

## Data Collection:

I have chosen CISI test collection. Documents, queries, and query judgments (qrels) were parsed from their respective files and converted into appropriate format to deal with them.

## Preprocessing:

To prepare the text data for further work, three functions were implemented

- 1. Function was implemented using the NLTK library to stem the text. Firstly, it tokenizes the passed text into parts, then stems each part.
- 2. Function was implemented using the regex module to ensure that there is no sort of (Special Characters, Tabs, Line Jumps, Extra White Spaces) in the passed text.
- 3. Function was implemented to traverse over the tokenized text and check if there is any tokenized text that exists in the stop words list defined from the NLTK. If it exists, it's removed immediately; else, it passes.

```
def Stem_text(text):
    tokens = word_tokenize(text)
    stemmed_tokens = [stemmer.stem(word) for word in tokens]
    return ' '.join(stemmed_tokens)

def clean(text):
    text = re.sub(r'[\.\,\#\_\|\:\;\!\@\&]', ' ', text) # remove special characters
    text = re.sub(r'\t', ' ', text) # remove tabs
    text = re.sub(r'\n', ' ', text) # remove line jump
    text = re.sub(r'\s+', ' ', text) # remove extra white space
    text = text.strip()
    return text

def remove_stopwords(text):
    tokens = word_tokenize(text)
    filtered_tokens = [word.lower() for word in tokens if word.lower() not in stop_words]
    return ' '.join(filtered_tokens)

def preprocess(sentence):
    sentence = clean(sentence)
    sentence = remove_stopwords(sentence)
    sentence = Stem_text(sentence)
    return sentence

[11] documents_df["Processed_text"] = documents_df["Text"].apply(preprocess).apply(clean)
documents_df["docno"] = documents_df["ID"].astype(str)
documents_df.drop(columns="ID", inplace=True)
```

# Indexing:

This stage focused on creating an efficient index for document retrieval using the DFIndexer. Based on requirements i have developed Data Structure that map each term with it's postings and beside each doc id there is an assignment (number ) this number refer to the frequency of the term inside this posting.

## Built In Method To Map The Inverted Index .

```
Indexing

[13] indexer = pt.DFIndexer("DatasetIndex", overwrite=True)
index_ref = indexer.index(documents_df[["Processed_text"], documents_df[["docno"]])
index = pt.IndexFactory.of(index_ref.toString())

Built in Method

[14] for kv in index.getKeyList(): print( kv.getKey(), " -> " , kv.getValue() )

woles -> term2075 doc=1 tf=1 maxtf=1 @[0 92248 5]
wolf -> term2555 doc=1 tf=1 maxtf=1 @[0 92231 2]
wolves -> term1004 doc=1 tf=1 maxtf=1 @[0 92225 1]
wols -> term827 doc=1 tf=1 maxtf=1 @[0 92237 1]
womens -> term211 doc=1 tf=1 maxtf=1 @[0 92238 1]
wonder -> term6815 doc=1 tf=1 maxtf=1 @[0 92246 2]
wook -> term2923 doc=1 tf=1 maxtf=1 @[0 92251 0]
wook -> term2379 doc=1 tf=1 maxtf=1 @[0 92251 2]
woodford -> term6777 doc=1 tf=1 maxtf=1 @[0 92268 6]
woodward -> term5134 doc=1 tf=1 maxtf=1 @[0 92273 4]
woolf -> term151 doc=1 tf=1 maxtf=1 @[0 92274 0]
woolgie -> term6017 doc=1 tf=1 maxtf=1 @[0 92276 4]
woosier -> term644 doc=1 tf=1 maxtf=1 @[0 92279 2]
woz -> term5441 doc=1 tf=1 maxtf=1 @[0 92280 2]
```

