# The Bank Project
## May, 13th 2023
### By Craig Schlachter

## Key Findings

I analyzed data on 9,480 customers to identify users that are about to churn. My findings suggest that the following customers are at a higher risk of churning:

- Customers who spend less
- Customers who transact infrequently
- Customers who are dormant
- Female customers
- Middle-aged customers (41-55)
- Customers with low card utilization (0 - 0.333)
- Customers making less than $40,000 a year

These findings have implications for our bank that wants to reduce customer churn. The bank should focus on retaining customers who are at risk of churning by offering them discounts, rewards, or other incentives to keep them engaged.

The limitations of this study includes the fact that it is third party data and we noticed an unusually high amount of customers have been on the books for 36 months. We can't assess if this is a strange outlier or maybe a product promotion the bank ran for

a card product.

For future research, I would recommend conducting a study with a larger sample size and collecting more data on the customers who are at risk of churning. This would help to better understand the factors that contribute to customer churn and develop more effective strategies for retaining customers.

. . .

# Step 1: Ask

Here, we clarify the problem we are trying to solve, and the objectives we are trying to meet.

## 1.1 Background

A local bank has contacted us to help find solutions for their customer relations department. The bank manager is concerned about customer churn and wants to predict who is most likely to leave so they can take action to retain them.

The bank manager is confident that an in-depth analysis of their 10,000 banking customers would reveal more opportunities for customer service initiatives.

## 1.2 Business Task

Analyze a banking dataset consisting of 10,000 unique records to gain insights into how customer demographics, spending habits, and risk affects customer churn rates. Create a customer segmentation model that can be used to predict future churn rates and develop strategies to retain them.

### 1.3 Business Objectives

- What are the average churn rates for customers with RFM scores from 1-5?
- Who are the most common customer segments?
- What is each customer segments' average monthly RFM values?
- What segments are most likely to churn, and what are their characterizations?

### 1.4 Deliverables

- A clear summary of the business task
- A description of all data sources used
- Documentation of any cleaning or manipulation of data
- A summary of analysis
- Supporting visualizations and key findings
- High-level content recommendations based on the analysis

### 1.5 Key Stakeholders

- The bank's CEO, who is the leader of the bank and in charge of making financial lives better for their customers.
- The bank manager, who oversees the day-to-day operations of the branch, supervises staff and works to keep and attract new customers.
- The customer relations department, who are in charge of building and maintaining customer relationships.

. . .

# Step 2: Prepare

In the Prepare phase, we assess the data and its limitations.

## 2.1 Information on Data Source

1. Data is publicly available on [Kaggle: Credit Card Customers](#) and stored in 1 csv file.
2. This dataset was collected from a website [Leaps Analyttica](#), a website dedicated to learning data science.
3. This dataset consists of 10,000 unique records.
4. Data collected includes age, gender, marital status, income, education, credit limit, average utilization ratio, and churn.

## 2.2 Limitations of Data Set

- The data was collected at least 2 years ago in 2021. The data covers customers' demographics, spending habits, credit limits, average utilization ratio, revolving balance, average accounts open, and etc.
- It is a small sample size. There are only 10,000 customers in the dataset. This is a small number compared to the millions of credit card customers in the United States.
- This data is from a third party so we cannot ascertain its integrity or accuracy.
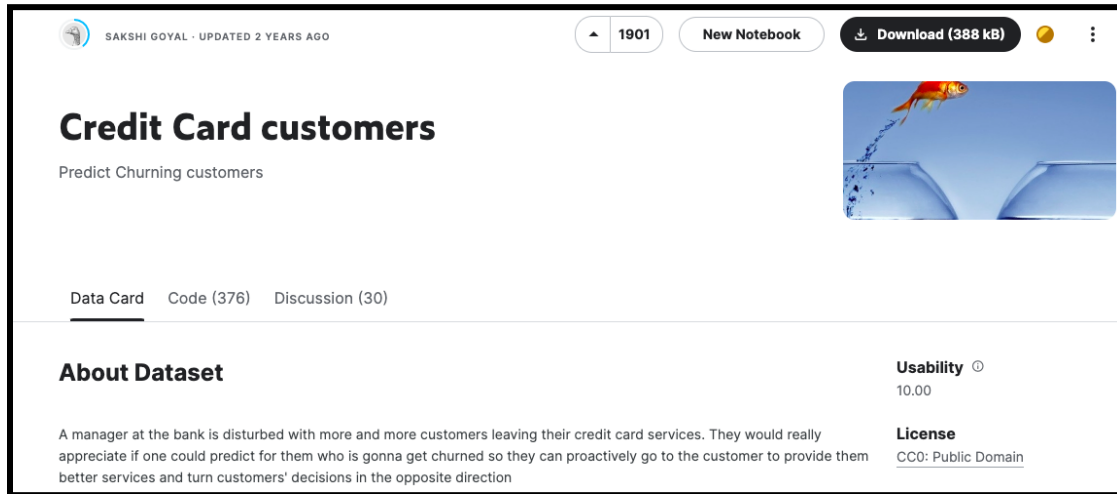
## 2.3 Is Data ROCCC?

A good data source is ROCCC which stands for Reliable, Original, Comprehensive, Current, and Cited.

- Reliable - MED - Not highly reliable as it only has 10,000 records
- Original - LOW - Third Party Provider ([Kaggle: Credit Card Customers](#))
- Comprehensive - MED - Parameters are sufficient to solve business task
- Current - MED - The data is 2 years old and may not be as relevant
- Cited - LOW - The data was collected from a third party, hence unknown

Overall, due to the limitations of this dataset, it is not recommended to use it to produce business intelligence recommendations.

## 2.4 Data Selection

The following file is downloaded and then imported into our created SQL table 'bank_cleaning'.



## 2.5 Tools

We are using SQL for data-wrangling and exploratory data analysis. Finally, we are using Tableau for visualizations.

. . .

