

DBMSS

27435A Pauletto Jacopo

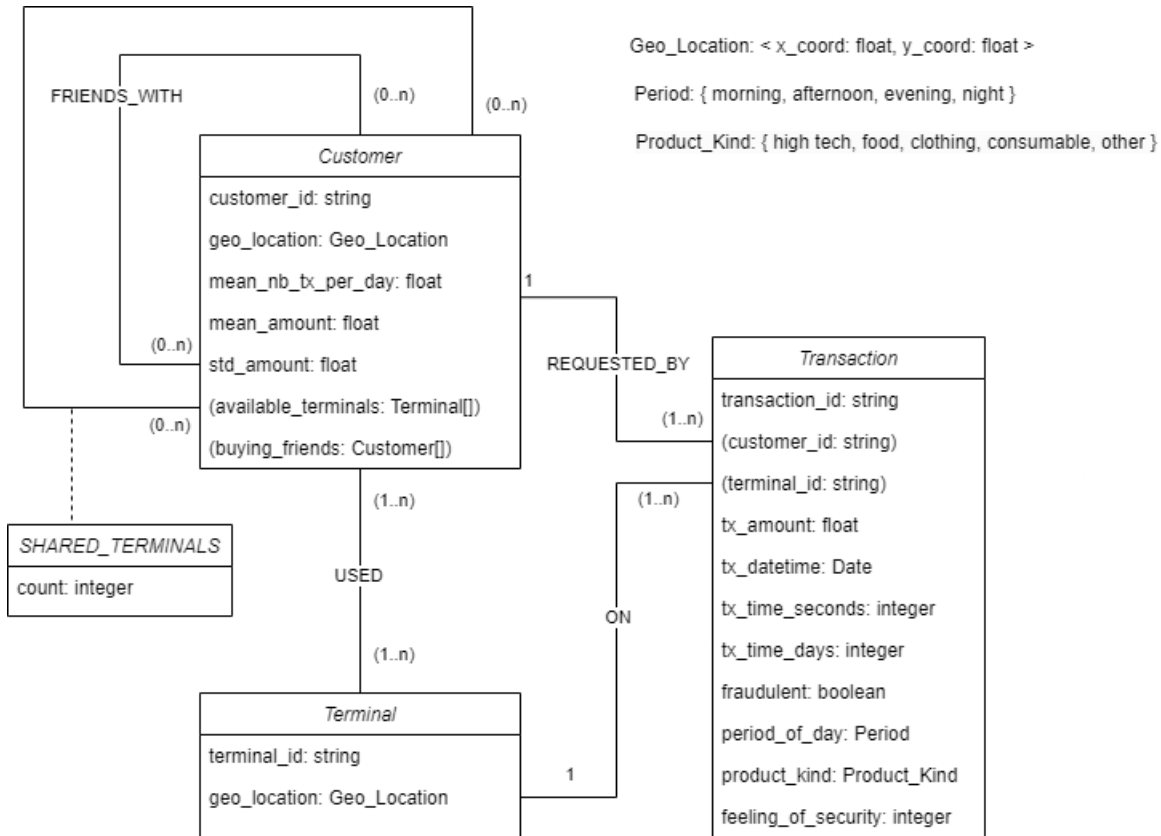
27436A Nava Stefano

2025

Contents

1	UML Class Diagram	2
1.1	Assumptions	2
1.2	Constraints	2
2	Logical Data Model	3
3	Import scripts	3
3.1	Queries for importing nodes	4
3.2	Queries for importing relationships	5
4	Operations scripts	6
4.1	Simple queries	6
4.2	Complex query C - Co-Customers	7
4.3	Complex query D - Extensions	8
5	Performances	9
5.1	Data import performances	9
5.2	Queries performances	10
5.3	Possible improvements	10

