

Regarding the data model, the new nature of data requires a de facto new database paradigm -continuously evolving databases- where data can be both stable and volatile. Even though evolving databases can be implemented according to any approach, graph databases seem especially well suited here [1, 2]. Indeed, the natural way to process evolving graphs as streams of edges gives insights on how to proceed in order to maintain dynamic graph databases. Hence, we consider that a suitable data model is a continuously evolving data graph, a graph having persistent (stable) as well as non persistent (volatile) relations. Stable relations correspond to edges occurring in standard graph databases while volatile relations are edges arriving in data streams during a set time interval. Once this time interval is over, the relations are not longer valid so that there is no need to store them in the (stable) graph database. However, when required -as for further legal or auditing purposes- timestamped occurrences of volatile relations can be kept in a log file. Volatile relations induce subgraphs that exist only while the relations are still valid. Without loss of generality, in this work we consider property graphs (PG) [3, 4] as the basic reference data model. As an example, Figure 1a depicts part of a schema of a PG database where stable relations correspond to the data that a bank typically gathers on its issued cards, ATMs (Automated Teller Machines) network, etc. Volatile relations model the interaction between cards and ATM entities

In the context of our work we could see the data we are considering to be both static and streaming data, as we are considering a bank system application that contains all the information related to it on the cards, clients..., and that it is receiving the streaming of transactions that happens on it. More specifically, the static data can be thought of the classical bank database data, that is, the data a bank typically gathers on its issued cards, clients, accounts, ATMs.... Whereas as the streaming data we can consider the transactions the clients of the bank produce with their cards on ATMs, PoS... that reach the bank system. Therefore, due to this nature of the data, we consider a *continuously evolving database* paradigm, where data can be both stable and volatile. Even though evolving databases can be implemented according to any approach, graph databases seem especially well suited here [1, 2].

The property graph data model consists of two sub property graphs: a stable and a volatile property graph. On the one hand, the stable is composed of the static part of the data that a bank typically gathers such as information about its clients, cards, ATMs (Automated Teller Machines). On the other hand, the volatile property graph models the transaction operations, which defines the most frequent and reiterative kind of interaction between entities of the data model.

The main difference and the main reason for this separation is the semantics with which we intentionally define each of the subgraphs: the stable will be understood like a fixed static bank database, whereas the volatile will be understood as the data model to define the transactions, as continuous interactions between the entities of the model, which will not be permanently saved, but instead, only for a certain window of time under the mission of detecting anomalous bank operations. Note that we will only model the transaction interaction in the volatile

subgraph, only letting them occur here. This separation will allow us to have a really simple and light property graph schema single-centered on the transactions with the minimal needed information (mostly identifiers of the entities a transaction links) and another, the stable, acting as a traditional bank database schema, from which to obtain the information details of the entities.

Due to the confidential and private nature of bank data, it was impossible to find a real bank dataset nor a real bank data model. Therefore, we developed our own proposal of a bank database model.

0.1 Design of the Data Model

In what follows we describe the design of our data model as a Property Graph data model, divided into the stable and volatile property graphs.

0.1.1 Stable Property Graph

We propose a simplified stable data model where the defined entities, relations and properties modeling the bank database are reduced to the essential ones (see Figure 1). Although it is obvious that a real bank data model is way more complex than the one we propose, we believe that ours is relevant and representative enough and therefore sufficient for the purpose of our work.

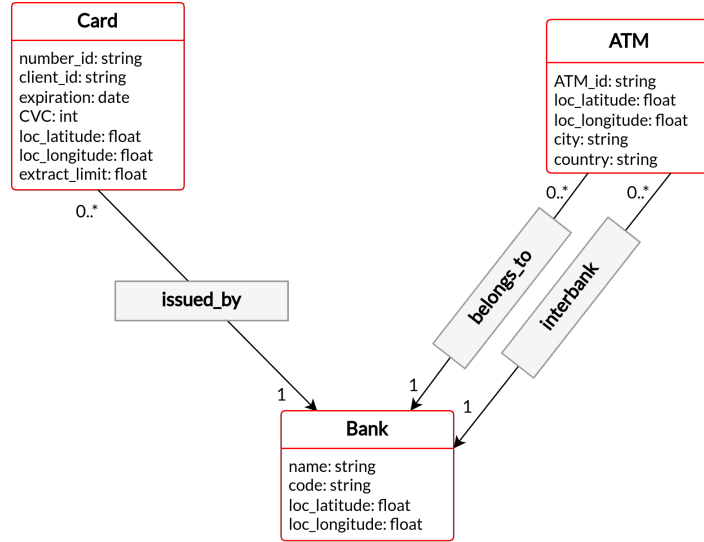


Figure 1: Definitive stable property graph model

Before obtaining this final model, a first model proposal was developed (see Figure 2), and then from it, we ended up reaching the final model version.

The first model idea was the initial attempt to capture the data that a bank system database typically gathers. It contains four entities: Bank, ATM, Client and Card with their respective properties, and the corresponding relationships between them.

