Queen Mary
**University of London**

# Experimental Search Engine Strategy and Demo
## Prepared for Information Retrieval

**Prepared by: Danielle Souza Da Silva, Georgia Dean, Adam Toth**

**14 April 2022**

# Agenda

**Overview of Original Proposal**
    Overview of Experiment
    Overview of Architecture
**Data**
    Data Sources
    Data Scraping Pipeline
    Data Structures
**Tooling and Equations**
    Tools
    Indexing Architecture
    Model Configurations
    Retrieval Architecture
**Demo**
**Results**
**Conclusion**

# Overview of Experiment

## Review of goals for experimental search engine

Our experimental search engine uses BM25F model to improve retrieval results of regulatory text.

**Why BM25F?** *BM25*, considered one of the most effective retrieval models, considers texts as unstructured, undifferentiated in any way. As an extension of this model, *BM25F* considers documents to be composed of several fields (title, abstract, body, etc.) with possibly different degrees of importance, term relevance saturation and length normalisation.
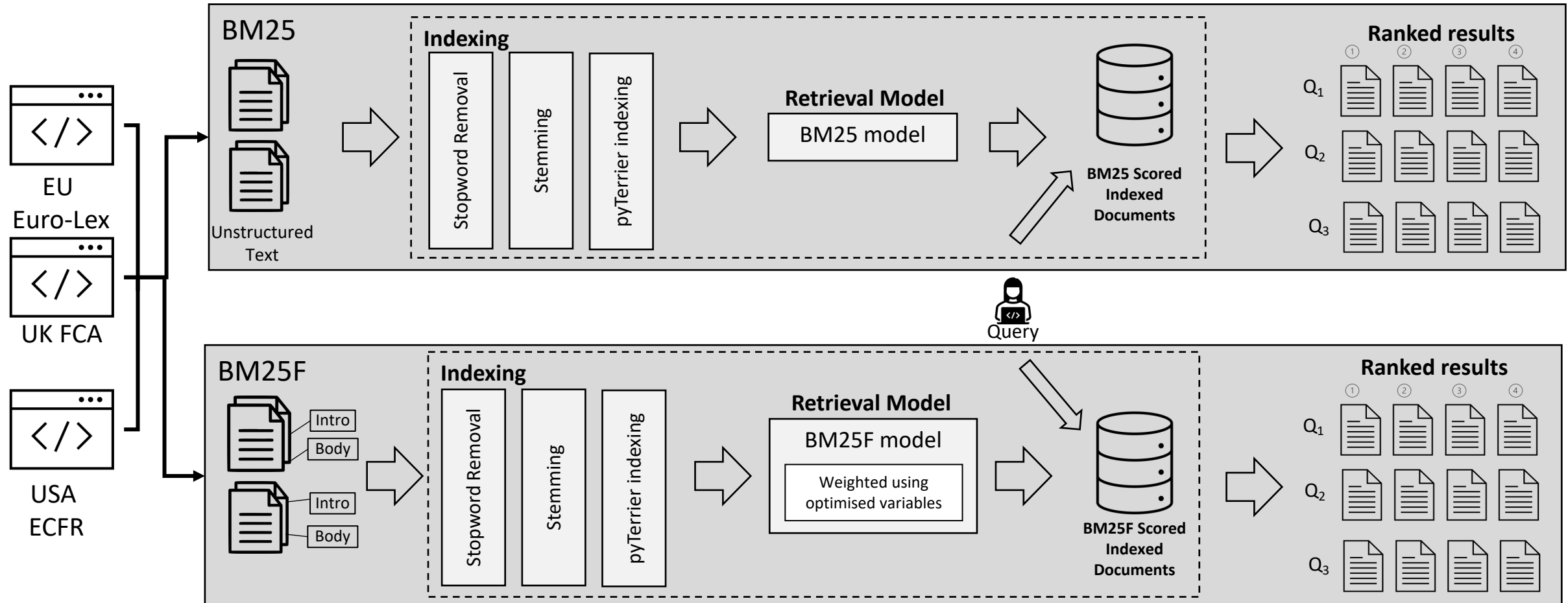
**Why regulatory texts?** Regulatory texts tend to have consistent structure and unique approach to topic inclusion - highly relevant concepts are mentioned in opening sections, often without repetition in the body. Therefore, some fields may be more predictive of relevance than others.

