# Retail Banking Data- CRM

Madhav Appaneni
*University at Buffalo*
Buffalo, New York, USA
madhavap@buffalo.edu

## I. INTRODUCTION

The dataset consists of retail banking demo data. Banking is one such sector where databases are most useful since the transactions happen at a faster pace and all the updates need to happen in real time. Since the transactions happen in real time, the data that gets generated will be huge which makes it difficult to store and process in excel sheets. On the other hand, when a customer performs transactions on his/her account, it would not be ideal for both customer and financial institution if they get erroneous results. Further, databases support concurrency. In this case, customers, bank officials and customer support agents can seamlessly ingest and access data. Databases also help in keeping the data secure in an organization with multiple access and role support. For example, the customer support agents must not be allowed to edit information about a branch or a bank. Using a spreadsheet would not be an ideal solution for this use. In this case, we have data coming from many input sources, if we use a spreadsheet, we will have to duplicate parts of records for better readability. But with databases, we get to join multiple relations and access data with better flexibility To tackle all these issues, we aim to design a database which handles all CRUD operations in real time.

## II. TARGET AUDIENCE

This database is used by three sets of people. 1. Bank officials: The bank officials use this database to store and retrieve information about the customers – their accounts, balances, addresses and other information 2. Customers: The customers use this database indirectly via a web interface to check their balance, transactions, addresses, card details and their customer support incidents. 3. Customer Support Agents: The customer support agents will use this database via a web interface or via CRM software to address customer queries about their accounts and raise any incidents or tickets if any customer faces any issues.

## III. DATABASE ADMINISTRATOR AND USERS

The main users of the CRM Retail Banking dashboard will be the Customer support agents. When customers call the customer care center, they take the complaint from the customers. This dashboard will help to maintain and keep track of the complaints using complaint IDs of different customers including past records. The database administrators will be the IT personnel of the bank. They would be responsible for maintaining the integrity, optimizing the run times, and ensuring the sanity of the database.

## IV. DATABASE STRUCTURE

### A. Intial ER Diagram

Entity relationship diagram is represented in Fig. 1 It shows the relation between different tables in the database.
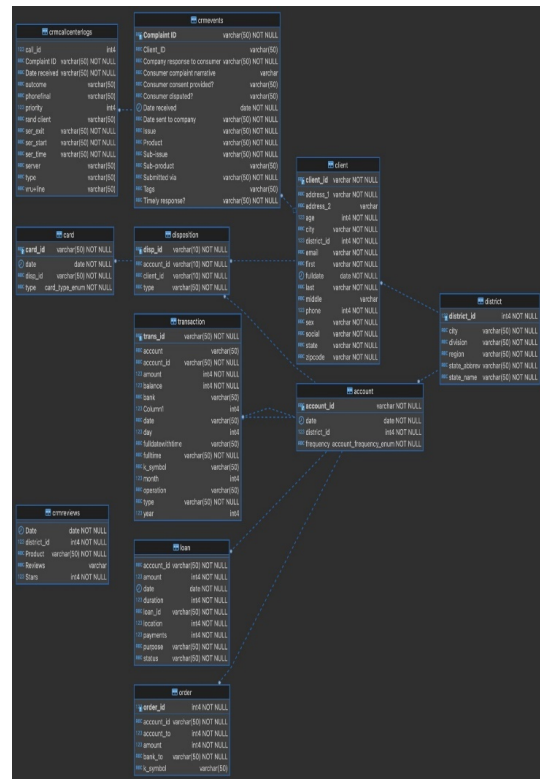


Fig 1: ER Diagram

### B. Relations

- we have 11 relations in the database.The details of the relations are described as below
- **Account:**

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| account_id | varchar | Account id of the account | NA | NOT NULL |
| district_id | integer | District if of the customer | NA | NOT NULL |
| frequency | enum | Account issuance frequency | NA | NOT NULL |
| date | date | Account opening date | NA | NOT NULL |

Fig 2: Schema of Account table

- **Card:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| card_id | varchar | Card if of the card | NA | NOT NULL |
| disp_id | varchar | Disposition Id of the card | NA | NOT NULL |
| type | enum | Type of card | NA | NOT NULL |
| date | date | Card issuance date | NA | NOT NULL |

Fig 3: Schema of Card Table

- **Client:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| sex | varchar | Gender of the client | NA | Not null |
| fulldate | date | Date on which client registered | NA | Not null |
| age | integer | Age of the client | NA | Not null |
| social | varchar | Social security number of the client | NA | Not null |
| first | varchar | First name of the client | NA | Not null |
| middle | varchar | Middle name of the client | NA | Nullable |
| last | varchar | Last name of the client | NA | Not null |
| phone | integer | Phone number of the client | NA | Not null |
| email | varchar | Email address of the client | NA | Not null |
| address_1 | varchar | Address Line 1 of the client | NA | Not null |
| address_2 | varchar | Address Line 2 of the client | NA | Nullable |
| city | varchar | City of the beneficiary | NA | Not null |
| state | varchar | State of the beneficiary | NA | Not null |
| zipcode | varchar | Zipcode of the city | NA | Not null |
| district_id | integer | District Id of where the beneficiary lives | NA | Not null |
| client_id | varchar | Client ID, unique ID of a client | NA | Not null |