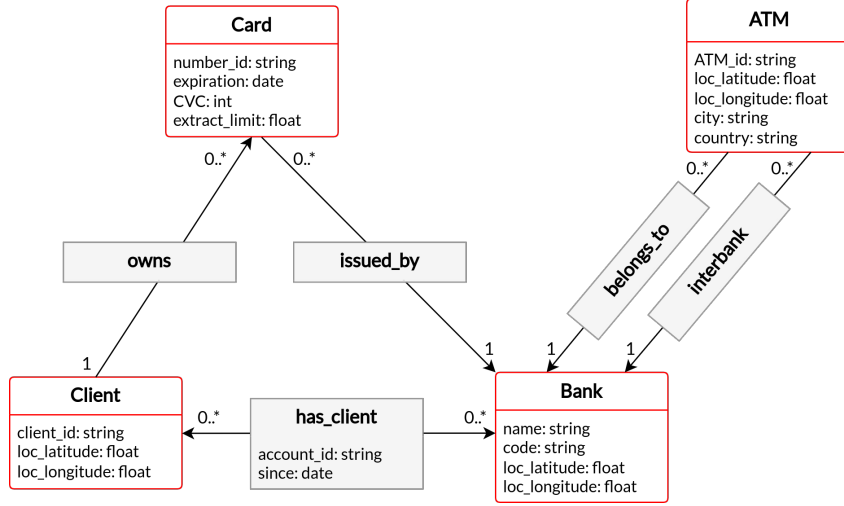


Figure 2: Initial bank stable property graph model

The relations are: a directed relationship from Client to Card **owns** representing that a client can own multiple credit cards and that a card is owned by a unique client, then a bidirectional relation **has_client** between Client and Bank; representing bank accounts of the clients in the bank. The relation between Card and Bank to represent that a card is **issued_by** the bank, and that the bank can have multiple cards issued. Finally, the relations **belongs_to** and **interbank** between the ATM and Bank entities, representing the two different kinds of ATMs depending on their relation with the bank; those ATMs owned and operated by the bank and those that, while not owned by the bank, are still accessible for the bank customers to perform transactions.

However, the final version of the model (Figure 1) was simplified to reduce it to the minimal needed entities. In particular, we decided to remove the Client entity and to merge it inside the Card entity. For this, all the Client properties were included in the Card entity. In the initial schema the Client entity was defined with three properties: the identifier of the client and the GPS coordinates representing the usual residence of the client. This change is done while preserving the restriction of a Card belonging to a unique client the same way it was previously done with the relation between Card and Client **owns** in the initial schema, which is now therefore removed.

Another derived consequence of this simplification is the removal of the other relation that the Client entity had with other entities: the **has_client** relation between Client and Bank, which was originally made with the intention of rep-

representing the bank accounts between clients and banks. Maintaining a bank account would imply having to consistently update the bank account state after each transaction of a client, complicating the model. Nevertheless, we eliminate the bank account relation, since its removal is considered negligible and at the same time helpful for the simplification of the model needed for the purposes of our work. However, for the sake of completeness the property *extract_limit* is introduced in the Card entity, representing a money amount limit a person can extract, which will be related with the amount of money a person owns. This will allow the detection of anomalies related with frequent or very high expenses. Other properties that are included with the purpose to allow the detection of some other kinds of anomalies are the GPS coordinates, which are added intentionally in the case of the ATM and Card entities, in the first case referring to the geolocation of each specific ATM and in the last case referring to each specific client address geolocation.

As a result, our stable property graph model contains three node entities: Bank, Card and ATM, and three relations: **issued_by** associating Card entities with the Bank entity, and **belongs_to** and **interbank** associating the ATM entities with the Bank entity.

The Bank entity represents the bank we are considering in our system. Its properties consist on the bank *name*, its identifier *code* and the location of the bank headquarters, expressed in terms of *latitude* and *longitude* coordinates, as seen in Table 1.

Name	Description and value
name	Bank name
code	Bank identifier code
loc_latitude	Bank headquarters GPS-location latitude
loc_longitude	Bank headquarters GPS-location longitude

Table 1: Bank node properties

Note that, from the beginning we were considering more than 1 bank entity. This lead to consider the creation of this entity, which now as only 1 bank is considered it may not be needed anymore, being able to reformulate and simplify the model. However, it is left since we considered it appropriate to be able to model the different kinds of ATMs a bank can have with different relation types instead of with different ATM types.

The ATM entity represents the Automated Teller Machines (ATM) that either belong to the bank’s network or that the bank can interact with. For the moment, this entity is understood as the classic ATM, however note that this entity could potentially be generalized to a Point Of Sale (POS), allowing a more general kind of interactions apart from the current Card-ATM interaction, where also online transactions could be included apart from the physical ones. We distinguish two different kinds of ATMs, depending on their relation with

the bank:

- Internal ATMs: ATMs owned and operated by the bank. They are fully integrated within the bank’s network. Modeled with the **belongs_to** relation.
- External ATMs: These ATMs, while not owned by the bank, are still accessible for the bank customers to perform transactions. Modeled with the **interbank** relation.

Both types of ATMs are considered to be of the same type of ATM node. Their difference is modeled as their relation with the bank instance: **belongs_to** for the internal ATMs and **interbank** for the external ATMs.

Name	Description and value
ATM_id	ATM unique identifier
loc_latitude	ATM GPS-location latitude
loc_longitude	ATM GPS-location longitude
city	ATM city location
country	ATM country location

Table 2: ATM node properties

The ATM node type properties consist on the ATM unique identifier *ATM_id*, its location, expressed in terms of *latitude* and *longitude* coordinates, and the *city* and *country* in which it is located, as seen in Table 2. Note that the last two properties are somehow redundant, considering that location coordinates are already included. In any case both properties are maintained since their inclusion provides a more explicit description of the location of the ATMs.

Finally, the Card node type represents the cards of the clients in the bank system. The Card node type properties, as depicted in Table 3, consist on the card unique identifier *number_id*, the associated client unique identifier *client_id*, as well as the coordinates of the associated client habitual residence address *loc_latitude* and *loc_longitude*. Additionally it contains the card validity expiration date *expiration*, the Card Verification Code, *CVC* and the *extract_limit* property, which represents the limit on the amount of money it can be extracted with the card on a single withdrawal.

Name	Description and value
<code>number_id</code>	Card unique identifier
<code>client_id</code>	Client unique identifier
<code>expiration</code>	Card validity expiration date
<code>CVC</code>	Card Verification Code
<code>extract_limit</code>	Card money amount extraction limit
<code>loc_latitude</code>	Client’s habitual address GPS-location latitude
<code>loc_longitude</code>	Client’s habitual address GPS-location longitude

Table 3: Card node properties

The client is completely anonymized in the system (no name, surname, age, or any other confidential details) by using only a *client_id*. Currently, *client_id* is included in the Card node type for completeness. However, it could be omitted for simplicity, as we assume a one-to-one relationship between card and client for the purposes of our work – each card is uniquely associated with a single client, and each client holds only one card. Thus, the *client_id* is not essential at this stage but is retained in case the database model is expanded to support clients with multiple cards or cards shared among different clients.