## Adham From Scratch Method to map the terms

```
DS Map From Scratch

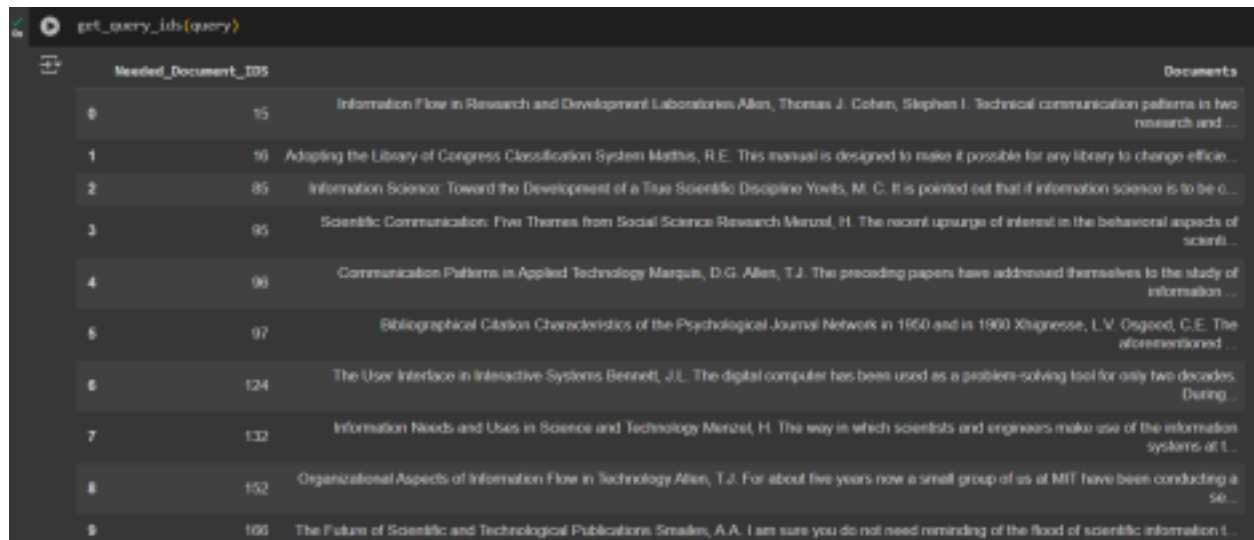
[15] TF = {}
for row in documents_df[["docno", "Processed_text"]].values:
    for term in row[1].split():
        if not term in TF: TF[term] = dict()
        if not row[0] in TF[term]: TF[term][row[0]] = 0
        TF[term][row[0]] += 1

for key in sorted(TF.keys()):
    print(f"{key} -> {TF[key]}")

store -> {'1004': 1, '1005': 1}
stop -> {'204': 1}
store -> {'14': 1, '44': 1, '67': 2, '97': 1, '120': 1, '129': 2, '131': 1, '135': 1, '137': 1, '174': 2, '175': 1, '254': 2, '267': 1, '289': 1, '31'}
store -> {'6': 1, '10': 1, '63': 1, '87': 1, '112': 1, '122': 1, '137': 1, '167': 1, '210': 1, '220': 2, '267': 1, '281': 1, '430': 1, '445': 2, '48'}
storehouse -> {'5': 1, '1329': 1}
store -> {'544': 1, '1330': 1}
store -> {'1': 1, '2': 1, '414': 1, '5051': 2}
store -> {'607': 1, '700': 1}
straight -> {'290': 1, '365': 1, '860': 1}
straightforward -> {'109': 1, '257': 1, '664': 1, '1387': 1}
strains -> {'189': 1, '170': 1, '366': 1, '843': 1, '937': 1, '1201': 2, '1442': 1}
strakhov -> {'1135': 1}
strong -> {'180': 1, '401': 1, '1420': 1}
strong -> {'314': 1, '700': 1}
stratagi -> {'46': 1, '49': 1, '304': 1, '700': 1, '458': 1, '502': 1, '500': 5, '510': 1, '505': 2, '555': 1, '700': 1, '722': 1, '730': 2, '773': 1}
stratifi -> {'80': 5, '93': 1, '954': 1, '1137': 2, '1166': 1}
stratifi -> {'1346': 1}
strains -> {'1310': 1}
```

