



Information Retrieval - Submission Report

Project N.14: Art For Sale

Jeferson Morales Mariciano, Filippo Piloni

Università della Svizzera italiana, Faculty of Informatics, Lugano, Switzerland

Contents

| | | |
|----------|--|----------|
| 1 | Introduction | 2 |
| 1.1 | Websites | 2 |
| 1.2 | Features | 2 |
| 2 | Architecture | 3 |
| 2.1 | Technology Stack | 3 |
| 2.1.1 | Scrapy | 3 |
| 2.1.2 | PyTerrier | 4 |
| 2.1.3 | Django & SQLite | 4 |
| 2.1.4 | Vue.js 3 & PrimeVue | 4 |
| 3 | Crawling | 4 |
| 4 | Indexing | 5 |
| 5 | Retrieval | 5 |
| 5.1 | Backend | 5 |
| 5.2 | Frontend | 6 |
| 5.2.1 | User Interface | 6 |
| 6 | Evaluation | 6 |
| 6.1 | User Evaluation | 7 |
| 6.1.1 | User Experience | 7 |
| 7 | Limitations and Future Improvements | 7 |
| 8 | Conclusions | 7 |
| 9 | References | 7 |

1 Introduction

The Information Retrieval course project "*ArtForSale*" aims to be a working prototype of an information retrieval system for a specific task: to display for-sale artworks from art selling related websites listed in Section 1.1. The system has around $\approx 8,500$ indexed documents for searching, browsing and presented through a user-friendly and intuitive interface, further details in UI/UX Section 5.2.1, 6.1.1. The source code repository of the project is linked in Section 9.

1.1 Websites

The following websites were chosen for the project based on availability of data, ease of scraping and low restrictions regarding policies and terms of use to scrape documents i.e. *robots.txt* rules:

- www.artsy.net
- www.artfinder.com
- www.saatchiart.com

The cumulative number of documents gathered from the above websites is around $\approx 8,500$ stored in *results.json* file containing the list of samples in JSON format. Further information concerning the scraping process and implementation is discussed in Crawler Section 3.

1.2 Features

The course project requires 2 additional categorized features to be implemented besides the retrieval system: 1 each from simple and complex group features.

The *Art For Sale* project implements 2 + 1 bonus additional simple features:

- Simple
 - **Result presentation:** results should be presented in a tabular format so that many results could be seen at the same time. Each table cell should contain appropriate information for your project.
 - **Filtering:** in addition to being able to search by title, an user should be able to filter the results based on at least 3 relevant attributes for your project.
- Complex
 - **Automatic recommendation:** In addition to the relevant search results pertaining to the user query, the user should also be suggested "similar" products based on, say, category, description, price etc. The ordering among the recommended items is not important. However, you should mention how did you arrive at the recommendations. For this purpose, you can use any open-source recommenders available.

Shown results are sorted in descending order of relevance from left to right in the tabular view. Implementation details of features are on Architecture Section 2.

Moreover, by exploiting the automatic recommendation feature implemented using category tags and the filtering feature to filter them, the **Results clustering** complex feature was tackled by allowing to group the results into topics based on category tags. Though, a negative filter to exclude documents based on categories was not added, thus partially implementing another extra feature.

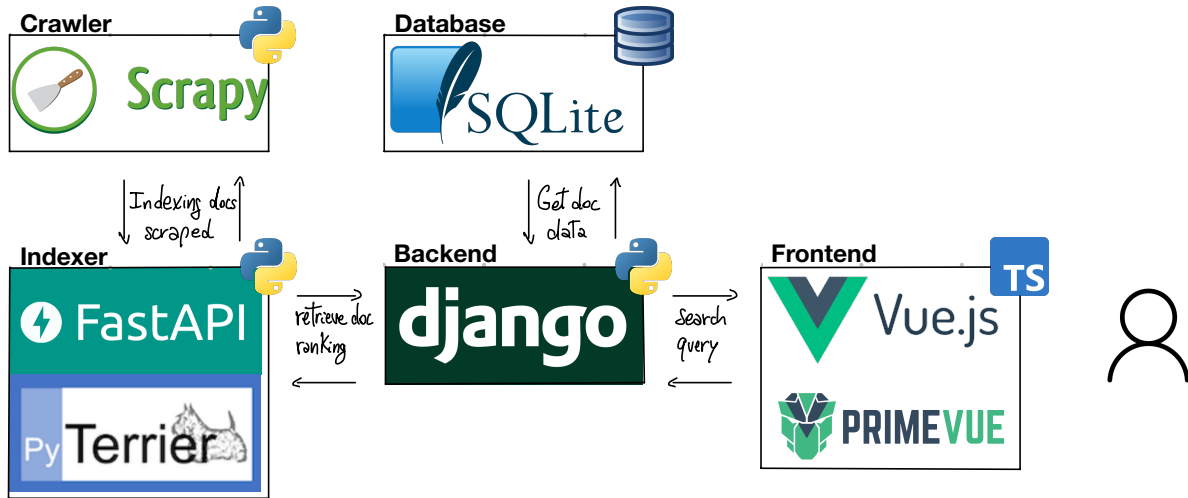
2 Architecture

The **Information Retrieval System** has capabilities limited to a **TIS¹** with **Text retrieval**.

