

Week2_Facts_and_Dimensional_Modeling

Facts and Dimensions

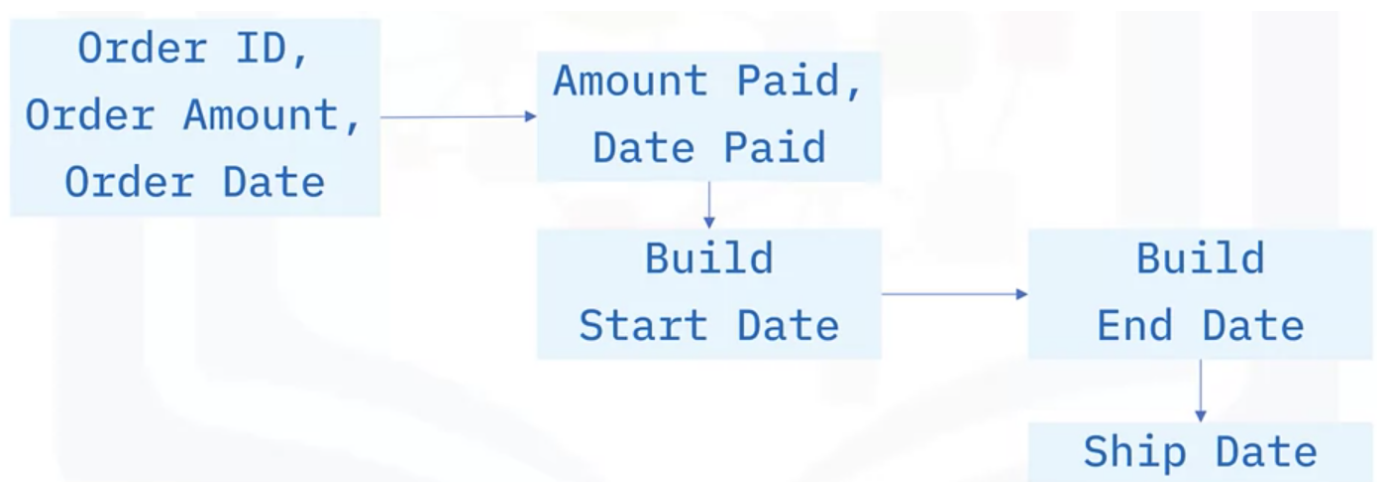
- Data can be categorized as facts and dimensions
- **Facts** are usually measured quantities, such as temperature, number of sales, or mm of rainfall
- Facts can also be qualitative
- Dimensions are attributes relating to facts
- Dimensions provide context to facts (location, date and time)

Fact table

- Facts of a business process, plus
- Foreign keys to dimension tables
 - Dollar amounts for sales transactions
- Can contain detail level facts, or
- Facts that have been aggregated
- Summary tables contain aggregated facts
 - "Quarterly Sales" summary table, with
 - "store_id" as foreign key

Accumulating snapshot fact tables

- Used to record events during a well-defined business process



Dimensions

- Dimensions categorize facts
- Called categorical variables in stats and machine learning

