**POLITECNICO**

MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

# Systems and Methods for Big and Unstructured Data Project

Author(s): **Tomaso Stefanizzi**

Personal Code: **10713635**

Academic Year: 2023-2024

# Contents

# 1 | Introduction

My idea for this project was to develop a dashboard in the context of an Amazon-style bookshop. The dataset at the core of this project comprises information about books, including details such as category, authors, user ratings, and other relevant attributes. To effectively model the relationships within this dataset, Neo4j, a graph database technology, was selected as the preferred database management system. In fact, the interconnected nature of books, authors, categories, users and reviews, aligns perfectly with the graph structure.

## 1.1.  Data Wrangling/Data Generation

To work with a relatively small dataset, the original data underwent significant filtering through Python code:

- The data was cleaned and filtered to remove null values

- Certain columns were excluded to focus on essential attributes

- Books with multiple authors have been excluded for simplicity

- One column have been added to represent the id of a Review

## 1.2.  Dataset

### 1.2.1.  Original Dataset

The original data is divided into two separated files, one for the Ratings and one for the Books. In the following tables I provide their features.

**Ratings**

| Feature | Description |
|---|---|
| id | The Id of Book |
| Title | Book Title |
| Price | The price of Book |
| User_id | Id of the user who rates the book |
| profileName | Name of the user who rates the book |
| review/helpfulness | helpfulness rating of the review, e.g. 2/3 |
| review/score | rating from 0 to 5 for the book |
| review/time | time of given the review |
| review/summary | the summary of a text review |
| review/text | the full text of a review |

**Books**

| Feature | Description |
|---|---|
| Title | Book Title |
| description | description of book |
| authors | Name of book authors |
| image | url for book cover |
| previewLink | link to access this book on google Books |
| publisher | name of the publisher |
| publishedDate | the date of publish |
| categories | genres of books |
| ratingsCount | averaging rating for book |

## 1.2.2.   Neo4j

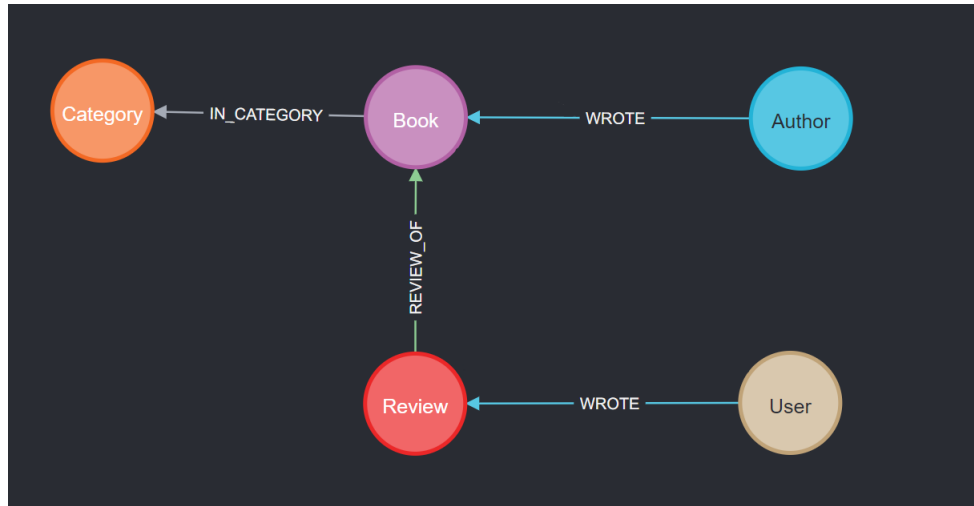After the Data Wrangling process, I created the different entities and relationships in Neo4j:



Figure 1.1: Database Schema

| Node Type | Properties |
|-----------|------------|
| Book | title, ratingsCount, publisher, publishedDate, description |
| Author | name |
| Category | name |
| Review | id, time, text, summary, score, price, helpfulness |

# 2 | Queries

In this section I list 10 of my queries for this project.