1 UML Class Diagram



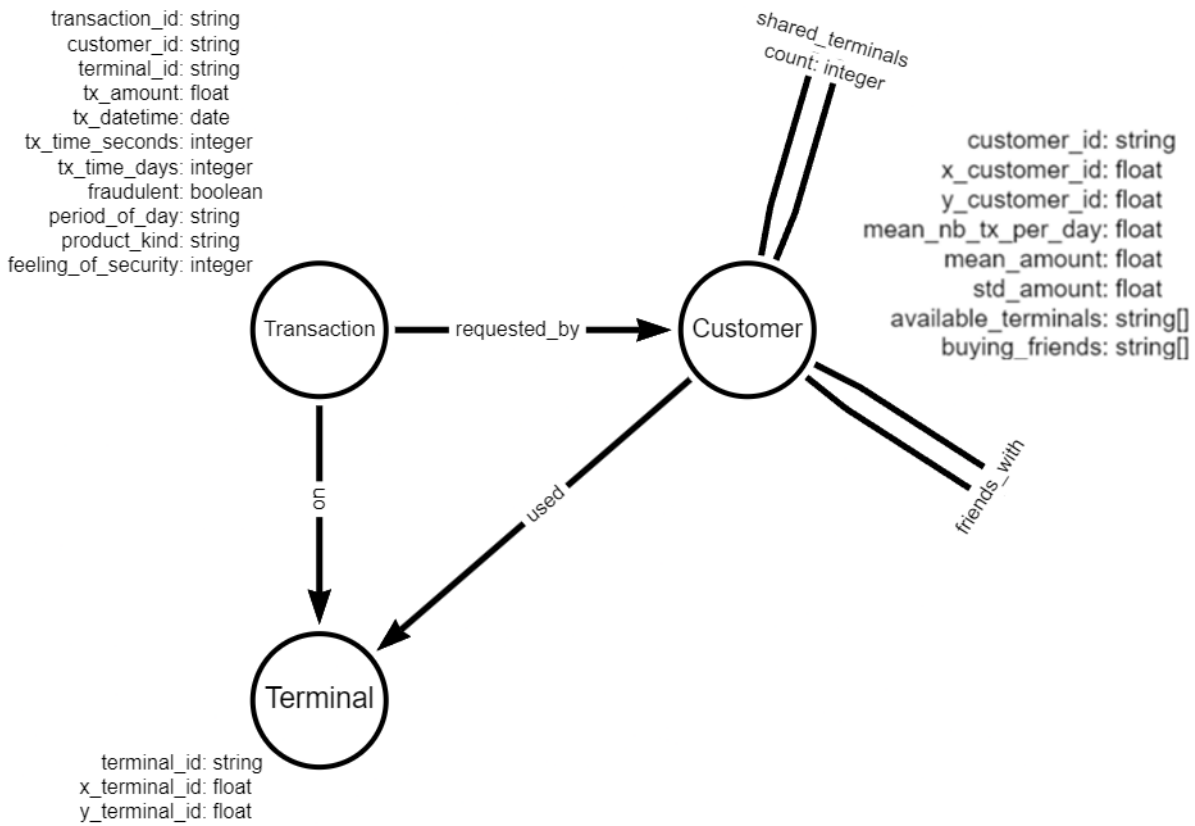
1.1 Assumptions

- We add an ID on Customer in order to univocally identify a customer.
- We add an ID on Terminal in order to univocally identify a terminal.
- We assume that there is a Geo_Location tuple with the x and y coordinates.
- We assume that the period of the day for a Transaction is one of the values contained in the Period enumerator.
- We assume that the product kind for a Transaction is one of the values contained in the Product Kind enumerator.

1.2 Constraints

- The x and y values of the Geo_Location must be between 0 and 100 (the min and max values of the grid).
- The mean_spending_amount and std_spending_amount of a customer can't be greater than the maximum amount_paid for the customer's transactions or lesser than 0.
- The feeling_of_security field of a transaction must be a value between 1 (low security) and 5 (high security).

2 Logical Data Model



We chose to develop this project on a Neo4J database because its main focus is on the relationships instead of the entities themselves. We created the **Transaction**, **Terminal** and **Customer** nodes and their respective properties according to the guidelines. The relationships between these three nodes have been generated based on the specifications and we added the two self-relationships on **Customer** following the queries' requests.