The **Text Access** part helping users to retrieve the right information at the right time, is structured by connecting the system with users in both modes:

- **Pull**: user takes initiative to fetch ad-hoc relevant information from the system, which offers **querying** and **browsing** capabilities. It satisfies **short-term needs**.
- **Push**: system takes initiative and has good knowledge of what are the user's **long-term needs** of information throughout the filtering and recommendation system using **browsing traces** to model the behavior of the user.

Fig 1 Search Engine Architecture of *ArtForSale* project



The software components are: crawler, indexer, database, backend, frontend.

RESTful APIs are used as interface for communication between indexer-backend and backend-frontend. Communication between database-backend is encapsulated by Django, while crawler-indexer needs to be run manually through CLI prompt and file writing.

2.1 Technology Stack

Mandatory technologies to use for the project implementation: scrapy, pyterrier. Here below is listed the tech stack used with a short description and motivation of choice.

2.1.1 Scrapy

Scrapy is a Python application framework for crawling web sites and extracting structured data. It is used to initialize the crawling & scraping process for a specified set of artworks. Crawling process is further discussed in section 3.

¹Text Information System

2.1.2 PyTerrier

PyTerrier is a Python declarative platform for information retrieval experiments. It uses the Java-based Terrier information retrieval platform internally to support indexing and retrieval operations. Indexer process is further discussed in section 4.

2.1.3 Django & SQLite

Django is a FOSS² high-level Python web framework for rapid development and clean, pragmatic design. It takes care of much of the hassle of web development.

By default, the Django configuration uses SQLite as database, which is the easiest choice of implementation already included in Python. Backend functionality and interaction with database is detailed in Sections 5.1.

2.1.4 Vue.js 3 & PrimeVue

The progressive javascript framework described as approachable, performant and versatile for building web user interfaces. It is an industry standard among the js component-based frameworks and it is used in courses since the 2nd year of the bachelor.

Together with PrimeVue, a complete UI suite for Vue.js, it elevates the web application with customizable, feature-rich UI components. Frontend design specifics are in Sections 5.2, 5.2.1, 6.1.1.

The search engine requirements are further described, measured and discussed in the Evaluation Section 6.

3 Crawling

It is part of **Indexing Process**, regarding **Text Acquisition**, in which the system identifies and stores documents for indexing.

The purpose is to harness (relatively) bit text data characterized with the 5 Vs³ located in previous mentioned web sites regarding for sale art. Among the 5 Vs, automating the Velocity attribute to keep track of the dynamic real-time data characterizing the art selling market was not possible: manual re-crawling of the spiders and inverted index re-creation is needed. Crawling through Scrapy parses text that is unstructured to group it in well-structured documents, allowing them to be easily indexed. The structure of each documents is:

- author
- title
- price
- description
- categories
- image url reference
- post url reference

²Free and Open Source

³Volume, Variety, Velocity, Value, Veracity

The extracted data is subsequently channeled into a *.json* file specified at start. After all spiders have been run, every result is manually merged into the *results.json* file. Our sample volume is around $\approx 8,500$. To run crawling refer to project's *README.md*.

4 Indexing

It is part of **Indexing Process**, regarding **Text Transformation**, in which the system transforms documents into index terms, **Index Creation**, where index terms are employed to create data structures to support fast searching, and **Ranking and Retrieval**, which uses query and index to generate a ranked list of docs.

The Pyterrier platform supports indexing and retrieval operations by employing shallow NLP⁴ techniques.

The indexer served through a FastAPI server in *crawler/index.py* file, returns a ranked list of documents in descending order of relevance probability $p \in [0, 1]$, which is the optimal standard strategy justified by the **Probability Ranking Principle** [Robertson 77]. The indexer creates the inverted index from the *results.json* file containing the scraped documents. Characteristics of index:

- Stemmer: Porter, algorithmic-based, Pyterrier's default
- Stopword: Pyterrier's default
- Tokenisation: English words, Pyterrier's default

The retrieval model is BM25: similarity-based model with ranked retrieval including TF⁵, IDF⁶ and sublinear TF transformation to avoid single term dominance. To run the indexer refer to project's *README.md*.

