

# **Banking Dataset: Predictive Marketing Insights**



**University of Maryland, Baltimore County**

**Department of Data Science**

**Date: 05/14/2024**

**Project Advisor: Mesfin Abate**

**Group Members:**

Lokesh Reddy Venna

Upendra Madaraboena

Charan Pathakamuri

Sashank Gandikota

# **Contents**

## **1. Business Process and Requirements**

### **1.1 Database System Overview**

#### **1.1.1 Customer Segmentation and Benefit**

#### **1.1.2 Key Components of the Database System**

#### **1.1.3 Benefits to Customers**

### **1.2 Data Storage and Capture Requirements**

#### **1.2.1 Data to be Stored**

#### **1.2.2 Data Types**

#### **1.2.3 Data Entry/Capture Methods**

#### **1.2.4 Reports Needed**

### **1.3 Project Risks**

## **2. Design Specifications**

### **2.1 Relational Concepts**

#### **2.1.1 ER Diagram and Technology**

### **2.2 Big Data Concepts**

#### **2.2.1 Data-Warehouse Structure**

#### **2.2.2 ETL**

#### **2.2.3 Power-BI**

### **2.3 Cost Estimates**

#### **2.3.1 Hardware**

#### **2.3.2 Software**

#### **2.3.3 Staff Resources**

### **2.4 Technologies Used**

### **2.5 Teammates contributions**

## **1.Business Process and Requirements**

### **1.1. Database System Overview:**

The database system for this project is designed to support the direct marketing campaigns of a Portuguese banking institution focused on selling term deposits. The primary goal is to predict whether a client will subscribe to a term deposit based on various customer attributes and interaction history.

#### **1.1.1. Customer Segmentation and Benefit:**

The database system categorizes customers based on demographics, financial history, and previous interactions with the bank. By analyzing this data, the system aims to identify potential customers who are most likely to subscribe to a term deposit. This segmentation allows the bank to tailor its marketing efforts more effectively, resulting in higher conversion rates and increased revenue.

#### **1.1.2. Key Components of the Database System:**

- ❖ **Customer Information:** Includes demographic details such as age, education, marital status, and job type. Financial information such as income level, and loans(housing and personal).
- ❖ **Interaction History:** Tracks previous interactions between the bank and customers, including past term deposit subscriptions, responses to marketing campaigns, number of campaigns per individual.
- ❖ **Predictive Analytics Models:** Utilizes the Relational database management system to analyze historical data and predict the likelihood of a customer subscribing to a term deposit. This system trained on features such as customer demographics, financial behavior, and response patterns.

#### **1.1.3. Benefits to Customers:**

- ❖ **Personalized Offers:** Customers receive targeted offers based on their financial needs and behavior, enhancing their banking experience.
- ❖ **Timely Communication:** Customers are contacted through preferred channels at optimal times, reducing unnecessary marketing outreach and improving engagement.

- ❖ Financial Planning: By considering term deposits, customers can make informed decisions about their savings and investments, potentially increasing their wealth over time.

## **1.2. Data Storage and Capture Requirements:**

### **1.2.1. Data to be Stored:**

- ❖ Bank client data including age (numeric), job type (categorical), marital status (categorical), education level (categorical), credit default status (binary), average yearly balance (numeric), housing loan status (binary), personal loan status (binary).
- ❖ Last contact details including communication type (categorical), last contact day of the month (numeric), last contact month of the year (categorical), last contact duration in seconds (numeric).
- ❖ Campaign-related attributes such as the number of contacts performed during this campaign and for this client (numeric), days since the client was last contacted from a previous campaign (numeric), number of contacts performed before this campaign and for this client (numeric), outcome of the previous marketing campaign (categorical).
- ❖ Output variable: whether the client subscribed to a term deposit (binary).

### **1.2.2. Data Types:**

Age: Numeric (Integer)

Job: Categorical (String)

Marital Status: Categorical (String)

Education: Categorical (String)

Default: Binary (String: "yes"/"no")

Balance: Numeric (Float)

Housing: Binary (String: "yes"/"no")

Loan: Binary (String: "yes"/"no")

Contact: Categorical (String)

Day: Numeric (Integer)

Month: Categorical (String)

Duration: Numeric (Float)

Campaign: Numeric (Integer)

Pdays: Numeric (Integer)

Previous: Numeric (Integer)

Poutcome: Categorical (String)

Y (Output): Binary (String: "yes"/"no")

#### 1.2.3. Data Entry/Capture Methods:

- ❖ Automated Data Capture: Utilize automated systems to capture data from forms, CRM systems, or other digital sources. This ensures accuracy and efficiency in data entry.
- ❖ Manual Data Entry: In cases where automated capture is not possible, trained personnel should manually enter data into the system. Double-checking and validation protocols should be in place to minimize errors.
- ❖ Real-time Updates: Ensure that data is updated in real-time or at regular intervals to reflect the latest customer interactions and campaign outcomes.