- **shared_terminals** identifies the number of terminals shared between two customers;
- **friends_with** identifies the concept of friendship based on the fact that two customers which made more than three transactions from the same terminal expressed a similar average feeling of security.

3 Import scripts

We built a generalized import function which considers the chosen dataset between the three options (50MB, 100MB, 200MB) and creates nodes and relationships based on the contents of the said dataset's CSV files.

```
1 with driver.session() as session:
2     print("Inizio caricamento nodi")
3     for query in nodes_queries:
4         session.execute_write(execute_query, query)
5
6     print("Fine caricamento nodi, inizio caricamento relazioni")
7
8     for query in relationships_queries:
9         session.execute_write(execute_query, query)
10
11     print("Fine caricamento relazioni")
12 driver.close()
```

The queries are collected in 2 arrays (one for the nodes, one for the relationships) which are iterated on by the generalized function.

3.1 Queries for importing nodes

We built the importing queries following the specifications of the Fraud Detection Handbook. Each one loads the CSV data (with its headers) from one of the 3 files and creates the related set of nodes.

```
1 LOAD CSV WITH HEADERS FROM 'file:///{{datasets_folders[dataset_number]}}/
  customers.csv' AS row
2 CREATE (:Customer {{
3     customer_id: row.CUSTOMER_ID,
4     x_customer_id: toFloat(row.x_customer_id),
5     y_customer_id: toFloat(row.y_customer_id),
6     mean_amount: toFloat(row.mean_amount),
7     std_amount: toFloat(row.std_amount),
8     mean_nb_tx_per_day: toFloat(row.mean_nb_tx_per_day),
9     available_terminals: row.available_terminals,
10    nb_terminals: toInteger(row.nb_terminals)
11 }});
```

```
1 LOAD CSV WITH HEADERS FROM 'file:///{{datasets_folders[dataset_number]}}/
  terminals.csv' AS row
2 CREATE (:Terminal {{
3     terminal_id: row.TERMINAL_ID,
4     x_terminal_id: toFloat(row.x_terminal_id),
5     y_terminal_id: toFloat(row.y_terminal_id)
6 }});
```

```
1 CALL apoc.periodic.iterate(
2     "LOAD CSV WITH HEADERS FROM 'file:///{{datasets_folders[dataset_number]}}/transactions.csv' AS row RETURN row",
3     "CREATE (:Transaction {{
4         transaction_id: row.TRANSACTION_ID,
5         customer_id: row.CUSTOMER_ID,
6         terminal_id: row.TERMINAL_ID,
7         tx_datetime: row.TX_DATETIME,
8         tx_amount: toFloat(row.TX_AMOUNT),
9         tx_time_seconds: toInteger(row.TX_TIME_SECONDS),
10        tx_time_days: toInteger(row.TX_TIME_DAYS)
11    }})",
12     {{ batchSize: 10000, batchMode: "BATCH" }}
13 );
```

For the transactions we used the APOC module in order to create the nodes in batches to avoid running out of memory with the 100mb and 200mb dataset.

3.2 Queries for importing relationships

Each query creates a relationship between nodes based on the requirements of the project.

As we did for the transactions nodes importing phase, we used the APOC module to perform batch imports in order to avoid running out of memory.

REQUESTED_BY

This relationship identifies the request of a transaction by a customer. The nodes are linked using the `customer_id`.

```
1 CALL apoc.periodic.iterate(  
2     "MATCH (c:Customer) MATCH (tr:Transaction) WHERE c.customer_id = tr.  
3         customer_id RETURN c, tr",  
4     "MERGE (tr)-[:REQUESTED_BY]->(c)",  
5     {{ batchSize: 10000, batchMode: "BATCH" }}  
);
```

ON

This relationship identifies the terminal on which the transaction is executed. The nodes are linked using the `terminal_id`.

```
1 CALL apoc.periodic.iterate(  
2     "MATCH (t:Terminal) MATCH (tr:Transaction) WHERE t.terminal_id = tr.  
3         terminal_id RETURN t, tr",  
4     "MERGE (tr)-[:ON]->(t)",  
5     {{ batchSize: 10000, batchMode: "BATCH" }}  
);
```