Fig 4: Schema of client Table

- **Transaction:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| trans_id | varchar | Transaction id | NA | NOT NULL |
| Account | varchar | Account number | NA | NOT NULL |
| Amount | Integer | Transaction amount | NA | NOT NULL |
| Balance | Integer | Balance amount before transaction | NA | NOT NULL |
| Bank | varchar | Name of the bank | NA | NULL |
| Date | Date | Transaction date time | NA | NOT NULL |
| type | varchar | Credit or Debit | NA | NOT NULL |
| Operation | varchar | Operation that initiated the transaction | NA | NOT NULL |

Fig 5: Schema of transaction table

- **order:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| order_id | integer | Uniquely identifies each order | N/A | NOT NULL |
| account_id | varchar(50) | identifies each bank account | N/A | NOT NULL |
| bank_to | varchar(50) | Identifies the bank to which money is transferred | N/A | NOT NULL |
| account_to | integer | Identifies the account to which money is transferred | N/A | NOT NULL |
| amount | integer | Identifies the amount of the order | N/A | NOT NULL |
| K_symbol | varchar(50) | Identifies the purpose of the payment | N/A | NULL |

Fig 6: Schema of order Table

- **Loan:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| loan_id | varchar(50) | Uniquely identifies each loan | N/A | NOT NULL |
| account_id | varchar(50) | identifies each bank account | N/A | NOT NULL |
| amount | integer | Identifies the loan amount | N/A | NOT NULL |
| duration | integer | Identifies the loan tenure | N/A | NOT NULL |
| payments | integer | Identifies the loan payments | N/A | NOT NULL |
| status | varchar(50) | Identifies the loan status | N/A | NOT NULL |
| date | varchar(50) | Identifies the loan starting date | N/A | NOT NULL |
| location | integer | Identifies the location number | N/A | NOT NULL |
| purpose | varchar(50) | Identifies the loan purpose | N/A | NOT NULL |

Fig 7: Schema of loan table

- **Call Center Logs:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| Date received | varchar(50) | Date when the complaint was received | NA | NOT NULL |
| Complaint Id | varchar(50) | Unique Identifier for the complaint | NA | NOT NULL |
| Rand client | varchar(50) | | NA | NULL |
| Phonefinal | varchar(50) | | NA | NULL |
| Vru_line | varchar(50) | | NA | NULL |
| Call_id | varchar(50) | Unique Identifier of the call log | NA | NULL |
| Priority | varchar(50) | Priority of the complaint | NA | NULL |
| Type | varchar(50) | Type of complaint | NA | NULL |
| Outcome | varchar(50) | Outcome resolution of the complaint | NA | NULL |
| Server | varchar(50) | Executive who attended the customer | NA | NULL |
| Ser_start | varchar(50) | When did the service start | NA | NOT NULL |
| Ser_exit | varchar(50) | When did the service end | NA | NOT NULL |
| Ser_time | varchar(50) | Total service time elapsed | NA | NOT NULL |

Fig 8: Schema of Call center logs Table

- **CRM Events:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| Date Received | Date | Date of the Consumer Relation Event | NA | NOT NULL |
| Product | Varchar(50) | Product involved in the event | NA | NOT NULL |
| Sub-product | Varchar(50) | Sub Product category | NA | NULL |
| Issue | Varchar(50) | Issue in Product, main category | NA | NOT NULL |
| Sub-Issue | Varchar(50) | Issue in Product, sub category | NA | NOT NULL |
| Consumer Complaint Narrative | Varchar | Narrative of issue as per customer | NA | NULL |
| Tags | Varchar(50) | Tags that are relevant to the nature of complaint | NA | NULL |
| Consumer Consent Provided | Varchar(50) | If consumer provided consent to call them back | NA | NOT NULL |
| Submitted via | Varchar(50) | Medium via which the event was registered | NA | NOT NULL |
| Date sent to company | Varchar(50) | Date the event was sent to the company | NA | NOT NULL |
| Company response to customer | Varchar(50) | Company resolution of the complaint | NA | NOT NULL |
| Timely response | Varchar(50) | Yes/No | NA | NOT NULL |
| Consumer disputed | Varchar(50) | If the consumer disputed against the resolution | NA | NULL |
| Complaint Id | Varchar(50) | Unique Id of the complaint | NA | NOT NULL |
| Client_ID | Varchar(50) | Id of the client who registered the complaint | NA | NULL |

Fig 9: Schema of CRM Events Table

- **CRM Reviews:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| Date | date | Date of the review | NA | NULL |
| Stars | integer | Number of stars rating the service | NA | NOT NULL |
| Reviews | varchar | The actual Review | NA | NULL |
| Product | varchar(50) | Product against which the review is made | NA | NOT NULL |
| district_id | integer | District ID at which the product was reviewed by customer | NA | NOT NULL |

Fig 11: Schema of CRM Reviews Table

- **Disposition:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| disp_id | varchar(10) | Uniquely identifies each disposition | N/A | NOT NULL |
| client_id | varchar(10) | Identifies each client | N/A | NOT NULL |
| account_id | varchar(10) | Identifies each account | N/A | NOT NULL |
| type | enum | Identifies whether its user or owner | N/A | NOT NULL |

Fig 12: Schema of Disposition Table

- **District:**

| Name | Data type | Description | Default value | Null constraints |
|------|-----------|-------------|---------------|------------------|
| district_id | integer | Uniquely identifies each district | N/A | NOT NULL |
| city | varchar(50) | Identifies the city of the client to corresponding district | N/A | NOT NULL |
| state_name | varchar(50) | Identifies the state name | N/A | NOT NULL |
| state_abbrev | varchar(50) | Identifies State's name by short abbreviation | N/A | NOT NULL |
| region | varchar(50) | Identifies the region to the corresponding district | N/A | NOT NULL |
| division | varchar(50) | Identifies the division to the corresponding district | N/A | NOT NULL |