#### 1.2.4. Reports Needed:

1. Which contact method has highest number of outcomes?
2. Which age group has taken highest number of housing loans?
3. Which age group has the highest number of personal loans?
4. Which educational backgrounds have the highest average bank balance?
5. Which age group has the more credit default percentage?
6. Which job-category is more present in the dataset?
7. Does the campaign has more success or failure?
8. Which loan holders are in credit default?
9. Does Average and total count of educational backgrounds are same?
10. Which campaign has more Success(previous or present)?

### **1.3. Project Risks:**

#### **Data Privacy and Security Risks:**

Risk: Unauthorized access or data breaches could compromise sensitive customer information, leading to legal and reputational consequences.

Mitigation: Implement robust data encryption, access controls, and regular security audits to safeguard customer data. Adhere to data protection regulations such as GDPR or CCPA.

#### **Data Quality Risks:**

Risk: Inaccurate or incomplete data could lead to biased predictions and unreliable insights, impacting the effectiveness of marketing campaigns.

Mitigation: Establish data validation and cleansing processes to ensure data accuracy and completeness. Regularly monitor data quality metrics and address any issues promptly.

#### **Resource Constraints:**

Risk: Insufficient resources, such as skilled personnel, technology infrastructure, or budgetary constraints, may hinder project execution and effectiveness.

Mitigation: Allocate adequate resources and invest in training for personnel involved in data analysis, marketing, and compliance. Prioritize key initiatives based on available resources and strategic goals.

#### **Technology Risks:**

Risk: Technical failures, system downtime, or compatibility issues with data management tools and platforms could disrupt project operations.

Mitigation: Implement robust IT infrastructure, backup systems, and contingency plans to mitigate technology-related risks. Regularly update and maintain software systems to ensure reliability and performance.

By identifying and addressing these risks proactively, the project can enhance its resilience, optimize outcomes, and achieve its objectives effectively.

## 2. Design Specifications:

### 2.1. Relational concepts:

Customer Entity: Contains information about customers.

Attributes: age, job, marital, education, default, balance, housing, loan

Contact Entity: Information about marketing contacts with customers.

Attributes: contact, day, month, duration, campaign, pdays, previous

Outcome Entity: The outcome of marketing contacts.

Attributes: poutcome, y

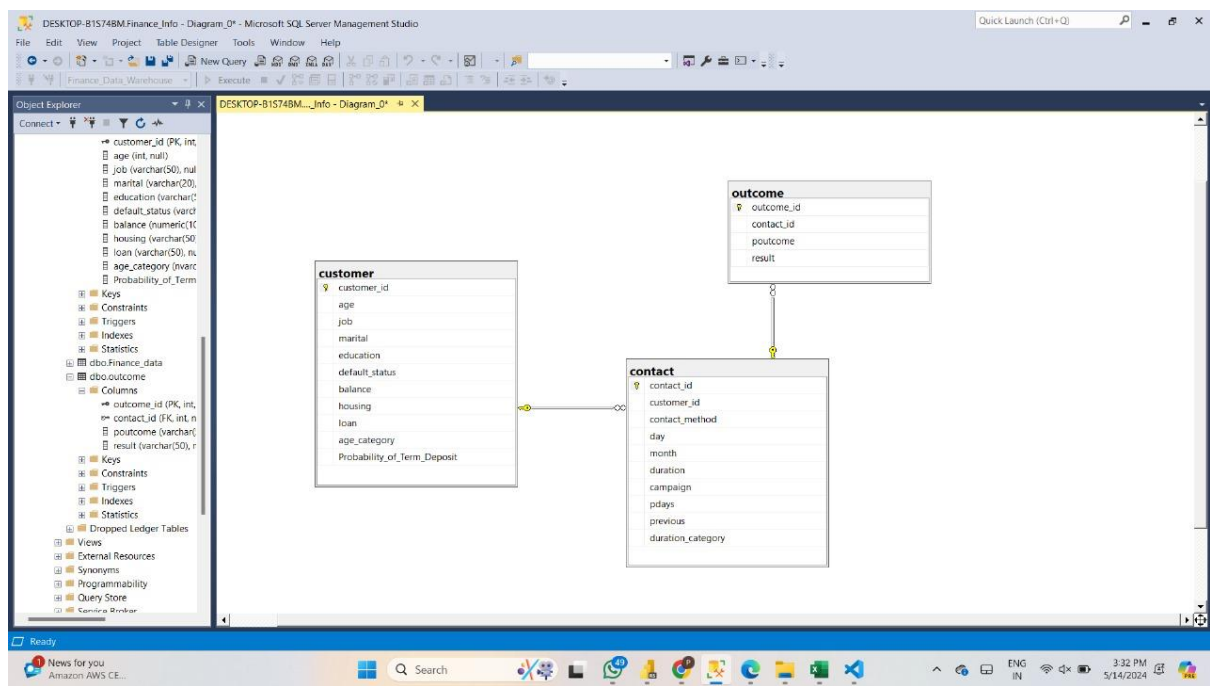
### Relationships:

Customer - Contact: A customer can have multiple contacts, but each contact is associated with one customer (One-to-Many relationship).

Contact - Outcome: Each contact can lead to an outcome (One-to-One or One-to-Many relationship).

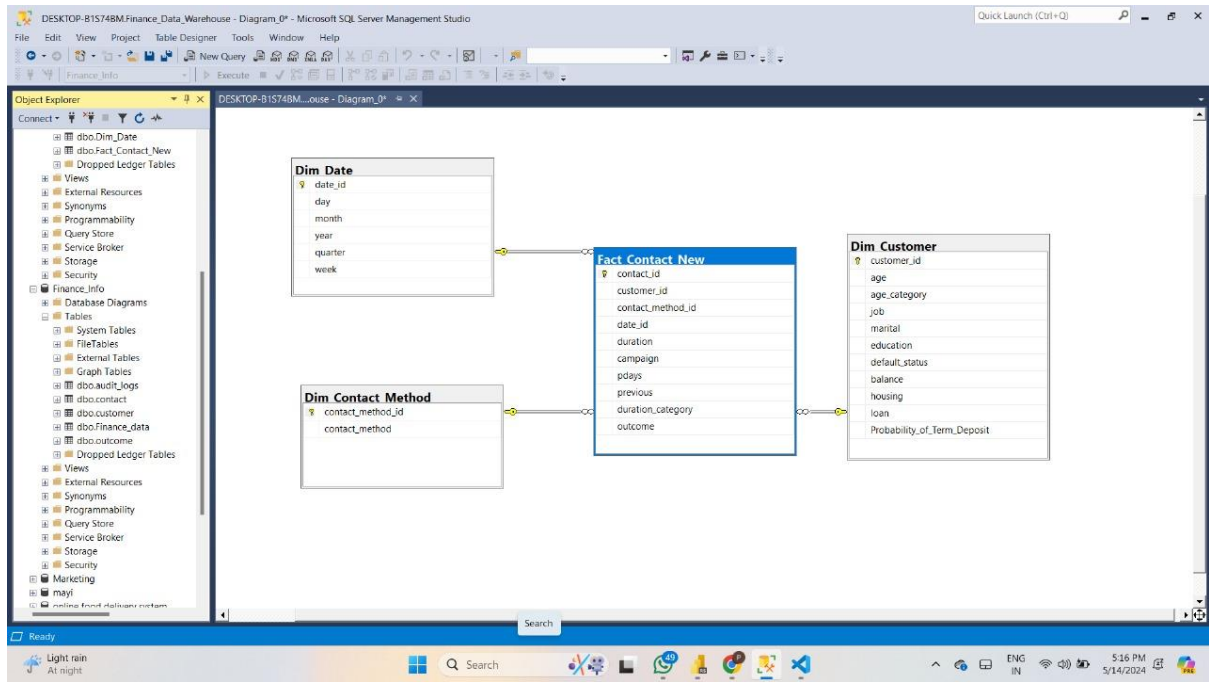
### 2.1.1. Technology to be used and ER Diagram:

Database Management System (DBMS): we used a relational DBMS, SQL Server management studio to create an ER diagram.