## Query processing:

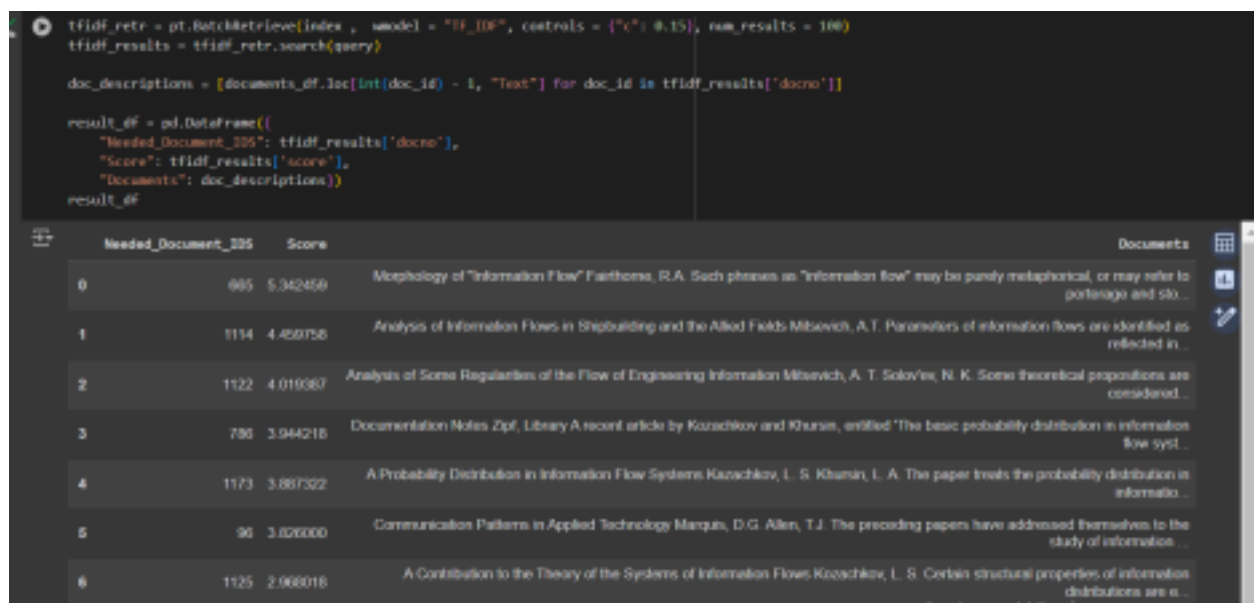
- Received query from the CLI then processed it and called builtin method from pyterrier to map the passed query with relevant documents



```
get_query_ids(query)
```

|   | Needed_Document_IDs | Documents   |
|---|---------------------|---|
| 0 | 15                  | Information Flow in Research and Development Laboratories Allen, Thomas J. Cohen, Stephen I. Technical communication patterns in two research and ... |
| 1 | 16                  | Adopting the Library of Congress Classification System Mathis, R.E. This manual is designed to make it possible for any library to change effice...   |
| 2 | 85                  | Information Science: Toward the Development of a True Scientific Discipline Yovits, M. C. It is pointed out that if information science is to be o... |
| 3 | 95                  | Scientific Communication: Five Themes from Social Science Research Merz, H. The recent upsurge of interest in the behavioral aspects of scienti...    |
| 4 | 96                  | Communication Patterns in Applied Technology Marquis, D.G. Allen, T.J. The preceding papers have addressed themselves to the study of information ... |
| 5 | 97                  | Bibliographical Citation Characteristics of the Psychological Journal Network in 1950 and in 1960 Xhignesse, L.V. Osgood, C.E. The aforementioned ... |
| 6 | 124                 | The User Interface in Interactive Systems Bennett, J.L. The digital computer has been used as a problem-solving tool for only two decades. During ... |
| 7 | 132                 | Information Needs and Uses in Science and Technology Menzel, H. The way in which scientists and engineers make use of the information systems of t... |
| 8 | 152                 | Organizational Aspects of Information Flow in Technology Allen, T.J. For about five years now a small group of us at MIT have been conducting a ...   |
| 9 | 166                 | The Future of Scientific and Technological Publications Smalen, A.A. I am sure you do not need reminding of the flood of scientific information t...  |

