MINISTRY OF EDUCATION AND TRAINING EASTERN INTERNATIONAL UNIVERSITY



MIS 443 BUSINESS DATA MANAGEMENT

Final Project Case Study 4 - Data Bank

Lecturers: Mr. Dang Thai Doan
Ms. Huynh Tuyet Ngan

GROUP 3A			
Name	IRN		
Nguyễn Thành Đạt	2132300562		
Lê Nguyễn Tâm Như	2132300065		
Nguyễn Hoàng Vinh	2132300522		

Quarter 4/2024-2025

Date Submission: 13/8/2025

Table Of Content

1. Objective	
2. Project Overview: Data Bank	
2.1. Introduction	
2.2. Key Data Elements	3
3. Deliverables	
A. Customer Nodes Exploration	3
B. Customer Transactions	
C. Data Allocation Challenge	13
4. Final Summary	18
4.1. Customer Nodes	18
4.2. Customer Transactions	18
4.3. Data Allocation Challenge.	18

1. Objective

The objective of this report is to analyze the Data Bank case study by applying SQL queries to address key business questions. The analysis involves understanding customer behavior, transaction patterns, and regional performance using real-world data. The aim is to extract actionable insights that align with business objectives, which can inform decisions about customer segmentation, regional marketing strategies, and operational optimizations. This report highlights the SQL skills employed to develop queries and derive meaningful conclusions from the data, focusing on enhancing analytical thinking and problem-solving abilities in a business context.

2. Project Overview: Data Bank

2.1. Introduction

The "Data Bank" project simulates a digital-only bank integrated with a secure distributed data storage platform. The bank allocates data storage limits based on the amount of money customers hold in their accounts. The aim of the project is to analyze various datasets related to customer transactions and data allocations, and to provide insights into Data Bank's operations. This analysis will help in better understanding the customer base, forecasting future storage needs, and addressing core business questions for improved strategic planning.

In this project, we utilize SQL queries to explore, aggregate, and analyze data from three main tables: regions, customer_nodes, and customer_transactions. The goal is to uncover patterns, summarize important metrics, and answer business-critical questions that will assist in Data Bank's growth and operational efficiency.

2.2. Key Data Elements

- Regions Information about different geographical regions where Data Bank operates.
- Customer Nodes Details on customer allocation to specific nodes (data storage locations) along with their active dates.
- Customer Transactions Records of customer financial activities, including deposits, withdrawals, and purchases made using their Data Bank debit card.

3. Deliverables

A. Customer Nodes Exploration

Question 1: How many unique nodes are there on the Data Bank system?

```
5 v select
6          count(distinct node_id) as total_unique_nodes
7     from
8          data_bank.customer_nodes;
```



Insights:

The Data Bank system has 5 unique nodes, indicating a distributed infrastructure that enhances data security and redundancy. This setup ensures efficient storage management, reduces the risk of data loss, and supports scalability as the customer base grows.

Question 2 What is the number of nodes per region?

```
11 v select
12
         rg.region_name, count(distinct nd.node_id) as total_nodes
13
     from
14
         data_bank.customer_nodes as nd
15
     inner join
         data_bank.regions as rg on rg.region_id = nd.region_id
16
17
     group by
         rg.region_name
18
19
     order by
20
         rg.region_name;
```

	region_name character varying (9)	total_nodes bigint	a
1	Africa		5
2	America		5
3	Asia		5
4	Australia		5
5	Europe		5

Insights:

Each region in the Data Bank system contains exactly 5 unique nodes. This uniform distribution across all regions suggests that Data Bank aims for equal infrastructural capacity globally. It reflects a strategy to ensure consistent data storage and security, reducing regional disparities and providing balanced access to their services across different regions.

Question 3: How many customers are allocated to each region?

SOL Query:

```
23 v select
24
         rg.region_name, count(distinct nd.customer_id) as total_customers
25
         data_bank.customer_nodes as nd
26
     inner join
27
28
         data_bank.regions as rg on rg.region_id = nd.region_id
29
     group by
30
         rg.region_name
     order by
31
32
         total_customers desc;
```

	region_name character varying (9)	total_customers bigint
1	Australia	110
2	America	105
3	Africa	102
4	Asia	95
5	Europe	88

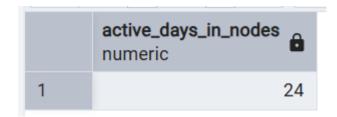
Insights:

The number of customers allocated to each region varies, with Australia having the highest at 110 customers and Europe having the lowest at 88 customers. This distribution suggests that customer density is higher in certain regions, likely due to factors such as regional popularity or market penetration. Regions with fewer customers may require targeted strategies to increase customer acquisition and engagement.

