

INFORMATION RETRIEVAL



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΔΥΤΙΚΗΣ ΑΤΤΙΚΗΣ
UNIVERSITY OF WEST ATTICA

DEPARTMENT OF INFORMATION AND COMPUTER ENGINEERING

INFORMATION RETRIEVAL PROJECT BUILDING AN ACADEMIC PAPER SEARCH ENGINE

WORK DETAILS

LABORATORY RESPONSIBILITY: TSELENTI PANAGIOTA

DELIVERY DATE: 1/29/2024

SUBMISSION DEADLINE: 1/29/2024

INFORMATION RETRIEVAL

STUDENT DETAILS 1

STUDENT PHOTO:



NAME: ATHANASIOU VASILIOS EVANGELOS

STUDENT ID: 19390005

STUDENT SEMESTER: 9th

STUDENT STATUS : UNDERGRADUATE

PROGRAM OF STUDY : UNIWA

RECOVERY OF INFORMATION

STUDENT DETAILS 2

STUDENT PHOTO:



NAME: TATSIS PANTELIS

STUDENT ID: 20390226

STUDENT SEMESTER: 7th

STUDENT STATUS : UNDERGRADUATE

PROGRAM OF STUDY : UNIWA

RECOVERY OF INFORMATION

CONTENTS

Steps	5
1. Web Crawler	5
1.a. Select target site	5
1.b. Web crawler implementation	9
1.c. Store data in a structured format	14
2. Text Processing (Text Processing)	16
2.a. Planning	16
2.b. Implementation	16
2.c. Evaluation	17
2.d. Application	18
3. Indexing	22
3.a. Creating the inverted index data structure	22
3.b. Storing the index in a data structure	23
4. Search Engine (Search Engine)	24
4.a. Develop user interface for job search	24
4.b. Implementation of recovery algorithms	25
4.c. Filtering search results by various criteria	40
Query Processing	43
Ranking of results (Ranking)	45
5. System evaluation	49
5.a. Data sets (Dataset)	49
5.b. Evaluation Scenarios	49
5.c. Python Libraries	49
5.d. Applications of metrics	49
5.e. Analysis and Improvements	49

RECOVERY OF INFORMATION

Steps

1. Web Crawler

1.a. Select target site

1.a.1 Planning

The target site is the academic papers repository [arXiv](#). The design model concerns the retrieval of the page with the results produced by the website's search engine from one or more random user queries, with the ultimate goal of creating the repository (dataset) of the local search engine, on which searches will be made with user queries.

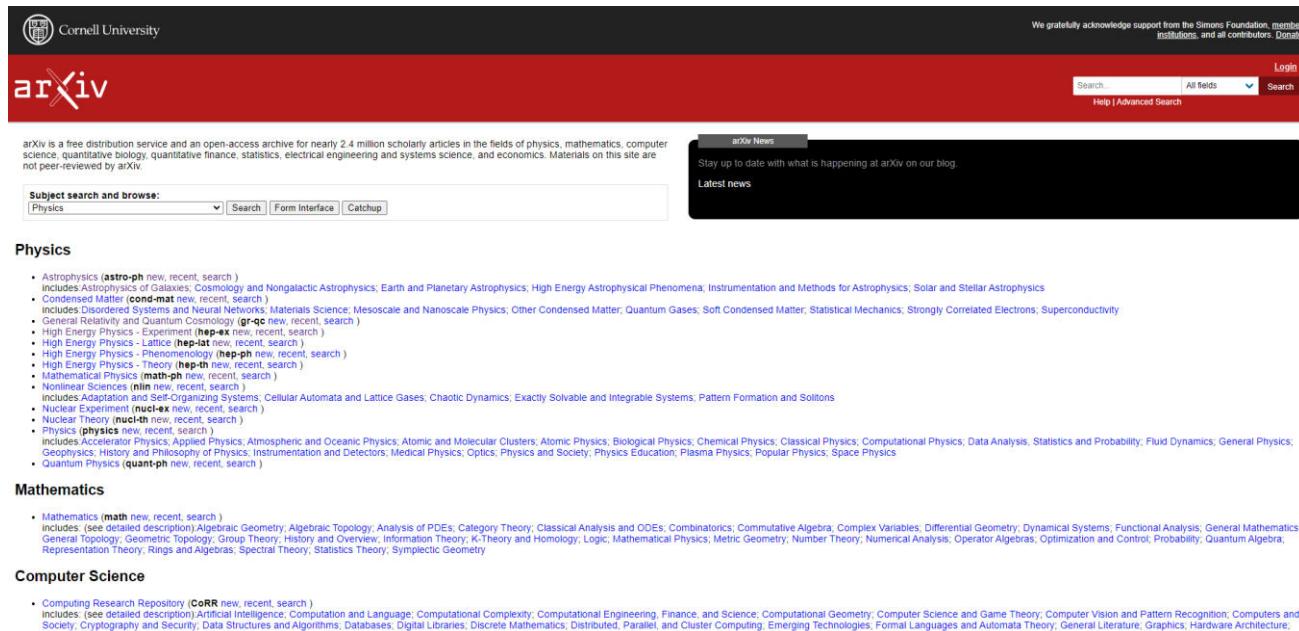


Figure 1.A.1 [The target site arXiv](#)

1.a.2 Implementation

Initially, the number of user queries and their selection is random in each execution of the local search engine (main.py). The questions are limited from 2 to 8 and specifically it has been chosen to be the titles of the basic courses on which academic papers have been posted. The implementation model involves sending an HTTP-GET request to the website's search engine [link](#) along with the corresponding random user query to display results. The process is iterative since the number of questions is more than 1. It is worth noting that for each question, the link that displays 100 results of academic papers has been selected. For example, for the user queries Physics, Statistics and Computer, 3 HTTP-GET requests are sent to the respective links:

https://arxiv.org/search/?query=Physics&searchtype=all&abstracts=show&order=-announced_date_first&size=100

RECOVERY OF INFORMATION

https://arxiv.org/search/?query=Statistics&searchtype=all&abstracts=show&order=-announced_date_first&size=100

https://arxiv.org/search/?query=Computer&searchtype=all&abstracts=show&order=-announced_date_first&size=100

The screenshot shows the arXiv search interface. At the top, there's a navigation bar with the Cornell University logo, the arXiv logo, and links for Help, Advanced Search, and Login. A red banner at the top right credits the Simons Foundation and member institutions. Below the banner is a search bar with a dropdown menu set to 'All fields' and a 'Search' button. To the left of the search bar is a 'Search term...' input field with 'Show abstracts' and 'Hide abstracts' checkboxes. Below the search bar are two sections: 'Searching by Author Name' and 'Searching by subcategory'. The 'Author Name' section contains tips for using the 'Author(s)' field, such as using quotes for exact matches and separating multiple authors with semicolons. The 'subcategory' section provides instructions for searching within a category or combining it with an author or keyword search. A 'Tips' link is located at the bottom left of these sections.

Figure 1.A.2 [The arXiv search engine](#)

The screenshot shows the search results for 'Physics'. The top navigation bar and banner are identical to Figure 1.A.2. The search bar now has 'Physics' selected in the dropdown. Below the search bar, there are buttons for '100 results per page' and 'Sort results by Announcement date (newest first)', along with a 'Go' button. The main content area displays three search results. Result 1 is titled 'Nucleosynthesis in magnetorotational supernovae: impact of the magnetic field configuration' by M. Reichert, M. Bugl, J. Guillet, M. Obergaulinger, M. A. Aloy, A. Arcones. Result 2 is titled 'pix2gestalt: Amodal Segmentation by Synthesizing Wholes' by Ege Ozguroglu, Ruoshi Liu, Didier Suris, Dian Chen, Adhal Dave, Pavel Tokmakov, Carl Vondrick. Result 3 is titled 'Entanglement entropy and deconfined criticality: emergent SO(5) symmetry and proper lattice bipartition'. Each result includes a link to the arXiv page, file formats (pdf, other), and subject categories (astro-ph.HE, cs.CV, cs.LG, cond-mat.str-el, hep-lat).

Figure 1.A.3 [The first 100 search results for Physics](#)

RECOVERY OF INFORMATION

The screenshot shows the arXiv search results for the query "Statistics". The top navigation bar includes the Cornell University logo, the arXiv logo, and a link to "We gratefully acknowledge support from the Simons Foundation and member institutions". The search bar has "Statistics" entered, with "All fields" selected and a "Search" button. Below the search bar are filters for "Show abstracts" (selected) and "Hide abstracts". The results page displays 100 results per page, sorted by "Announcement date (newest first)". The first result is a paper titled "Clustering-based spatial interpolation of parametric post-processing models" by Sándor Baran and Mária Lakatos, with links to PDF, other formats, and categories stat.AP and stat.ME.

Showing 1-100 of 221,264 results for all: Statistics

Search v0.5.6 released 2020-02-24 | Feedback?

All fields | Advanced Search | Login

Statistics

All fields

Search

Show abstracts | Hide abstracts

Advanced Search

100 results per page. Sort results by Announcement date (newest first) Go

1 2 3 4 5 ...

Next

1. arXiv:2401.14393 [pdf, other] stat.AP stat.ME

Clustering-based spatial interpolation of parametric post-processing models

Authors: Sándor Baran, Mária Lakatos

Abstract: ...development of the last three decades, ensemble forecasts still often suffer from the lack of calibration and might exhibit systematic bias, which calls for some form of statistical post-processing. Nowadays, one can choose from a large variety of post-processing approaches, where parametric methods provide full predictive distributions of the investigated w... [More](#)

Submitted 25 January 2024; originally announced January 2024.

Comments: 19 pages, 6 figures

2. arXiv:2401.14378 [pdf, ps, other] cond-mat.stat-mech

Single-file dynamics with general charge measures

Authors: Ziga Krajnik

Abstract: ...transformation acting on the finite-time distribution of particle fluctuations. The transformation is mapped to a simple substitution rule for corresponding full-counting statistics. By taking the asymptotics of the dressing transformation we analyze typical and large scale charge fluctuations.

Typical charge fluctuations in equilibrium states with vanishing... [More](#)

Submitted 25 January 2024; originally announced January 2024.

Comments: 18+9 pages

3. arXiv:2401.14372 [pdf, other] physics.geo-ph

Figure 1.A.4 The first 100 search results of the query Statistics

The screenshot shows the arXiv search results for the query "Computer". The top navigation bar includes the Cornell University logo, the arXiv logo, and a link to "We gratefully acknowledge support from the Simons Foundation and member institutions". The search bar has "Computer" entered, with "All fields" selected and a "Search" button. Below the search bar are filters for "Show abstracts" (selected) and "Hide abstracts". The results page displays 100 results per page, sorted by "Announcement date (newest first)". The first result is a paper titled "Multimodal Pathway: Improve Transformers with Irrelevant Data from Other Modalities" by Yiyuan Zhang, Xiaohan Ding, Kaixiong Gong, Yuxiao Ge, Ying Shan, Xiangyu Yue, with links to PDF, other formats, and categories cs.CV, cs.AI, and cs.LG.

Showing 1-100 of 735,886 results for all: Computer

Search v0.5.6 released 2020-02-24 | Feedback?

All fields | Advanced Search | Login

Computer

All fields

Search

Show abstracts | Hide abstracts

Advanced Search

100 results per page. Sort results by Announcement date (newest first) Go

1 2 3 4 5 ...

Next

1. arXiv:2401.14405 [pdf, other] cs.CV cs.AI cs.LG

Multimodal Pathway: Improve Transformers with Irrelevant Data from Other Modalities

Authors: Yiyuan Zhang, Xiaohan Ding, Kaixiong Gong, Yuxiao Ge, Ying Shan, Xiangyu Yue

Abstract: We propose to improve transformers of a specific modality with irrelevant data from other modalities, e.g., improve an ImageNet model with audio or point cloud datasets. We would like to highlight that the data samples of the target modality are irrelevant to the other modalities, which distinguishes our method from other works utilizing paired (e.g., CLIP) or interleaved data of different modality... [More](#)

Submitted 25 January 2024; originally announced January 2024.

Comments: The code and models are available at <https://github.com/AI-Lab-CV/M2PT>

2. arXiv:2401.14404 [pdf, other] cs.CV cs.LG

Deconstructing Denoising Diffusion Models for Self-Supervised Learning

Authors: Xinlei Chen, Zhuang Liu, Saining Xie, Kaiming He

Abstract: In this study, we examine the representation learning abilities of Denoising Diffusion Models (DDM) that were originally purposed for image generation. Our philosophy is to deconstruct a DDM, gradually transforming it into a classical Denoising Autoencoder (DAE). This deconstructive procedure allows us to explore how various components of modern DDMs influence self-supervised representation learning... [More](#)

Submitted 25 January 2024; originally announced January 2024.

Comments: Technical report, 10 pages

3. arXiv:2401.14403 [pdf, other] cs.RO cs.AI cs.CV cs.LG eess.IV

Figure 1.A.5 The first 100 search results for the query Computer

1.a.3 Evaluation

The local engine supports exception handling for the case that the return code of an HTTP-GET request is other than 200. Strictly speaking, a code of 200 means that the request was served and the page on which the request was made was successfully retrieved. It is worth noting that the selection of random user queries for the creation of the dataset is from 2 to 8, so that the size of the dataset is from 200 to 800 tasks.

RECOVERY OF INFORMATION

1.b. Web crawler implementation

1.b.1 Planning

The web crawler is designed to collect the metadata of the academic papers (web scrape) from the page displayed by the website search engine by inputting one of the random user queries. The process is iterative from the fact documented in the [implementation](#) of step 1.a Select a target site.

1.b.2 Implementation

The implementation follows the Python web-crawler BeautifulSoup and takes place in a separate module (web_crawler.py), where it returns the page requested for retrieval through a successful HTTP-GET request, in unstructured HTML format. The data collected from the tasks are as follows:

- Title
- Authors
- Courses and related sub-courses
- Summary
- Comments
- Publication date
- Link to download the work in pdf format

In the implementation, an identifier integer (doc_id), which is unique, is added to the metadata of each task. This technique aims to easily retrieve a task from the local search engine without having to have all its metadata as an argument. The metadata is collected exactly as it is embedded in the unstructured HTML format and with few variations to store the main content of each field. It is worth noting that the authors and the courses and sub-courses related to the work, have been stored in a structure instead of a single textual content for the better organization of the data. Also, some tasks have no content in the comments field, so it is replaced with a blank. Finally, each task's data is stored in a dictionary structure.

The screenshot shows a web browser displaying the arXiv search results for the query "Statistics". The results page shows 1-100 of 221,264 results. The developer tools (Elements tab) are open, showing the HTML structure of a search result item. The HTML includes various CSS classes like .title, .mathjax, .arxiv-result, and .arxiv-message. The developer tools also show the network requests, styles, and other developer-related information.

RECOVERY OF INFORMATION

Figure 1.B.1 Retrieving the title of the work

The screenshot shows the arXiv search interface for the 'Statistics' category. The search bar at the top contains the query 'Statistics'. Below the search bar, there are filters for 'p.authors', 'Color', 'Font', 'Margin', 'Role', and 'Keyboard-focusable'. The main search results list the paper 'On of parametric post-processing models' by Sándor Baran and Mária Lakatos. The title is displayed in a large blue box. The abstract discusses ensemble-based probabilistic forecasting and proposes a general clustering-based interpolation technique. The paper was submitted on January 25, 2024, and originally announced on January 2024. It has 19 pages and 6 figures.

Figure 1.B.2 Retrieving the authors of the paper

This screenshot is identical to Figure 1.B.1, showing the arXiv search results for 'Statistics'. The paper 'On of parametric post-processing models' by Sándor Baran and Mária Lakatos is highlighted. The authors' names are listed in the title section of the search results page.

Figure 1.B.3 Retrieval of work-related courses

RECOVERY OF INFORMATION

The screenshot shows the arXiv search interface. At the top, there's a red header bar with the Cornell University logo, the word "arXiv", and a search bar containing "Statistics". Below the header, the main content area has a title "Showing 1–100 of 221,264 results for all: Statistics". A sub-header indicates "Search v0.5.6 released 2020-02-24 Feedback?". The search results are displayed in a grid format. Each result includes a thumbnail, the title, authors, date, and a "View" button. A sidebar on the left shows navigation links like "Statistics", "All fields", and "Search". On the right, there's a detailed view of the first result, showing its abstract and full text. The bottom right corner features a blue "arXiv" logo.

Figure 1.B.4 Retrieve the job summary

The screenshot shows the arXiv search interface for the 'Statistics' category. At the top, there's a navigation bar with Cornell University and the arXiv logo, followed by a search bar and a 'Login' button. The main title 'Showing 1–100 of 221,264 results for all: Statistics' is displayed above a search results table. The table has columns for 'Statistics', 'All fields', and 'Search'. Below the table, there are filters for 'Show abstracts' and 'Hide abstracts', and a link to 'Advanced Search'. A pagination bar shows pages 1 through 5, with a 'Next' button. The first result is a paper titled 'Clustering-based spatial interpolation of parametric post-processing models' by Sándor Baran and Mária Lakatos. The abstract discusses the use of ensemble prediction systems for probabilistic forecasting and the challenges of calibration and systematic bias. It proposes a clustering-based method for post-processing. The paper is available as arXiv:2401.14393 [pdf, other] and includes links to stat.AP and stat.ME. The right side of the screen shows the browser's developer tools, specifically the Elements and Styles tabs, which are being used to inspect the page's CSS.

