

# Data Modeling documentation - Yalo Case

## Business Overview

Local retail company focused in beverage sales, localized in Iowa-US, which wants to use data analytics to understand financial performance and get insights for next actions, for both finance and sales department (main stakeholders)

## Tech Stack

### Data Source

Considering a paid service that is an API Endpoint with daily information about all liquor sales in Iowa.

### Tools

Considering a rather small company, the data stack will follow a full open source scenario

- Linux Server
- Python
- PostgreSQL
- dbt
- Airflow
- Power BI / Looker Studio

### Ingestion Phase

- Python Script + Airflow

A python script that will use a GET request into the API to gather the **previous day** data, every morning at 8AM.

To avoid risks of losing data or any eventual issue in the script, the code will check the last updated date on the PostgreSQL bronze layer and will get any day that is not present. So, in case of a failed run, it may get the missing days as well

Everything orchestrated by airflow.

At the end of the code, the python script will return a .parquet file for each day in a folder, that would be the raw layer

## Data Modeling

### Bronze Layer

The bronze layer will have a single table called *t\_bronze\_iowa\_sales*, that would be inserted by a Python script in PostgreSQL database db\_bronze\_iowa.

The ingestion type would be **incremental**. Only adding new days into the *t\_bronze\_iowa\_sales*.

At the beginning of dbt model run, it will perform quality tests into this table to check

- Duplicity
- NULL values

Stopping the pipeline in case of any of these issues appears, sending a mail/message to the team to verify what happened.

This table will only be accessed by Data & Analytics engineering team.

<i>t_bronze_iowa_sales</i>	
<i>invoice_and_item_number</i>	TEXT
<i>date</i>	DATE
<i>store_number</i>	TEXT
<i>store_name</i>	TEXT
<i>address</i>	TEXT
<i>city</i>	TEXT
<i>zip_code</i>	TEXT
<i>store_location</i>	TEXT
<i>county_number</i>	TEXT
<i>county</i>	TEXT
<i>category</i>	TEXT
<i>category_name</i>	TEXT
<i>vendor_name</i>	TEXT
<i>vendor_number</i>	TEXT
<i>item_number</i>	TEXT
<i>item_description</i>	TEXT
<i>pack</i>	INTEGER
<i>bottle_volume_ml</i>	INTEGER
<i>state_bottle_cost</i>	NUMERIC(10,2)
<i>state_bottle_retail</i>	NUMERIC(10,2)
<i>bottles_sold</i>	INTEGER
<i>sale_dollars</i>	NUMERIC(12,2)
<i>volume_sold_liters</i>	NUMERIC(12,2)
<i>volume_sold_gallons</i>	NUMERIC(12,2)

## Silver Layer

Here we will apply the Star Schema methodology to model this data, based on the business needs:

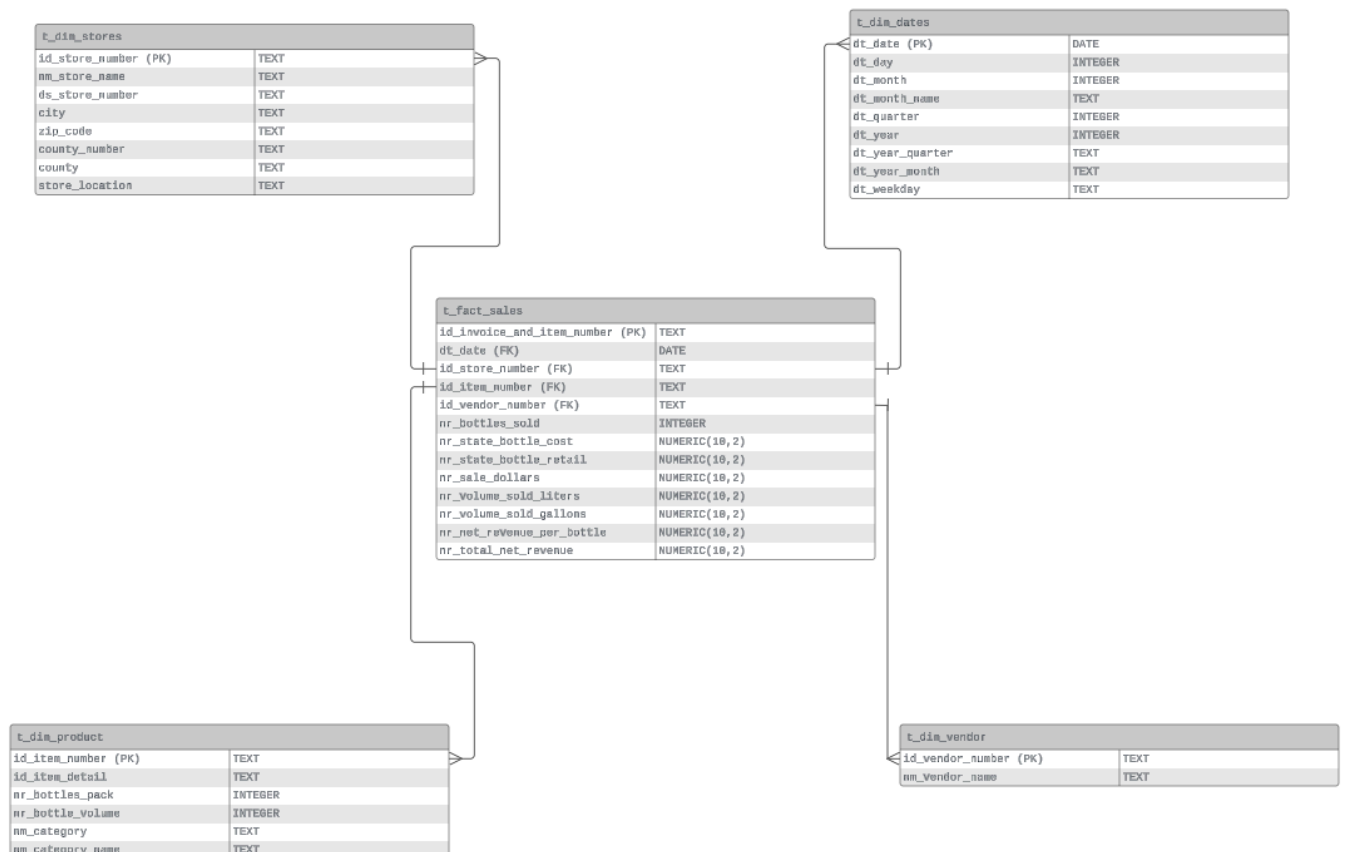
fact Table: `t_fact_sales`

- A materialized table with informations about sales, in a daily granularity

dimension tables:

separated by domain

- dates: `t_dim_dates`
  - with information about weekday, year, quarter, etc
- stores: `t_dim_stores`
  - detailed view of stores and its locations
- products: `t_dim_products`
  - details such as pack size, volume per bottle, etc
- vendors: `t_dim_vendors`
  - vendor names



Everything would be orchestrated by dbt + Airflow, with a daily model run.  
At the end of each run, there will be tests as well.

## Gold Layer

Here the modeling is focused in 2 approaches

### ABT (Analytical Base Table)

A single source of truth for BI teams and Data Scientists, with tested and clean data to link to a BI tool (Looker Studio or Power BI). And for Data Scientists, allowing them to prepare Machine Learning Models (such as Time Series forecasting). In case of a MLOps team, this would be the beginning of their workflow.

### Performance Tables (Materialized Views)

Aggregated materialized views, which could provide status of sales per vendor, product, store and county, allowing business teams to have fast answers and also use it on BI tools, as shown in the chart below

Analytical Base Table  
ABT

t_sales_transactions	
id_invoice_and_item_number	TEXT
dt_date	DATE
dt_day	TEXT
dt_month	TEXT
dt_quarter	TEXT
dt_year	DATE
dt_year_and_quarter	TEXT
dt_year_and_month	TEXT
dt_weekday	INTEGER
sa_store_number	TEXT
sa_store_name	TEXT
sa_address	TEXT
sa_city	TEXT
sa_state	TEXT
sa_region	TEXT
sa_zip_code	TEXT
sa_county	TEXT
sa_item_number	TEXT
sa_item_description	TEXT
pr_bottles_pack	INTEGER
pr_per_bottle_volume	INTEGER
sa_category	TEXT
sa_category_name	TEXT
sa_vendor_number	TEXT
sa_vendor_name	TEXT
pr_bottles_sold	INTEGER
pr_unit_price_per_bottle	NUMERIC(10,2)
pr_unit_price_per_bottle_retail	NUMERIC(10,2)
pr_sale_dollars	NUMERIC(10,2)
pr_revenue_per_bottle	NUMERIC(10,2)
pr_total_revenue	NUMERIC(10,2)
pr_volume_sold_liters	NUMERIC(10,2)
pr_volume_sold_gallons	NUMERIC(10,2)

Aggregate Tables -  
Sales Performance

mv_store_sales_performance	
dt_date	DATE
sa_store_number	TEXT
pr_total_bottles_sold	INTEGER
pr_total_sale_dollars	NUMERIC(10,2)
pr_total_revenue	NUMERIC(10,2)
pr_total_volume_sold_liters	NUMERIC(10,2)
pr_total_volume_sold_gallons	NUMERIC(10,2)

mv_product_sales_performance	
dt_date	DATE
sa_item_number	TEXT
sa_item_description	TEXT
sa_category	TEXT
sa_category_name	TEXT
pr_total_bottles_sold	INTEGER
pr_total_sale_dollars	NUMERIC(10,2)
pr_total_revenue	NUMERIC(10,2)
pr_total_volume_sold_liters	NUMERIC(10,2)
pr_total_volume_sold_gallons	NUMERIC(10,2)
pr_avg_retail_price	NUMERIC(10,2)

mv_vendor_sales_performance	
dt_date	DATE
sa_vendor_number	TEXT
sa_vendor_name	TEXT
pr_total_sale_dollars	NUMERIC(10,2)
pr_total_revenue	NUMERIC(10,2)
pr_total_bottles_sold	INTEGER

mv_county_sales_performance	
dt_date	DATE
sa_county	TEXT
sa_county_number	TEXT
pr_total_bottles_sold	INTEGER
pr_total_sale_dollars	NUMERIC(10,2)
pr_total_revenue	NUMERIC(10,2)
pr_total_volume_sold_liters	NUMERIC(10,2)
pr_total_volume_sold_gallons	NUMERIC(10,2)

At the end of silver layer processing, there will be several tests before creating those tables and materialized views, with actions in case of trigger

- Not null - Stop flow
- Duplicates - Stop flow
- Sales standard deviation - Warning
- Outliers - Warning

## Overall Architecture view:

