



POLITECNICO
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

Systems and Methods for Big and Unstructured Data Project

Author(s): **Samuele Forner 10726370**

Lara Longhi 11098675

Francesco Cavalieri 11020855

Group Number: **49**

Academic Year: 2024-2025

Contents

Contents	i
1 Introduction	1
2 Data Wrangling	3
3 Datasets	5
3.1 FoodMart dataset	5
3.2 Career Changes dataset	8
4 Queries	9
4.1 Neo4j Queries	9
4.2 MongoDB Queries	18
5 Application	37
List of Figures	41
List of Tables	43

1 | Introduction

In this report we will present an analysis of two totally unrelated datasets:

- The first one contains data on food stores across the United States. Since the data points (which were available in csv format) presented a lot of objects belonging to clearly distinct classes (like Customers, Region, ... see more details later) and several different relationships with their own attributes, the most natural choice for the Database technology to use was Neo4j, in order to represent the different entities as nodes and correctly model the relationships between them.

All the queries are somehow useful to improve the business in different ways, like finding the best spot for a new warehouse, analyzing the effectiveness of a promotion or even designing some new ones;

- The other one is about people's job experiences:

Do they like their current occupation? Are they satisfied about their salary? Are they willing to change their career path?

Also in this case we found the data in a csv format, where every different row of the file represented a different person. Moreover, there are little to no links between the data points.

For this two reasons the perfect way to work with these data was to store them in a json-like format, and there is no better Database technology to work with this kind of data format, as far as we know, than MongoDB.

In the next pages we will provide a better insight about the datasets, as well as how we manipulated them in order to make them suitable for our purpose and which query we performed on them.

2 | Data Wrangling

The MongoDB dataset presented all the values of the fields of the documents in string format, thus preventing us from performing any kind of mathematical operation or comparison on them. In order to perform useful queries, we had to convert all the fields that contained numbers, from strings to integer values.

This goal has been achieved by running the following javascript script on the MongoDB shell:

```

1 db.Carrer.updateMany(
2   {},
3   [
4     {
5       $set: {
6         // Convert any field of the document
7         allFields: {
8           $arrayToObject: {
9             $map: {
10              input: { $objectToArray: "$$ROOT" }, // Convert any
                  document in key-values' arrays'
11              as: "field",
12              in: {
13                k: "$$field.k", // Keep the key unchanged
14                v: {
15                  $cond: [
16                    {
17                      // Verify if the value is a numerical string
18                      $and: [
19                        { $eq: [{ $type: "$$field.v" }, "string"] },
20                        { $regexMatch: { input: "$$field.v", regex:
21                          /^[0-9]+$ / } }
22                      ]
23                    },
24                    { $toInt: "$$field.v" }, // Convert to number if the
                  value is a numerical string
25                    "$$field.v" // Keep the original values otherwise
26                  ]
27                }
28              }
29            }
30          }
31        }
32      },
33      {
34        $replaceRoot: { newRoot: "$allFields" } // Restore the documents
                  to the original format
35      }
36    ]
37  )

```


3 | Datasets

3.1. FoodMart dataset

A dataset used to manage the classic Food Marts and their sales, located in different regions, but with a more flexible and scalable NoSql approach, in this case Neo4j.

Nodes are linked among them with relationships, here below is provided a comprehensive list of those relationships:

Node Labels	Relationships
[Customer]	[IN_REGION, PURCHASED_BY]
[Region]	[IN_REGION]
[ProductSubCategory]	[IN_FAMILY, IN_CATEGORY, IN_DEPARTMENT]
[ProductFamily]	[IN_FAMILY]
[ProductCategory]	[IN_CATEGORY]
[Department]	[IN_DEPARTMENT]
[Product]	[IN_CATEGORY, FROM_BRAND, LINE_ITEM]
[Brand]	[FROM_BRAND]
[City]	[IN_REGION, IN_CITY]
[Store]	[OF_TYPE, IN_CITY, IN_STORE]
[StoreType]	[OF_TYPE]
[Year]	[HAS_YEAR]
[Month]	[HAS_YEAR, HAS_DAY]
[Date]	[HAS_DAY, NEXT_DAY, STARTS_ON, ENDS_ON, ON_DATE]
[Promotion]	[STARTS_ON, ENDS_ON, APPLIED_PROMOTION]
[Sale]	[ON_DATE, IN_STORE, PURCHASED_BY, LINE_ITEM, APPLIED_PROMOTION]

An "unique" constraint is applied to each node in order to avoid to add duplicated elements:

```

1 CREATE CONSTRAINT FOR (r:Region) REQUIRE r.id IS UNIQUE;
2 CREATE CONSTRAINT FOR (r:Region) REQUIRE r.name IS UNIQUE;
3 CREATE CONSTRAINT FOR (c:Customer) REQUIRE c.id IS UNIQUE;
4 CREATE CONSTRAINT FOR (c:Customer) REQUIRE c.account_num IS UNIQUE;
5 CREATE CONSTRAINT FOR (pf:ProductFamily) REQUIRE pf.name IS UNIQUE;
6 CREATE CONSTRAINT FOR (d:Department) REQUIRE d.name IS UNIQUE;
7 CREATE CONSTRAINT FOR (pc:ProductCategory) REQUIRE pc.name IS UNIQUE;
8 CREATE CONSTRAINT FOR (ps:ProductSubCategory) REQUIRE ps.id IS UNIQUE;
9 CREATE CONSTRAINT FOR (b:Brand) REQUIRE b.name IS UNIQUE;
10 CREATE CONSTRAINT FOR (p:Product) REQUIRE p.id IS UNIQUE;
11 CREATE CONSTRAINT FOR (st:StoreType) REQUIRE st.name IS UNIQUE;
12 CREATE CONSTRAINT FOR (s:Store) REQUIRE s.id IS UNIQUE;
13 CREATE CONSTRAINT FOR (y:Year) REQUIRE y.year IS UNIQUE;
14 CREATE CONSTRAINT FOR (m:Month) REQUIRE m.id IS UNIQUE;
15 CREATE CONSTRAINT FOR (d:Date) REQUIRE d.id IS UNIQUE;
16 CREATE CONSTRAINT FOR (d:Date) REQUIRE d.date IS UNIQUE;
17 CREATE CONSTRAINT FOR (d:Date) REQUIRE d.day IS UNIQUE;
18 CREATE CONSTRAINT FOR (p:Promotion) REQUIRE p.id IS UNIQUE;
19 CREATE CONSTRAINT FOR (s:Sale) REQUIRE s.id IS UNIQUE;