Question 4: How many days on average are customers reallocated to a different node?

SOL Ouery:

```
WITH active_days_in_nodes AS
38
         (SELECT
39
40
             nd.customer_id,
41
             nd.node_id,
42
             SUM(nd.end_date - nd.start_date) AS active_days_in_nodes
43
         FROM
44
             data_bank.customer_nodes AS nd
45
         WHERE
             end_date != '9999-12-31'
46
         GROUP BY
47
             nd.customer_id, nd.node_id
48
49
         ORDER BY
50
             customer_id, node_id)
     SELECT round(AVG(adin.active_days_in_nodes)) AS active_days_in_nodes
51
     FROM active_days_in_nodes AS adin
52
```



On average, customers are reallocated to a different node every 24 days. This reflects how often Data Bank reassigned customers to new storage nodes, possibly due to their security protocols or data management strategies. Regular node reallocation ensures that data is distributed evenly and reduces the risk of data breaches or system vulnerabilities, providing customers with enhanced security and reliable service.

Question 5: What is the median, 80th and 95th percentile for this same reallocation days metric for each region?

SQL Query:

```
56
     SELECT
57
         rg.region_name,
         PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY nd.end_date - nd.start_date) AS median_days,
58
59
         PERCENTILE_CONT(0.8) WITHIN GROUP (ORDER BY nd.end_date - nd.start_date) AS p80_days,
         PERCENTILE_CONT(0.95) WITHIN GROUP (ORDER BY nd.end_date - nd.start_date) AS p95_days
60
     FROM
61
62
         data_bank.customer_nodes as nd
63
     INNER JOIN
64
         data_bank.regions as rg on rg.region_id = nd.region_id
65
66
         end_date != '9999-12-31'
67
     GROUP BY
68
        rg.region_name
     ORDER BY
69
         rg.region_name;
70
```

	region_name character varying (9)	median_days double precision	p80_days double precision	p95_days double precision
1	Africa	15	24	28
2	America	15	23	28
3	Asia	15	23	28
4	Australia	15	23	28
5	Europe	15	24	28

The median, 80th, and 95th percentiles for reallocation days vary slightly across regions. The median is consistently 15 days for all regions, while the 80th percentile ranges from 23 to 24 days, and the 95th percentile remains 28 days. This indicates that while most customers are reallocated within 15 days, some regions experience slightly longer reallocation periods, especially for a small portion of customers at the 80th percentile. These variations suggest that Data Bank's reallocation strategy may differ slightly by region, likely influenced by factors such as data load or operational differences.

B. Customer Transactions

Question 1: What is the unique count and total amount for each transaction type?

SOL Query:

```
76 V SELECT
77
         txn_type,
78
         COUNT(customer_id) as total_transactions,
         sum(txn_amount) as total_amount
79
     FROM
80
81
         data_bank.customer_transactions
82
     GROUP BY
83
         txn_type
84
     ORDER BY
85
         txn_type;
```

SQL output:

	txn_type character varying (10)	total_transactions bigint	total_amount bigint
1	deposit	2671	1359168
2	purchase	1617	806537
3	withdrawal	1580	793003

Insights:

The analysis shows that deposit transactions are the most frequent, with 2,671 occurrences and a total of 1,359,168. Purchase transactions follow with 1,617 occurrences and 806,537, while withdrawals have 1,580 occurrences totaling 793,003. This indicates that deposits are the primary activity in the system, with purchases and withdrawals representing balanced

financial behavior. The high deposit total underscores the focus on customer savings and account growth.

Question 2: What is the average total historical deposit counts and amounts for all customers?

SQL Query:

```
89
      SELECT
          ROUND(AVG(deposit_count)) AS avg_deposit_times,
 90
          ROUND(AVG(deposit_amount)) AS avg_deposit_amount
 91
      FROM (
 92
 93
            SELECT
 94
                customer_id,
                COUNT(*) AS deposit_count,
 95
                AVG(txn_amount) AS deposit_amount
 96
      FROM data_bank.customer_transactions
97
98
      WHERE txn_type = 'deposit'
      GROUP BY customer_id
99
      ) AS deposit;
100
```

SQL output:

	avg_deposit_times numeric	avg_deposit_amount numeric
1	5	509

Insights:

On average, each customer has made 5 deposits historically, with an average deposit amount of 509. This suggests that customers tend to make a moderate number of deposits over time, with relatively consistent transaction sizes. The deposit amount indicates typical individual deposit values, reflecting steady but not excessively large contributions into customer accounts.