Figure 1.B.5 Retrieving the comments of the work

RECOVERY OF INFORMATION

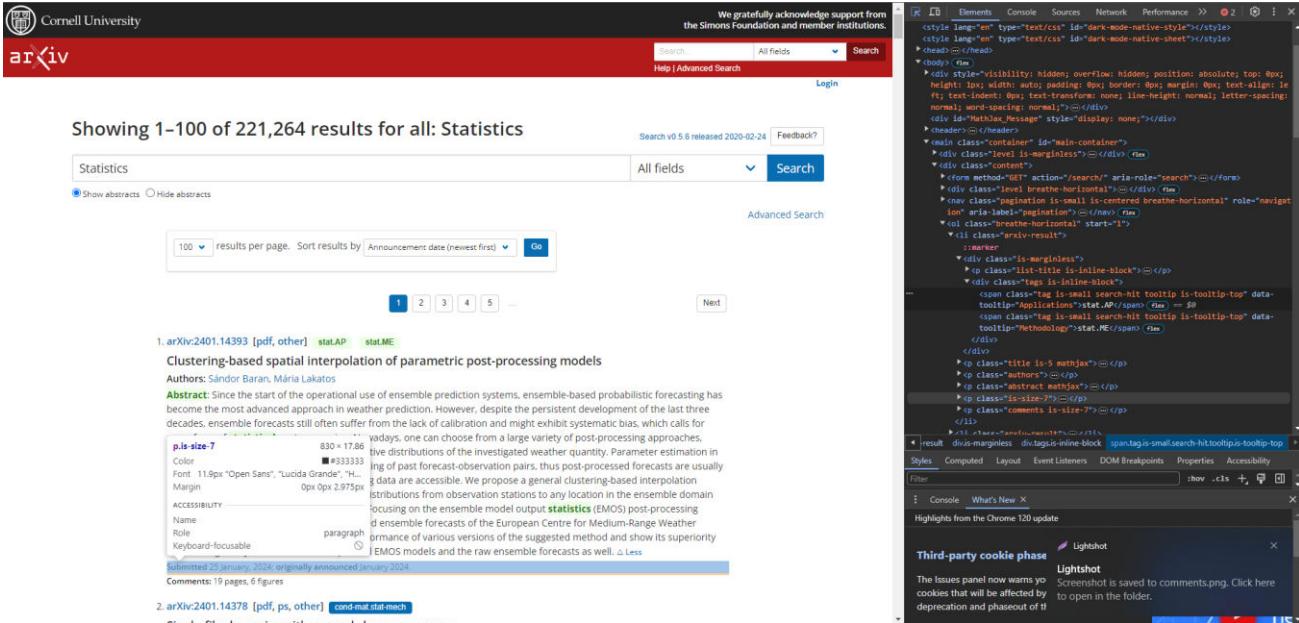


Figure 1.B.6 Retrieving the publication date of the work

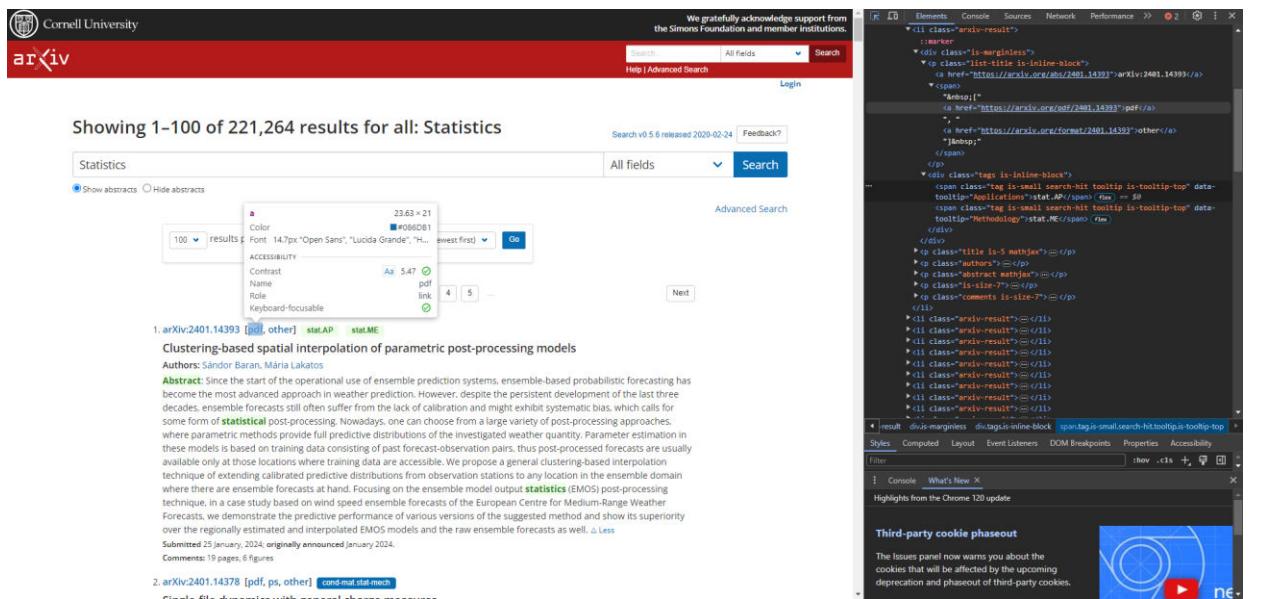


Figure 1.B.7 Retrieving the link to download the work in pdf

1.b.3 Evaluation

The original implementation of the scraper was more complicated, as, for each paper, the collection of metadata (web scrape) took place on the separate pages of each paper with a link arxiv.org/abs/XXX.XXXX (eg <https://arxiv.org/abs/2401.14393>). This created significant delays in the harvester's performance, as collecting data from 200-800 jobs meant sending HTTP-GET requests to 200-800 job links to display each job's respective data pages in unstructured HTML. The final [implementation](#) reduces page retrieval from 200-800 links to 2-8 links, where job data is collected as-is as embedded in unstructured HTML and stored cleanly according to the necessary techniques, in a data structure, achieving this way to reliable dataset size and gleaner performance.

RECOVERY OF INFORMATION

The screenshot shows a detailed view of a research paper on arXiv. The title is "Clustering-based spatial interpolation of parametric post-processing models" by Sándor Baran and Mária Lakatos. The abstract discusses ensemble-based probabilistic forecasting and proposes a clustering-based interpolation technique. The page includes sections for "Statistics > Applications", "Submission history", and "Bibliographic Tools". On the right, there are links for "Access Paper" (PDF, HTML, Other Formats), "Current browse context" (stat.AP), "References & Citations" (NASAADS, Google Scholar, Semantic Scholar), and "Bookmark". The footer of the arXiv interface is visible at the bottom.

Figure 1.B.8 [The job page with all the data](#)

1.c. Store data in a structured format

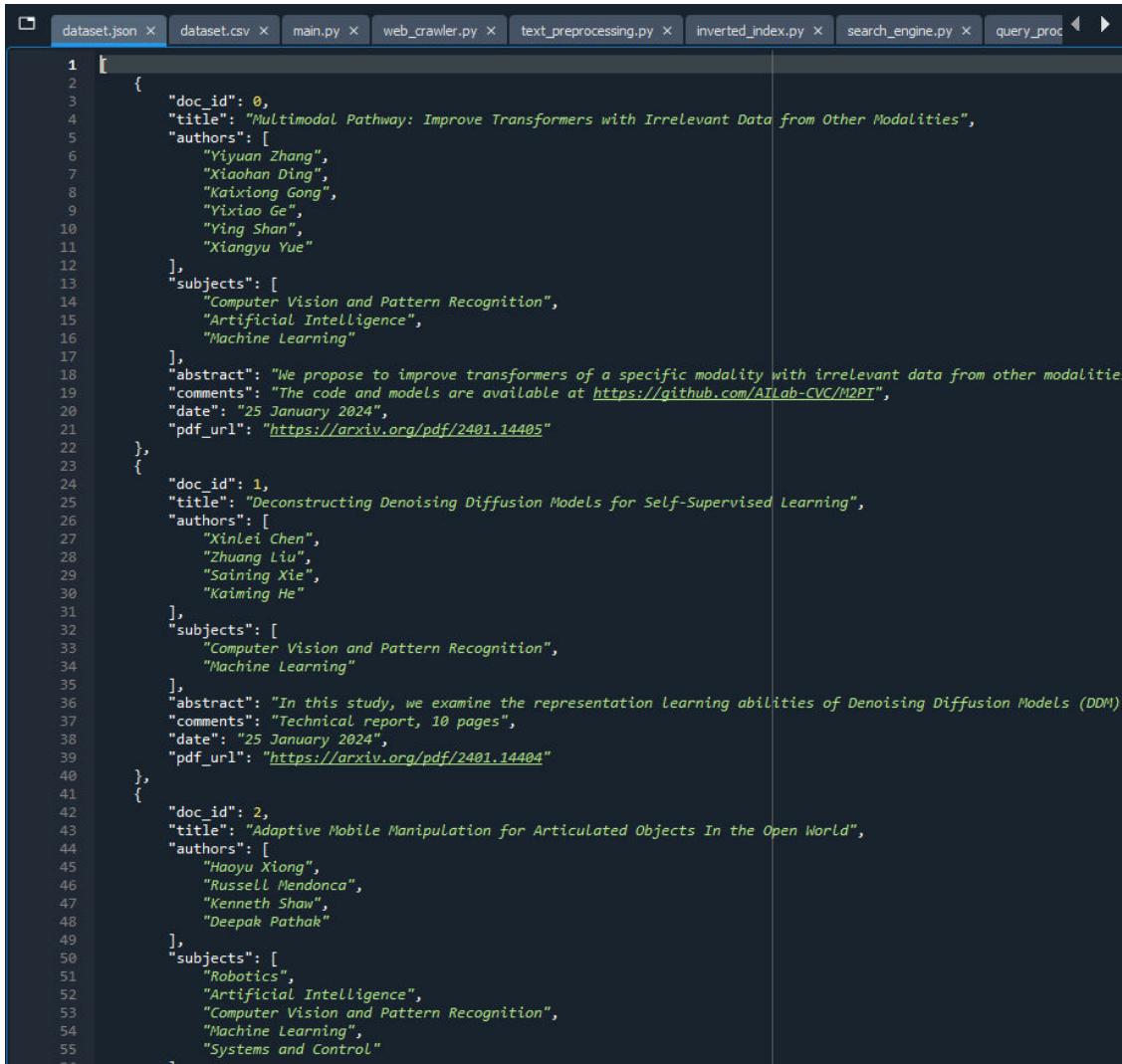
1.c.1 Planning

The design adopts the structured JSON format instead of the CSV format to store the task data, and by performing this step, the creation of the local search engine dataset is completed

1.c.2 Implementation

The collection of task metadata from the unstructured HTML format is completed in the [implementation](#) of step 1.b. Web crawler implementation, so what remains is to write the data in a .json file and the repository (dataset) of the local search engine is this. The implementation takes place in the main program (main.py), where the results returned by the gleaner in dictionary format are stored in JSON format.

RECOVERY OF INFORMATION



The screenshot shows a code editor with multiple tabs at the top: dataset.json, dataset.csv, main.py, web_crawler.py, text_preprocessing.py, inverted_index.py, search_engine.py, and query_proc. The dataset.json tab is active, displaying a JSON object with three documents (doc_id 0, 1, 2). Each document has fields: doc_id, title, authors, subjects, abstract, comments, date, and pdf_url. The titles and abstracts are in English, mentioning topics like "Multimodal Pathway", "Computer Vision and Pattern Recognition", and "Machine Learning". The code editor interface includes line numbers on the left and standard navigation buttons at the top right.

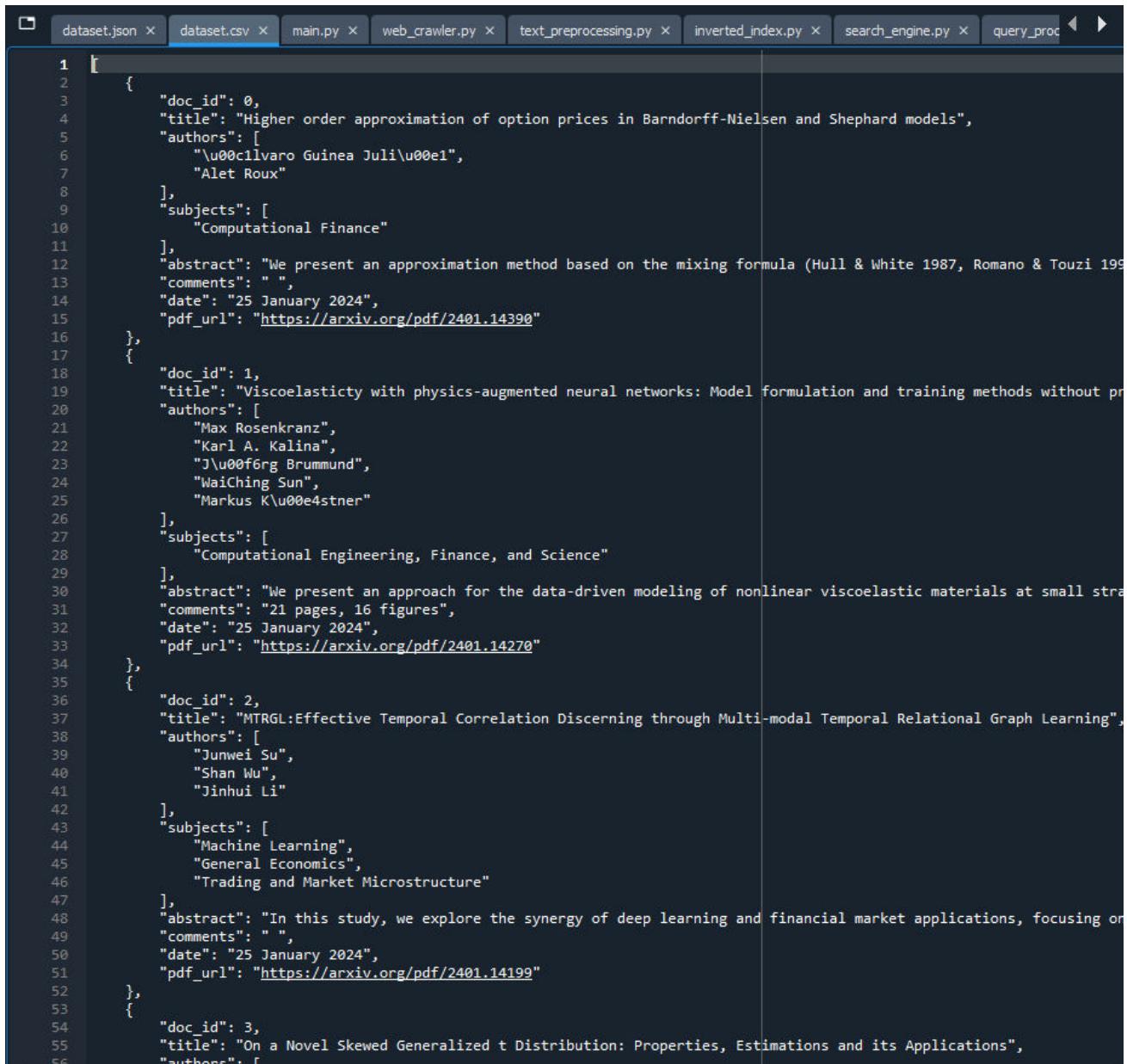
```
1  {
2      "doc_id": 0,
3          "title": "Multimodal Pathway: Improve Transformers with Irrelevant Data from Other Modalities",
4          "authors": [
5              "Yiyuan Zhang",
6              "Xiaohan Ding",
7              "Kaixiong Gong",
8              "Yixiao Ge",
9              "Ying Shan",
10             "Xiangyu Yue"
11         ],
12         "subjects": [
13             "Computer Vision and Pattern Recognition",
14             "Artificial Intelligence",
15             "Machine Learning"
16         ],
17         "abstract": "We propose to improve transformers of a specific modality with irrelevant data from other modalities",
18         "comments": "The code and models are available at https://github.com/AI-Lab-CVC/M2PT",
19         "date": "25 January 2024",
20         "pdf_url": "https://arxiv.org/pdf/2401.14405.pdf"
21     },
22     {
23         "doc_id": 1,
24         "title": "Deconstructing Denoising Diffusion Models for Self-Supervised Learning",
25         "authors": [
26             "Xinlei Chen",
27             "Zhuang Liu",
28             "Saining Xie",
29             "Kaiming He"
30         ],
31         "subjects": [
32             "Computer Vision and Pattern Recognition",
33             "Machine Learning"
34         ],
35         "abstract": "In this study, we examine the representation learning abilities of Denoising Diffusion Models (DDMs)",
36         "comments": "Technical report, 10 pages",
37         "date": "25 January 2024",
38         "pdf_url": "https://arxiv.org/pdf/2401.14404.pdf"
39     },
40     {
41         "doc_id": 2,
42         "title": "Adaptive Mobile Manipulation for Articulated Objects In the Open World",
43         "authors": [
44             "Haoyu Xiong",
45             "Russell Mendonca",
46             "Kenneth Shaw",
47             "Deepak Pathak"
48         ],
49         "subjects": [
50             "Robotics",
51             "Artificial Intelligence",
52             "Computer Vision and Pattern Recognition",
53             "Machine Learning",
54             "Systems and Control"
55         ]
56     }
57 }
```

Figure 1.C.1 The repository in JSON format

1.c.3 Evaluation

The choice of JSON format instead of CSV to store task data in a structured format is justified by the fact that there are data (authors, courses) that are more readable in a structure (e.g. list) that supports JSON rather than a single comma-separated text (CSV). The volume of information ranges from 200 to 800 tasks, where each task includes 8 fields of data (documentation in the [implementation](#) of step 1.b. Web-crawler implementation). Therefore, leveraging large amounts of data is more easily manageable in JSON format than in simple CSV format.

RECOVERY OF INFORMATION



The screenshot shows a code editor with multiple tabs at the top: dataset.json, dataset.csv, main.py, web_crawler.py, text_preprocessing.py, inverted_index.py, search_engine.py, and query_proc. The dataset.json tab is active, displaying a JSON object with five documents. Each document is represented by a key-value pair where the key is the document ID (e.g., 0, 1, 2, 3) and the value is a JSON object containing fields like title, authors, subjects, abstract, comments, date, and pdf_url. The JSON is formatted with line numbers on the left.

```
1 {
2     "doc_id": 0,
3     "title": "Higher order approximation of option prices in Barndorff-Nielsen and Shephard models",
4     "authors": [
5         "\u00c1lvaro Guinea Juli\u00e1",
6         "Alet Roux"
7     ],
8     "subjects": [
9         "Computational Finance"
10    ],
11    "abstract": "We present an approximation method based on the mixing formula (Hull & White 1987, Romano & Touzi 1990).",
12    "comments": " ",
13    "date": "25 January 2024",
14    "pdf_url": "https://arxiv.org/pdf/2401.14390"
15 },
16 {
17     "doc_id": 1,
18     "title": "Viscoelasticity with physics-augmented neural networks: Model formulation and training methods without prior knowledge",
19     "authors": [
20         "Max Rosenkranz",
21         "Karl A. Kalina",
22         "J\u00f6rg Brummund",
23         "WaiChing Sun",
24         "Markus K\u00f6stner"
25     ],
26     "subjects": [
27         "Computational Engineering, Finance, and Science"
28     ],
29     "abstract": "We present an approach for the data-driven modeling of nonlinear viscoelastic materials at small strains and large deformations using physics-augmented neural networks. The proposed model is able to capture complex material behavior from experimental data and predict it under different loading conditions without prior knowledge of the underlying physical processes. The model is trained using a combination of experimental data and numerical simulations, and its performance is evaluated using various validation criteria. The results show that the proposed model is able to predict the material behavior accurately and robustly, even for cases where the data is limited or noisy. The model is also able to handle complex loading paths and predict the material behavior under different loading conditions. The proposed model is a promising tool for the development of new materials and their applications in various engineering fields.",",
30     "comments": "21 pages, 16 figures",
31     "date": "25 January 2024",
32     "pdf_url": "https://arxiv.org/pdf/2401.14270"
33 },
34 {
35     "doc_id": 2,
36     "title": "MTRGL:Effective Temporal Correlation Discerning through Multi-modal Temporal Relational Graph Learning",
37     "authors": [
38         "Junwei Su",
39         "Shan Wu",
40         "Jinhui Li"
41     ],
42     "subjects": [
43         "Machine Learning",
44         "General Economics",
45         "Trading and Market Microstructure"
46     ],
47     "abstract": "In this study, we explore the synergy of deep learning and financial market applications, focusing on the MTRGL framework. This framework leverages multi-modal temporal relational graph learning to discern effective temporal correlations between financial time series. By integrating various types of data (e.g., price, volume, and sentiment) and their temporal dependencies, MTRGL can capture complex patterns that are often overlooked by traditional statistical methods. The proposed framework is demonstrated to be highly accurate and efficient in predicting future market movements, outperforming several baseline models. The results show that MTRGL can provide valuable insights for financial decision-making, such as risk management and portfolio optimization. The framework is also flexible and can be applied to other domains where temporal dependencies are important.",",
48     "comments": " ",
49     "date": "25 January 2024",
50     "pdf_url": "https://arxiv.org/pdf/2401.14199"
51 },
52 {
53     "doc_id": 3,
54     "title": "On a Novel Skewed Generalized t Distribution: Properties, Estimations and its Applications",
55     "authors": [
56 
```

Figure 1.C.2 The repository in CSV format

2. Text Processing

2.a. Planning

The design follows the preprocessing of the textual content (Abstract field) of the repository created in [step 1](#), with various text preprocessing techniques. Text preprocessing is one of the most basic natural language processing (NLP) techniques, where it achieves the matching of data terms with the search query, so that the ambiguity that is apparently present in the wording of the user query, does not affect the search result .

