



# Data Engineer Bootcamp



## Capstone Project 1



Presented by Bador Alsharif and Sami Alfifi



# Agenda

## Capstone Project 1

01 Project Objectives

02 Data Overview

03 Project Steps

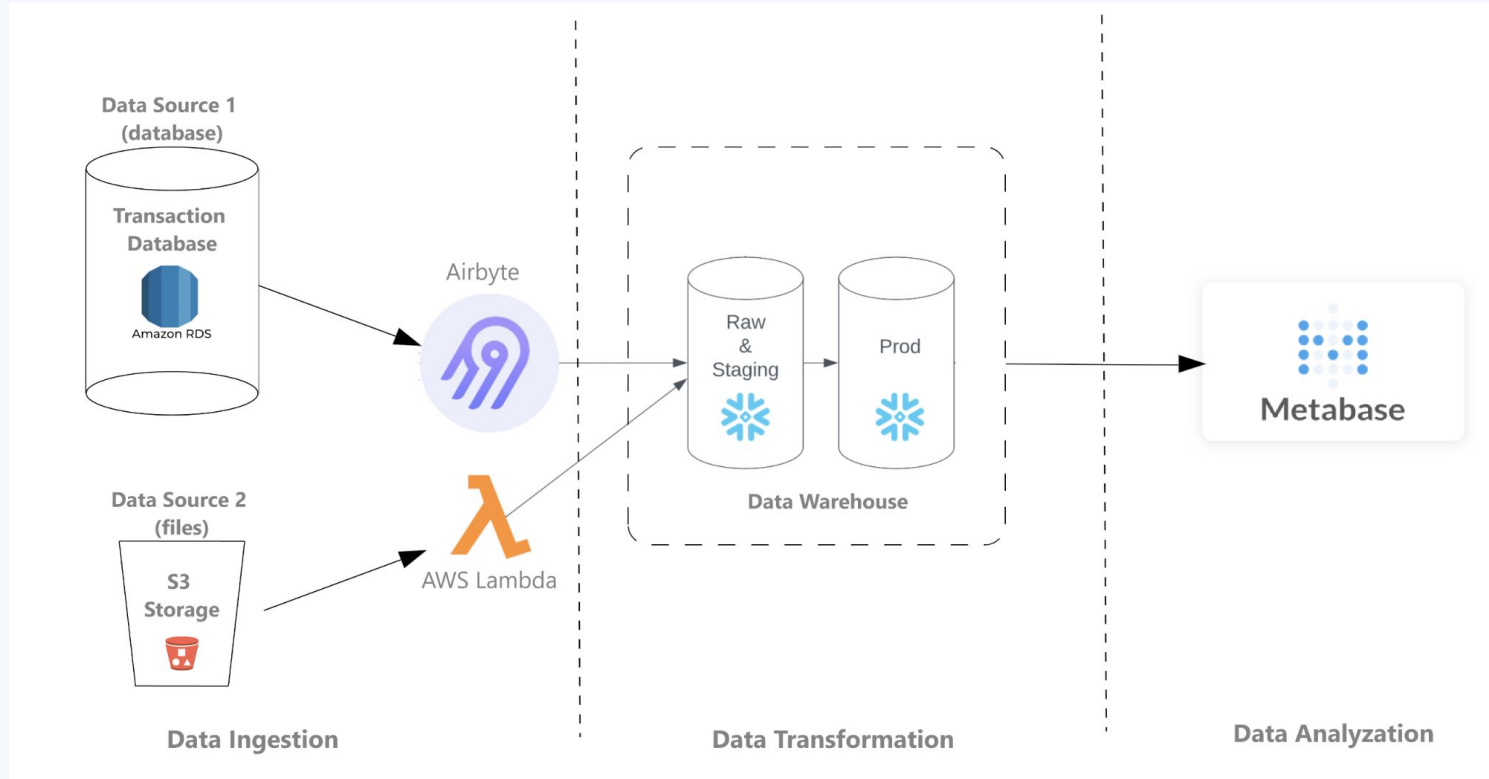
04 Project Achievement

05 Future



WeCloudData

# Project Objectives



# Data Overview

- TPCDS a well-known dataset designed for database testing
- In TPC-DS, fact tables store quantitative data like sales, while dimension tables contain descriptive attributes like customer details.