⇒? Include in the card properties the properties related with the gathered behavior for the card: *withdrawal_day*, *transfer_day*, *withdrawal_avg...* or just in the CSV to use them for the creation of the synthetic transactions, but do not store them in the stable bank database.

Finally, as some remarks to complete the description, note that for both the ATM and the Card entities we have the GPS coordinates information. In the first case referring to the geolocation of each specific ATM and in the last case referring to each specific client address geolocation. These are included since they are needed for some of the transactions potential frauds that we want to detect. Another property that was intentionally included with this purpose is the `extract_limit` property on the Card node entity, it represents the money amount limit a person can extract, which will be related with the amount of money a person owns. This will allow the detection of anomalies related with frequent or very high expenses.

0.1.2 Volatile Property Graph

The volatile subgraph model describes the most *variable* part of our model, the continuous interactions between the client’s cards and the ATMs. These interactions represent the transactions that are continuously occurring and arrive to our system as a continuous data stream. This subgraph model contains the minimal information needed to identify the Card and ATM entities – `number_id` and `ATM_id` Card and ATM identifiers – between which the interaction occurs, along with additional details related to the interaction. The proposed data model can be seen in the Figure 3. On it we define the Card and ATM node entities with only the identifier properties, `number_id` and `ATM_id`, respectively. These identifiers are enough to be able to recover, if needed, the whole information

about the specific Card or ATM entity in the stable subgraph. Finally we define the *interaction* relationship between the Card and the ATM nodes. The *interaction* relation contains as properties: **id** as the interaction unique identifier, **type** which describes the type of the interaction (withdrawal, deposit, balance inquiry or transfer), **amount** describing the amount of money involved in the interaction in the local currency considered, and finally, **start** and **end** which define the interaction *datetime* start and end moments, respectively.

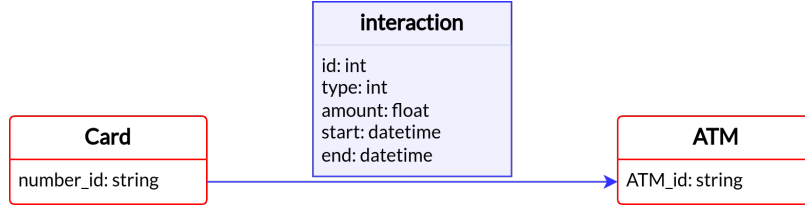


Figure 3: Volatile property graph model

As a whole, the proposed property graph data model is represented in Figure 4. where both the stable and volatile property subgraphs are merged to give a full view on the final property graph model.

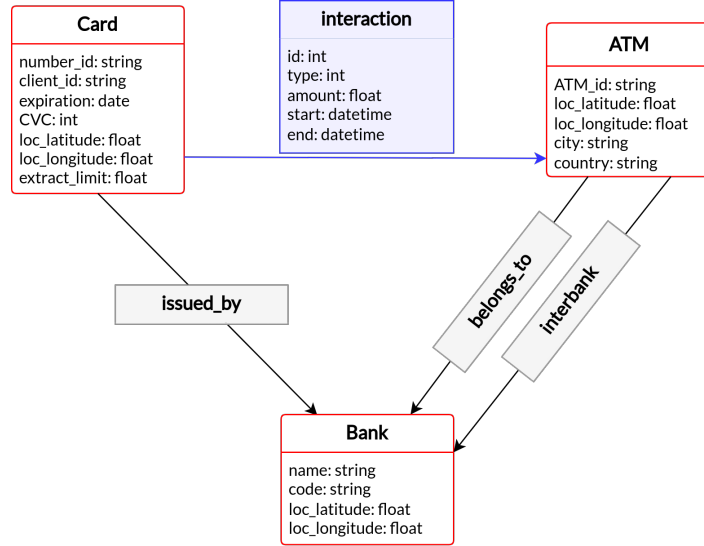


Figure 4: Complete property graph model

0.2 Creation of the Synthetic Dataset

As mentioned, given the confidential and private nature of bank data, it was not possible to find any real bank datasets. In this regard, and based on the described data model, we built a synthetic dataset consisting on a stable property

graph bank database and a set of volatile interactions (transactions). For this we used the *Wisabi Bank Dataset*¹ as a base to do the generation of the synthetic dataset. The *Wisabi Bank Dataset* is a fictional banking dataset that was made publicly available in the Kaggle platform. We considered it of interest as a base for the synthetic bank database that we wanted to develop. The interest to use this bank dataset as a base was mainly because of its size: it contains 8819 different customers, 50 different ATM locations and 2143838 transactions records of the different customers during a full year (2022). Additionally, it provides good heterogeneity on the different kind of transactions: withdrawals, deposits, balance inquiries and transfers. The main uses of this bank dataset are the obtention of a geographical distribution for the locations of our generated ATMs and the construction of a card/client *behavior*, which we will use for the generation of the synthetic transactions.

In particular, we divide the creation of the synthetic dataset in two. On the one hand on the creation of the stable bank database and on the other hand on the creation of the set of synthetic transactions and anomalous transactions that will conform the stream of data reaching our system.

Details of the *Wisabi Bank Dataset*

The *Wisabi Bank Dataset* consists on ten CSV tables. Five of them are of transaction records of five different states of Nigeria (Federal Capital Territory, Lagos, Kano, Enugu and Rivers State) that refers to transactions of cardholders in ATMs. In particular they contain 2143838 transactions records done during the year 2022, of which 350251 are in Enugu, 159652 in Federal Capital Territory, 458764 in Kano, 755073 in Lagos and 420098 in Rivers. Then, the rest of the tables are: a customers table ('customers_lookup') where the data of 8819 different cardholders is gathered, an ATM table ('atm_location lookup') with information of each of the 50 different locations of the ATMs, and then three remaining tables as complement of the previous ones ('calendar lookup', 'hour lookup' and 'transaction_type lookup') (tables summary).

