small amounts of data extracted from the *OLAP server*, the *Data server* or an operational DBMS. The fact that a subset of a data cube is transferred on a user's machine makes it a good choice for those who travel and move extensively, such as sales people, by using portable computers, or who do not regularly perform such complex queries that a faster server is preferred to the speed of the client. The main product of this type is Business Objects.

Among the DOLAP systems there are those specialized for interactive multidimensional analysis, with some limitations as regards the functionalities of the *OLAP server*. For example, Business Objects and MicroStrategy allows the definition of interactive reports with operations such as drill down and roll up. The system does not maintain aggregates in the local memory, but only the results of recent operations. The aggregates are calculated by the *OLAP server*, by the *Data server* or by the operational DBMS.

In all the solutions, the *metadata*, with information on the structure of the fact table, dimensions and hierarchies, are created and maintained by the *OLAP or Data server* and are imported from the *OLAP client* using the standard *CWM (Common Warehouse Metamodel)* or proprietary formats.

Finally, in all the various solutions, the systems are supported by ETL (*Extract, Transform, Load*) tools to load data from operational databases and other external sources.

4.2 Data Analysis Using SQL

In the following, several examples are presented to show how to write SQL queries to produce reports for commonly asked business questions. The examples are based on the following table with attributes without null values.

Sales(Customer, Product, Brand, Date, City, Region, Area, Quantity, Revenue, Margin)

Only in some cases will a graphical representation of the result also be shown, but this is, in general, *essential* to make the results understandable and useful to those who need information for decision support, the main motivation of the multidimensional analysis. Sometimes patterns can be seen in visual data that cannot be seen in numerical data. All reporting tools allow us to perform both an analysis of data without writing the query in SQL, and to produce a graphical representation of the result. Some DBMS, such as Oracle, can produce a graphical representation of the result with analysis expressed in SQL too.

Some of the more commonly used business reports follow.

Simple Reports.

Many kinds of commonly requested business reports can be readily expressed as SQL queries.

- What were the total revenue and margin (in value terms and as a percentage of the revenue) of sales for the month of January 2009, by brand and by product?
- What were the total revenue and margin (in value terms and as a percentage of the revenue) of sales for the month of January 2009, by brand and by product, and the brand subtotals for all products?

Moderately Difficult Reports.

However, many other commonly requested reports cannot be expressed so easily. Reports that require comparisons often challenge both the query writers and SQL itself.

- Revenues for 2009 by brand and product, with the percentage change from the previous year?
- How did product revenues in 2009 compare by geographic area, in a readable spreadsheet, or “cross-tab”, format?
- Which suppliers charge the most for bulk tea products?
- What was the most successful promotion last December in Rome?

Very Difficult Reports Without Analytic SQL.

Reports that require sequential processing are very difficult to express as SQL queries, for example deriving a simple running total. Data analysts typically run several queries with a program, then paste the results together. This approach is awkward because it requires a sophisticated user.

The standard SQL analytic functions provide a better solution because they are easy to use, and perform a broad range of calculations that execute quickly on the server.

- What were the cumulative totals (*running totals*) for Best coffee sales during each month of last year?
- What were the ratios of monthly sales to total sales (expressed as percentages) for Best coffee during the same period?
- Which ten cities had the worst coffee sales in 2010 with regard to dollar sales and quantities sold?
- Which supermarket falls into the top 25% in terms of sales revenue for the first quarter of 2010?
- What products fell into the top 20%, middle 60%, and bottom 20% of sales margins totals for the second week of 2011, at stores in the Center area?

4.3 Simple Reports with SQL

A simple kind of query involves *grouping* and *aggregation* of the data.

To write such kind of **SELECT** the **GROUP BY** clause must be used with the following version of the command syntax.

```
SELECT      DISTINCT  $S_A, S_{AF}$ 
FROM        $T$ 
WHERE       $W_C$ 
GROUP BY    $G_A$ 
HAVING      $H_C$ 
ORDER BY    $O_A$ ;
```

where (a) S_A are the **SELECT** attributes and S_{AF} are the **SELECT** aggregation functions; (b) T are the **FROM** tables; (c) W_C is the **WHERE** condition; (d) G_A are the grouping attributes, with $S_A \subseteq G_A$; (e) H_C is the **HAVING** condition with aggregation functions H_{AF} ; (f) O_A are the **ORDER BY** attributes; (g) the **DISTINCT**, **WHERE**, **HAVING** and **ORDER BY** clauses are optional.

The command semantics with tables R and S , and all the optional clauses specified, in terms of the extended relational algebra is shown in Figure 4.2.

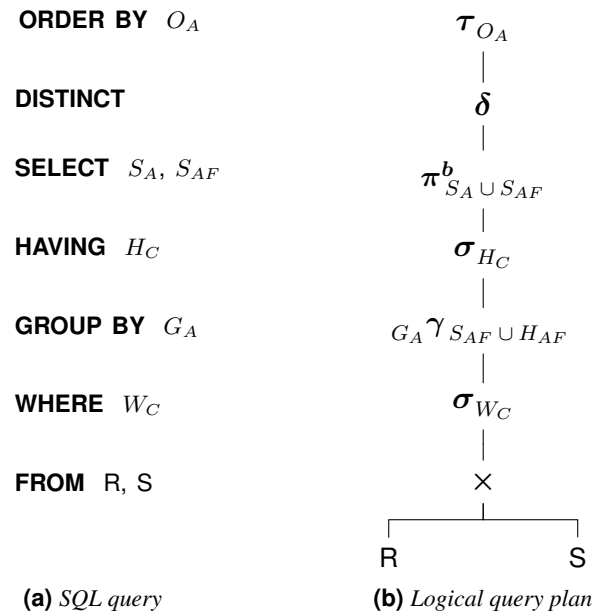


Figure 4.2: SQL query with **GROUP BY** semantics

Figure 4.3 shows an example of the analysis report for “total revenue and margin (in value terms and as a percentage of the revenue) of sales for the year 2009, by brand and by product.”

Margin by Brand and by Product Year 2009				
Brand	Product	Revenue (€)	Margin (€)	Margin% (%)
B1	P1	2 100	273	13
	P2	3 720	624	17
	P3	15 300	1 803	12
B2	P4	12 600	756	6
	P5	22 500	2 196	10
	P6	48 300	4 496	9

Figure 4.3: A Simple Report

The following query with a **GROUP BY** produces the desired result:¹

```

SELECT      Brand, Product, SUM(Revenue) AS Revenue,
              SUM(Margin) AS Margin,
              ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
GROUP BY   Brand, Product
ORDER BY   Brand, Product;
  
```

1. In all the examples in this chapter, the SQL queries produce the data needed for reports, but not their graphic representation. The YEAR, QUARTER, MONTH functions retrieve subfields from DATE values.

Note that the operations *Slice* and *Dice* are expressed by a selection and projection, while *Roll-up* and *Drill-down* require a **GROUP BY**. For example, a *roll-up* on Brand, to find, for the year 2009, the total revenue, total margin and margin percentage by Brand, is expressed with the following query:

```

SELECT      Brand, SUM(Revenue) AS Revenue,
              SUM(Margin) AS Margin,
              ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
GROUP BY   Brand
ORDER BY   Brand;

```

while a *drill-down* on Customer, on the previous result, to find, for the year 2009, the total revenue, total margin and margin percentage by Brand and Customer, is expressed with the following query:

```

SELECT      Brand, Customer, SUM(Revenue) AS Revenue,
              SUM(Margin) AS Margin,
              ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
GROUP BY   Brand, Customer
ORDER BY   Brand, Customer;

```

In general, adding an attribute to the **GROUP BY** and **SELECT**, a *drill-down* is made, while dropping an attribute, a *roll-up* is made.

4.3.1 The Operator *ROLLUP*

Many OLAP queries use a **GROUP BY** to partition data into groups that are reduced to a single row of aggregates and grouping columns. However, standard SQL limit the types of OLAP queries that can be easily expressed. One extension is the **ROLLUP** clause. Suppose that we want to obtain the report of Figure 4.4. Any report that contains a metric is likely to contain a “total” at the end. If the report has more than one dimensional attribute, the metric may also be subtotaled.

Margin by Brand and by Product Year 2009				
Brand	Product	Revenue (€)	Margin (€)	Margin% (%)
B1	P1	2 100	273	13
	P2	3 720	624	17
	P3	15 300	1 803	12
B1	Total	21 120	2 700	13
B2	P4	12 600	756	6
	P5	22 500	2 196	10
	P6	48 300	4 496	9
B2	Total	83 400	7 448	9
Total		104 520	10 148	10

Figure 4.4: A report with some subtotals

A possible solution with standard SQL would be:

```

SELECT      Brand, Product, SUM(Revenue) AS Revenue,
              SUM(Margin) AS Margin,
              ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
GROUP BY   Brand, Product

UNION ALL

SELECT      Brand, NULL AS Product, SUM(Revenue) AS Revenue,
              SUM(Margin) AS Margin,
              ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
GROUP BY   Brand

UNION ALL

SELECT      NULL AS Brand, NULL AS Product, SUM(Revenue) AS Revenue,
              SUM(Margin) AS Margin,
              ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
ORDER BY   Brand, Product;

```

Three statements are required because the report requires three aggregations applied to groups of values produced by a different **GROUP BY** clause. Computing all of these queries independently is time consuming, and this is the main motivation for the **ROLLUP** clause which is included in the SQL of several DBMSs:

```

SELECT      Brand, Product, SUM(Revenue) AS Revenue,
              SUM(Margin) AS Margin,
              ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
GROUP BY   ROLLUP (Brand, Product)
ORDER BY   Brand, Product;

```

■ Definition 4.1 *ROLLUP in SQL:1999*

A **ROLLUP** group is an extension of the **GROUP BY** clause that produces a result that contains subtotal records in addition to the regular grouped records, whose aggregate values are derived by applying the same functions. A **ROLLUP**($A_1, A_2, \dots, A_{n-1}, A_n$) group is equivalent to the union of the $n + 1$ grouping results on the attributes $(A_1, A_2, \dots, A_{n-1}, A_n)$, $(A_1, A_2, \dots, A_{n-1})$, \dots , (A_1, A_2) , (A_1) , and $()$. Notice that each grouping result is created by eliminating an attribute from the list specified in the **ROLLUP** clause, by moving from right to left. Therefore, the order in which the attributes are specified is significant for the **ROLLUP** result. The operator produces its results with just one table access.

For example, the rows of the table in Figure 4.4 are calculated first grouping on Brand, Product, and then the subtotals are calculated progressively moving from right to left through the list of grouping columns: first grouping on Brand, and then on $()$ (super-aggregate rows).

4.3.2 The Operator *CUBE*

Suppose now that we want to obtain a table such as that in Figure 4.5, similar to the table in Figure 4.4 except that, in addition, it has totals for each row and each column.

Margin by Brand and by Product Year 2009				
Brand	Product	Revenue (€)	Margin (€)	Margin% (%)
B1	P1	2 100	273	13
	P2	3 720	624	17
	P3	15 300	1 803	12
Total B1		21 120	2 700	13
B2	P4	12 600	756	6
	P5	22 500	2 196	10
	P6	48 300	4 496	9
Total B2		83 400	7 448	9
	Total P1	2 100	273	13
	Total P2	3 720	624	17
	Total P3	15 300	1 803	12
	Total P4	12 600	756	6
	Total P5	22 500	2 196	10
	Total P6	48 300	4 496	9
Total		104 520	10 148	10

Figure 4.5: Report with subtotals

Again, computing all of these queries independently is time consuming, and this is the main motivation for the **CUBE** clause which is included in the SQL of several DBMSs:

```

SELECT      Brand, Product, SUM(Revenue) AS Revenue,
            SUM(Margin) AS Margin,
            ROUND(100*SUM(Margin)/SUM(Revenue)) AS Margin%
FROM        Sales
WHERE       YEAR(Date) = 2009
GROUP BY   CUBE (Brand, Product)
ORDER BY   Brand, Product;

```

■ **Definition 4.2** *CUBE in SQL:1999*

A **CUBE** group is an extension of the **GROUP BY** clause that produces a result that contains subtotal records in addition to the regular grouped records, whose aggregate values are derived by applying the same functions. A **CUBE** $(A_1, A_2, \dots, A_{n-1}, A_n)$ group is equivalent to the union of the 2^n grouping results on the attributes of all possible subsets of the attributes specified in the **CUBE** clause. Unlike **ROLLUP**, the order in which the attributes are specified doesn't matter for **CUBE**. The operator produces its results with just one table access.

Some systems also provide the operator **GROUPING SETS** to group only for certain combinations of attributes. For example, replacing in the previous query **GROUP BY CUBE**(Brand, Product) with **GROUP BY GROUPING SETS**((Brand, Product), (Brand)) data are grouped only for the two combinations listed.

4.3.3 Observations

In general, in the **GROUP BY** clause both attributes and different **ROLLUP** and **CUBE** can be used. For example, the following query


```

SELECT   Date, Brand, Product, SUM(Revenue) AS Revenue
FROM     Sales
GROUP BY Date, ROLLUP(Brand, Product);

```

generates the following groupings: (Date, Brand, Product), (Date, Brand) and (Date), but not (). The result is justified by recalling that Date generates the set of groupings {(Date)}, ROLLUP generates the set of groupings {(Brand, Product), (Brand), ()} and their combination generates the cartesian product of two sets. The following query

```

SELECT   Date, Brand, Product, SUM(Revenue) AS Revenue
FROM     Sales
GROUP BY CUBE(Date), ROLLUP(Brand, Product);

```

generates the following groupings instead: (Date, Brand, Product), (Date, Brand), (Date), (Brand, Product), (Brand) and ().

Usually, in relational systems when using operators **ROLLUP** and **CUBE**, the result shows the value NULL to indicate a running total, creating ambiguity because the value might be present in the data and not a result of the operators **ROLLUP** and **CUBE**. To correctly interpret the meaning of a record the function **GROUPING** is used with an attribute parameter in the **GROUP BY**: the function returns 1 if the value NULL has been created by **ROLLUP** or **CUBE**, and returns 0 otherwise. For example, the result of the query

```

SELECT   Brand, Product, SUM(Revenue) AS Revenue,
          GROUPING(Brand), GROUPING(Product)
FROM     Sales
WHERE    YEAR(Date) = 2009
GROUP BY ROLLUP(Brand, Product);

```

produces the result of the query without **GROUPING**, extended with two more columns that have the value 1 when the record has a field NULL, which corresponds to a total, as shown in Figure 4.6.