RECOVERY OF INFORMATION

Preprocessing was designed with the ultimate goal of grouping together different wording versions of a user query so that there is not much variance in search results when the same desired result is sought but with different wording. The design follows various text preprocessing techniques to preprocess the repository job data. The techniques designed for implementation are as follows:

- Tokenization: The separation of text into verbal units
- Punctuation characters removal: Removal of all punctuation marks
- Special characters removal: Removal of all special characters
- Normalization: The replacement of uppercase letters with their lowercase counterparts
- Stopwords removal: Removal of "forbidden" words
- Stemming: Removing endings and keeping the main stem of the term

2.b. Implementation

The implementations of the text preprocessing techniques take place in a separate module (text_preprocessing.py) and are as follows:

- Tokenization: Use nltk.word_tokenize() to split text into word units and separate punctuation
- Punctuation characters removal: Removal of all terms belonging to the set string.punctuation, that is, punctuation marks.
- Special characters removal: Removal of all terms that are not numbers or letters
- Normalization: Replace uppercase letters with their lowercase counterparts using the .lower() routine
- Stopwords removal: Remove all terms belonging to the set stopwords.words('english') from nltk.corpus. Usually, these are terms that appear frequently in texts such as articles, pronouns, adverbs, etc.
- Stemming: Remove endings and keep the main term stem using nltk.stem's PorterStemmer() algorithm

Given that a search query (query) will also be needed to search for jobs, this should also be pre-processed (details in [step 4 Search engine](#)). In search queries that include Boolean operations, the pre-processing technique "Stopwords removal" is not performed, as this will also remove the logical operators AND, NOT, OR from the query (details in [step 4.b.1 Boolean Retrieval](#)).

2.c. Evaluation

The pre-processing is done in the order presented in the [design](#) and [implementation](#). The evaluation is based on the choice of text preprocessing techniques, as well as the order in which they take place. The evaluations of the functions of each technique are as follows:

- Tokenization: The separation of the text into verbal units contributes to the management of words and punctuation marks and is therefore applied 1st.
- Punctuation characters removal: The removal of punctuation marks helps to clean the text from unnecessary information and is applied 2nd ^{for} this reason. However, this technique may affect the understanding of a text that uses punctuation to distinguish terms (eg by removing the punctuation on a date 19/03/2001, the preprocessed date is displayed as 3 numbers 19 03 2001 with resulting in a loss of information)

RECOVERY OF INFORMATION

- Special characters removal: Special characters removal also helps to clean the text of unnecessary information. However, as with the punctuation characters removal technique, it may affect the understanding of the text (e.g. by removing the special characters in an email jondoe@gmail.com, the preprocessed email is rendered as jondoe gmail com, resulting in loss of information)
- Normalization: Normalization helps bring terms into a common form by converting all uppercase letters to lowercase.
- Stopwords removal: Stopwords removal cleans the text of words that appear too often and do not provide important information. However, it may affect the understanding of the text in terms accompanied by forbidden words (eg, by removing the forbidden words in King of Denmark, the preprocessed text is rendered as king denmark, thereby losing information).
- Stemming: The process of 'stemming' removes the endings from words and keeps the main stem of the word, thus succeeding in grouping different versions of a word into a common stem. However, here too there is a risk of information loss (e.g. in the boolean search queries operational AND research, operating AND system, operative AND dentistry, the stemming process groups the terms operational, operating, operative into the common stem oper, resulting in a loss of information)

The lemmatization technique offers more accurate results, as it takes into account the grammar of the word to convert it to its basic form in contrast to the stemming technique which simply removes endings from words and presents a potential problem of information loss. However, given that the dataset is small compared to large search engine repositories (consisting of millions of documents), the stemming technique offers simplicity and speed in search engine performance, and this is why it was preferred over the lemmatization technique.

2.d. Application

For the needs of the text preprocessing application, the "abstract" field of the work was used with the following data:

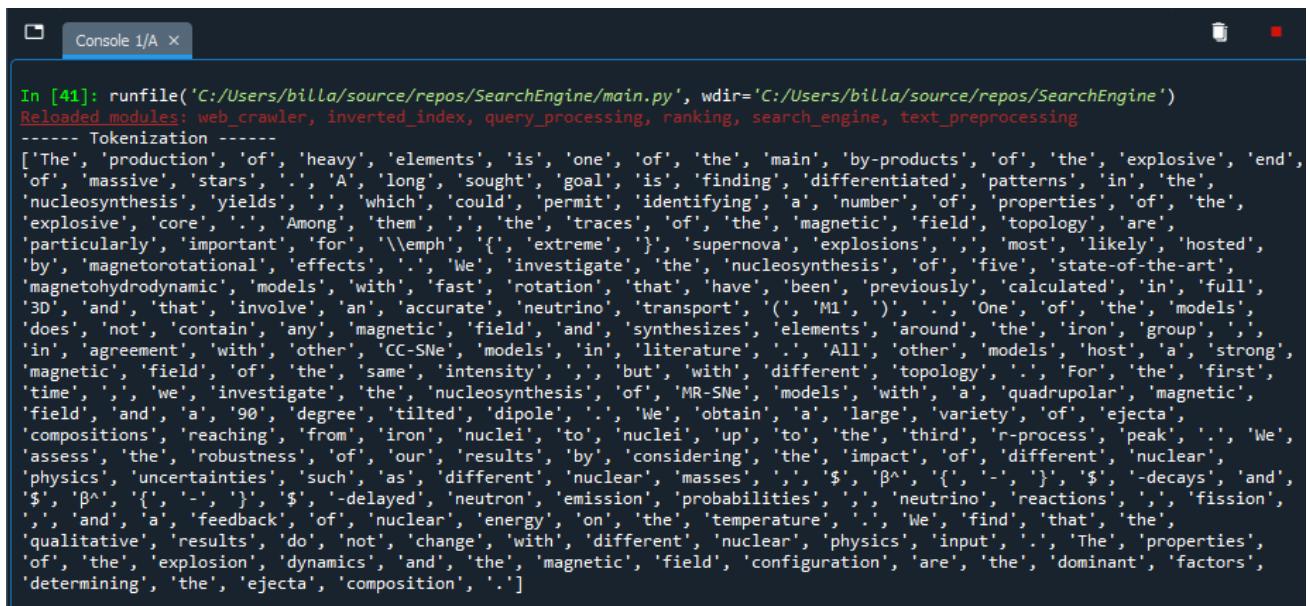
```
{  
  "doc_id": 0,  
  "title": "Nucleosynthesis in magnetorotational supernovae: impact of the magnetic field configuration",  
  "authors": [  
    "M. Reichert",  
    "M. Bugli",  
    "J. Guilet",  
    "M. Obergaulinger",  
    "M. \u00c1loy",  
    "A. Arcones"
```

RECOVERY OF INFORMATION

```
[],  
"subjects": [  
    "High Energy Astrophysical Phenomena"  
],
```

"abstract": "The production of heavy elements is one of the main by-products of the explosive end of massive stars. A long sought goal is finding differentiated patterns in the nucleosynthesis yields, which could permit identifying a number of properties of the explosive core. Among them, the traces of the magnetic field topology are particularly important for \emph{extreme} supernova explosions. We investigate the nucleosynthesis of five state-of-the-art magnetohydrodynamic models that have been previously calculated in full 3D and that involve an accurate neutrino transport (M1). One of the models does not contain any magnetic field and synthesizes elements around the iron group, in agreement with other CC-SNe models in literature models host a strong magnetic field of the same intensity, but with different topology. For the first time, we investigate the nucleosynthesis of MR-SNe models with a quadrupolar magnetic field and a 90 degree tilted dipole. We obtain a large variety of ejecta compositions reaching from iron nuclei to nuclei up to the third r-process peak. We assess the robustness of our results by considering the impact of different nuclear physics uncertainties such as different nuclear masses, \$\u03b2^{\{-\}}\$-decays and \$\u03b2^{\{-\}}\$-delayed neutron emission probabilities, neutrino reactions, fission, and a feedback of nuclear energy on the temperature. We find that the qualitative results do not change with different nuclear physics input. The properties of the explosion dynamics and the magnetic field configuration are the dominant factors determining the ejecta composition.",

```
"comments": " ",  
"date": "25 January 2024",  
"pdf_url": https://arxiv.org/pdf/2401.14402  
}
```



The screenshot shows a Jupyter Notebook cell with the following content:

```
In [41]: runfile('C:/Users/billa/source/repos/SearchEngine/main.py', wdir='C:/Users/billa/source/repos/SearchEngine')  
Reloaded modules: web_crawler, inverted_index, query_processing, ranking, search_engine, text_preprocessing  
----- Tokenization -----  
['The', 'production', 'of', 'heavy', 'elements', 'is', 'one', 'of', 'the', 'main', 'by-products', 'of', 'the', 'explosive', 'end',  
'of', 'massive', 'stars', '.', 'A', 'long', 'sought', 'goal', 'is', 'finding', 'differentiated', 'patterns', 'in', 'the',  
'nucleosynthesis', 'yields', ',', 'which', 'could', 'permit', 'identifying', 'a', 'number', 'of', 'properties', 'of', 'the',  
'explosive', 'core', '.', 'Among', 'them', ',', 'the', 'traces', 'of', 'the', 'magnetic', 'field', 'topology', 'are',  
'particularly', 'important', 'for', '\emph', '{}', 'extreme', '{}', 'supernova', 'explosions', ',', 'most', 'likely', 'hosted',  
'by', 'magnetorotational', 'effects', '.', 'We', 'investigate', 'the', 'nucleosynthesis', 'of', 'five', 'state-of-the-art',  
'magnetohydrodynamic', 'models', 'with', 'fast', 'rotation', 'that', 'have', 'been', 'previously', 'calculated', 'in', 'full',  
'3D', 'and', 'that', 'involve', 'an', 'accurate', 'neutrino', 'transport', '(', 'M1', ')', ',', 'One', 'of', 'the', 'models',  
'does', 'not', 'contain', 'any', 'magnetic', 'field', 'and', 'synthesizes', 'elements', 'around', 'the', 'iron', 'group', ',',  
'in', 'agreement', 'with', 'other', 'CC-SNe', 'models', 'in', 'literature', ',', 'All', 'other', 'models', 'host', 'a', 'strong',  
'magnetic', 'field', 'of', 'the', 'same', 'intensity', ',', 'but', 'with', 'different', 'topology', ',', 'For', 'the', 'first',  
'time', ',', 'we', 'investigate', 'the', 'nucleosynthesis', 'of', 'MR-SNe', 'models', 'with', 'a', 'quadrupolar', 'magnetic',  
'field', 'and', 'a', '90', 'degree', 'tilted', 'dipole', ',', 'We', 'obtain', 'a', 'large', 'variety', 'of', 'ejecta',  
'compositions', 'reaching', 'from', 'iron', 'nuclei', 'to', 'nuclei', 'up', 'to', 'the', 'third', 'r-process', 'peak', '.', 'We',  
'assess', 'the', 'robustness', 'of', 'our', 'results', 'by', 'considering', 'the', 'impact', 'of', 'different', 'nuclear',  
'physics', 'uncertainties', 'such', 'as', 'different', 'nuclear', 'masses', ',', '$', 'B^{\{-\}}', '{}', '$', '-decays', 'and',  
'$', 'B^{\{-\}}', '{}', '$', '-delayed', 'neutron', 'emission', 'probabilities', ',', 'neutrino', 'reactions', ',', 'fission',  
, 'and', 'a', 'feedback', 'of', 'nuclear', 'energy', 'on', 'the', 'temperature', ',', 'We', 'find', 'that', 'the',  
'qualitative', 'results', 'do', 'not', 'change', 'with', 'different', 'nuclear', 'physics', 'input', '.', 'The', 'properties',  
'of', 'the', 'explosion', 'dynamics', 'and', 'the', 'magnetic', 'field', 'configuration', 'are', 'the', 'dominant', 'factors',  
'determining', 'the', 'ejecta', 'composition', '.']
```

Figure 2.D.1 Tokenization

RECOVERY OF INFORMATION

In [43]: runfile('C:/Users/billa/source/repos/SearchEngine/main.py', wdir='C:/Users/billa/source/repos/SearchEngine')
Reloaded modules: web_crawler, inverted_index, query_processing, ranking, search_engine, text_preprocessing
----- Punctuation characters removal -----
['The', 'production', 'of', 'heavy', 'elements', 'is', 'one', 'of', 'the', 'main', 'by-products', 'of', 'the', 'explosive', 'end', 'of', 'massive', 'stars', 'A', 'long', 'sought', 'goal', 'is', 'finding', 'differentiated', 'patterns', 'in', 'the', 'nucleosynthesis', 'yields', 'which', 'could', 'permit', 'identifying', 'a', 'number', 'of', 'properties', 'of', 'the', 'explosive', 'core', 'Among', 'them', 'the', 'traces', 'of', 'the', 'magnetic', 'field', 'topology', 'are', 'particularly', 'important', 'for', '\u2192', 'extreme', 'supernova', 'explosions', 'most', 'likely', 'hosted', 'by', 'magnetorotational', 'effects', 'We', 'investigate', 'the', 'nucleosynthesis', 'of', 'five', 'state-of-the-art', 'magnetohydrodynamic', 'models', 'with', 'fast', 'rotation', 'that', 'have', 'been', 'previously', 'calculated', 'in', 'full', '3D', 'and', 'that', 'involve', 'an', 'accurate', 'neutrino', 'transport', 'M1', 'One', 'of', 'the', 'models', 'does', 'not', 'contain', 'any', 'magnetic', 'field', 'and', 'synthesizes', 'elements', 'around', 'the', 'iron', 'group', 'in', 'agreement', 'with', 'other', 'CC-SNe', 'models', 'in', 'literature', 'All', 'other', 'models', 'host', 'a', 'strong', 'magnetic', 'field', 'of', 'the', 'same', 'intensity', 'but', 'with', 'different', 'topology', 'For', 'the', 'first', 'time', 'we', 'investigate', 'the', 'nucleosynthesis', 'of', 'MR-SNe', 'models', 'with', 'a', 'quadrupolar', 'magnetic', 'field', 'and', 'a', '90', 'degree', 'tilted', 'dipole', 'We', 'obtain', 'a', 'large', 'variety', 'of', 'ejecta', 'compositions', 'reaching', 'from', 'iron', 'nuclei', 'to', 'nuclei', 'up', 'to', 'the', 'third', 'r-process', 'peak', 'We', 'assess', 'the', 'robustness', 'of', 'our', 'results', 'by', 'considering', 'the', 'impact', 'of', 'different', 'nuclear', 'physics', 'uncertainties', 'such', 'as', 'different', 'nuclear', 'masses', ' β^+ ', '-decays', 'and', ' β^- ', '-delayed', 'neutron', 'emission', 'probabilities', 'neutrino', 'reactions', 'fission', 'and', 'a', 'feedback', 'of', 'nuclear', 'energy', 'on', 'the', 'temperature', 'We', 'find', 'that', 'the', 'qualitative', 'results', 'do', 'not', 'change', 'with', 'different', 'nuclear', 'physics', 'input', 'The', 'properties', 'of', 'the', 'explosion', 'dynamics', 'and', 'the', 'magnetic', 'field', 'configuration', 'are', 'the', 'dominant', 'factors', 'determining', 'the', 'ejecta', 'composition']

Figure 2.D.2 Punctuation characters removal

In [45]: runfile('C:/Users/billa/source/repos/SearchEngine/main.py', wdir='C:/Users/billa/source/repos/SearchEngine')
Reloaded modules: web_crawler, inverted_index, query_processing, ranking, search_engine, text_preprocessing
----- Special characters removal -----
['The', 'production', 'of', 'heavy', 'elements', 'is', 'one', 'of', 'the', 'main', 'byproducts', 'of', 'the', 'explosive', 'end', 'of', 'massive', 'stars', 'A', 'long', 'sought', 'goal', 'is', 'finding', 'differentiated', 'patterns', 'in', 'the', 'nucleosynthesis', 'yields', 'which', 'could', 'permit', 'identifying', 'a', 'number', 'of', 'properties', 'of', 'the', 'explosive', 'core', 'Among', 'them', 'the', 'traces', 'of', 'the', 'magnetic', 'field', 'topology', 'are', 'particularly', 'important', 'for', '\u2192', 'extreme', 'supernova', 'explosions', 'most', 'likely', 'hosted', 'by', 'magnetorotational', 'effects', 'We', 'investigate', 'the', 'nucleosynthesis', 'of', 'five', 'stateoftheart', 'magnetohydrodynamic', 'models', 'with', 'fast', 'rotation', 'that', 'have', 'been', 'previously', 'calculated', 'in', 'full', '3D', 'and', 'that', 'involve', 'an', 'accurate', 'neutrino', 'transport', 'M1', 'One', 'of', 'the', 'models', 'does', 'not', 'contain', 'any', 'magnetic', 'field', 'and', 'synthesizes', 'elements', 'around', 'the', 'iron', 'group', 'in', 'agreement', 'with', 'other', 'CCSNe', 'models', 'in', 'literature', 'All', 'other', 'models', 'host', 'a', 'strong', 'magnetic', 'field', 'of', 'the', 'same', 'intensity', 'but', 'with', 'different', 'topology', 'For', 'the', 'first', 'time', 'we', 'investigate', 'the', 'nucleosynthesis', 'of', 'MRSNe', 'models', 'with', 'a', 'quadrupolar', 'magnetic', 'field', 'and', 'a', '90', 'degree', 'tilted', 'dipole', 'We', 'obtain', 'a', 'large', 'variety', 'of', 'ejecta', 'compositions', 'reaching', 'from', 'iron', 'nuclei', 'to', 'nuclei', 'up', 'to', 'the', 'third', 'rprocess', 'peak', 'We', 'assess', 'the', 'robustness', 'of', 'our', 'results', 'by', 'considering', 'the', 'impact', 'of', 'different', 'nuclear', 'physics', 'uncertainties', 'such', 'as', 'different', 'nuclear', 'masses', 'decays', 'and', 'delayed', 'neutron', 'emission', 'probabilities', 'neutrino', 'reactions', 'fission', 'and', 'a', 'feedback', 'of', 'nuclear', 'energy', 'on', 'the', 'temperature', 'We', 'find', 'that', 'the', 'qualitative', 'results', 'do', 'not', 'change', 'with', 'different', 'nuclear', 'physics', 'input', 'The', 'properties', 'of', 'the', 'explosion', 'dynamics', 'and', 'the', 'magnetic', 'field', 'configuration', 'are', 'the', 'dominant', 'factors', 'determining', 'the', 'ejecta', 'composition']

