

Open in app ↗

Medium

Search

Write



★ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



# Dimensional Modeling for Data Warehousing



Sanjay Kumar PhD · [Follow](#)

6 min read · Sep 25, 2024



95



1



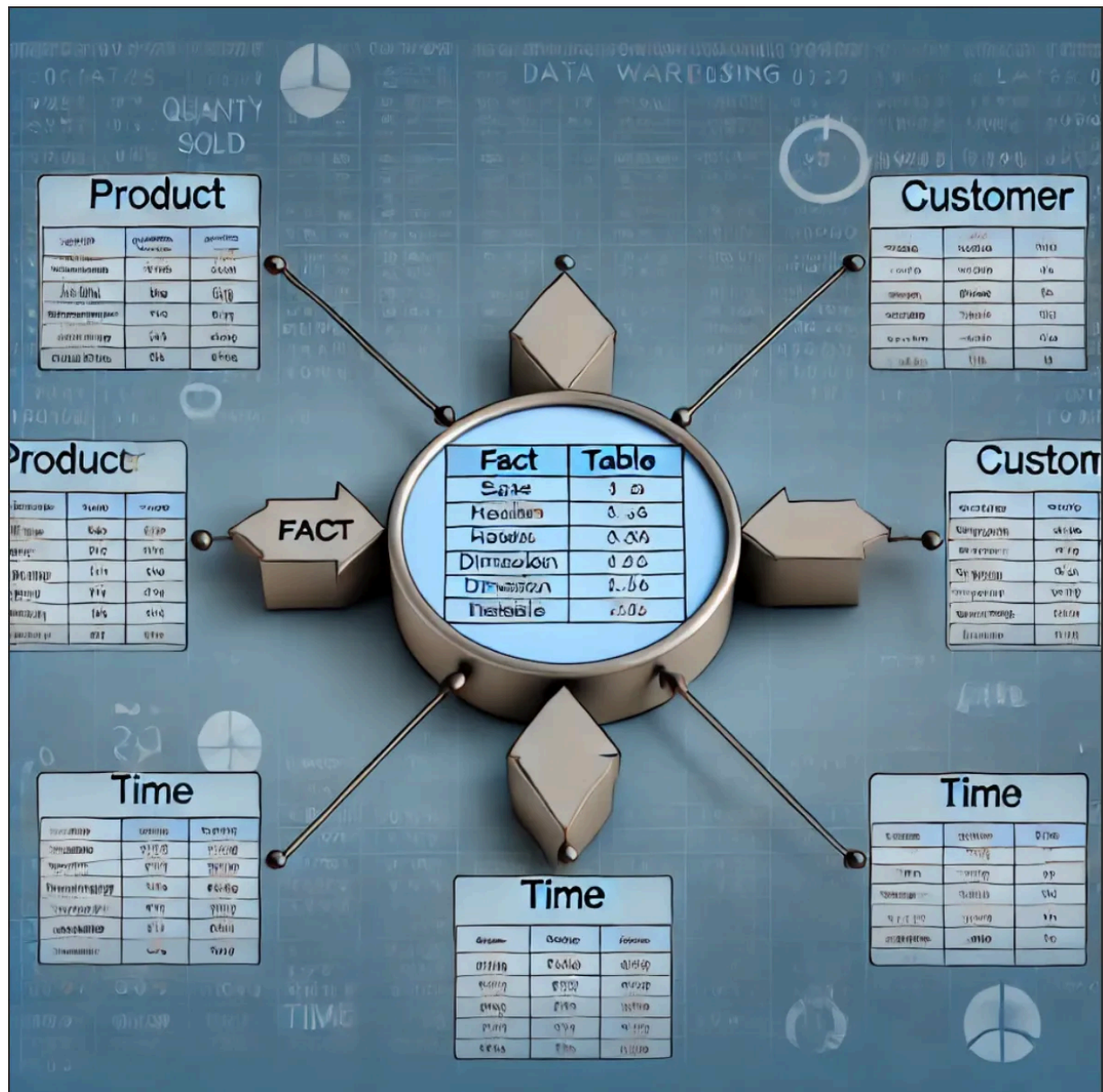


Image Credit : DALL E

Dimensional modeling is a powerful and widely-used approach for organizing data in data warehouses, optimizing it for fast and easy querying. The main goal is to create a structure that simplifies data retrieval and supports business decision-making. In this detailed blog post, we will break down the key concepts of dimensional modeling, explain its importance, and demonstrate how it can be applied effectively in real-world scenarios.

## What is Dimensional Modeling?

Dimensional modeling is a data structure design technique optimized for data warehouses and decision support systems. It aims to simplify complex data sets for non-technical users by creating a design that's easy to navigate, query, and understand. This technique revolves around the use of fact and dimension tables and is integral to structuring data warehouses.

Data warehouses are meant to store large volumes of historical data from various sources (like transactional databases) for analysis and reporting. Dimensional modeling makes it easier for business users to ask questions, run reports, and extract insights from this stored data without having to deal with the complexities of underlying database structures.

At its core, dimensional modeling is built on the following concepts:

- **Fact Tables:** These hold the measurable, quantitative data points (like sales figures, revenue, or quantities sold).
- **Dimension Tables:** These contain descriptive attributes related to the facts, such as time, geography, customer information, or product details.

Together, fact and dimension tables enable users to query and analyze data in ways that reveal meaningful patterns, trends, and insights.

## Star Schema and Snowflake Schema: The Backbone of Dimensional Modeling

Dimensional modeling typically uses one of two types of schemas: **Star Schema** or **Snowflake Schema**. These schemas define how fact and dimension tables are organized and related to each other in a data warehouse.

1. **Star Schema:** The star schema is the simplest and most commonly used schema in dimensional modeling. In this design, a central fact table (representing the “facts” or measurable data) is surrounded by several