Brand	Product	Revenue	GROUPING (Brand)	GROUPING (Product)
B1	P1	2 100	0	0
B1	P2	3 720	0	0
B1	P3	15 300	0	0
B1		21 120	0	1
B2	P4	12 600	0	0
B2	P5	22 500	0	0
B2	P6	48 300	0	0
B2		83 400	0	1
		104 520	1	1

Figure 4.6: Report with **ROLLUP** and **GROUPING**

To get the result without additional columns, but with the value Total when necessary, we write:

```

SELECT   CASE WHEN GROUPING(Brand) = 1 THEN 'Total' ELSE Brand
          END AS Brand,
          CASE WHEN GROUPING(Product) = 1 THEN 'Total' ELSE Product
          END AS Product,
          SUM(Revenue) AS Revenue
FROM     Sales
WHERE    YEAR(Date) = 2009
GROUP BY ROLLUP(Brand, Product);

```

where, with the first **CASE**, if the value of Brand is a NULL, then the string Total will appear (any string can be chosen). Otherwise, its actual value will appear, as shown in Figure 4.7.

Brand	Product	Revenue
B1	P1	2 100
B1	P2	3 720
B1	P3	15 300
B1	Total	21 120
B2	P4	12 600
B2	P5	22 500
B2	P6	48 300
B2	Total	83 400
Total	Total	104 520

Figure 4.7: Displaying the ALL values with Total

The function **GROUPING**, like any other aggregate function, can be used in **HAVING** to select only some of the records produced by **ROLLUP** or **CUBE**. For example, the following query finds only the record with totals:

```

SELECT  CASE WHEN GROUPING(Brand) = 1 THEN 'Total' ELSE Brand
        END AS Brand,
        CASE WHEN GROUPING(Product) = 1 THEN 'Total' ELSE Product
        END AS Product,
        SUM(Revenue) AS Revenue
FROM    Sales
WHERE   YEAR(Date) = 2009
GROUP BY ROLLUP (Brand, Product)
HAVING  GROUPING(Brand) = 1 OR GROUPING(Product) = 1;

```

4.4 Moderately Difficult Reports with SQL

Let us show examples of queries to present results in a spreadsheet-type cross-tab format, rather than the form of lists of values, or to produce reports with metrics to be calculated by comparison with others.

Example 4.1

Let us produce a report with the total revenue in 2009 by product and by geographical areas. The following simple SQL query

```
SELECT Product, Area, SUM(Revenue) AS TotalRevenue
FROM Sales
WHERE YEAR(Date) = 2009
GROUP BY Product, Area
ORDER BY Product, Area;
```

produces a vertically ordered result set that makes it difficult to compare the revenues by product and by geographical area.

Total Revenue by Product and by Geographical area Year 2009		
Product	Area	Revenue
P1	Center	600
P1	Islands	300
P1	North	900
P1	South	300
P2	Center	1 200
P2	Islands	360
P2	North	1 800
P2	South	360
P3	Center	4 680
P3	Islands	1 980
P3	North	7 020
P3	South	1 620
P4	Center	3 600
P4	Islands	1 800
P4	North	5 400
P4	South	1 800
P5	Center	6 300
P5	Islands	3 150
P5	North	9 450
P5	South	3 600
P6	Center	15 000
P6	Islands	5 100
P6	North	22 500
P6	South	5 700

This kind of data is much easier to compare when it is formatted like a spreadsheet-type cross-tab or *pivot table*:

Comparison between Revenues by Product and by Area Year 2009				
Product	North	Center	South	Islands
P1	900	600	300	300
P2	1 800	1 200	360	360
P3	7 020	4 680	1 620	1 980
P4	5 400	3 600	1 800	1 800
P5	9 450	6 300	3 600	3 150
P6	22 500	15 000	5 700	5 100

The result is obtained by grouping the data by Product and by using in the **SELECT** the aggregate function **SUM** with argument a **CASE** expression:

```

SELECT Product,
SUM(CASE
  WHEN Area = 'North' THEN Revenue ELSE 0 END) AS North,
SUM(CASE
  WHEN Area = 'Center' THEN Revenue ELSE 0 END) AS Center,
SUM(CASE
  WHEN Area = 'South' THEN Revenue ELSE 0 END) AS South,
SUM(CASE
  WHEN Area = 'Islands' THEN Revenue ELSE 0 END) AS Islands
FROM Sales
WHERE YEAR(Date) = 2009
GROUP BY Product
ORDER BY Product;

```

DBMSs such as Oracle 11g and SQL Server 2005 provide an extension to SQL to create the cross-tab with a **PIVOT** clause.

Figure 4.8 shows a graphical representation of product revenues by geographic area with multiple groups of stacked bars, while Figure 4.9 shows another graphical representation often used to show the percentage revenue mix of product by geographic area, but the SQL query to produce the revenue percentage is more complex, and we will see later how to write it.

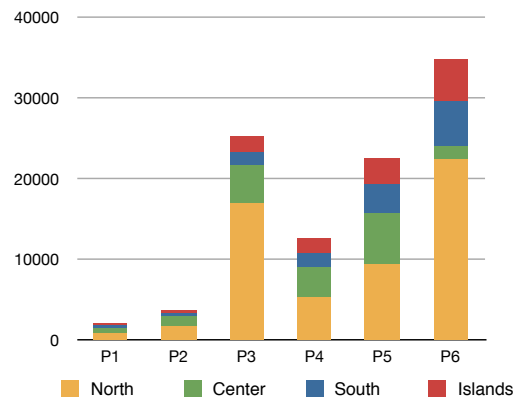


Figure 4.8: An example of a stacked bar report

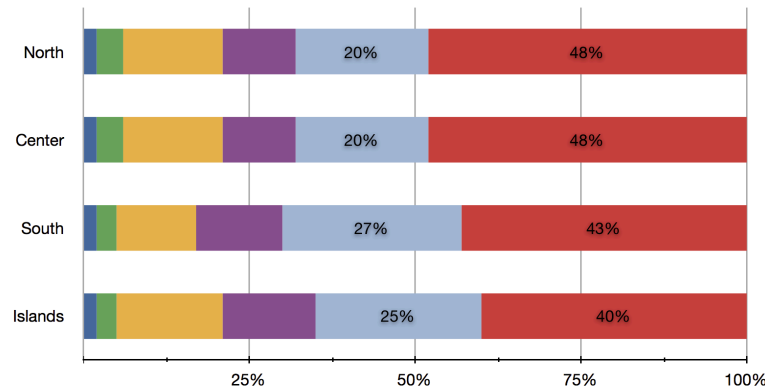


Figure 4.9: Another example of a stacked bar report

Example 4.2

Another very common type of analysis requires reports that compare data columns that refer to different periods (*variance report*). For example, “Revenues for 2009 by brand and by product, with the percentage change from the previous year ($Delta = 100 \times ((Revenue_{2009} - Revenue_{2008}) / Revenue_{2009})$)”.

Comparison between Revenue by Brand and by Product 2009 – 2008				
Brand	Product	Revenue (€) 2009	Revenue (€) 2008	Delta (%)
B1	P1	2 100	13 560	–546
	P2	3 720	23 640	–535
	P3	15 300	20 340	–33
B2	P4	12 600	1 440	89
	P5	22 500	2 100	91
	P6	48 300		100

The annual revenues for 2009 and for 2008, by brand and by product, are obtained with the following SQL queries:

```
Revenue09 = SELECT Brand, Product, SUM(Revenue) AS Revenue2009
             FROM Sales
             WHERE YEAR(Date) = 2009
             GROUP BY Brand, Product ;
```

```
Revenue08 = SELECT Brand, Product, SUM(Revenue) AS Revenue2008
             FROM Sales
             WHERE YEAR(Date) = 2008
             GROUP BY Brand, Product ;
```

If the same products were sold in both 2009 and 2008, the final result would be obtained with a *natural join* of Revenue09 and Revenue08.

Instead, to take into account that not necessarily the same products were sold in both 2009 and 2008, the analysis in SQL requires a *full join* of Revenue09 and Revenue08.

```

SELECT Revenue09.Brand AS Brand, Revenue09.Product AS Product
, Revenue2009
, Revenue2008
, CASE
  WHEN Revenue2009 IS NULL THEN -100
  WHEN Revenue2008 IS NULL THEN 100
  ELSE ROUND(100*(Revenue2009 - Revenue2008) / Revenue2009)
END AS Delta
FROM ( SELECT Brand, Product, SUM(Revenue) AS Revenue2009
FROM Sales
WHERE YEAR(Date) = 2009
GROUP BY Brand, Product
) AS Revenue09

FULL JOIN

( SELECT Brand, Product, SUM(Revenue) AS Revenue2008
FROM Sales
WHERE YEAR(Date) = 2008
GROUP BY Brand, Product
) AS Revenue08

USING (Brand, Product)
ORDER BY Brand, Product;

```

In Figure 4.10 there is shown a graphical representation of the result with a histogram, useful for comparing metrics.

Another very useful graph is the comparison of revenues for the months of different years (Figure 4.11), calculated with an SQL query similar to the previous one.

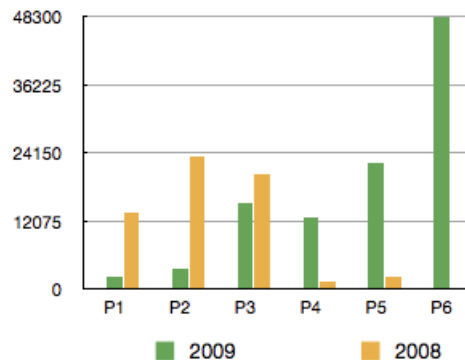


Figure 4.10: The histogram of the revenues variation by product

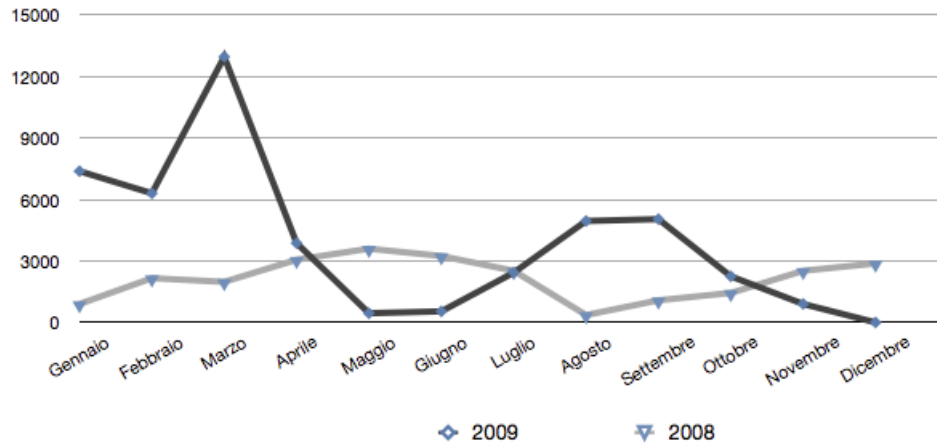


Figure 4.11: The trend in revenues per month of different years

4.4.1 The WITH Clause in SQL

The above example shows how SQL queries for complex analysis generally require the use of subqueries in the **FROM** clause, usually handled by relational systems as a temporary view definition. To make it easier to understand these types of queries, it is useful to write them using the **WITH** clause to break down complex queries with subqueries into simpler parts.

For example, the previous query is written as follows.

```

WITH      Revenue09 AS
          ( SELECT Brand, Product, SUM(Revenue) AS Revenue2009
            FROM   Sales
            WHERE  YEAR(Date) = 2009
            GROUP BY Brand, Product
          )
, Revenue08 AS
          ( SELECT Brand, Product, SUM(Revenue) AS Revenue2008
            FROM   Sales
            WHERE  YEAR(Date) = 2008
            GROUP BY Brand, Product
          )

SELECT    Revenue09.Brand AS Brand, Revenue09.Product AS Product
, Revenue2009
, Revenue2008
, CASE
  WHEN Revenue2009 IS NULL THEN -100
  WHEN Revenue2008 IS NULL THEN 100
  ELSE ROUND(100*(Revenue2009 - Revenue2008) / Revenue2009)
END AS Delta
FROM      Revenue09 FULL JOIN Revenue08 USING (Brand, Product)
ORDER BY Brand, Product;
```

A temporary view defined with **WITH** can use one of the previous defined views or be recursive, if defined with **WITH RECURSIVE**. Recursive queries are typically used to deal with hierarchical or tree-structured fact tables. For example, suppose we have the relation

Flights		
Code	From	To
10	MI	PI
11	MI	TO
12	MI	RM
13	PI	FI
14	TO	RM
15	RM	VE
16	NA	BA
17	TO	PA

The result of the query

WITH RECURSIVE

```

CitiesReachableFrom AS
( SELECT From, To
  FROM Flights
   UNION
  SELECT Flights.From AS From, CitiesReachableFrom.To AS To
  FROM Flights, CitiesReachableFrom
  WHERE Flights.To = CitiesReachableFrom.From
)

SELECT *
FROM CitiesReachableFrom
WHERE From = 'MI';

```

is a relation with pairs of cities reachable from Milano by taking one or more flights.

From	To
MI	PI
MI	TO
MI	RM
MI	FI
MI	VE
MI	PA

The definition of the SQL temporary recursive view `CitiesReachableFrom` has the structure *BasisSelect UNION RecursiveSelect*, with a *linear recursion*, that is the **FROM** of the *RecursiveSelect* contains one occurrence only of the temporary recursive view. The *FinalSelect* of the **WITH** clause is executed with the `CitiesReachableFrom` relation value that might be computed as follows:

- 1) `CitiesReachableFrom := BasisSelect;`
- 2) **WHILE** (changes to `CitiesReachableFrom`) **DO**
- 3) `CitiesReachableFrom := CitiesReachableFrom UNION RecursiveSelect;`

`CitiesReachableFrom` is initially empty and with rule (1) its records become those of the *BasisSelect*, the cities reachable with direct flights, because the *RecursiveSelect* produces an empty relation.