Figure 2.D.3 Special characters removal

RECOVERY OF INFORMATION

```
In [47]: runfile('C:/Users/billa/source/repos/SearchEngine/main.py', wdir='C:/Users/billa/source/repos/SearchEngine')
Reloaded modules: web_crawler, inverted_index, query_processing, ranking, search_engine, text_preprocessing
----- Normalization -----
['the', 'production', 'of', 'heavy', 'elements', 'is', 'one', 'of', 'the', 'main', 'byproducts', 'of', 'the', 'explosive', 'end',
'of', 'massive', 'stars', 'a', 'long', 'sought', 'goal', 'is', 'finding', 'differentiated', 'patterns', 'in', 'the',
'nucleosynthesis', 'yields', 'which', 'could', 'permit', 'identifying', 'a', 'number', 'of', 'properties', 'of', 'the',
'explosive', 'core', 'among', 'them', 'the', 'traces', 'of', 'the', 'magnetic', 'field', 'topology', 'are', 'particularly',
'important', 'for', 'emph', 'extreme', 'supernova', 'explosions', 'most', 'likely', 'hosted', 'by', 'magnetorotational', 'effects',
've', 'investigate', 'the', 'nucleosynthesis', 'of', 'five', 'stateoftheheart', 'magnetohydrodynamic', 'models', 'with', 'fast',
'rotation', 'that', 'have', 'been', 'previously', 'calculated', 'in', 'full', '3d', 'and', 'that', 'involve', 'an', 'accurate',
'neutrino', 'transport', 'm1', 'one', 'of', 'the', 'models', 'does', 'not', 'contain', 'any', 'magnetic', 'field', 'and',
'synthesizes', 'elements', 'around', 'the', 'iron', 'group', 'in', 'agreement', 'with', 'other', 'ccsne', 'models', 'in',
'literature', 'all', 'other', 'models', 'host', 'a', 'strong', 'magnetic', 'field', 'of', 'the', 'same', 'intensity', 'but',
'with', 'different', 'topology', 'for', 'the', 'first', 'time', 'we', 'investigate', 'the', 'nucleosynthesis', 'of', 'mrsne',
'models', 'with', 'a', 'quadrupolar', 'magnetic', 'field', 'and', 'a', '90', 'degree', 'tilted', 'dipole', 'we', 'obtain', 'a',
'large', 'variety', 'of', 'ejecta', 'compositions', 'reaching', 'from', 'iron', 'nuclei', 'to', 'nuclei', 'up', 'to', 'the',
'third', 'rprocess', 'peak', 'we', 'assess', 'the', 'robustness', 'of', 'our', 'results', 'by', 'considering', 'the', 'impact',
'of', 'different', 'nuclear', 'physics', 'uncertainties', 'such', 'as', 'different', 'nuclear', 'masses', 'decays', 'and',
'delayed', 'neutron', 'emission', 'probabilities', 'neutrino', 'reactions', 'fission', 'and', 'a', 'feedback', 'of', 'nuclear',
'energy', 'on', 'the', 'temperature', 'we', 'find', 'that', 'the', 'qualitative', 'results', 'do', 'not', 'change', 'with',
'different', 'nuclear', 'physics', 'input', 'the', 'properties', 'of', 'the', 'explosion', 'dynamics', 'and', 'the', 'magnetic',
'field', 'configuration', 'are', 'the', 'dominant', 'factors', 'determining', 'the', 'ejecta', 'composition']
```

Figure 2.D.4 Normalization

```
In [49]: runfile('C:/Users/billa/source/repos/SearchEngine/main.py', wdir='C:/Users/billa/source/repos/SearchEngine')
Reloaded modules: web_crawler, inverted_index, query_processing, ranking, search_engine, text_preprocessing
----- Stop-words removal -----
['production', 'heavy', 'elements', 'one', 'main', 'byproducts', 'explosive', 'end', 'massive', 'stars', 'long', 'sought', 'goal',
'finding', 'differentiated', 'patterns', 'nucleosynthesis', 'yields', 'could', 'permit', 'identifying', 'number', 'properties',
'explosive', 'core', 'among', 'traces', 'magnetic', 'field', 'topology', 'particularly', 'important', 'emph', 'extreme',
'supernova', 'explosions', 'likely', 'hosted', 'magnetorotational', 'effects', 'investigate', 'nucleosynthesis', 'five',
'stateoftheheart', 'magnetohydrodynamic', 'models', 'fast', 'rotation', 'previously', 'calculated', 'full', '3d', 'involve',
'accurate', 'neutrino', 'transport', 'm1', 'one', 'models', 'contain', 'magnetic', 'field', 'synthesizes', 'elements', 'around',
'iron', 'group', 'agreement', 'ccsne', 'models', 'literature', 'models', 'host', 'strong', 'magnetic', 'field', 'intensity',
'different', 'topology', 'first', 'time', 'investigate', 'nucleosynthesis', 'mrsne', 'models', 'quadrupolar', 'magnetic', 'field',
'90', 'degree', 'tilted', 'dipole', 'obtain', 'large', 'variety', 'ejecta', 'compositions', 'reaching', 'iron', 'nuclei', 'nuclei',
'third', 'rprocess', 'peak', 'assess', 'robustness', 'results', 'considering', 'impact', 'different', 'nuclear', 'physics',
'uncertainties', 'different', 'nuclear', 'masses', 'decays', 'delayed', 'neutron', 'emission', 'probabilities', 'neutrino',
'reactions', 'fission', 'feedback', 'nuclear', 'energy', 'temperature', 'find', 'qualitative', 'results', 'change', 'different',
'nuclear', 'physics', 'input', 'properties', 'explosion', 'dynamics', 'magnetic', 'field', 'configuration', 'dominant', 'factors',
'determining', 'ejecta', 'composition']
```

Figure 2.D.5 Stop-words removal

```
In [51]: runfile('C:/Users/billa/source/repos/SearchEngine/main.py', wdir='C:/Users/billa/source/repos/SearchEngine')
Reloaded modules: web_crawler, inverted_index, query_processing, ranking, search_engine, text_preprocessing
----- Stemming -----
['product', 'heavi', 'element', 'one', 'main', 'byproduct', 'explos', 'end', 'massiv', 'star', 'long', 'sought', 'goal', 'find',
'differenti', 'pattern', 'nucleosynthesi', 'yield', 'could', 'permit', 'identifi', 'number', 'properti', 'explos', 'core', 'among',
'trace', 'magnet', 'field', 'topolog', 'particularli', 'import', 'emph', 'extrem', 'supernova', 'explos', 'like', 'host',
'magnetrot', 'effect', 'investig', 'nucleosynthesi', 'five', 'stateoftheheart', 'magnetohydrodynan', 'model', 'fast', 'rotat',
'previous', 'calcul', 'full', '3d', 'involv', 'accur', 'neutrino', 'transport', 'm1', 'one', 'model', 'contain', 'magnet', 'field',
'synthes', 'element', 'around', 'iron', 'group', 'agreement', 'ccsne', 'model', 'literatur', 'model', 'host', 'strong', 'magnet',
'field', 'intens', 'differ', 'topolog', 'first', 'time', 'investig', 'nucleosynthesi', 'mrsne', 'model', 'quadrupolar', 'magnet',
'field', '90', 'degre', 'tilt', 'dipol', 'obtain', 'larg', 'varieti', 'ejecta', 'composit', 'reach', 'iron', 'nuclei', 'nuclei',
'third', 'rprocess', 'peak', 'assess', 'robust', 'result', 'consid', 'impact', 'differ', 'nuclear', 'physic', 'uncertainiti',
'differ', 'nuclear', 'mass', 'decay', 'delay', 'neutron', 'emiss', 'probabl', 'neutrino', 'reaction', 'fission', 'feedback',
'nuclear', 'energi', 'temperatur', 'find', 'qualit', 'result', 'chang', 'differ', 'nuclear', 'physic', 'input', 'properti',
'explos', 'dynam', 'magnet', 'field', 'configur', 'domin', 'factor', 'determin', 'ejecta', 'composit']
```

Figure 2.D.6 Stemming (instead of Lemmatization)

RECOVERY OF INFORMATION



The screenshot shows a Jupyter Notebook cell output. The code runs a file named main.py with parameters wdir='C:/Users/billa/source/repos/SearchEngine' and reloads modules like web_crawler, inverted_index, query_processing, ranking, search_engine, and text_preprocessing. It then prints a list of words under the heading 'Lemmatization'. The list contains numerous words from scientific and technical fields, such as 'production', 'heavy', 'element', 'one', 'main', 'byproduct', 'explosive', 'end', 'massive', 'star', 'long', 'sought', 'goal', 'finding', 'differentiated', 'pattern', 'nucleosynthesis', 'yield', 'could', 'permit', 'identifying', 'number', 'property', 'explosive', 'core', 'among', 'trace', 'magnetic', 'field', 'topology', 'particularly', 'important', 'emph', 'extreme', 'supernova', 'explosion', 'likely', 'hosted', 'magnetorotational', 'effect', 'investigate', 'nucleosynthesis', 'five', 'stateoftheart', 'magnetohydrodynamic', 'model', 'fast', 'rotation', 'previously', 'calculated', 'full', '3d', 'involve', 'accurate', 'neutrino', 'transport', 'ml', 'one', 'model', 'contain', 'magnetic', 'field', 'synthesizes', 'element', 'around', 'iron', 'group', 'agreement', 'ccsne', 'model', 'literature', 'model', 'host', 'strong', 'magnetic', 'field', 'intensity', 'different', 'topology', 'first', 'time', 'investigate', 'nucleosynthesis', 'mrsne', 'model', 'quadrupolar', 'magnetic', 'field', '90', 'degree', 'tilted', 'dipole', 'obtain', 'large', 'variety', 'ejecta', 'composition', 'reaching', 'iron', 'nucleus', 'nucleus', 'third', 'rprocess', 'peak', 'ass', 'robustness', 'result', 'considering', 'impact', 'different', 'nuclear', 'physic', 'uncertainty', 'different', 'nuclear', 'mass', 'decay', 'delayed', 'neutron', 'emission', 'probability', 'neutrino', 'reaction', 'fission', 'feedback', 'nuclear', 'energy', 'temperature', 'find', 'qualitative', 'result', 'change', 'different', 'nuclear', 'physic', 'input', 'property', 'explosion', 'dynamic', 'magnetic', 'field', 'configuration', 'dominant', 'factor', 'determining', 'ejecta', 'composition']

Figure 2.D.7 Lemmatization (instead of Stemming)

3. Indexing

3.a. Creating the inverted index data structure

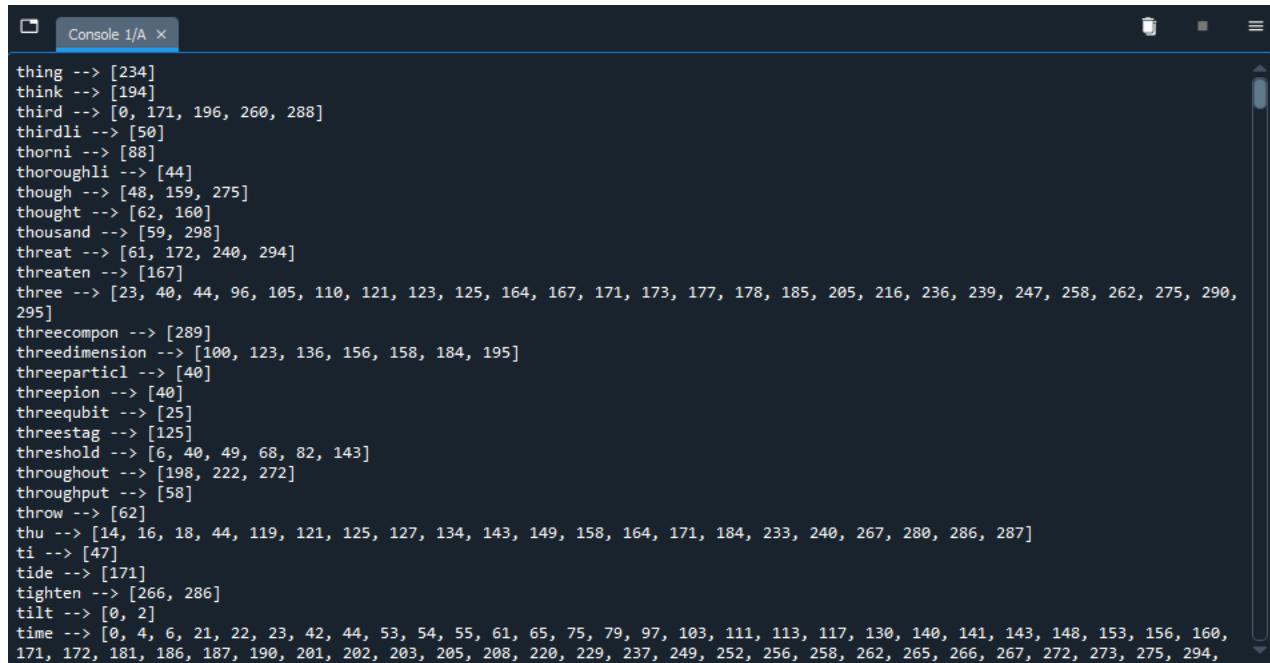
3.a.1 Planning

The design of the inverted index data structure consists of the collection of tasks for indexing ([step 1](#)), the conversion of the textual content of each task into a list of elements (tokenization) as well as their preprocessing ([step 2](#)). Finally, it consists of the indexing of the works containing the terms, that is, it matches the terms (tokens) that have been collected with the specific texts in which they have been found. Virtually every word corresponds to a number which in turn shows the text in which the word is found. For example, if a word corresponding to the numbers 1,3,7 we know that it is in the texts with doc_id 1,3,7.

RECOVERY OF INFORMATION

3.a.2 Implementation

To implement the inverted index data structure it was necessary to initialize an empty inverted index, then to separate the words from the summaries for each summation separately into word units (tokenization). All preprocessed words are then added to the index, and if a word already exists in it, the id of the task in which the word is found is simply added to the index keyword, thus avoiding duplicate concepts within the index. Finally, the keys of the dictionary and its elements are sorted.



```
Console 1/A ×

thing --> [234]
think --> [194]
third --> [0, 171, 196, 260, 288]
thirdli --> [50]
thorni --> [88]
thoroughli --> [44]
though --> [48, 159, 275]
thought --> [62, 160]
thousand --> [59, 298]
threat --> [61, 172, 240, 294]
threaten --> [167]
three --> [23, 40, 44, 96, 105, 110, 121, 123, 125, 164, 167, 171, 173, 177, 178, 185, 205, 216, 236, 239, 247, 258, 262, 275, 290, 295]
threecompon --> [289]
threedimension --> [100, 123, 136, 156, 158, 184, 195]
threeparticl --> [40]
threepion --> [40]
threequbit --> [25]
threestag --> [125]
threshold --> [6, 40, 49, 68, 82, 143]
throughout --> [198, 222, 272]
throughput --> [58]
throw --> [62]
thu --> [14, 16, 18, 44, 119, 121, 125, 127, 134, 143, 149, 158, 164, 171, 184, 233, 240, 267, 280, 286, 287]
ti --> [47]
tide --> [171]
tighten --> [266, 286]
tilt --> [0, 2]
time --> [0, 4, 6, 21, 22, 23, 42, 44, 53, 54, 55, 61, 65, 75, 79, 97, 103, 111, 113, 117, 130, 140, 141, 143, 148, 153, 156, 160, 171, 172, 181, 186, 187, 190, 201, 202, 203, 205, 208, 220, 229, 237, 249, 252, 256, 258, 262, 265, 266, 267, 272, 273, 275, 294,
```

Figure 3.A.1 The inverted index data structure

3.a.3 Evaluation

The inverted index data structure greatly improves search terms in a text and is therefore a powerful tool for search engines. Matching the terms to the texts in which they appear creates an important picture of the frequency of the terms in the collection with the papers, which reinforces how important information a term is to the collection. The optimal indexing in terms of efficiency would be to list the position(s) of the term within the text, so that searches are more efficient and accurate, as in this way more information is provided to the user. However, this methodology was not chosen when [implementing](#) the local search engine index, as the repository is small and does not need such an implementation that spends memory space and search time due to its large size.

3.b. Store the index in a data structure

3.b.1 Planning

RECOVERY OF INFORMATION

The creation of the inverted index data structure has been completed in [step 3.a](#), so what remains is to store it in a Python language data structure that will be efficient in terms of searching. The design model boils down to the choice of dictionary structure to store the index.

3.b.2 Implementation

The structure of the dictionary includes two fields, the keys (keys) and the contents (values). In the local search engine implementation, keywords are all the unique terms of the collection of tasks (terms), while contents are lists (posting lists) with the unique integers (doc_id) that denote the texts in which the term appears-key.

3.b.3 Evaluation

The choice of the dictionary structure for storing the index is evaluated by the fact that it yields search efficiency by matching key terms to the lists of texts in which they appear. However, there may be significant memory consumption due to the large index size resulting from the number of unique keywords. However, the requirements of the local search engine are more about the accuracy and efficiency of the search than about saving space. By choosing another structure (eg table, list, tuple, etc.) there would be a significant improvement in space, but there would be a significant impact on the search result. In closing, the repository is small so memory consumption would not be so much that it would cause a significant problem in search engine performance.

4. Search Engine

4.a. Develop user interface for job search

4.a.1 Planning

In order to make the search engine more user friendly a user interface (GUI) was designed to search the tasks. More specifically, this interface contains:

- User interface window
- User query input field
- Retrieval Algorithm Selection Field (Boolean Retrieval, Vector Space Model, Probabilistic Retrieval Model)
- Search recovery algorithm results button
- Filter criteria input field
- Filter selection field

RECOVERY OF INFORMATION

- Search filter results button

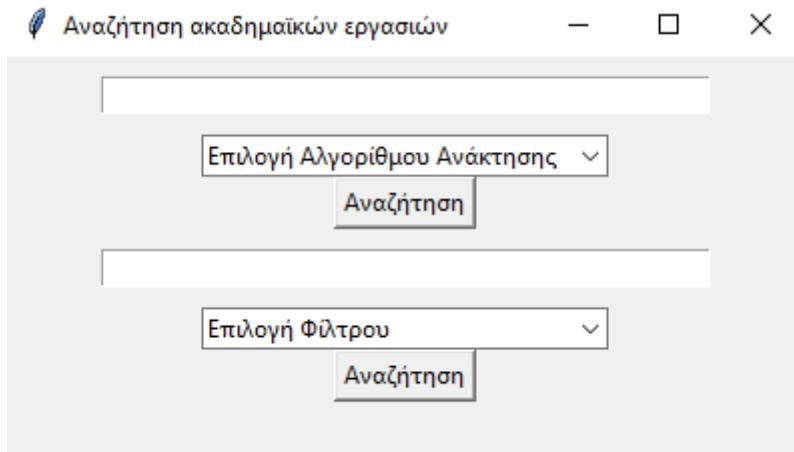


Figure 4.A.1 The user interface

4.a.2 Implementation

The implementation of the user interface was defined in a separate module (`search_engine.py`) and specifically the `tkinter` library was used. The interface serves the two basic functions of searching for works based on a user query and filtering results based on specific criteria, such as author or publication date. The interface and task search functionality are implemented in a class with properties the repository with preprocessed fields "abstract" ([step 2](#)) and the inverted index ([step 3](#)).