Fig 13: Schema of District Table

### C. Database Constraints

- *Primary Keys*:
  - Attribute account_id – uniquely identifies each account and is a primary key for account table
  - Attribute card_id – uniquely identifies each card and is the primary key for card table.
  - Attribute client_id – uniquely identifies each client and is a primary key for client table.
  - Attribute disp_id – uniquely identifies each disposition and is a primary key for disposition table
  - Attribute district_id – uniquely identifies each district and is a primary key for district table
  - Attribute loan_id – uniquely identifies each loan and is a primary key for loan table
  - Attribute order_id – uniquely identifies each order and is a primary key for order table
  - Attribute trans_id – uniquely identifies each transaction and is a primary key for transaction table
  - Attribute call_id – uniquely identifies each call and is a primary key for callcenterlogs table
  - Attribute complaint_id – uniquely identifies each complaint and is a primary key for crmevents table
  - Attribute review_id – uniquely identifies each review and is a primary key for crmreviews table
- *Foreign Keys*:

  - Attribute district_id – in the account table references attribute distict_id in the district table.
  - Attribute disp_id – in the card table references attribute disp_id in the disposition table.
  - Attribute district_id – in the client table references attribute district_id in the district table.
  - Attribute client_id and account_id in the disposition table references attribute client_id and account_id in the client and account tables respectively.
  - Attribute account_id – in the loan table references attribute account_id in the account table.
  - Attribute account_id – in the order table references attribute account_id in the account table.
  - Attribute account_id – in the transaction table references attribute account_id in the account table.
  - Attribute complaint_id – in the callcenterlogs table references attribute complaint_id in the crmreviews table.
  - Attribute client_id in the crmevents table references attribute client_id in the client table
  - Attribute district_id – in the crmreviews table references attribute district_id in the crmreviews table.
- *NOT NULL* and *DEFAULT*: Banking data is very prone to the information supplied by the customer and other banking officials. For this reason, no information in the tables is considered a default. This is in accordance with our use case.

  - Account - all attributes must be supplied and are not null.
  - Card - all attributes must be supplied and are not null.
  - Client - all attributes must be supplied and are not null.
  - Transaction - all attributes must be supplied and are not null.
  - order - all attributes except k_symbol must be supplied and are not null. K_symbol is the purpose of the payment and can be null.
  - loan - all attributes must be supplied and are not null.
  - call center logs - rand_client, phonefinal, vru_line, call_id and priority are null, rest have to be supplied and not null.
  - crm events - consumer complaint narrative, tags, consumerdisputed, and call_id are null, rest have to be supplied and not null.
  - crm reviews - except date and review all fields are not null.
  - disposition - all attributes must be supplied and are not null.
  - district - all attributes must be supplied and are not null.
  - *On Delete Cascades*
    * On transaction, deleting account_id deletes account(account_id).
    * On loan, deleting account_id deletes

account(account_id).
- ∗ On CRMEvents, deleting Client_ID deletes client(client_id).
- ∗ On CRMCallCenterLogs, deleting Complaint ID deletes account( Complaint ID).
- ∗ On order, deleting account_id deletes account(account_id).
- ∗ On account, deleting district_id deletes district(district_id).
- ∗ On card, deleting disp_id deletes disposition(disp_id).
- ∗ On disposition, deleting client_id deletes client(client_id).
- ∗ On disposition, deleting client_id deletes clientaccount(client_id).
- ∗ On client, deleting district_id deletes district(district_id).
- ∗ On CRMReviews, deleting district_id deletes district(district_id).
- ∗ On district, deleting state_name deletes account(state_name).

## V. Data Migration

The dataset contains 11 CSV files. We carefully examined each CSV file and removed the duplicate fields present in each CSV. We devised efficient data types, default values, null constraints, and on-delete constraints for each column in a relation and created the DDL scripts. We further identified the relationships between each relation and included these relations in our create.sql script. We then imported the raw CSV into the database client, trimmed the table according to our analysis using SQL queries, and achieved the resulting schema for each relation. We exported the SQL insert statement files for each relation as a separate file. To automate the process of creating the database schema and inserting the data, we made use of a shell script that internally uses the PSQL command line tool to run each .sql file over the selected database.