On the next iteration, with rule (3) possible, other cities reachable with more flights are added to the relation *RecursiveSelect*. If the new value is equal to the previous one, the final result has been obtained and the loop ends. Otherwise, the value obtained is used to find other cities reachable.

The *FinalSelect* is executed with the final value of `CitiesReachableFrom` and the result is the subset of records with `From = 'MI'`.

4.5 Very Difficult Reports Without Analytic SQL

The SQL has been extended to allow the use of several analytic functions, also known as *Windows Functions*, to use them in new ways in order to facilitate the development of complex data analysis.

A **SELECT**, for simplicity with a single analytic function, without **DISTINCT**, with a fact table and one dimensional table only, has the following structure:

SELECT Select Attributes (S_A), Select Aggregation Functions (S_{AF}),
 Analytic Function (A_F) **OVER**(
 [**PARTITION BY** <attribute list>]
 [**ORDER BY** <sort attribute list>
 [<window clause>]])
FROM Fact table (F) and a dimension table (D1)
WHERE Where condition (W_C)
GROUP BY Grouping Attributes (G_A)
HAVING Having condition (H_C) with aggregation functions (H_{AF})
ORDER BY Sorting attributes (O_A);

The command semantics in terms of the extended relational algebra is shown in Figure 4.12.

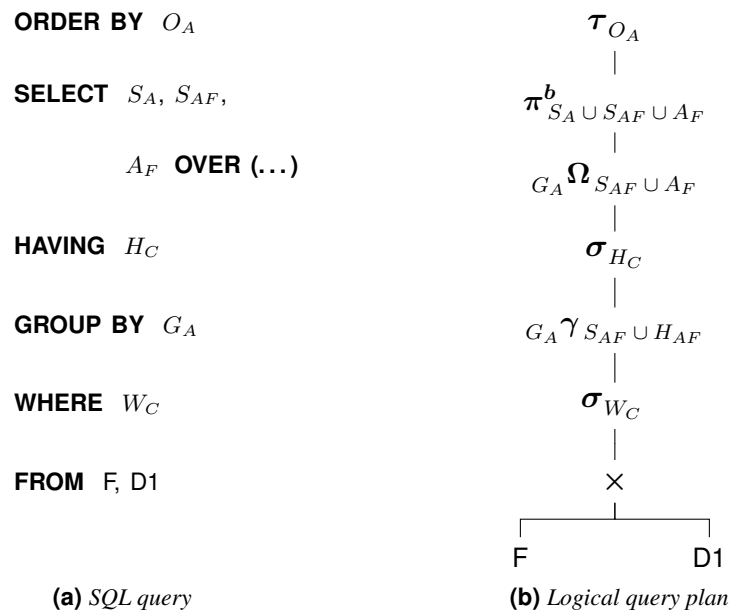


Figure 4.12: Analytic SQL query semantics

The query is processed in the following steps:

1. Perform the operations specified in **FROM**, **WHERE**, **GROUP BY** and **HAVING** clauses to compute the set of records resulting from the subtree rooted at σ_{H_C} .
2. Apply the specified analytic functions to the result of the subtree rooted at σ_{H_C} , to produce a new set of record that differs from the previous ones only for new attributes calculated using the analytic functions.
3. Apply any projection, and then the **ORDER BY** to produce the query result.

The analytic functions can be used only in **SELECT**, with the **OVER** clause, and usually they are applied to the entire set of records produced in the first phase, but may also

be applied separately to disjoint subsets obtained by partitioning the records by the value of an expression defined on the attributes of the record (option **PARTITION BY**). The aggregate functions can also be applied to non-disjoint subsets of records in a partition defined by the notion of *moving window*: for each record r of a partition, aggregate functions apply to the data identified by a “window” placed on r . They are useful for analysis of data, such as: “What is the moving average of weekly sales?”.

The result of the traditional aggregate functions SUM, COUNT, AVG, MIN, MAX does not depend on the order of records in the collection on which they operate. Instead, the result of the new analytic functions, which we will see that later on, may depend on the order of the data specified with the **ORDER BY** clause in the **OVER** of the **SELECT**.

The partitioning operation is like that for the calculation of a **GROUP BY**, but **PARTITION BY** does not produce a record for each group as with the **GROUP BY**, but rather produces as many records as there are elements of the group, which will then be extended with new attributes calculated using the analytic functions. When the **PARTITION BY** clause is not present, the set of records behaves as a single group.

Figure 4.13 shows the main analytic functions available in some DBMS.

Function	Oracle	DB2	SQL Server	PostgreSQL	MySQL
COVAR_POP	x	x	x	x	
CUBE	x	x	x		
CUME_DIST	x		x	x	
DENSE_RANK	x	x	x	x	
LAG, LEAD	x	x	x	x	
NTILE	x	x	x	x	x
PERCENT_RANK	x		x		
RANK	x	x	x	x	
RATIO_TO_REPORT	x	x			
REGR_Functions	x	x		x	
ROLLUP	x	x	x		x
ROW_NUMBER	x	x	x	x	
STDDEV_POP	x	x	x	x	x
VAR_POP	x	x	x	x	x
Window.Clause	x	x	x		

Figure 4.13: Analytic SQL in some DBMS

4.5.1 Premise

Note the difference between the result of a query with an aggregate function, and the traditional **GROUP BY**, and an analytic function.

Let us consider the relation

R	
P	...
P1	...
P1	...
P2	...
P2	...
P2	...
P2	...
P2	...

The query

```

SELECT   P, COUNT(*) AS No
FROM     R
GROUP BY P;

```

returns the relation

P	No
P1	2
P2	5

while the query

```

SELECT   P,
           COUNT(*) OVER (PARTITION BY P) AS No
FROM     R
ORDER BY P;

```

returns the relation

P	No
P1	2
P1	2
P2	5
P2	5
P2	5
P2	5
P2	5

While the **GROUP BY** groups a set of records and a record for each group with two attributes is obtained – with the value of the grouping attribute and the value of the aggregate function **COUNT** – with the analytic function, for each record of the set the value of the aggregate function **COUNT** is computed when applied to subsets obtained with the **PARTITION**.

If the query was without **PARTITION**, the analytic function would be applied to the whole record set:

```

SELECT   P, COUNT(*) OVER() AS No
FROM     R
ORDER BY P;

```

to get the table

P	No
P1	7
P1	7
P2	7
P2	7
P2	7
P2	7
P2	7

Compare this result with that of

```

SELECT   COUNT(*) AS No
FROM     R

```

No
7

The next example shows the usefulness of the **OVER** clause to solve a nontrivial problem left unresolved in a previous example.

Example 4.3

Let us reconsider the example of the report with the total sales revenue in 2009, by product and by geographical area:

Total revenue by Product and by Geographical area Year 2009				
Product	North	Center	South	Islands
P1	900	600	300	300
P2	1 800	1 200	360	360
P3	7 020	4 680	1 620	1 980
P4	5 400	3 600	1 800	1 800
P5	9 450	6 300	3 600	3 150
P6	22 500	15 000	5 700	5 100

to change it by replacing the total revenue by product and by geographical area, with the revenue percentage of total revenue for the area.

Percentage of total revenue by area by Product and by Geographical area Year 2009				
Product	PctNorth	PctCenter	PctSouth	PctIslands
P1	2	2	2	2
P2	4	4	3	3
P3	15	15	12	16
P4	11	11	13	14
P5	20	20	27	25
P6	48	48	43	40

The result is obtained with the query

```

SELECT Product
, ROUND(100*SUM(CASE
                WHEN Area = 'North'
                THEN Revenue ELSE 0 END)
        / SUM(SUM(CASE
                WHEN Area = 'North'
                THEN Revenue ELSE 0 END))
        OVER ()
        ) AS PctNorth
, ROUND(100*SUM(CASE
                WHEN Area = 'Center'
                THEN Revenue ELSE 0 END)
        / SUM(SUM(CASE
                WHEN Area = 'Center'
                THEN Revenue ELSE 0 END))
        OVER ()
        ) AS PctCenter
, ROUND(100*SUM(CASE
                WHEN Area = 'South'
                THEN Revenue ELSE 0 END)
        / SUM(SUM(CASE
                WHEN Area = 'South'
                THEN Revenue ELSE 0 END))
        OVER ()
        ) AS PctSouth
, ROUND(100*SUM(CASE
                WHEN Area = 'Islands'
                THEN Revenue ELSE 0 END)
        / SUM(SUM(CASE
                WHEN Area = 'Islands'
                THEN Revenue ELSE 0 END))
        OVER ()
        ) AS PctIslands
FROM Sales
WHERE YEAR(Date) = 2009
GROUP BY Product
ORDER BY Product;

```

GROUP BY with CASE

We have seen how the construct **CASE** is useful in the **SELECT** for certain analyses, and other examples will be seen later. Let us see how it can also be useful in the **GROUP BY**, for grouping the data of a report not on the values of certain attributes, but on values calculated from those of other attributes.

Let us consider the relation

S		
P	Prc	...
P1	10	...
P2	20	...
P3	30	...
P4	40	...
P5	50	...
P6	60	...
P7	70	...

Suppose we want a report to display the products classified in 3 categories: *Cheap* ($Prc \leq 20$), *Medium* ($20 < Prc \leq 50$) e *Expensive* ($50 < Prc \leq 100$).

P	Prc	Category
P1	10	Cheap
P2	20	Cheap
P3	30	Medium
P4	40	Medium
P5	50	Medium
P6	60	Expensive
P7	70	Expensive

The result is obtained with the query

```

SELECT P, Prc
      , CASE
        WHEN Prc <= 20 THEN 'Cheap'
        WHEN Prc > 20 AND Pz <= 50 THEN 'Medium'
        ELSE 'Expensive'
      END AS Category
FROM S;

```

Now suppose we want a report showing the average price of products by category.

Category	AvgPrice
Cheap	15
Medium	40
Expensive	65

```

WITH CategoryAndPrice AS
( SELECT CASE
      WHEN Prc <= 20 THEN 'Cheap'
      WHEN Prc > 20 AND Prc <= 50 THEN 'Medium'
      ELSE 'Expensive'
    END AS Category
  , Price
  FROM S
)

SELECT Category, AVG(Price) AS AvgPrice
FROM CategoryAndPrice
GROUP BY Category
ORDER BY AvgPrice;

```

Without the **WITH** the query becomes

```

SELECT CASE
      WHEN Prc <= 20 THEN 'Cheap'
      WHEN Prc > 20 AND Prc <= 50 THEN 'Medium'
      ELSE 'Expensive'
    END AS Category
  , AVG(Prc) AS AvgPrice
FROM S
GROUP BY CASE
      WHEN Prc <= 20 THEN 'Cheap'
      WHEN Prc > 20 AND Prc <= 50 THEN 'Medium'
      ELSE 'Expensive'
    END
ORDER BY AvgPrice;

```

The query must be written with two *syntactically equal* **CASE** expressions, one in the **GROUP BY** and another in the **SELECT**.

Some relational systems, such as PostgreSQL, allow the use in the **GROUP BY** of the **CASE** expression label in the **SELECT** to avoid rewriting the expression in the **GROUP BY** and, therefore, the query can be written as

```

SELECT   CASE
          WHEN Prc <= 20 THEN 'Cheap'
          WHEN Prc > 20 AND Prc <= 50 THEN 'Medium'
          ELSE 'Expensive'
        END AS Category
        ,AVG(Prc) AS AvgPrice
FROM     S
GROUP BY Category
ORDER BY AvgPrice;

```

We shall see later the use of **GROUP BY** with **CASE** in a case of more complex analysis.

4.5.2 Analytic Functions with the Use of Partitions

RANK and DENSE_RANK

These functions are used to sort the records out in a set based on the value of an attribute or of an expression (aggregate function), and to assign to each record its position (*rank*) in the set. The standard record order is ascending, and so the records with rank 1 have the minimum value of the attribute, but the descending order can be specified. The result is sorted by the rank value, unless otherwise specified. A ranking function is specified in the **SELECT** clause with the following syntax:

```

<RankFunction>()
OVER(
  [PARTITION BY <attribute list>]
  ORDER BY <sort attribute list>
) [ AS lde ]

```

RANK and **DENSE_RANK** require an **ORDER BY** clause, because to determine the values rank the data must be ordered. If no partitioning is specified, the entire set of records composes a single partition.

The functions **RANK** and **DENSE_RANK** produce different results when the values to be ranked are not different. The rank of a value a_i is defined as 1 plus the number of values that *strictly* precede a_i . If $k > 1$ values are equal, they are assigned the same value rank p , and the next value has the rank $p + k$. Therefore there will be a gap in the sequential rank numbering. Instead, with **DENSE_RANK** there will be no gaps in the sequential rank numbering, with ties being assigned the same rank. The rank of a value a_i is defined as 1 plus the number of *distinct* values that *strictly* precede a_i . For example, the values in the ascending order (10, 20, 20, 30, 30, 40) have the ranks (1, 2, 2, 4, 4, 6) and the dense ranks (1, 2, 2, 3, 3, 4).

Example 4.4

“Show for the year 2009, and the regions of Tuscany and Lazio, the total revenue by region and product, the rank of the products for total revenue in each region and for total revenue.”

Revenues and Ranks in the 2009 by Region and by Product				
Region	Product	Total Revenue	Product Rank by Region	Product Rank Global
Lazio	P3	2880	3	4
	P2	960	5	8
	P4	2700	4	5
	P1	480	6	10
	P5	4800	2	2
	P6	11400	1	1
Toscana	P1	120	6	12
	P6	3600	1	3
	P3	1800	2	6
	P5	1500	3	7
	P4	900	4	9
	P2	240	5	11

```

SELECT Region, Product
, SUM(Revenue) AS TotalRevenue
, RANK() OVER (PARTITION BY Region ORDER BY SUM(Revenue) DESC)
AS ProductRankByRegion
, RANK() OVER (ORDER BY SUM(Revenue) DESC)
AS ProductRankGlobal

FROM Sales
WHERE YEAR(Date) = 2009
AND Region IN ('Toscana', 'Lazio')

GROUP BY Region, Product
ORDER BY Region;

```

Example 4.5

“Show for the year 2009 total revenue, the rank and dense rank of total revenue for clients.”

