# Comparative Analysis of Performance Metrics of NoSQL and Relational Databases

Author:
Angela Mladenovska

June 2024

# Table of Contents

# 1. Introduction

Data is an extremely significant component in the modern world. Various conclusions are drawn from their manipulation and analysis. They are arranged and kept in a database to ensure effective organization and productive work. Although there are many other kinds of databases, relational and non-relational databases are the most well-known.

The foundation of relational databases is the tabular arrangement of data, which consists of rows that display individual entries and columns that represent features that are present in the table. Only structured data is managed by this kind of database. It is possible to establish relationships that correspond between several tables. Relationships can be formed between different tables, as the name depicts. Most well known relational database management systems are MySQL, Oracle, Microsoft SQL Server, and PostgreSQL. PostgreSQL's performances will be actually taken into consideration throughout this project.

On the other hand, non-relational databases or NoSQL have a different approach from the relational ones: they do not have a classic tabular structure. They can be graph-based, document-based, key-value and column family [1]. The non-relational database management system MongoDB, which stores data in JSON format, a representation of a document-oriented database will be introduced to us in this project.

To clarify, document-oriented databases offer efficient data management, mainly in a semi-structured format.

Different types of databases produce different results when running the same queries. In this project assignment we will compare the timing performance of PostgreSQL and MongoDB.

# 2. Dataset Selection and Analysis

For the purpose of this project, secondary data analysis[2] is performed, which is analysis of data that has been collected from various sources for a particular purpose.

This process consists of several steps: defining a research topic, selecting a dataset, reviewing the selected dataset. We will later discuss the steps that were taken in detail.

It was necessary to select a larger dataset, with several hundreds to a million records. On the link[3] where the selected dataset is located, there are two datasets: books_data.csv and books_rating.csv. For our purposes, we only use the books_rating dataset.

The Books rating dataset contains about 3 million records and consists of 10 columns: Id, Title, Price, User_id, profileName, review/helpfulness, review/score, review/time, review/summary, review/text.

After the selection was defined, the dataset was reviewed, and since it initially contained null values and unnecessarily formatted values, it needed to be processed and analyzed accordingly. The data analysis and data preparation was performed in the Google Colab tool.



Picture 1: Initial Overview of the Dataset

Additionally, a check was made on how many null (nan/null) values exist in each of the columns. This was done in order to properly handle them, because in the general case such values would be unnecessary.



Picture 2: Overview of the presence of null values in the dataset

A specific method for handling the null values is selected for each column. For the Title, review/summary, review/text and Price columns we dropped the rows where the values are missing. Null values in the User_id and profileName columns are replaced by the value 'Unknown'.

In order to maintain consistency when importing the data, those values from the columns Title, profileName, review/summary, review/text, which contained the value ';'(semicolon), were deleted, because the separator when saving the dataset after its processing was done to be ';'. This is due to the presence of commas in the values represented in this dataset, and by default the comma is the separator in a csv file.

The values from the Id column were not unique, and in order for it to serve as the primary key in the database table, some changes were made to make those values unique.

Another change that has been made is the removal of unnecessary quotes on values from the profileName column.

Also, in order to have a richer choice for data manipulation, a new date_rated column has been added that represents the dates when each book was rated. The column is filled with random dates using Python functions.

After completing the data cleaning, 385 902 records remained in the dataset.

# 3. Creating Relational and Aggregation Models for Data Importing

## 3.1 Relational model

For the relational model, the DataGrip IDE was used, mainly because of the friendly user interface that it offers, for data manipulation. A new Postgres database was created in PgAdmin, where a table was generated for importing the dataset. In DataGrip, a connection was made with the database.

| id | title | price | user_id | profile_name | review_helpfulnes |
|---|---|---|---|---|---|
| 0829814000_1 | Wonderful Worship in Smaller Churches | 19.4 | A373VVEU6Z9M0N | Dr. Terry W. Dorsett | 1/1 |
| 0829814000_2 | Wonderful Worship in Smaller Churches | 19.4 | AGKGOH65VTRR4 | Cynthia L. Lajoy  Cindy La Joy | 1/1 |
| 0829814000_3 | Wonderful Worship in Smaller Churches | 19.4 | A30QWLU31BU1Y | Maxwell Grant | 1/1 |
| 0595344550_1 | Whispers of the Wicked Saints | 10.95 | A1E9M6APK30ZAU | V. Powell | 1/2 |
| 0595344550_2 | Whispers of the Wicked Saints | 10.95 | AUR0VA5H0C66C | LoveToRead  Actually Read Books | 1/2 |

| review_score | review_time | review_summary | review_text | date_rating |
|---|---|---|---|---|
| 5 | 1291766400 | Small Churches CAN Have Wonderful Worship | "Many small churches feel like they can not hav… | 2024-02-23 |
| 5 | 1248307200 | Not Just for Pastors! | I just finished reading this amazing book and c… | 2022-04-22 |
| 5 | 1222560000 | Small church pastor? This is the book on worship | "I hadn't been a small church pastor very long … | 2022-09-13 |
| 4 | 1119571200 | Here is my opinion | "I have to admit, I am not one to write reviews… | 2020-02-25 |
| 1 | 1119225600 | Buyer beware | "This is a self-published book, and if you want… | 2022-01-23 |

Picture 3: Overview of the imported dataset in DataGrip

## 3.2 Aggregation model

The MongoDB non-relational database management system was chosen for working with the non-relational database. MongoDB was installed locally on the machine, additionally MongoDBCompass IDE was installed where a connection was made with the database. The csv dataset was imported there as well, but its structure was modified in JSON format.