```

The desired structure is represented in following figure:

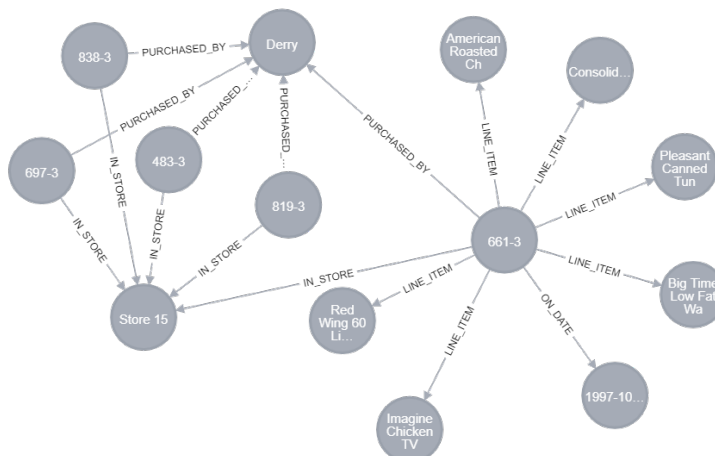


Figure 3.1: Example of nodes relationship

Table 3.1: Node Labels and Their Properties

Node Labels	Minimum Properties
[ProductFamily]	[<u>name</u>]
[ProductCategory]	[<u>name</u>]
[Department]	[<u>name</u>]
[Brand]	[<u>name</u>]
[Sale]	[<u>id</u>]
[StoreType]	[<u>name</u>]
[Year]	[<u>year</u>]
[ProductSubCategory]	[<u>id</u> , <u>name</u>]
[Month]	[<u>id</u> , <u>month</u>]
[City]	[<u>name</u> , <u>country</u> , <u>state</u>]
[Date]	[<u>id</u> , <u>date</u> , <u>day</u> , <u>day_of_month</u> , <u>day_of_week</u>]
[Promotion]	[<u>id</u> , <u>name</u> , <u>media_type</u> , <u>cost</u> , <u>start_date</u> , <u>end_date</u>]
[Region]	[<u>id</u> , <u>sales_city</u> , <u>sales_state_province</u> , <u>sales_district</u> , <u>sales_region</u> , <u>sales_country</u>]
[Product]	[<u>id</u> , <u>name</u> , <u>SKU</u> , <u>SRP</u> , <u>gross_weight</u> , <u>net_weight</u> , <u>recyclable_package</u> , <u>low_fat</u> , <u>units_per_case</u> , <u>cases_per_pallet</u> , <u>shelf_width</u> , <u>shelf_height</u> , <u>shelf_depth</u>]
[Store]	[<u>id</u> , <u>name</u> , <u>city</u> , <u>postal_code</u> , <u>country</u> , <u>state</u> , <u>number</u> , <u>street_address</u> , <u>coffee_bar</u> , <u>video_store</u> , <u>salad_bar</u> , <u>prepared_food</u> , <u>florist</u>]
[Customer]	[<u>id</u> , <u>account_num</u> , <u>lname</u> , <u>fname</u> , <u>address1</u> , <u>city</u> , <u>state_province</u> , <u>postal_code</u> , <u>country</u> , <u>phone1</u> , <u>phone2</u> , <u>birthdate</u> , <u>marital_status</u> , <u>yearly_income</u> , <u>gender</u> , <u>total_children</u> , <u>num_children_at_home</u> , <u>education</u> , <u>date_accnt_opened</u> , <u>member_card</u> , <u>occupation</u> , <u>houseowner</u> , <u>num_cars_owned</u> , <u>fullname</u>]

3.2. Career Changes dataset

A dataset with several details about different people's career, like their level of instruction, current occupation, previous experiences and other demographic factors.

Each person is represented by a single json document with the following structure (after the data wrangling process):

```

1 {
2     // field : tipe
3     _id: <ObjectId>,           //unique for every document
4     Field of Study: String,
5     Current Occupation: String,
6     Age: Int,
7     Gender: String,
8     Years of Experience: Int,
9     Education Level: String, //highest level of education
10    Industry Growth Rate: String, //"High", "Medium", "Low"
11    Job Satisfaction: Int, //A rating from 1 to 10
12    Work-Life Balance: Int, //A rating from 1 to 10
13    Job Opportunities: Int, //number of available job
                                //opportunities in the field
14    Salary: Int,
15    Job Security: Int, //A rating from 1 to 10
16    Career Change Interest: Int, //1 for yes, 0 for no
17    Skills Gap: Int, //A measure of how well the individual's
                                //current skills match their job requirements
18    Family Influence: String, //degree of influence the family
                                //has on their career choice
19    Mentorship Available: Int, //1 for yes, 0 for no
20    Certifications: Int, //1 for yes, 0 for no
21    Freelancing Experience: Int, //1 for yes, 0 for no
22    Geographic Mobility: Int, //1 for yes, 0 for no
23    Professional Networks: Int, //how strong the individual's
                                //professional network is (1-10)
24    Career Change Events: Int, //number of career changes
                                //made in the past
25    Technology Adoption: Int, //A measure of the comfort level
                                //with adopting new technologies (1-10)
26    Likely to Change Occupation: Int //1 for yes, 0 for no
27 }
```

4 | Queries

4.1. Neo4j Queries

Query.1 - Selling Hubs

Retrieve the cities with most sold units from their stores.