## VI. BCNF PROVEMENT

- The relation account has the following functional dependencies account_id district_id, frequency, date. Here account_id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation account is in BCNF.
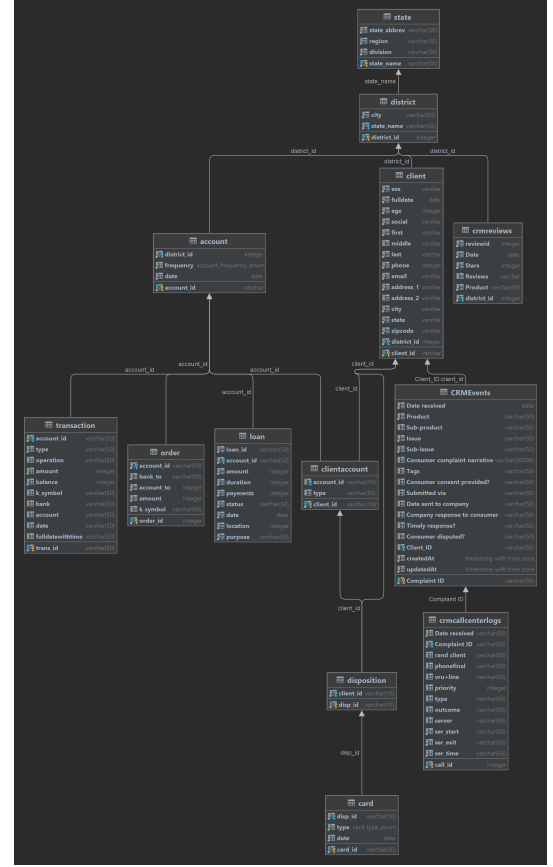- The relation card has the following functional dependencies card_id ⟶ disposition_id, type, date Here card_id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation card is in BCNF

- The relation client has the following functional dependencies client_id ⟶ sex, fulldate, age, social, first, middle, last, phone, email, address_1, address_2, zipcode, district_id. Here client_id is determining all the other attributes and it is the super key.Also there is no non prime attribute that is determining other non-prime attribute(3NF) and the above functional dependency is non-trivial. Therefore,the relations client is in BCNF.
  When we tried to split the geolocation part of the client table into a separate table, we ran into issues like one zipcode being present in multiple districts and states. We found that a zip code can span multiple cities, districts, and states https://gis.stackexchange.com/questions/53918/determining-which-us-zipcodes-map-to-more-than-one-state-or-more-than-one-city. Further, the dataset documentation mentions that the zip codes are randomly generated with values between 4000 and 6000. So, we proved that the functional dependencies we came up with were invalid. Moreover, since this issue is technically valid, we left the table as is without splitting.
- The relation disposition has the following functional dependencies disp_id ⟶ client_id, account_id, type. Here disp_id is the super key, since it determines all the other attributes.But there are non-prime attributes determining another non-prime attributes (3NF fails) i.e client_id ⟶ account_id,type and disp_id. Clearly, the FD violates the BCNF. So we decompose the relation into two tables disposition(disp_id,client_id) and client_account_map(client_id,account_id,type). Both the new relations are now in BCNF.
- The relation district has following functional dependencies district_id ⟶ city, state_name, state_abbrev, region, division. state_name ⟶ state_abbrev, region, division; state_abbrev ⟶ state_name, region, division. Here in the first FD the district_id is the super key, since it determines all the other attributes. Also, the above functional dependency is non-trivial. It satisfies the BCNF conditions. But in the second and third FD's, the state_name and state_abbrev are not super keys. Therefore, the relation disposition is not in BCNF. So, it is decomposed into two relations district (district_id, city, state_name) and state(state_abbrev, state_name, region, division). Both the new relations are now in BCNF.
- The relation loan has the following functional dependencies loan_id ⟶ account_id, amount, duration, payments, status, date, location, purpose. Here loan_id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation loan is in BCNF.

- The relation order has the following functional dependencies order_id ⟶ account_id, bank_to, account_to, amount, k_symbol. Here order_id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation order is in BCNF
- The relation transaction has the following functional dependencies trans_id ⟶ account, amount, balance, bank string, date, type, operation. Here trans_id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation transaction is in BCNF
- The relation callcenterlogs has the following functional dependencies call_id ⟶ date_received, complaint_id, rand_client, phonefinal, vru_line, priority, type, outcome, server, ser_start, ser_exit, ser_time. Here call_id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation callcenterlogs is in BCNF.
- The relation crmevents has the following functional dependencies complaint_id ⟶ date received, product, sub-product, issue, sub-issue, consumer complaint narrative, tags, consumer consent provided, submitted via, date sent to company, company response to consumer, timely response, consumer disputed, client_id. Here complaint id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation crmevents is in BCNF.
- The relation crmreviews has the following functional dependencies review_id ⟶ date, stars, reviews, product, district_id. Here review_id is the super key, since it determines all the other attributes. There is no other non-prime attribute determining any other non-prime attribute(3NF). Also, the above functional dependency is non-trivial. Therefore, the relation crmreviews is in BCNF.

## VII. TRANSFORMATION FROM INITIAL SCHEMA TO THE FINAL SCHEMA

### A. Final ER Diagram

After decomposing district and disposition tables using BCNF, the final Entity relationship diagram is represented in Fig. 14 It shows the relation between different tables in the database.



Fig 14: Final ER Diagram

### B. Relations after BCNF Provement

- We now have 13 relations in the database after BCNF decomposition.The details of the relations are described below:

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| account_id | varchar | Account id of the account | NA | NOT NULL |
| district_id | integer | District if of the customer | NA | NOT NULL |
| frequency | enum | Account issuance frequency | NA | NOT NULL |
| date | date | Account opening date | NA | NOT NULL |

Fig 15: Schema of Account table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| card_id | varchar | Card if of the card | NA | NOT NULL |
| disp_id | varchar | Disposition Id of the card | NA | NOT NULL |
| type | enum | Type of card | NA | NOT NULL |
| date | date | Card issuance date | NA | NOT NULL |