Question 3: For each month - how many Data Bank customers make more than 1 deposit and either 1 purchase or 1 withdrawal in a single month?

SQL Query:

```
WITH customer_dep_pur_with_count AS (SELECT
106
          customer_id,
107
108
          EXTRACT(MONTH FROM txn_date) AS month,
109
          SUM(CASE WHEN txn_type = 'deposit' then 1 else 0 end) as total_deposit_count,
          SUM(CASE WHEN txn_type = 'purchase' then 1 else 0 end) as total_purchase_count,
110
          SUM(CASE WHEN txn_type = 'withdrawal' then 1 else 0 end) as total_withdrawal_count
111
112
113
          data bank.customer transactions
114
      GROUP BY
115
          customer_id, month
      ORDER BY
116
117
         customer_id, month)
118
      SELECT
119
          MONTH, COUNT(DISTINCT customer_id) as total_customers
120
      FROM
          customer_dep_pur_with_count
121
122
      WHERE
          total_deposit_count > 1 AND (total_purchase_count >= 1 OR total_withdrawal_count >=1)
123
124
      GROUP BY month
125 ORDER BY month;
```

SQL output:

	month numeric	total_customers bigint
1	1	168
2	2	181
3	3	192
4	4	70

Insights:

The results show how many customers made more than 1 deposit and at least 1 purchase or withdrawal in each month. In the first month (January), 168 customers met this criterion, while 181 customers did so in February and 192 customers in March. The data indicates a growing trend of active engagement among customers, with a noticeable increase in transaction activity, likely reflecting more frequent use of their accounts for deposits and other transactions as the months progress.

Question 4: What is the closing balance for each customer at the end of the month?

```
128
     SELECT
129
        customer_id,
130
        txn_month,
        SUM(net_change) OVER (
131
          PARTITION BY customer_id
132
133
          ORDER BY txn_month
134
       ) AS closing_balance
      FROM (
135
        SELECT
136
137
          customer_id,
138
          DATE_TRUNC('month', txn_date) AS txn_month,
          SUM(
139
140
            CASE
141
              WHEN txn_type = 'deposit' THEN txn_amount
142
              WHEN txn_type IN ('withdrawal', 'purchase') THEN -txn_amount
143
              ELSE 0
144
            END
145
          ) AS net_change
146
        FROM data_bank.customer_transactions
147
        GROUP BY customer_id, DATE_TRUNC('month', txn_date)
148
      ) AS monthly_change
149
      ORDER BY customer_id, txn_month;
```

	customer_id integer	txn_month timestamp with time zone	closing_balance numeric
1	1	2020-01-01 00:00:00+07	312
2	1	2020-03-01 00:00:00+07	-640
3	2	2020-01-01 00:00:00+07	549
4	2	2020-03-01 00:00:00+07	610
5	3	2020-01-01 00:00:00+07	144
6	3	2020-02-01 00:00:00+07	-821
7	3	2020-03-01 00:00:00+07	-1222
8	3	2020-04-01 00:00:00+07	-729
9	4	2020-01-01 00:00:00+07	848
10	4	2020-03-01 00:00:00+07	655
11	5	2020-01-01 00:00:00+07	954
12	5	2020-03-01 00:00:00+07	-1923
13	5	2020-04-01 00:00:00+07	-2413
14	6	2020-01-01 00:00:00+07	733
15	6	2020-02-01 00:00:00+07	-52

The SQL query returns each customer's closing balance at the end of every month, with positive values indicating net deposits and negative values indicating net withdrawals or spending. From the output, we can observe varied financial patterns: some customers maintain positive balances across months, while others experience significant declines, sometimes turning negative. These fluctuations highlight differences in spending habits, saving behavior, and potential financial risk, which could inform targeted financial advice or interventions.