It was explored how MongoDB works, before we moved on to further work. [5]

```
▾ {
▸   "_id": {...},
    "Id": "0829814000_1",
    "Title": "Wonderful Worship in Smaller Churches",
    "Price": 19.4,
    "User_id": "A373VVEU6Z9M0N",
    "profileName": "Dr. Terry W. Dorsett",
    "review/helpfulness": "1/1",
    "review/score": 5,
    "review/time": 1291766400,
    "review/summary": "Small Churches CAN Have Wonderful Worship",
    "review/text": "Many small churches feel like they can not have great worship because they lack the resources that larger churches take fo…
    "Date": "2024-02-23"
  }
```

Picture 4: Overview of a single JSON object in MongoDBCompass

Also, a check was made whether the number of the records is equal in both of the databases.



Picture 5: Total number of rows in the relational model



Picture 6: Total number of records in NoSQL

With this, we ensured that the dataset was imported appropriately. This indicates that there were not any issues during the importing of the datasets.

# 4. Defining queries

After the models were created, the next step was to define queries of various types: aggregation, filtering, retrieving.

## 4.1 Queries in PostgreSQL

After the successful importing of the data in the database, 10 queries were written which purpose is to be executed in the database. Pictures of the queries can be found in this section, as well as description of the query can be found in the comments on the pictures.

```sql
--Query 1:List the books with their prices, whose title contains less than 5 words and the price is not zero.
-- Order the result in descending order
select distinct title, price
from books_ratings
where length(title) - length(replace(title,' ','')) + 1 < 5
and price!=0
order by price DESC;
```

```sql
--Query 2:List the titles of the books that have the best rating(5), the rating's text contains more than 5 words,
-- the summarization text has more than 3 words and the profile name is not Unknown
select title
from books_ratings
where review_score = 5
and length(review_text) - length(replace(review_text,' ','')) + 1 > 5
and length(review_summary) - length(replace(review_summary,' ','')) + 1 > 3
and profile_name not like '%Unknown%';
```

```sql
--Query 3:List the books which price is between 2 and 4, are given rating between 1 and 3, have a valid rating date format
-- (yyyy-mm-dd) and the name of the user that has given the rating must not be Unknown or an email address.
select date_rating,profile_name
from books_ratings
where price between 2 and 4
and review_score between 1 and 3
AND (date_rating LIKE '____-__-__')
and (profile_name not like 'Unknown' and profile_name not like '%.com');
```

```sql
--Query 4:According to the words found in the review_summary, for every user calculate how many reviews have positive and negative reactions
select profile_name,
    count(*) as total_reviews,
    sum(case when review_summary ilike '%excellent%' or review_summary ilike '%good%' or review_summary ilike '%great%'
        or review_summary ilike '%awesome%' or review_summary ilike '%would recommend%' then 1 else 0 end) as good_reviews,
    sum(case when review_summary ilike '%did not enjoy%' or review_summary ilike '%did not like%' or review_summary ilike '%boring%'
        or review_summary ilike '%would not recommend%' or (review_summary ilike '%too bad%' and review_summary not ilike '%not%')
        then 1 else 0 end) as bad_reviews
from books_ratings
where profile_name not ilike 'Unknown' and review_score > 3
group by profile_name
order by total_reviews desc;
```

```sql
--Query 5:List the users' information who have given reviews with helpfulness greater than 50%, including the average
--result of the given ratings for every user
select user_id,profile_name,
        count(*) as total_reviews,
        avg(review_score) as average_score
from books_ratings
where cast(substring(review_helpfulness,1,position('/' in review_helpfulness)-1) as double precision) /
      nullif(cast(substring(review_helpfulness,position('/' in review_helpfulness)+1) as double precision),0) > 0.5
group by user_id, profile_name
having count(*) >= 5
order by average_score desc;
```

```sql
--Query 6:List the details of the books with highest average rating from the users, including the number of ratings and the
-- time of the most recent rating for the book.
select rated_books.title, average_score, total_reviews,
       most_recent_review_time, review_summary, review_text
from(
    select title,
           avg(review_score)as average_score,
           count(*) as total_reviews,
           max(review_time)as most_recent_review_time
    from books_ratings
    group by title) as rated_books
join books_ratings on rated_books.title = books_ratings.title
order by average_score desc
limit 5;
```

```sql
--Query 7:List the books with greatest number of ratings from anonymous users
select title,
       count(case when profile_name = 'Unknown' then 1 else 0 end)as number_of_anonymous_reviews,
       avg(review_score)as average_rating_for_anonymous_reviews,
       case
           when avg(review_score) >= 4 then 'Positive'
           when avg(review_score) < 2 then 'Negative'
           else 'Mixed'
       end as average_sentiment_for_book
from books_ratings
group by title
having count(case when profile_name='Unknown' then 1 else 0 end) > 5
order by number_of_anonymous_reviews desc;
```

```sql
--Query 8:List the top 10 best rated books with more than 3 reviews according to the most recent reviews, sorted by
--the average rating and the number of reviews
select br.title, br.price, br.profile_name as reviewer_name,
       avg(review_score)as average_rating,
       count(id)as number_of_reviews
from books_ratings as br,
    (
        select title,max(review_time)as latest_review_time
        from books_ratings
        group by title
    ) as latest_reviews
where br.title = latest_reviews.title and br.review_time = latest_reviews.latest_review_time
group by br.title, br.price, br.profile_name
having count(br.id) > 3
order by average_rating desc, number_of_reviews desc
limit 10;
```

```sql
--Query 9: List the number of books that start on "The" and that have been given rating during
--the 4 trimesters with review_score greater than 4
SELECT
  (SELECT COUNT(*)
    FROM books_ratings
    WHERE review_score > 4
      AND EXTRACT(MONTH FROM TO_DATE(date_rating, 'YYYY-MM-DD')) BETWEEN 1 AND 3
      AND title LIKE 'The%') AS num_books_above_4_first_quarter,

  (SELECT COUNT(*)
    FROM books_ratings
    WHERE review_score > 4
      AND EXTRACT(MONTH FROM TO_DATE(date_rating, 'YYYY-MM-DD')) BETWEEN 4 AND 6
      AND title LIKE 'The%') AS num_books_above_4_second_quarter,

  (SELECT COUNT(*)
    FROM books_ratings
    WHERE review_score > 4
      AND EXTRACT(MONTH FROM TO_DATE(date_rating, 'YYYY-MM-DD')) BETWEEN 7 AND 9
      AND title LIKE 'The%') AS num_books_above_4_third_quarter,

  (SELECT COUNT(*)
    FROM books_ratings
    WHERE review_score > 4
      AND EXTRACT(MONTH FROM TO_DATE(date_rating, 'YYYY-MM-DD')) BETWEEN 10 AND 12
      AND title LIKE 'The%') AS num_books_above_4_fourth_quarter;
```

```sql
--Query 10:List all the books that were rated in the past 2 years and whose review_text contain more than 60 words
select title,price,review_text,date_rating
from books_ratings
where TO_DATE(date_rating, 'YYYY-MM-DD') BETWEEN CURRENT_DATE - INTERVAL '2 year' AND CURRENT_DATE
and length(review_text) - length(replace(review_text,' ','')) + 1 > 60
order by date_rating asc;
```

