# Implementing Data Versioning & Testing in dbt
by Fariz Wakan

April 24th, 2024

# Data versioning
# with dbt

# Dimensional modeling

- Dimensional modeling is a **data modeling technique** where you break data up into **"facts"** and **"dimensions"** to organize and describe entities within your data warehouse
- The purpose of dimensional modeling is to **optimize the database for faster retrieval of data**

# Dimensional model vs Relational model

- A dimensional model is **designed to read, summarize, analyze numeric information** like values, balances, counts, weights, etc. in a data warehouse. In contrast, relational model is **optimized for addition, updating and deletion of data** in an OLTP system.
- In the relational model, **normalization and ER models reduce redundancy in data**. On the contrary, dimensional model arranges data in such a way that it is easier to retrieve information and generate reports, *denormalized and redundant*.
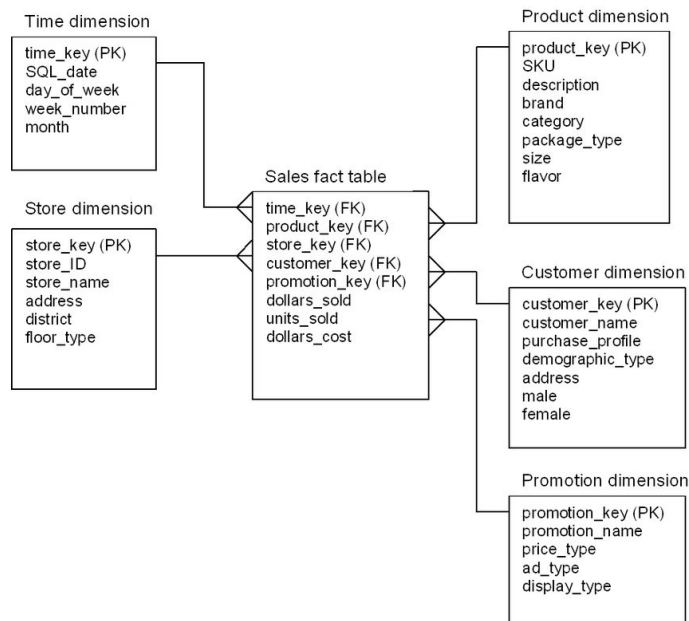
# Facts & Dimensions

- **Facts**
  The **measurements/metrics** or facts of the business process
- **Dimensions**
  The **context** surrounding a business process event. It give who, what, where of a fact.
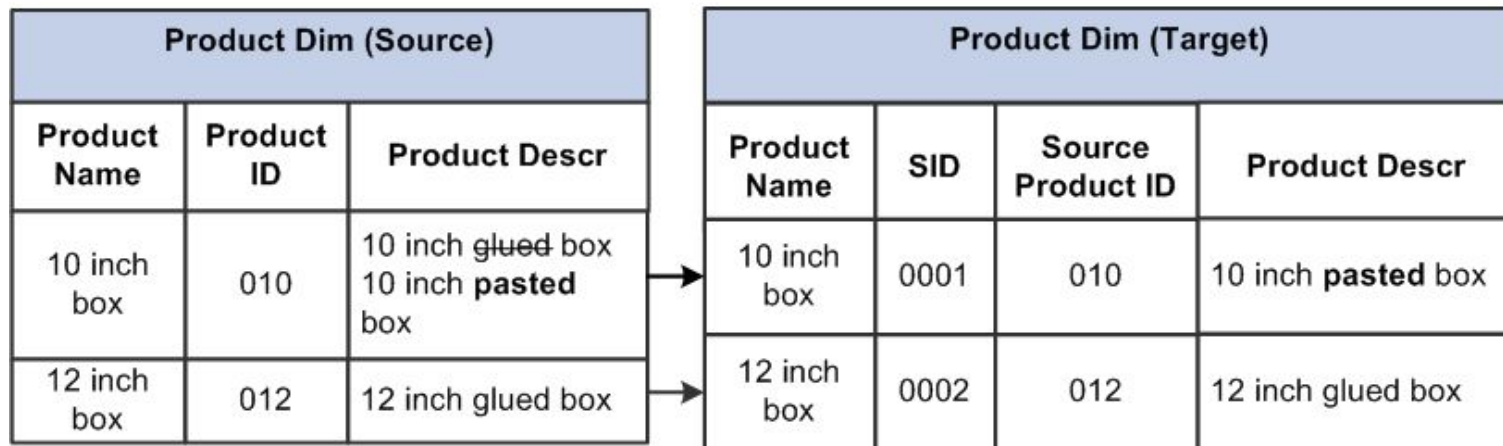
# Fact & Dimension Tables

# Slowly Changing Dimension (SCD)

- SCD is a dimension that **stores and manages both current and historical data over time** in a data warehouse
- The three types of SCDs,
    - Type 1 SCDs - *Overwriting*
    - Type 2 SCDs - *Creating another dimension record*
    - Type 3 SCDs - *Creating a current value field*

# Type 1 SCD

**Type 1 Slowly Changing Dimension**

| Product Dim (Source) | | | Product Dim (Target) | | | |
|---|---|---|---|---|---|---|
| **Product Name** | **Product ID** | **Product Descr** | **Product Name** | **SID** | **Source Product ID** | **Product Descr** |
| 10 inch box | 010 | 10 inch ~~glued~~ box 10 inch **pasted** box | 10 inch box | 0001 | 010 | 10 inch **pasted** box |
| 12 inch box | 012 | 12 inch glued box | 12 inch box | 0002 | 012 | 12 inch glued box |

# Type 2 SCD

## Type 2 Slowly Changing Dimension

| Product Dim (Source) | | | Product Dim (Target) | | | | | |
|---|---|---|---|---|---|---|---|---|
| Product Name | Product ID | Product Descr | SID | Source Product ID | Product Name | Product Descr | EFF_ START_DT | EFF_ END_DT |
| 12 inch box | 012 | 12 inch glued box | 0001 | 012 | 12 inch box | 12 inch glued box | Jan-01-1753 | Dec-31-9999 |
| 10 inch box | 010 | 10 inch glued box 10 inch **pasted** box | 0002 | 010 | 10 inch box | 10 inch glued box | Jan-01-1753 | May-12-06 |
| | | | 0003 | 010 | **10 inch box** | **10 inch pasted box** | **May-12-06** | **Dec-31-9999** |

# Type 3 SCD

| cust_id | customer_number | first_name_pre | fist_name_curr | last_name | effective_date |
|---------|-----------------|----------------|----------------|-----------|----------------|
| 1 | 10001 | NULL | Kontext | Wonderful | 2022-01-01 |

SCD Type 3 Merge

| cust_id | customer_number | first_name_pre | fist_name_curr | last_name | effective_date |
|---------|-----------------|----------------|----------------|-----------|----------------|
| 1 | 10001 | Kontext | Context | Wonderful | 2022-07-01 |

# dbt supports SCD ..

# dbt snapshot

Demo Session

# Data testing
# with dbt

Data
Fellowship

. . .

# dbt test

Demo Session

# References

- https://docs.getdbt.com/docs/introduction

Thank you!