Fact tables	Dimension tables
Catalog_Sales	Date_Dim
Web_Sales	Customer
Inventory	Item
	Promotion
	Customer_Demographics
	Call_Center
	Customer_Address
	Catalog_Page
	Warehouse
	Time_Dim
	Ship_Mode
	Household_Demographics
	Income_Band
	Web_page
	Web_Site



# Project Steps

01

AWS Lambda Function

02

Airbyte Installation

03

Data Modeling

04

ETL and Data Loading

05

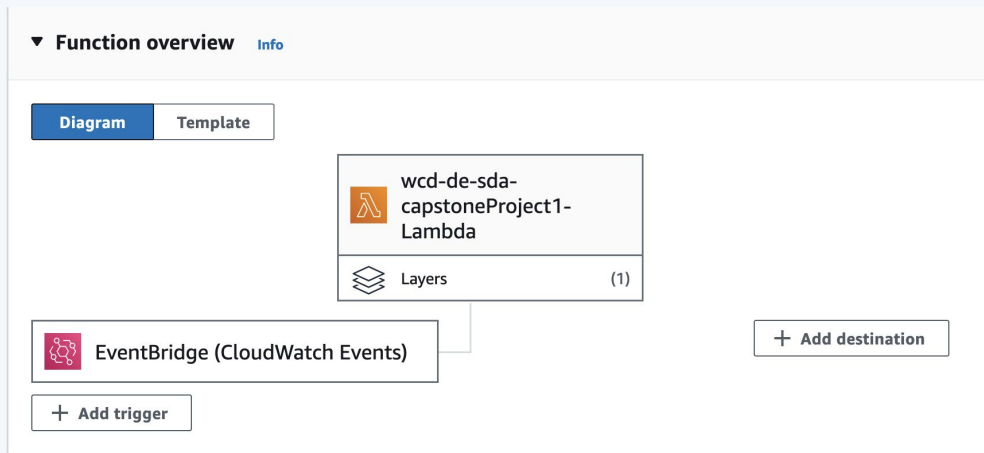
Metabase




WeCloudData

# AWS Lambda Function

1. download 'inventory.csv' from a public AWS S3 URL.
2. upload the 'inventory.csv' file to the Snowflake table.
3. create a schedule for every night 2 am (Riyadh time).



Trigger


 **EventBridge (CloudWatch Events):** [runs\\_at\\_2am\\_Riyadh\\_time](#)  
arn:aws:events:us-east-1:058264394404:rule/runs\_at\_2am\_Riyadh\_time  
Rule state: **ENABLED**

▼ Details

Description: **create a schedule for every night 2 am (Riyadh time).**  
Event bus: **default**  
isComplexStatement: **No**  
name: **runs\_at\_2am\_Riyadh\_time**  
Schedule expression: **cron(0 23 \* \* ? \*)**  
Service principal: **events.amazonaws.com**  
Statement ID: **lambda-0a9067e8-b9cb-4c18-ba74-2f237bd181bd**  
url: **events/home#/rules/runs\_at\_2am\_Riyadh\_time**

# Airbyte Installation

1. Install Airbyte on EC2 instance.
2. Connect to the Postgres database in AWS RDS.
3. Create connections to data sources and set up replication pipelines.