## 4.2 Queries in MongoDB

The same queries were written for MongoDB. Pictures of the queries can be found in the following section.

Query 1

```
 1 ▼ [
 2 ▼   {
 3        $match:
 4 ▼        {
 5 ▼          $and: [
 6 ▼            {
 7 ▼              Title: {
 8                  $exists: true,
 9                  $ne: ""
10              }
11            },
12 ▼            {
13 ▼              Price: {
14                  $ne: 0
15              }
16            },
17 ▼            {
18 ▼              $expr: {
19 ▼                $lt: [
20 ▼                  {
21 ▼                    $subtract: [
22 ▼                      {
23                          $strLenCP: "$Title"
24                      },
25 ▼                      {
26 ▼                        $strLenCP: {
27 ▼                          $replaceAll: {
28                              input: "$Title",
```

```
                             input:  $Title ,
29                            find: " ",
30                            replacement: ""
31                          }
32                        }
33                      }
34                    ]
35                  },
36                  5
37                ]
38              }
39            }
40          ]
41        }
42    },
43 ▼  {
44      $group:
45 ▼      {
46          _id: "$Title",
47          // Group by title
48 ▼        price: {
49            $first: "$Price"
50          }
51        }
52    },
53 ▼  {
54      $sort:
55 ▼      {
56          price: -1
```

Query 2

```
 2 ▼   {
 3        $match:
 4 ▼        {
 5            "review/score": 5,
 6 ▼          "review/text": {
 7              $exists: true,
 8              $type: "string",
 9              $gt: 5
10            },
11 ▼          "review/summary": {
12              $exists: true,
13              $type: "string",
14              $gt: 3
15            },
16 ▼          profileName: {
17 ▼            $not: {
18                $regex: /Unknown/
19              }
20            }
21          }
22    },
23 ▼  {
24      $project:
25 ▼      {
26          _id: 0,
27          Title: 1
28        }
29      }
30    ]
```