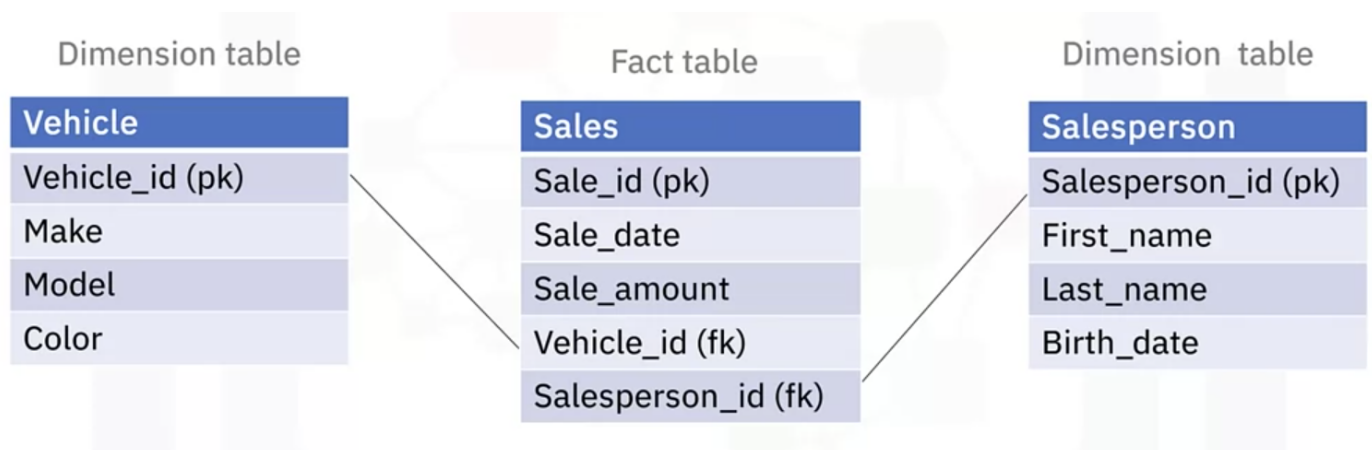
- Used to answer business questions
- Used for filtering, grouping, and labelling
 - People, product, and place names, and date or time stamps
- Dimension table stores dimensions of a fact
 - Joined to fact table via foreign key

Dimension table examples

- Product tables:
 - Make, model, color, size
- Employee tables:
 - Name, title, department
- Temporal tables:
 - Date/time at granularity of recorded events
- Geography tables:
 - Country, state, city, region

Example schema with fact & dimension tables

Each fact table typically has multiple dimension tables related to it.



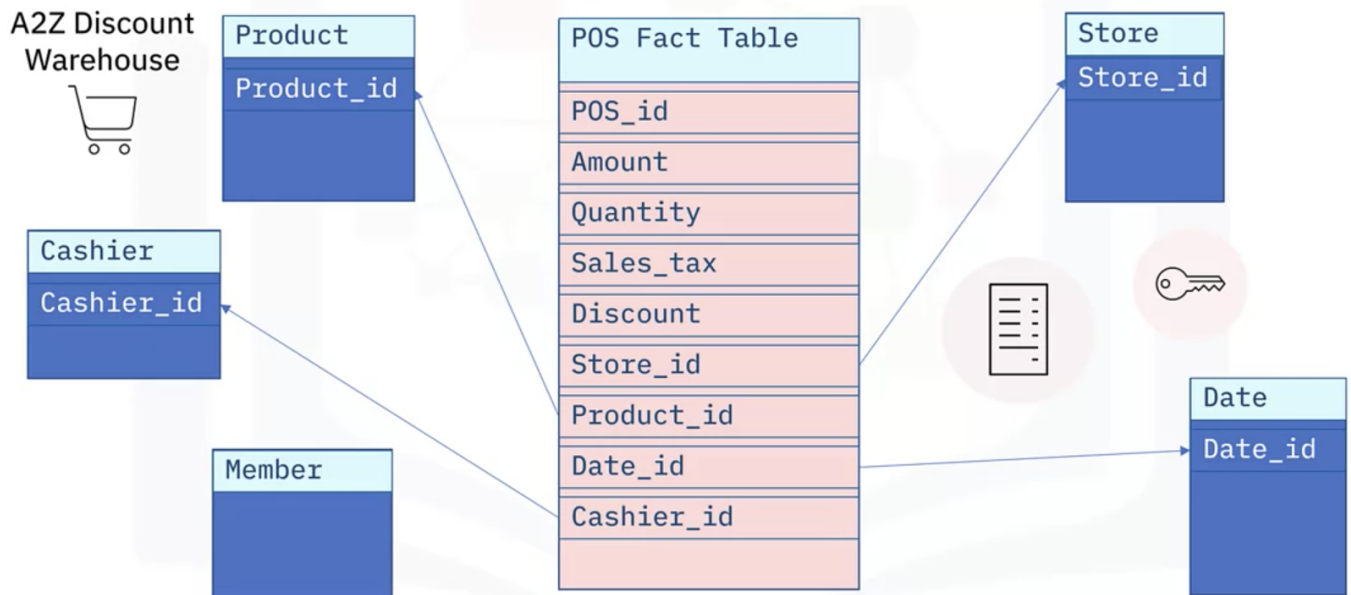
Data Modeling using Star and Snowflake Schemas