dimension tables (representing the context or descriptive information for the facts), forming a star-like shape. This structure is easy to understand and fast to query because there are fewer joins between tables.

- **Example:** In a retail business, a sales fact table might contain data such as sales amount, quantity sold, and discount applied. Dimension tables would include attributes like product, customer, time, and store location, providing additional context for the sales data.

The star schema's simplicity makes it highly performant for querying, as most queries only involve joining the fact table with a few dimension tables.

1. **Snowflake Schema:** The snowflake schema is a more normalized version of the star schema. In this design, dimension tables are further broken down into sub-dimension tables, creating a more complex structure resembling a snowflake. This schema reduces data redundancy by normalizing dimension tables, but at the cost of making queries more complex and sometimes slower due to additional joins.

- **Example:** In a snowflake schema, the product dimension table might be split into several related tables like product category, product subcategory, and product details. This structure saves storage space but requires more complex SQL queries to retrieve the same information.

## Fact Tables: The Core of Business Data

The **Fact Table** is the centerpiece of dimensional modeling. It holds the numerical, measurable data that businesses want to analyze. Fact tables typically contain **metrics** (e.g., sales revenue, number of units sold, or profit margin) and **keys** that link to dimension tables.

Fact tables can be classified into three types:

1. **Transaction Fact Table:** Captures a single transaction or event, such as a purchase made by a customer. Each row represents a specific occurrence, making it the most granular type of fact table.
2. **Snapshot Fact Table:** Captures the state of a process at a particular point in time. For example, a snapshot table might track inventory levels at the end of each month.
3. **Accumulating Snapshot Fact Table:** Tracks the progress of a process over time by updating rows as the process progresses. A typical example would be tracking an order's journey from placement to shipment and delivery.

## Dimension Tables: Adding Context to Facts

**Dimension Tables** provide the context needed to understand the data stored in fact tables. These tables are descriptive and typically contain textual attributes like customer names, geographic regions, or product categories. Each dimension table has a **primary key** that uniquely identifies each record and links it to the fact table.