## Query 3

```
1  ▾ [
2  ▾   {
3        $match:
4  ▾       {
5  ▾         Price: {
6             $gte: 2,
7             $lte: 4
8           },
9  ▾         "review/score": {
10            $gte: 1,
11            $lte: 3
12          },
13 ▾        Date: {
14            $regex: /^[0-9]{4}-[0-9]{2}-[0-9]{2}$/
15          },
16 ▾        $nor: [
17 ▾          {
18              profile_name: "Unknown"
19            },
20 ▾          {
21 ▾            profile_name: {
22                $regex: /\.com$/
23              }
24            }
25          ]
26        }
27      },
28 ▾    {
29        $project:
30 ▾        {
31            _id: 0,
32            Date: 1,
33            profileName: 1
34          }
35      }
36    ]
```

## Query 4

```
1  ▾ [
2  ▾   {
3        $match:
4  ▾       {
5  ▾         profileName: {
6  ▾           $not: {
7                $regex: /Unknown/
8              }
9            },
10 ▾         "review/score": {
11             $gt: 3
12           }
13         }
14      },
15 ▾    {
16        $group:
17 ▾        {
18            _id: "$profileName",
19 ▾          total_reviews: {
20              $sum: 1
21            },
```

10

```
29 ▾          good_reviews: {
30 ▾            $sum: {
31 ▾              $cond: [
32 ▾                {
33 ▾                  $or: [
34 ▾                    {
35 ▾                      $regexMatch: {
36                          input: "$review/summary",
37                          regex:
38                            /excellent|good|great|awesome|would recommend/i
39                        }
40                      }
41                    ]
42                  },
43                  1,
44                  0
45                ]
46              }
47            },
48 ▾          bad_reviews: {
49 ▾            $sum: {
50 ▾              $cond: [
51 ▾                {
52 ▾                  $or: [
53 ▾                    {
54 ▾                      $regexMatch: {
55                          input: "$review/summary",
```

```
43 ▾              $or: [
44 ▾                {
45 ▾                  $regexMatch: {
46                      input: "$review/summary",
47                      regex:
48                        /did not enjoy|did not like|boring|would not recommend|too bad(?! not)
49                    }
50                  }
51                ]
52              },
53              1,
54              0
55            ]
56          }
57        }
58      }
59    },
60 ▾  {
61      $sort:
62 ▾      {
63        total_reviews: -1
64      }
65    }
66  ]
```

## Query 5

```
1 ▾ [
2 ▾   {
3 ▾     $addFields: {
4 ▾       helpfulness_numerator: {
5 ▾         $toDouble: {
6 ▾           $arrayElemAt: [
7               { $split: ["$review/helpfulness", "/"] },
8               0
9             ]
10          }
11        },
12 ▾      helpfulness_denominator: {
13 ▾        $toDouble: {
14 ▾          $arrayElemAt: [
15              { $split: ["$review/helpfulness", "/"] },
16              1
17            ]
18          }
19        }
20      }
21    },
22 ▾  {
23 ▾    $addFields: {
24 ▾      helpfulness_ratio: {
25 ▾        $cond: {
26            if: { $eq: ["$helpfulness_denominator", 0] },
27            then: 0,
28            else: { $divide: ["$helpfulness_numerator", "$helpfulness_denominator"] }
```

11

```
33 ▾    {
34 ▾      $match: {
35          helpfulness_ratio: { $gt: 0.5 }
36        }
37      },
38 ▾    {
39 ▾      $group: {
40          _id: { user_id: "$User_id", profile_name: "$profileName" },
41          total_reviews: { $count: {} },
42          average_score: { $avg: "$review/score" }
43        }
44      },
45 ▾    {
46 ▾      $match: {
47          total_reviews: { $gte: 5 }
48        }
49      },
50 ▾    {
51 ▾      $sort: {
52          average_score: -1
53        }
54      },
55 ▾    {
56 ▾      $project: {
57          _id: 0,
58          user_id: "$_id.user_id",
59          profile_name: "$_id.profile_name",
60          total_reviews: 1,
61          average_score: 1
```

Query 6

```
 1 ▾  [
 2 ▾    {
 3 ▾      $group: {
 4          _id: "$Title",
 5          average_score: { $avg: "$review/score" },
 6          total_reviews: { $count: {} },
 7          most_recent_review_time: { $max: "$review/time" }
 8        }
 9      },
10 ▾    {
11        $sort: { average_score: -1 }
12      },
13 ▾    {
14        $limit: 5
15      },
16 ▾    {
17 ▾      $lookup: {
18          from: "books_ratings",
19          let: { title: "$_id", most_recent_time: "$most_recent_review_time" },
20 ▾        pipeline: [
21 ▾          {
22 ▾            $match: {
23 ▾              $expr: {
24 ▾                $and: [
25                    { $eq: ["$Title", "$$title"] },
26                    { $eq: ["$review/time", "$$most_recent_time"] }
27                  ]
28                }
```

```
29                }
30              },
31 ▾          {
32 ▾            $project: {
33                _id: 0,
34                review_summary: "$review/summary",
35                review_text: "$review/text"
36              }
37            }
38          ],
39          as: "recent_review"
40        }
41      },
42 ▾    {
43        $unwind: "$recent_review"
44      },
45 ▾    {
46 ▾      $project: {
47          _id: 0,
48          title: "$_id",
49          average_score: 1,
50          total_reviews: 1,
51          most_recent_review_time: 1,
52          review_summary: "$recent_review.review_summary",
53          review_text: "$recent_review.review_text"
54        }
55      }
56    ]
```