### 2.0.1. Top 50 prolific Writers

Returns the list of the top 50 Authors with the highest number of written books

```
1 MATCH (a:Author)-[:WROTE]->(b:Book)
2 RETURN a.name AS Author, count(b) AS NumberOfBooks
3 ORDER BY NumberOfBooks DESC LIMIT 50
```



### 2.0.2. Top 50 Active Users and their Average Rating

Returns the top 50 Users with the highest number of written reviews

```
1 MATCH (u:User)-[:WROTE]->(r:Review)
2 WITH u.id as ID, u.profileName as User, avg(r.score) AS
    AverageRating, count(r) AS NumberOfReviews
3 RETURN User, NumberOfReviews,AverageRating
```

```
4 ORDER BY NumberOfReviews DESC LIMIT 50
```



### 2.0.3.   Average Rating of the 10 most populated Categories

Returns the 10 most populated categories in the dataset, with the average rating of all the books of that specific category
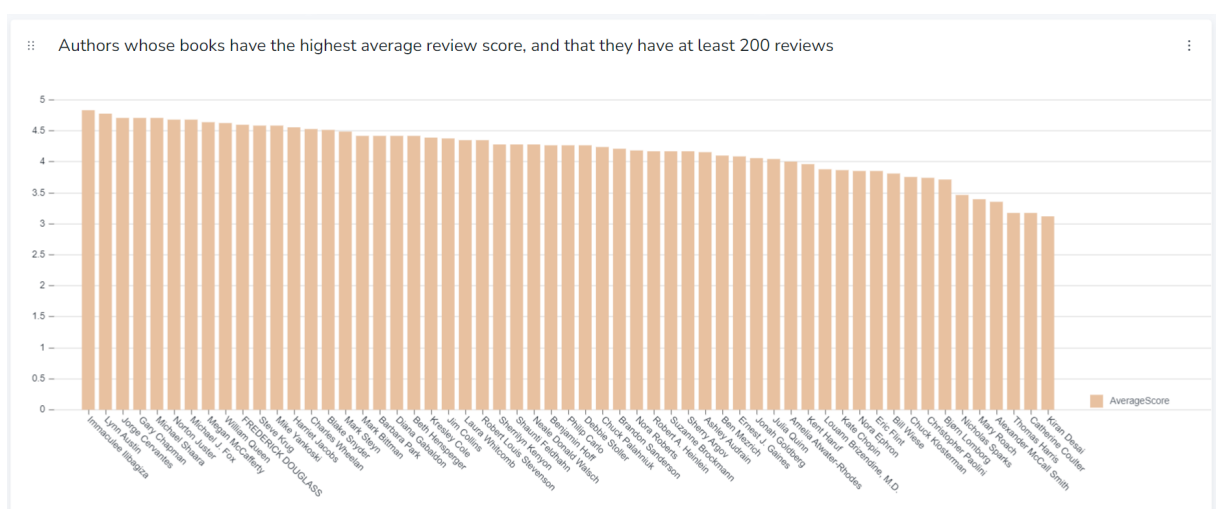
```
1 MATCH (c:Category)<-[:IN_CATEGORY]-(b:Book)
2 WITH c, count(b) AS NumberOfBooks
3 ORDER BY NumberOfBooks DESC
4 LIMIT 10
5 WITH collect(c.name) AS TopCategories
6 UNWIND TopCategories AS Category
7 MATCH (c:Category {name: Category})<-[:IN_CATEGORY]-(b:Book)<-[:
    REVIEW_OF]-(r:Review)
8 RETURN c.name AS Category, avg(r.score) AS AverageRating
```