Useful to track where to deliver most of the items and detect the potential position of warehouses

```

1 MATCH (store:Store) - [:IN_CITY] - (citta:City)
2 WITH citta
3 MATCH (citta:City) - [ : IN_CITY] - (:Store) - [:IN_STORE] - (s:Sale) -[
    purchase:LINE_ITEM] - (:Product)
4 WITH citta, s , SUM(toInteger(purchase.quantity)) AS total_units
5 WITH citta , SUM(total_units) AS unit_per_city
6 RETURN citta,unit_per_city
7 ORDER BY unit_per_city DESC

```

City	Units per City
(:City {country: Mexico, name: Hidalgo, state: Zacatecas})	108,356
(:City {country: USA, name: Salem, state: OR})	80,762
(:City {country: USA, name: Tacoma, state: WA})	74,339
(:City {country: USA, name: Seattle, state: WA})	54,865
(:City {country: USA, name: Portland, state: OR})	54,521
(:City {country: USA, name: Los Angeles, state: CA})	52,630
(:City {country: USA, name: Spokane, state: WA})	52,477
(:City {country: USA, name: San Diego, state: CA})	52,413
...	...

Table 4.1: Neo4j query 1

Query.2 - Unforgettable day

Retrieve the most profitable day for each store

```

1 MATCH (store:Store) - [:IN_STORE] - (s:Sale) -[purchase:LINE_ITEM] - (:
   Product)
2 WITH store, s , SUM(toFloat(purchase.price - purchase.cost_price)/10000
   * purchase.quantity) AS total_sale_profit
3 WITH store, s, round(total_sale_profit * 100) / 100 AS
   rounded_total_sale_profit
4 MATCH (s:Sale) - [:ON_DATE] - (day:Date)
5 WITH store , day , SUM(rounded_total_sale_profit) AS total_profit_day
6 WITH store , day , round(total_profit_day * 100) / 100 AS
   rounded_total_profit_day
7 ORDER BY rounded_total_profit_day DESC
8 WITH store, collect({day: day, total: rounded_total_profit_day})[0] AS
   top_day_profit
9 RETURN store.name, top_day_profit.day.date AS day, top_day_profit.total
   AS max_total_per_day
10 ORDER BY store ASC

```

Store Name	Date	Max Total Per Day
Store 1	1998-12-20	3,336.88
Store 2	1998-12-23	196.65
Store 3	1998-11-07	3,164.02
Store 4	1998-01-02	2,845.41
Store 5	1998-12-14	160.41
Store 6	1998-11-05	2,775.15
Store 7	1998-11-05	3,859.08
Store 8	1998-12-03	4,254.03
Store 9	1998-06-03	1,404.08
Store 10	1998-08-29	2,880.64
Store 11	1998-12-23	3,285.34
Store 12	1998-09-16	4,383.52
Store 13	1997-07-27	9,282.36
Store 14	1998-12-29	181.26
.....

Table 4.2: Neo4j query 2

Query.3 - What's popular in Acapulco?

*Retrieve the 10 most popular couple of sold item in a sale at **Acapulco's** stores, useful in order of starting new promotions*

```

1 MATCH (store:Store{city:"Acapulco"}) - [:IN_STORE] - (s:Sale) , (p1:
    Product) - [:LINE_ITEM] - (s:Sale) - [:LINE_ITEM] - (p2:Product)
2 WITH p1, p2
3 WHERE id(p1) < id(p2)
4 WITH p1, p2, COUNT(*) AS pair_count
5 RETURN p1, p2, pair_count
6 ORDER BY pair_count DESC
7 LIMIT 10

```

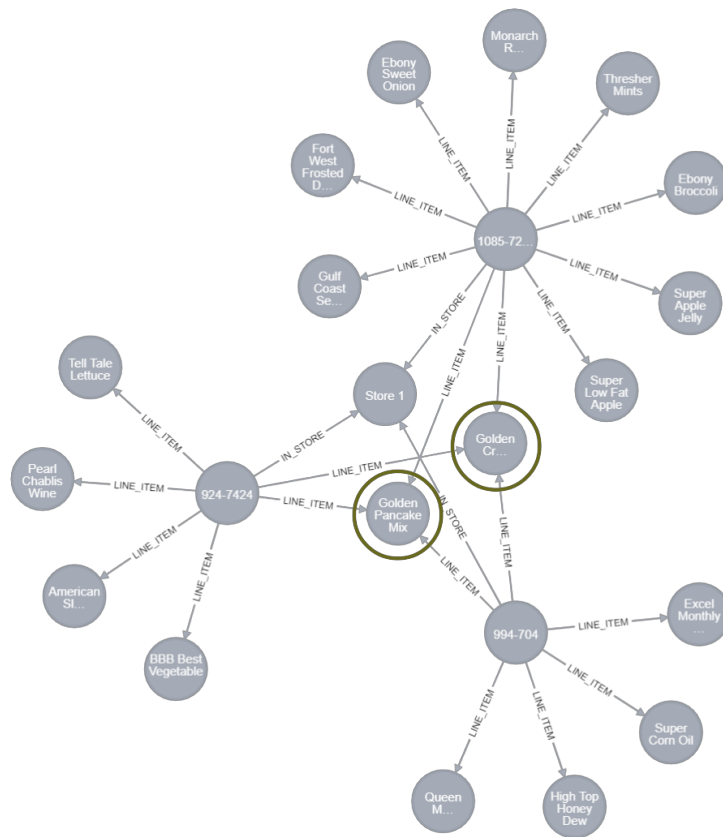


Figure 4.1: Example of sold items

Query.4 - Have you drank enough today?

*Trend of **Beers** sold over time at **USA**'s stores, useful to spot trendings and eventually to manage the pricing*

```

1 MATCH ( store : Store { country : "USA" }) - [: IN_STORE ] - ( s : Sale
    ) - [sell:LINE_ITEM] -(beer:Product) - [:IN_CATEGORY] - (beer_cat:
      ProductSubCategory{name:"Beer"}),
2 (s:Sale) - [:ON_DATE] - (day:Date)
3 WITH beer , day ,SUM(sell.quantity) as sold_beers
4 WITH day , SUM(sold_beers) as total_sold_beers
5 RETURN day.date AS date , day.day_of_week , total_sold_beers
6 ORDER BY day.day ASC

```

Date	Day of Week	Total Sold Beers
1997-01-02	Thursday	2
1997-01-03	Friday	3
1997-01-05	Sunday	8
1997-01-06	Monday	3
1997-01-07	Tuesday	10
1997-01-09	Thursday	4
1997-01-10	Friday	2
1997-01-11	Saturday	12
1997-01-13	Monday	9
.....

Table 4.3: Neo4j query 4

Query.5 - Best clients

For each year the top 3 customers who spent more.

Fidelity is the most valuable trait of our clients, so let's find who reward

```

1 MATCH (customer:Customer) - [:PURCHASED_BY] - (s:Sale) - [sell:LINE_ITEM
    ] - (p:Product)
2 WITH customer, s, SUM(toFloat(sell.price ) / 10000 * sell.quantity) AS
    money_spent
3 WITH customer, s, round(money_spent * 100) / 100 AS rounded_money_spent
4 MATCH (s) - [:ON_DATE] - (day:Date)
5 WITH customer, substring(day.date, 0, 4) AS year, SUM(toFloat(
    rounded_money_spent)) AS Y_spent
6 WITH customer, year, round(Y_spent *100) /100 AS Y_spent_rounded
7 ORDER BY Y_spent_rounded DESC
8 WITH year, COLLECT({customer: customer.fullname, money_spent:
    Y_spent_rounded})[0..3] AS customers_per_year
9 RETURN year , customers_per_year

```

Year	Customers Per Year
1998	{money_spent: 6757.93, customer: Scott Littleford}, {money_spent: 6503.93, customer: Isaiah Heymsfield}, {money_spent: 6259.62, customer: Craig Blackwell}
1997	{money_spent: 3954.59, customer: Mary Francis Benigar}, {money_spent: 3882.0, customer: Wildon Cameron}, {money_spent: 3540.47, customer: Ida Rodriguez}

Query.6 - Top Stores

Find the top 5 stores that earned more

```
1 MATCH (s:Store)<-[:IN_STORE]-(sa:Sale)-[l:LINE_ITEM]->(:Product)
2 WITH s, sum(l.price*l.quantity)/10000 AS sales
3 ORDER BY sales DESC
4 RETURN s.name, sales LIMIT 5
```

Store Name	Sales
Store 13	565,242
Store 17	527,800
Store 15	387,822
Store 11	387,257
Store 24	375,108

Query.7 - Cutting Time

Find, for each store, which department made less money, that's usefull in order to design new promotion for that store

```

1 MATCH (st:Store)<-[:IN_STORE]-(sa:Sale)-[l:LINE_ITEM]->(pr:Product)-[:
   IN_CATEGORY]->(sc:ProductSubCategory)-[:IN_DEPARTMENT]->(dep:
   Department)
2 WITH st, dep, sum(l.price*l.quantity)/10000 AS sales
3 ORDER BY st.name ASC, sales ASC
4 WITH st, collect([dep, sales]) as dep_sales
5 RETURN st.name, dep_sales[0][0].name AS department, dep_sales[0][1] AS
   profit

```

Store Name	Department	Profit
Store 1	Carousel	558
Store 2	Carousel	37
Store 3	Carousel	967
Store 4	Carousel	538
Store 5	Canned Products	18
Store 6	Carousel	935
Store 7	Carousel	1198
Store 8	Carousel	761
Store 9	Carousel	122
Store 10	Carousel	567
Store 11	Carousel	1336
Store 12	Carousel	647
Store 13	Carousel	1422
Store 14	Checkout	41
Store 15	Carousel	1063
Store 16	Carousel	1127
Store 17	Carousel	1609
Store 18	Carousel	173
Store 19	Carousel	679
Store 20	Carousel	157
Store 21	Carousel	876
Store 22	Carousel	52
Store 23	Carousel	266
Store 24	Carousel	823

Table 4.4: Neo4j query 7

Query.8 - Promotions Effectiveness

Find the number of purchase made with and without using a promotion, to analyse their effectiveness:

```

1 MATCH (s:Sale)-[:APPLIED_PROMOTION]->(pr:Promotion)
2 WITH collect(s) AS cs, COUNT(DISTINCT s) AS with_promo
3 MATCH (ns:Sale)
4 WHERE NOT ns IN cs
5 WITH with_promo, COUNT(DISTINCT ns) AS no_promo
6 RETURN with_promo, no_promo

```

With promo	Without promo
15060	43248

Query.9 - Best Customers

For every member card type, find the customer who spent more money, in order to reward them in future context:

```

1 MATCH (c:Customer)<-[:PURCHASED_BY]-(s:Sale)-[l:LINE_ITEM]->(p:Product)
2 WITH c.fullname AS name, c.member_card AS card, sum(l.price*l.quantity)
   /10000 AS rank
3 ORDER BY card, rank DESC
4 WITH card, collect([rank, name]) AS ranking
5 RETURN card, ranking[0][1] AS winner, ranking[0][0] AS rank

```

Card	Winner	Rank
Bronze	Ida Rodriguez	8337
Gold	James Horvat	7702
Normal	Kristin Miller	6247
Silver	Wildon Cameron	7818

Query.10 - How Many Customers

Find how many customers each type of store has

```
1 MATCH (c:Customer) <-[:PURCHASED_BY]-(sa:Sale)-[:IN_STORE]->(s:Store)-[:  
    OF_TYPE]->(st:StoreType)  
2 WITH st.name AS name, COUNT(DISTINCT c) AS customers  
3 ORDER BY customers  
4 RETURN name, customers
```

Name	Customers
Small Grocery	977
Mid-Size Grocery	1064
Gourmet Supermarket	1535
Deluxe Supermarket	2283
Supermarket	4767

4.2. MongoDB Queries

Query.1 - Satisfaction

Compute the Satisfaction index (average of each Job Satisfaction) for each job, then rank them

```
1 SMBUD_project.Carrer.aggregate([
2   {
3     $group: {
4       _id: "$Current Occupation",
5       Satisfaction: {
6         $avg: "$Job Satisfaction"
7       }
8     }
9   },
10  {
11    $match: {
12      _id: {
13        $ne: null
14      }
15    }
16  },
17  {
18    $project: {
19      Occupation: "$_id",
20      _id: 0,
21      Satisfaction: 1
22    }
23  },
24  {
25    $sort: {
26      Satisfaction: -1
27    }
28  }
29 ])
```