Fig 16: Schema of Card Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| sex | varchar | Gender of the client | NA | Not null |
| fulldate | date | Date on which client registered | NA | Not null |
| age | integer | Age of the client | NA | Not null |
| social | varchar | Social security number of the client | NA | Not null |
| first | varchar | First name of the client | NA | Not null |
| middle | varchar | Middle name of the client | NA | Nullable |
| last | varchar | Last name of the client | NA | Not null |
| phone | integer | Phone number of the client | NA | Not null |
| email | varchar | Email address of the client | NA | Not null |
| address_1 | varchar | Address Line 1 of the client | NA | Not null |
| address_2 | varchar | Address Line 2 of the client | NA | Nullable |
| city | varchar | City of the beneficiary | NA | Not null |
| state | varchar | State of the beneficiary | NA | Not null |
| zipcode | varchar | Zipcode of the city | NA | Not null |
| district_id | integer | District Id of where the beneficiary lives | NA | Not null |
| client_id | varchar | Client ID, unique ID of a client | NA | Not null |

Fig 17: Schema of client Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| trans_id | varchar | Transaction id | NA | NOT NULL |
| Account | varchar | Account number | NA | NOT NULL |
| Amount | Integer | Transaction amount | NA | NOT NULL |
| Balance | Integer | Balance amount before transaction | NA | NOT NULL |
| Bank | varchar | Name of the bank | NA | NULL |
| Date | Date | Transaction date time | NA | NOT NULL |
| type | varchar | Credit or Debit | NA | NOT NULL |
| Operation | varchar | Operation that initiated the transaction | NA | NOT NULL |

Fig 18: Schema of transaction table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| order_id | integer | Uniquely identifies each order | N/A | NOT NULL |
| account_id | varchar(50) | identifies each bank account | N/A | NOT NULL |
| bank_to | varchar(50) | Identifies the bank to which money is transferred | N/A | NOT NULL |
| account_to | integer | Identifies the account to which money is transferred | N/A | NOT NULL |
| amount | integer | Identifies the amount of the order | N/A | NOT NULL |
| K_symbol | varchar(50) | Identifies the purpose of the payment | N/A | NULL |

Fig 19: Schema of order Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| loan_id | varchar(50) | Uniquely identifies each loan | N/A | NOT NULL |
| account_id | varchar(50) | identifies each bank account | N/A | NOT NULL |
| amount | integer | Identifies the loan amount | N/A | NOT NULL |
| duration | integer | Identifies the loan tenure | N/A | NOT NULL |
| payments | integer | Identifies the loan payments | N/A | NOT NULL |
| status | varchar(50) | Identifies the loan status | N/A | NOT NULL |
| date | varchar(50) | Identifies the loan starting date | N/A | NOT NULL |
| location | integer | Identifies the location number | N/A | NOT NULL |
| purpose | varchar(50) | Identifies the loan purpose | N/A | NOT NULL |

Fig 20: Schema of loan table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| Date received | varchar(50) | Date when the complaint was received | NA | NOT NULL |
| Complaint Id | varchar(50) | Unique Identifier for the complaint | NA | NOT NULL |
| Rand client | varchar(50) | | NA | NULL |
| Phonefinal | varchar(50) | | NA | NULL |
| Vru_line | varchar(50) | | NA | NULL |
| Call_id | varchar(50) | Unique Identifier of the call log | NA | NULL |
| Priority | varchar(50) | Priority of the complaint | NA | NULL |
| Type | varchar(50) | Type of complaint | NA | NULL |
| Outcome | varchar(50) | Outcome resolution of the complaint | NA | NULL |
| Server | varchar(50) | Executive who attended the customer | NA | NULL |
| Ser_start | varchar(50) | When did the service start | NA | NOT NULL |
| Ser_exit | varchar(50) | When did the service end | NA | NOT NULL |
| Ser_time | varchar(50) | Total service time elapsed | NA | NOT NULL |

Fig 21: Schema of Call center logs Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| Date Received | Date | Date of the Consumer Relation Event | NA | NOT NULL |
| Product | Varchar(50) | Product involved in the event | NA | NOT NULL |
| Sub-product | Varchar(50) | Sub Product category | NA | NULL |
| Issue | Varchar(50) | Issue in Product, main category | NA | NOT NULL |
| Sub-Issue | Varchar(50) | Issue in Product, sub category | NA | NOT NULL |
| Consumer Complaint Narrative | Varchar | Narrative of issue as per customer | NA | NULL |
| Tags | Varchar(50) | Tags that are relevant to the nature of complaint | NA | NULL |
| Consumer Consent Provided | Varchar(50) | If consumer provided consent to call them back | NA | NOT NULL |
| Submitted via | Varchar(50) | Medium via which the event was registered | NA | NOT NULL |
| Date sent to company | Varchar(50) | Date the event was sent to the company | NA | NOT NULL |
| Company response to customer | Varchar(50) | Company resolution of the complaint | NA | NOT NULL |
| Timely response | Varchar(50) | Yes/No | NA | NOT NULL |
| Consumer disputed | Varchar(50) | If the consumer disputed against the resolution | NA | NULL |
| Complaint Id | Varchar(50) | Unique Id of the complaint | NA | NOT NULL |
| Client_ID | Varchar(50) | Id of the client who registered the complaint | NA | NULL |

Fig 22: Schema of CRM Events Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| Date | date | Date of the review | NA | NULL |
| Stars | integer | Number of stars rating the service | NA | NOT NULL |
| Reviews | varchar | The actual Review | NA | NULL |
| Product | varchar(50) | Product against which the review is made | NA | NOT NULL |
| district_id | integer | District ID at which the product was reviewed by customer | NA | NOT NULL |