# Step 3: Process

Here, we will process the data to ensure it is clean, correct, relevant, complete and free of errors and outliers by performing:

- Explore and observe the dataset
- Check for and handle any missing values
- Check for and remove any duplicate rows

- Ensure data is input and formatted correctly
- Check for and handle any outlier values
- Save cleaned data to a new file

# 3.1 Preparing the environment

The SQL table is created, columns are named, and data types are set.
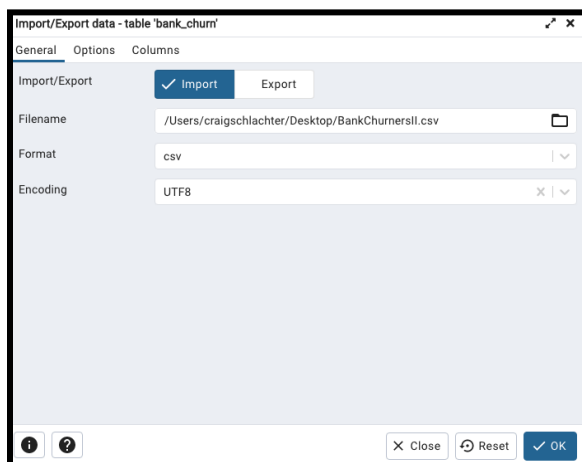
```
1   -- Table: public.bank_churn
2
3   -- DROP TABLE IF EXISTS public.bank_churn;
4
5   CREATE TABLE IF NOT EXISTS public.bank_churn
6   (
7       clietnum integer,
8       attrition_flag text COLLATE pg_catalog."default",
9       customer_age integer,
10      gender text COLLATE pg_catalog."default",
11      dependent_count integer,
12      education_level text COLLATE pg_catalog."default",
13      marital_status text COLLATE pg_catalog."default",
14      income_category text COLLATE pg_catalog."default",
15      card_category text COLLATE pg_catalog."default",
16      months_on_book integer,
17      total_relationship_count integer,
18      months_inactive integer,
19      contacts_count integer,
20      credit_limit numeric,
21      total_revolving_bal numeric,
22      avg_open_to_buy numeric,
23      total_amt_change_q4_q1 numeric,
24      total_trans_amt numeric,
25      total_trans_ct integer,
26      total_ct_change_q4_q1 numeric,
27      avg_utilization_ratio numeric
28  )
29
30  TABLESPACE pg_default;
31
32  ALTER TABLE IF EXISTS public.bank_churn
33      OWNER to postgres;
```

# 3.2 Importing Dataset

Reading in the selected file.

## 3.2 Data Cleaning and Manipulation

1.  Observe and familiarize with data

2.  Check for nulls or missing values

3.  Check for and remove duplicates

4.  Check for and remove outliers

5.  Perform validation checks of data

Previewing the first 100 rows to familiarize with the data.

```sql
/*

Cleaning Data For Bank Customer Segmentation Project

*/

-- Preview the data
SELECT *
FROM bank_churn
LIMIT 100
```

| clietnum integer | attrition_flag text | customer_age integer | gender text | dependent_count integer | education_level text | marital_status text | income_category text | card_category text |
|---|---|---|---|---|---|---|---|---|
| 713061558 | Existing Customer | 44 | M | 2 | Graduate | Married | $40K - $60K | Blue |
| 708790833 | Existing Customer | 42 | M | 5 | Uneducated | Unknown | $120K + | Blue |
| 710821833 | Existing Customer | 65 | M | 1 | Unknown | Married | $40K - $60K | Blue |
| 712396908 | Existing Customer | 57 | F | 2 | Graduate | Married | Less than $40K | Blue |
| 709327383 | Existing Customer | 45 | F | 2 | Graduate | Married | Unknown | Blue |
| 806165208 | Existing Customer | 47 | M | 1 | Doctorate | Divorced | $60K - $80K | Blue |
| 708508758 | Attrited Customer | 62 | F | 0 | Graduate | Married | Less than $40K | Blue |
| 811604133 | Existing Customer | 47 | F | 4 | Unknown | Single | Less than $40K | Blue |
| 789124683 | Existing Customer | 54 | M | 2 | Unknown | Married | $80K - $120K | Blue |
| 771071958 | Existing Customer | 41 | F | 3 | Graduate | Single | Less than $40K | Blue |
| 720466383 | Existing Customer | 59 | M | 1 | High School | Unknown | $40K - $60K | Blue |
| 804424383 | Existing Customer | 63 | M | 1 | Unknown | Married | $60K - $80K | Blue |
| 806624208 | Existing Customer | 47 | M | 4 | High School | Married | $40K - $60K | Blue |
| 787937058 | Existing Customer | 58 | M | 0 | Graduate | Married | $80K - $120K | Blue |

Our initial impression of the dataset is that it is well-structured, organized, and formatted in a way that is suitable for the business task at hand.

Now, we will check for nulls or missing values.

```sql
-- 1. Handle Missing Values
-- Check for the missing values in each column
SELECT
    SUM(CASE WHEN clietnum IS NULL THEN 1 ELSE 0 END) as clietnum_nulls
    ,SUM(CASE WHEN attrition_flag IS NULL THEN 1 ELSE 0 END) as attrition_flag_nulls
    ,SUM(CASE WHEN customer_age IS NULL THEN 1 ELSE 0 END) as customer_age_nulls
    ,SUM(CASE WHEN gender IS NULL THEN 1 ELSE 0 END) as gender_nulls
    ,SUM(CASE WHEN dependent_count IS NULL THEN 1 ELSE 0 END) as dependent_count_nulls
    ,SUM(CASE WHEN education_level IS NULL THEN 1 ELSE 0 END) as education_level_nulls
    ,SUM(CASE WHEN marital_status IS NULL THEN 1 ELSE 0 END) as marital_status_nulls
    ,SUM(CASE WHEN income_category IS NULL THEN 1 ELSE 0 END) as income_category_nulls
    ,SUM(CASE WHEN card_category IS NULL THEN 1 ELSE 0 END) as card_category_nulls
    ,SUM(CASE WHEN months_on_book IS NULL THEN 1 ELSE 0 END) as months_on_book_nulls
    ,SUM(CASE WHEN total_relationship_count IS NULL THEN 1 ELSE 0 END) as total_relationship_count_nulls
    ,SUM(CASE WHEN months_inactive IS NULL THEN 1 ELSE 0 END) as months_inactive_nulls
    ,SUM(CASE WHEN contacts_count IS NULL THEN 1 ELSE 0 END) as contacts_count_nulls
    ,SUM(CASE WHEN credit_limit IS NULL THEN 1 ELSE 0 END) as credit_limit_nulls
    ,SUM(CASE WHEN total_revolving_bal IS NULL THEN 1 ELSE 0 END) as total_revolving_bal_nulls
    ,SUM(CASE WHEN avg_open_to_buy IS NULL THEN 1 ELSE 0 END) as avg_open_to_buy_nulls
    ,SUM(CASE WHEN total_amt_change_q4_q1 IS NULL THEN 1 ELSE 0 END) as total_amt_change_q4_q1_nulls
    ,SUM(CASE WHEN total_trans_amt IS NULL THEN 1 ELSE 0 END) as total_trans_amt_nulls
    ,SUM(CASE WHEN total_trans_ct IS NULL THEN 1 ELSE 0 END) as total_trans_ct_nulls
    ,SUM(CASE WHEN total_ct_change_q4_q1 IS NULL THEN 1 ELSE 0 END) as total_ct_change_q4_q1_nulls
    ,SUM(CASE WHEN avg_utilization_ratio IS NULL THEN 1 ELSE 0 END) as avg_utilization_ratio_nulls

FROM bank_churn
```

| clietnum_nulls bigint | attrition_flag_nulls bigint | customer_age_nulls bigint | gender_nulls bigint | dependent_count_nulls bigint | education_level_nulls bigint | marital_status_nulls bigint |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

All rows were absent of nulls or missing values. Now, we will validate each individual column for missing values or nulls to validate our above query.

```
45    -- Going through each column to validate my first query's results
46    SELECT *
47    FROM bank_churn
48    --WHERE COALESCE(card_category,'') = ''
49    WHERE avg_utilization_ratio IS NULL --COALESCE(avg_utilization_ratio,0) = 0
50
```

Data Output    Messages    Notifications

| clietnum integer | attrition_flag text | customer_age integer | gender text | dependent_count integer | education_level text | marital_status text |
|---|---|---|---|---|---|---|

Our validation process confirms our initial query's findings.