The implementations of the interface components mentioned in the [design](#) are as follows:

- User query input field: The user enters the query
- Retrieval Algorithm Selection Field: The user selects from a menu one of the 3 retrieval algorithms (Boolean Retrieval, Vector Space Model, Probabilistic Retrieval Model) to be applied to retrieve tasks
- Search Retrieval Algorithm Results Button: Clicking the Search button displays the query results in the terminal window, based on the retrieval algorithm selected by the user
- Filter criterion selection field: The user selects from a menu (if desired) one of the 2 result filtering criteria (Author, Publication date)
- Filter criteria input field: User enters the filter (author name or publication date)
- Filter Results Search Button: Pressing the "Search" button displays the filtering results in the terminal window, based on the filter criteria selected by the user

4.a.3 Evaluation

The interface, although simple, meets the criteria so that the user can easily search for jobs based on certain words that he sets as a search query. It just as easily allows filtering based on specific criteria (published date, author), so there's nothing else in the interface that isn't there as a user-facing function. The interface takes care of the correct entry of data into the search engine, data processing and possible filtering. The usefulness of this class is catalytic, since closing the interface also terminates the program. Perhaps, it would be preferable

RECOVERY OF INFORMATION

not to have so many functions concentrated in the interface class purely as a matter of code flexibility and readability .

4.b. Implementation of recovery algorithms

4.b.1 Planning

Boolean Retrieval

The Boolean Retrieval algorithm implemented in the search engine is based on Boolean algebra and set theory, using three basic keywords: OR, AND and NOT. It follows the process of [Query Processing](#). These functions affect the search results according to their position in the query.

- The keyword NOT, when placed before a word, searches for texts that do not contain that word.
- The AND keyword, when placed between two words, searches for texts that contain both words.
- The OR keyword, when placed between two words, searches for text that contains either word.

In addition, the algorithm takes into account the priority of operations in complex Boolean user queries, performing the operations in this order of priority:

1. Operations with the logical NOT operators in parentheses (from left to right)
2. Operations with the AND or OR operators inside parentheses (from left to right)
3. Operations with the logical NOT operators (from left to right)
4. Operations with the AND or OR operators (from left to right)

In this way, the algorithm can be adapted to respond to complex search queries.

Vector Space Model

The Vector Space Model algorithm implemented in the search engine represents texts and user queries as vectors in an n-dimensional space. From the [ranking algorithms](#) it follows TF-IDF and Cosine Similarity to calculate the similarity between the user query and the collection texts. Due to the problem with large Euclidean distances of vectors, calculating the cosine of the angle that the query vector forms with the vector of some text in the collection is more efficient. The closer the vectors are, the closer the query is to the text, that is, the cosine of the angle tends to 1, where 1 means identical.

Probabilistic Retrieval Model

The Probabilistic Retrieval Model algorithm implemented in the search engine searches for jobs based on the probabilistic model. It essentially calculates the probability that a text is relevant to the user query. From the [ranking algorithms](#) it follows TF-IDF and Okapi BM25 to calculate the probabilities.

RECOVERY OF INFORMATION

4.b.2 Implementation

Boolean Retrieval

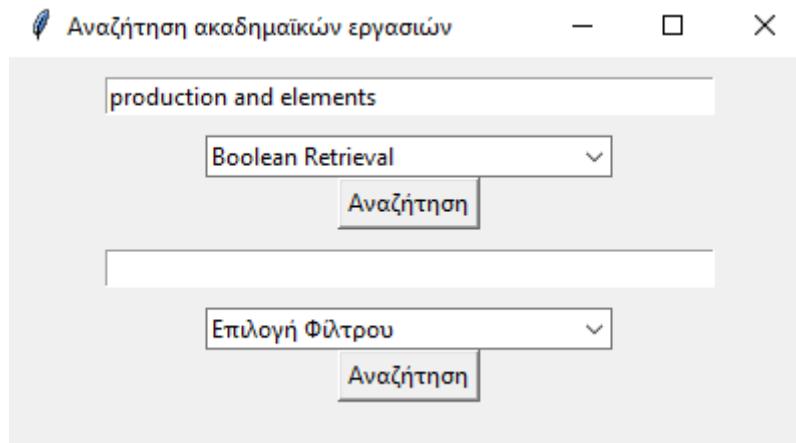
The implementation of Boolean Retrieval was carried out in the same module (`search_engine.py`) as the [user interface](#) and specifically in the same class. As input, the algorithm takes the following data:

- The Boolean query that the user will provide to perform the search
- The [inverted index](#) that will do the necessary matching between words and texts

The function of the algorithm is to create a stack of simple Boolean sub-queries of the complex or simple Boolean user query. The position of each sub-query in the stack is determined by the priority order of the operations mentioned in the [design](#). For each sub-query, the operation from the [query processing process is performed](#). The first concern of the algorithm is whether there are parentheses. In the case of parentheses, the algorithm allocates the logical operations within the parentheses to the stack first, and then any remaining operations. If there are no parentheses in the user's query, the search is performed normally from left to right.

Each word in the query must be properly examined to check if it corresponds to a keyword (and, or, not) or a word found within the index. For this reason, the terms of the query must become tokens, through which it will be possible to identify possible keywords that will allow the correct matching between words of the query and words found in the index.

Some illustrative case studies from user queries are as follows:



RECOVERY OF INFORMATION

```
In [17]: runfile('C:/Users/billa/source/repos/SearchEngine/main.py', wdir='C:/Users/billa/source/repos/SearchEngine')
Reloaded modules: web_crawler, text_preprocessing, inverted_index, query_processing, ranking, search_engine

===== Ερύτημα αναζήτησης : production and elements =====
===== Αλγόριθμος ανάκτησης : Boolean Retrieval =====
-----
#1

Document ID : 0
Title      : Nucleosynthesis in magnetorotational supernovae: impact of the magnetic field configuration
Authors    : M. Reichert, M. Bugli, J. Guilet, M. Obergaulinger, M. Á. Aloy, A. Arcones
Subjects   : High Energy Astrophysical Phenomena
Abstract    : The production of heavy elements is one of the main by-products of the explosive end of massive stars. A long sought goal is finding differentiated patterns in the nucleosynthesis yields, which could permit identifying a number of properties of the explosive core. Among them, the traces of the magnetic field topology are particularly important for extreme supernova explosions, most likely hosted by magnetorotational effects. We investigate the nucleosynthesis of five state-of-the-art magnetohydrodynamic models with fast rotation that have been previously calculated in full 3D and that involve an accurate neutrino transport (M1). One of the models does not contain any magnetic field and synthesizes elements around the iron group, in agreement with other CC-SNe models in literature. All other models host a strong magnetic field of the same intensity, but with different topology. For the first time, we investigate the nucleosynthesis of MR-SNe models with a quadrupolar magnetic field and a 90 degree tilted dipole. We obtain a large variety of ejecta compositions reaching from iron nuclei to nuclei up to the third r-process peak. We assess the robustness of our results by considering the impact of different nuclear physics uncertainties such as different nuclear masses,  $\beta^{-}$ -decays and  $\beta^{\{-\}}$ -delayed neutron emission probabilities, neutrino reactions, fission, and a feedback of nuclear energy on the temperature. We find that the qualitative results do not change with different nuclear physics input. The properties of the explosion dynamics and the magnetic field configuration are the dominant factors determining the ejecta composition.
Comments   :
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14402
```

Figure 4.B.1 Simple Boolean query with AND

RECOVERY OF INFORMATION

Αναζήτηση ακαδημαϊκών εργασιών

— □ ×

production or elements

Boolean Retrieval

Αναζήτηση

Επιλογή Φίλτρου

Αναζήτηση

```
===== Ερώτημα αναζήτησης : production or elements =====
===== Άλγορίθμος ανάκτησης : Boolean Retrieval =====
#1
-----
Document ID : 0
Title : Nucleosynthesis in magnetorotational supernovae: impact of the magnetic field configuration
Authors : M. Reichert, M. Bugli, J. Guilet, M. Obergaulinger, M. Á. Aloy, A. Arcones
Subjects : High Energy Astrophysical Phenomena
Abstract : The production of heavy elements is one of the main by-products of the explosive end of massive stars. A long sought goal is finding differentiated patterns in the nucleosynthesis yields, which could permit identifying a number of properties of the explosive core. Among them, the traces of the magnetic field topology are particularly important for \emph{extreme} supernova explosions, most likely hosted by magnetorotational effects. We investigate the nucleosynthesis of five state-of-the-art magnetohydrodynamic models with fast rotation that have been previously calculated in full 3D and that involve an accurate neutrino transport (M1). One of the models does not contain any magnetic field and synthesizes elements around the iron group, in agreement with other CC-SNe models in literature. All other models host a strong magnetic field of the same intensity, but with different topology. For the first time, we investigate the nucleosynthesis of MR-SNe models with a quadrupolar magnetic field and a 90 degree tilted dipole. We obtain a large variety of ejecta compositions reaching from iron nuclei to nuclei up to the third r-process peak. We assess the robustness of our results by considering the impact of different nuclear physics uncertainties such as different nuclear masses, $\beta^{-}$-decays and $\beta^{+}$-delayed neutron emission probabilities, neutrino reactions, fission, and a feedback of nuclear energy on the temperature. We find that the qualitative results do not change with different nuclear physics input. The properties of the explosion dynamics and the magnetic field configuration are the dominant factors determining the ejecta composition.
Comments :
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14402
#2
-----
Document ID : 42
Title : How far can we see back in time in high-energy collisions using charm quarks?
Authors : Laszlo Gyulai, Gabor Biro, Robert Vertesi, Gergely Gabor Barnafoldi
Subjects : High Energy Physics - Phenomenology, Nuclear Theory
Abstract : We use open charm production to estimate how far we can see back in time in high-energy hadron-hadron collisions. We analyze the transverse momentum distributions of the identified D mesons from pp, p-Pb and A-A collisions at the ALICE and STAR experiments covering the energy range from  $\sqrt{s_{\text{NN}}} = 200$  GeV up to 7 TeV. Within a non-extensive statistical framework, the common Tsallis parameters for D mesons represent higher temperature and more degrees of freedom than that of light-flavour hadrons. The production of D mesons corresponds to a significantly earlier proper time,  $\tau_{\text{D}} = (0.18 \pm 0.06) \tau_{\text{LF}}$ .
Comments : 18 pages, 6 figures, 1 table
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14282
#3
-----
Document ID : 45
Title : Phenomenology of TMD parton distributions in Drell-Yan and $Z^0$ boson production in a hadron structure oriented
```

IPython Console History

RECOVERY OF INFORMATION

taken from the 2014 Global Adult Population Survey or the Global Entrepreneurship Monitor project. The propensity for innovation amongst tourism entrepreneurs has a statistically significant relationship to gender, age, level of education and informal investments in previous businesses.

Comments : Journal ref: Sustainability 12:5003 (2020)
Date : 5 December 2023
PDF_URL : <https://arxiv.org/pdf/2401.13679>

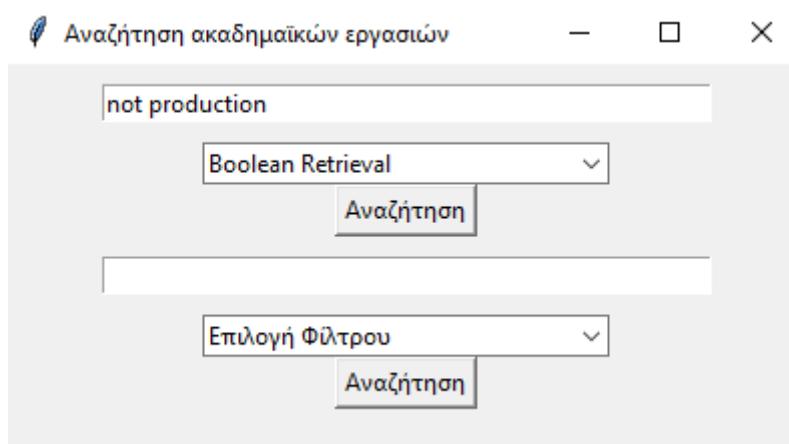
#19

Document ID : 236
Title : Realized Stochastic Volatility Model with Skew-t Distributions for Improved Volatility and Quantile Forecasting
Authors : Makoto Takahashi, Yuta Yamauchi, Toshiaki Watanabe, Yasuhiro Omori
Subjects : Econometrics
Abstract : Forecasting volatility and quantiles of financial returns is essential for accurately measuring financial tail risks, such as value-at-risk and expected shortfall. The critical elements in these forecasts involve understanding the distribution of financial returns and accurately estimating volatility. This paper introduces an advancement to the traditional stochastic volatility model, termed the realized stochastic volatility model, which integrates realized volatility as a precise estimator of volatility. To capture the well-known characteristics of return distribution, namely skewness and heavy tails, we incorporate three types of skew-t distributions. Among these, two distributions include the skew-normal feature, offering enhanced flexibility in modeling the return distribution. We employ a Bayesian estimation approach using the Markov chain Monte Carlo method and apply it to major stock indices. Our empirical analysis, utilizing data from US and Japanese stock indices, indicates that the inclusion of both skewness and heavy tails in daily returns significantly improves the accuracy of volatility and quantile forecasts.
Comments :
Date : 23 January 2024
PDF_URL : <https://arxiv.org/pdf/2401.13179>

#20

Document ID : 248
Title : Optimal design of a local renewable electricity supply system for power-intensive production processes with demand response
Authors : Sonja H. M. Germscheid, Benedikt Nilges, Niklas von der Assen, Alexander Mitsos, Manuel Dahmen
Subjects : Optimization and Control
Abstract : This work studies synergies arising from combining industrial demand response and local renewable electricity supply. To this end, we optimize the design of a local electricity generation and storage system with an integrated demand response scheduling of a continuous power-intensive production process in a multi-stage problem. We optimize both total annualized cost and global warming impact and consider local photovoltaic and wind electricity generation, an electric battery, and electricity trading on day-ahead and intraday market. We find that installing a battery can reduce emissions and enable large trading volumes on the electricity markets, but significantly increases cost. Economic and ecologic process and battery operation are driven primarily by the electricity price and grid emission factor, respectively, rather than locally generated electricity. A parameter study reveals that economic savings from the local system and flexibilizing the process behave almost additive.
Comments : manuscript (32 pages, 9 figures, 6 tables), supporting materials (11 pages, 9 figures, 2 tables)
Date : 23 January 2024
PDF_URL : <https://arxiv.org/pdf/2401.12759>

Figure 4.B.2 Simple Boolean query with OR



RECOVERY OF INFORMATION

```
=====
  Ερώτημα αναζήτησης : not production =====
  Αλγόριθμος ανάκτησης : Boolean Retrieval =====
-----

#1
-----
Document ID : 1
Title      : pix2gestalt: Amodal Segmentation by Synthesizing Wholes
Authors    : Ege Ozguroglu, Ruoshi Liu, Didac Suris, Dian Chen, Achal Dave, Pavel Tokmakov, Carl Vondrick
Subjects   : Computer Vision and Pattern Recognition, Machine Learning
Abstract   : We introduce pix2gestalt, a framework for zero-shot amodal segmentation, which learns to estimate the shape and appearance of whole objects that are only partially visible behind occlusions. By capitalizing on large-scale diffusion models and transferring their representations to this task, we learn a conditional diffusion model for reconstructing whole objects in challenging zero-shot cases, including examples that break natural and physical priors, such as art. As training data, we use a synthetically curated dataset containing occluded objects paired with their whole counterparts. Experiments show that our approach outperforms supervised baselines on established benchmarks. Our model can furthermore be used to significantly improve the performance of existing object recognition and 3D reconstruction methods in the presence of occlusions.
Comments   : Website: https://gestalt.cs.columbia.edu/
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14398.pdf
-----

#2
-----
Document ID : 2
Title      : Entanglement entropy and deconfined criticality: emergent SO(5) symmetry and proper lattice bipartition
Authors    : Jonathan D'Emidio, Anders W. Sandvik
Subjects   : Strongly Correlated Electrons, High Energy Physics - Lattice
Abstract   : We study the Rényi entanglement entropy (EE) of the two-dimensional  $J=Q$  model, the emblematic quantum spin model of deconfined criticality at the phase transition between antiferromagnetic and valence-bond-solid ground states. Quantum Monte Carlo simulations with an improved EE scheme reveal critical corner contributions that scale logarithmically with the system size, with a coefficient in remarkable agreement with the form expected from a large- $N$  conformal field theory with  $SO(N=5)$  symmetry. However, details of the bipartition of the lattice are crucial in order to observe this behavior. If the subsystem for the reduced density matrix does not properly accommodate valence-bond fluctuations, logarithmic contributions appear even for corner-less bipartitions. We here use a  $45^\circ$  tilted cut on the square lattice. Beyond supporting an  $SO(5)$  deconfined quantum critical point, our results for both the regular and tilted cuts demonstrate important microscopic aspects of the EE that are not captured by conformal field theory.
Comments   : 5 pages, 3 figures
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14396.pdf
-----

#3
-----
Document ID : 3
Title      : Summing up perturbation series around superintegrable point
Authors    : A. Mironov, A. Morozov, A. Popolitov, Sh. Shakirov
Subjects   : High Energy Physics - Theory, Mathematical Physics
Abstract   : We work out explicit formulas for correlators in the Gaussian matrix model perturbed by a logarithmic potential, i.e. by inserting Miwa variables. In this paper, we concentrate on the example of a single Miwa variable. The ordinary Gaussian
-----
```