USED

This relationship identifies the fact that a terminal has been used by a customer, based on his transactions exploiting the `REQUESTED_BY` relationship.

```
1 CALL apoc.periodic.iterate(  
2     "MATCH (tr:Transaction)-[:REQUESTED_BY]->(c:Customer) MATCH (tr)-[:ON  
3         ]->(t:Terminal) RETURN tr, c, t",  
4     "MERGE (c)-[:USED]->(t)",  
5     {{ batchSize: 10000, batchMode: "BATCH" }}  
);
```

SHARED_TERMINALS

This relationship identifies the fact that a customer has shared one or more terminals with another customer, based on his previous `USED` relationships between himself and a terminal. The relationship keeps track of the number of shared terminals with the `count` property on the relationship itself. We decided to add this relationship following the 3a request.

```
1 CALL apoc.periodic.iterate(  
2     "MATCH (c1:Customer)-[:USED]->(t:Terminal)<-[:USED]-(c2:Customer) WHERE  
3         c1 <> c2 WITH c1, c2, COUNT(t) AS sharedTerminalCount RETURN c1,  
4         c2, sharedTerminalCount",  
5     "MERGE (c1)-[rel:SHARED_TERMINALS]-(c2) SET rel.count =  
6         sharedTerminalCount",  
7     {{ batchSize: 10000, batchMode: "BATCH" }}  
);
```

4 Operations scripts

We decided to organize the query files based on their complexity and necessary operations. For the queries "a", "b", "e" we used a generalized function because they don't need parameters and can be directly run. The query "c" requires building some clauses based on the selected number of levels for the co-customer relationship. The query "d" has to generate some random values in order to correctly populate the db so we treated it as a completely different operation.

4.1 Simple queries

a - Shared terminals

This query exploits the SHARED_TERMINALS relationship that we build during the data import in order to count the number of shared terminals between customers.

```
1 MATCH (X:Customer)-[st:SHARED_TERMINALS]-(Y:Customer)
2 WITH X, Y, st.count AS sharedTerminals
3 WHERE
4     sharedTerminals >= 3 AND
5     abs(Y.mean_amount - X.mean_amount) <= X.mean_amount * 0.1
6 RETURN
7     X.customer_id AS Customer_X,
8     Y.customer_id AS Customer_Y,
9     X.mean_amount AS SpendingAmount_X,
10    Y.mean_amount AS SpendingAmount_Y;
```

b - Monthly fraudulent transactions

For this query, we start by retrieving the current month and year and the previous month and year. Then we calculate the average amount spent for the previous month and then we compare it with the current month's amount for each transaction. If the transaction appears to be fraudulent, we set a `fraudulent` property as true.

```
1 WITH date() AS current_date
2 WITH current_date.year AS current_year, current_date.month AS current_month
3 WITH current_year, current_month, CASE
4     WHEN current_month = 1 THEN 12
5     ELSE current_month - 1
6 END AS previous_month,
7 CASE
8     WHEN current_month = 1 THEN current_year - 1
9     ELSE current_year
10 END AS previous_year
11 MATCH (t_prev:Transaction)-[:ON]->(term:Terminal)
12 WHERE toInteger(SPLIT(t_prev.tx_datetime, "-")[1]) = previous_month
13 AND toInteger(SPLIT(t_prev.tx_datetime, "-")[0]) = previous_year
14 WITH term, AVG(t_prev.tx_amount) AS avg_prev_month_import, current_month,
15     current_year
16 WITH term, avg_prev_month_import, avg_prev_month_import * 1.2 AS threshold,
17     current_month, current_year
18 MATCH (t_curr:Transaction)-[:ON]->(term)
19 WHERE toInteger(SPLIT(t_curr.tx_datetime, "-")[1]) = current_month
20 AND toInteger(SPLIT(t_curr.tx_datetime, "-")[0]) = current_year
21 AND t_curr.tx_amount > threshold
22 SET t_curr.fraudulent = true
23 RETURN t_curr
```

e - Period of day transactions

This query counts the transactions and checks for the fraudulent ones exploiting the added `period_of_day` property.

```
1 MATCH (t:Transaction)
2 WITH t.period_of_day AS period, COUNT(t) AS total_transactions,
3     AVG(CASE WHEN t.fraudulent = true THEN 1 ELSE 0 END) AS
4     avg_fraudulent_transactions
5 RETURN period, total_transactions, avg_fraudulent_transactions
ORDER BY period
```