Next, I am checking for and removing any duplicates column by column.

```
53    -- 2. Removing Duplicates
54    -- Check for duplicate rows in the table
55
56    SELECT avg_utilization_ratio
57          ,COUNT(*) AS amount
58    FROM bank_churn
59    GROUP BY 1
60    HAVING COUNT(*) > 1
61    ORDER BY 2 DESC
```

Data Output    Messages    Notifications

| | avg_utilization_ratio numeric | amount bigint |
|---|---|---|
| 1 | 0 | 2317 |
| 2 | 0.073 | 41 |
| 3 | 0.057 | 32 |
| 4 | 0.06 | 30 |
| 5 | 0.048 | 30 |
| 6 | 0.069 | 27 |
| 7 | 0.045 | 27 |
| 8 | 0.059 | 27 |
| 9 | 0.061 | 27 |
| 10 | 0.039 | 25 |
| 11 | 0.053 | 25 |

After reviewing the results for each column, I did not find any duplicates that appeared to be unnatural for this type of dataset.

Now, I am checking for and removing outliers.

```
94   -- 6. Handle Outlier Values
95   -- Looking at descriptive statistics to find outliers
96
97
98   -- Calculate the average, median, standard deviation, minimum, and maximum age of customers.
99   SELECT ROUND(avg(customer_age),2) AS avg_age
100         ,PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY customer_age) AS median_age
101         ,ROUND(STDDEV(customer_age),2) AS stdv_of_age
102         ,MIN(customer_age) AS min_age
103         ,MAX(customer_age) AS max_age
104  FROM bank_churn
105
106  -- There is one customer who is 73 years old, which is 3 standard deviations above the mean.
107  -- This customer is an outlier and may be skewing the results.
108
109
110  -- This code then identifies the rows that are outliers
111  SELECT customer_age
112         ,NTILE(100) OVER (ORDER BY customer_age) AS percentile
113  FROM bank_churn
114
115
116  -- This code then deletes the rows that are outliers
117  DELETE FROM bank_churn
118  WHERE customer_age = 73
```

We identified an outlier in the "customer_age" column and removed it.

```
148  -- Calculate the average, median, standard deviation, minimum, and maximum months inactive of customers
149  SELECT ROUND(avg(months_inactive),2) AS avg_months_inactive
150         ,PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY months_inactive) AS median_months_inactive
151         ,ROUND(STDDEV(months_inactive),2) AS stdv_of_months_inactive
152         ,MIN(months_inactive) AS min_months_inactive
153         ,MAX(months_inactive) AS max_months_inactive
154  FROM bank_churn
155
156  -- There is one customer who has been inactive for 3 standard deviations above the mean.
157  -- This customer is an outlier and may be skewing the results.
158
159
160  -- This code then identifies the rows that are outliers
161  SELECT months_inactive
162         ,NTILE(100) OVER (ORDER BY months_inactive) AS percentile
163  FROM bank_churn
164
165  -- This code then deletes the rows that are outliers
166  DELETE FROM bank_churn
167  WHERE months_inactive IN(
168                      WITH q1 AS (
169                              SELECT months_inactive
170                              ,NTILE(100) OVER (ORDER BY months_inactive) AS percentile
171                              FROM bank_churn
172                              )
173
174                      SELECT months_inactive
175                      FROM q1
176                      WHERE percentile = 100)
177                      RETURNING *;
```

We identified an outlier in the "months_inactive" column and removed it.

```
180  -- Calculate the average, median, standard deviation, minimum, and maximum contacts count of customers.
181  SELECT ROUND(avg(contacts_count),2) AS avg_contacts_count
182        ,PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY contacts_count) AS median_contacts_count
183        ,ROUND(STDDEV(contacts_count),2) AS stdv_of_contacts_count
184        ,MIN(contacts_count) AS min_contacts_count
185        ,MAX(contacts_count) AS max_contacts_count
186  FROM bank_churn
187
188  -- There is one customer who has contacts which are 3 standard deviations above the mean.
189  -- This customer is an outlier and may be skewing the results.
190
191
192  -- This code then identifies the rows that are outliers
193  SELECT contacts_count
194        ,NTILE(100) OVER (ORDER BY contacts_count) AS percentile
195  FROM bank_churn
196
197
198  -- This code then deletes the rows that are outliers
199  DELETE FROM bank_churn
200  WHERE contacts_count IN(
201                     WITH q1 AS (
202                            SELECT contacts_count
203                            ,NTILE(100) OVER (ORDER BY contacts_count) AS percentile
204                            FROM bank_churn
205                            )
206
207                     SELECT contacts_count
208                     FROM q1
209                     WHERE percentile = 100)
210                     RETURNING *;
```

We identified an outlier in the "contacts_count" column and removed it.

```
241  -- Calculate the average, median, standard deviation, minimum, and maximum values for total_amt_change_q4_q1
242  SELECT ROUND(avg(total_amt_change_q4_q1),2) AS avg_total_amt_change_q4_q1
243        ,PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY total_amt_change_q4_q1) AS median_total_amt_change_q4_q1
244        ,ROUND(STDDEV(total_amt_change_q4_q1),2) AS stdv_of_total_amt_change_q4_q1
245        ,MIN(total_amt_change_q4_q1) AS min_total_amt_change_q4_q1
246        ,MAX(total_amt_change_q4_q1) AS max_total_amt_change_q4_q1 -- Outlier of 3 STDDEV
247  FROM bank_churn
248
249
250  -- This code then identifies the rows that are outliers
251  SELECT total_amt_change_q4_q1
252        ,NTILE(100) OVER (ORDER BY total_amt_change_q4_q1) AS percentile
253  FROM bank_churn
254
255
256  -- This code then deletes the rows that are outliers
257  DELETE FROM bank_churn
258  WHERE total_amt_change_q4_q1 IN(
259                     WITH q1 AS (
260                            SELECT total_amt_change_q4_q1
261                            ,NTILE(100) OVER (ORDER BY total_amt_change_q4_q1) AS percentile
262                            FROM bank_churn
263                            )
264
265                     SELECT total_amt_change_q4_q1
266                     FROM q1
267                     WHERE percentile = 100)
268                     RETURNING *;
```