Query 7

```
1 ▼  [
2 ▼    {
3 ▼      $group: {
4          _id: "$Title",
5 ▼        number_of_anonymous_reviews: {
6 ▼          $sum: {
7              $cond: [{ $eq: ["$profileName", "Unknown"] }, 1, 0]
8            }
9          },
10 ▼       average_rating_for_anonymous_reviews: {
11 ▼         $avg: {
12             $cond: [{ $eq: ["$profileName", "Unknown"] }, "$review/score", null]
13           }
14         },
15         average_score: { $avg: "$review/score" }
16       }
17     },
18 ▼   {
19 ▼     $addFields: {
20 ▼       average_sentiment_for_book: {
21 ▼         $switch: {
22 ▼           branches: [
23 ▼             {
24                 case: { $gte: ["$average_score", 4] },
25                 then: "Positive"
26               },
27 ▼             {
28                 case: { $lt: ["$average_score", 2] },
```

```
28                 case: { $lt: ["$average_score", 2] },
29                 then: "Negative"
30               }
31             ],
32             default: "Mixed"
33           }
34         }
35       }
36     },
37 ▼   {
38 ▼     $match: {
39         number_of_anonymous_reviews: { $gt: 5 }
40       }
41     },
42 ▼   {
43       $sort: { number_of_anonymous_reviews: -1 }
44     },
45 ▼   {
46 ▼     $project: {
47         _id: 0,
48         title: "$_id",
49         number_of_anonymous_reviews: 1,
50         average_rating_for_anonymous_reviews: 1,
51         average_sentiment_for_book: 1
52       }
53     }
54   ]
```

## Query 8

```
1 ▼  [
2 ▼    {
3 ▼      $group: {
4          _id: "$Title",
5 ▼        latest_review_time: {
6            $max: "$review/time"
7          }
8        }
9      },
10 ▼   {
11 ▼     $lookup: {
12         from: "books_ratings",
13 ▼       let: {
14           title: "$_id",
15           latest_time: "$latest_review_time"
16         },
17 ▼       pipeline: [
18 ▼         {
19 ▼           $match: {
20 ▼             $expr: {
21 ▼               $and: [
22 ▼                 {
23                     $eq: ["$Title", "$$title"]
24                   },
25 ▼                 {
26 ▼                   $eq: [
27                       "$review/time",
28                       "$$latest_time"
```

```
29                        ]
30                      }
31                    ]
32                  }
33                }
34              }
35            ],
36            as: "latest_reviews"
37          }
38        },
39      {
40          $unwind: "$latest_reviews"
41        },
42      {
43          $group: {
44            _id: {
45              title: "$latest_reviews.Title",
46              price: "$latest_reviews.Price",
47              profile_name:
48                "$latest_reviews.profileName"
49            },
50            average_rating: {
51              $avg: "$latest_reviews.review/score"
52            },
53            number_of_reviews: {
54              $sum: 1
55            },
```

```
56          reviews: {
57            $push: "$latest_reviews"
58          }
59        }
60      },
61      {
62        $match: {
63          number_of_reviews: {
64            $gt: 3
65          }
66        }
67      },
68      {
69        $sort: {
70          average_rating: -1,
71          number_of_reviews: -1
72        }
73      },
74      {
75        $limit: 10
76      },
77      {
78        $project: {
79          _id: 0,
80          title: "$_id.title",
81          price: "$_id.price",
82          reviewer_name: "$_id.profile_name",
83          average_rating: 1,
```

```
84          number_of_reviews: 1
85        }
86      }
87    ]
```

## Query 9

```
1   [
2     {
3       $match: {
4         "review/score": { $gt: 4 },
5         Title: { $regex: "^The" }
6       }
7     },
8     {
9       $project: {
10        month: { $month: { $dateFromString: { dateString: "$Date", format: "%Y-%m-%d" } } }
11      }
12    },
13    {
14      $facet: {
15        num_books_above_4_first_quarter: [
16          { $match: { month: { $in: [1, 2, 3] } } },
17          { $count: "count" }
18        ],
19        num_books_above_4_second_quarter: [
20          { $match: { month: { $in: [4, 5, 6] } } },
21          { $count: "count" }
22        ],
23        num_books_above_4_third_quarter: [
24          { $match: { month: { $in: [7, 8, 9] } } },
25          { $count: "count" }
26        ],
27        num_books_above_4_fourth_quarter: [
28          { $match: { month: { $in: [10, 11, 12] } } },
```

```
27 ▾        num_books_above_4_fourth_quarter: [
28              { $match: { month: { $in: [10, 11, 12] } } },
29              { $count: "count" }
30            ]
31          }
32        },
33 ▾    {
34 ▾      $project: {
35            num_books_above_4_first_quarter: { $arrayElemAt: ["$num_books_above_4_first_quarter.cou
36            num_books_above_4_second_quarter: { $arrayElemAt: ["$num_books_above_4_second_quarter.c
37            num_books_above_4_third_quarter: { $arrayElemAt: ["$num_books_above_4_third_quarter.cou
38            num_books_above_4_fourth_quarter: { $arrayElemAt: ["$num_books_above_4_fourth_quarter.c
39          }
40        }
41      ]
```