Star schemas

- Keys connect facts with dimensions
- Dimensions radiate from a central fact
- Graph whose edges are relations between facts and dimensions

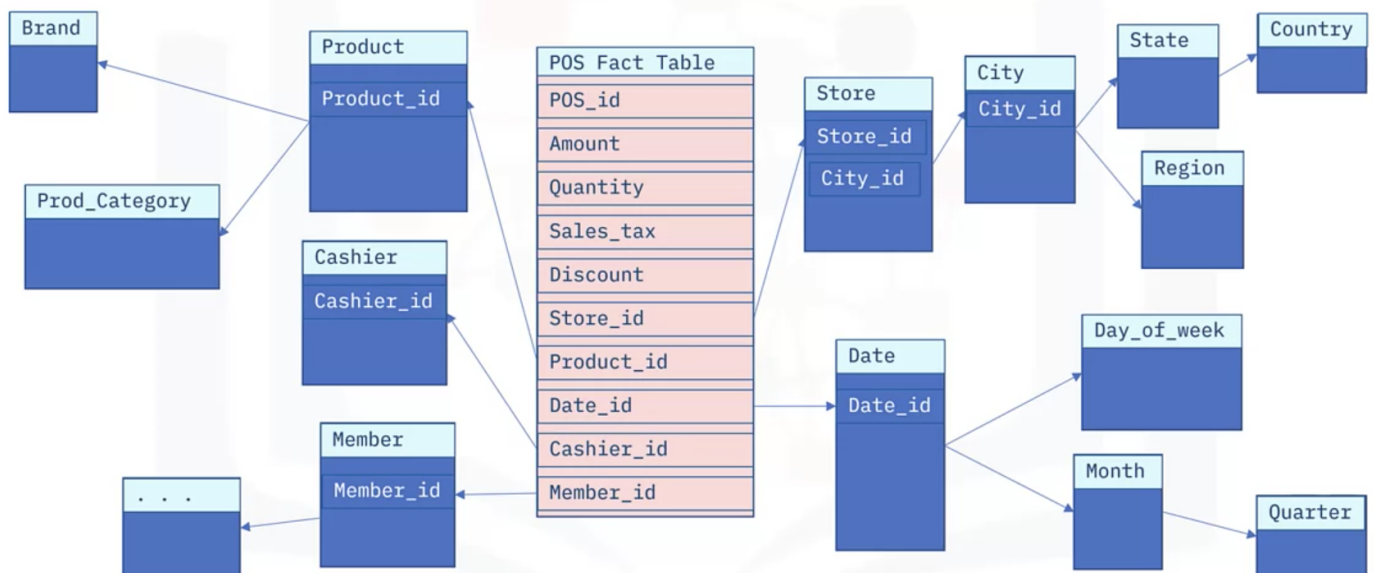
Design considerations

1. Select a business process (sales, manufacturing, supply chain logistics)
2. Choose level of detail/granularity (annual, regional sales number, monthly sales by person)
3. Identify the dimensions (date and time, people)
4. Identify the facts



Snowflake schemas

- Are normalized star schemas
- Dimension tables split into child tables
- Snowflakes don't need to be fully normalized



Star schemas are optimized for reads and are widely used for designing data marts, whereas snowflake schemas are optimized for writes and are widely used for transactional data warehousing. A star schema is a special case of a snowflake schema in which all hierarchical dimensions have been denormalized, or flattened.

Attribute	Star schema	Snowflake schema
Read speed	Fast	Moderate
Write speed	Moderate	Fast
Storage space	Moderate to high	Low to moderate
Data integrity risk	Low to moderate	Low
Query complexity	Simple to moderate	Moderate to complex
Schema complexity	Simple to moderate	Moderate to complex
Dimension hierarchies	Denormalized single tables	Normalized over multiple tables
Joins per dimension hierarchy	One	One per level
Ideal use	OLAP systems, Data Marts	OLTP systems

Both star and snowflake schemas benefit from the application of normalization. “Normalization reduces redundancy” is an idiom that points to a key advantage leveraged by both schemas.