Revenues, Rank and Dense Rank in the 2009 by Customer			
Customer	Total Revenue	Customer Rank	Customer Dense Rank
C15	21120	1	1
C03	13650	2	2
C08	11400	3	3
C09	9240	4	4
C04	8640	5	5
C11	7560	6	6
C02	7560	6	6
C12	7200	8	7
C16	4680	9	8
C14	4200	10	9
C06	3150	11	10
C10	2520	12	11
C01	1920	13	12
C13	1680	14	13


```

SELECT Customer
, SUM(Revenue) AS TotalRevenue
, RANK() OVER (ORDER BY SUM(Revenue) DESC)
AS CustomerRank
, DENSE_RANK() OVER (ORDER BY SUM(Revenue) DESC)
AS CustomerDenseRank
FROM Sales
WHERE YEAR(Date) = 2009
GROUP BY Customer
ORDER BY TotalRevenue DESC;

```

Example 4.6

The ranking functions can be used to find only the first N records with the value of higher or lower ranks (**TOP_N**, **BOTTOM_N**): (a) there is the **SELECT** with the function of rank in the **FROM** clause and then (b) an appropriate selection of the result of the external **SELECT** on the rank values. Let us see an example.

“Show for the year 2009, and the Islands geographic area, the region and the customers with the two highest revenues (TOP 2 customers).”

Top 2 customers in the islands Year 2009			
Region	Customer	Total Revenue	Top 2
Sardegna	C15	2 640	1
Sicilia	C11	1 080	1
Sicilia	C04	1 080	1
Sardegna	C03	1 560	2

```

WITH SalesWithRank AS
( SELECT Region, Customer
, SUM(Revenue) AS TotalRevenue
, RANK() OVER (PARTITION BY Region
ORDER BY SUM(Revenue) DESC )
AS Rank
FROM Sales
WHERE YEAR(Date) = 2009 AND Area = 'Islands'
GROUP BY Region, Customer
)

SELECT Region, Customer, TotalRevenue, Rank AS Top2
FROM SalesWithRank
WHERE Rank <= 2
ORDER BY Top2;

```

NTILE(n)

A set of sorted records is partitioned into n groups with the same number of records (plus or minus 1) and the group number to which it belongs is assigned to each record.

Example 4.7

Customers are divided into 4 groups on the basis of total revenue, and their rank is calculated in each quartile.

Revenue and rank in quartiles by 4 customer groups			
Customer	Total Revenue	Quartile	Rank by Quartile
C03	22 890	1	1
C15	21 120	1	2
C08	16 440	1	3
C04	14 400	1	4
C02	12 600	2	1
C12	11 760	2	2
C09	9 240	2	3
C05	9 240	2	3
C16	9 240	3	1
C14	8 820	3	2
C11	7 560	3	3
C07	7 200	3	4
C01	4 800	4	1
C06	4 410	4	2
C13	3 360	4	3
C10	2 520	4	4

```

WITH CustomersQuartiles AS
( SELECT Customer, SUM(Revenue) AS TotalRevenue
  , NTILE(4) OVER (ORDER BY SUM(Revenue) DESC)
    AS Quartile
  FROM Sales
  GROUP BY Customer
)
SELECT Customer, TotalRevenue, Quartile
  , RANK() OVER (PARTITION BY Quartile ORDER BY TotalRevenue DESC)
    AS RankByQuartile
FROM CustomersQuartiles;

```

ROW_NUMBER() and CUME_DIST()

On a set of sorted records

- **ROW_NUMBER()** assigns a sequence number to each record and
- **CUME_DIST()** assigns a value between 0 and 1 to each record of a sorted set according to the number of records that precede it. For a record r in a set with n elements sorted in increasing order, if k is the number of records that precede it, the **CUME_DIST()** of r is $0 < (k + 1)/n \leq 1$. To equal values **CUME_DIST()** assigns equal values.

Example 4.8

Consider the customers ordered by the sum of all their purchases (sales revenue) in descending order.

1. We want to see if the Pareto rule holds: 80 percent of revenue comes from 20 percent of customers. These customers are important because the business

depends on their loyalty. In particular, of the Top20% of customers, that is the 20% of customers with the highest sales, we want to know their name and the sales revenue, total revenue of all sales, their position n in the Top20% and the percentage of the sum of their revenue compared to total revenue of all sales.

Customers Top20% by Revenue					
Customer	Revenue by Customer	Total Revenue	n	Percent of Total Revenue	Percent of Running Totals
C03	22 890	165 600	1	14	6
C15	21 120	165 600	2	13	12
C08	16 440	165 600	3	10	19

```

WITH RowNumberCustomer AS
( SELECT Customer, SUM(Revenue) AS RevenueByCustomer
, SUM(SUM(Revenue)) OVER() AS TotalRevenue
, ROW_NUMBER()
OVER(ORDER BY SUM(Revenue) DESC) AS n
FROM Sales
GROUP BY Customer ORDER BY RevenueByCustomer
)
, RowNumberCustomerExtended AS
( SELECT Customer, RevenueByCustomer, TotalRevenue, n
, ROUND(100*RevenueByCustomer / TotalRevenue)
AS PercentOfTotalRevenue
, ROUND(100*CUME_DIST() OVER(ORDER BY n))
AS PercentOfRunningTotals
FROM RowNumberCustomer
)
SELECT *
FROM RowNumberCustomerExtended
WHERE PercentOfRunningTotals <= 20;

```

In this case the Pareto rule does not hold: only 37% of sales are related to 19% of the customers.

2. We want to partition the customers into four groups:
 - Top5%, with 5% of customers with the highest amount of revenues.
 - Next15%, with 15% of other customers with the highest amount of revenues.
 - Middle30%, with 30% of other customers with the highest amount of revenues.
 - Bottom50%, with 50 % of the other customers with the lowest amount of revenues.

For each customer group we want to know their number, and the percentage of the sum of their revenues compared to total revenue of all sales.

Customers by Revenue Top5%, Next15%, Middle30% and Bottom50%		
Group	Number of Customers	Percent of TotalRevenue
Next15%	2	27
Middle30%	2	19
Bottom50%	12	55

```

WITH RowNumberCustomers AS
( SELECT Customer, SUM(Revenue) AS RevenueByCustomer
  ,SUM(SUM(Revenue)) OVER() AS TotalRevenue
  ,ROW_NUMBER()
    OVER (ORDER BY SUM(Revenue) DESC) AS n
  FROM Sales
  GROUP BY Customer
)
,RowNumberCustomersExtended AS
( SELECT Customer, RevenueByCustomer, TotalRevenue
  ,ROUND(100*CUME_DIST() OVER (ORDER BY n))
    AS PercentOfCustomers
  FROM RowNumberCustomers
)
,RowNumberCustomersExtendedWithGroup AS
( SELECT Customer, RevenueByCustomer, TotalRevenue
  ,(CASE
    WHEN PercentOfCustomers <= 6
      THEN 'Top6%'
    WHEN PercentOfCustomers > 6 AND
      PercentOfCustomers <= 14
      THEN 'Next14%'
    WHEN PercentOfCustomers > 14 AND
      PercentOfCustomers <= 30
      THEN 'Middle30%'
    ELSE 'Bottom50%' END
  ) AS Group
  FROM RowNumberCustomersExtended
)

SELECT Group
  ,COUNT(Customer) AS NumberOfCustomers
  ,ROUND(100*SUM(RevenueByCustomer) / TotalRevenue)
    AS PercentOfTotalRevenue
FROM RowNumberCustomersExtendedWithGroup
GROUP BY Group, TotalRevenue
ORDER BY Group DESC;

```

4.5.3 Analytic Functions with the Use of Windows

For each record in a set, called the *current record*, a *window* can be defined on the data to determine the record set ‘nearby’ to be taken into account for the calculation of the new fields to be added to the record. The window size can be determined in a physical way (option **ROWS**), based on the number of records, or in a logical way (option **RANGE**), using a condition usually based on an attribute of type DATE.

The current record of a set (or a partition) is both the one of reference for the calculation of an aggregate function, and that for which the window size is defined by specifying the start and the end record, which can then change when the next current record is selected.

A window can include all records of the set on which it is defined, or include only the current record. For example, to calculate a cumulative sum function, the first record is fixed and the end of the set moves from first to last record, while to calculate a moving average both the first and last record move.

For each current record, the records of the specified window are considered, and with them the value of an aggregate function is computed. The general format of the window clause is:

```

<AggregateFunction>(<expr>)
OVER(
  [PARTITION BY <attribute list>]
  [ORDER BY <sort attribute list>]
  [<ROWS or RANGE> <window size specification>]]
) [ AS lde ]

```

Example 4.9

Consider the following table with the transactions data of bank accounts:

BankAccount(AccountNumber, TransactionDate, TransactionType)

We want to find the balance of the accounts sorted by date of transactions.

```

SELECT  AccountNumber, TransactionDate, TransactionType
        ,SUM(TransactionType) OVER
          ( PARTITION BY AccountNumber ORDER BY TransactionDate
            ROWS UNBOUNDED PRECEDING) AS Balance
FROM    BankAccount
ORDER BY AccountNumber, TransactionDate;

```

where **ROWS UNBOUNDED PRECEDING** specifies that the window begins with the first record of the partition and ends with the current record.

Account Number	Transaction Date	Transaction Type	Balance
1234	2009-11-01	113.00	113.00
1234	2009-11-05	-52.00	61.00
1234	2009-11-13	36.00	97.00
4321	2009-11-01	10.00	10.00
4321	2009-11-21	32.00	42.00
4321	2009-11-29	-5.00	37.00

Example 4.10

The cumulative total (running totals) is another piece of useful information to highlight the features of the report data and to facilitate analysis. In the figure an example is shown to know about the “cumulative monthly revenue by quarter (*Quarter-to-Date*) and year (*Year-to-Date*) for the product P1 in 2009”.²

Product P1 Revenue by Quarter and Month Year 2009				
Quarter	Month	Revenue (€)	Revenue QtoD (€)	Revenue YtoD (€)
Q1	January	16 500	16 500	16 500
Q1	February	14 220	30 720	30 720
Q1	March	27 480	58 200	58 200
Q2	April	7 920	7 920	66 120
Q2	May	1 200	9 120	67 320
Q2	June	1 260	10 380	68 580
Q3	July	5 400	5 400	73 980
Q3	August	11 730	17 130	85 710
Q3	September	10 860	27 990	96 570
Q4	October	5 850	5 850	102 420
Q4	November	2 100	7 950	104 520
Q4	December			

```

SELECT Quarter_Name(QUARTER(Date)) AS Quarter
, Month_Name(MONTH(Date)) AS Month
, SUM(Revenue) AS Revenue
, SUM(SUM(Revenue)) OVER
(PARTITION BY QUARTER(Date)
ORDER BY MONTH(Date)
ROWS UNBOUNDED PRECEDING) AS RevenueQToD
, SUM(SUM(Revenue)) OVER
(ORDER BY MONTH(Date)
ROWS UNBOUNDED PRECEDING) AS RevenueYToD
FROM Sales
WHERE YEAR(Date) = 2009
GROUP BY QUARTER(Date), MONTH(Date)
ORDER BY Quarter, Month;

```

Finally, an example is given of a moving window with a size defined in a physical way. Assuming that the total revenues from product sales vary greatly during the year, their values do not make a clear global trend, but it might be useful to know how to predict future revenues. For this reason, an interesting report is the one that shows a moving average of total revenues over three months – the current one, the one that precedes it and the other that follows it.

```

SELECT MONTH(Date) AS Month
, ROUND(AVG(SUM(Revenue))
OVER (ORDER BY MONTH(Date)
ROWS BETWEEN 1 PRECEDING AND 1 FOLLOWING), 2)
AS MovingAverageRevenue
FROM Sales
GROUP BY MONTH(Date)
ORDER BY Month;

```

In Oracle, when you define a query in analytic SQL you can choose to display a graphical representation of the result. For example, Figure 4.14 shows the trends of the moving average of total revenue, with a moving window of 3 or 5 months.

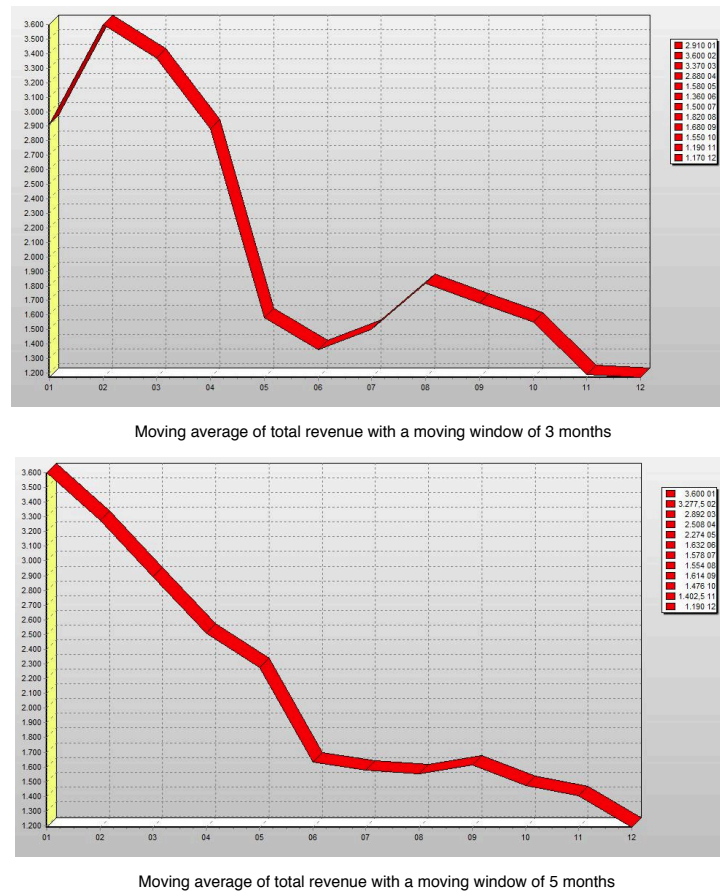


Figure 4.14: Two graphical representations of the moving average of total revenue by product and month

4.6 Summary

- The main commercial solutions for OLAP systems have been presented. OLAP servers are implemented using either a multidimensional storage engine (MOLAP), a relational DBMS engine (ROLAP) as the backend, or a hybrid combination called HOLAP. Changes in the hardware technology will change how the backend of large data warehouses are organized, and as cloud data services take root, more changes are expected.
- The standard SQL language supports relatively simple data analysis, and so it has been extended with new operators and analytic functions to allow complex data analysis.
- Data analysis in SQL is useful for understanding both the functionality of tools that provide graphical interfaces to formulate the queries, and their limitations in expressive power.