## Query 10

```
 1 ▾  [
 2 ▾    {
 3 ▾      $addFields: {
 4 ▾        date_rating_date: {
 5 ▾          $dateFromString: {
 6              dateString: "$Date",
 7              format: "%Y-%m-%d"
 8            }
 9          },
10 ▾        word_count: {
11 ▾          $size: {
12              $split: ["$review/text", " "]
13            }
14          }
15        }
16      },
17 ▾    {
18 ▾      $match: {
19 ▾        date_rating_date: {
20 ▾          $gte: {
21 ▾            $dateSubtract: {
22                startDate: "$$NOW",
23                unit: "year",
24                amount: 2
25              }
26            },
27            $lte: "$$NOW"
28          },
```

```
29 ▾          word_count: {
30                $gt: 60
31              }
32            }
33          },
34 ▾    {
35 ▾      $sort: {
36            date_rating_date: 1
37          }
38        },
39 ▾    {
40 ▾      $project: {
41            _id: 0,
42            title: "$Title",
43            price: "$Price",
44            review_text: "$review/text",
45            date_rating: "$Date"
46          }
47        }
48      ]
```

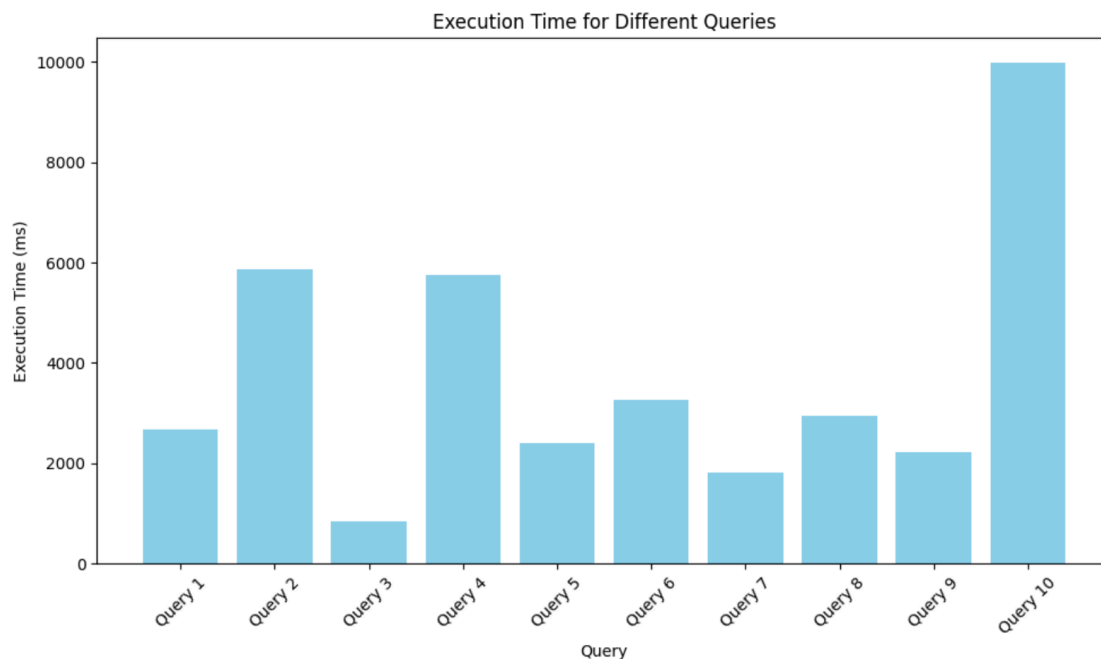# 5. Generating reports about the time performances for the queries

The queries were executed on the following machine: MacBook Air 2019 Intel core i5 8GB RAM. In this section a detailed overview of the time performances are given for each query, in both databases.