#19

```
-----
Document ID : 19
Title      : Uncovering Heterogeneity of Solar Flare Mechanism With Mixture Models
Authors    : Bach Viet Do, Yang Chen, XuanLong Nguyen, Ward Manchester
Subjects   : Solar and Stellar Astrophysics, Applications, Methodology
Abstract   : The physics of solar flares occurring on the Sun is highly complex and far from fully understood. However, observations show that solar eruptions are associated with the intense kilogauss fields of active regions, where free energies are stored with field-aligned electric currents. With the advent of high-quality data sources such as the Geostationary Operational Environmental Satellites (GOES) and Solar Dynamics Observatory (SDO)/Helioseismic and Magnetic Imager (HMI), recent works on solar flare forecasting have been focusing on data-driven methods. In particular, black box machine learning and deep learning models are increasingly adopted in which underlying data structures are not modeled explicitly. If the active regions indeed follow the same laws of physics, there should be similar patterns shared among them, reflected by the observations. Yet, these black box models currently used in the literature do not explicitly characterize the heterogeneous nature of the solar flare data, within and between active regions. In this paper, we propose two finite mixture models designed to capture the heterogeneous patterns of active regions and their associated solar flare events. With extensive numerical studies, we demonstrate the usefulness of our proposed method for both resolving the sample imbalance issue and modeling the heterogeneity for rare energetic solar flare events.
Comments   :
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14345.pdf
-----

#20
-----
Document ID : 20
Title      : From the Choi Formalism in Infinite Dimensions to Unique Decompositions of Generators of Completely Positive Dynamical Semigroups
Authors    : Frederik vom Ende
Subjects   : Functional Analysis, Mathematical Physics, Quantum Physics
Abstract   : Given any separable complex Hilbert space, any trace-class operator  $B$  which does not have purely imaginary trace, and any generator  $L$  of a norm-continuous one-parameter semigroup of completely positive maps we prove that there exists a unique bounded operator  $K$  and a unique completely positive map  $S$  such that (i)  $L = K(\cdot) + (\cdot)K^*$ , (ii) the superoperator  $\Phi(B^*(\cdot)B)$  is trace class and has vanishing trace, and (iii)  $\|\mathrm{tr}(B^*K)\|$  is a real number. Central to our proof is a modified version of the Choi formalism which relates completely positive maps to positive semi-definite operators. We characterize when this correspondence is injective and surjective, respectively, which in turn explains why the proof idea of our main result cannot extend to non-separable Hilbert spaces. In particular, we find examples of positive semi-definite operators which have empty pre-image under the Choi formalism as soon as the underlying Hilbert space is infinite-dimensional.
Comments   : 25+3 pages. Generalizes arXiv:2310.04037 to infinite dimensions. To be submitted to J. Funct. Anal
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14344.pdf
```

Figure 4.B.3 Simple Boolean query with NOT

RECOVERY OF INFORMATION

```
===== Ερώτημα αναζήτησης : production or not elements or (framework or not entropy) =====
===== Άλγοριθμος ανάκτησης : Boolean Retrieval =====
-----
#1
-----
Document ID : 0
Title      : Nucleosynthesis in magnetorotational supernovae: impact of the magnetic field configuration
Authors    : M. Reichert, M. Bugli, J. Guilet, M. Obergaulinger, M. Á. Aloy, A. Arcones
Subjects   : High Energy Astrophysical Phenomena
Abstract    : The production of heavy elements is one of the main by-products of the explosive end of massive stars. A long sought goal is finding differentiated patterns in the nucleosynthesis yields, which could permit identifying a number of properties of the explosive core. Among them, the traces of the magnetic field topology are particularly important for \emph{extreme} supernova explosions, most likely hosted by magnetorotational effects. We investigate the nucleosynthesis of five state-of-the-art magnetohydrodynamic models with fast rotation that have been previously calculated in full 3D and that involve an accurate neutrino transport (M1). One of the models does not contain any magnetic field and synthesizes elements around the iron group, in agreement with other CC-SNe models in literature. All other models host a strong magnetic field of the same intensity, but with different topology. For the first time, we investigate the nucleosynthesis of MR-SNe models with a quadrupolar magnetic field and a 90 degree tilted dipole. We obtain a large variety of ejecta compositions reaching from iron nuclei to nuclei up to the third r-process peak. We assess the robustness of our results by considering the impact of different nuclear physics uncertainties such as different nuclear masses, $\beta^{\{-\}}$-decays and $\beta^{\{-\}}$-delayed neutron emission probabilities, neutrino reactions, fission, and a feedback of nuclear energy on the temperature. We find that the qualitative results do not change with different nuclear physics input. The properties of the explosion dynamics and the magnetic field configuration are the dominant factors determining the ejecta composition.
Comments   :
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14402
-----
#2
-----
Document ID : 1
Title      : pix2gestalt: Amodal Segmentation by Synthesizing Wholes
Authors    : Ege Ozguroglu, Ruoshi Liu, Dídac Surís, Dian Chen, Achal Dave, Pavel Tokmakov, Carl Vondrick
Subjects   : Computer Vision and Pattern Recognition, Machine Learning
Abstract    : We introduce pix2gestalt, a framework for zero-shot amodal segmentation, which learns to estimate the shape and appearance of whole objects that are only partially visible behind occlusions. By capitalizing on large-scale diffusion models and transferring their representations to this task, we learn a conditional diffusion model for reconstructing whole objects in challenging zero-shot cases, including examples that break natural and physical priors, such as art. As training data, we use a synthetically curated dataset containing occluded objects paired with their whole counterparts. Experiments show that our approach outperforms supervised baselines on established benchmarks. Our model can furthermore be used to significantly improve the performance of existing object recognition and 3D reconstruction methods in the presence of occlusions.
Comments   : Website: https://gestalt.cs.columbia.edu/
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14398
-----
#3
-----
Document ID : 2
Title      : Entanglement entropy and deconfined criticality: emergent SO(5) symmetry and proper lattice binartition
```

RECOVERY OF INFORMATION

```
model with unperturbed boundary magnetic fields based on the off-diagonal Bethe ansatz solution. At the initial setting, we obtain the different patterns of Bethe roots of the reduced Bethe ansatz equations for the different boundary parameters. According to them, we obtain the densities of states, ground state energy density and surface energy. Our results show that the system has the stable boundary bound states when the boundary magnetic fields satisfy some constraints.
Comments : 13 pages, 3 figures
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14356
-----
#19
-----
Document ID : 18
Title : Initial data for Minkowski stability with arbitrary decay
Authors : Allen Juntao Fang, Jérémie Szeftel, Arthur Touati
Subjects : Analysis of PDEs, General Relativity and Quantum Cosmology, Mathematical Physics
Abstract : We construct and parametrize solutions to the constraint equations of general relativity in a neighborhood of Minkowski spacetime with arbitrary prescribed decay properties at infinity. We thus provide a large class of initial data for the results on stability of Minkowski which include a mass term in the asymptotics. Due to the symmetries of Minkowski, a naive linear perturbation fails. Our construction is based on a simplified conformal method, a reduction to transverse traceless perturbations and a nonlinear fixed point argument where we face linear obstructions coming from the cokernels of both the linearized constraint operator and the Laplace operator. To tackle these obstructions, we introduce a well-chosen truncated black hole around which to perturb. The control of the parameters of the truncated black hole is the most technical part of the proof, since its center of mass and angular momentum could be arbitrarily large.
Comments : 86 pages
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14353
-----
#20
-----
Document ID : 19
Title : Uncovering Heterogeneity of Solar Flare Mechanism With Mixture Models
Authors : Bach Viet Do, Yang Chen, XuanLong Nguyen, Ward Manchester
Subjects : Solar and Stellar Astrophysics, Applications, Methodology
Abstract : The physics of solar flares occurring on the Sun is highly complex and far from fully understood. However, observations show that solar eruptions are associated with the intense kilogauss fields of active regions, where free energies are stored with field-aligned electric currents. With the advent of high-quality data sources such as the Geostationary Operational Environmental Satellites (GOES) and Solar Dynamics Observatory (SDO)/Helioseismic and Magnetic Imager (HMI), recent works on solar flare forecasting have been focusing on data-driven methods. In particular, black box machine learning and deep learning models are increasingly adopted in which underlying data structures are not modeled explicitly. If the active regions indeed follow the same laws of physics, there should be similar patterns shared among them, reflected by the observations. Yet, these black box models currently used in the literature do not explicitly characterize the heterogeneous nature of the solar flare data, within and between active regions. In this paper, we propose two finite mixture models designed to capture the heterogeneous patterns of active regions and their associated solar flare events. With extensive numerical studies, we demonstrate the usefulness of our proposed method for both resolving the sample imbalance issue and modeling the heterogeneity for rare energetic solar flare events.
Comments :
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14345
```

Figure 4.B.4 Complex Boolean query

Vector Space Model

The implementation of the Vector Space Model was carried out in the same module (`search_engine.py`) as the [user interface](#) and specifically in the same class. As input, the algorithm takes the user query and the implementations that take place are as follows:

- [Text Preprocessing: Step 2](#)
- [TF-IDF calculation of text terms: Ranking of results \(Ranking\)](#)
- [TF-IDF calculation of query terms: Ranking of results \(Ranking\)](#)
- [Calculation of the cosines between the texts and the query: Ranking of results \(Ranking\)](#)
- [Ranking of results \(Ranking\)](#)

Some illustrative case studies from user queries are as follows:

RECOVERY OF INFORMATION

Αναζήτηση ακαδημαϊκών εργασιών

the production of heavy elements

Vector Space Model

Αναζήτηση

Επιλογή Φίλτρου

Αναζήτηση

```
Console 1/A
=====
Eρώτησμα αναζήτησης : the production of heavy elements =====
Αλγόριθμος ανάκτησης : Vector Space Model =====
#1 Cosine Similarity: 0.1566
-----
Document ID : 0
Title      : Nucleosynthesis in magnetorotational supernovae: impact of the magnetic field configuration
Authors    : M. Reichert, M. Bugli, J. Guilet, M. Obergaulinger, M. Á. Aloy, A. Arcones
Subjects   : High Energy Astrophysical Phenomena
Abstract   : The production of heavy elements is one of the main by-products of the explosive end of massive stars. A long sought goal is finding differentiated patterns in the nucleosynthesis yields, which could permit identifying a number of properties of the explosive core. Among them, the traces of the magnetic field topology are particularly important for \text{extreme} supernova explosions, most likely hosted by magnetorotational effects. We investigate the nucleosynthesis of five state-of-the-art magnetohydrodynamic models with fast rotation that have been previously calculated in full 3D and that involve an accurate neutrino transport (M1). One of the models does not contain any magnetic field and synthesizes elements around the iron group, in agreement with other CC-SNe models in literature. All other models host a strong magnetic field of the same intensity, but with different topology. For the first time, we investigate the nucleosynthesis of MR-SNe models with a quadrupolar magnetic field and a 90 degree tilted dipole. We obtain a large variety of ejecta compositions reaching from iron nuclei to nuclei up to the third r-process peak. We assess the robustness of our results by considering the impact of different nuclear physics uncertainties such as different nuclear masses, $\beta^{-}$-decays and $\beta^{+}$-$\gamma$-delayed neutron emission probabilities, neutrino reactions, fission, and a feedback of nuclear energy on the temperature. We find that the qualitative results do not change with different nuclear physics input. The properties of the explosion dynamics and the magnetic field configuration are the dominant factors determining the ejecta composition.
Comments   :
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14402
#2 Cosine Similarity: 0.1539
-----
Document ID : 236
Title      : Realized Stochastic Volatility Model with Skew-t Distributions for Improved Volatility and Quantile Forecasting
Authors    : Makoto Takahashi, Yuta Yamauchi, Toshiaki Watanabe, Yasuhiro Omori
Subjects   : Econometrics
Abstract   : Forecasting volatility and quantiles of financial returns is essential for accurately measuring financial tail risks, such as value-at-risk and expected shortfall. The critical elements in these forecasts involve understanding the distribution of financial returns and accurately estimating volatility. This paper introduces an advancement to the traditional stochastic volatility model, termed the realized stochastic volatility model, which integrates realized volatility as a precise estimator of volatility. To capture the well-known characteristics of return distribution, namely skewness and heavy tails, we incorporate three types of skew-t distributions. Among these, two distributions include the skew-normal feature, offering enhanced flexibility in modeling the return distribution. We employ a Bayesian estimation approach using the Markov chain Monte Carlo method and apply it to major stock indices. Our empirical analysis, utilizing data from US and Japanese stock indices, indicates that the inclusion of both skewness and heavy tails in daily returns significantly improves the accuracy of volatility and quantile forecasts.
Comments   :
Date       : 23 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.13179
#3 Cosine Similarity: 0.0002
```

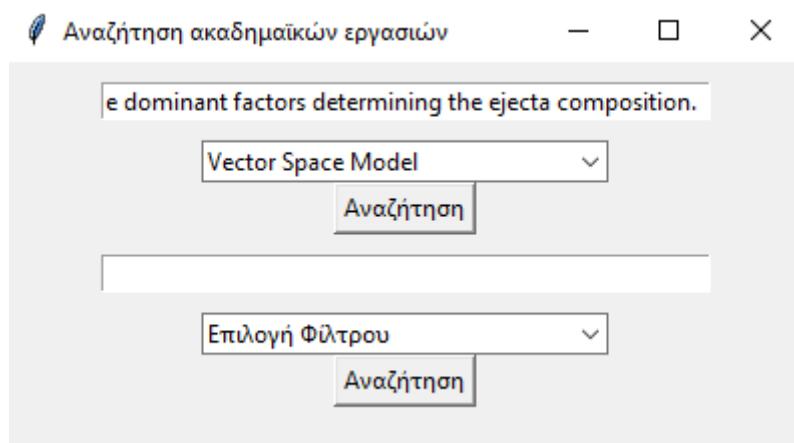
RECOVERY OF INFORMATION

```
Console 1/A X
stepwise statistical multi-regression model with leave-one-out cross-validation. Under specific management conditions (e.g., three annual cuts) and from one to five months in advance, the generated model successfully provided a p-value<0.01 in correlation (t-test), a root mean square error percentage (%RMSE) of 14.6% and a 71.43% hit rate predicting above/below average years in terms of forage yield collection.
Comments : Journal ref: Gomara I, Bellocchi G, Martin R, Rodriguez-Fonseca B, Ruiz-Ramos M (2020) Agricultural and Forest Meteorology, 280, 107768
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14053
-----
#19 Cosine Similarity: 0.0262

Document ID : 165
Title : FIMBA: Evaluating the Robustness of AI in Genomics via Feature Importance Adversarial Attacks
Authors : Heorhi Skovorodnikov, Hoda Alkhzaimi
Subjects : Machine Learning, Cryptography and Security, Genomics
Abstract : With the steady rise of the use of AI in bio-technical applications and the widespread adoption of genomics sequencing, an increasing amount of AI-based algorithms and tools is entering the research and production stage affecting critical decision-making streams like drug discovery and clinical outcomes. This paper demonstrates the vulnerability of AI models often utilized downstream tasks on recognized public genomics datasets. We undermine model robustness by deploying an attack that focuses on input transformation while mimicking the real data and confusing the model decision-making, ultimately yielding a pronounced deterioration in model performance. Further, we enhance our approach by generating poisoned data using a variational autoencoder-based model. Our empirical findings unequivocally demonstrate a decline in model performance, underscored by diminished accuracy and an upswing in false positives and false negatives. Furthermore, we analyze the resulting adversarial samples via spectral analysis yielding conclusions for countermeasures against such attacks.
Comments : 15 pages, core code available at: https://github.com/Heorhiis/fimba-attack
Date : 19 January 2024
PDF_URL : https://arxiv.org/pdf/2401.10657
-----
#20 Cosine Similarity: 0.0229

Document ID : 248
Title : Optimal design of a local renewable electricity supply system for power-intensive production processes with demand response
Authors : Sonja H. M. Germscheid, Benedikt Nilges, Niklas von der Assen, Alexander Mitsos, Manuel Dahmen
Subjects : Optimization and Control
Abstract : This work studies synergies arising from combining industrial demand response and local renewable electricity supply. To this end, we optimize the design of a local electricity generation and storage system with an integrated demand response scheduling of a continuous power-intensive production process in a multi-stage problem. We optimize both total annualized cost and global warming impact and consider local photovoltaic and wind electricity generation, an electric battery, and electricity trading on day-ahead and intraday market. We find that installing a battery can reduce emissions and enable large trading volumes on the electricity markets, but significantly increases cost. Economic and ecologic process and battery operation are driven primarily by the electricity price and grid emission factor, respectively, rather than locally generated electricity. A parameter study reveals that economic savings from the local system and flexibilizing the process behave almost additive.
Comments : manuscript (32 pages, 9 figures, 6 tables), supporting materials (11 pages, 9 figures, 2 tables)
Date : 23 January 2024
PDF_URL : https://arxiv.org/pdf/2401.12759
```

Figure 4.B.5 Small User Query (VSM)