**DE-RDS → Snowflake** Enabled 

DE-RDS → Snowflake

[Status](#) [Job History](#) [Replication](#) [Transformation](#) [Settings](#)

**Configuration** [>](#)


Replication frequency:  
Cron - 0 6 \* \* \* ? UTC


Destination Namespace:  
Destination default

Destination Stream Prefix:  
Mirror source name

Detect and propagate schema changes:  
Ignore

**Sync History** Reset your data Sync now

 Sync Succeeded  
1.29 GB | 2,152,946 emitted records | 2,152,946 committed records | Job id: 9 | 6m 45s

11:54AM 05/27/2024 

# Data Modeling

1. Create an analytics schema within the Snowflake database to hold the data model.
2. Develop an integrated customer dimension table from multiple related tables.
3. Design a weekly fact table containing sales and inventory metrics.

## TPCDS

### ANALYTICS

#### Tables

CUSTOMER\_DIM

DATE\_DIM

ITEM\_DIM

WAREHOUSE\_DIM

WEEKLY\_SALES\_INVENTORY

#### CUSTOMER\_DIM

100K Rows

A	C_SALUTATION	VARCHAR(16777216)
A	C_PREFERRED_CUST_F...	VARCHAR(16777216)
#	C_FIRST_SALES_DATE_SK	NUMBER(38,0)
#	C_CUSTOMER_SK	NUMBER(38,0)
A	C_LOGIN	VARCHAR(16777216)
#	C_CURRENT_CDEMO_SK	NUMBER(38,0)
A	C_FIRST_NAME	VARCHAR(16777216)
#	C_CURRENT_HDEMO_SK	NUMBER(38,0)
#	C_CURRENT_ADDR_SK	NUMBER(38,0)
A	C_LAST_NAME	VARCHAR(16777216)
A	C_CUSTOMER_ID	VARCHAR(16777216)
#	C_LAST_REVIEW_DATE_SK	NUMBER(38,0)