0.2.1 Stable Bank Database

To do the generation of a stable bank database we provide the Python program `bankDataGenerator.py`, in which it is needed to enter the bank properties' values, as well as the parameters on the number of the bank ATMs (internal and external) and Cards: `n` and `m`, respectively. In what follows we give the details on the generation of the instances of our static database entities. For simplicity and to do it in a more stepwise manner, we are going to first create all the CSV data tables for the nodes and for the relations in the corresponding format and then we will populate the Neo4j GDB with them.

Bank Since a unique bank instance is considered, the values of the properties of the bank node are manually assigned, leaving them completely customisable.

¹Wisabi bank dataset on kaggle

For the bank, we will generate n ATM and m Card entities. Note that apart from the generation of the ATM and Card node types we will also need to generate the relationships between the ATM and Bank entities (`belongs_to` and `external`) and the Card and Bank entities (`issued_by`).

ATM We generate $n = n_internal + n_external$ ATMs, where $n_internal$ is the number of internal ATMs owned by the bank and $n_external$ is the number of external ATMs that are accesible to the bank. The generation of n ATMs for the bank is done following the geographical distribution of the locations of the ATMs in the *Wisabi Bank Dataset*. On this dataset there are 50 ATMs locations distributed along Nigerian cities. Note that for each of these ATMs locations, there can be more than one ATM. However, this is not taken into account and only one ATM per location is assumed for the distribution.

⇒ Put a plot of the distribution of the ATM locations

This distribution of the ATMs matches the relevance of the location in terms of its population, since the number of ATM locations is larger in the most populated Nigerian cities (30% of the ATM locations are in the city of Lagos, then the 20% in Kano...). Therefore, for the generation of the location of each of the n ATMs, the location/city of an ATM selected uniformly at random from the *Wisabi Bank Dataset* is assigned as *city* and *country*. Then, new random geolocation coordinates inside a bounding box of this city location are set as the *loc_latitude* and *loc_longitude* exact coordinates of the ATM.

Finally, as the ATM unique identifier *ATM_id* it is assigned a different code depending on the ATM internal or external category:

$$ATM_id = \begin{cases} bank_code + "-" + i & 0 \leq i < n_internal \text{ if internal ATM} \\ EXT + "-" + i & 0 \leq i < n_external \text{ if external ATM} \end{cases}$$

Card We generate a total of m cards that the bank manages, for each of them the assignment of the different properties is done as follows:

- Card and client identifiers:

$$\begin{cases} number_id = c-bank_code-i & 0 \leq i < m \\ client_id = i \end{cases}$$

- **Expiration** and **CVC** properties: they are not relevant, could be empty value properties indeed or a same toy value for all the cards. For completeness the same values are given for all the cards: **Expiration** = 2050-01-17, **CVC** = 999.
- Client's habitual address location (*loc_latitude*, *loc_longitude*): two possible options were designed to define the client habitual residence address. In both cases they are random coordinates drawn from a bounding

box of a location/city. The difference is on how the selection of the location/city is done:

1. Wisabi customers selection: Take the city/location of the habitual ATM of a random selected *Wisabi* database customer. Note that in the *Wisabi Bank Dataset* customers contain an identifier of their usual ATM, more in particular, the dataset is designed in such a way that customers only perform operations in the same ATM. With this approach, we maintain the geographical distribution of the *Wisabi* customers.
 2. Generated ATMs selection: Take the city/location of a random ATM of the n generated ATMs. This method is the one utilized so far.
- **Behavior:** It contains relevant attributes that will be of special interest when performing the generation of the synthetic transactions of each of the cards. The defined *behavior* parameters are shown in Table 4.

Behavior parameter	Description
amount_avg_withdrawal	Withdrawal amount mean
amount_std_withdrawal	Withdrawal amount standard deviation
amount_avg_deposit	Deposit amount mean
amount_std_deposit	Deposit amount standard deviation
amount_avg_transfer	Transfer amount mean
amount_std_transfer	Transfer amount standard deviation
withdrawal_day	Average number of withdrawal operations per day
deposit_day	Average number of deposit operations per day
transfer_day	Average number of transfer operations per day
inquiry_day	Average number of inquiry operations per day

Table 4: *Behavior* parameters

For each card, its *behavior* parameters are gathered from the transactions record of a randomly selected customer on the *Wisabi Bank Dataset*, from which we can access the transactions record of 8819 different customers for one year time interval. On it, there are four different types of operations that a customer can perform: withdrawal, deposit, balance inquiry and transaction. The parameters for the *behavior* gather information about these four different types of operations. Note that all these *behavior* parameters are added as additional fields of the CSV generated card instances, so, as mentioned, they can later be utilized for the generation of the synthetic transactions.

Another possible way to assign the *behavior* parameters could be the assignation of the same behavior to all of the card instances. However, this method will provide less variability in the generation of the synthetic

transactions than the aforementioned method. Nevertheless, other tailored generation methods to generate different *behavior* for each the cards could also be considered to similarly obtain this variability.

- `extract_limit: amount_avg_withdrawal * 5` Other possible ways could be chosen for assigning a value to this property.

0.2.2 Transactions Set

The transaction set constitutes the simulated data stream continuously arriving into our system. Each transaction of the set has the form of a *interaction* edge/relation from the volatile subgraph (see 0.1.2) matching a Card with an ATM entity from the stable bank database.

Regarding the purposes of our work, where we want our system to detect the presence of anomalous ATM transactions, we divide the transaction set in two subsets: the regular transaction set and the anomalous transaction set. The regular transaction set consists of *ordinary/correct* transactions, whereas the anomalous transaction set consists of *irregular/incorrect* transactions that are intentionally created to produce anomalous scenarios. The main objective under this division is the ability to have the control on the amount of anomalous ATM transactions that we generate so that later we are able to measure the efficiency of our system detecting them, as we have the exact number of anomalous transactions created to compare with, represented by size of the anomalous subset. Otherwise, if we were generating all the transactions together at the same time it would be more difficult to have the control on the amount of anomalous scenarios created, so that later, it will not be possible to measure the efficiency of our system detecting them, since we will not know their amount.