We identified an outlier in the "total_amt_change_q4_q1" column and removed it.

```
271  -- Calculate the average, median, standard deviation, minimum, and maximum values for total_trans_amt
272  SELECT ROUND(avg(total_trans_amt),2) AS avg_total_trans_amt
273       ,PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY total_trans_amt) AS median_total_trans_amt
274       ,ROUND(STDDEV(total_trans_amt),2) AS stdv_of_total_trans_amt
275       ,MIN(total_trans_amt) AS min_total_trans_amt
276       ,MAX(total_trans_amt) AS max_total_trans_amt -- Outlier of 3 STDDEV
277  FROM bank_churn
278
279
280  -- This code then identifies the rows that are outliers
281  SELECT total_trans_amt
282       ,NTILE(100) OVER (ORDER BY total_trans_amt) AS percentile
283  FROM bank_churn
284
285  -- This code then deletes the rows that are outliers
286  DELETE FROM bank_churn
287  WHERE total_trans_amt IN(
288                      WITH q1 AS (
289                              SELECT total_trans_amt
290                              ,NTILE(100) OVER (ORDER BY total_trans_amt) AS percentile
291                              FROM bank_churn
292                              )
293
294                      SELECT total_trans_amt
295                      FROM q1
296                      WHERE percentile = 100)
297                      RETURNING *;
```

We identified an outlier in the "total_trans_amt" column and removed it.

```
300  -- Calculate the average, median, standard deviation, minimum, and maximum values for total_trans_ct
301  SELECT ROUND(avg(total_trans_ct),2) AS avg_total_trans_ct
302       ,PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY total_trans_ct) AS median_total_trans_ct
303       ,ROUND(STDDEV(total_trans_ct),2) AS stdv_of_total_trans_ct
304       ,MIN(total_trans_ct) AS min_total_trans_ct
305       ,MAX(total_trans_ct) AS max_total_trans_ct -- Outlier of 3 STDDEV
306  FROM bank_churn
307
308
309  -- This code then deletes the rows that are outliers
310  DELETE FROM bank_churn
311  WHERE total_trans_ct > 133
```

We identified an outlier in the "total_trans_ct" column and removed it.

```
315    -- Calculate the average, median, standard deviation, minimum, and maximum values for total_ct_change_q4_q1
316    SELECT ROUND(avg(total_ct_change_q4_q1),2) AS avg_total_ct_change_q4_q1
317          ,PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY total_ct_change_q4_q1) AS median_total_ct_change_q4_q1
318          ,ROUND(STDDEV(total_ct_change_q4_q1),2) AS stdv_of_total_ct_change_q4_q1
319          ,MIN(total_ct_change_q4_q1) AS min_total_ct_change_q4_q1
320          ,MAX(total_ct_change_q4_q1) AS max_total_ct_change_q4_q1 -- Outlier of 3 STDDEV
321    FROM bank_churn
322
323
324
325    -- This code then identifies the rows that are outliers
326    SELECT total_ct_change_q4_q1
327          ,NTILE(100) OVER (ORDER BY total_ct_change_q4_q1 ASC) AS percentile
328    FROM bank_churn
329
330
331    -- This code then deletes the rows that are outliers
332    DELETE FROM bank_churn
333    WHERE total_ct_change_q4_q1 IN(
334                        WITH q1 AS (
335                                    SELECT total_ct_change_q4_q1
336                                    ,NTILE(100) OVER (ORDER BY total_ct_change_q4_q1) AS percentile
337                                    FROM bank_churn
338                                    )
339
340                        SELECT total_ct_change_q4_q1
341                        FROM q1
342                        WHERE percentile = 100)
343                        RETURNING *;
344
345    DELETE FROM bank_churn
346    WHERE total_ct_change_q4_q1 = 0
347    RETURNING *;
```

We identified an outlier in the "total__ct_change_q4_q1" column and removed it.

```
360    -- 7. Validate the Data
361    /* Months_on_book(36) has nearly ten times the count of next record (37)
362    It may be caused by a bank policy or deal offered to customers. Will be
363    looking at how sensitive data is to this metric in EDA phase. */
364
```

We identified another outlier in the "months_on_book" column while we were validating the data and will determine its sensitivity during the EDA phase.

Now that our data is cleaned and saved, we will move onto the next phase to explore and analyze the data.

# Step 4: Analyze

## 4.1 Performing Calculations

Pulling statistics for analysis:

1. Examine and select the most significant metrics for bank customers
2. Create customer segments to identify customers at risk of churning
3. Build an RFM model to improve customer retention
4. Compute the average recency, frequency, and monetary value for each customer segment by month.
5. Find the percentage of customers who fall into each demographic segment and have certain spending habits within each RFM customer segment.
6. Calculate the average churn rate for each unique combination of recency and frequency/monetary.

**[Step 1's Results]**

```
1  /*
2
3
4    Performing Exploratory Data Analysis for Bank Segmentation Project
5
6
7  */
8
9
10
11  -- 0. Looking at the key metrics of the credit card users
12
13  SELECT ROUND(COUNT(clietnum)/1000.0,2)|| 'K' AS total_customers
14      ,'$'||ROUND(SUM(total_trans_amt)/1000000,2)|| 'M' AS total_spend
15      ,'$'||ROUND((SUM(total_trans_amt)/12)/1000000,2)|| 'M' AS total_monthly_spend
16      ,ROUND((SUM(total_trans_ct)/12)/1000.0,2)||'K' AS total_monthly_transactions
17  FROM bank_churn
18
```

Data Output   Messages   Notifications

| total_customers text | total_spend text | total_monthly_spend text | total_monthly_transactions text |
|---|---|---|---|
| 9.48K | $41.18M | $3.43M | 51.44K |

**[Step 2's Results]**

```
25  -- 2. Looking at the prominent segment groups
26  WITH prfm AS (
27  /* Create a Common Table Expression (CTE) called 'prfm' to compute the Recency,
28  Frequency,and Monetary (RFM) values for each customer in the dataset */
29      SELECT  clietnum
30          ,months_inactive AS Recency
31          ,total_trans_ct AS Frequency
32          ,total_trans_amt AS Monetary
33          ,NTILE(5) OVER (ORDER BY months_inactive DESC) AS r
34          ,NTILE(5) OVER (ORDER BY total_trans_ct ASC) AS f
35          ,NTILE(5) OVER (ORDER BY total_trans_amt ASC) AS m
36      FROM bank_churn
37      ORDER BY 2 DESC
38  ),
39
40      segment_atr AS(
41      -- Create another CTE called 'segment_atr' to combine RFM values and other attributes of each customer
42      SELECT bc.clietnum
43          ,bc.customer_age
44          ,bc.marital_status
45          ,bc.education_level
46          ,bc.income_category
47          ,prfm.recency  -- Computing rfm by month
48          ,ROUND(prfm.frequency::numeric/12,2) AS frequency
49          ,ROUND(prfm.monetary::numeric/12,2) AS monetary
50          ,prfm.r
51          ,prfm.f
52          ,prfm.m
53          ,ROUND((prfm.f + prfm.m)/2,0) AS fm
54          ,SUM(CASE WHEN bc.attrition_flag = 'Attrited Customer' THEN 1 ELSE 0 END) AS churned_customers
55
56      FROM bank_churn bc
57      JOIN prfm ON (bc.clietnum = prfm.clietnum)
58      GROUP BY 1,2,3,4,5,6,7,8,9,10,11,12
59      ORDER BY 1 ASC
```