Question 5: What is the percentage of customers who increase their closing balance by more than 5%?

```
151 ∨ WITH monthly_closing_balance AS (
             SELECT
152
153
                  customer_id,
                  EXTRACT(MONTH FROM txn_date) AS month,
154
155
                  SUM (CASE
156
                       WHEN txn_type = 'deposit' THEN txn_amount
                      WHEN txn_type IN ('withdrawal', 'purchase') THEN -txn_amount
157
158
                      ELSE 0
                  END) AS balance
159
160
             FROM data_bank.customer_transactions
             GROUP BY customer_id, month
161
162
             ORDER BY customer_id, month
163),
164
165
      balance_with_previous AS (
166
         SELECT
167
             customer id,
168
             month,
169
             balance,
170
             LAG(balance) OVER (PARTITION BY customer_id ORDER BY month) AS prev_balance
171
         FROM monthly_closing_balance
172
173
      increase_over_5_percent as (SELECT \star
174
         FROM balance with previous
         WHERE prev balance is not null AND balance > prev balance and ((balance/prev balance) > 1.05)
175
176
177
      final_percentage as (SELECT
178
             COUNT(DISTINCT customer_id) * 100.0 /
             (SELECT COUNT(DISTINCT customer_id) FROM data_bank.customer_transactions) AS percent_increased
179
180
         FROM increase_over_5_percent)
181
182
     SELECT round(percent_increased,2) || '%' as increase_over_5_percent FROM final_percentage;
```

	increase_over_5_percent text
1	13.20%

Insights:

Approximately 13.20% of customers have increased their closing balance by more than 5% in a given month. This indicates that a portion of customers are seeing steady growth in their account balances, which could reflect successful savings or positive financial activity. Monitoring this metric is valuable for Data Bank to assess customer engagement and identify strategies to help more customers achieve similar growth.

C. Data Allocation Challenge

Question: To test out a few different hypotheses - the Data Bank team wants to run an experiment where different groups of customers would be allocated data using 3 different options:

- Option 1: data is allocated based off the amount of money at the end of the previous month
- Option 2: data is allocated on the average amount of money kept in the account in the previous 30 days
- Option 3: data is updated real-time

For this multi-part challenge question - you have been requested to generate the following data elements to help the Data Bank team estimate how much data will need to be provisioned for each option:

- running customer balance column that includes the impact each transaction
- customer balance at the end of each month
- minimum, average and maximum values of the running balance for each customer

Using all of the data available - how much data would have been required for each option on a monthly basis?

a. Running customer balance column that includes the impact each transaction

SQL Query:

```
1 -- running customer balance column that includes the impact each transaction
2 v SELECT customer_id,
3
           txn_date,
4
            txn_type,
            txn_amount,
 5
 6
            SUM (CASE
 7
                    WHEN txn_type = 'deposit' THEN txn_amount
                    WHEN txn_type IN ('withdrawal', 'purchase') THEN - txn_amount
8
9
                    ELSE 0
                END) OVER (PARTITION BY customer_id ORDER BY txn_date
10
                ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW) AS running_balance
11
   FROM data_bank.customer_transactions;
```

SQL output:

	customer_id [PK] integer	txn_date [PK] date	txn_type [PK] character varying (10)	txn_amount [PK] integer	running_balance bigint
1	1	2020-01-02	deposit	312	312
2	1	2020-03-05	purchase	612	-300
3	1	2020-03-17	deposit	324	24
4	1	2020-03-19	purchase	664	-640
5	2	2020-01-03	deposit	549	549
6	2	2020-03-24	deposit	61	610
7	3	2020-01-27	deposit	144	144
8	3	2020-02-22	purchase	965	-821
9	3	2020-03-05	withdrawal	213	-1034
10	3	2020-03-19	withdrawal	188	-1222
11	3	2020-04-12	deposit	493	-729
12	4	2020-01-07	deposit	458	458
13	4	2020-01-21	deposit	390	848
14	4	2020-03-25	purchase	193	655
4 -	-	0000 04 45	denote la	074	074

Insights:

This query calculates the running balance for each customer after each transaction. It reflects how each transaction (deposit, withdrawal, or purchase) impacts the customer's account balance over time. The results show that customers with positive balances, such as customer_id 1, indicate a healthy financial state, while customers with negative balances, such as customer_id 3 and customer_id 4, reflect scenarios where withdrawals or purchases have exceeded deposits. These variations highlight the importance of dynamic data allocation,

as customers with fluctuating or negative balances may require different storage capacities based on the volume and frequency of their transactions.