To do the generation of the synthetic set of transactions we provide the Python program `transactionGenerator.py`. On it we need to enter some parameters needed to customise the generation of the set of transactions.

Regular Transaction Set

- **By card:** Generation of the transactions for each of the cards independently. **We have control** to avoid anomalous scenarios when selecting the ATMs and distributing the transactions along time.
- **In general:** Linking ATM and client composing (card,ATM) pairs, and distributing these pairs along time according to a certain distribution. **No control / More difficult to control the possible derived anomalous scenarios produced among same card pairs.**

Therefore, the generation by card option it is considered to be the best so to be able to have the control on the possible anomalous scenarios for the generated transactions of each of the cards. Some ideas to explore:

- Selection of ATMs:
 - Neighborhood / Closed ATM subset.
 - Random walk. To do the selection of the sequence of ATMs for the generated transactions.
- Distribution of the transactions along time:
 - Uniform distribution.
 - Poisson process distribution.
- Other options:
 - Random walk for both the ATM and the transaction time selection, in the same algorithm together.

0.2.3 ATM closed subset + Uniform time distribution

- **ATM selection:** Closed ATM subset.
- **Time distribution:** Uniform distribution of the *num_tx* transactions for a card on a certain day (*num_tx* is drawn from a Poisson distribution of $\lambda = \text{withdrawal_day}$, based on the behavior of the *Wisabi* clients).

More detailed explanation follows.

The transaction generator is done to be able to generate transactions for each of the cards based on the gathered client transaction behavior of each of the cards for a customisable **d** number of days starting in a **start_date**.

For a card, the idea is to create a certain number of transactions per day, by linking the card to a certain ATM that is no farther than **max_distance** kms from the residence location of the client of the card. Also, we will limit the time distance between two consecutive transactions so that the final set of created transactions can not produce a potential fraud related with having two transactions in different ATM locations with an insufficient feasible time distance.

For each card:

- **ATM subset:** Create a subset of ATMs that are considered to be *usual* for the card client, so that they are all at a distance inferior or equal to **max_distance** kms to the residence location of the client of the card. We limit the size of this subset to be of **max_size_atm_subset**, so that we take only a maximum of **max_size_atm_subset** of the closest ATMs. In addition, among the ATMs in this subset, preference is given to those that are the closest to the residence location of the client and those that belong to the same bank company as the client's card. The transactions

generated for this card will be linked only to ATMs of this subset called Neighborhood.

$$\text{Neighborhood} = \{\text{ATM} \mid \text{dist}(\text{ATM}, \text{residence_loc}) \leq \text{max_distance}\}$$

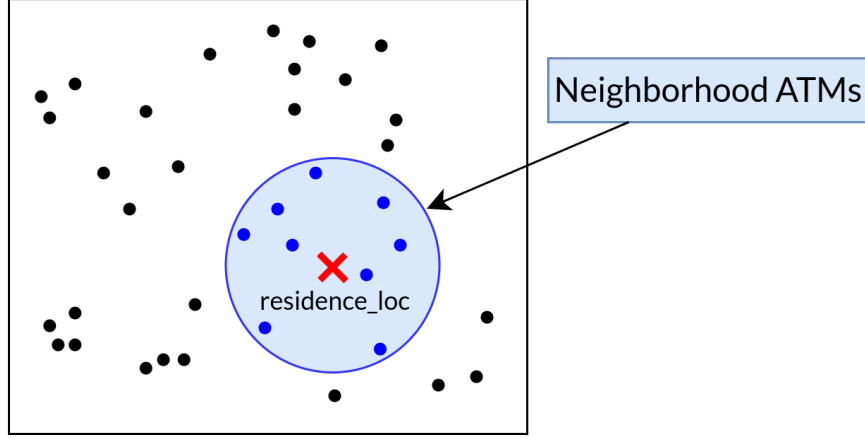


Figure 5: Neighborhood ATM subset

- **t_min**: Minimum threshold time between any two consecutive transactions of the client. That is, the minimum time distance between the end of a transaction and the start of the next consecutive transaction of a card.

$$t_{min} = \frac{2 * \text{max_distance}}{\text{max_speed}}$$

For the calculation of this minimum time distance:

- **max_distance**: two options:
 - * In general: taking $2 * \text{max_distance}$ kms as the upper bound on the maximum distance between 2 ATMs of the ATM subset, set the **t_min** to be the time needed to traverse that distance at a selected **max_speed**. (* So far done like this).

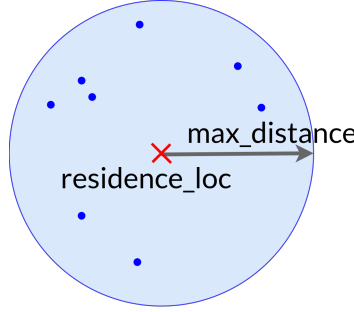


Figure 6: `max_distance` upper bound for the calculation of `t_min`

- * Specifically: Taking the specific maximum distance between all the ATM pairs in the ATM subset to get the `t_min`.
- `max_speed`: maximum speed at which it is possible to travel (by any possible means of transport) between any pair of ATMs \in `Neighborhood`.
- For each day generate `num_tx` transactions, random number drawn from a Poisson distribution of $\lambda = \text{withdrawal_day}$.
- Distribution of the `num_tx` transaction times (`transaction_start`), doing a uniform distribution of the `num_tx` for each of the days:

We create an ordered list of `num_tx` start moments in seconds in a day (in the range $[\text{t_min}/2, 86400 - (\text{t_min}/2) - \text{max_duration}]$) so that all of them are at a minimum time distance of $\text{t_min} + \text{max_duration}$. See Figure 7. *max_duration* limits the maximum duration of an ordinary transaction. For the moment it was set to be of 10 minutes (600s).