## 2.0.4. Authors whose books have the highest average review score

Returns the Authors with the their average ratings. Select only the authors with at least 200 reviews.

```
1  MATCH (a:Author)-[:WROTE]->(b:Book)<-[:REVIEW_OF]-(r:Review)
2  WITH a, avg(r.score) AS AverageScore, count(r) AS NumberOfReviews
3  WHERE NumberOfReviews >= 200
4  RETURN a.name AS Author, AverageScore
5  ORDER BY AverageScore DESC
```

### 2.0.5.   Top 50 publishers

Returns the 50 Publishers that published the highest number of books

```
1 MATCH (b: Book)
2 WITH b.publisher as Publisher, count(*) as PublishedBooks
3 RETURN Publisher, PublishedBooks
4 ORDER BY PublishedBooks DESC limit 50
```

| ⠿ Top 50 publishers | ⋮ |
| --- | --- |
| Publisher | PublishedBooks |
| Penguin | 536 |
| Simon and Schuster | 393 |
| Harper Collins | 326 |
| Vintage | 182 |
| Random House | 148 |
| | 1–5 of 50    ‹   › |

### 2.0.6.   Books with the highest average review score

Returns the books ordered by their average score. Select only the books with at least 200 reviews.

```
1 MATCH (b:Book)<-[:REVIEW_OF]-(r:Review)
2 WITH b, avg(r.score) AS AverageScore, count(r) AS NumberOfReviews
3 WHERE NumberOfReviews >= 200
4 RETURN b.title AS Book, AverageScore
5 ORDER BY AverageScore DESC
```

Top Books that have the highest average review score, with at least 200 reviews

## 2.0.7. Ratings distribution

Returns each score and its count to plot the rating distribution

```
1 MATCH (r:Review)
2 RETURN r.score AS Rating, count(*) AS Count
3 ORDER BY Rating
```

Distribution of the Ratings

## 2.0.8. Top 10 Books by Revenue

Returns the top 10 Books ordered by their revenue.

```
1 MATCH (b:Book)<-[:REVIEW_OF]-(r:Review)
2 WITH b, sum(r.price) AS Revenue
3 RETURN b.title as Book, round(Revenue,2) as 'Revenue in $'
```

```
4 ORDER BY Revenue DESC
5 LIMIT 10
```



### 2.0.9.   Top 10 Publishers by Revenue

Returns the top 10 Publishers ordered by their revenue.

```
1 MATCH (b:Book)<-[:REVIEW_OF]-(r:Review)
2 WITH b.publisher AS Publisher, sum(r.price) AS 'Revenue in $'
3 RETURN Publisher, 'Revenue in $'
4 ORDER BY 'Revenue in $' DESC
5 LIMIT 10
```

## 2.0.10.   Price Popularity

Returns the prices that produced more revenue.

```
1 MATCH (b:Book)<-[:REVIEW_OF]-(r:Review)
2 WITH b.publisher AS Publisher, sum(r.price) AS 'Revenue in $'
3 RETURN Publisher, 'Revenue in $'
4 ORDER BY 'Revenue in $' DESC
5 LIMIT 10
```

| ⠿ Rank of the best Prices by Revenue | ⋮ |
| --- | --- |
| Price in $ | Total Revenue in $ |
| 33,97 | 43.243,81 |
| 26,95 | 32.178,3 |
| 54 | 31.806 |
| 32,95 | 29.127,8 |
| 29,95 | 27.164,65 |
| | 1–5 of 419   ‹ › |

# 3 | Dashboard

In order to provide an overview of the system, I created a Dashboard with useful data visualizations.

## 3.1. Structure

I created 3 pages:

- **Main Page**: It contains general informations of the system.

- **Data Insights**: It contains some tables/graphs of the data distribution. On the bottom, I added an interactive plot: for a selected Author, there is a plot of all their books in a graph and in a list on the side.