## 5.1 Time performances in PostgreSQL

There are 2 times in the PostgreSQL execution: execution time and fetching time.
Execution time refers to the amount of time it takes for the database server to process a query, while fetching time refers to the time taken to send the result set from the database server to the client application. The total execution time is a sum from both of these times.

| Query(Q) | Execution time of the query(ms) |
|----------|--------------------------------|
| Q1 | 2673 |
| Q2 | 5859 |
| Q3 | 836 |
| Q4 | 5742 |
| Q5 | 2401 |
| Q6 | 3263 |
| Q7 | 1823 |
| Q8 | 2939 |
| Q9 | 2215 |
| Q10 | 9969 |

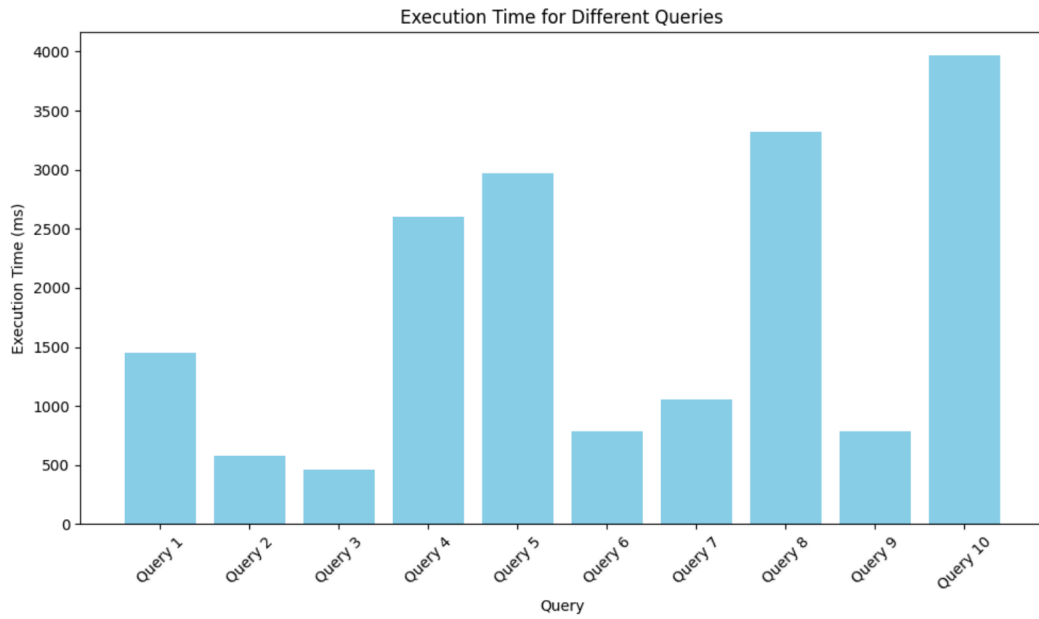Table 1: Overview of the execution time of each query in PostgreSQL

Picture 7: Graphic overview of the execution time of each query in PostgreSQL
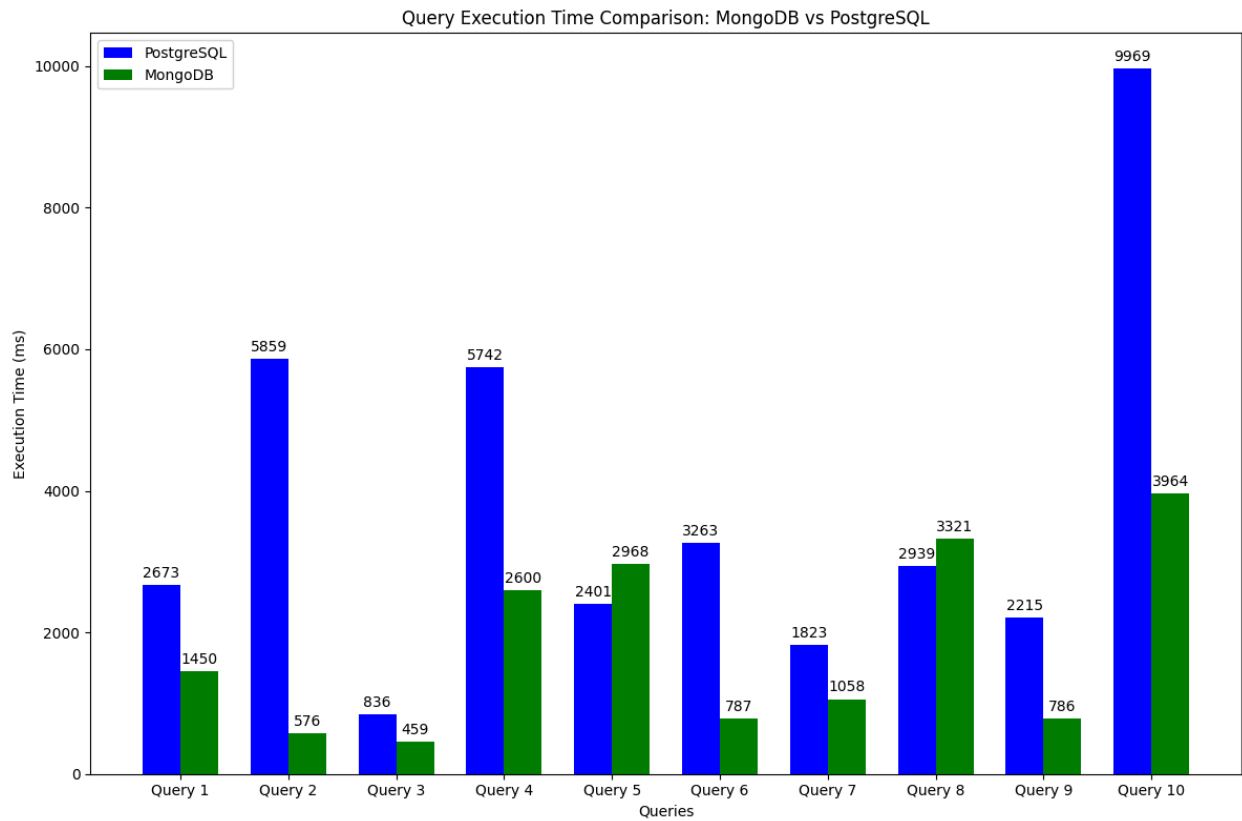
## 5.2 Time performances in MongoDB