## 2.2. Big data concepts:

### 2.2.1. Data-Warehouse Structure:

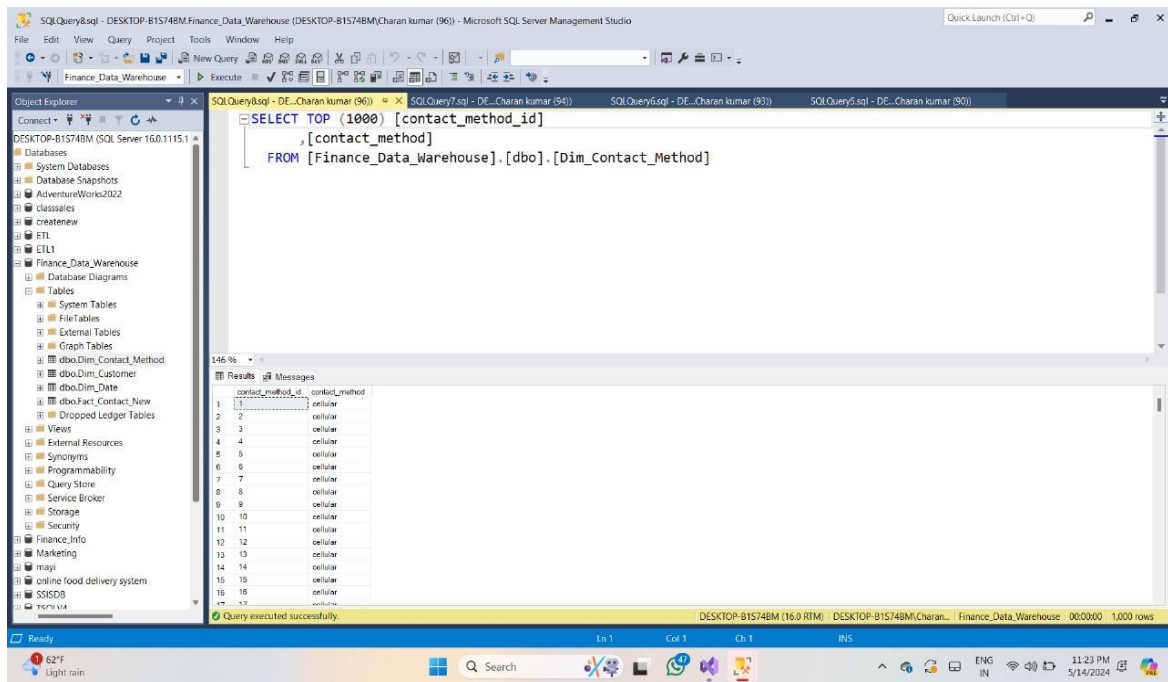


There are 3 dimensions tables and 1 Fact table.

Dimension Tables:

1. Dim\_Date:  
Primary Key: Date\_id  
Other attributes: Day, Month, year, quarter, Week
2. Dim\_Contact\_method:  
Primary Key: Contact\_method\_id  
Other attributes: Contact\_method