Fig 23: Schema of CRM Reviews Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| district_id | integer | Uniquely identifies each district | N/A | NOT NULL |
| city | varchar(50) | Identifies the city of the client to corresponding district | N/A | NOT NULL |
| state_name | varchar(50) | Identifies the state name | N/A | NOT NULL |

Fig 24: Schema of District Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| state_name | varchar(50) | Identifies the state name | N/A | NOT NULL |
| state_abbrev | varchar(50) | Identifies State's name by short abbreviation | N/A | NOT NULL |
| region | varchar(50) | Identifies the region to the corresponding district | N/A | NOT NULL |
| division | varchar(50) | Identifies the division to the corresponding district | N/A | NOT NULL |

Fig 25: Schema of State Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| disp_id | varchar(10) | Uniquely identifies each disposition | N/A | NOT NULL |
| client_id | varchar(10) | Identifies each client | N/A | NOT NULL |

Fig 26: Schema of Disposition Table

| Name | Data type | Description | Default value | Null constraints |
|---|---|---|---|---|
| client_id | varchar(10) | Identifies each client | N/A | NOT NULL |
| account_id | varchar(10) | Identifies each account | N/A | NOT NULL |
| type | enum | Identifies whether its user or owner | N/A | NOT NULL |

Fig 27: Schema of Client Account Table

*C. Database Constraints for tables formed after decomposition*

– *Primary Keys*:
  * Attribute state_name – uniquely identifies each state and is a primary key for state table
  * Attribute client_id and account_id – uniquely identifies each client account and is a composite primary key for clientaccount table
  * Attribute disp_id – uniquely identifies each disposition and is a primary key for the disposition table
  * Attribute disrict_id – uniquely identifies each district and is a primary key for the district table

– *Nullable and Default Constraints*:
  for all the tables there is no default constraint and all the attributes are non-nullable.

## VIII. PROBLEMS WHILE HANDLING LARGER DATA SETS

One of the problems that we have observed using larger data sets is taking longer execution times.one such scenario is when we want to query transaction table on date columns. Since the transaction table is huge, it is taking longer time to retrieve data. So, we have adopted indexing concept on the date column and have observed the execution time reduced drastically from 360.413 ms to 0.186 ms.



Fig 25: Before - Without Indexing



Fig 26: After - Indexing on date

## IX. BASIC QUERIES

1) **Find the average rating of each product of Eagle National Bank for each state**

```
select "Product", state_name, avg("Stars") from crmreviews c
join district d on c.district_id = d.district_id
group by "Product", state_name order by "Product", state_name;
```

2) **Find the distribution of clients in each state that has VISA signature Card.**

```sql
select c.state,count(c.client_id) from client c
join disposition ON disposition.client_id = c.client_id
join card c1 on c1.disp_id = disposition.disp_id
where c1."type"='VISA Signature' group by c.state;
```



3) **For a particular sub-product, let's say "Savings Account", find how the customer support performed in terms of how many queries by customers were closed/open per region.**
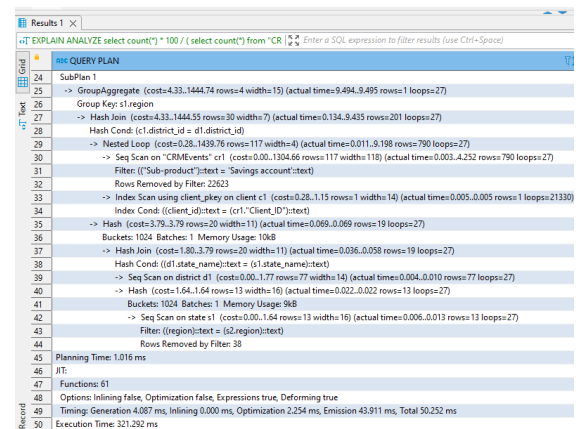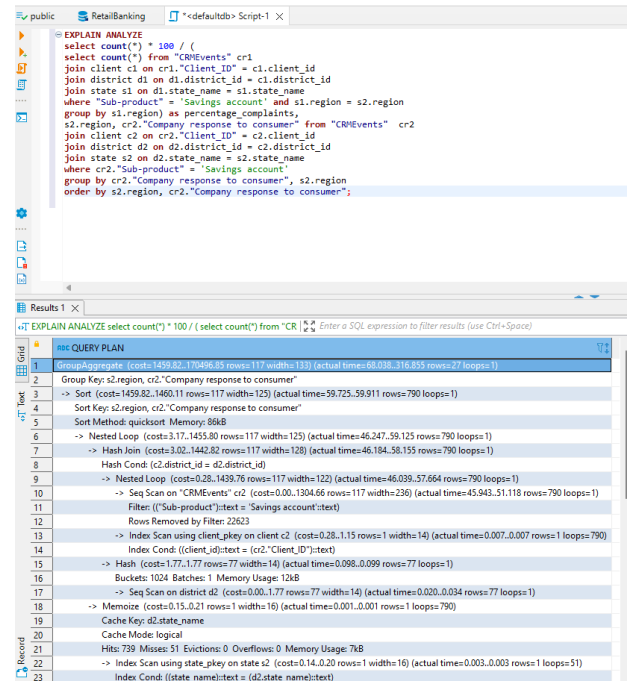
```sql
select count(*) * 100 / (
select count(*) from crmevents cr1
join client c1 on cr1."Client_ID" = c1.client_id
join district d1 on d1.district_id = c1.district_id
where "Sub-product" = 'Savings account' and d1.region = d2.region
group by d1.region) as percentage_complaints,
d2.region, cr2."Company response to consumer" from crmevents cr2
join client c2 on cr2."Client_ID" = c2.client_id
join district d2 on d2.district_id = c2.district_id
where cr2."Sub-product" = 'Savings account'
group by d2.region, cr2."Company response to consumer"
order by d2.region, cr2."Company response to consumer";
```