| Прашалник(Query-Q) | Време на извршување на прашалникот(ms) |
|---|---|
| Q1 | 1450 |
| Q2 | 576 |
| Q3 | 459 |
| Q4 | 2600 |
| Q5 | 2968 |
| Q6 | 787 |
| Q7 | 1058 |
| Q8 | 3321 |
| Q9 | 786 |
| Q10 | 3964 |

Table 2: Overview of the execution time of each query in MongoDB

Picture 8: Graphic overview of the execution time of each query in MongoDB



Picture 9: Comparison of the execution times for each query in both of the databases using a graph

According to the received results, we can conclude that the non-relational database gives better performances than the relational. Eight out of ten queries have shorter execution time in MongoDB.

The possible reasons why the two remaining queries have larger execution time in MongoDB, can be the fact that MongoDB is not very suitable for joining data from multiple instances, as well as working with functions like cast and substring that take more execution time in this type of database.

# 6. Conclusions

There is often a challenge and dilemma in choosing an appropriate database. A major factor influencing the choice is the nature of the data itself.

Relational databases are known to have schema and take data from fewer sources. They are used when the integrity of the data is preserved and they are normalized. This means that if any data does not match the schema and breaks the constraints, it will not be able to be inserted. This would affect the performance of the database.

On the other hand, non-relational databases take a different approach. They arise mainly as a need for improved performance due to the emergence of huge amounts of data in recent times. Apart from managing unstructured data, they can also manage structured and semi-structured data. They do not maintain a schema like relational databases and have eventual consistency. Specifically in our case, the document-oriented Mongo database creates dynamic schemas, which means that records can be inserted without the structure being created beforehand. MongoDB also supports indexing, which enables fast searching through documents. It is also worth mentioning that this type of database has horizontal scaling, which means that data is distributed across multiple nodes, giving faster performance.[6]

In this project, MongoDB gave better performance. The reasons for this phenomenon are the mentioned characteristics. MongoDB is known to provide better performance for queries involving filtering and aggregations. On the other hand, its weaknesses are working with data joins and nested queries, which is actually inherent to relational databases.

Depending on what we want to achieve and what queries we want to write, strengths and weaknesses should be considered before choosing an appropriate database to work with.

# 7. References

[1] M. Sharma, V. D. Sharma and M. M. Bundele, "Performance Analysis of RDBMS and No SQL Databases: PostgreSQL, MongoDB and Neo4j," 2018 3rd International Conference and Workshops on Recent Advances and Innovations in Engineering (ICRAIE), Jaipur, India, 2018, pp. 1-5, doi: 10.1109/ICRAIE.2018.8710439.

[2] Smith, A.K., Ayanian, J.Z., Covinsky, K.E. *et al. Conducting High-Value Secondary Dataset Analysis: An Introductory Guide and Resources*. *J GEN INTERN MED* **26**, 920–929 (2011). https://doi.org/10.1007/s11606-010-1621-5

[3] *Amazon Books Reviews*. (2022, September 13). Kaggle. https://www.kaggle.com/datasets/mohamedbakhet/amazon-books-reviews/data?fbclid=IwAR07 NneLx8rHJAWxBfAmC2iMHPwXabX5KoAt9Ad6IaGyQ75W8itBlIexIb8

[4] Alex Sington. (2022, October 29). *Import CSV files into PostgreSQL | pgAdmin | Data Analyst Skill Tutorial #3* [Video]. YouTube. https://www.youtube.com/watch?v=UdMj8iKSCoc

[5] *MonGoDB: The Definitive Guide*. (n.d.). Google Books. https://books.google.mk/books?hl=en&lr=&id=pIrCDwAAQBAJ&oi=fnd&pg=PR2&dq=installing+mongodb&ots=bnNziBEsEd&sig=-IzwpZWmOYMWDARjsvpemm9pggQ&redir_esc=y#v=onepage&q=installing%20mongodb&f=false

[6] Dipina Damodaran, B., Salim, S. and Vargese, S.M., 2016. Performance evaluation of MySQL and MongoDB databases. *Int. J. Cybern. Inform.(IJCI)*, *5*, pp.387-394.
Link to paper