Dimension tables allow users to filter and group data in a meaningful way. For instance, a time dimension table might include fields for year, month, day, and quarter, enabling users to analyze sales trends over different periods.

## Slowly Changing Dimensions (SCD): Tracking Changes Over Time

Handling changes in dimension data over time is a critical challenge in dimensional modeling, particularly when tracking historical data. This is where the concept of **Slowly Changing Dimensions (SCD)** comes in. SCDs help track changes in attributes of dimension tables over time while maintaining historical accuracy.

There are three common types of SCDs:

1. **Type 1 SCD:** This is the simplest approach, where the old data is overwritten with the new data. There is no historical record maintained. For example, if a customer changes their address, the old address is replaced with the new one.
2. **Type 2 SCD:** In this approach, a new record is added every time a change occurs, preserving historical data. Each record has a unique key, and metadata (such as start and end dates) is used to identify the current and previous versions of the data.
3. **Type 3 SCD:** This method adds a new column to the dimension table to track the change. For instance, a customer dimension table might include both a “previous address” and “current address” column. However, this method is less flexible as it can only track a limited number of historical changes.

## Performance Optimization Techniques in Dimensional Modeling

As the volume of data in a warehouse grows, query performance can degrade. Therefore, optimizing the data warehouse for performance is critical. Some common techniques for performance optimization include:

1. **Fact Table Aggregation:** Aggregating data at higher levels of granularity (e.g., monthly or yearly instead of daily) can improve query performance by reducing the amount of data processed. For instance, you could maintain a summary table that stores total monthly sales for each product category.
2. **Indexing:** Creating indexes on frequently queried columns (such as foreign keys in fact tables) can significantly speed up queries.
3. **Partitioning:** Breaking large fact tables into smaller, more manageable pieces (partitions) based on date ranges or other criteria can help improve query performance.

4. **Derived Tables:** Creating derived tables or materialized views that precompute complex joins or calculations can simplify queries and improve performance. For example, a derived table could store the average sales per customer segment, eliminating the need to compute this on-the-fly.

## Key Design Considerations for Dimensional Modeling

When designing a dimensional model, there are several key considerations to keep in mind:

- **Grain of the Fact Table:** Deciding the level of detail or granularity in the fact table is one of the most important design decisions. For example, should the fact table store data at the transaction level (each individual sale) or the summary level (daily or monthly sales totals)? The grain will determine how much data is stored and how it can be analyzed.
- **Conformed Dimensions:** These are shared dimensions that can be used across multiple fact tables or data marts, ensuring consistency and uniformity. For example, a date dimension used in both the sales and inventory data marts ensures that both have a consistent view of time.
- **Junk Dimensions:** These combine low-cardinality flags and indicators into a single dimension table to avoid clutter. For example, flags such as “is promotional” and “is returned” could be combined into a single junk dimension.

## Conclusion

Dimensional modeling is a cornerstone of data warehousing, providing a simple and intuitive framework for organizing and querying large volumes of business data. By understanding key concepts such as star and snowflake schemas, fact and dimension tables, slowly changing dimensions, and

performance tuning techniques, businesses can design data warehouses that are efficient, scalable, and optimized for decision support.

- Data
- Data Warehouse
- Data Science
- Data Visualization
- Data Analysis



Written by Sanjay Kumar PhD

Follow

257 Followers · 433 Following

AI Product | Data Science| GenAI | Machine Learning | LLM | AI Agents | NLP|  
Data Analytics | Data Engineering | Deep Learning | Statistics

Responses (1)



What are your thoughts?

Respond



Easy Marketing Tips  
4 days ago



Dimensional modeling isn't just about organizing data—it's about making insights accessible. A well-structured star schema transforms raw numbers into business decisions, enabling faster queries and smarter analytics.



Reply