RECOVERY OF INFORMATION

```
Console 1/A ×
=====
Ερώτηση ανάζήτησης : The production of heavy elements is one of the main by-products of the explosive end of massive stars. A long sought goal is finding differentiated patterns in the nucleosynthesis yields, which could permit identifying a number of properties of the explosive core. Among them, the traces of the magnetic field topology are particularly important for extreme supernova explosions, most likely hosted by magnetorotational effects. We investigate the nucleosynthesis of five state-of-the-art magnetohydrodynamic models with fast rotation that have been previously calculated in full 3D and that involve an accurate neutrino transport (M1). One of the models does not contain any magnetic field and synthesizes elements around the iron group, in agreement with other CC-SNe models in literature. All other models host a strong magnetic field of the same intensity, but with different topology. For the first time, we investigate the nucleosynthesis of MR-SNe models with a quadrupolar magnetic field and a 90 degree tilted dipole. We obtain a large variety of ejecta compositions reaching from iron nuclei to nuclei up to the third r-process peak. We assess the robustness of our results by considering the impact of different nuclear physics uncertainties such as different nuclear masses,  $\beta^{(-)}$ -decays and  $\beta^{(+)}$ -delayed neutron emission probabilities, neutrino reactions, fission, and a feedback of nuclear energy on the temperature. We find that the qualitative results do not change with different nuclear physics input. The properties of the explosion dynamics and the magnetic field configuration are the dominant factors determining the ejecta composition.
=====
Αλγόριθμος ανάκτησης : Vector Space Model =====
#1 Cosine Similarity: 1.0000
-----
Document ID : 0
Title : Nucleosynthesis in magnetorotational supernovae: impact of the magnetic field configuration
Authors : M. Reichert, M. Bugli, J. Guilet, M. Obergaulinger, M. Á. Aloy, A. Arcones
Subjects : High Energy Astrophysical Phenomena
Abstract : The production of heavy elements is one of the main by-products of the explosive end of massive stars. A long sought goal is finding differentiated patterns in the nucleosynthesis yields, which could permit identifying a number of properties of the explosive core. Among them, the traces of the magnetic field topology are particularly important for extreme supernova explosions, most likely hosted by magnetorotational effects. We investigate the nucleosynthesis of five state-of-the-art magnetohydrodynamic models with fast rotation that have been previously calculated in full 3D and that involve an accurate neutrino transport (M1). One of the models does not contain any magnetic field and synthesizes elements around the iron group, in agreement with other CC-SNe models in literature. All other models host a strong magnetic field of the same intensity, but with different topology. For the first time, we investigate the nucleosynthesis of MR-SNe models with a quadrupolar magnetic field and a 90 degree tilted dipole. We obtain a large variety of ejecta compositions reaching from iron nuclei to nuclei up to the third r-process peak. We assess the robustness of our results by considering the impact of different nuclear physics uncertainties such as different nuclear masses,  $\beta^{(-)}$ -decays and  $\beta^{(+)}$ -delayed neutron emission probabilities, neutrino reactions, fission, and a feedback of nuclear energy on the temperature. We find that the qualitative results do not change with different nuclear physics input. The properties of the explosion dynamics and the magnetic field configuration are the dominant factors determining the ejecta composition.
Comments :
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14402
-----
#2 Cosine Similarity: 0.1291
-----
Document ID : 70
Title : Magnetic fields of protoplanetary disks
Authors : Sergey A. Khabibrakhmanov
Subjects : Solar and Stellar Astrophysics, Earth and Planetary Astrophysics, Plasma Physics
Abstract : We review the current status of studies on accretion and protoplanetary disks of young stars with large-scale
-----
Console 1/A ×
=====
electric and magnetic charges has been proved. Using conformal positive energy theorem, as well as, the positive mass theorem and adequate conformal transformations, we envisage the two alternative ways of proving that the exterior region of a certain radius of the studied static (i.e photon sphere), is characterized by ADM mass, electric and magnetic charges.
Comments : 22 pages, RevTex, to be published in Phys.Rev.D15
Date : 25 January 2024
PDF_URL : https://arxiv.org/pdf/2401.14116
-----
#19 Cosine Similarity: 0.0453
-----
Document ID : 133
Title : A distribution-guided Mapper algorithm
Authors : Yuyang Tao, Shufei Ge
Subjects : Algebraic Topology, Machine Learning, Quantitative Methods
Abstract : Motivation: The Mapper algorithm is an essential tool to explore shape of data in topology data analysis. With a dataset as an input, the Mapper algorithm outputs a graph representing the topological features of the whole dataset. This graph is often regarded as an approximation of a reeb graph of data. The classic Mapper algorithm uses fixed interval lengths and overlapping ratios, which might fail to reveal subtle features of data, especially when the underlying structure is complex.
Results: In this work, we introduce a distribution guided Mapper algorithm named D-Mapper, that utilizes the property of the probability model and data intrinsic characteristics to generate density guided covers and provides enhanced topological features. Our proposed algorithm is a probabilistic model-based approach, which could serve as an alternative to non-probabilistic ones. Moreover, we introduce a metric accounting for both the quality of overlap clustering and extended persistence homology to measure the performance of Mapper type algorithm. Our numerical experiments indicate that the D-Mapper outperforms the classical Mapper algorithm in various scenarios. We also apply the D-Mapper to a SARS-CoV-2 coronavirus RNA sequences dataset to explore the topological structure of different virus variants. The results indicate that the D-Mapper algorithm can reveal both vertical and horizontal evolution processes of the viruses.
Availability: Our package is available at https://github.com/ShufeiGe/D-Mapper.
Comments :
Date : 19 January 2024
PDF_URL : https://arxiv.org/pdf/2401.12237
-----
#20 Cosine Similarity: 0.0429
-----
Document ID : 138
Title : Approximating a linear dynamical system from non-sequential data
Authors : Cliff Stein, Pratik Worah
Subjects : Genomics
Abstract : Given non-sequential snapshots from instances of a dynamical system, we design a compressed sensing based algorithm that reconstructs the dynamical system. We formally prove that successful reconstruction is possible under the assumption that we can construct an approximate clock from a subset of the coordinates of the underlying system.
As an application, we argue that our assumption is likely true for genomic datasets, and we recover the underlying nuclear receptor networks and predict pathways, as opposed to genes, that may differentiate phenotypes in some publicly available datasets.
Comments :
Date : 22 January 2024
PDF_URL : https://arxiv.org/pdf/2401.11858
```

Figure 4.B.6 Identical user query with some text (VSM)

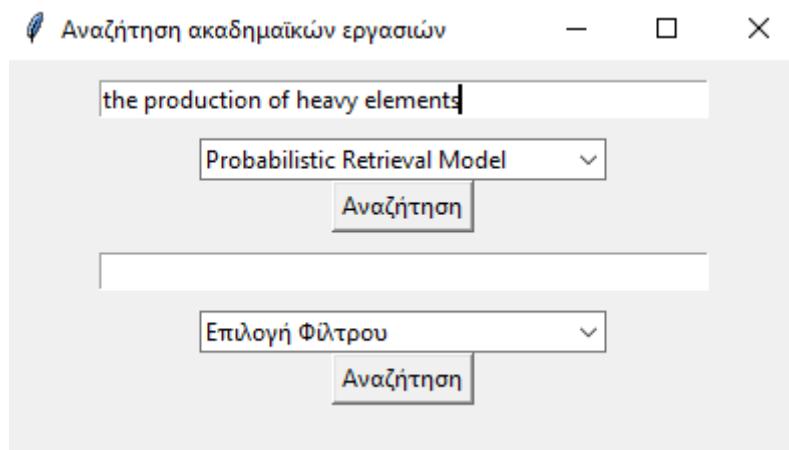
RECOVERY OF INFORMATION

Probabilistic Retrieval Model

The implementation of the Probabilistic Retrieval Model was carried out in the same module (search_engine.py) as the [user interface](#) and specifically in the same class. As input, the algorithm takes the user query and the inverted index, and the implementations that take place are as follows:

- [Text Preprocessing: Step 2](#)
- [TF-IDF calculation of text terms: Ranking of results \(Ranking\)](#)
- [TF-IDF calculation of query terms: Ranking of results \(Ranking\)](#)
- [Calculation of the Okapi BM25 score between the texts and the query: Ranking of results \(Ranking\)](#)
- [Ranking of results \(Ranking\)](#)

Some illustrative case studies from user queries are as follows:



RECOVERY OF INFORMATION

```

Console 1/A ×

===== Ερώτημα αναζήτησης : the production of heavy elements =====
===== Αλγόριθμος ανάκτησης : Probabilistic Retrieval Model =====

#1 BM25 Score: 41.4798
-----
Document ID : 211
Title      : Optimal Queueing Regimes
Authors    : Marco Scarsini, Eran Shmaya
Subjects   : Theoretical Economics, Computer Science and Game Theory, Probability
Abstract   : We consider an M/M/1 queueing model where customers can strategically decide whether to join the queue or balk and when to renege. We characterize the class of queueing regimes such that, for any parameters of the model, the socially efficient behavior is an equilibrium outcome.
Comments   : MSC Class: 91A40; 60J28
Date       : 24 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.13812

#2 BM25 Score: 41.0244
-----
Document ID : 194
Title      : The Interplay Between Logical Phenomena and the Cognitive System of the Mind
Authors    : Kazem Haghnejad Azar
Subjects   : Neurons and Cognition
Abstract   : In this article, we employ mathematical concepts as a tool to examine the phenomenon of consciousness experience and logical phenomena. Through our investigation, we aim to demonstrate that our experiences, while not confined to limitations, cannot be neatly encapsulated within a singular collection. Our conscious experience emerges as a result of the developmental and augmentative trajectory of our cognitive system. As our cognitive abilities undergo refinement and advancement, our capacity for logical thinking likewise evolves, thereby manifesting a heightened level of conscious experience. The primary objective of this article is to embark upon a profound exploration of the concept of logical experience, delving into the intricate process by which these experiences are derived from our mind.
Comments   :
Date       : 21 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.09465

#3 BM25 Score: 40.8304
-----
Document ID : 256
Title      : The outcomes of generative AI are exactly the Nash equilibria of a non-potential game
Authors    : Boualem Djehiche, Hamidou Tembine
Subjects   : Computer Science and Game Theory
Abstract   : In this article we show that the asymptotic outcomes of both shallow and deep neural networks such as those used in BloombergGPT to generate economic time series are exactly the Nash equilibria of a non-potential game. We then design and analyze deep neural network algorithms that converge to these equilibria. The methodology is extended to federated deep neural networks between clusters of regional servers and on-device clients. Finally, the variational inequalities behind large language models including encoder-decoder related transformers are established.
Comments   : 24 pages. Accepted and to appear in: International Econometric Conference of Vietnam
Date       : 22 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.12321

Console 1/A ×
PDF_URL   : https://arxiv.org/pdf/2401.14009v1.pdf
#19 BM25 Score: 39.5568
-----
Document ID : 189
Title      : Dimensional Neuroimaging Endophenotypes: Neurobiological Representations of Disease Heterogeneity Through Machine Learning
Authors    : Junhao Wen, Mathilde Antoniades, Zhijian Yang, Gyujoon Hwang, Ioanna Skampardonis, Rongguang Wang, Christos Davatzikos
Subjects   : Machine Learning, Image and Video Processing, Quantitative Methods
Abstract   : Machine learning has been increasingly used to obtain individualized neuroimaging signatures for disease diagnosis, prognosis, and response to treatment in neuropsychiatric and neurodegenerative disorders. Therefore, it has contributed to a better understanding of disease heterogeneity by identifying disease subtypes that present significant differences in various brain phenotypic measures. In this review, we first present a systematic literature overview of studies using machine learning and multimodal MRI to unravel disease heterogeneity in various neuropsychiatric and neurodegenerative disorders, including Alzheimer disease, schizophrenia, major depressive disorder, autism spectrum disorder, multiple sclerosis, as well as their potential in transdiagnostic settings. Subsequently, we summarize relevant machine learning methodologies and discuss an emerging paradigm which we call dimensional neuroimaging endophenotype (DNE). DNE dissects the neurobiological heterogeneity of neuropsychiatric and neurodegenerative disorders into a low dimensional yet informative, quantitative brain phenotypic representation, serving as a robust intermediate phenotype (i.e., endophenotype) largely reflecting underlying genetics and etiology. Finally, we discuss the potential clinical implications of the current findings and envision future research avenues.
Comments   :
Date       : 17 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.09517

#20 BM25 Score: 39.5309
-----
Document ID : 19
Title      : Uncovering Heterogeneity of Solar Flare Mechanism With Mixture Models
Authors    : Bach Viet Do, Yang Chen, XuanLong Nguyen, Ward Manchester
Subjects   : Solar and Stellar Astrophysics, Applications, Methodology
Abstract   : The physics of solar flares occurring on the Sun is highly complex and far from fully understood. However, observations show that solar eruptions are associated with the intense kilogauss fields of active regions, where free energies are stored with field-aligned electric currents. With the advent of high-quality data sources such as the Geostationary Operational Environmental Satellites (GOES) and Solar Dynamics Observatory (SDO)/Helioseismic and Magnetic Imager (HMI), recent works on solar flare forecasting have been focusing on data-driven methods. In particular, black box machine learning and deep learning models are increasingly adopted in which underlying data structures are not modeled explicitly. If the active regions indeed follow the same laws of physics, there should be similar patterns shared among them, reflected by the observations. Yet, these black box models currently used in the literature do not explicitly characterize the heterogeneous nature of the solar flare data, within and between active regions. In this paper, we propose two finite mixture models designed to capture the heterogeneous patterns of active regions and their associated solar flare events. With extensive numerical studies, we demonstrate the usefulness of our proposed method for both resolving the sample imbalance issue and modeling the heterogeneity for rare energetic solar flare events.
Comments   :
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14345

```

Figure 4.B.7 Small User Query (PRM)

RECOVERY OF INFORMATION

Αναζήτηση ακαδημαϊκών εργασιών

s by which these experiences are derived from our mind.

Probabilistic Retrieval Model

Αναζήτηση

Επιλογή Φίλτρου

Αναζήτηση

Console 1/A

```
===== Ερώτημα αναζήτησης : In this article, we employ mathematical concepts as a tool to examine the phenomenon of consciousness experience and logical phenomena. Through our investigation, we aim to demonstrate that our experiences, while not confined to limitations, cannot be neatly encapsulated within a singular collection. Our conscious experience emerges as a result of the developmental and augmentative trajectory of our cognitive system. As our cognitive abilities undergo refinement and advancement, our capacity for logical thinking likewise evolves, thereby manifesting a heightened level of conscious experience. The primary objective of this article is to embark upon a profound exploration of the concept of logical experience, delving into the intricate process by which these experiences are derived from our mind. =====
===== Αλγόριθμος ανάκτησης : Probabilistic Retrieval Model =====
#1 BM25 Score: 719.5806
Document ID : 194
Title      : The Interplay Between Logical Phenomena and the Cognitive System of the Mind
Authors    : Kazem Haghnejad Azar
Subjects   : Neurons and Cognition
Abstract   : In this article, we employ mathematical concepts as a tool to examine the phenomenon of consciousness experience and logical phenomena. Through our investigation, we aim to demonstrate that our experiences, while not confined to limitations, cannot be neatly encapsulated within a singular collection. Our conscious experience emerges as a result of the developmental and augmentative trajectory of our cognitive system. As our cognitive abilities undergo refinement and advancement, our capacity for logical thinking likewise evolves, thereby manifesting a heightened level of conscious experience. The primary objective of this article is to embark upon a profound exploration of the concept of logical experience, delving into the intricate process by which these experiences are derived from our mind.
Comments   :
Date       : 21 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.09465
#2 BM25 Score: 711.1157
Document ID : 144
Title      : Accelerating Seed Location Filtering in DNA Read Mapping Using a Commercial Compute-in-SRAM Architecture
Authors    : Courtney Golden, Dan Ilan, Nicholas Cebry, Christopher Batten
Subjects   : Hardware Architecture, Genomics
Abstract   : DNA sequence alignment is an important workload in computational genomics. Reference-guided DNA assembly involves aligning many read sequences against candidate locations in a long reference genome. To reduce the computational load of this alignment, candidate locations can be pre-filtered using simpler alignment algorithms like edit distance. Prior work has explored accelerating filtering on simulated compute-in-DRAM, due to the massive parallelism of compute-in-memory architectures. In this paper, we present work-in-progress on accelerating filtering using a commercial compute-in-SRAM accelerator. We leverage the recently released Gemini accelerator platform from GSI Technology, which is the first, to our knowledge, commercial-scale compute-in-SRAM system. We accelerate the Myers' bit-parallel edit distance algorithm, producing average speedups of 14.1x over single-core CPU performance. Individual query/candidate alignments produce speedups of up to 24.1x. These early results suggest this novel architecture is well-suited to accelerating the filtering step of sequence-to-sequence DNA alignment.
Comments   : Journal ref: 5th Workshop on Accelerator Architecture in Computational Biology and Bioinformatics (AACBB), June 2023
Date       : 21 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.11685
```

RECOVERY OF INFORMATION

```
Console 1/A ×
#18 BM25 Score: 687.6884
-----
Document ID : 10
Title      : Efficient Optimisation of Physical Reservoir Computers using only a Delayed Input
Authors     : Enrico Picco, Lina Jaurigue, Kathy Lüdge, Serge Massar
Subjects    : Emerging Technologies, Artificial Intelligence, Neural and Evolutionary Computing, Optics
Abstract    : We present an experimental validation of a recently proposed optimization technique for reservoir computing, using an optoelectronic setup. Reservoir computing is a robust framework for signal processing applications, and the development of efficient optimization approaches remains a key challenge. The technique we address leverages solely a delayed version of the input signal to identify the optimal operational region of the reservoir, simplifying the traditionally time-consuming task of hyperparameter tuning. We verify the effectiveness of this approach on different benchmark tasks and reservoir operating conditions.
Comments   :
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14371
-----
#19 BM25 Score: 687.0624
-----
Document ID : 251
Title      : Moen Meets Rotemberg: An Earthly Model of the Divine Coincidence
Authors     : Pascal Michaillat, Emmanuel Saez
Subjects    : Theoretical Economics
Abstract    : This paper proposes a model of the divine coincidence, explaining its recent appearance in US data. The divine coincidence matters because it helps explain the behavior of inflation after the pandemic, and it guarantees that the full-employment and price-stability mandates of the Federal Reserve coincide. In the model, a Phillips curve relating unemployment to inflation arises from Moen's (1997) directed search. The Phillips curve is nonvertical thanks to Rotemberg's (1982) price-adjustment costs. The model's Phillips curve guarantees that the rate of inflation is on target whenever the rate of unemployment is efficient, generating the divine coincidence. If we assume that wage decreases -- which reduce workers' morale -- are more costly to producers than price increases -- which upset customers -- the Phillips curve also displays a kink at the point of divine coincidence.
Comments   :
Date       : 22 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.12475
-----
#20 BM25 Score: 685.8325
-----
Document ID : 211
Title      : Optimal Queueing Regimes
Authors     : Marco Scarsini, Eran Shmaya
Subjects    : Theoretical Economics, Computer Science and Game Theory, Probability
Abstract    : We consider an M/M/1 queueing model where customers can strategically decide whether to join the queue or balk and when to renege. We characterize the class of queueing regimes such that, for any parameters of the model, the socially efficient behavior is an equilibrium outcome.
Comments   : MSC Class: 91A40; 60J28
Date       : 24 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.13812
```

Figure 4.B.8 Identical user query with some text (PRM)

4.b.3 Evaluation

Boolean Retrieval

The Boolean Retrieval code is implemented in such a way that it can search with both simple and complex expressions that will be given to it by the user through a search query. However, the specific keywords necessary for the operation of the algorithm (and, or, not) must be spelled correctly, because whether the letters entered by the user are in lowercase or uppercase does not affect in search of something. In the event that a word given through a query by the user does not exist in the index, the code will handle the error as in other search engines if a word is not in the indexes, there are no relevant results.

Vector Space Model

RECOVERY OF INFORMATION

The Vector Space Model code performs for the small search engine repository and is an important tool for retrieving tasks. The cosine calculation is an advanced ranking technique, so the retrieval results could be limited to the simple TF-IDF ranking algorithm. It is worth noting that the Python library, scikit-learn, which provides ready-made methods that generate the vectors and calculate the cosines, was not used.

Probabilistic Retrieval Model

The Probabilistic Retrieval Model code also performs for large search engine repositories and is an important tool for retrieving tasks. The calculation of the bm25 score is an advanced ranking technique, where it offers efficiency in the recovery of tasks, removing with the constant saturation coefficients b and k, the terms that appear very often in the collection and are usually of little importance information.

4.c. Filter search results by various criteria

4.c.1 Planning

Filtering the results is an option that the user has after entering the query he wants to search for and the algorithm he chooses to do the search. Essentially it offers the user the choice of the results he got to filter them either based on their author, or based on the date they were published. This is achieved through a second query in which the user must enter an author's name (in the case that the "Authors" field has been selected as filtering) so that only those in which he/she is included in the authors are displayed from the results obtained who asks Accordingly, he must do the same if he wants to filter by date, but the field must be filtered by "Date".

4.c.2 Implementation

To implement the filtering, we needed access to all the data we have in json format as well as the algorithm used for the search. Knowing the algorithm that was used, the filtering function knows from which table it should get the results (boolean_results if boolean retrieval was used etc). These results are the doc_ids of the texts printed by each algorithm. The function that deals with filtering takes the filter selected by the user (date, author) and also a query that contains what the user wants to search for in the results. Later, the necessary check is done with the author or date fields and the texts that have the same information as the user's query are printed.