#### WEEKLY\_SALES\_INVENTORY

1.2M Rows

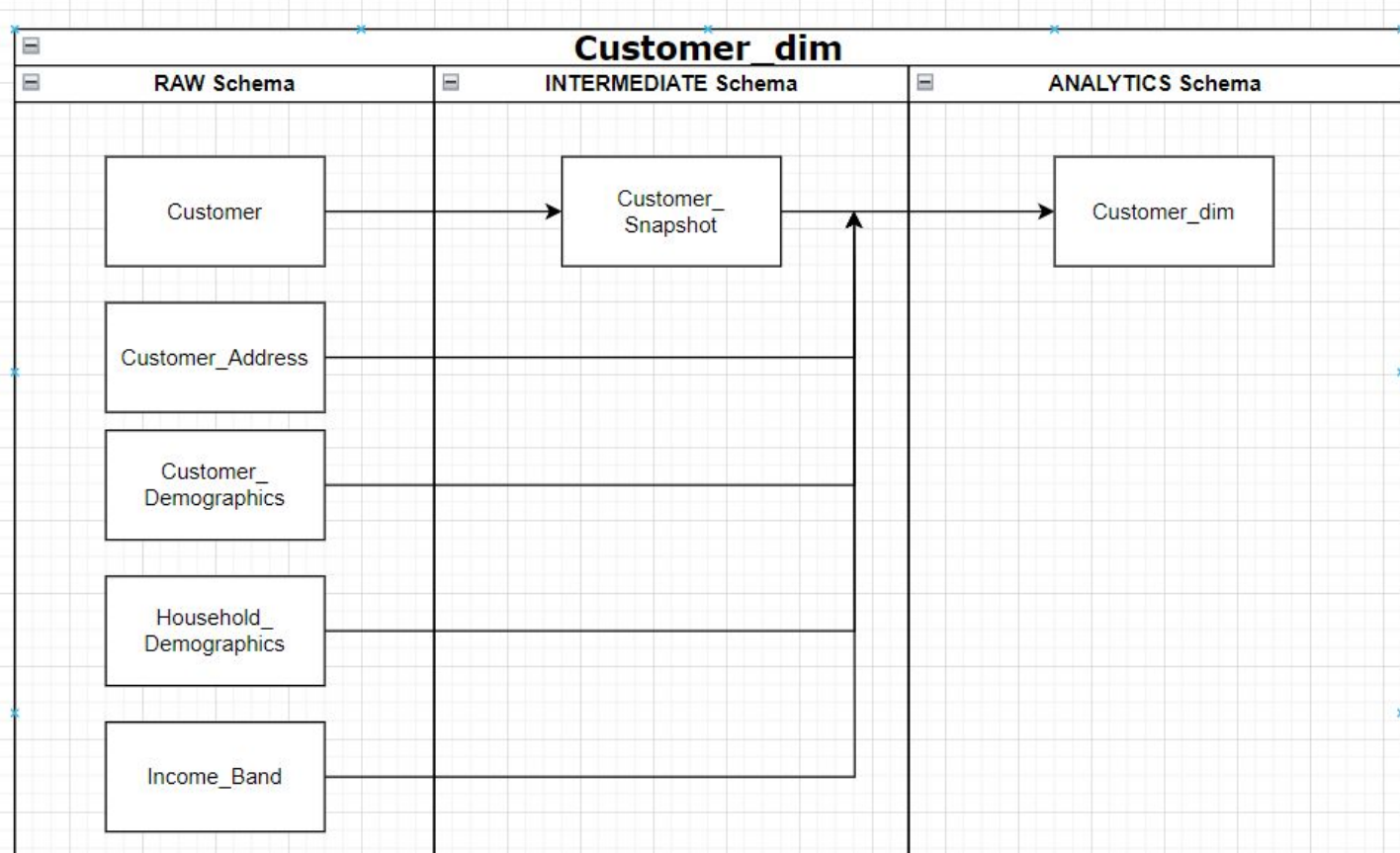
#	WAREHOUSE_SK	NUMBER(38,0)
#	ITEM_SK	NUMBER(38,0)
#	SOLD_WK_SK	NUMBER(38,0)
#	SOLD_WK_NUM	NUMBER(38,0)
#	SOLD_YR_NUM	NUMBER(38,0)
#	SUM_QTY_WK	NUMBER(38,0)
#	SUM_AMT_WK	FLOAT
#	SUM_PROFIT_WK	FLOAT
#	AVG_QTY_DY	NUMBER(38,6)
#	INV_QTY_WK	NUMBER(38,0)
#	WKS_SPLY	NUMBER(38,6)
0 1	LOW_STOCK_FLG_WK	BOOLEAN