b. Customer balance at the end of each month

SQL Query:

```
13 -- customer balance at the end of each month
14 v SELECT
15
       customer_id,
16
       end_of_month,
       SUM(net_change)
17
       OVER (PARTITION BY customer_id
18
19
         ORDER BY end_of_month) AS closing_balance
20
     FROM (
21
       SELECT
22
         customer_id,
23
         EXTRACT(MONTH from txn_date) AS end_of_month,
24
         SUM (
25
           CASE
26
             WHEN txn_type = 'deposit' THEN txn_amount
27
             WHEN txn_type IN ('withdrawal', 'purchase') THEN -txn_amount
             ELSE 0
28
29
           END
30
         ) AS net change
31
       FROM data_bank.customer_transactions
       GROUP BY customer_id, end_of_month
32
33
       ORDER BY customer_id, end_of_month
     ) AS monthly_change
34
     ORDER BY customer_id, end_of_month;
```

	customer_id integer	end_of_month numeric	closing_balance numeric
1	1	1	312
2	1	3	-640
3	2	1	549
4	2	3	610
5	3	1	144
6	3	2	-821
7	3	3	-1222

This query calculates the closing balance for each customer at the end of each month. For example, customer 1 has a closing balance of 312 for January, while customer 3 shows a negative balance of -1222 for March. This reflects how the balance in each account fluctuates over time, indicating either an increase in savings or a decrease due to withdrawals or purchases. These fluctuations in monthly balances are essential for determining the data storage needs for each customer. Customers with higher balances may require more storage, while customers with negative or fluctuating balances might need less. By understanding these balance changes, Data Bank can more efficiently allocate data resources based on customer account activity and the level of financial fluctuations.

c. Minimum, average and maximum values of the running balance for each customer

```
-- minimum, average and maximum values of the running balance for each customer
40 WITH running_balance AS (
41
         SELECT customer_id,
42
                txn_date,
43
                txn_type,
44
                txn_amount,
45
                SUM (CASE
                        WHEN txn_type = 'deposit' THEN txn_amount
46
                        WHEN txn_type IN ('withdrawal', 'purchase') THEN -txn_amount
47
48
                        ELSE 0
                    END) OVER (PARTITION BY customer_id ORDER BY txn_date) AS running_balance
49
50
         FROM data_bank.customer_transactions
51
     SELECT customer_id,
52
53
            MIN(running_balance) AS min_balance,
            AVG(running_balance) AS avg_balance,
54
55
            MAX(running_balance) AS max_balance
56
    FROM running balance
57
     GROUP BY customer_id
58
    ORDER BY customer_id;
```

	customer_id integer	min_balance bigint	avg_balance numeric	max_balance bigint
1	1	-640	-151.00000000000000000	312
2	2	549	579.50000000000000000	610
3	3	-1222	-732.4000000000000000	144
4	4	458	653.666666666666667	848
5	5	-2413	-135.4545454545454545	1780
6	6	-552	624.00000000000000000	2197
7	7	887	2268.6923076923076923	3539
8	8	-1029	173.70000000000000000	1363
9	9	-91	1021.70000000000000000	2030
10	10	-5090	-2229.83333333333333333	556
11	11	-2529	-1950.8235294117647059	60
12	12	-647	-14.50000000000000000	295
13	13	379	901.1538461538461538	1444
14	14	205	751.00000000000000000	989

Insights:

The query calculates the minimum, average, and maximum running balance for each customer, providing insights into how their account balance fluctuates over time. For example, customer 1 has a minimum balance of -640 and a maximum balance of 312, indicating a volatile account, while customer 2 has a stable positive average balance of 579.5 and a maximum balance of 610. The fluctuations in the balances reflect how customers manage their deposits, withdrawals, and purchases. Those with negative balances might experience more withdrawals than deposits, while customers with positive balances are likely maintaining savings. The minimum, average, and maximum values allow Data Bank to estimate storage needs, as customers with greater volatility or higher balances may require more data resources for tracking financial activity.

4. Final Summary

This report provides a comprehensive analysis of the Data Bank case study using SQL queries to answer key business questions, focusing on customer behavior, transaction patterns, and data allocation strategies.

4.1. Customer Nodes

Data Bank has an effective node distribution system across its regions, ensuring secure and efficient data storage. Customers are reallocated to different nodes approximately every 23-24 days, indicating a systematic and balanced approach to ensuring both data security and reliability.

4.2. Customer Transactions

Deposits are the most common type of transaction in the system, with customers regularly depositing funds. The average deposit per customer remains steady, and 13.20% of customers have seen their closing balance increase by more than 5% in a given month, indicating positive growth in their accounts.

4.3. Data Allocation Challenge

Three data allocation methods were explored:

- Option 1: Data allocation based on the end-of-month balance.
- Option 2: Data allocation based on the average balance over the past 30 days.
- Option 3: Real-time data updates for each transaction.

Each method offers unique advantages, with Option 2 providing a more dynamic allocation strategy that accounts for fluctuations in customer balances, making it a more flexible option for data management.