4.2 Complex query C - Co-Customers

The code constructs a Cypher query dynamically. The process involves:

- Starting with a `MATCH` clause to find the initial customer node (`u1`).
- Iteratively adding `MATCH` clauses to connect customers over `num_levels`.
- Adding a `WHERE` clause to ensure that all connected customers in the chain are unique.
- Appending a `RETURN` statement to retrieve the ID of the customer at the final level.

```
1 c_query = f"MATCH (u1:Customer {{customer_id: {customer_id}}})"
2
3 for i in range(1, num_levels):
4     c_query += f" MATCH (u{i})-[:SHARES_TERMINAL]-(u{i+1}:Customer)"
5
6 where_clause = f" WHERE"
7 for i in range(1, num_levels + 1):
8     for j in range(i+1, num_levels + 1):
9         where_clause += f" u{i}.customer_id <> u{j}.customer_id AND"
10 where_clause = where_clause[:-4] + " RETURN DISTINCT u" + str(num_levels) +
11     ".customer_id AS co_customer"
12
13 c_query += where_clause
14 print(c_query)
```

Generated Query Example

Given the inputs:

- `customer_id = 1`
- `num_levels = 3`

The generated Cypher query would look like:

```
1 MATCH (u1:Customer {customer_id: 1})
2 MATCH (u1)-[:SHARES_TERMINAL]-(u2:Customer)
3 MATCH (u2)-[:SHARES_TERMINAL]-(u3:Customer)
4 WHERE u1.customer_id <> u2.customer_id AND
5       u1.customer_id <> u3.customer_id AND
6       u2.customer_id <> u3.customer_id
7 RETURN DISTINCT u3.customer_id AS co_customer
```

This query retrieves the IDs of all customers (`u3`) who are three levels away from the initial customer (`u1`) via shared terminals, ensuring no duplicates in the connections.

4.3 Complex query D - Extensions

The implementation consists of two core functions: `add_fields` and `create_friends`.

Extending Transactions

The `add_fields` function iteratively processes transactions and enriches them with additional attributes. It determines the period of the day based on the transaction's timestamp and assigns random values for `product_kind` and `feeling_of_security` as follows:

```
1 def add_fields(tx):
2     query = f"""
3         MATCH (t:Transaction{{transaction_id: $transaction_id}})
4         WITH t, toInteger(substring(SPLIT(t.tx_datetime, " ")[1], 0, 2)) as
5             intDateTime
6         SET t.period_of_day = CASE
7             WHEN intDateTime >= 6 AND
8                 intDateTime < 12 THEN "morning"
9             WHEN intDateTime >= 12 AND
10                intDateTime < 18 THEN "afternoon"
11            WHEN intDateTime >= 18 AND
12                intDateTime < 24 THEN "evening"
13            ELSE "night"
14        END,
15        t.product_kind = $product_kind,
16        t.feeling_of_security = $feeling_of_security
17    """
18    for record in tx.run("MATCH (t:Transaction) WHERE t.product_kind IS
19        NULL RETURN t.tx_datetime AS tx_datetime, t.transaction_id AS
20        transaction_id LIMIT 10000"):
21        product_kind = random.choice(["high tech", "food", "clothing", "
22            consumable", "other"])
23        feeling_of_security = random.randint(1, 5)
24        tx.run(
25            query,
26            product_kind=product_kind,
27            feeling_of_security=feeling_of_security,
28            transaction_id=record["transaction_id"]
29        )
```