4) **Find the account details of the person that has taken the highest loan amount**

```sql
select a.account_id,c.first,c.last,max(l.amount) as HighestLoanAmount from
client c
join disposition d on d.client_id=c.client_id
join account a on a.account_id = d.account_id
join loan l on l.account_id = a.account_id
group by a.account_id,c.first,c.last
having max(l.amount) in (select max(amount) from loan l);
```



5) **In any particular region, how much money was transferred from Eagle National Bank to other banks?**

```sql
select count(amount), bank_to, region from "order"
natural join account
natural join district
group by bank_to, region
order by region, bank_to;
```

```sql
select count(amount), bank_to, region from "order" natural join account natural join district
    group by bank_to, region
    order by region, bank_to;
```

**order(+) 1**

select count(amount), bank_to, region from "order" nat... | Enter a SQL expression to filter results (use Ctrl+Space)

| | count | bank_to | region |
|---|---|---|---|
| 1 | 43 | AB | Midwest |
| 2 | 41 | CD | Midwest |
| 3 | 41 | EF | Midwest |
| 4 | 49 | GH | Midwest |
| 5 | 47 | IJ | Midwest |
| 6 | 51 | KL | Midwest |
| 7 | 45 | MN | Midwest |
| 8 | 52 | OP | Midwest |
| 9 | 48 | QR | Midwest |
| 10 | 47 | ST | Midwest |
| 11 | 52 | UV | Midwest |
| 12 | 49 | WX | Midwest |
| 13 | 60 | YZ | Midwest |
| 14 | 272 | AB | Northeast |
| 15 | 236 | CD | Northeast |
| 16 | 245 | EF | Northeast |
| 17 | 265 | GH | Northeast |
| 18 | 261 | IJ | Northeast |
| 19 | 252 | KL | Northeast |

6) **Find the maximum transaction of each account using Credit as transaction type**

```sql
EXPLAIN ANALYZE
select acc.account_id, max(trans.amount) from account acc
join transaction trans on trans.account_id=acc.account_id
where type='Credit'
group by acc.account_id;
```

public | RetailBanking | *<defaultdb> Script-1 ×

```sql
select acc.account_id, max(trans.amount) from account acc
join transaction trans on trans.account_id=acc.account_id
where type='Credit'
group by acc.account_id;
```

**account 1**

select acc.account_id, max(trans.amount) from account acc join tra... | Enter a SQL expression to filter r

| | account_id | max |
|---|---|---|
| 1 | A00000001 | 12,600 |
| 2 | A00000002 | 30,354 |
| 3 | A00000003 | 11,253 |
| 4 | A00000004 | 5,553 |
| 5 | A00000005 | 5,017 |
| 6 | A00000006 | 6,669 |
| 7 | A00000007 | 33,975 |
| 8 | A00000008 | 30,712 |
| 9 | A00000009 | 45,691 |
| 10 | A00000010 | 26,529 |
| 11 | A00000011 | 3,494 |
| 12 | A00000012 | 5,938 |
| 13 | A00000013 | 6,803 |
| 14 | A00000014 | 22,137 |
| 15 | A00000015 | 25,598 |
| 16 | A00000016 | 31,262 |
| 17 | A00000017 | 42,988 |
| 18 | A00000018 | 47,679 |
| 19 | A00000019 | 22,708 |
| 20 | A00000020 | 3,005 |
| 21 | A00000021 | 14,122 |
| 22 | A00000022 | 25,310 |
| 23 | A00000023 | 33,132 |
| 24 | A00000024 | 18,925 |
| 25 | A00000025 | 49,734 |
| 26 | A00000026 | 37,203 |
| 27 | A00000027 | 12,940 |
| 28 | A00000029 | 26,667 |
| 29 | A00000030 | 49,290 |
| 30 | A00000031 | 49,537 |
| 31 | A00000032 | 24,822 |
| 32 | A00000033 | 35,384 |
| 33 | A00000034 | 47,768 |
| 34 | A00000035 | 5,446 |
| 35 | A00000036 | 29,778 |

7) **Find the count of number of events registered for each Issue in crmevents which are not null**

```sql
select "CRMEvents"."Issue",count(*) from "CRMEvents" where "CRMEvents"."Issue" is not null
    group by "CRMEvents"."Issue"
```

public | RetailBanking | *<defaultdb> Script-1 ×

```sql
select "CRMEvents"."Issue",count(*) from "CRMEvents" where "CRMEvents"."Issue" is not null
    group by "CRMEvents"."Issue"
```