- Re ranked the same query using TFIDF



```
tfidf_retr = pt.BatchRetrieve(index, model = "tfidf", controls = {'c': 9.15}, num_results = 100)
tfidf_results = tfidf_retr.search(query)

doc_descriptions = [documents_df.loc[int(doc_id) - 1, "Text"] for doc_id in tfidf_results['docno']]

result_df = pd.DataFrame({
    "Needed_Document_IDs": tfidf_results['docno'],
    "Score": tfidf_results['score'],
    "Documents": doc_descriptions})
result_df
```

|   | Needed_Document_IDs | Score    | Documents  |
|---|---------------------|----------|--|
| 0 | 665                 | 5.342459 | Morphology of "Information Flow" Parthasarathy, R.A. Such phrases as "information flow" may be purely metaphorical, or may refer to portage and sto... |
| 1 | 1114                | 4.459756 | Analysis of Information Flows in Shipbuilding and the Allied Fields Mitsevich, A.T. Parameters of information flows are identified as reflected in ... |
| 2 | 1122                | 4.019387 | Analysis of Some Regularities of the Flow of Engineering Information Mitsevich, A. T. Solov'ev, N. K. Some theoretical propositions are considered ... |
| 3 | 786                 | 3.944218 | Documentation Notes Zpf, Library A recent article by Kozachkov and Khramin, entitled "The basic probability distribution in information flow syst...   |
| 4 | 1173                | 3.887322 | A Probability Distribution in Information Flow Systems Kozachkov, L. S. Khramin, L. A. The paper treats the probability distribution in informatio...  |
| 5 | 96                  | 3.626000 | Communication Patterns in Applied Technology Marquis, D.G. Allen, T.J. The preceding papers have addressed themselves to the study of information ...  |
| 6 | 1125                | 2.968018 | A Contribution to the Theory of the Systems of Information Flows Kozachkov, L. S. Certain structural properties of information distributions are a ... |

# Query Expansion

Here we can find easily by focusing on the score of the screenshot and the above one the influence of the query expansion

```
[25] tfidf_results = tfidf_retr.search(expanded_q)
results_documents = documents_df[documents_df["docno"].isin(tfidf_results["docno"])]
doc_descriptions = [documents_df.loc[int(doc_id) - 1, "Text"] for doc_id in tfidf_results["docno"]]

result_df = pd.DataFrame({
    "Needed_Document_IDs": tfidf_results['docno'],
    "Score": tfidf_results['score'],
    "Documents": doc_descriptions})
result_df
```

|     | Needed_Document_IDs | Score     | Documents   |
|-----|---------------------|-----------|---|
| 0   | 1122                | 36.285919 | Analysis of Some Regularities of the Flow of Engineering Information Mitsevich, A. T. Solov'ev, N. K. Some theoretical propositions are considered... |
| 1   | 1340                | 10.147523 | Social Theory and Social Structure Merlon, R.K. Of the four chapters, added to this edition, two come from published symposia, one of which is out... |
| 2   | 1114                | 9.731981  | Analysis of Information Flows in Shipbuilding and the Allied Fields Mitsevich, A.T. Parameters of information flows are identified as reflected in... |
| 3   | 735                 | 9.000457  | Journals Most Cited by Chemists and Chemical Engineers Barrett, R.L. Barrett, M.A. The purpose of this paper is to present up-to-date material to...  |
| 4   | 788                 | 9.241198  | Documentation Notes Zipf, Library A recent article by Kozachkov and Khursin, entitled 'The basic probability distribution in information flow syst... |
| ... | ...                 | ...       | ...   |
| 95  | 1200                | 3.330770  | On the Statistics of Individual Variations of Productivity in Research Laboratories Shockley, W. In the following pages a co-winner of the 1956 No... |

Then Reranked the docs using Elmo Embeddings

|   | score    | docs  |
|---|----------|---|
| 0 | 0.581055 | Analysis of Some Regularities of the Flow of Engineering Information Mitsevich, A. T. Solov'ev, N. K. Some theoretical propositions are considered... |
| 1 | 0.332099 | Documentation Notes Zipf, Library A recent article by Kozachkov and Khursin, entitled 'The basic probability distribution in information flow syst... |
| 2 | 0.242589 | Analysis of Information Flows in Shipbuilding and the Allied Fields Mitsevich, A.T. Parameters of information flows are identified as reflected in... |
| 3 | 0.176654 | A Probability Distribution in Information Flow Systems Kazachkov, L. S. Khursin, L. A. The paper treats the probability distribution in informatio... |
| 4 | 0.041194 | Morphology of "Information Flow" Fairthorne, R.A. Such phrases as "information flow" may be purely metaphorical, or may refer to portorage and sto... |