Normalizing a table means to create, for each dimension:

1. A surrogate key to replace the natural key, that is, the unique values of the given column, and
2. A lookup table to store the surrogate and natural key pairs.

Each surrogate key’s values are repeated exactly as many times within the normalized table as the natural key was before moving the natural key to its new lookup table. Thus, you did nothing to reduce the redundancy of the original table.

However, dimensions typically contain groups of items that appear frequently, such as a “city name” or “product category”. Since you only need one instance from each group to build your lookup table, your lookup table will have many fewer rows than your fact table. If there are child dimensions involved, then the lookup table may still have some redundancy in the child dimension columns. In other words, if you have a hierarchical dimension, such as “Country”, “State”, and “City”, you can repeat the process on each level to further reduce the redundancy.

Notice that further normalizing your hierarchical dimensions has no effect on the size or content of your fact table - star and snowflake schema data models share identical fact tables.

Normalization reduces data size

When you normalize a table, you typically reduce its data size, because in the process you likely replace expensive data types, such as strings, with much smaller integer types. But to preserve the information content, you also need to create a new lookup table that contains the original objects.

The question is, does this new table use less storage than the savings you just gained in the normalized table?

For small data, this question is probably not worth considering, but for big data, or just data that is growing rapidly, the answer is yes, it is inevitable. Indeed, your fact table will grow much more quickly than your dimension tables, so normalizing your fact table, at least to the minimum degree of a star schema is likely warranted. Now the question is about which is better – star or snowflake?

Comparing benefits: snowflake vs. star data warehouses

The snowflake, being completely normalized, offers the least redundancy and the smallest storage footprint. If the data ever changes, this minimal redundancy means the snowflaked data needs to be changed in fewer places than would be required for a star schema. In other words, writes are faster, and changes are easier to implement.

However, due to the additional joins required in querying the data, the snowflake design can have an adverse impact on read speeds. By denormalizing to a star schema, you can boost your query efficiency.

You can also choose a middle path in designing your data warehouse. You could opt for a partially normalized schema. You could deploy a snowflake schema as your basis and create views or even materialized views of denormalized data. You could for example simulate a star schema on top of a snowflake schema. At the cost of some additional complexity, you can select from the best of both worlds to craft an optimal solution to meet your requirements.

Practical differences

Most queries you apply to the dataset, regardless of your schema choice, go through the fact table. Your fact table serves as a portal to your dimension tables.

The main practical difference between star and snowflake schema from the perspective of an analyst has to do with querying the data. You need more joins for a snowflake schema to gain access to the deeper levels of the hierarchical dimensions, which can reduce query performance over a star schema. Thus, data analysts and data scientists tend to prefer the simpler star schema.

Snowflake schemas are generally good for designing data warehouses and in particular, transaction processing systems, while star schemas are better for serving data marts, or data warehouses that have simple fact-dimension relationships. For example, suppose you have point-of-sale records accumulating in an Online Transaction Processing System (OLTP) which are copied as a daily batch ETL process to one or more Online Analytics Processing (OLAP) systems where subsequent analysis of large volumes of historical data is carried out. The OLTP source might use a snowflake schema to optimize performance for frequent writes, while the OLAP system uses a star schema to optimize for frequent reads. The ETL pipeline that moves the data between systems includes a denormalization step which collapses each hierarchy of dimension tables into a unified parent dimension table.

Too much of a good thing?

There is always a tradeoff between storage and compute that should factor into your data warehouse design choices. For example, do your end-users or applications need to have precomputed, stored dimensions such as 'day of week', 'month of year', or 'quarter' of the year? Columns or tables which are rarely required are occupying otherwise usable disk space. It might be better to compute such dimensions within your SQL statements only when they are needed. For example, given a star schema with a date dimension table, you could apply the SQL 'MONTH' function as `MONTH(dim_date.date_column)` on demand instead of joining the precomputed month column from the MONTH table in a snowflake schema.

Scenario

Suppose you are handed a small sample of data from a very large dataset in the form of a table by your client who would like you to take a look at the data and consider potential schemas for a data warehouse based on the sample. Putting aside gathering specific requirements for the moment, you start by exploring the table and find that there are exactly two types of columns in the dataset - facts and dimensions. There are no foreign keys although there is an index. You think of this table as being a completely denormalized, or flattened dataset.