5 Retrieval

browser traces collected in cookies for every search query user issues, part of recommender system.

5.1 Backend

The backend relies on a SQLite database to store the main information about the documents, making possible to retrieve the relevant ones following a call to the index.

Django serves as the foundational framework for constructing the primary structure, bringing with it a multitude of advantages such as rapid development, a robust and secure architecture, built-in administrative features, and a thriving community support ecosystem.

The backend handles the two main fetches of the application:

- `/api/documents/get-documents/"query"/`: This endpoint facilitates the retrieval of relevant documents based on a provided query string.
- `/api/recommendation/get-recommended/`: This endpoint is tasked with retrieving recommended artworks. By inputting a set of tags, it conducts a search within the index, returning all pertinent documents associated with the specified criteria.

⁴Natural Language Processing

⁵Term Frequency

⁶Inverse Document Frequency

The backend handles the main fetches from the frontend. The backend is also the joint between the frontend and the index, making calls to a FASTApi that returns relevant documents given a query. To return to the frontend the main information about the documents, the backend iterates through a database in SQLite that contains the main information for each artwork.

recommender system use-based information filtering and content-based information filtering

5.2 Frontend

Part of **Query and Retrieval Process**, regarding **User Interaction**, in which the system supports the creation of a query and display results.

transforms documents into index terms or features, and **Index Creation**, where the system takes index terms and creates data structures to support fast searching.

The frontend is the main part of the application, where all the features are implemented. For the interface, it has been chosen a simple layout, consisting of only a search-bar and a search button. Once the user performs the search, the aspect of the page changes, adding new components. We will examine each of them:

- First, under the search-bar will appear a new set of control to filter the obtained results according to four main characteristics: site of origin, minimum price, maximum price and tag.
- Second, there will be shown some carousels for the recommended results based on the most frequents tags of the resulted artworks. Essentially, once all the artworks related to the query have been retrieved, the frontend extracts from all these artworks the related tags, it chooses the ones with the higher frequency and, using the `/api/recommendation/get-recommended/` fetch, it shows a maximum of 5 results related to each tag in a moving carousel
- Finally, the artworks retrieved from the query are displayed in a card-like style, showing an insight of the artwork, the author, the title, the price, a short description and the main tags, with a Visit button that will bring the user to the page of that artwork in the original site

Vue.js allows us to create a reactive and dynamic user interface,

5.2.1 User Interface

Emotional design: key in Apple, Google alike

Jacobs heuristics respected

standard filtering system

6 Evaluation

Part of the **Query and Retrieval Process**, regarding **Evaluation**.

Employing of **query expansion** was tested by always adding the *art* keyword at the end of each search to get more results. This led to under-constrained query giving back not relevant results. The reason was due to the fact that for our specific set of concerns, hence documents, the word "*art*" is an high frequency word, thus useless to achieve pertinency as explained by **Zipf's Law**.

At start, every field was in a *pandas.Series* and then given to the indexer but this lead to unsatisfying results: it seemed like only the first field (title) got indexed, giving an **over-constrained query** feeling where either no documents or no relevant ones were returned, Hence we decided to put every field into one big "blob" of text and then indexed that.

Effectiveness: quality of results. **Efficiency:** response time and throughput.

System evaluation: Recall and precision paragraph

6.1 User Evaluation

Part of the **Relevance Feedback Process**, regarding **User Evaluation**.

SUS includes pia, quote are designers users too?

6.1.1 User Experience

Sentiment analysis

7 Limitations and Future Improvements

The project is a *Text Retrieval* system at its current state. Since the type of data provided by artworks is not descriptive nor verbose, implementing *Text Analysis* features becomes troublesome because lacking the large amounts of text data in order to discover interesting patterns buried in text. Nevertheless, *Text Organization* implementation feasibility is high due to short well-structured categorized-based descriptions usually artwork documents have, allowing to connect and navigate scattered information.

8 Conclusions

We had fun with deciding the technology stack to use and the architecture structure. We enjoyed the project freedom that allowed us to put in practice all the methodologies learnt from this course and all the others we had during the bachelor.

9 References

- *ArtForSale* GitHub public repository: <https://github.com/JekxDevil/IR-ArtForSale>
- Course slides on iCorsi: <https://www.icorsi.ch/course/view.php?id=16928>
- *Scrapy*: <https://docs.scrapy.org>
- *Pyterrier*: <https://pyterrier.readthedocs.io>
- *FastAPI*: <https://fastapi.tiangolo.com/>
- *Django*: <https://www.djangoproject.com/>
- *Vue.js 3*: <https://vuejs.org/>
- *PrimeVue*: <https://primevue.org/>