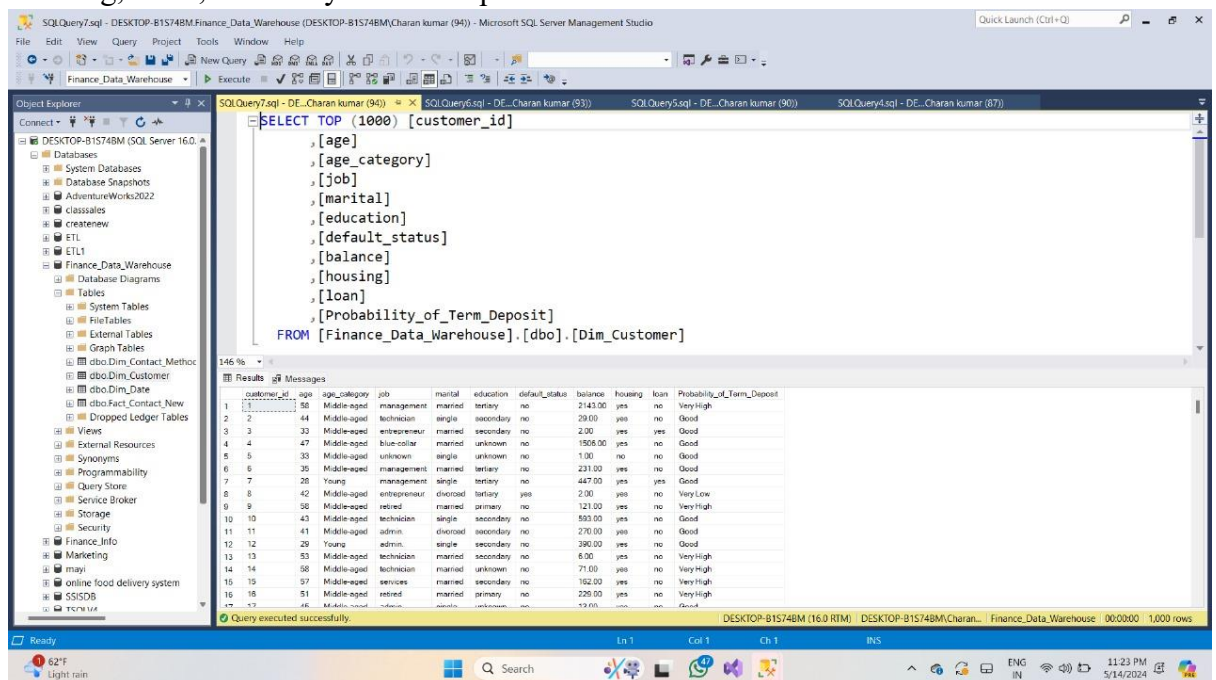




### 3. Dim\_Customer:

Primary Key: Customer\_id

Other attributes: age, age\_category, job, marital, education, Default Status, Balance, Housing, Loan, Probability of term deposit



Fact table:

Comntact\_id, Customer\_id, Contact\_method\_Id, date\_id, duration, campaign, pdays, previous, Duration\_Category, Outcome.

### 2.2.2 ETL Development:

In ETL, we have derived three columns from already existing columns.

i) Age Category:

From age column, age category based on following criteria:

Below 30: young Age

Above 30 and below 60: Middle-Age

Above 60: Elderly

ii) Probability of Term Deposit:

We have calculated this value based on age and credit default.

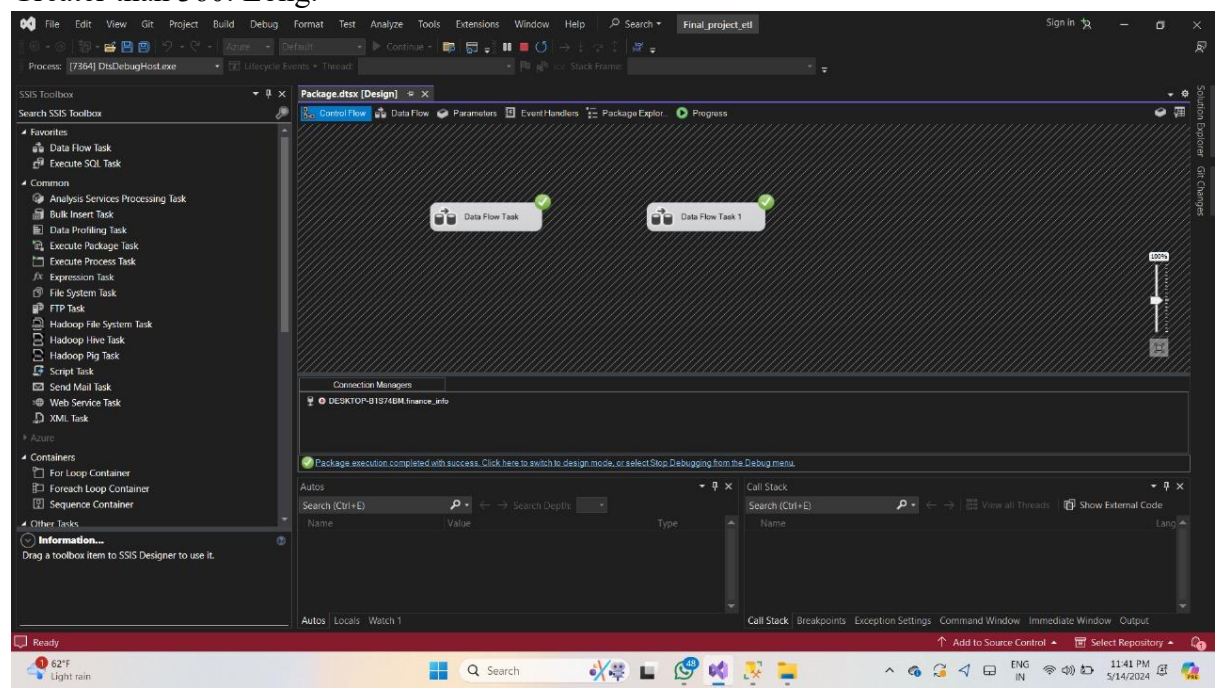
iii) Duration Category:

This has been categorized based on duration column.