You also notice that amongst the dimensions are columns with relatively expensive data types in terms of storage size, such as strings for names of people and places.

At this stage you already know you could equally well apply either a star or snowflake schema to the dataset, thereby normalizing to the degree you wish. Whether you choose star or snowflake, the total data size of the central fact table will be dramatically reduced. This is because instead of using dimensions directly in the main fact table, you use surrogate keys, which are typically integers; and you move the natural dimensions to their own tables or hierarchy of tables which are referenced by the surrogate keys. Even a 32-bit integer is small compared to say a 10-character string ($8 \times 10 = 80$ bits).

Now it's a matter of gathering requirements and finding some optimal normalization scheme for your schema.

Staging Areas for Data Warehouses

A staging area is:

- Intermediate storage for ETL processing
- Bridge between data sources and the target system
- Sometimes transient
- Sometimes held for archiving or troubleshooting
- Used to monitor and monitor ETL jobs

Objectives:

- Integration (consolidation of data from multiple source)
 - Change detection (manage extraction of new and modified data as needed)
 - Scheduling
 - Cleansing data and validation (missing data & duplicates)
 - Aggregation (summarize data)
 - Normalization (enforce consistency of data types and naming conventions)
 - Separate location where data from source system is extracted to
 - Decouples operations such as validation, cleansing and other processes
 - Minimizes corruption risk
 - Simplifies ETL workflows
 - Simplifies recovery
-

Verify Data Quality

Data verification includes checking data for:

- **Accuracy**
 - Source and destination records match
 - Duplicated records
 - Typos
 - Out-of-range values
 - Spelling errors
 - Mass misalignment
 - Misinterpretation of a comma as a separator
- **Completeness**
 - Locating missing values
 - Void or null fields
 - Placeholders like "999" or "-1"
 - Locating missing data records due to system failures
- **Consistency**
 - Non-conformance to standard terms
 - Date formatting
 - Inconsistent data entry
 - "Mr John Doe" and "John Doe"

- Inconsistent units
 - metric and imperial
- **Currency**
 - Avoiding outdated information
 - Update addresses
 - Manage name changes

How to manage and resolve these issues:

1. Write SQL queries to detect and test for these conditions
2. Create rules for correcting and managing those conditions
3. Create a script that runs the data validation sql queries every night
4. Automate the script created in 3 that performs detection and correction processes
5. Review the automated report and address issues that could not be resolved automatically

Data Quality Solutions

OpenRefine

Populating Data Warehouse

Loading frequency

Populating the warehouse is an ongoing process:

- Initial load + periodic incremental loads
- Full refreshes due to schema changes or catastrophic failure are rare
- Fact tables need more frequent updating than dimension tables
- City and store lists change slowly, unlike sales transactions

Typical ways of loading data

Many tools can help you automate the ongoing loading process:

- You can automate loading as part of your ETL pipeline using tools like Airflow and kafka
- tools like bash, python and SQL to build your data pipeline and schedule it with cron

Populating your data warehouse

Preparation:

- Your schema has been modeled

- Your data has been staged
- You have verified the data quality

Setup and initial load:

- Instantiate the data warehouse and its schema
- Create production tables
- Establish relations between tables
- Load your transformed and cleaned data

Ongoing loads:

- Automate incremental loads using a script as part of your ETL pipeline
- Incremental loading requires you to have a method to correctly detect new and changed data
- Normally, you detect changes in the source system
- Many RDBMSs have change tracking to identify new, changed or deleted records
- You might have timestamps identifying when data was first written and when it was modified
- Otherwise, you may need to use brute-force comparison to your staged data

Periodic maintenance

- Perform monthly or yearly data archiving
- Script both the deletion of older data and its archiving to slower, less costly storage

Querying the Data

- `CUBE` and `ROLLUP` provide summary reports
- Easier to implement than alternative sql QUERIES
- Materialized views are stored tables
- Useful for reducing load of complex, frequently requested views on large data sets
- Querying materialized views can be much faster than querying the underlying tables
- Combining cubes or rollups with materialized views can enhance performance