**CRMEvents 1**

select "CRMEvents"."Issue",count(*) from "CRMEvents" where "CRM | Enter a SQL expression to filter results (use Ctrl+Space)

| | Issue | count |
|---|---|---|
| 1 | Balance transfer fee | 32 |
| 2 | Other | 978 |
| 3 | Identity theft / Fraud / Embezzlement | 841 |
| 4 | Deposits and withdrawals | 3,786 |
| 5 | Account opening, closing, or management | 5,882 |
| 6 | Delinquent account | 314 |
| 7 | Arbitration | 40 |
| 8 | Privacy | 48 |
| 9 | Sale of account | 31 |
| 10 | Bankruptcy | 43 |
| 11 | Advertising and marketing | 192 |
| 12 | Payoff process | 217 |
| 13 | Application processing delay | 34 |
| 14 | Cash advance fee | 36 |
| 15 | Customer service / Customer relations | 373 |
| 16 | Credit reporting | 270 |
| 17 | Problems caused by my funds being low | 2,041 |
| 18 | Collection practices | 92 |
| 19 | Convenience checks | 26 |
| 20 | Forbearance / Workout plans | 61 |
| 21 | Making/receiving payments, sending money | 1,246 |
| 22 | Credit line increase/decrease | 252 |
| 23 | Other fee | 347 |
| 24 | Billing statement | 262 |
| 25 | Balance transfer | 139 |
| 26 | Rewards | 290 |
| 27 | Transaction issue | 346 |
| 28 | Closing/Cancelling account | 834 |
| 29 | Using a debit or ATM card | 1,132 |
| 30 | Cash advance | 39 |
| 31 | Billing disputes | 1,560 |
| 32 | Credit determination | 203 |
| 33 | Unsolicited issuance of credit card | 119 |
| 34 | APR or interest rate | 493 |
| 35 | Late fee | 344 |
| 36 | Collection debt dispute | 135 |
| 37 | Overlimit fee | 17 |
| 38 | Credit card protection / Debt protection | 318 |

*1) Insert / Update Queries:*

*<postgres> Script-33 ×

```sql
INSERT INTO card (card_id, disp_id, type, date)
    VALUES ('V00001635', 'D00001234', 'VISA Infinite', '2018-12-24');
```

**Statistics 1 ×**

| Name | Value |
|---|---|
| Updated Rows | 1 |
| Query | INSERT INTO card (card_id, disp_id, type, date) |
| | VALUES ('V00001635', 'D00001234', 'VISA Infinite', '2018-12-24') |
| Finish time | Thu Dec 01 20:02:13 EST 2022 |

Fig 36: Insert Card details

*<postgres> Script-33 ×

```sql
INSERT INTO district ("district_id", city, state_name)
    VALUES (78,'Buffalo', 'New York');
```

**Statistics 1 ×**

| Name | Value |
|---|---|
| Updated Rows | 1 |
| Query | INSERT INTO district ("district_id", city, state_name) |
| | VALUES (78,'Buffalo', 'New York') |
| Finish time | Thu Dec 01 20:01:57 EST 2022 |

Fig 37: Add a district

Fig 38: Update phone number of a client



Fig 39: Update priority of a complaint



Fig 40: delete account based on account_id



Fig 41: Delete order with order id and empty K_symbol

## X. QUERY EXECUTION ANALYSIS

In this part, three of the queries has been executed and analyzed using 'EXPLAIN ANALYZE'.

1) **For a particular sub-product, let's say "Savings Account", find how the customer support performed in terms of how many queries by customers were closed/open per region.**





Fig 27: Before - Without Indexing

Here cost=1459.82 and Execution Time: 321.292 ms Lets now index on Sub-product column in CRMEvents table and region column on state as they are executing in sequential scan which is making the query execution time and cost more.

Fig 28: Before Indexing

cost=16593.46 and Execution Time: 472.276 ms

Here as the transaction type is taking sequential scan and is causing the execution time of the query longer, we tried to index on transaction type column and see the results



Fig 28: After Indexing

cost=3325.53 and Execution Time: 232.250 ms
We can see there is significant decrease in the cost and execution time after Indexing.

3) **Find the count of number of events registered for each Issue in crmevents which are not null**



Fig 28: Before Indexing

cost=1362.61 and Execution Time: 10.869 ms

Here as the crmevents issue is taking sequential scan and is causing the execution time of the query longer, we tried to index on crmevents "issue" column and see the results



Fig 28: After Indexing

cost=90.33 and Execution Time: 151.332 ms
We can see there is significant decrease in the cost and execution time after Indexing.

2) **Find the maximum transaction of each account using Credit as transaction type:**

Fig 28: After Indexing

cost=0.29 and Execution Time: 4.570 ms
We can see there is significant decrease in the cost and execution time after Indexing.

## XI. DASHBOARD

The web application made as a part of the bonus task is up and available at https://lionfish-app-y3ij5.ondigitalocean.app/



Fig 29: Home Page



Fig 30: Search Clients



Fig 31: Search Loans



Fig 32: Search Transactions



Fig 33: Complaint submission Portal



Fig 34: Complaint submission Portal — Submitted Complaint



Fig 35: Search Complaints

### FUTURE SCOPE

* **Website and Public Portal** - We have already built a CRM dashboard that would help the bank officials manage behind-the-scenes. A future prospect is to add a public-facing customer portal to help audience at a broader scale interact with Eagle National Bank. Our database and design would ensure that there is no redundant data across the tables and maintain database's integrity.
* **Login Credentials** - Personalized dashboards can be secured with login credentials and distributed to officials to prevent fraudulent activities.
* **Trigger** - We can add Triggers at suitable places to perform routine cleanup activities to maintain the sanity of the database.

### REFERENCES

[1] "Retail Banking Demo Data - Dataset by Lpetrocelli." data.world, May 3, 2022. https://data.world/lpetrocelli/retail-banking-demo-data.