If the call duration is less than 180: small.

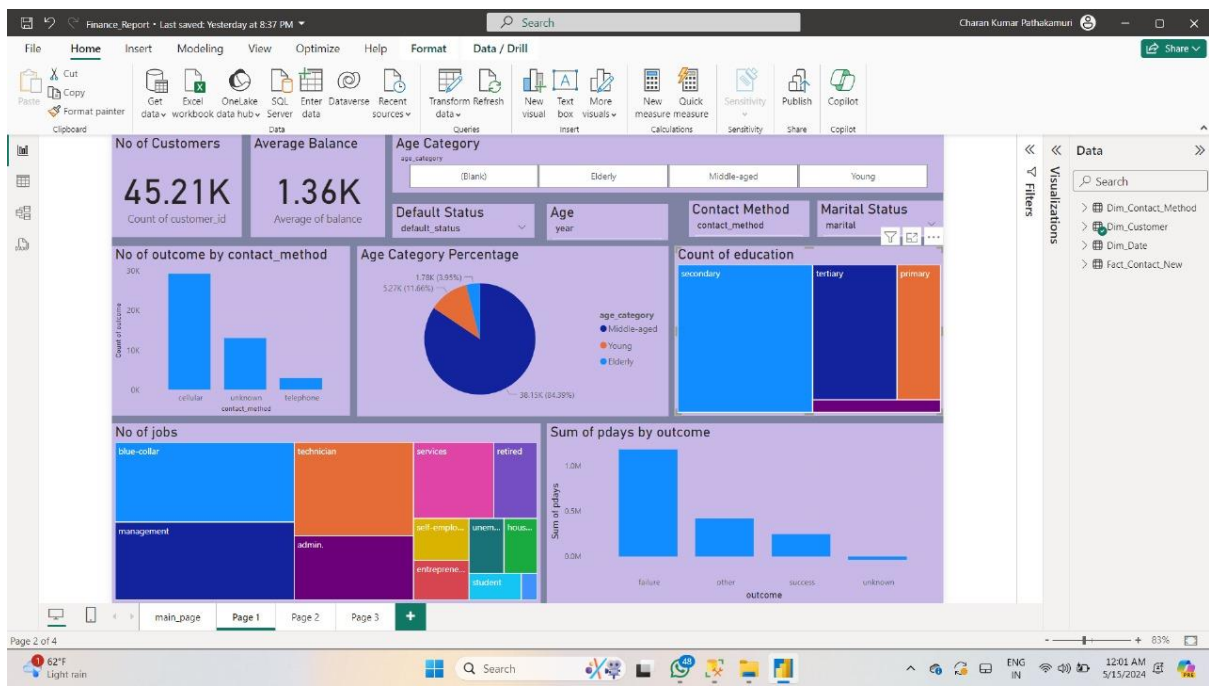
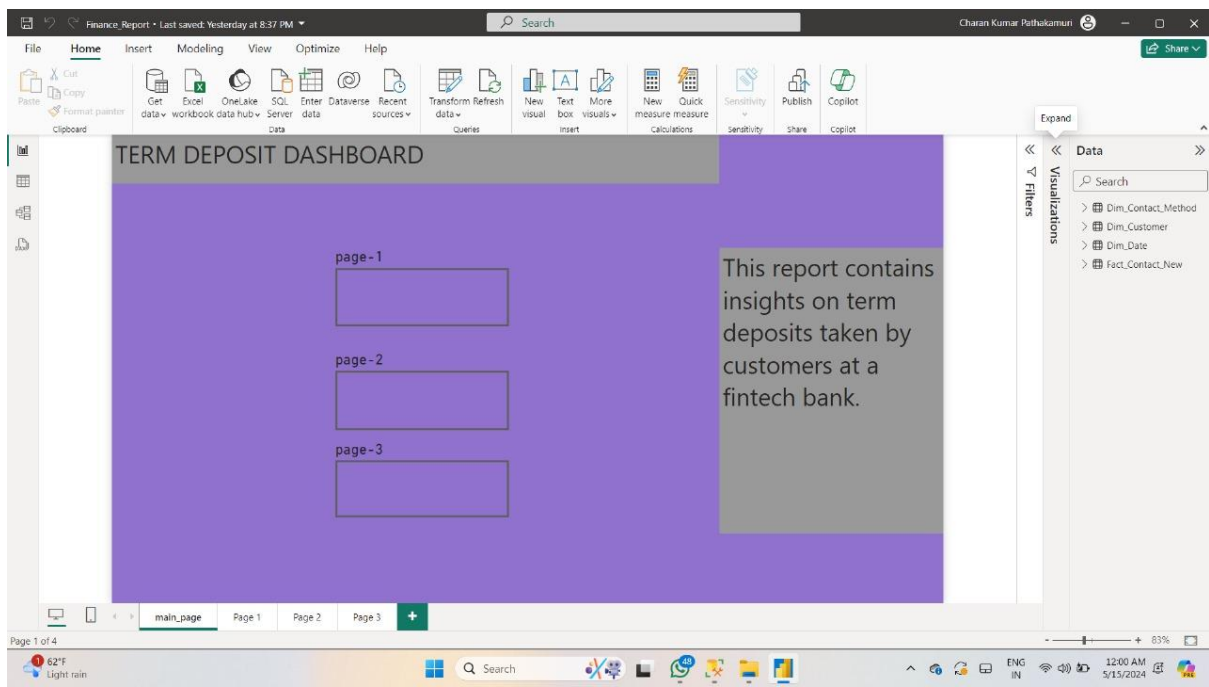
If the call higher than 180 and lesser than 360: Medium