**[Step 2's Results]**

```
62      segments AS(
63      -- Create a CTE called 'segments' to assign RFM segments based on the quantiles
64      SELECT *
65          ,CASE WHEN (r = 5 AND fm = 5) OR (r = 5 AND fm = 4) OR (r = 4 AND fm = 5) THEN 'Champions'
66          WHEN (r = 5 AND fm = 3) OR (r = 4 AND fm = 4) OR (r = 3 AND fm = 5) OR (r = 3 AND fm = 4) THEN 'Engaged'
67          WHEN (r = 5 AND fm = 2) OR (r = 4 AND fm = 2) OR (r = 3 AND fm = 3) OR (r = 4 AND fm = 3) THEN 'Potential'
68          WHEN (r = 5 AND fm = 1) THEN 'Recent Customers'
69          WHEN (r = 4 AND fm = 1) OR (r = 3 AND fm = 1) THEN 'Promising'
70          WHEN (r = 3 AND fm = 2) OR (r = 2 AND fm = 3) OR (r = 2 AND fm = 2) THEN 'Need Attention'
71          WHEN (r = 2 AND fm = 1) THEN 'About to sleep'
72          WHEN (r = 2 AND fm = 5) OR (r = 2 AND fm = 4) OR (r = 1 AND fm = 3) THEN 'At risk'
73          WHEN (r = 1 AND fm = 5) OR (r = 1 AND fm = 4) THEN 'Can't lose them'
74          WHEN r = 1 AND fm = 2 THEN 'Hibernating'
75          WHEN r = 1 AND fm = 1 THEN 'Lost' END AS rfm_segment
76      FROM segment_atr
77      ORDER BY 1 ASC
78          )
79
80  SELECT rfm_segment
81      ,COUNT(*)
82      ,COUNT(*) FILTER(WHERE churned_customers = 1) AS total_churned
83      ,COALESCE(ROUND((SUM(churned_customers) FILTER(WHERE churned_customers = 1)/COUNT(*))*100,2),0) AS churn_rate
84  FROM segments
85  GROUP BY 1
86  ORDER BY 4 DESC
87
88  /*This query computes the RFM (Recency, Frequency, Monetary) values for each customer and categorizes them
89  into five segments using the NTILE function. It also computes the number of churned customers per segment. */
```

**[Step 2's Results]**

| | rfm_segment text | count bigint | total_churned bigint | churn_rate numeric |
|---|---|---|---|---|
| 1 | Hibernating | 478 | 286 | 59.83 |
| 2 | Lost | 1049 | 420 | 40.04 |
| 3 | Promising | 841 | 245 | 29.13 |
| 4 | Need Attention | 1359 | 192 | 14.13 |
| 5 | At risk | 909 | 128 | 14.08 |
| 6 | Recent Customers | 417 | 50 | 11.99 |
| 7 | Can't lose them | 204 | 14 | 6.86 |
| 8 | Potential | 1432 | 73 | 5.10 |
| 9 | Engaged | 1754 | 58 | 3.31 |
| 10 | Champions | 1015 | 22 | 2.17 |
| 11 | About to sleep | 21 | 0 | 0 |

## [Step 3's Results]

```sql
92    -- 3. Compute RFM Reports for all customers
93    WITH crfm AS (
94        -- Create a Common Table Expression (CTE) called 'crfm'
95        -- Calculate Recency, Frequency and Monetary (RFM) values for each client
96        -- by dividing the data into 5 equal segments (quintiles) using NTILE
97        SELECT  clietnum
98                ,months_inactive AS Recency
99                ,total_trans_ct AS Frequency
100               ,total_trans_amt AS Monetary
101               ,NTILE(5) OVER (ORDER BY months_inactive DESC) AS r
102               ,NTILE(5) OVER (ORDER BY total_trans_ct ASC) AS f
103               ,NTILE(5) OVER (ORDER BY total_trans_amt ASC) AS m
104       FROM bank_churn
105       ORDER BY 2 DESC
```

## [Step 4's Results]

```sql
93    WITH crfm AS (
94        -- Create a Common Table Expression (CTE) called 'crfm'
95        -- Calculate Recency, Frequency and Monetary (RFM) values for each client
96        -- by dividing the data into 5 equal segments (quintiles) using NTILE
97        SELECT  clietnum
98                ,months_inactive AS Recency
99                ,total_trans_ct AS Frequency
00                ,total_trans_amt AS Monetary
01                ,NTILE(5) OVER (ORDER BY months_inactive DESC) AS r
02                ,NTILE(5) OVER (ORDER BY total_trans_ct ASC) AS f
03                ,NTILE(5) OVER (ORDER BY total_trans_amt ASC) AS m
04        FROM bank_churn
05        ORDER BY 2 DESC
06    ),
07
08    segment_atr AS(
09        -- Create another CTE called 'segment_atr'
10        -- Join the 'crfm' CTE with the 'bank_churn' table on client number
11        -- Calculate RFM values per month
12        SELECT bc.clietnum
13                ,bc.customer_age
14                ,bc.marital_status
15                ,bc.education_level
16                ,bc.income_category
17                ,crfm.recency  -- Computing rfm by month
18                ,ROUND(crfm.frequency::numeric/12,2) AS frequency
19                ,ROUND(crfm.monetary::numeric/12,2) AS monetary
20                ,crfm.r
21                ,crfm.f
22                ,crfm.m
23                ,ROUND((crfm.f + crfm.m)/2,0) AS fm
24                ,SUM(CASE WHEN bc.attrition_flag = 'Attrited Customer' THEN 1 ELSE 0 END) AS churned_customers
25
26        FROM bank_churn bc
27        JOIN crfm ON (bc.clietnum = crfm.clietnum)
```

## [Step 4's Results]

```sql
126           FROM bank_churn bc
127           JOIN crfm ON (bc.clietnum = crfm.clietnum)
128           GROUP BY 1,2,3,4,5,6,7,8,9,10,11,12
129           ORDER BY 1 ASC
130       ),
131
132       segments AS(
133           -- Create a third CTE called 'segments'
134           -- Classify clients into segments based on their RFM scores
135           SELECT *
136               ,CASE WHEN (r = 5 AND fm = 5) OR (r = 5 AND fm = 4) OR (r = 4 AND fm = 5) THEN 'Champions'
137                   WHEN (r = 5 AND fm = 3) OR (r = 4 AND fm = 4) OR (r = 3 AND fm = 5) OR (r = 3 AND fm = 4) THEN 'Engaged'
138                   WHEN (r = 5 AND fm = 2) OR (r = 4 AND fm = 2) OR (r = 3 AND fm = 3) OR (r = 4 AND fm = 3) THEN 'Potential'
139                   WHEN (r = 5 AND fm = 1) THEN 'Recent Customers'
140                   WHEN (r = 4 AND fm = 1) OR (r = 3 AND fm = 1) THEN 'Promising'
141                   WHEN (r = 3 AND fm = 2) OR (r = 2 AND fm = 3) OR (r = 2 AND fm = 2) THEN 'Need Attention'
142                   WHEN (r = 2 AND fm = 1) THEN 'About to sleep'
143                   WHEN (r = 2 AND fm = 5) OR (r = 2 AND fm = 4) OR (r = 1 AND fm = 3) THEN 'At risk'
144                   WHEN (r = 1 AND fm = 5) OR (r = 1 AND fm = 4) THEN 'Can''t lose them'
145                   WHEN r = 1 AND fm = 2 THEN 'Hibernating'
146                   WHEN r = 1 AND fm = 1 THEN 'Lost' END AS rfm_segment
147           FROM segment_atr
148           ORDER BY 1 ASC
149       )
150
151   SELECT rfm_segment
152           ,ROUND(AVG(recency),2) AS average_recency
153           ,ROUND(AVG(frequency),2) AS average_frequency
154           ,ROUND(AVG(monetary),2) AS average_monetary
155   FROM segments
156   GROUP BY 1
157   ORDER BY 2 DESC, 3 ASC, 4 ASC
158
```

## [Step 4's Results]

| | rfm_segment<br>text 🔒 | average_recency<br>numeric 🔒 | average_frequency<br>numeric 🔒 | average_monetary<br>numeric 🔒 |
|---|---|---|---|---|
| 1 | Can't lose them | 4.34 | 7.06 | 519.66 |
| 2 | Hibernating | 3.27 | 4.06 | 210.99 |
| 3 | At risk | 3.20 | 6.42 | 437.38 |
| 4 | Lost | 3.16 | 2.95 | 138.75 |
| 5 | About to sleep | 3.00 | 4.78 | 148.72 |
| 6 | Need Attention | 2.83 | 5.38 | 298.12 |
| 7 | Engaged | 1.95 | 7.10 | 560.78 |
| 8 | Promising | 1.91 | 2.93 | 136.66 |
| 9 | Potential | 1.73 | 5.46 | 297.91 |
| 10 | Champions | 1.23 | 7.61 | 684.25 |
| 11 | Recent Customers | 0.97 | 2.98 | 128.68 |

## [Step 5's Results]

```sql
161   -- 4. Showing pct. of gender, age_ranges, spend utilization
162   --ranges, income segments, and educational segments.
163   WITH crfm AS (
164               SELECT  clietnum
165                       ,months_inactive AS Recency
166                       ,total_trans_ct AS Frequency
167                       ,total_trans_amt AS Monetary
168                       ,NTILE(5) OVER (ORDER BY months_inactive DESC) AS r
169                       ,NTILE(5) OVER (ORDER BY total_trans_ct ASC) AS f
170                       ,NTILE(5) OVER (ORDER BY total_trans_amt ASC) AS m
171               FROM bank_churn
172               ORDER BY 2 DESC
173           ),
174
175       segment_atr AS(
176               SELECT bc.clietnum
177                       ,bc.customer_age
178                       ,bc.gender
179                       ,bc.marital_status
180                       ,bc.education_level
181                       ,bc.income_category
182                       ,bc.avg_utilization_ratio
183                       ,crfm.recency  -- Computing rfm by month
184                       ,ROUND(crfm.frequency::numeric/12,2) AS frequency
185                       ,ROUND(crfm.monetary::numeric/12,2) AS monetary
186                       ,crfm.r
187                       ,crfm.f
188                       ,crfm.m
189                       ,ROUND((crfm.f + crfm.m)/2,0) AS fm
190                       ,SUM(CASE WHEN bc.attrition_flag = 'Attrited Customer' THEN 1 ELSE 0 END) AS churned_customers
191
192               FROM bank_churn bc
193               JOIN crfm ON (bc.clietnum = crfm.clietnum)
194               GROUP BY 1,2,3,4,5,6,7,8,9,10,11,12,13,14
195               ORDER BY 1 ASC
```

## [Step 5's Results]

```sql
198   segments AS(
199       SELECT *
200           ,CASE WHEN (r = 5 AND fm = 5) OR (r = 5 AND fm = 4) OR (r = 4 AND fm = 5) THEN 'Champions'
201               WHEN (r = 5 AND fm = 3) OR (r = 4 AND fm = 4) OR (r = 3 AND fm = 5) OR (r = 3 AND fm = 4) THEN 'Engaged'
202               WHEN (r = 5 AND fm = 2) OR (r = 4 AND fm = 2) OR (r = 3 AND fm = 3) OR (r = 4 AND fm = 3) THEN 'Potential'
203               WHEN (r = 5 AND fm = 1) THEN 'Recent Customers'
204               WHEN (r = 4 AND fm = 1) OR (r = 3 AND fm = 1) THEN 'Promising'
205               WHEN (r = 3 AND fm = 2) OR (r = 2 AND fm = 3) OR (r = 2 AND fm = 2) THEN 'Need Attention'
206               WHEN (r = 2 AND fm = 1) THEN 'About to sleep'
207               WHEN (r = 2 AND fm = 5) OR (r = 2 AND fm = 4) OR (r = 1 AND fm = 3) THEN 'At risk'
208               WHEN (r = 1 AND fm = 5) OR (r = 1 AND fm = 4) THEN 'Can''t lose them'
209               WHEN r = 1 AND fm = 2 THEN 'Hibernating'
210               WHEN r = 1 AND fm = 1 THEN 'Lost' END AS rfm_segment
211           ,CASE WHEN customer_age BETWEEN 26 AND 40 THEN 'Younger'
212               WHEN customer_age BETWEEN 41 AND 55 THEN 'Middle Aged'
213               WHEN customer_age BETWEEN 56 AND 70 THEN 'Older' END AS age_group
214           ,CASE WHEN avg_utilization_ratio BETWEEN 0.000 AND 0.333 THEN 'Low Use'
215               WHEN avg_utilization_ratio BETWEEN 0.334 AND 0.667 THEN 'Medium Use'
216               WHEN avg_utilization_ratio BETWEEN 0.668 AND 0.999 THEN 'High Use'
217               END AS utilization_category
218       FROM segment_atr
219       ORDER BY 1 ASC
220       ),
221
222   seg_vit AS (
223       SELECT rfm_segment
224               ,age_group
225               ,customer_age
226               ,gender
227               ,education_level
228               ,income_category
229               ,utilization_category
230               ,COUNT(*) OVER (PARTITION BY rfm_segment, age_group, utilization_category) AS total_segment_count
231
232       FROM segments
```