Output

```
1 [{
2   "Satisfaction": 5.558755464129597,
3   "Occupation": "Psychologist"
4 },
5 {
6   "Satisfaction": 5.556703762586116,
7   "Occupation": "Biologist"
8 },
9 {
10  "Satisfaction": 5.486328643482878,
11  "Occupation": "Economist"
12 },
13 {
14  "Satisfaction": 5.486003110419906,
15  "Occupation": "Business Analyst"
16 },
17 {
18  "Satisfaction": 5.483205501190161,
19  "Occupation": "Lawyer"
20 },
21 {
22  "Satisfaction": 5.473792394655704,
23  "Occupation": "Software Developer"
24 },
25 {
26  "Satisfaction": 5.4702396289616075,
27  "Occupation": "Artist"
28 },
29 {
30  "Satisfaction": 5.465002573340196,
31  "Occupation": "Teacher"
32 },
33 {
34  "Satisfaction": 5.46023166023166,
35  "Occupation": "Doctor"
36 },
37 {
38  "Satisfaction": 5.452835118892082,
39  "Occupation": "Mechanical Engineer"
40 }]
```

Query.2 - sameDegreeOccupation

Find how many people are currently working in the same field they have a degree in:

```
1 SMBUD_project.Carrer.aggregate([
2   {
3     $match: {
4       $or: [
5         {$and: [
6           { "Field of Study": "Psychology" },
7           { "Current Occupation": "Psychologist" }
8         ]
9       },
10      {$and: [
11        { "Field of Study": "Medicine" },
12        { "Current Occupation": "Doctor" }
13      ]
14    },
15    {$and: [
16      { "Field of Study": "Education" },
17      { "Current Occupation": "Teacher" }
18    ]
19  },
20  {$and: [
21    { "Field of Study": "Arts" },
22    { "Current Occupation": "Artist" }
23  ]
24 },
25 {$and: [
26   { "Field of Study": "Computer Science" },
27   { "Current Occupation": "Software Developer" }
28 ]
29 },
30 {$and: [
31   { "Field of Study": "Business" },
32   { "Current Occupation": "Business Analyst" }
33 ]
34 },
35 {$and: [
36   { "Field of Study": "Mechanical Engineering" },
37   { "Current Occupation": "Mechanical Engineer" }
38 ]
39 },
40
41 ]
```



```
42     {$and: [  
43         { "Field of Study": "Biology" },  
44         { "Current Occupation": "Biologist" }  
45     ]  
46 },  
47     {$and: [  
48         { "Field of Study": "Law" },  
49         { "Current Occupation": "Lawyer" }  
50     ]  
51 },  
52     {$and: [  
53         { "Field of Study": "Economics" },  
54         { "Current Occupation": "Economist" }  
55     ]  
56 }  
57 ]  
58 }  
59 },  
60 {$group: {  
61     _id: {  
62         Field: "$Field of Study",  
63         Occupation: "$Current Occupation"  
64     },  
65     tot: { $sum: 1 }  
66 }  
67 },  
68 {$project: {  
69     Field: "$_id.Field",  
70     Occupation: "$_id.Occupation",  
71     tot: 1,  
72     _id: 0  
73 }  
74 }  
75 ]);
```

Output - Partial

```
1 [{
2   "tot": 379,
3   "Field": "Law",
4   "Occupation": "Lawyer"
5 },
6 {
7   "tot": 389,
8   "Field": "Economics",
9   "Occupation": "Economist "
10 },
11 {
12   "tot": 395,
13   "Field": "Psychology",
14   "Occupation": "Psychologist "
15 },
16 {
17   "tot": 415,
18   "Field": "Computer Science",
19   "Occupation": "Software Developer "
20 },
21 {
22   "tot": 354,
23   "Field": "Business",
24   "Occupation": "Business Analyst "
25 },
26 {
27   "tot": 352,
28   "Field": "Arts",
29   "Occupation": "Artist "
30 },
31 {
32   "tot": 397,
33   "Field": "Biology",
34   "Occupation": "Biologist "
35 },
36 {
37   "tot": 405,
38   "Field": "Medicine",
39   "Occupation": "Doctor "
40 },
41 ...
42 }]
```

Query.3 - changeOccupation

For each job, find how many people would like to change carrer path and how many would not:

```
1
2 SMBUD_project.Carrer.aggregate([
3   {
4     $group: {
5       _id: "$Current Occupation",
6       change: {
7         $sum: {
8           $cond: [
9             {$eq: [
10              "$Likely to Change Occupation",1
11            ]},
12            1,0
13          ]}},
14       not_change: {
15         $sum: {
16           $cond: [
17             {$eq: [
18              "$Likely to Change Occupation",0
19            ]},
20            1,0
21          ]}}
22       },
23   {
24     $project:
25     {
26       Occupation: "$_id",
27       _id: 0,
28       change: 1,
29       not_change: 1
30     }
31   }
32 ])
```

Output - Partial

```
1 [{
2   "change": 2161,
3   "not_change": 1606,
4   "Occupation": "Economist"
5 },
6 {
7   "change": 2275,
8   "not_change": 1610,
9   "Occupation": "Doctor"
10 },
11 {
12   "change": 0,
13   "not_change": 0,
14   "Occupation": null
15 },
16 {
17   "change": 2213,
18   "not_change": 1676,
19   "Occupation": "Psychologist"
20 },
21 {
22   "change": 2212,
23   "not_change": 1680,
24   "Occupation": "Software Developer"
25 },
26 {
27   "change": 2232,
28   "not_change": 1654,
29   "Occupation": "Teacher"
30 },
31 {
32   "change": 2209,
33   "not_change": 1618,
34   "Occupation": "Mechanical Engineer"
35 },
36 {
37   "change": 2206,
38   "not_change": 1568,
39   "Occupation": "Biologist"
40 },
41 ...
42 ]
```

Query.4 - Education vs. Salary Analysis

*Find which education level have a better **average** salary*

```
1 SMBUD_project.Carrer.aggregate([
2   {
3     $group: {
4       _id: "$Education Level",
5       averageSalary: {
6         $avg: "$Salary"
7       }
8     }
9   },
10  {
11    $match: {
12      _id: {
13        $ne: null
14      }
15    }
16  },
17  {
18    $project: {
19      educationLevel: "$_id",
20      _id: 0,
21      averageSalary: 1
22    }
23  },
24  {
25    $sort: {
26      averageSalary: -1
27    }
28  }
29 ])
```