Greater than 360: Long.

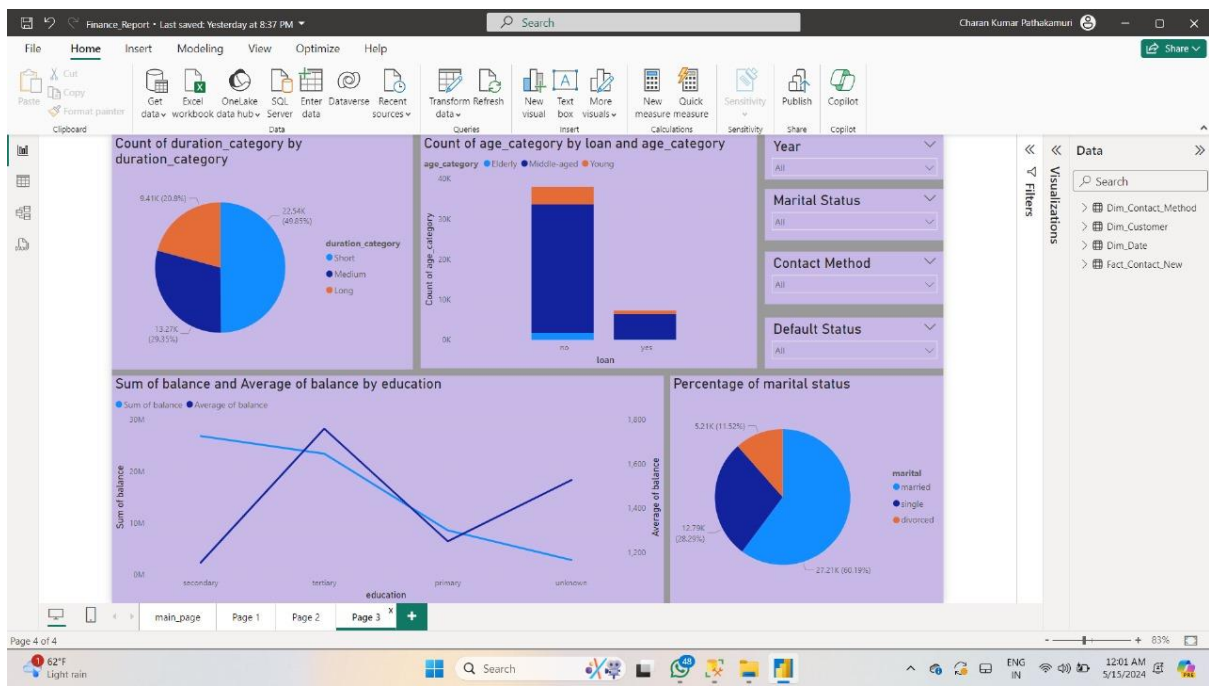
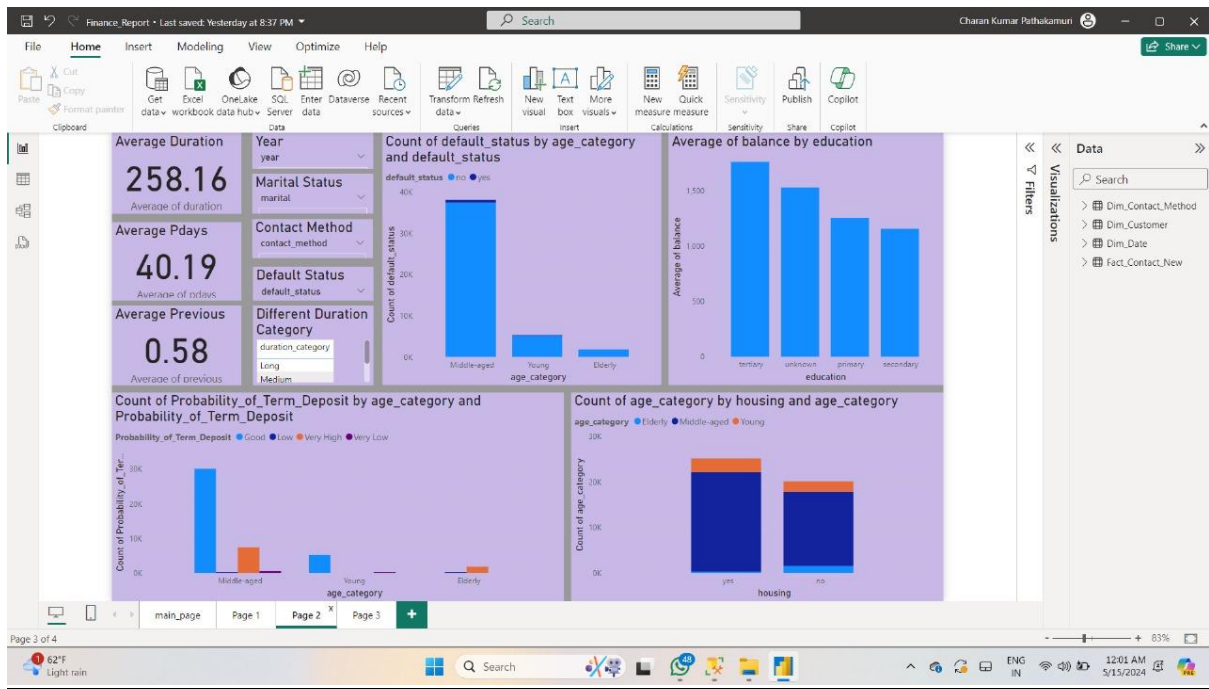