Note that the interval bounds for the transaction start moments are designed in such a way that the `t_min` minimum time distance between the end of a transaction and the start of the next transaction is also respected between transactions belonging to different consecutive days. See Figure 8. *Note that this way of generation (day by day) we are not allowing transactions to occur in the interval marked among the two red dotted lines of the Figure 8. NOTE/TODO: This could be fixed by looking to the start time of the last transaction of the previous day...* The start moments are therefore taken to define the `transaction_start` of each of the transactions of that day. `transaction_end` is assigned a shifted time difference with the respective `transaction_start`, in particular the difference is drawn from a normal distribution $\mathcal{N}(300, 120)$ that defines the duration of a transaction to be of mean of 5 minutes (300s) and a standard deviation of 2 minutes (120s), bounding it to be of a maximum time of `max_duration` of 10 minutes (600s) and setting it to the mean if the distribution sample was negative.

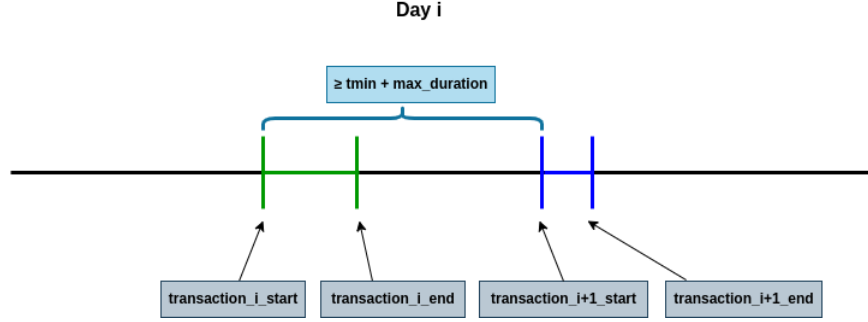


Figure 7: Time distance limit between two consecutive transactions

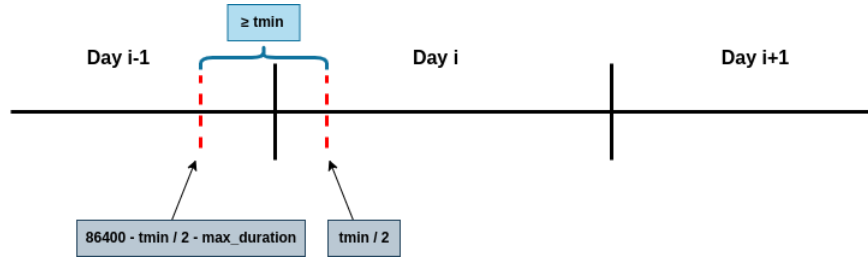


Figure 8: Time distance limit between two consecutive transactions between two days

Note that a limit was set on the duration of the transaction so that we can have the control avoiding time overlapping transactions, which will be producing irregular undesired fraud pattern alerts.

- **transaction_amount**: based on card behavior parameters, it is drawn from a normal distribution $\mathcal{N}(\text{amount_avg}, \text{amount_std})$. If negative amount, drawn from a uniform distribution $\mathcal{U}(0, \text{amount_avg} * 2)$.

Note that this approach is done with the focus on avoiding the production of potential fraud scenarios. By taking into account the previous generated transaction: both for the linked ATM of the new transaction and its transaction time (to avoid transactions that are overlapped or that come one directly after the other, since this may be fraudulent!).

0.2.4 ATM closed subset + Poisson Process

To do the generation of the synthetic set of transactions we provide the Python program `transactionGenerator.py`. On it we need to enter some parameters needed to customise the generation of the set of transactions.

Generate `num_tx` transactions for a selected period of time `t`. Distribution following a Poisson process distribution along $[0, t]$.

- $\lambda = \text{avg_tx}$ on a day for the client if $t = 24\text{h}$. Otherwise, decide a specific λ for the considered t .
- Inter-arrival times X_1, X_2, \dots are distributed following an exponential distribution: $X_i \sim \text{Exp}(\lambda)$. They represent the time elapsed between two of these events, e.g. X_2 . To do the generation of the synthetic set of transactions we provide the Python program `transactionGenerator.py`. On it we need to enter some parameters needed to customise the generation of the set of transactions. To avoid introducing anomalous undesired scenarios, we have to impose that: $X_i \geq \Delta_{i-1} + t_{min}, \forall X_i$, where:
 - X_i : time elapsed between the $i-1$ and i -th transaction.
 - Δ_{i-1} : time duration of the $i-1$ transaction. This duration will be upper bounded by Δ_{max} , which will be considered the maximum possible duration of a transaction.
 - t_{min} : minimum calculated time distance between any 2 consecutive transactions of the client.

With these interarrival times, we obtain the arrival times T_i . Note that they are not independent, in particular: $T_1 \leq T_2 \leq T_3 \dots$.

Therefore, to generate the Poisson process with rate λ :

1. Generate i.i.d. random variables X_1, X_2, \dots where $X_i \sim \text{Exp}(\lambda)$.
2. Obtain the arrival times as:
 - $T_1 = X_1$
 - $T_2 = X_1 + X_2$
 - $T_3 = X_1 + X_2 + X_3$
 - \dots

Note that, having imposed the previous we will have that:

$$\begin{cases} T_i = T_{i-1} + X_i \\ X_i \geq \Delta_{max} + t_{min} \end{cases} \quad \forall X_i \quad (1)$$

which implies that:

$$T_i \geq T_{i-1} + \Delta_{max} + t_{min}, \forall X_i \quad (2)$$

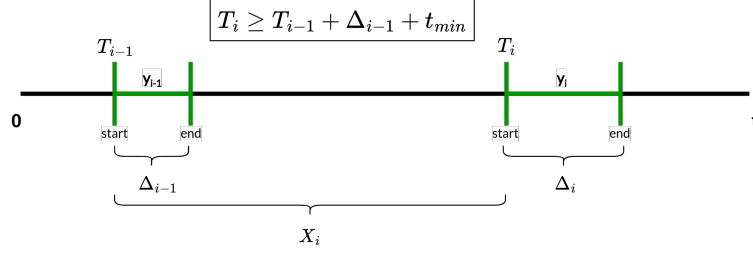


Figure 9: Schema of the interarrival and arrival event times on the poisson process.