CASE STUDIES

A.1 Hospital

A hospital uses the database shown in Figure A.1 to store the following information about inpatients' treatments.

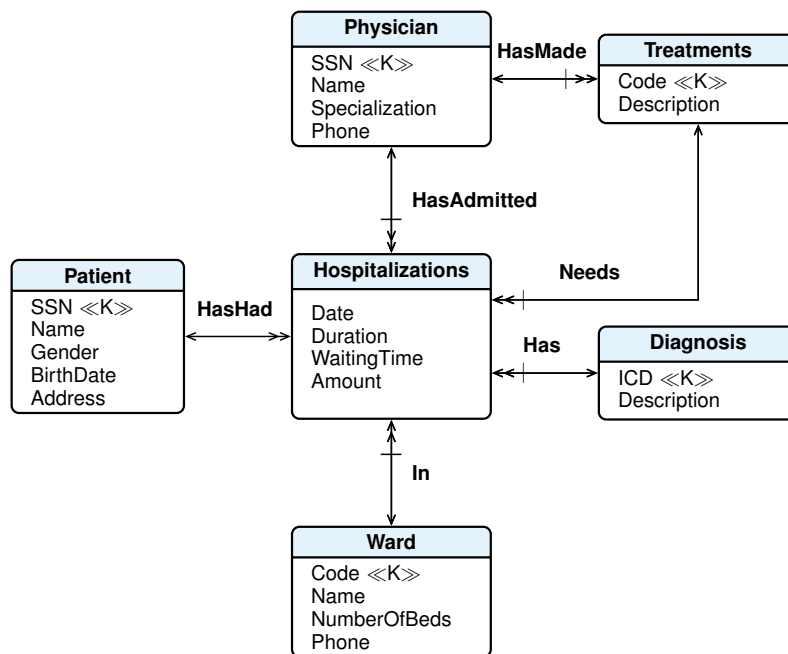


Figure A.1: Conceptual design of a database for inpatients' treatments

For each patient the information of interest is the SSN (*Social Security Number*), which is unique, the name and the address.

A patient may be hospitalized several times, and each time the information of interest is the date, the physician who admitted the patient, the ward assigned, the diagnosis, the duration in days of hospitalization, the number of days of waiting time for the hospitalization, the received treatment, and the billed amount. For simplicity, let us assume that each patient may receive only one treatment from a given physician on a given hospitalization.

For each hospital ward the information of interest is the code, which is unique, the name, the number of beds and the phone.

For each treatment the information of interest is the code, which is unique, the description, and the physician who carried it out.

For each diagnosis the information of interest is the ICD code, *International Classification of Diseases*, which is unique, and its description.

Give a conceptual and logical designs of a data mart assuming that the following examples of business questions have been collected during the user interviews:

1. Total billed amount for hospitalizations by diagnosis code and description, by month (year).
2. Total number of hospitalizations and billed amount by ward, by patient gender (age at date of admission, city, region).
3. Total billed amount, average length of stay and average waiting time by diagnosis code and description, by name (specialization) of the physician who admitted the patient.
4. Total billed amount, and average waiting time of admission by patient age (region), by treatment code (description).

A.2 Airline Companies

We want to analyze airline companies' flights to compare them from the point of view of their ability to fly with occupied seats and therefore to make profits.

For each flight the information of interest is the company name, the departure and the destination cities, the departure time (hour, day, month, year), the number of unoccupied seats in each class (economic, business, first), the revenue of each class.

A flight code (a combination of the *ICAO airline designator* with the flight number) identifies a flight of an airline company from a departure airport to a destination airport (e.g. AP2701 is an Alitalia flight from Malpensa to Fiumicino, available on certain days a week).

A flight is identified by the flight code and the departure time.

For each city the information of interest is the city's name, the country and the continent.

For each company the information of interest is the name and the type (private or national).

Give a conceptual and logical data mart designs assuming that the following examples of business questions have been collected during the user interviews:

1. Number of unoccupied seats in a given year, by flight code, by company name (or type), by class, by departure time (hour, day, month, year)
2. Number of unoccupied seats in a given class and year, by flight code, by company name, by class, by departure (destination) city (country, continent)
3. Number of unoccupied seats and income of the Alitalia company, by year, by month, by destination country.

A.3 Airline Flights

An airline uses the database shown in Figure A.2 to store the following traffic information on passengers from its flights.

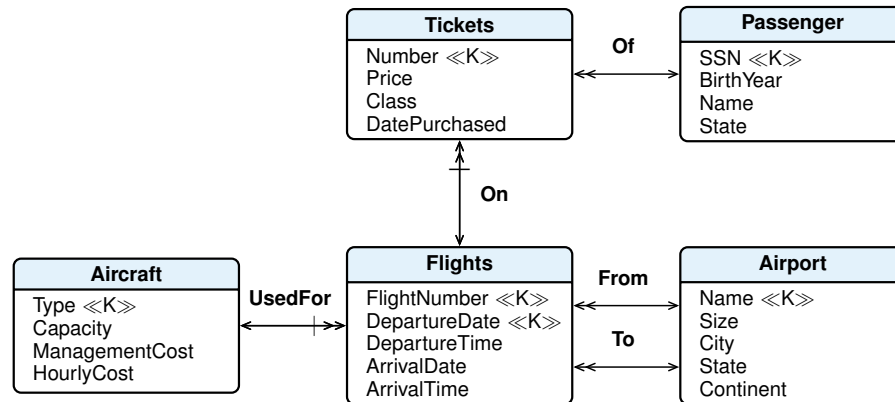


Figure A.2: The conceptual design of a database for airline flights

For each flight the information of interest is the flight code, the date and time of departure, the date and time of arrival, the departure and destination airport, and the aircraft used. The flight code identifies a possible flight of an airline from a departure airport to a destination airport (e.g. AP2701 is an Alitalia flight from Malpensa to Fiumicino, available on certain days a week). A specific flight is identified by a flight code and the departure date.

A flight is used by a set of passengers who bought a ticket. For each ticket the information of interest is the number, which is unique, the price, the class and the date of purchase. For each passenger the information of interest is the SSN, which is unique, the name and country.

A flight is made with an aircraft for which the information of interest is the type, which is unique, the capacity, the monthly management cost (management cost for brevity) and the hourly operating cost (fuel and crew).

For each airport the information of interest is the name, which is unique, the country, the continent, and the size, with values “small”, “medium”, “large”.

Give a conceptual and logical data mart designs and the SQL queries for the following business questions collected during the user interviews:

1. Number of first-class passengers in a given month and year, by country and by age range of passengers.
2. Number of passengers from Europe to the U.S. in a given month and year, and the total revenue, by country and by age range of passengers.
3. Number of flights, by departure city, by destination city.
4. Average number of airline passengers, by month, by aircraft type, by country of destination.
5. Average number of airline passengers, by class, by holiday date.
6. Number of passengers, by year, by size of the destination airport.
7. Number of flights to airports in Germany from the October to December quarter of a given year, and total management cost of the aircraft, by aircraft type.
8. Average profit of all flights, by country of departure, by destination country. The profit of a flight is the total passenger price minus the total flight cost.

9. Total revenue in a given year of flights by month, by destination country. The total revenue by month, total revenue by destination country, and the total revenue are also of interest.

A.4 Inventory

A company has a set of warehouses in different cities containing some of the products for sale. Proper inventory management is very important for a company and requires reconciling two conflicting demands: keeping up their level to meet customer demands, and minimizing their level to reduce the capital investment and space required for storage.

The quantity of a product that is added to the warehouse is called QtyAcquired (QtyA), that which is removed from the warehouse and shipped to customers is called QtyShipped (QtyS) and that present on a given day is called QtyOnHand (QtyOH). For example:

Date	Description	QtyA	QtyS	QtyOH	Days	QtyOH × Days	Average monthly QtyOH
01/01/2008	Initial QtyOnHand			100	14	1 400	
15/01/2008	QtyReceived	120		220	6	1 320	
21/01/2008	QtyShipped		80	140	4	560	
25/01/2008	QtyShipped		60	80	7	560	
Totals January		120	140		31	3 840	123.87
01/02/2008				80	6	480	
07/02/2008	QtyShipped		20	60	8	480	
15/02/2008	QtyReceived	150		210	10	2 100	
25/02/2008	QtyShipped		50	160	4	640	
29/02/2008	QtyShipped		100	60	1	60	
Totals February		150	170		29	3 760	129.66

Among the possible models for the analysis of inventory products, and their handling, the *model for periodic snapshots* is considered, simplified as follows: at the end of each month, for each deposit, the following quantities are considered: (a) the monthly average QtyOnHand of each product, (b) the total quantity of each product acquired in the month, and (c) the total quantity of each product shipped during the month. The monthly average QtyOnHand is calculated as the monthly arithmetic average of the various values of existing stocks for one month, weighted by their durations:

$$\frac{\sum_{i=1}^n q_i \times d_i}{T}$$

where q_i is the value of the quantity on hand for d_i days and T is the number of days of the month.

For the sake of simplicity, we will use the terms QtyOnHand, QtyAcquired and QtyShipped instead of Monthly average QtyOnHand, Monthly QtyAcquired and Monthly QtyShipped.

The company is interested in analyzing the QtyOnHand on a volume basis, and not on a financial basis, by considering the following metrics in a given *time period* T_m measured in months:

- *Inventory Turns (Inventory Turnover Ratio or Turns)* is the top inventory metric used by any business. This metric measures how fast a product moves in and out of the warehouse, and is calculated with the following ratio:

$$\frac{\text{Total QtyShipped}}{\text{Average QtyOnHand}}$$

A high value of the *Inventory Turns* means that products are more frequently sold, and so their sale allow the purchase of new quantities on hand.

- *Days in Inventory* of a product (*Days Sales in Inventory*, *Average Turnover Period* or *Days Inventory Outstanding*) measures the average time (in days) of the product in stock, and is calculated with the following ratio:

$$\frac{\text{Number of days}}{\text{Inventory Turns}}$$

A low value of the *Days in Inventory* means that the recovery of capital invested in stocks is more rapid.¹

The following are some examples of business questions collected during the user interviews. For each report there are others similar with QtyAcquired and QtyShipped, by city or region of warehouses, by product category, by quarter or year.

Report 1. Total of *Quantity on Hand* (QtyOnHand) in January 2010, by product (SKU and Product Name), by region. The subtotal, by all regions, is also of interest.

Quantity on Hand January 2010			
SKU Product	Product Name	Region	Total of QtyOnHand
1	P1	North	200
		South	150
		East	50
		West	100
		All	...
2	P2	North	400
...

Report 2. Total of *Quantity on Hand* in the first quarter of 2010, by product category, by month name.

Product Category Quantity on Hand First Quarter 2010		
Product Category	Month	Total of QtyOnHand
C1	January	900
	February	300
	March	500
C2	January	400
...

Report 3. Values of the *Inventory Turns* and *Days in Inventory* in the year 2010, by product category, by quarter name.

1. For simplicity, we assume that a month is 30 days, a quarter is 90 days and one year is 365 days.

Inventory Turns and Days in Inventory Year 2010			
Product Category	Quarter	Inventory Turns	Days in Inventory
C1	Q1
	Q2
	Q3
	Q4
C2	Q1

Give a conceptual and logical data mart designs. For each measure, specify if it is additive, semi-additive or non-additive. Give the SQL queries to produce the data of the reports.

A.5 Hotels

The managers of an international hotel chain are interested in analyzing the degree of use of different types of hotel rooms to determine how to price them.

Every day the rooms may be *occupied*, *vacant* or *unavailable* for maintenance reasons. There are different types of rooms on the basis of the following properties: the type (standard, suite, deluxe), the number of beds, the maximum number of occupants, and optional features, such as Minibar, satellite TV, Internet, whirlpool bath, kitchenette, suite.

For each hotel the information of interest is the name, the location and the category (5 star, 4 star, . . . , 1 star).

The managers are interested in analyzing the daily capacity utilization (*occupancy rate*) of each *room type* using the following metrics:

- The *room occupancy rate*, defined as the ratio of the number of rooms occupied to the total number of rooms (occupied, free and unavailable).
- The *average room revenue*, defined as the ratio of the total revenue for rooms occupied to the number of rooms occupied.
- The *revenue per available room*, defined as the ratio of the total revenue for rooms occupied to the number of rooms available, equivalent to the *average room revenue* \times *room occupancy rate*.

The following are some examples of business questions collected during the user interviews, of interest also by category, region and country of the hotel, and by month, year and day of a holiday date.