Output

```
1 [{
2   "averageSalary": 115348.93935604877,
3   "educationLevel": "Master's"
4 },
5 {
6   "averageSalary": 115320.96336686077,
7   "educationLevel": "Bachelor's"
8 },
9 {
10  "averageSalary": 114665.00756878388,
11  "educationLevel": "PhD"
12 },
13 {
14  "averageSalary": 114579.17316287289,
15  "educationLevel": "High School"
16 }]
```

Query.5 - Gender Distribution Analysis

This query calculates the percentage distribution of genders within each field of study, providing insights into gender representation across disciplines.

```

1 SMBUD_project.Carrer.aggregate([
2   {
3     $group: {
4       _id: {
5         fieldOfStudy: "$Field of Study",
6         gender: "$Gender"
7       },
8       count: {$sum: 1}
9     }
10  },
11  {
12    $group: {
13      _id: "$_id.fieldOfStudy",
14      total: {$sum: "$count"},
15      genders: {$push: {gender: "$_id.gender", count: "$count"}}
16    }
17  },
18  {
19    $match: {_id: {$ne: null}}
20  },
21  {
22    $project: {
23      fieldOfStudy: "$_id",
24      _id: 0,
25      genderPercentages: {
26        $map: {
27          input: "$genders",
28          as: "gender",
29          in: {
30            gender: "$$gender.gender",
31            percentage: {
32              $multiply: [
33                {$divide: ["$$gender.count", "$total"]},
34                100
35              ]
36            }
37          }
38        }
39      }
40    })

```

Output - Partial

```
1 [{
2   "fieldOfStudy": "Biology",
3   "genderPercentages": [
4     {
5       "gender": "Female",
6       "percentage": 50.14027033919919
7     },
8     {
9       "gender": "Male",
10      "percentage": 49.85972966080082
11    }
12  ]
13 },
14 {
15   "fieldOfStudy": "Medicine",
16   "genderPercentages": [
17     {
18       "gender": "Female",
19       "percentage": 50.797670296277545
20     },
21     {
22       "gender": "Male",
23       "percentage": 49.20232970372246
24     }
25  ]
26 },
27 {
28   "fieldOfStudy": "Mechanical Engineering",
29   "genderPercentages": [
30     {
31       "gender": "Female",
32       "percentage": 49.40629839958699
33     },
34     {
35       "gender": "Male",
36       "percentage": 50.59370160041301
37     }
38  ]
39 },
40 ...
41 ]
```


Query.6 - Relation between experience and job opportunities

Analyse the relationship between years of experience and job opportunities, ranking them by highest average opportunity

```
1 SMBUD_project.Carrer.aggregate([
2   {
3     $group: {
4       _id: {
5         yearsOfExperience: "$Years of Experience"
6       },
7       avgJobOpportunities: {$avg: "$Job Opportunities"}
8     }
9   },
10  {
11    $project: {
12      yearsOfExperience: "$_id.yearsOfExperience",
13      avgJobOpportunities: 1,
14      _id: 0
15    }
16  },
17  {
18    $sort: {avgJobOpportunities: -1}
19  }
20 ])
```

Output - Partial

```
1 [{
2   "avgJobOpportunities": 52.05675954592363,
3   "yearsOfExperience": 20
4 },
5 {
6   "avgJobOpportunities": 51.99171842650104,
7   "yearsOfExperience": 33
8 },
9 {
10  "avgJobOpportunities": 51.48221757322176,
11  "yearsOfExperience": 32
12 },
13 ...
14 ]
```

Query.7 - Total Certifications Analysis

Calculates the total number of certifications within each field of study and education level

```

1 SMBUD_project.Carrer.aggregate([
2   {
3     $group: {
4       _id: {
5         fieldOfStudy: "$Field of Study",
6         educationLevel: "$Education Level"
7       },
8       totalCertifications: {$sum: "$Certifications"}
9     }
10  },
11  {
12    $project: {
13      fieldOfStudy: "$_id.fieldOfStudy",
14      educationLevel: "$_id.educationLevel",
15      totalCertifications: 1,
16      _id: 0
17    }
18  },
19  {$sort: {totalCertifications: -1}}
20 ])
```

Output - Partial

```

1 [{
2   "totalCertifications": 318,
3   "fieldOfStudy": "Arts",
4   "educationLevel": "PhD"
5 },
6 {
7   "totalCertifications": 314,
8   "fieldOfStudy": "Medicine",
9   "educationLevel": "High School"
10 },
11 {
12   "totalCertifications": 310,
13   "fieldOfStudy": "Biology",
14   "educationLevel": "High School"
15 },
16 ...
17 ]
```

Query.8 - Skill gap Analysis

This query calculates the mean skill gap by the employee field of study, for each field of study, it also collects the cases where the skill gap is too high (greater or equal 5)

```

1 SMBUD_project.Carrer.aggregate([
2   {
3     $group: {
4       _id: "$Field of Study",
5       averageSkillGap: { $avg: "$Skills Gap" },
6       documents: {
7         $push: {
8           docId: "$_id",
9           skillsGap: "$Skills Gap",
10          currentOccupation: "$Current Occupation",
11          age: "$Age",
12          salary: "$Salary"
13        }
14      }
15    }
16  },
17  {
18    $project: {
19      fieldOfStudy: "$_id",
20      _id: 0,
21      averageSkillGap: 1,
22      documents: {
23        $filter: {
24          input: "$documents",
25          as: "doc",
26          cond: { $gte: ["$$doc.skillsGap", 5] }
27        }
28      },
29      howMany: { // same as before, just counting
30        $size: {
31          $filter: {
32            input: "$documents",
33            as: "doc",
34            cond: { $gte: ["$$doc.skillsGap", 5] }
35          }
36        }
37      }
38    }
39  }
40 ])
```

Output - Partial

```
1 [{
2   "averageSkillGap": 5.4410492157923205,
3   "fieldOfStudy": "Psychology",
4   "documents": [
5     {
6       "docId": {
7         "$oid": "6750823f8a036fa9e4893c1d"
8       },
9       "skillsGap": 9,
10      "currentOccupation": "Psychologist",
11      "age": 34,
12      "salary": 80448
13    },
14    ...
15  ],
16  "howMany": 2175
17 },
18 {
19   "averageSkillGap": 5.537269569863746,
20   "fieldOfStudy": "Economics",
21   "documents": [
22     {
23       "docId": {
24         "$oid": "6750823e8a036fa9e4893c09"
25       },
26       "skillsGap": 10,
27       "currentOccupation": "Business Analyst",
28       "age": 49,
29       "salary": 116672
30     },
31     ...
32   ],
33   "howMany": 2294
34 },
35 ...
36 ]
```