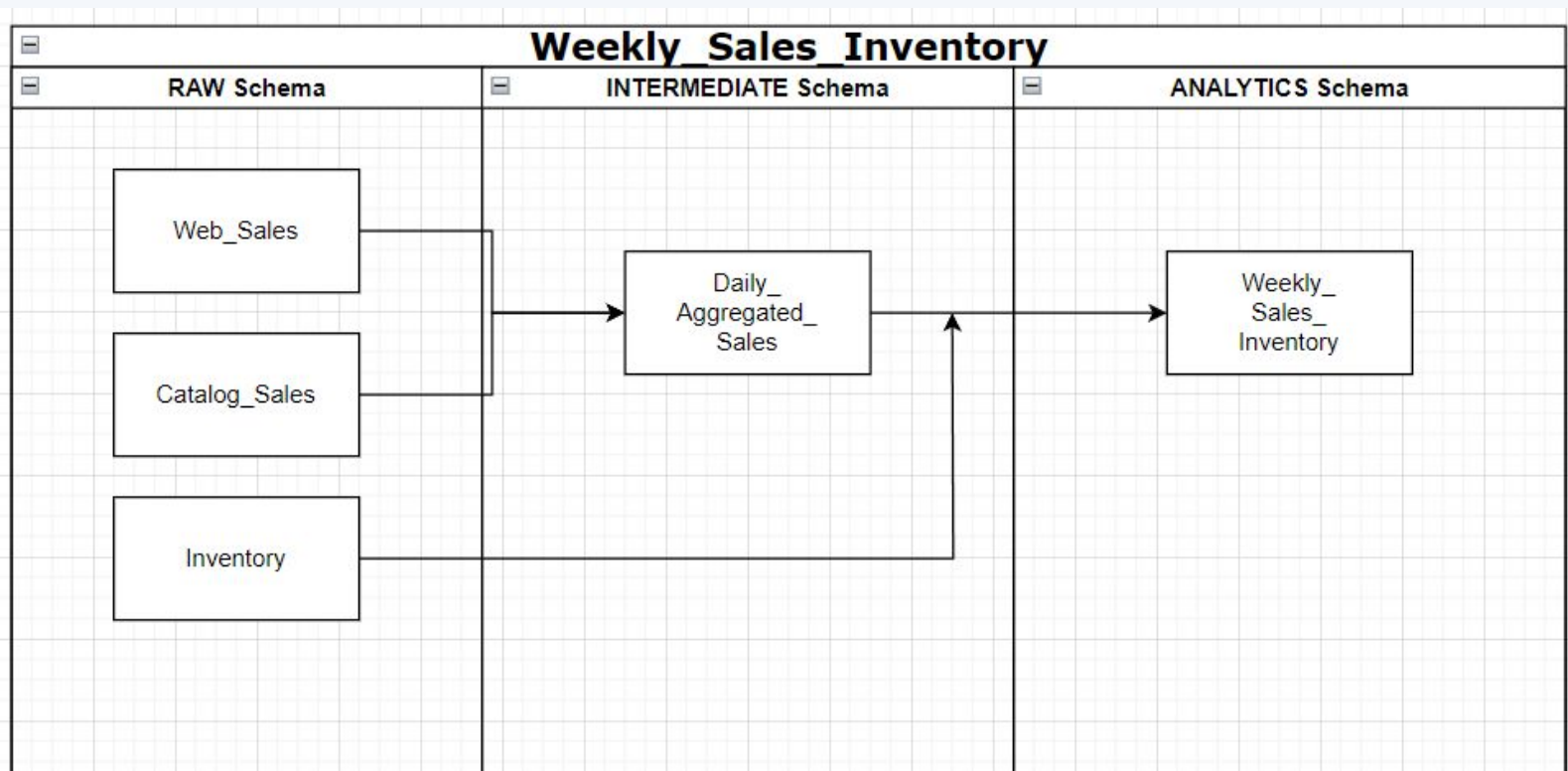
WeCloudData



# Data Modeling - (Customer Dimension)

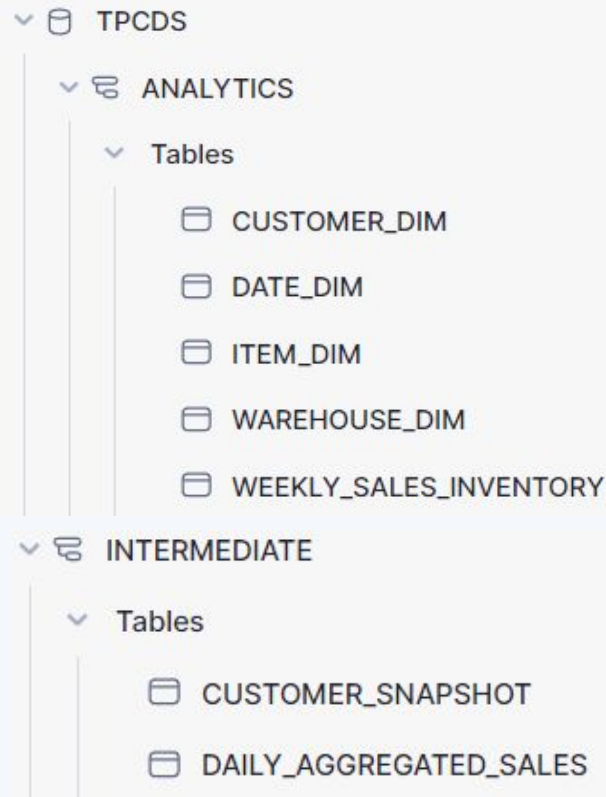


# Data Modeling - (weekly\_sales\_inventory)

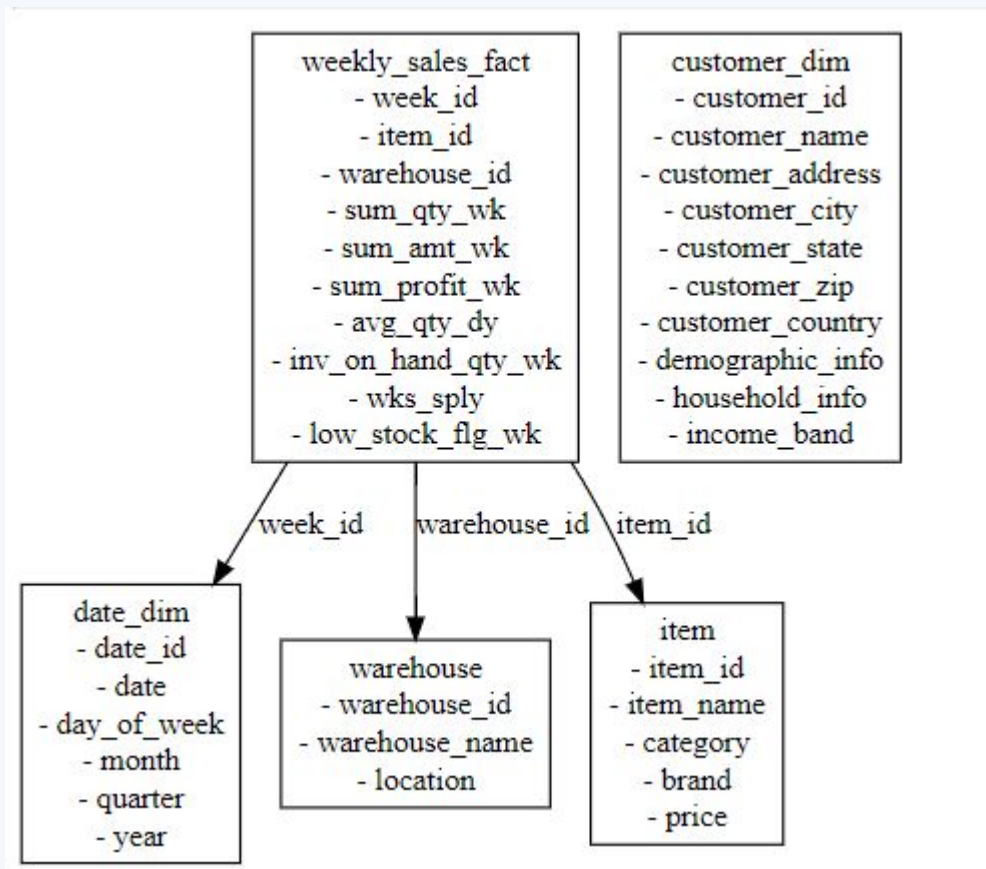


# ETL and Data Loading

1. Implement merge scripts for integrating customer dimension and daily sales records.
2. Create SQL queries for joining daily sales and updated inventory tables to generate the weekly sales and inventory fact table.
3. in order to realize the requirements from the Metabase BI purpose, we also need the Date\_dim, Warehouse, and Item tables in the Prod schema.



# ETL and Data Loading



# Testing

1. Not\_null, Unique, Relationship, Accepted Value, Value Range, Length Tests.
2. Foreign Key Integrity Test to ensuring foreign key relationships are valid.