1. The *room occupancy rate* of hotels of a given city and day, by hotel.

Occupancy Rate Hotel Best, Florence July 17, 2010	
Hotel	Occupancy Rate
Best 1	47%
Best 2	53%
Best 3	19%

2. The *room occupancy rate* of hotels of a given region and day, by room type.

Occupancy Rate Hotel Best, Tuscany July 17, 2010	
Room Type	Occupancy Rate
Standard	47%
Suite	53%
Deluxe	61%

3. The *room occupancy rate* at a given month and year, by hotel of a given city.

Occupancy Rate Hotel Best, Florence July 2010	
Hotel	Occupancy Rate
Best 1	74%
Best 2	79%
Best 3	60%

4. The *room occupancy rate* and *average room revenue* of hotels in a given city, at a given month and year, by hotel.

Occupancy Rate and Average Room Revenue Hotel Best, Milan July 2010		
Hotel	Occupancy Rate	Average Room Revenue
Best 1	74%	145
Best 2	79%	60
Best 3	60%	75

5. The monthly revenue and the cumulative revenue of 4-star hotels in a given year, by country and by month.
6. In a given year, the total revenue, and the cumulative revenue, of the rooms with the maximum number of occupants and whirlpool bath, by hotel.

Give a conceptual and logical data mart designs to analyze the room type utilization. For each measure, specify if it is additive, semi-additive or non-additive. Give the SQL queries to produce the data of the reports.

CASE STUDIES: SOLUTIONS

It is likely that the solutions shown here will turn out to be not perfect. If you disagree with an answer, please feel free to mail us.

B.1 Hospital

Requirements specification

Each business question is analyzed to identify the dimensions and the measures used, and the aggregations to compute (*metrics*):

			Hospitalization
Requirements analysis	Dimensions	Measures	Metrics
Total billed amount for hospitalizations, by diagnosis code and description, by month (year).	Diagnosis (ICD, Description), Date (Month, Year)	Amount	Total Amount
Total number of hospitalizations and billed amount, by ward, by patient gender (age at date of admission, city, region).	Ward, Patient (Gender, Age, City, Region)	Amount	Total number Total Amount
Total billed amount, average length of stay and average waiting time by diagnosis code and description, by name (specialization) of the physician who admitted the patient.	Diagnosis (ICD code, Description), Physician (Name, Specialization)	Amount, Duration, WaitingTime	Total Amount Average Duration Average WaitingTime
Total billed amount, and average waiting time for admission by patient age (region), by treatment code (description).	Patient (Age, Region), Treatment (Code, Description)	Amount, Duration, WaitingTime	Total Amount Average WaitingTime

From the requirements specification the following fact granularity arises:

	Fact granularity
Description	A fact is a hospitalization of a patient, assuming that they may require one treatment only
Preliminary dimensions	Patient, Date, Ward, Diagnosis, Treatment, Physician
Preliminary measures	Duration, WaitingTime, Amount

The measure **Amount** is **additive**. The measures **Duration** and **WaitingTime** are **non-additive**.

Conceptual Design

The data mart conceptual design is shown in Figure B.1.

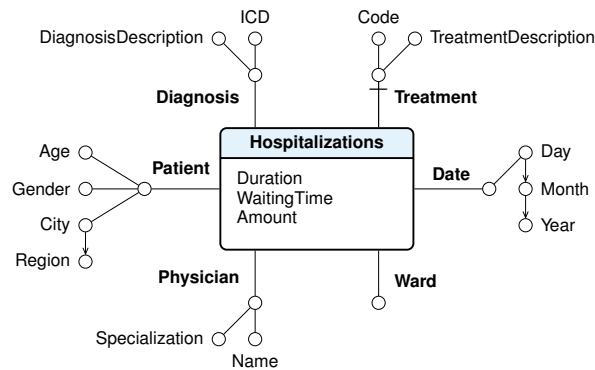


Figure B.1: The conceptual design of a data mart for the hospitalizations

Logical design

In the logical design, the facts are stored in the relation *Hospitalizations*, with the measures, the degenerate dimension *Ward* and a foreign key for each dimension table, with its own surrogate primary key (Figure B.2). The surrogate primary key for the *Date* dimension is a day, an integer of the form YYYYMMDD.

This solution is correct, assuming that if a patient is hospitalized several times with different values of age, its value in the dimension *Patient* is that of the last hospitalization. If we are interested in storing the value of a patient age at each hospitalization, as desired by the requirements, with the admission of a patient with an age different from the one already present in *Patient*, a new record is created in the table *Patient* with a different surrogate primary key (changes dealt with mode Type 2). To find out which data refer to the same patient hospitalizations (for example, to count the different patients hospitalized), *InitialPatientKey* is added as the attribute in the fact table, with the first surrogate key value assigned to a patient (Figure B.3). This solution also allows us to deal with cases in which, at each new hospitalization, the patient also changes the city and region of residence.

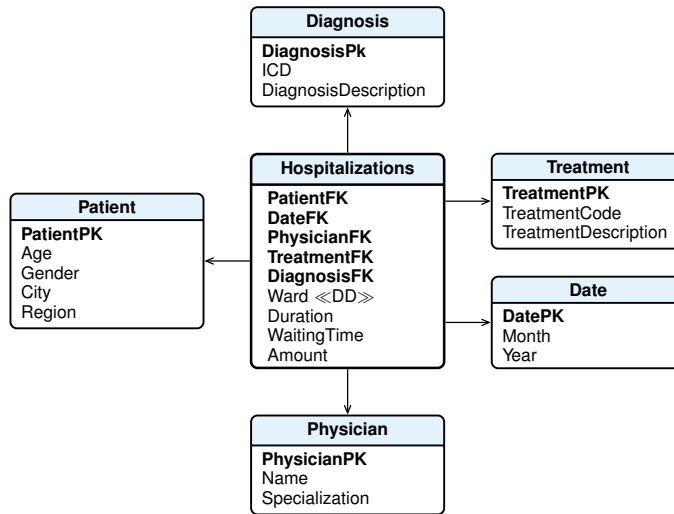


Figure B.2: The initial logical design of a data mart for the hospitalizations

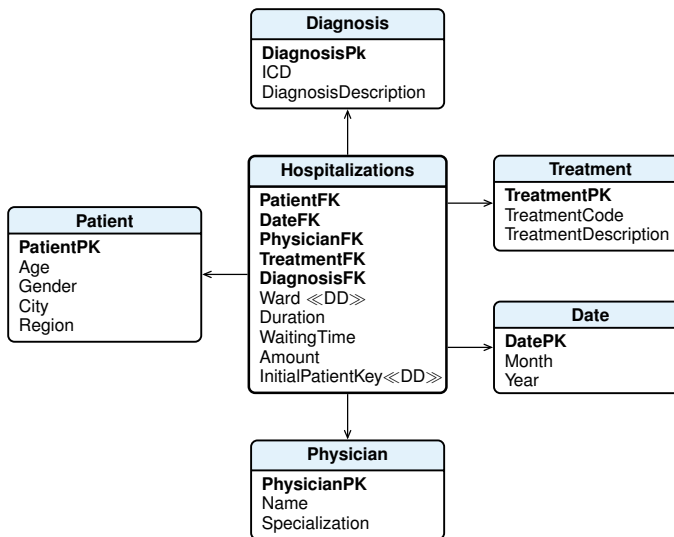


Figure B.3: The final logical design of a data mart for the hospitalizations

B.2 Airline Companies

Requirements specification

Each business requirement analysis is analyzed to identify the dimensions and the measures used, and the aggregations to compute (*metrics*):

Airline companies			
Requirements analysis	Dimensions	Measure	Metrics
Number of unoccupied seats in a given year, by flight code, by company name (or type), by class, by departure time (time, day, month, year)	FlightCode, Class, Company(Name, Type), DepartureTime (Time, Day, Month, Year)	UnoccupiedSeats	Total UnoccupiedSeats
Number of unoccupied seats in a given class and year, by flight code, by company name, by class, by departure (destination) city (country, continent).	FlightCode, Class, Company(Name), DepartureCity (Country, Continent), DestinationCity (Country, Continent)	UnoccupiedSeats	Total UnoccupiedSeats
Number of unoccupied seats and revenue of the Alitalia company, by year, by month, by destination country.	Company(Name), DepartureTime (Month, Year), DepartureCity(Country)	UnoccupiedSeats Revenue	Total UnoccupiedSeats, Revenue

From the requirements specification the following fact granularity arises:

Fact granularity	
Description	A fact is the information on the number of unoccupied seats on a flight of a class of a company
Preliminary dimensions	Class, FlightCode, Company, Departure time, Departure city, Destination city
Preliminary measures	UnoccupiedSeats, Revenue

The measures are **additive**.

Conceptual Design

The data mart conceptual design is shown in Figure B.4.

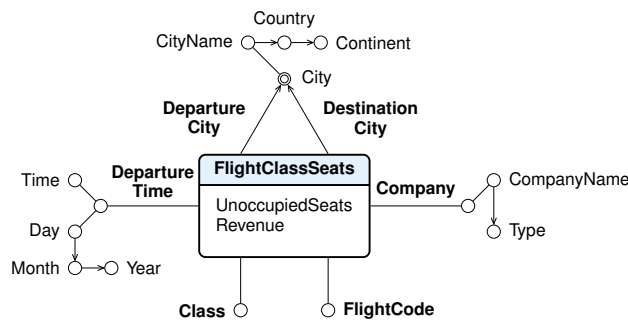


Figure B.4: The conceptual design of a data mart for the airline companies

B.3 Airline Flights

Requirements specification

Each business requirement analysis is analyzed to identify the dimensions and the measures used, and the aggregations to compute (*metrics*):

			Flight Process
Requirements analysis	Dimensions	Measures	Metrics
Number of first-class passengers in a given month and year, by country, by age range of passengers.	Passenger (Nationality, AgeRange), Class, DepartureDate(Month, Year)		Number of passengers
Number of passengers from Europe to the U.S. in a given month and year, and the total revenue, by country, by age range of passengers.	Passenger (Nationality, AgeBand), DepartureDate(Month, Year), DepartureAirport(Continent), DestinationAirport(Country)	Price	Number of passengers, Total price
Number of flights by departure city, by destination city.	DepartureAirport(City), DestinationAirport(City)		Number of flights
Average number of airline passengers by month, by aircraft type, by destination country.	DepartureDate(Month), Aircraft(Type), DestinationAirport(Country)		Average number of passengers
Average number of airline passengers by class, by holiday date.	Class, DepartureDate(HolidayFlag)		Average number of passengers
Number of passengers per year, by size of the destination airport.	DestinationAirport(Size), DepartureDate(Year)		Number of passengers
Number of flights to airports in Germany from the October to December quarter of a given year, and total management cost of the aircraft, by aircraft type.	DepartureDate(Month, Year), DestinationAirport(Country), Aircraft(Type, ManagementCost)		Number of flights, Total management cost
Average profit of all flights, by country of departure, by destination country. The profit of a flight is the total passenger price minus the total flight cost.	DepartureAirport(Country), DestinationAirport(Country), Flight(Duration), Aircraft(ManagementCost, Hourly-OperatingCost)	Price	Average profit
Total revenue in a given year of flights, by month, by destination country. The total revenue by month, total revenue by destination country, and the total revenue are also of interest.	DestinationAirport(Country), DepartureDate(Month, Year)	Price	Total Price

From the requirements specification the following fact granularity arises:

		Fact granularity
Description	A fact is the information on the ticket of a passenger flight	
Preliminary dimensions	Passenger, Flight, Class, Aircraft, Departure Airport, Destination Airport	
Preliminary measures	Price	

The measure **Price** is **additive**.

Conceptual Design

The data mart conceptual design is shown in Figure B.6:

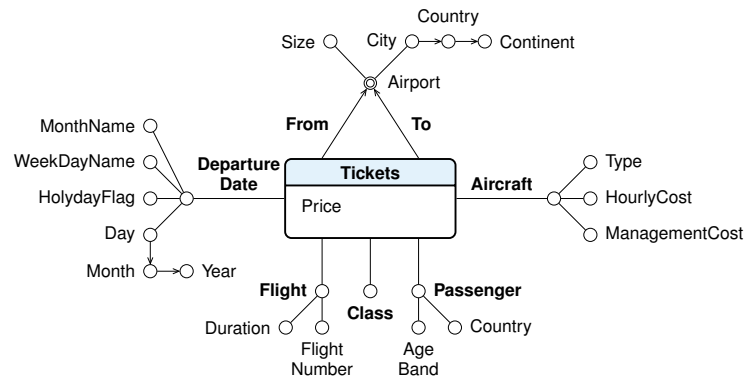


Figure B.6: The conceptual design of a data mart for the airline flights

Logical design

In the logical design, the facts are stored in the relation *Tickets*, with the measures, the degenerate dimension *Class* and a foreign key for each dimension table, with its own surrogate primary key (Figure B.7). The surrogate primary key for the *DepartureDate* dimension is a day, an integer of the form YYYYMMDD.

To simplify the SQL analysis, the degenerate dimension *FlightID* has been added to the fact table to identify the flight of a ticket, with a value the chaining together of the *FlightFK* and *DepartureDateFK* values. If *FlightID* is not used, in the SQL analysis it will be substituted by the expression:

`(CAST (FlightFK AS varchar) + CAST (DepartureDateFK AS varchar))`

The table *Passenger* has as many elements as are the different combinations of *Nationality* and *AgeBand* values.

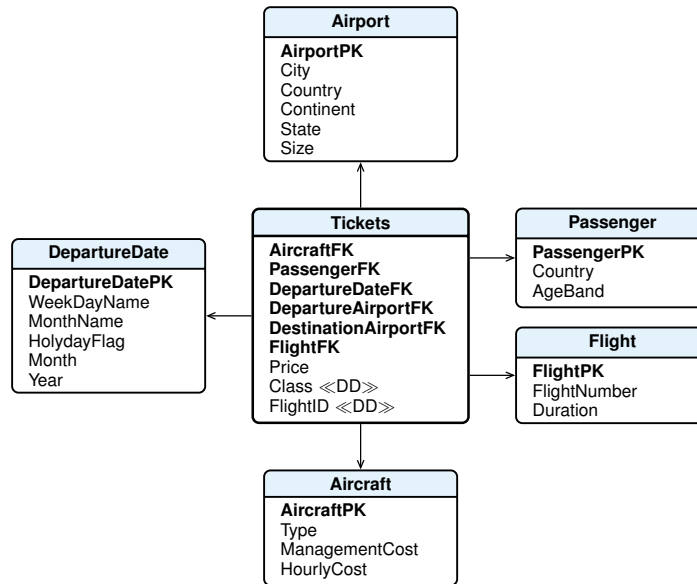


Figure B.7: The logical design of a data mart for the airline flights

Data Analysis

Let us assume that a month is represented as the integer YYYYMM, and a holiday date has the HolidayFlag = true.