Helper Function: `get_period_of_day`

This function determines the period of the day based on the hour extracted from the transaction's timestamp:

```
1 def get_period_of_day(tx_datetime):
2     hour = int(tx_datetime.split(" ")[1].split(":")[0])
3     if 6 <= hour < 12:
4         return "morning"
5     elif 12 <= hour < 18:
6         return "afternoon"
7     elif 18 <= hour < 24:
8         return "evening"
9     else:
10        return "night"
```


Creating Buying Friends Relationships

The `create_friends` function identifies customers who have executed more than three transactions at the same terminal and calculates their average `feeling_of_security`. If the average feelings of two customers differ by less than 1, a `FRIENDS_WITH` relationship is established between them:

```
1 def create_friends(tx):
2     query_avg_feeling = """
3         MATCH (c:Customer)-[:REQUESTED_BY]-(t:Transaction)-[:ON]->(terminal
4             :Terminal)
5         WITH c, terminal, AVG(t.security_feeling) AS avg_feeling, COUNT(t)
6             AS count
7         WHERE count > 3
8         RETURN c.customer_id AS customer_id, terminal.terminal_id AS
9             terminal_id, avg_feeling
10
11     """
12     results = tx.run(query_avg_feeling)
13     customers = list(results)
14
15     query_create_friends = """
16         MATCH (c1:Customer), (c2:Customer)
17         WHERE c1.customer_id = $customer1_id AND c2.customer_id =
18             $customer2_id
19         MERGE (c1)-[:FRIENDS_WITH {
20             terminal_id: $terminal_id,
21             avg_feeling_diff: $avg_feeling_diff
22         }]->(c2)
23     """
24
25     for c1 in customers:
26         for c2 in customers:
27             if c1["customer_id"] != c2["customer_id"] and c1["terminal_id"]
28                 == c2["terminal_id"]:
29                 diff = abs(c1["avg_feeling"] - c2["avg_feeling"])
30                 if diff < 1:
31                     tx.run(query_create_friends,
32                         customer1_id=c1["customer_id"],
33                         customer2_id=c2["customer_id"],
34                         terminal_id=c1["terminal_id"],
35                         avg_feeling_diff=diff)
```

5 Performances

5.1 Data import performances

Import type	50mb	100mb	200mb
Time to generate customer profiles table	0.057s	0.067s	0.052s
Time to generate terminal profiles table	0.0059s	0.0059s	0.0059s
Time to associate terminals to customers	1.8s	1.8s	1.8s
Time to generate transactions	58s	100s	190s

We measured these times while importing the different datasets in a Neo4J local instance. Since the number of customers and terminals doesn't vary between the various dataset, their import times are similar. On the other hand, transactions massively increase and their import times are different.

5.2 Queries performances

Query	50mb	100mb	200mb
a	2139 ms	2211 ms	2067 ms
b	2453 ms	2361 ms	2090 ms
c	2125 ms	2156 ms	2104 ms
d	69.42 h	352.8 h	1192.6 h
e	2100 ms	2096 ms	2059 ms

Observations:

- For the query `c` we used a co-customer degree of 3 with `customer_id` set to 0.
- For the query `d` we estimated the total time based on a subset of data (1000 for the 50mb dataset, 100 for the 100mb and 200mb datasets) because of its execution length.

5.3 Possible improvements

We observed we could add some indexes on `customer_id`, `terminal_id` and `transaction_id` regarding the execution of queries `c` and `d`:

- About the query `c`, the caller provides an ID for the customer for which the co-customer relationship is retrieved. If we create an index on the ID, the number of customers scanned in order to find that specific ID is lower than without an index.
- About the query `d`, we have several parameters used in more steps:
 - When adding the `period_of_day`, `product_kind` and `feeling_of_security` properties, we pass a `transaction_id` in order to set the new properties for the right transaction node.
 - When creating friendships, we have a similar situation in which we need to find the right customers ("friends") after verifying that they have a similar feeling of security. Using an index we can avoid scanning more customers than needed.

We also observed we could use an anchor node when creating friends exploiting `customer_id`, as we did in the step about adding the `period_of_day`, `product_kind` and `feeling_of_security` properties for the transactions. This would result in having a better starting point for the query execution, possibly reducing the times.