# User Interface

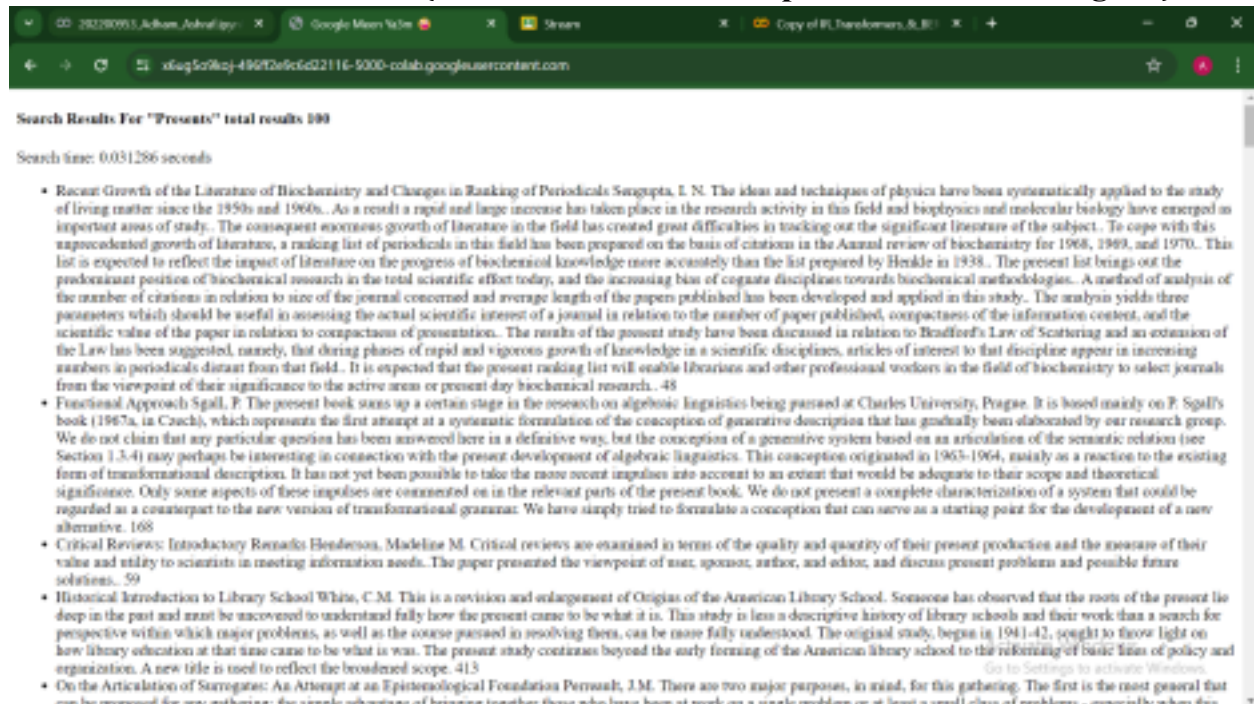
Using Flask I have decided to built very basic web application to show off my work :

**In the Search Bar I have Typed “ Presents “**



Activate Windows  
Go to Settings to activate Windows.

**As Shown In the screenshot {Total Num of results , Speed of the search engine }**



25220095\_Ashwin\_Ashwin@y... Google Meet Na... Stream Copy of PLTransformers\_6... x6q5o7kzj-499T2e5c6d2116-5000-colab.googleusercontent.com