- **Revenue Report**: It contains some informations related to the revenue. Also here, I provided an interactive plot: for a selected User, we can see his overall purchases (in $), a graph with all his Reviews and also a plot with the price profiling (it's the distribution of all the prices on his purchases).
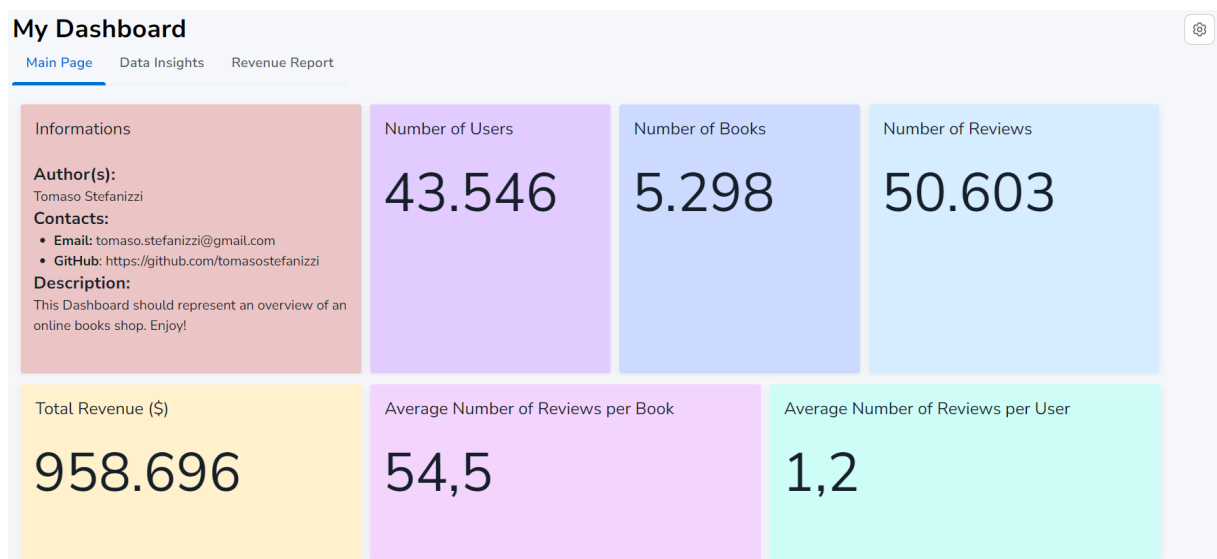


Figure 3.1: Main Page

**My Dashboard** ⚙

Main Page    Data Insights    Revenue Report

| Top 50 prolific Writers | |
| --- | --- |
| Author | NumberOfBooks |
| Georgette Heyer | 33 |
| Agatha Christie | 30 |
| Robert A. Heinlein | 21 |
| John Steinbeck | 20 |
| Terry Pratchett | 20 |
| 1–5 of 50   ‹   › | |

| Top 50 Active Users and their Average Rating | | |
| --- | --- | --- |
| User | NumberOfReviews | AverageRating |
| Harriet Klausner | 86 | 4,709 |
| Midwest Book Review | 55 | 5 |
| Gail Cooke | 29 | 4,483 |
| Blue Tyson "- Research Fir | 25 | 3,48 |
| E. A Solinas "ea_solinas" | 25 | 3 |
| 1–5 of 50   ‹   › | | |

| Average Rating of the 10 most populated Categories | |
| --- | --- |
| Category | AverageRating |
| Fiction | 4,076 |
| Juvenile Fiction | 4,301 |
| Religion | 4,281 |
| Biography & Autobiography | 4,557 |
| Young Adult Fiction | 4,311 |
| 1–5 of 10   ‹   › | |

Authors whose books have the highest average review score, and that they have at least 200 reviews

Distribution of the Ratings

| Top 50 publishers | |
| --- | --- |
| Publisher | PublishedBooks |
| Penguin | 536 |
| Simon and Schuster | 393 |
| Harper Collins | 326 |
| Vintage | 182 |
| Random House | 148 |
| 1–5 of 50   ‹   › | |