**[Step 5's Results]**

```
232         FROM segments
233         WHERE churned_customers = 1
234         AND rfm_segment IN ('Hibernating','Lost','Promising','Need Attention','At risk')
235         ORDER BY 8 DESC
236         )
237
238 SELECT rfm_segment
239     ,ROUND((COUNT(*) FILTER(WHERE gender = 'M')/COUNT(*)::numeric)*100.0,2) AS percent_males
240     ,ROUND((COUNT(*) FILTER(WHERE gender = 'F')/COUNT(*)::numeric)*100.0,2) AS percent_females
241     ,ROUND((COUNT(*) FILTER(WHERE age_group = 'Younger')/COUNT(*)::numeric)*100.0,2) AS percent_young
242     ,ROUND((COUNT(*) FILTER(WHERE age_group = 'Middle Aged')/COUNT(*)::numeric)*100.0,2) AS percent_middle_aged
243     ,ROUND((COUNT(*) FILTER(WHERE age_group = 'Older')/COUNT(*)::numeric)*100.0,2) AS percent_old
244     ,ROUND((COUNT(*) FILTER(WHERE utilization_category = 'Low Use')/COUNT(*)::numeric)*100.0,2) AS percent_low_use
245     ,ROUND((COUNT(*) FILTER(WHERE utilization_category = 'Medium Use')/COUNT(*)::numeric)*100.0,2) AS percent_med_use
246     ,ROUND((COUNT(*) FILTER(WHERE utilization_category = 'High Use')/COUNT(*)::numeric)*100.0,2) AS percent_high_use
247     ,ROUND((COUNT(*) FILTER(WHERE income_category = 'Less than $40K')/COUNT(*)::numeric)*100.0,2) AS percent_under_$40K
248     ,ROUND((COUNT(*) FILTER(WHERE income_category = '$40K - $60K')/COUNT(*)::numeric)*100.0,2) AS percent_$40K_$60K
249     ,ROUND((COUNT(*) FILTER(WHERE income_category = '$60K - $80K')/COUNT(*)::numeric)*100.0,2) AS percent_$60K_$80K
250     ,ROUND((COUNT(*) FILTER(WHERE income_category = '$80K - $120K')/COUNT(*)::numeric)*100.0,2) AS percent_$80K_$120K
251     ,ROUND((COUNT(*) FILTER(WHERE income_category = '$120K +')/COUNT(*)::numeric)*100.0,2) AS percent_over_$120K
252     ,ROUND((COUNT(*) FILTER(WHERE income_category = 'Unknown')/COUNT(*)::numeric)*100.0,2) AS percent_unknown
253     ,ROUND((COUNT(*) FILTER(WHERE education_level = 'Uneducated')/COUNT(*)::numeric)*100.0,2) AS percent_uneducated
254     ,ROUND((COUNT(*) FILTER(WHERE education_level = 'High School')/COUNT(*)::numeric)*100.0,2) AS percent_high_school
255     ,ROUND((COUNT(*) FILTER(WHERE education_level = 'College')/COUNT(*)::numeric)*100.0,2) AS percent_college
256     ,ROUND((COUNT(*) FILTER(WHERE education_level = 'Graduate')/COUNT(*)::numeric)*100.0,2) AS percent_graduate
257     ,ROUND((COUNT(*) FILTER(WHERE education_level = 'Post-Graduate')/COUNT(*)::numeric)*100.0,2) AS percent_post_graduate
258     ,ROUND((COUNT(*) FILTER(WHERE education_level = 'Doctorate')/COUNT(*)::numeric)*100.0,2) AS percent_doctorate
259     ,ROUND((COUNT(*) FILTER(WHERE education_level = 'Unknown')/COUNT(*)::numeric)*100.0,2) AS percent_unknown
260 FROM seg_vit
261 GROUP BY 1
```

**[Step 5's Results(Summary of Results)]**

| rfm_segment text | percent_males numeric | percent_females numeric | percent_young numeric | percent_middle_aged numeric | percent_old numeric | percent_low_use numeric | percent_med_use numeric |
|---|---|---|---|---|---|---|---|
| 1 At risk | 64.84 | 35.16 | 15.63 | 75.78 | 8.59 | 91.41 | 7.03 |
| 2 Hibernating | 30.42 | 69.58 | 13.64 | 74.83 | 11.54 | 73.78 | 14.69 |
| 3 Lost | 43.57 | 56.43 | 14.05 | 70.24 | 15.71 | 79.29 | 10.71 |
| 4 Need Attention | 38.54 | 61.46 | 28.13 | 64.58 | 7.29 | 79.17 | 7.81 |
| 5 Promising | 37.96 | 62.04 | 22.86 | 67.35 | 9.80 | 77.14 | 11.84 |

**[Step 6's Results]**

```
265 -- 5. Showing a matrix for recency, frequency/monetary and the resulting churn rate
266 WITH matrix AS (
267         SELECT  clietnum
268             ,months_inactive AS Recency
269             ,total_trans_ct AS Frequency
270             ,total_trans_amt AS Monetary
271             ,NTILE(5) OVER (ORDER BY months_inactive DESC) AS r
272             ,NTILE(5) OVER (ORDER BY total_trans_ct ASC) AS f
273             ,NTILE(5) OVER (ORDER BY total_trans_amt ASC) AS m
274         FROM bank_churn
275         ORDER BY 2 DESC
276         )
277
278
279 SELECT r
280     ,ROUND((f + m)/2,0) AS fm
281     ,ROUND(COUNT(*) FILTER(WHERE bc.attrition_flag = 'Attrited Customer')/COUNT(*)::numeric*100,2) churned_percent
282
283 FROM matrix mt
284 JOIN bank_churn bc ON (mt.clietnum=bc.clietnum)
285 GROUP BY 1,2
286 ORDER BY churned_percent DESC
```

**[Step 6's Results(Summary of Results)]**

| | r integer | fm numeric | churned_percent numeric |
|---|---|---|---|
| 1 | 1 | 2 | 59.83 |
| 2 | 3 | 2 | 59.21 |
| 3 | 1 | 1 | 40.04 |
| 4 | 3 | 1 | 32.65 |
| 5 | 4 | 1 | 21.09 |
| 6 | 1 | 3 | 16.97 |
| 7 | 2 | 4 | 13.88 |

Interpreting statistical findings:

1. Our data analysis showed that our total customer base was **9,480 customers**. They **spent $41.18M** over the **year**, with an average of **51.44K transactions** per **month**. Unfortunately, **15.70%** of them **churned**.

2. We created **11 different segments** for our customers based on their RFM scores. This allowed us to **identify customers who are at risk of churning** and target them with specific marketing campaigns.

3. We analyzed our data and found that customers who were **inactive for 3 months** or more were more likely to churn. We also found that customers who **spent less than $300 per month** and who **transacted less than 5.5 times per month** were at a higher risk of churning.