	TEST_NAME	RESULT
1	Not Null Test	Pass
2	Unique Test	Pass
3	Relationship Test	Pass
4	Accepted Value Test	Pass
5	Foreign Key Integrity Test	Pass
6	Value Range Test	Pass
7	Length Test	Pass



# Automating

1. Set up automated processes for daily and weekly data aggregations.
2. Create SQL queries for joining daily sales and updated inventory tables to generate the weekly sales and inventory fact table.
3. Configure cron jobs using Snowflake Tasks to execute daily and weekly aggregation scripts at specified times.

## ANALYTICS

### Tables

### Tasks

🔄 CREATING\_CUSTOMER\_DIME...

🔄 CREATING\_WEEKLY\_AGGREG...

### Procedures

🔄 POPULATING\_CUSTOMER\_DI...

🔄 POPULATING\_WEEKLY\_AGGR...

## INTERMEDIATE

### Tables

### Tasks

🔄 CREATING\_DAILY\_AGGREGAT...

### Procedures

🔄 POPULATING\_DAILY\_AGGREG...



# Metabase (Business Requirements)

1. Install and configure Metabase for data visualization and report generation.
2. Connect Metabase to the Snowflake database.
3. Fix field names, and assign types to the columns to ensure proper visualization.
4. Generate reports based on business requirements using dashboards.



Top Selling Items

CAL_WK	ITEM_SK	SUM_AMT_WK	SUM_QTY_WK	RANKING
June 12, 2022	12,748	39,543.32	288	1
June 12, 2022	15,976	27,140.85	244	2
June 12, 2022	14,959	26,416.16	172	3
June 12, 2022	13,734	26,333.56	97	4
June 12, 2022	9,700	25,858.86	173	5
June 12, 2022	2,257	25,597.49	221	6

Rows 1-6 of first 2000 < >

Bottom Selling Items

CAL_WK	ITEM_SK	SUM_AMT_WK	SUM_QTY_WK	RANKING
September 11, 2022	14,467	0	1	1
September 11, 2022	625	0	3	2
September 11, 2022	9,038	0	5	3
September 11, 2022	17,617	0	5	3
September 11, 2022	4,477	0	6	5
September 11, 2022	9,430	0	9	6
September 11, 2022	7,231	0	13	7
September 11, 2022	9,266	0	15	7

Rows 1-8 of first 2000 < >



# Metabase (Additional Requirements)

Monthly Sales Amount

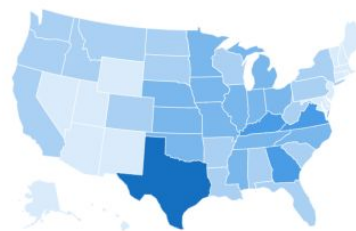


Customer Countries



Customer States

- 26 - 1,003
- 1,181 - 2,320
- 2,435 - 3,303
- 3,849 - 4,955
- 7,799 +





# Challenges

1. **Data Quality and Consistency:**  
Implement data validation and cleansing routines to ensure data quality before aggregation.
2. **Complex Transformations:**  
Ensuring that the logic for these calculations is robust and handles edge cases properly.
3. **Deployment and Migration:**  
Migrating existing data to the new schema and ensuring that historical data is accurately transformed to include the new metrics.



# Project Achievement

1. Exploring the data and understanding business requirements.
2. Merging data in temporary tables to avoid affecting tables in use.
3. Testing tables to validate the results before automating the pipelines.



# Future

1. Enhanced Data Quality Measures:

Implementing robust data cleansing and validation techniques to ensure high data quality and accuracy.

2. Advanced Analytics Capabilities:

Incorporating advanced analytics techniques such as predictive modeling or machine learning to derive deeper insights from the data.

3. Security Enhancements:

Implementing additional security measures to protect sensitive data throughout the entire data lifecycle.





Thank you!