## 2.2.3 Power BI:







## 2.3 Cost Estimations:

### 2.3.1 Software Costs:

Microsoft SQL Server:

- Licensing: SQL Server enterprise edition will cost \$15,123
- ODBC Driver: Which is free

Python Libraries:

- a) Pandas: Free and Open source
- b) SQL Alchemy: free and open source.
- c) pyodbc: free and open-source.
- d) Openpyxl: free and open source.

Visual Studio 2022: Enterprise Edition : \$1,700

SSIS Package: \$3,400

Microsoft Fabric Power BI : \$5000

Visual Studio Code : Free and Open Source

Total software Costs: \$25,223

### 2.3.2 Hardware Costs:

On-premises:

Server Chassis: \$1,000

CPU : \$2000

Ram: 200

Storage: \$200

Network Interface: \$500

Cooling Systems: \$250

RAID Controllers: \$200

Backup Solutions: \$500

Total: 4,850

### 2.3.3 Staff Resources:

Developer: 10,000

DBA: 15,000

Data Engineer: \$5000

Power BI developer: \$3000

Total: \$33,000

Miscellaneous:

Training and Documentation: \$1000

Total cost: \$63,850

#### **2.4 Technologies Used:**

- ❖ Python
- ❖ Pandas.
- ❖ SQL Alchemy Engine
- ❖ Openpyxl
- ❖ SQL Server
- ❖ Pyodbc for data Migration
- ❖ SSIS used for ETL.
- ❖ Power Bi
- ❖ Visual Studio 2022

#### **2.5.Teamates Contributions:**

Lokesh Reddy: Data Base Designing, Modelling and normalizing

Upendra: ETL, Data Quality

Charan: Power BI, Data Exploring

Sashank: Data Cleaning, Loading