4. Our data analysis shows that the majority of customers who are **most likely to churn** are **female**, **middle-aged**, with l**ow card utilization**, **low income**, and a

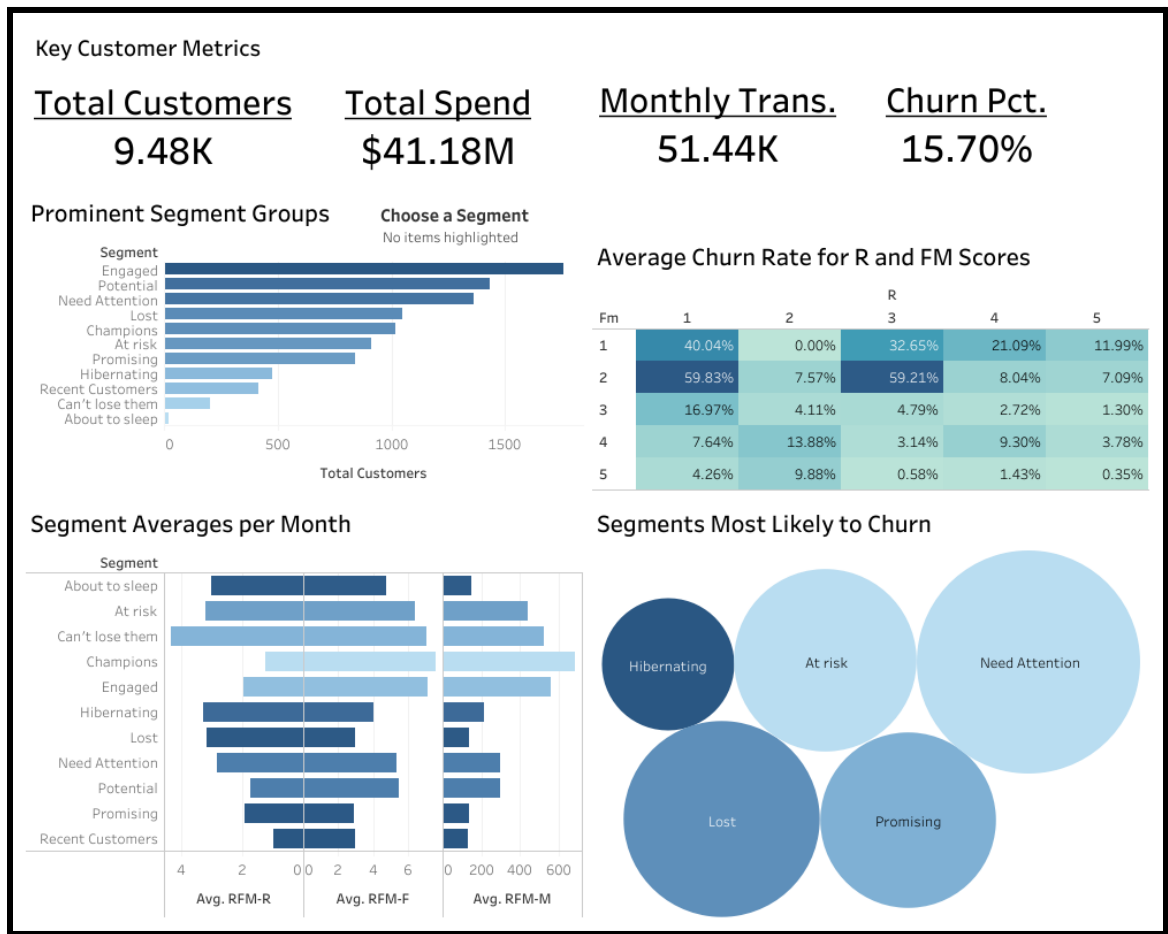**graduate level of education**. These customers are in the Hibernating, Lost, and Promising segments.

5. We found that customers who had recent activity and **low frequency and monetary scores of 1 or 2** were more likely to churn. This was also true for customers who had a **recency score of 3 or 4** and low frequency and monetary scores.

# Step 5: Share

In this step, we are creating visualizations and communicating our findings based on our analysis.

## 5.1 Data Visualizations and Findings

**Customer Segmentation by RFM, Demographics, and Spending**
(Visit my tableau profile for interactivity with this dashboard)

This dashboard analyzes the most prominent segment groups, their RFM averages, and identifies trends about segments that are most likely to churn.

1.  We have identified three customer segments that are most likely to churn: **Hibernating**, **Lost**, and **Promising**. These segments represent **24.97% of our customer base** and their behavior and demographics can help us develop strategies to limit churn.
2.  The Hibernating, Lost, and Promising segments had a higher churn rate, **40.16%**, than other segments. This is likely due to their behavior, which is characterized by less recent, less frequent, and lower spending.
3.  We believe that there is a **correlation** between **frequency** of purchase and **retention rate**. Customers who purchase more frequently are more likely to be satisfied with our products and services, and they are less likely to switch to a competitor.

.   .   .

# Step 6: Act

In the final step, we will be delivering our insights and providing recommendations based on our analysis.

Here, we revisit our business questions and share with you our high-level business recommendations.

1. **What are the average churn rates for customers with RFM scores from 1-5?**

   ● Our analysis found that customers with low recency, frequency, and monetary scores had up to a **59.83% churn rate**. Customers with a recency score of 3 but low frequency and monetary scores had up to a **59.21% churn rate**. We should focus on increasing customer engagement by encouraging the use of our products and services more often. We can accomplish this by offering special promotions and discounts, or by providing more personalized customer service.

2. **Who are our most common customer segments?**

   ● Our analysis found that the majority of our customers are happy with our service, but two segments, **Need Attention** and **Lost**, are at risk of churning. These segments **represent 25.40% of our customer base**. Based on the data we see between frequency and customer retention, I recommend that we focus on offering special promotions and discounts to help these customers engage and enjoy their card service more.

3. **What is each customer segments' average monthly RFM values?**

   ● Our analysis of monthly averages revealed patterns among the different segment groups. **The segments most likely to churn** were inactive for an average of 2.87 months, transacted 4.35 times per month, and spent $244.40 per month. **In contrast**, the rest of the group had an average of 2.20 months of inactivity, 5.83 transactions per month, and $390.02 in spending. I recommend personalized customer interactions which can lead to increased customer satisfaction, loyalty, and repeat business.

4. **What segments are most likely to churn, and what are their characterizations?**

   ● In summary, our analysis revealed that the following demographic and spending habits are common among customers who belong to the **Hibernating, Lost, Promising, Need Attention, and At risk** segments:

- **Female**: 56.93%
- **Middle-aged**: 70.56%
- **Low card utilization**: 80.16%
- **Income under $40,000**: 37.71%
- **Graduate-level education**: 30.34%

These findings can be used to develop targeted interventions to reduce customer churn.