1. Number of first-class passengers in a given month and year, by country, by age range of passengers.

```

SELECT    Country, AgeBand, COUNT(*) AS NoOfPassengers
FROM      Tickets, DepartureDate, Passenger
WHERE     PassengerFK = PassengerPK AND DepartureDateFK = DepartureDatePK
           AND Month = 200812 AND Class = 1
GROUP BY Country, AgeBand;

```

2. Number of passengers from Europe to the U.S. in a given month and year, and the total revenue, by country, by age range of passengers.

```

SELECT    Country, AgeBand
           , COUNT(*) AS NoOfPassengers, SUM(Price) AS Revenue
FROM      Tickets, Airport FRM, Airport TO, DepartureDate, Passenger
WHERE     DepartureAirportFK = FRM.AirportPK
           AND DestinationAirportFK = TO.AirportPK
           AND PassengerFK = PassengerPK AND DepartureDateFK = DepartureDatePK
           AND Month = 200812 AND FRM.Continent = 'Europa' AND TO.Country = 'USA'
GROUP BY Country, AgeBand;

```

3. Number of flights, by departure city, by destination city.

```

SELECT   FRM.City AS DepartureCity, TO.City AS DestinationCity
           , COUNT(DISTINCT FlightID) AS NoOfFlights
FROM     Tickets, Airport FRM, Airport TO
WHERE    DepartureAirportFK = FRM.AirportPK
           AND DestinationAirportFK = TO.AirportPK
GROUP BY FRM.City, TO.City;

```

4. Average number of airline passengers by month, by aircraft type, by destination country.

```

SELECT   MonthName , Type AS AircraftType, Country AS DestinationCountry
           , COUNT(*) / COUNT(DISTINCT FlightID) AS AvgNoOfPassengers
FROM     Tickets, Airport, DepartureDate, Aircraft
WHERE    DestinationAirportFK = AirportPK
           AND DepartureDateFK = DepartureDatePK AND AircraftFK = AircraftPK
GROUP BY MonthName, Type, Country;

```

5. Average number of airline passengers, by class, by holiday date.

```

SELECT   Class, DepartureDateFK AS HolydayDate
           , COUNT(*) / COUNT(DISTINCT FlightID) AS AvgNoOfPassengers
FROM     Tickets, DepartureDate
WHERE    DepartureDateFK = DepartureDatePK AND HolydayFlag
GROUP BY Class, DepartureDateFK;

```

6. Number of passengers, by year, by size of the destination airport.

```

SELECT   Year, Size AS SizeDestinationAirport
           , COUNT(*) AS NoOfPassengers
FROM     Tickets, Airport, DepartureDate
WHERE    DestinationAirportFK = AirportPK
           AND DepartureDateFK = DepartureDatePK
GROUP BY Year, Size;

```

7. Number of flights to airports in Germany from the October to December quarter of a given year, and total management cost of the aircraft, by aircraft type.

```

SELECT   Type
           , COUNT(DISTINCT FlightID) AS NoOfFlights
           , COUNT(DISTINCT Month)*ManagementCost AS TotalManagementCost
FROM     Tickets, Airport, DepartureDate, Aircraft
WHERE    DestinationAirportFK = AirportPK
           AND DepartureDateFK = DepartureDatePK AND AircraftFK = AircraftPK
           AND Country = 'Germania' AND Month IN (200710 , 200711 , 200712)
GROUP BY Type, ManagementCost;

```

8. Average profit of all flights by country of departure and by destination country. The profit of a flight is the total passenger price minus the total flight cost.

```

WITH Price-FlightHourlyCost-FlighManagementCost AS
( SELECT Type
  , FRM.Country AS DepartureCountry
  , TO.Country AS DestinationCountry
  , SUM(Price) AS TotalPrice
  , HourlyCost*Duration*COUNT(DISTINCT FlightID)
    AS FlightHourlyCost
  , ManagementCost*COUNT(DISTINCT Month)
    AS FlighManagementCost
FROM Tickets, Airport FRM, Airport TO, Flight, Aircraft, DepartureDate
WHERE DepartureAirportFK = FRM.AirportPK
AND DestinationAirportFK = TO.AirportPK
AND FlightFK = FlightPK
AND AircraftFK = AircraftPK
AND DepartureDateFK = DepartureDatePK
GROUP BY FlightFK, Type, FRM.Country, TO.Country, HourlyCost,
ManagementCost, Duration
)

SELECT DepartureCountry
, DestinationCountry
, (SUM(TotalPrice) –
  SUM(FlightHourlyCost) – SUM(FlighManagementCost))/COUNT(*)
  AS FlightsAvgProfit
FROM Price-FlightHourlyCost-FlighManagementCost
GROUP BY DepartureCountry, DestinationCountry;

```

9. Total revenue in a given year of flights by month and by destination country. The total revenue by month, total revenue by destination country, and the total revenue are also of interest.

```

SELECT MonthName, Country AS DestinationCountry
, SUM(Price) AS TotalRevenue,
FROM Tickets, Airport, DepartureDate
WHERE DestinationAirportFK = AirportPK AND DepartureDateFK = DepartureDatePK
AND Year = 2008
GROUP BY CUBE (MonthName, Country);

```

B.4 Inventory

Requirements specification

From the examples of business questions the following fact granularity arises:

	Fact granularity
Description	A fact is the information on the monthly values of the quantities of products on hand, acquired and shipped
Preliminary dimensions	Product (SKUProduct, Name, Category), Date (Month, Quarter, Year) Warehouse (Name, City, Region, Area)
Preliminary measures	Quantity on hand, Quantity acquired, Quantity shipped

The measures **Quantity acquired** and **Quantity shipped** are **semi-additive** with respect to the dimension **Product**.

The measure **Quantity on hand** is **semi-additive** with respect to both the dimension **Date**, and the dimension **Product**.

The metrics **Inventory Turns** and **Days in Inventory**, defined with a ratio, are **non-additive** and cannot be considered as measures.

Conceptual Design

The data mart conceptual design is shown in Figure B.8.

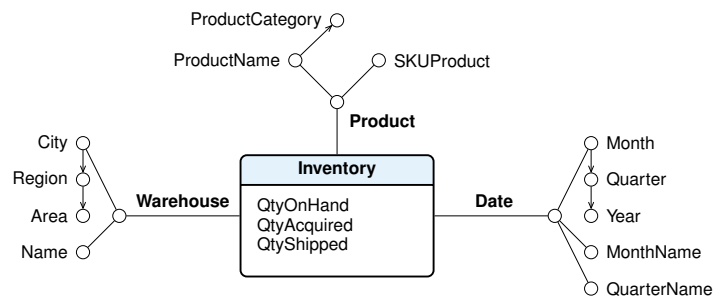


Figure B.8: The conceptual design of a data mart for the Inventory

Logical design

In the logical design, the facts are stored in the relation Inventory, with the measures, and a foreign key for each dimension table, with its own surrogate primary key (Figure B.9). The surrogate primary key for the Date dimension is a month, an integer of the form YYYYMM.

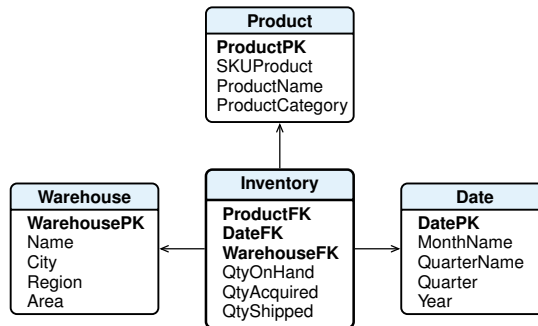


Figure B.9: The logical design of a data mart for the Inventory

Data Analysis

- Report 1.** Total of *Quantity on Hand* in January 2010, by product (SKU and Product Name), by region. The subtotal by all regions is also of interest.

```

SELECT  SKUProduct, ProductName, Region
        , SUM(QtyOnHand) AS TotalQtyOnHand
FROM    Inventory, Product, Warehouse
WHERE   ProductFK = ProductPK AND WarehouseFK = WarehousePK
        AND DateFK = 201001
GROUP BY SKUProduct, ProductName, ROLLUP(Region);
  
```

- Report 2.** Total of *Quantity on Hand* in the first quarter 2010, by product category, by month name.

A value of the attribute Quarter is an integer of the form YYYYQ.

This report has no subtotals, as the previous one, because a subtotal for each category would require totalling the quantities over time, which is meaningless. Adding together the month-end quantities for January, February, and March produces a number that has no meaning. It does not represent the quantity on hand at the end of the period; the March value alone tells us that.

When summing a semi-additive measure such as Quantity on Hand, the dimension across which it is not additive (time) must be used to constrain the query, as in **Report 1**, or the semi-additive measure must be grouped by the dimension in question, as in this report, without a further total or subtotal.

As “subtotals” we can compute the average of the values, but attention is required in correctly computing the average of a set of values as the sum of the values divided by the number of values. In this example the standard SQL average function will not perform this calculation correctly because it assumes as cardinality of a set of values the number of elements of a group of records. For example, if we have two products of the same category available in two warehouses every month of a quarter,

QtyOnHand of product category C1 First Quarter 2010					
Product Category	DateFK	WarehouseFK	ProductFK	Month Name	QtyOnHand
C1	201001	1	1	January	500
C1	201001	2	2	January	400
C1	201002	1	1	February	100
C1	201002	2	2	February	100
C1	201003	1	1	March	200
C1	201003	2	2	March	300

grouping the data on ProductCategory, C1 will appear in 6 records and so

$$\text{AVG}(\text{QtyOnHand}) = \frac{\text{SUM}(\text{QtyOnHand})}{6}$$

while the correct value is

$$\text{AVG}(\text{QtyOnHand}) = \frac{\text{SUM}(\text{QtyOnHand})}{3 \text{ (months of the quarter)}}$$

The problem is avoided by computing the average without using the SQL average function, as follows.

```

SELECT ProductCategory, MonthName AS Month
, SUM(QtyOnHand) / COUNT(DISTINCT DateFK) AS TotalQtyOnHand
FROM Inventory, Product, Date
WHERE ProductFK = ProductPK AND DateFK = DatePK AND Quarter = 20101
GROUP BY ProductCategory, ROLLUP(MonthName);

```

3. **Report 3.** Values of the *Inventory Turns* and *Days in Inventory* in the year 2010, by product category, by quarter name.

The **non-additive** metrics *Inventory Turns* and *Days in Inventory*, must be computed by a “*ratio of sum* and not by a *sum of ratio*”.

```

SELECT ProductCategory, QuarterName AS Quarter
, SUM(QtyShipped) / (SUM(QtyOnHand) / COUNT(DISTINCT DateFK))
AS InventoryTurns
, 90 * (SUM(QtyOnHand) / COUNT(DISTINCT DateFK))
/ SUM(QtyShipped)
AS DaysInInventory
FROM Inventory, Product, Date
WHERE ProductFK = ProductPK AND DateFK = DatePK AND Year = 2010
GROUP BY ProductCategory, QuarterName;

```

B.5 Hotels

Requirements specification

From the requirements the following fact granularity arises :

	Fact granularity
Description	A fact is the information on the daily room type utilization and revenue of each hotel
Preliminary dimensions	Room type, Date, Hotel
Preliminary measures	NOccupiedRooms, NVacantRooms, NUnavailableRooms, NOccupants, Revenue

The dimension **Room type** has as many attributes as the properties of a room, with the attributes for the optional features available with values 'Y' or 'N'.

The measures **NOccupants** and **Reveue** are **additive**.

The measures **NoOccupiedRooms**, **NVacantRooms** and **NUnavailableRooms** are **semi-additive** with respect to **Date**.

The metrics **Occupancy Rate**, **Average Available Room Revenue** and **Average Room Revenue** are **non-additive** and must not be defined as measures.

Conceptual Design

The conceptual design of a data mart is shown in Figure B.10.

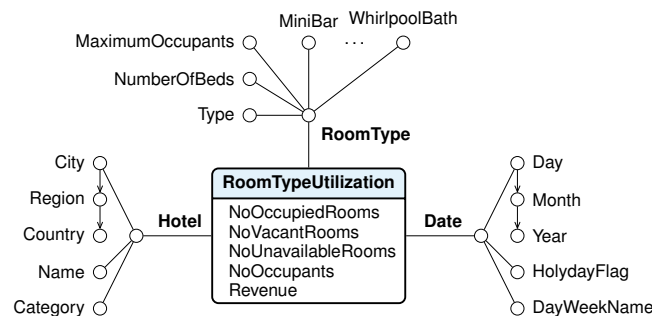


Figure B.10: The conceptual design of a data mart for the hotel room type utilization

Logical design

In the logical design, the facts are stored in the relation `RoomTypeUtilization`, with the measures, and a foreign key for each dimension table, with its own surrogate primary key (Figure B.11). The surrogate primary key for the `Date` dimension is a day, an integer of the form `YYYYMMDD`.