Some illustrative case studies from user queries are as follows:

RECOVERY OF INFORMATION

The screenshot shows a search interface with the following components:

- A top bar with the text "Αναζήτηση ακαδημαϊκών εργασιών" and three icons: a minus sign, a square, and a cross.
- A search input field containing "the product of heavy elements".
- A dropdown menu set to "Vector Space Model".
- A search button labeled "Αναζήτηση".
- An intermediate search result showing "Jonathan D'Emidio".
- A second dropdown menu set to "Συγγραφείς".
- A second search button labeled "Αναζήτηση".
- A bottom panel titled "Console 1/A" displaying the following text:

```
Φιλτράρισμα αποτελεσμάτων κατά: Συγγραφείς
-----
#1 Cosine Similarity: 0.0000
-----
Document ID : 2
Title      : Entanglement entropy and deconfined criticality: emergent SO(5) symmetry and proper lattice bipartition
Authors    : Jonathan D'Emidio, Anders W. Sandvik
Subjects   : Strongly Correlated Electrons, High Energy Physics - Lattice
Abstract   : We study the Rényi entanglement entropy (EE) of the two-dimensional  $J\text{-}Q$  model, the emblematic quantum spin model of deconfined criticality at the phase transition between antiferromagnetic and valence-bond-solid ground states. Quantum Monte Carlo simulations with an improved EE scheme reveal critical corner contributions that scale logarithmically with the system size, with a coefficient in remarkable agreement with the form expected from a large- $N$  conformal field theory with  $SO(N=5)$  symmetry. However, details of the bipartition of the lattice are crucial in order to observe this behavior. If the subsystem for the reduced density matrix does not properly accommodate valence-bond fluctuations, logarithmic contributions appear even for corner-less bipartitions. We here use a  $45^\circ$  tilted cut on the square lattice. Beyond supporting an  $SO(5)$  deconfined quantum critical point, our results for both the regular and tilted cuts demonstrate important microscopic aspects of the EE that are not captured by conformal field theory.
Comments   : 5 pages, 3 figures
Date       : 25 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.14396
```

Figure 4.C.1 Filtering by Author

The screenshot shows a search interface with the following components:

- A top bar with the text "Αναζήτηση ακαδημαϊκών εργασιών" and three icons: a minus sign, a square, and a cross.
- A search input field containing "the product of heavy elements".
- A dropdown menu set to "Probabilistic Retrieval Model".
- A search button labeled "Αναζήτηση".
- An intermediate search result showing "17 January 2024".
- A second dropdown menu set to "Ημερομηνία".
- A second search button labeled "Αναζήτηση".

RECOVERY OF INFORMATION

```
Console 1/A ×
Φιλτράρισμα αποτελεσμάτων κατά: Ημερομηνία
#1 BM25 Score: 39.8124
-----
Document ID : 190
Title      : Is the Emergence of Life an Expected Phase Transition in the Evolving Universe?
Authors    : Stuart Kauffman, Andrea Roli
Subjects   : Populations and Evolution, Biological Physics
Abstract   : We propose a novel definition of life in terms of which its emergence in the universe is expected, and its ever-creative open-ended evolution is entailed by no law. Living organisms are Kantian Wholes that achieve Catalytic Closure, Constraint Closure, and Spatial Closure. We here unite for the first time two established mathematical theories, namely Collectively Autocatalytic Sets and the Theory of the Adjacent Possible. The former establishes that a first-order phase transition to molecular reproduction is expected in the chemical evolution of the universe where the diversity and complexity of molecules increases; the latter posits that, under loose hypotheses, if the system starts with a small number of beginning molecules, each of which can combine with copies of itself or other molecules to make new molecules, over time the number of kinds of molecules increases slowly but then explodes upward hyperbolically. Together these theories imply that life is expected as a phase transition in the evolving universe. The familiar distinction between software and hardware loses its meaning in living cells. We propose new ways to study the phylogeny of metabolisms, new astronomical ways to search for life on exoplanets, new experiments to seek the emergence of the most rudimentary life, and the hint of a coherent testable pathway to prokaryotes with template replication and coding.
Comments   :
Date       : 17 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.09514
-----
#2 BM25 Score: 39.5568
-----
Document ID : 189
Title      : Dimensional Neuroimaging Endophenotypes: Neurobiological Representations of Disease Heterogeneity Through Machine Learning
Authors    : Junhao Wen, Mathilde Antoniades, Zhijian Yang, Gyujoon Hwang, Ioanna Skampardonis, Rongguang Wang, Christos Davatzikos
Subjects   : Machine Learning, Image and Video Processing, Quantitative Methods
Abstract   : Machine learning has been increasingly used to obtain individualized neuroimaging signatures for disease diagnosis, prognosis, and response to treatment in neuropsychiatric and neurodegenerative disorders. Therefore, it has contributed to a better understanding of disease heterogeneity by identifying disease subtypes that present significant differences in various brain phenotypic measures. In this review, we first present a systematic literature overview of studies using machine learning and multimodal MRI to unravel disease heterogeneity in various neuropsychiatric and neurodegenerative disorders, including Alzheimer disease, schizophrenia, major depressive disorder, autism spectrum disorder, multiple sclerosis, as well as their potential in transdiagnostic settings. Subsequently, we summarize relevant machine learning methodologies and discuss an emerging paradigm which we call dimensional neuroimaging endophenotype (DNE). DNE dissects the neurobiological heterogeneity of neuropsychiatric and neurodegenerative disorders into a low dimensional yet informative, quantitative brain phenotypic representation, serving as a robust intermediate phenotype (i.e., endophenotype) largely reflecting underlying genetics and etiology. Finally, we discuss the potential clinical implications of the current findings and envision future research avenues.
Comments   :
Date       : 17 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.09517
-----
Console 1/A ×
PDF_URL   : https://arxiv.org/pdf/2401.09255
-----
#16 BM25 Score: 35.4404
-----
Document ID : 295
Title      : Early Prediction of Geomagnetic Storms by Machine Learning Algorithms
Authors    : Iris Yan
Subjects   : Machine Learning
Abstract   : Geomagnetic storms (GS) occur when solar winds disrupt Earth's magnetosphere. GS can cause severe damages to satellites, power grids, and communication infrastructures. Estimate of direct economic impacts of a large scale GS exceeds $40 billion a day in the US. Early prediction is critical in preventing and minimizing the hazards. However, current methods either predict several hours ahead but fail to identify all types of GS, or make predictions within short time, e.g., one hour ahead of the occurrence. This work aims to predict all types of geomagnetic storms reliably and as early as possible using big data and machine learning algorithms. By fusing big data collected from multiple ground stations in the world on different aspects of solar measurements and using Random Forests regression with feature selection and downsampling on minor geomagnetic storm instances (which carry majority of the data), we are able to achieve an accuracy of 82.55% on data collected in 2021 when making early predictions three hours in advance. Given that important predictive features such as historic Kp indices are measured every 3 hours and their importance decay quickly with the amount of time in advance, an early prediction of 3 hours ahead of time is believed to be close to the practical limit.
Comments   : 14 pages, 7 figures
Date       : 17 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.10290
-----
#17 BM25 Score: 34.9005
-----
Document ID : 187
Title      : Functional Linear Non-Gaussian Acyclic Model for Causal Discovery
Authors    : Tian-Le Yang, Kuang-Yao Lee, Kun Zhang, Joe Suzuki
Subjects   : Machine Learning, Statistics Theory, Neurons and Cognition, Methodology
Abstract   : In causal discovery, non-Gaussianity has been used to characterize the complete configuration of a Linear Non-Gaussian Acyclic Model (LiNGAM), encompassing both the causal ordering of variables and their respective connection strengths. However, LiNGAM can only deal with the finite-dimensional case. To expand this concept, we extend the notion of variables to encompass vectors and even functions, leading to the Functional Linear Non-Gaussian Acyclic Model (Func-LiNGAM). Our motivation stems from the desire to identify causal relationships in brain-effective connectivity tasks involving, for example, fMRI and EEG datasets. We demonstrate why the original LiNGAM fails to handle these inherently infinite-dimensional datasets and explain the availability of functional data analysis from both empirical and theoretical perspectives. {We establish theoretical guarantees of the identifiability of the causal relationship among non-Gaussian random vectors and even random functions in infinite-dimensional Hilbert spaces.} To address the issue of sparsity in discrete time points within intrinsic infinite-dimensional functional data, we propose optimizing the coordinates of the vectors using functional principal component analysis. Experimental results on synthetic data verify the ability of the proposed framework to identify causal relationships among multivariate functions using the observed samples. For real data, we focus on analyzing the brain connectivity patterns derived from fMRI data.
Comments   :
Date       : 17 January 2024
PDF_URL   : https://arxiv.org/pdf/2401.09641
```

Figure 4.C.2 Filtering by Publication Date

RECOVERY OF INFORMATION

4.c.3 Evaluation

The function to filter the results works well, however it is a rather simplistic function that requires the user to precisely enter the items they are interested in looking for in the algorithm results. The only case where the function does not work is when the user enters more than 1 author as the input to the filter. However, for an author the function satisfies the filter logic.

Query Processing

Planning

The process of query processing (Query Processing) followed by the [Boolean Retrieval algorithm](#) is based on Boolean algebra and set theory, using three basic Boolean operations: OR, AND and NOT. These functions affect the search results according to their position in the query received by the user. The relevant documents are retrieved using the inverted index. The theoretical/mathematical background of the design model is as follows:

AND

$$AND_RES = Q_0 \cap Q_1 = \{doc_{id} \mid doc_{id} \in Q_0 \text{ and } doc_{id} \in Q_1\}$$

AND_RES □ The texts containing both terms Q_0 , Q_1

Q_0 □ The texts containing the term Q_0

Q_1 □ The texts containing the term Q_1

OR

$$OR_RES = Q_0 \cup Q_1 = \{doc_{id} \mid doc_{id} \in Q_0 \text{ or } doc_{id} \in Q_1\}$$

RES □ The texts containing either the term Q_0 , or the term Q_1 or both

Q_0 □ The texts containing the term Q_0

Q_1 □ The texts containing the term Q_1

NOT

RECOVERY OF INFORMATION

$$NOT_RES = \neg Q_0 = \{doc_{id} \mid doc_{id} \in D \setminus Q_0\}$$

$RES \sqsubseteq$ Texts that do not contain the term Q_0

$D \sqsubseteq$ The texts of the collection

$Q_0 \sqsubseteq$ The texts containing the term Q_0

The process is carried out for simple Boolean user queries e.g. black OR white, NOT black, black AND WHITE. For complex queries, the [Boolean Retrieval](#) algorithm determines the priority order of operations, where it breaks the complex query into simple sub-queries and sends them in order of priority for processing.

Implementation

The query processing implementation was performed in a separate module (query_processing.py). Initially, the terms of the Boolean query are separated into verbal units and using the inverted index from a routine, all the terms of the Boolean query are replaced with the posting lists with the texts (doc_id) in which they appear. Boolean operations are performed on these lists. If the term does not exist in the index, then the routine handles the error.

Retrieving the relevant documents for each logical operator is implemented as follows:

AND

If the current term is the AND operator, the repository doc_ids that belong to the list of doc_ids displayed by the previous term and the corresponding list of the next term are retrieved.

OR

If the current term is the OR operator, the repository doc_ids that belong to the list of doc_ids displayed by the previous term and the doc_ids that belong to the list of the next term are retrieved.

RECOVERY OF INFORMATION

NOT

If the current term is the 'NOT' operator, the repository doc_ids that are not in the list of doc_ids that the next term appears in are retrieved

The implementation focuses on handling the Boolean functional operators AND, OR, NOT for query processing, using sets for the set operations.

Evaluation

[Boolean Retrieval](#) algorithm . The philosophy "belongs to x text, does not belong to y text" provides an efficient way to retrieve texts based on the query posed by the user. The [design](#) presents the theoretical background based on set theory, while the [implementation](#) presents the practical background based on Boolean algebra. The replacement of the term-words with the lists of the texts of the collection they belong to, according to the inverted index, speeds up the process of retrieving the texts that the user poses as a query. Representing the process with set theory and Boolean algebra make text retrieval efficient.

Ranking of results (Ranking)

Planning

The design model of ranking results (Ranking) firstly follows the application of the simple TF-IDF (Term Frequency-Inverse Document Frequency) ranking algorithm and then the application of advanced Cosine Similarity ranking techniques followed by the [Vector Space Model](#) and Okapi BM25 recovery algorithm that the [Probabilistic Retrieval Model](#) retrieval algorithm follows . The aforementioned applications contribute to the estimation of the semantic importance of the terms in a set of texts, where, based on a user query, the cosine similarity coefficients or the BM25 coefficients (BM25 Score) are calculated between the terms of the user query and the terms of the texts in the collection. The TF-IDF model follows the calculation of the $TF \times IDF$ product for each term of the collection and for each term of the user query, while the Cosine Similarity and Okapi BM25 models follow the calculation of the cosine similarity coefficient and a score respectively based on these two events. Finally, the search results are sorted by the highest coefficient. The theoretical/mathematical background of the design model is as follows:

RECOVERY OF INFORMATION

TF-IDF

$$TF_{t,d} = \frac{N_{t,d}}{\sum_k N_{k,d}}$$

$TF_{t,d}$ □ The frequency of occurrence of the term t in text d

$N_{t,d}$ □ The number of occurrences of the term t in text d

$\sum_k N_{k,d}$ □ The number of terms in the text d

$$DF_t = |D_t|$$

DF_t □ The frequency of occurrence of texts where the term t appears

$|D_t|$ □ The number of texts where the term t appears

$$IDF_t = \log \log \frac{|D|}{DF_t}$$

IDF_t □ The inverse frequency of texts where the term t occurs. Terms with a high frequency of occurrence in texts are of little importance information and this is due to the fact that the logarithm tends to 0.

$|D|$ □ The number of tasks collection

DF_t □ The frequency of occurrence of texts where the term t appears

Cosine Similarity

$$\cos \cos (q, d) = \frac{q \times d}{\|q\| \times \|d\|} \rightarrow \cos \cos (q, d) = \frac{\sum_{t=1}^n (q_t \times d_t)}{\sqrt{\sum_{t=1}^n (q_t)^2} \times \sqrt{\sum_{t=1}^n (d_t)^2}}$$

$\sum_{t=1}^n (q_t \times d_t)$ □ The inner product of the user query vector with the vector of each text in the collection

$\sqrt{\sum_{t=1}^n (q_t)^2}$ □ The Euclidean distance of the user query vector

$\sqrt{\sum_{t=1}^n (d_t)^2}$ □ The Euclidean distance of the vector of each text in the collection

RECOVERY OF INFORMATION

Okapi BM25

$$(q, d) = \sum_{t=1}^n IDF(q_t) \times \frac{TF(q_t, d) \times (k + 1)}{TF(q_t, d) + k \times (1 - b + b \times \frac{|d|}{avgdlen})}$$

$\sum_{t=1}^n IDF(q_t)$ □ The inverse frequency of occurrence of the texts where the query terms appear

$TF(q_t, d)$ □ The frequency of occurrence of the query terms in the text d

k □ Positive parameter ($k = 1.2$) which controls the terms that appear very often in a text (TF) and usually, this is information of little importance for this text

b □ Positive parameter ($b = 0.75$) which controls the terms that appear very frequently in texts (DF) and usually, this is information of little importance for the collection

$|d|$ □ The size of the text d

$avgdlen$ □ The average size of the collection

Implementation

The ranking algorithms are implemented in a separate module (ranking.py) and include the following routines:

TF-IDF

- TF-IDF calculation of the terms of the tasks:

For each task collection term (from the abstract fields):

1. Text preprocessing ([step 2](#)) and separation into verbal units (tokenization)
2. Calculation of the frequency of each term in the text (Term Frequency)
3. Calculation of the frequency of occurrence of the texts in which each term appears (Document Frequency)
4. Calculation of the inverse frequency of the texts that each term appears (Inverse Document Frequency)
5. Calculation of the frequency of each term in the text (TF × IDF)

- Compute TF-IDF of user query terms:

For each term of the user query (query):

1. Preprocessing of the query (routine of [step 2](#)) and separation into verbal units (tokenization)
2. Calculation of the frequency of each query term in the query (Term Frequency)

RECOVERY OF INFORMATION

3. Calculation of the inverse frequency of occurrence of the texts in which each term of the text appears (Inverse Document Frequency)
4. Calculation of the frequency of each query term in each text ($TF \times IDF$)

Cosine Similarity

- Calculation of the cosine similarity ([Vector Space Model](#)) between the user query and the tasks (Cosine Similarity):
For each text in the collection, its similarity to the user query is calculated as follows:
 1. Compute the inner product of the user query vectors with the vector of each text in the collection
 2. Computing the Euclidean distance of the user query vector
 3. Calculation of the Euclidean distance of the vector of each text in the collection
 4. Calculate the similarity between the user query and the text
- Ranking tasks based on cosine similarity

Okapi BM25

- Calculation of the BM25 Score ([Probabilistic retrieval model](#)) coefficient between the user query and the tasks:
For each text in the collection, its similarity to the user query is calculated as follows:
 1. Query preprocessing ([step 2 routine](#))
 2. Calculation of the size of the task
 3. Calculate average task collection size
 4. Calculation of the frequency of occurrence of the texts in which each term appears (Document Frequency)
 5. Calculation of the inverse frequency of the texts that each term appears (Inverse Document Frequency)
 6. Calculation of frequency of appearance of the query term in each text (Term Frequency)
 7. Calculation of BM25 coefficient for each term of the query
- Ranking of documents based on the BM25 Score

Evaluation

Ranking algorithms are important tools for search engines in document retrieval efficiency. The simple TF-IDF ranking algorithm calculates how important a term is in the collection, while the advanced Cosine Similarity and Okapi BM25 algorithms focus on the similarity between the user query and the texts in the collection. The philosophy of "how close some texts are to the user's query" provides an efficient way to retrieve texts based on the user's query.

5. System evaluation

RECOVERY OF INFORMATION

5.a. Datasets

The data is extracted from arxiv.org and with various scientific fields such as Physics, Mathematics, etc. As a repository of scientific works, it provides a large number of documents and words for the index.

5.b. Evaluation Scenarios

The user can through the search engine search to make any search he wants in contents from a large index, however the final results that will be displayed are limited in size since they fit in the console of each IDE.

5.c. Python libraries

The libraries that were used to create the search engine are several, some of them are json which allows data to be stored in a specific format, request which allows retrieving a page through an HTTP-GET request and bs4 which returns the page in unstructured HTML format (parsing)

5. d. Applications of metrics

The search engine to start in an average test takes about 12 seconds and fetching the relevant data from arxiv.org takes about 12 seconds on average. The papers it collects range from 200 to 800 because as mentioned in the [implementation](#) of the web crawler, it addresses 2-8 corresponding links to retrieve the data that will be stored in the repository (dataset). Indicative assessment measures are:

Precision is the percentage of retrieved texts that are relevant

$$P = \#relevant\ recovered\ texts / \#recovered\ texts$$

Recall is the percentage of relevant texts that have been retrieved

$$R = \#relevant\ recovered\ texts / \#recovered\ texts$$

F1 Score is the percentage balancing precision and recall

$$F1 = 1/2 * (P + R) / (P * R)$$

5.e. Analysis and Improvements

The search engine cannot filter multiple authors and this is something that would make it even easier for the user. The GUI is also quite simple and could be more user-friendly and visually nicer in terms of its layout. Certainly, the efficiency of the algorithms is not the best possible for everyone, but this does not significantly affect the accuracy, let alone the use of the application.

RECOVERY OF INFORMATION

RECOVERY OF INFORMATION



Thank you for your attention.