Query.9 - Gender distribution Analysis pt.2

For each pair of positions and education level, it calculates the distribution of men and women

```

1 SMBUD_project.Carrer.aggregate([
2   // Stage 1 group occupation,gender,educationLevel
3   {
4     $group: {
5       _id: {
6         occupation: "$Current Occupation",
7         gender: "$Gender",
8         educationLevel: "$Education Level"
9       },
10      count: { $sum: 1 }
11    }
12  },
13  // Stage 2: Explode and count female / males
14  {
15    $group: {
16      _id: {
17        occupation: "$_id.occupation",
18        educationLevel: "$_id.educationLevel"
19      },
20      countsByGender: {
21        $push: {
22          gender: "$_id.gender",
23          count: "$count"
24        }
25      }
26    }
27  },
28  // Stage 3: Pivot countsByGender into separate fields for men and
    women
29  {
30    $project: {
31      _id: 0, // Exclude default _id
32      occupation: "$_id.occupation",
33      educationLevel: "$_id.educationLevel",
34      menCount: {
35        $reduce: {
36          input: "$countsByGender",
37          initialValue: 0,
38          in: {
39            $cond: [

```

```

40         { $eq: ["$$this.gender", "Male"] }, // Check if gender is
           Male
41         { $add: ["$$value", "$$this.count"] },
42         "$$value" // Keep the value as-is if not Male
43     ]
44 }
45 }
46 },
47 womenCount: {
48     $reduce: {
49         input: "$countsByGender",
50         initialValue: 0,
51         in: {
52             $cond: [
53                 { $eq: ["$$this.gender", "Female"] }, // Check if gender
                    is Female
54                 { $add: ["$$value", "$$this.count"] },
55                 "$$value" // Keep the value as-is if not Female
56             ]
57         }
58     }
59 }
60 }
61 },
62 {$sort: { occupation: 1 }}
63 ])

```

Output - Partial

```

1  [{
2    "occupation": null,
3    "educationLevel": null,
4    "menCount": 0,
5    "womenCount": 0
6  },
7  {
8    "occupation": "Artist",
9    "educationLevel": "PhD",
10   "menCount": 476,
11   "womenCount": 492
12  },
13  ...
14  ]

```

Query.10 - Carrer Salary trend

it calculates the salary trend for each position by number of experience

```

1 SMBUD_project.Carrer.aggregate([
2
3   // Stage 1: Group by occupation and years of experience
4   {
5     $group: {
6       _id: {
7         occupation: "$Current Occupation",    // Group by occupation
8         yearsOfExperience: "$Years of Experience" // Group by years of
          experience
9       },
10      avgSalary: { $avg: "$Salary" }
11    },
12  },
13  // Stage 2: Sort by occupation and years of experience
14  {
15    $sort: {
16      "_id.occupation": 1, // Sort by occupation alphabetically
17      "_id.yearsOfExperience": 1 // Sort by years of experience in
        ascending order
18    },
19  },
20  // Stage 3: Group by occupation and create an array of average
    salaries (progression)
21  {
22    $group: {
23      _id: "$_id.occupation", // Group by occupation
24      progression: {
25        $push: "$avgSalary"
26      }
27    },
28  },
29  // Stage 4: Project the final result to match the desired format
30  {
31    $project: {
32      _id: 0, // Exclude the default _id field
33      occupation: "$_id", // Include the occupation field
34      progression: "$progression" // Include the progression array
35    }
36  }
37 ])
```

Output - Partial

```
1 [{
2   "occupation": "Doctor",
3   "progression": [...],
4 },
5 {
6   "occupation": "Lawyer",
7   "progression": [...],
8 },
9 {
10  "occupation": "Economist",
11  "progression": [
12    111829.96551724138,
13    111638.35643564357,
14    119633.24509803922,
15    115185.45192307692,
16    113446.07142857143,
17    123441.76923076923,
18    110961.87912087912,
19    115924.43564356436,
20    ...
21  ]
22 },
23 ...
24 ]
```


5 | Application

As extra work, we designed and developed an Android application. This application can be run both on physical Android devices and using the online simulator *Appetize.com*, ensuring greater flexibility of use.

Application Description

The application was developed to display queries and their respective results in a clear and interactive manner. The interface is simple and intuitive, allowing the user to enter a query and view the results directly on the app. Moreover, the application is compatible with various versions of Android, ensuring broad accessibility.

Main Features

- Executing queries via a clear and responsive user interface.
- Immediate display of query results.
- Compatibility with physical Android devices and the <https://appetize.io/> simulator.
- Responsive layout, optimized for different screen sizes.

Application Screenshots

The following figures show some screenshots of the application.