Top Books that have the highest average review score, with at least 200 reviews

Select an Author

Author name
Agatha Christie

Autor graph

Author name    Book title

Selected Author Books

Books

CROOKED HOUSE

Murder on the Orient Express

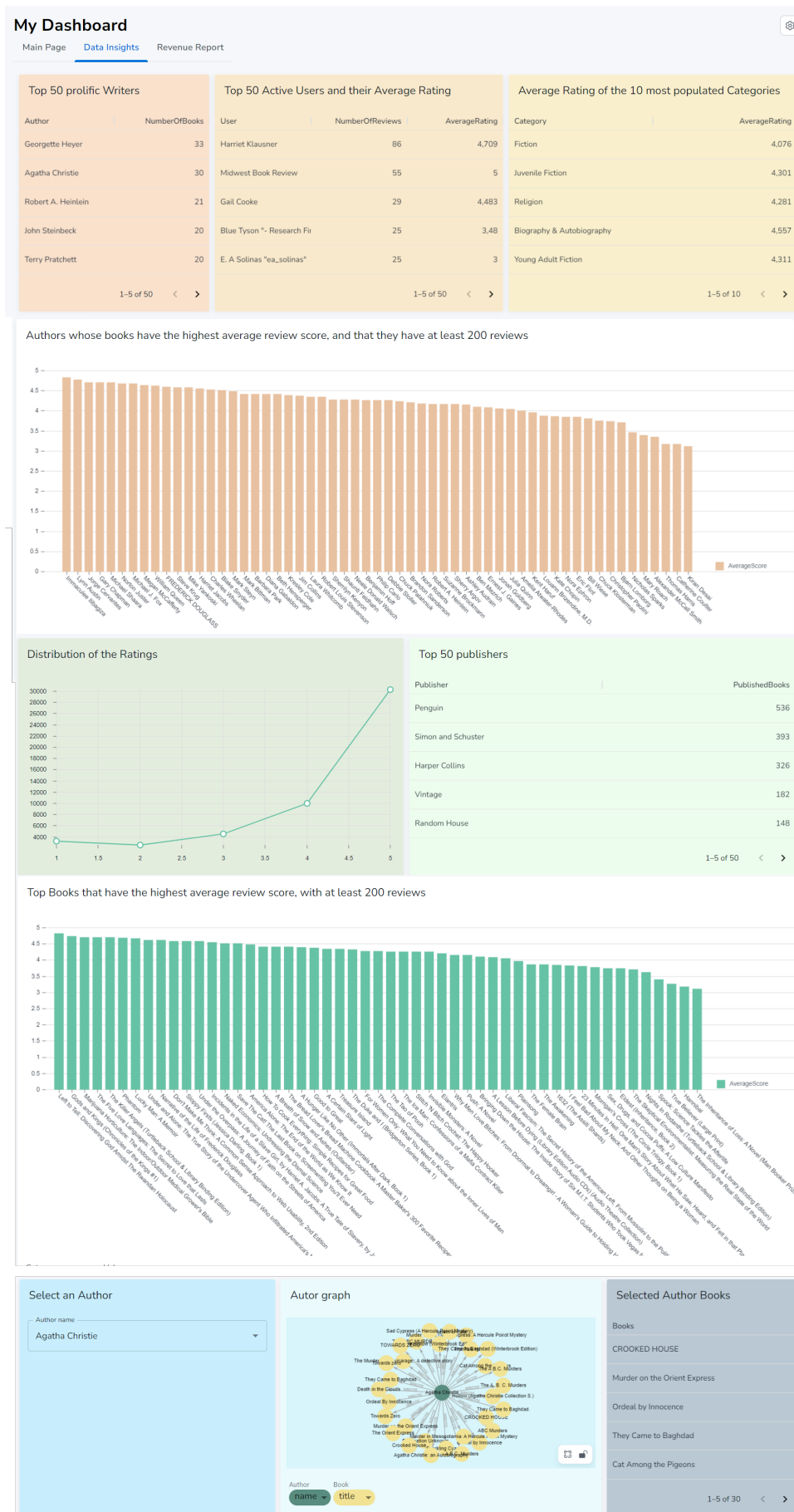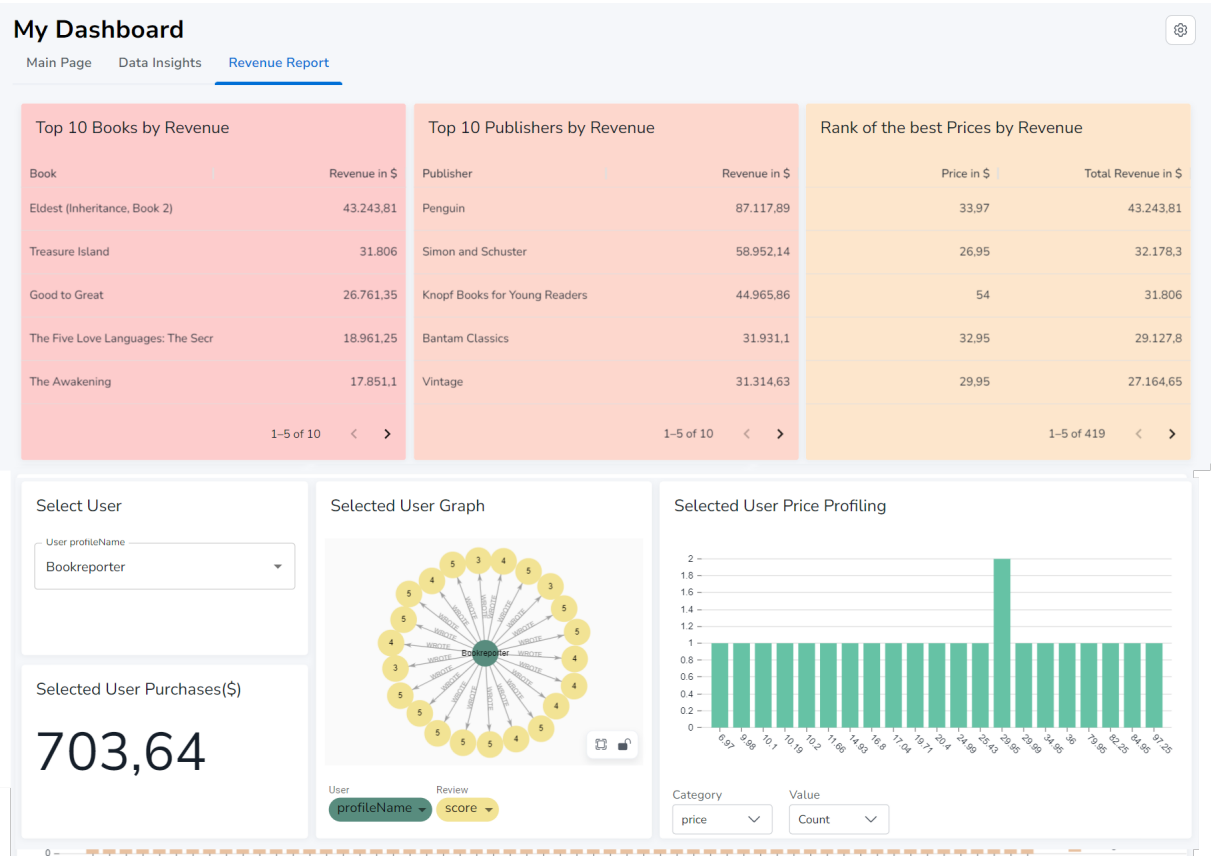Ordeal by Innocence

They Came to Baghdad

Cat Among the Pigeons

1–5 of 30   ‹   ›

Figure 3.2: Data Insights

Figure 3.3: Revenue Report