References:

- Poisson process modeling - Theoric description
- Poisson process modeling - Python implementation
- Less formal explanation with Python code
- More theory
- Wikipedia entry

0.2.5 Random Walk / Markov Chain

Idea is to model the sequence of transactions for each of the clients as a markov chain through the network of ATMs. The network is a fully-connected graph with ATMs as nodes and edges are weighted with the distance between each pair of ATMs as weights.

The idea is to obtain the sequence of transactions for each client as a markov chain, in which, for each step a transaction is generated in that specific ATM node, at a certain datetime respecting the constraint of the minimum time distance with the previous transaction, so that no undesired anomalous fraud scenarios are produced.

Some considerations:

- Initially, we will compute the transition matrix obtaining all the respective transition probabilities. The probability of transition to another ATM-node will be inversely proportional to the distance to the considered ATM-node.
- Transition matrix P_t : containing, for each state, the probability of transitioning between states. For example, the entry i, j : $(P_t)_{i,j} = \mathcal{P}(X_{t+1} = j | X_t = i)$ contains the probability of transitioning to state j when being in state i . Note that since we assume a complete graph, and also the possibility to transition to the same state, all entries of this matrix will be distinct to 0: $(P_t)_{i,j} \neq 0 \forall i, j$.

- Finally, once P_t is computed, we perform the simulation to obtain the sequence of transactions for the specific card/client.

References:

- Markov Chains
 - MC - Theory
 - Simulation of Markov Chains - Theory
 - MC - Implementation example
 - MC - Implementation video example
- Random Walks
 - Theory:
 - * RWs: A Review of Algorithms and Applications
 - * RWs on Graphs
 - * RW - More theory
 - Implementation examples:
 - * Simple RW on a Graph - Implementation example
 - * A RW on a graph - Implementation
 - * Simulate Random Walks Through Markov Chain - Matlab Implementation example

0.3 Population of the Graph Database

0.3.1 Neo4j graph database creation

→ TODO: 0. Describe on how to set up the database → TODO: Explanation of the versions of both Neo4j instances used - local and UPC VM cluster.

Prior to the population of the Neo4j graph database, a Neo4j graph database instance needs to be created. This was done both locally and in a **Virtual Machine of the UPC cluster**.

Version: Neo4j 5.21.0 Community edition.

- Accessing it: by default it runs on localhost port 7474: `http://localhost:7474`.
Start the neo4j service locally by: `sudo systemctl start neo4j`
It can be also be accessed by the internal utility `cypher-shell`. Username: `neo4j` and password: `bisaurin`.

0.3.2 Neo4j graph database population - CSV to PG

Once we created the Neo4j graph database instance and all the CSV data files representing all the nodes and relations, we populate the Property Graph instance in Neo4j. Before performing the population of the GDB, we create uniqueness constraints on the properties of the nodes that we use as our *de facto* IDs for the ATM and Card IDs: `ATM_id` and `number_id`, respectively. The reason to do this is to avoid having duplicated nodes of these types with the same ID in the database. Therefore, as an example, when adding a new ATM node that has the same `ATM_id` as another ATM already existing in the database, we are aware of this and we do not let this insertion to happen. ID uniqueness constraints are created with the following cypher directives:

```
CREATE CONSTRAINT ATM_id IF NOT EXISTS
FOR (a:ATM) REQUIRE a.ATM_id IS UNIQUE

CREATE CONSTRAINT number_id IF NOT EXISTS
FOR (c:Card) REQUIRE c.number_id IS UNIQUE

CREATE CONSTRAINT code IF NOT EXISTS
FOR (b:Bank) REQUIRE b.code IS UNIQUE
```

Listing 1: Uniqueness ID constraints

Once we created the Neo4j graph database instance and all the CSV data files representing all the nodes and relations, we populate the Property Graph instance in Neo4j. For this, we propose two different methods. The first does it by directly importing the CSV files using the Cypher's `LOAD CSV` command, while the second method does it by parsing the CSV data and running the creation of the nodes and relationships using Cypher. Both methods can be found and employed using the `populate` module golang module. In this module we can find the two subdirectories where each of the methods can be run. In detail, the module tree structure is depicted in Figure 10. On it, the `cmd` subdirectory contains the scripts to run each of the populating methods: the first method script on `csvimport` and the second on the `cypherimport`, while the `internal` subdirectory is a library of the files with the specific functions used by these methods.

Prior to run any of these methods we need to first set up correctly the `.env` file located inside the desired method directory, where we have to define the corresponding Neo4j URI, username and password to access the Neo4j graph database instance.

- **Method 1: Cypher's `LOAD CSV`**

The Cypher's `LOAD CSV` clause allows to load CSV into Neo4j, creating the nodes and relations expressed on the CSV files (see *load-csv cypher manual*). To use it simply follow these steps:

1. Place all the CSVs (`atm.csv`, `bank.csv`, `card.csv`, `atm-bank-internal.csv`, `atm-bank-external.csv` and `card-bank.csv`) under the `/var/lib/neo4j/import`

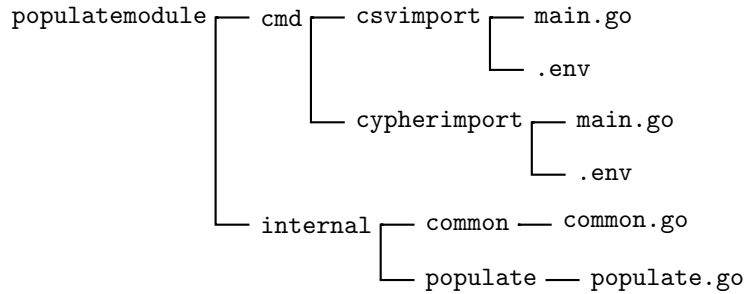


Figure 10: populatemodule file structure

directory of the machine containing the Neo4j graph database instance.

2. Run `$ go run populatemodule/cmd/csvimport/main.go`

Process description: Then the different CSV files containing all the data tables of our data set, were loaded into the GDB with the following cypher directives.