Search time: 0.020258 seconds

Search Results For "Cost Accounting" total results 100

Search time: 0.014243 seconds

- **Managerial Cost Accounting for a Technical Information Center** Helton, John G. The purpose of this paper is to describe a research project conducted at a technical center to test the hypothesis that: A theoretically-sound managerial cost-accounting system can be designed to meet the specific characteristics of a technical information center by reusing and innovating systems utilized by other enterprises. A computerized cost system was developed and operated for a three-month period to test this hypothesis. The results of the study indicate that effective managerial cost accounting is possible for a technical information center. Relevant cost information was generated periodically to measure the operating performance of the center's production process. A summary of the data that were reported regularly to management is presented in this paper. 24
- **Cost Accounting and Analysis for University Libraries** Leinhardt, Ferdinand F. Cooper, Michael D. The approach to library planning studies in this paper is the use of accounting models to measure library costs and implement program budgets. A cost-flow model for a university library is developed and tested with historical data from the General Library at the University of California, Berkeley. Various comparisons of an exploratory nature are made of the unit costs and total costs for different parts of the Berkeley system. 5
- **Library Cost Analysis: A Recipe** Kowitz, J. Unfortunately, time has passed since the days when the library's patron was the local monarch and cost was no deterrent. Time's passage has replaced the monarch with taxpayers or stockholders, and, consequently, sensitivity to cost has attained stellar importance. The causes for being unaware of costs may stem from a variety of reasons, but they cannot, in all fairness to the profession, belie an inability to perform the simple arithmetic of cost accounting. What is suspected is a lack of the few simple ground rules and the logical operations that bind them together, in short - a recipe for cost accounting and analysis. In the following is outlined one such set of ground rules and their related procedural requirements, which have evolved and been applied with success over the past few years. It is stressed that since this set represents the findings of one library, it may not fully satisfy the specific requirements of your own shop. Therefore, feel free to adapt the ground rules to your immediate requirements. With regard to discipline, it is pretty much summed up in the six steps and five resource requirements which follow. In addition to identifying steps, requirements, and the mysterious ways of cost analysis, these ingredients are blended together in a manner which will be meaningful for your internal operations and may be significant for your library's future. 90
- **Cost Comparison of Manual and On-Line Computerized Literature Searching** Elman, S.A. Cost and searching time comparisons are made between manual and on-line literature searches. The formula  $\text{Costal} = (T \times \text{Costm}) + P$  is presented which captures all on-line cost factors. A minimum cost of \$1.00 per minute of on-line searching is derived. Average searching time for manual searching is 22 hours at a total cost of \$230; for on-line it is 45 minutes at total cost of \$47.00. It is pointed out that most reported low-on-line search costs fail to account for all cost factors. Figures are those prevailing at the time of writing. 124
- **Library Participation in a Biomedical Communication and Information Network** Bridgman, Willis E., Jr. Meyerhoff, Erich The experience of two libraries participating in the SUNY Biomedical Communication Network is described. The history of the Network is briefly given together with its original aims and their current status. Use of the terminals and formulation of queries are explained. Figures are given for total costs, number of searches performed, and cost per search. There is a account of the internal structure of the administration of the Network. 65
- **Cost Accounting for the Library** Bratcher, C. Cleveland, G. Rinford, E. Increasingly, librarians have felt the need for more accurate cost data. The prime reason for this need has been in the development and presentation of the budget which is the instrument used to determine and obtain the funds for the library's forthcoming fiscal period. Since libraries do not charge for the service they render their users, they must derive the funds necessary for their operations and growth from supporting bodies such as federal, state, or local governments, private institutions, and industrial firms. 27
- **The Shared Cataloging Systems of the Ohio College Library Center** Kilgore, Frederick G. Long, Philip L. Londgraf, Alan L. Wyckoff, John A. Development and implementation of an off-line catalog card production system and on-line shared cataloging system are described. In off-line production, average cost per card for \$29.89 catalog cards in finished form and

# Evaluation

Here We can find the evaluation for the search engine

```
[34] evaluation = pt.Evaluate(tfidf_retr.transform(topics), qrels, metrics=["map", "recall", "P"], perquery=True)
      evaluation_df = pd.DataFrame(evaluation).T
      evaluation_df
```

|   | map      | pg5 | pg10 | pg15 | pg20 | pg10 | pg100 | pg200 | pg500 | pg1000 | pg5      | pg10     | pg15     | pg20    | pg50     | pg100   | pg200   | pg500   | pg1000  |
|---|----------|-----|------|------|------|------|-------|-------|-------|--------|----------|----------|----------|---------|----------|---------|---------|---------|---------|
| 1 | 0.373060 | 0.6 | 0.5  | 0.4  | 0.55 | 0.5  | 0.34  | 0.170 | 0.000 | 0.034  | 0.005217 | 0.100666 | 0.130425 | 0.23913 | 0.320087 | 0.73813 | 0.73613 | 0.73913 | 0.73913 |

*Screenshots Purpose For Illustrating My steps and Showing My results*  
*You can find the same results clearly in the [Note Book](#)*