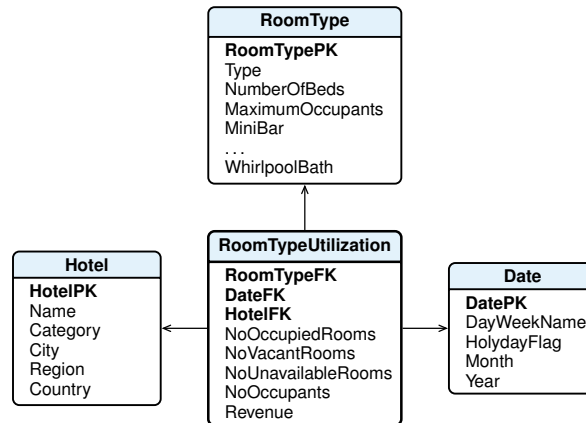


Figure B.11: The logical design of a data mart for the hotel room type utilization

Data Analysis

1. The *room occupancy rate* of hotels of a given city and day, by hotel.

```

SELECT    H.Name
          , SUM(F.NOccupiedRooms) / (SUM(F.NOccupiedRooms) +
                                     SUM(F.NVacantRooms) +
                                     SUM(F.NUnavailableRooms) )
          AS OccupancyRate
FROM      RoomTypeUtilization F, Hotel H
WHERE     F.HotelFK = H.HotelPK AND F.DateFK = 20100717
          AND H.City = 'Florence'
GROUP BY F.HotelFK, H.Name;
  
```

2. The *room occupancy rate* of hotels of a given region and day, by room type.

```

SELECT    R.Type
          , SUM(F.NOccupiedRooms) / (SUM(F.NOccupiedRooms) +
                                     SUM(F.NVacantRooms) +
                                     SUM(F.NUnavailableRooms) )
          AS OccupancyRate
FROM      RoomTypeUtilization F, RoomType R, Hotel H
WHERE     F.HotelFK = H.HotelPK AND F.RoomTypeFK = R.RoomTypePK
          AND H.Region = 'Tuscany' AND F.DateFK = 20100717
GROUP BY R.Type;
  
```

3. The *room occupancy rate* at a given month and year, by hotel in a given city.

```

SELECT      H.Name
           , SUM(F.NOccupiedRooms) / (SUM(F.NOccupiedRooms) +
                                     SUM(F.NVacantRooms) +
                                     SUM(F.NUnavailableRooms) )
           AS OccupancyRate,
FROM        RoomTypeUtilization F, Date D, Hotel H
WHERE       F.HotelFK = H.HotelPK AND F.DateFK = D.DatePK
           AND D.Month = 201007 AND H.City = 'Florence'
GROUP BY   F.HotelFK, H.Name;

```

4. The *room occupancy rate* and *average room revenue* of hotels in a given city, at a given month and year, by hotel.

```

SELECT      H.Name
           , SUM(F.NOccupiedRooms) / (SUM(F.NOccupiedRooms) +
                                     SUM(F.NVacantRooms) +
                                     SUM(F.NUnavailableRooms) )
           AS OccupancyRate
           , SUM(F.Revenue)/SUM(F.NOccupiedRooms) AS AvgRevenueByRoom
FROM        RoomTypeUtilization F, Date D, Hotel H
WHERE       F.HotelFK = H.HotelPK AND F.DateFK = D.DatePK
           AND D.Month = 201007 AND H.City = 'Milan'
GROUP BY   F.HotelFK, H.Name;

```

5. The monthly revenue and the cumulative revenue of 4-star hotels in a given year, by country and by month.

```

SELECT      H.Country, D.Month
           , SUM(F.Revenue) AS MonthlyRevenue
           , SUM(SUM(F.Revenue)) OVER
             (PARTITION BY H.Country ORDER BY D.Month
              ROWS UNBOUND PRECEDING)
           AS CumulativeRevenue
FROM        RoomTypeUtilization F, Date D, Hotel H
WHERE       F.HotelFK = H.HotelPK AND F.DateFK = D.DatePK AND D.Year = 2010
GROUP BY   H.Country, D.Month;

```

6. In a given year, the total revenue, and the cumulative revenue, of the rooms with the maximum number of occupants and whirlpool bath, by hotel.

```

SELECT      F.HotelFK, H.Name, SUM(F.Revenue) AS TotalRevenue
           , SUM(SUM(F.Revenue)) OVER
             (ROWS UNBOUND PRECEDING)
           AS CumulativeRevenue
FROM        RoomTypeUtilization F, RoomType R, Date D, Hotel H
WHERE       F.HotelFK = H.HotelPK AND F.DateFK = D.DatePK
           AND F.RoomTypeFK = R.RoomTypePK
           AND D.Year = 2010 AND R.WhirlpoolBath = 'Y'
           AND F.NOccupants = R.MaximumOccupants
GROUP BY   F.HotelFK, H.Name;

```

GLOSSARY

Aggregation

The result of an aggregate function (sum, count, average, minimum, maximum, etc.) applied to a bag of values.

Business Intelligence

A set of methods and tools for interactive data analysis used primarily by business administrative staff to understand and analyze business performance in order to obtain useful information to support unstructured decision making.

The term *intelligence* is used with the meaning of investigating to find out something interesting, like in *Intelligence Service*.

The *business intelligence* methods and tools are of the following types:

- *Reports*. Reporting is considered the basic level of decision support.
- *Multidimensional data analysis*. Data analysis is usually accomplished interactively with some kind of data analysis tool.
- *Exploratory data analysis*. This data analysis technique is very different from reports and multidimensional analysis: it uses what is called a *discovery technique of useful data models* with *data mining* algorithms.

Computerized Information System

A subset of an information system that use a variety of technologies to process information.

Conformed Dimension

A dimension shared by several fact tables.

Constellation schema

The relational schema of a data warehouse with several fact tables that share dimensional tables.

Cube

A multidimensional cube model (data cube) represent facts with n dimensions by points (a cell) in an n -dimensional space. The cells of the cube contain data measures and the edges of the cube represent the data dimensions.

Although a cube implies only 3 dimensions in geometry, a data cube may represent any number of dimensions.

Some vendors provide OLAP servers that implement the fact table as a data cube using a specialized data structure. Such implementations are referred to as MOLAP (Multidimensional OLAP).

Cuboid

Let us assume that each dimension domain is extended with an additional value “*”. This value has the intuitive meaning “all”, and it represents summarization along the dimension in which it appears, called *cuboid*. A cube can be extended with new “borders” made of cells containing the value of aggregate functions. The extended cube is a generalization of a cross-tabulation, which is 2-dimensional, to n dimensions.

To speed up data analysis, commercial data cube systems precompute all or some of the cuboids and store them as *materialized views* of the data cube.

Data, Information, Knowledge

Data is the representation of certain facts that a computer records, stores and processes. Data, or a condensed form of them, become information when are interpreted in a certain context. Information becomes knowledge when it provides insight upon which the recipient, on the base of his experience, competence, and attitude, can make informed and effective decisions and take proper actions.

Data Mart

A database that has the same characteristics of a data warehouse, but it is focused on a single measurable business process to analyze, and so it has only one fact.

Data Mining

An exploratory data analysis technique to discovery useful data models with specialized algorithms.

Data Warehouse

A decision support database with historical, nonvolatile data, pulled together primarily from operational business systems, structured and tuned to facilitate analysis of the performance of key business processes, worthy of improvement.

The first and still now the most widely cited definition of data warehouse was provided by William Inmon in 1990: “A data warehouse is a subject-oriented, integrated, nonvolatile, and time-varying collection of data in support of management’s decisions.”

A fundamental axiom of the data warehouse is that data is both read-only and non-volatile. As the amount of data within the data warehouse grows, the value of the data increases, allowing a user to perform longer-term analyses of the data. Whereas the operational data is generally real-time or near real-time, data within the data warehouse is historical, since the data warehouse is used primarily for reporting and analyzing relatively large volumes of historical data in an effort to decide what to do in the future.

Data Warehousing

The process used to organize data in a data warehouse and then allow users to analyze them with business intelligence tools.

Data Warehouse Management System (DWMS)

A specialized software for creating and managing large amount of nonvolatile data efficiently and allowing it to be analyzed with OLAP queries. There are three broad directions that have been taken to develop this specialized systems: Relational OLAP (ROLAP), Multidimensional OLAP (MOLAP), Column-Oriented OLAP.

DSS, Decision Support System

A software system used to support decision-making processes within an organization. While an operational system is for performing the business, a decision support systems is for analyzing the business.

Dice

An operator to selects a subcube of a given cube with a selection on two or more

dimensions. The operator does not make aggregations on the data cube.

Dimensional Data Model

A data model that represents measurements of a process and the independent variables that may affect that process. In a dimensional model, data are organized into multiple dimensions and each dimension contains multiple levels of abstraction defined by concept hierarchies. This organization provides the users with the flexibility to view data from different perspectives.

Dimensional Fact Model (DFM)

A conceptual dimensional data model.

Dimension

One of the perspectives that can be used to analyze the data in a data warehouse.

Dimensional Table

The table of a relational database which contains the data for one of the dimensions. The dimensional attributes describe individual characteristics of a dimension.

The dimension table has a primary key (usually a surrogate one) which is used to connect it to the fact table.

The dimension tables in a star schema are intentionally de-normalized.

DOLAP (Desktop OLAP)

A system which manage on a personal computer small amount of data extracted from a multidimensional OLAP server, a DW or an operational DBMS.

Drill-down or Roll-down

An operator to have an aggregated view of the data to a higher level of detail in two ways: by moving down along a dimensional hierarchy level or by adding a dimension of analysis.

ERP (Enterprise Resource Planning)

The meaning of the acronym ERP does not explain the purpose of these systems, which is not the enterprise resource planning, but the integration of business processes in a single software system that can meet all the information requirements of the company using a centralized database .

ETL (Extract, Transform, Load)

A set of back-end data staging steps that are used to (1) obtain data from operational sources (i.e. the extraction step), (2) cleanse and prepare data for import into the data warehouse (i.e. the transformation step), and (3) actually importing the transformed data into the data warehouse (i.e. the loading step).

Fact

A collection of related data items, consisting of measures and context data. Each fact typically represent a business item, a business transaction, or an event that can be used in analyzing the business or key business processes. The most useful data items are indeed numeric and often additive.

Fact Table

The table of a relational database which contains the individual facts being stored in the data warehouse.

There are two types of fields in a fact table: a) The fields storing the foreign keys which connect each particular fact to the appropriate value in each dimension; b)

The fields storing the individual fact measures, such as number, amount, or price. The granularity of the fact table is one of the most significant design decisions in creating a data warehouse. The facts should be as detailed as possible to allow for the data to be viewed from the greatest number of perspectives.

Granularity

The level of detail of the facts stored in a data warehouse, and so the meaning of

a single record in a fact table.

Hierarchy

Dimensional attributes can be arranged into one or more logical structures to analyze data at various levels of detail.

For example, the hierarchy among the attributes City and Region of the dimension Location, states that each city belongs to one region and a region generally contains several cities. The multidimensional analysis usually exploits the hierarchy among the dimensional attributes to perform aggregations of the measures at various levels of detail along the dimensions of the data warehouse. For example, a typical hierarchy is the dimension of time to analyze the facts by year, by quarter, by month or by day.

HOLAP (Hybrid On Line Analytical Processing)

A combined use of Relational OLAP (ROLAP) and Multidimensional OLAP (MOLAP).

Information System

A system whose purpose is to store, process, and communicate information.

Key Business Process

A business process that can be clearly defined, is measurable, and is worthy of improvement.

Measure

A numerical property of a fact useful for evaluating the performance of the processes to be analyzed.

Materialized View

The results of a query stored and automatically used to facilitate the execution of other more complex queries.

Metadata

It is referred to as being the data about data, which defines all aspects of the data contained in a data warehouse including where it originally comes from, its type, what transformations it has been subjected to, where it has been used and what it means from a business perspective.

MOLAP (Multidimensional On Line Analytical Processing)

OLAP systems that store cuboids in a specialized data structures.

OLAP (On Line Analytical Processing)

A category of database software systems that primarily involves aggregating large amounts of data from a data warehouse. The term was introduced to distinguish the activities of data analysis from daily activities on business data organized in databases.

OLAP Client

A system that provides interactive tools for multi-dimensional analysis.

OLAP Server

A system that provides a vision of data to be analyzed as a cube.

OLTP (On Line Transaction Processing)

A category of database software systems that typically involves processing transactions in real time.

Operational Systems

A transaction processing systems to process operational data.

ROLAP (Relational On Line Analytical Processing)

An OLAP system that store data and materialized views in a relational DBMS.

Roll-up

The operator performs aggregation on a data cube, either climbing up a concept

hierarchy for a dimension or by dimension reduction.

Schema

The definition of the logical structure of a database or a data warehouse.

Slice

An operator to selects a cross section that cut across a cube with a selection on one dimension. The result is a subcube, and so the operator does not make aggregations on a data cube.

Snowflake Schema

A variant of the star schema, where some dimension tables are normalized, thereby further splitting the data into additional tables.

Star Schema

The relational schema of a data warehouse with (1) a large central table (fact table) containing the bulk of the data without redundancy, and (2) a set of smaller attendant tables (dimension tables), one for each dimension.

BIBLIOGRAPHY

- Adamson, C. and Venerable, M. (1998). *Data Warehouse Design Solutions*. J. Wiley & Sons, New York. 18
- Artz, J. M. (2005). Data driven versus metric driven data warehouse design. In Wang, J., editor, *Encyclopedia of Data Warehousing and Mining*, pages 223–227. IDEA Group Reference, Hershey, PA, USA. 9, 36
- Ballard, C., Farrell, D. M., Gupta, A., Mazuela, C., and Vohnik, S. (2006). *Dimensional Modeling: In a Business Intelligence Environment*. IBM, <http://www.redbooks.ibm.com/redbooks/pdfs/sg247138.pdf>. 36
- Golfarelli, M., Maio, D., and Rizzi, S. (1998). Conceptual design of data warehouses from E/R schemes. In *Proc. Hawaii Int. Conf. on System Sciences, vol. VII*, pages 334–343, Kona, Hawaii. 15
- Kimball, R. and Ross, M. (2002a). *Data warehouse. La guida completa*. Hoepli Informatica, Milano. 18
- Kimball, R. and Ross, M. (2002b). *The Data Warehouse Toolkit: How to Design Dimensional Data Warehouses. Second Edition*. J. Wiley & Sons, New York. 44
- Moody, D. L. and Kortink, M. A. R. (2000). From enterprise models to dimensional models: A methodology for Data Warehouse and Data Mart design. In *Proceedings of the International Workshop on Design and Management of Data Warehouses (DMDW'2000)*, pages 1–12, Stockholm, Sweden. 41
- Song, I., Rowe, W., Mesker, C., and Ewen, E. (2001). An analysis of many-to-many relationships between fact and dimension table in dimensional modeling. In *Proceedings of the International Workshop on Design and Management of Data Warehouses (DMDW 2001)*, pages 6.1–6.13, Interlaken, Switzerland. 49