ATM (atm.csv)

```

LOAD CSV WITH HEADERS FROM 'file:///csv/atm.csv' AS row
MERGE (a:ATM {
  ATM_id: row.ATM_id,
  loc_latitude: toFloat(row.loc_latitude),
  loc_longitude: toFloat(row.loc_longitude),
  city: row.city,
  country: row.country
});

```

Listing 2: atm.csv

Some remarks:

- ATM is the node label, the rest are the properties of this kind of node.
- Latitude and longitude are stored as float values; note that they could also be stored as cypher *Point* data type. However for the moment it is left like this. In the future it could be converted when querying or directly be set as cypher point data type as property.

Bank (bank.csv)

```

LOAD CSV WITH HEADERS FROM 'file:///csv/bank.csv' AS row
MERGE (b:Bank {
    name: row.name,
    code: row.code,
    loc_latitude: toFloat(row.loc_latitude),
    loc_longitude: toFloat(row.loc_longitude)
});

```

Listing 3: bank.csv

Note that the code is stored as a string and not as an integer, since to make it more clear it was already generated as a string code name.

ATM-Bank relationships (atm-bank-internal.csv and atm-bank-external.csv)

```

LOAD CSV WITH HEADERS FROM 'file:///csv/atm-bank-internal.csv' AS row
MATCH (a:ATM {ATM_id: row.ATM_id})
MATCH (b:Bank {code: row.code})
MERGE (a)-[r:BELONGS_TO]->(b);

```

Listing 4: atm-bank-internal.csv

```

LOAD CSV WITH HEADERS FROM 'file:///csv/atm-bank-external.csv' AS row
MATCH (a:ATM {ATM_id: row.ATM_id})
MATCH (b:Bank {code: row.code})
MERGE (a)-[r:INTERBANK]->(b);

```

Listing 5: atm-bank-external.csv

Card (card.csv)

```

LOAD CSV WITH HEADERS FROM 'file:///csv/card.csv' AS row
MERGE (c:Card {
    number_id: row.number_id,
    client_id: row.client_id,
    expiration: date(row.expiration),
    CVC: toInteger(row.CVC),
    extract_limit: toFloat(row.extract_limit),
    loc_latitude: toFloat(row.loc_latitude),
    loc_longitude: toFloat(row.loc_longitude)});

```

Listing 6: card.csv

Notes:

- We do not include the fields that were generated to define the behavior of the card. They are only used for the generation of the transactions.
- expiration: set as *date* data type.

Card-Bank relationships (card-bank.csv)

```
LOAD CSV WITH HEADERS FROM 'file:///csv/card-bank.csv' AS
row
MATCH (c:Card {number_id: row.number_id})
MATCH (b:Bank {code: row.code})
MERGE (c)-[r:ISSUED_BY]->(b);
```

Listing 7: card-bank.csv

- Method 2:

→ **TODO: Describe the other population method**

1 Indexing

Useful for ensuring efficient lookups and obtaining a better performance as the database scales.

→ indexes will be created on those properties of the entities on which the lookups are going to be mostly performed; specifically in our case:

- Bank: `code` ?
- ATM: `ATM_id`
- Card: `number_id`

Why on these ones?

→ Basically the volatile relations / transactions only contain this information, which is the minimal information to define the transaction. This is the only information that the engine receives from a transaction, and it is the one used to retrieve additional information - the complete information details of the ATM and Card nodes on the complete stable bank database. Therefore these parameters/fields (look for the specific correct word on the PG world) are the ones used to retrieve / query the PG.

By indexing or applying a unique constraint on the node properties, queries related to these entities can be optimized, ensuring efficient lookups and better performance as the database scales.

From Neo4j documentation:

An index is a copy of specified primary data in a Neo4j database, such as nodes, relationships, or properties. The data stored in the index provides an access path to the data in the primary storage and allows users to evaluate query filters more efficiently (and, in some cases, semantically interpret query filters). In short, much like indexes in a book, their function in a Neo4j graph database is to make data retrieval more efficient.

Some references on indexing:

- Search-performance indexes
- The impact of indexes on query performance
- Create, show, and delete indexes

Okay... but before diving deeper...:

To Index or Not to Index?

When Neo4j creates an index, it creates a redundant copy of the data in the database. Therefore using an index will result in more disk space being utilized, plus slower writes to the disk.

Therefore, you need to weigh up these factors when deciding which data/properties to index.

Generally, it's a good idea to create an index when you know there's going to be a lot of data on certain nodes. Also, if you find queries are taking too long to return, adding an index may help.

From another tutorial on indexing in neo4j

→ Apparently, there are *Token lookup indexes* which are a default type of node indexes, from the Neo4j documentation: *"Two token lookup indexes are created by default when creating a Neo4j database (one node label lookup index and one relationship type lookup index). Only one node label and one relationship type lookup index can exist at the same time."* That is, two indexes are created by default, they are on the node label and on the relationship type.

→ There are different type of indexes depending on what we want to use them for, the type of property that they index...: see here for more details.

→ An example on when to use indexes: example case.

→ **Decision:** Create indexes on the ATM and Card node properties which serve as our node identifiers in practice: on `ATM.id` and on `number.id` node properties. The reason is that, so far, most of the queries are performed looking at a single node using its identifier; in particular to retrieve ATM nodes location coordinates to calculate the distance between two ATM nodes.

POSSIBLE IMPROVEMENT → make Neo4j obtain/calculate/return this distance as the result of a query.

Description of the creation on indexes

1. Select the type of index: range index VS text index... Some references on this election:
 - Use range index (the default), if we do not need to use `CONTAINS`, `ENDS WITH` or the value size is not above 8k.

2. Note that for the text index we (may) need to guarantee that the property type is `STRING`: create a single property type constraint.

References

- [1] Renzo Angles and Claudio Gutierrez. Survey of graph database models. *ACM Computing Surveys (CSUR)*, 40(1):1–39, 2008.
- [2] Rohit Kumar Kaliyar. Graph databases: A survey. In *International Conference on Computing, Communication & Automation*, pages 785–790. IEEE, 2015.