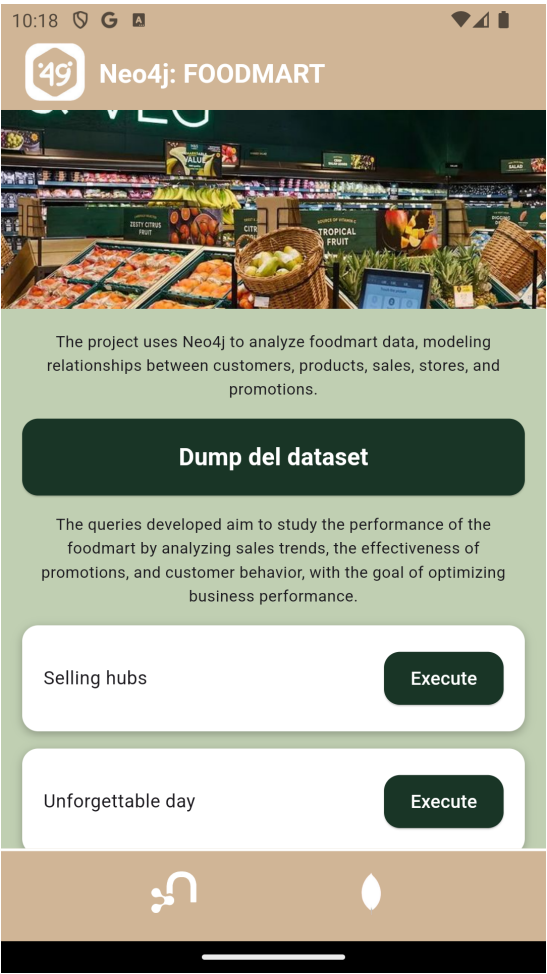


Figure 5.1: Neo4j app section

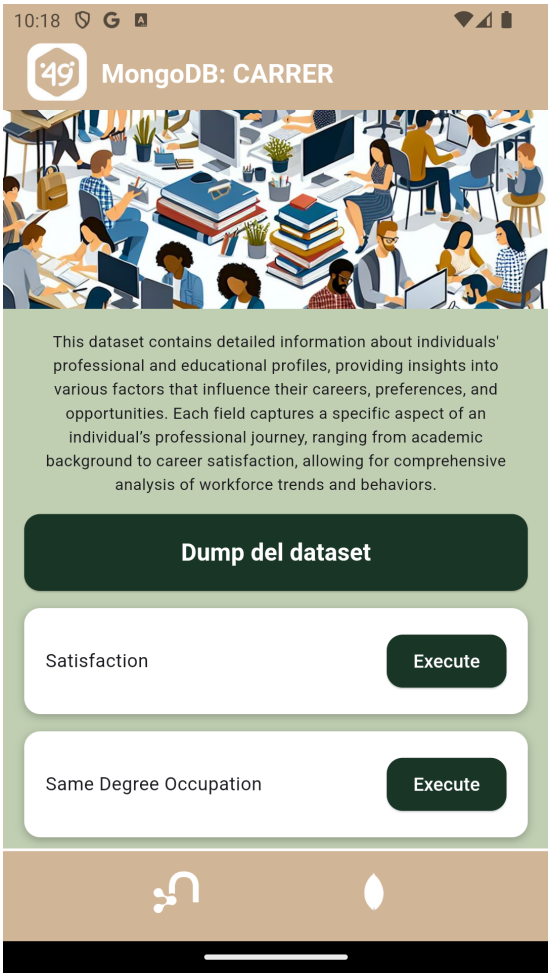


Figure 5.2: MongoDB app section

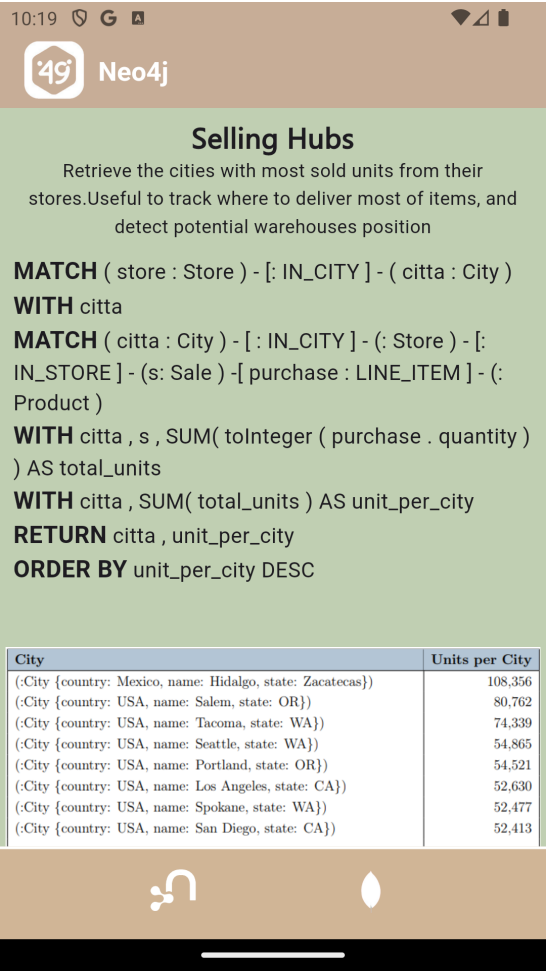


Figure 5.3: Example view of a query

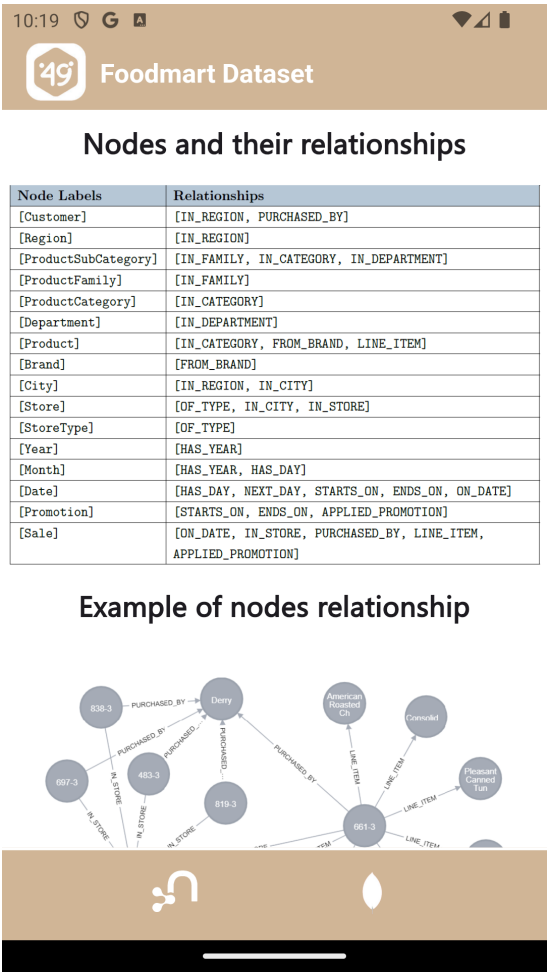


Figure 5.4: Example of a dataset dump view

List of Figures

3.1	Example of nodes relationship	6
4.1	Example of sold items	11
5.1	Neo4j app section	38
5.2	MongoDB app section	38
5.3	Example view of a query	39
5.4	Example of a dataset dump view	39

List of Tables

3.1	Node Labels and Their Properties	7
4.1	Neo4j query 1	9
4.2	Neo4j query 2	10
4.3	Neo4j query 4	12
4.4	Neo4j query 7	15