**Goal:** by computing and assigning relative weights to pre-specified fields of the documents, the goal is to improve upon the retrieval results of the benchmark model  BM25.

**Evaluation:** performance is measured using mean average precision and recall.

# Data Sources

**Dataset**

The dataset used for this experimental model was comprised regulatory text that was scraped by the researchers from various regulatory bodies' websites The sources are regulatory texts from the United States, United Kingdom and European Union focusing specifically on financial regulation for a total of **2387** documents.

| Jurisdiction | Regulations | Number of documents in corpus |
| --- | --- | --- |
| United States | E-CFR (Title 12, Title 17) | 478 |
| United Kingdom | FCA Handbook | 504 |
| European Union | Eur-Lex Financial Regulations | 1405 |

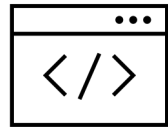**Query relevant document pairs**

The query relevant document pairs were completed by domain experts in the regulatory compliance field.

Total number of queries: 12

Range of relevant documents per query: 575

# Data Scraping Pipeline

## Process used to scrape data

| | | | | |
|---|---|---|---|---|
| HTML scraped from web pages and XML from APIs | HTML and XML cleaned and parsed in Python Azure ML to extract data fields | Processed data stored in Azure SQL Server Database for rapid retrieval | Relevant metadata and text rendered as dictionary for use in experimental search engine model | Model build in Google Colab Python notebook using PyTerrier library |

# Data Structures

## Explanation of how data was structured

### ECFR (USA)



### Euro-Lex (EU)



### FCA Handbook(UK)

# Tools

## Tools utilised for the search engine architecture proposed.

Python was used to implement the search engine components. The libraries that were utilised to scrape and prepare the data, implement the search engine and run experiments with results are listed below.

| Python libraryName | Description of use in proposed implementation |
|---|---|
| Beautiful Soup | It was used to scrape HTML from the web containing the relevant documents. |
| Requests | It was used to apply HTTP request and get URLs desired for the dataset. |
| Pickle | This library was used to move data. |
| Pandas | It was used to read the csv file and transform it to a data frame. |
| Pyterrier | It was used to index the data, apply the retrieval models, tune the parameters (field weights) for the BM25F model and run experiments. |

The main library used to implement each step of the architecture proposed was PyTerrier.

**PyTerrier** is a platform used for information retrieval experiments in Python. It uses Java-based Terrier information retrieval platform to support indexing and retrieval operations.

References: [pyterrier.readthedocs.io.], [terrier.org.],

# Indexing Architecture

## Method used for the search engine architecture proposed.

**Indexing**

1. A corpus is represented in the form of a collection object. Documents are provided with an instance of a Tokeniser class that breaks pieces of text into single indexing tokens.

2. The indexer manages the indexing process. It iterates over the documents of the collection and sends each term found through a TermPipeline component, which where stemming and stop word removal takes place (default PorterStemmer).

3. Once the terms have been processed through the TermPipeline, they are aggregated and data structures are created by their corresponding DocumentBuilders.

**Indexing**

The graphic below gives an overview of the interaction between the main components involved in the indexing process.

Corpus

Indexing

Collection

public boolean nextDocument()
public Document getDocument()

Document

Tokeniser
public TokenStream tokenise(Reader reader)
public String[] getTokens(Reader reader)

public String getNextTerm()
public Set<String> getFields()

TermPipeline

public void processTerm(String term)

Indexer

Data Structures

Reference: [GitHub. (2022). *terrier-org/terrier-core*]

# Model Configurations

## How the models are configured in PyTerrier

**Retrieval Models**

- **BM25**: PyTerrier uses a java class to implement the Okapi BM25 weighting model.

$$w_i^{BM25}(tf) = \frac{tf}{k_1\left((1-b) + b\frac{dl}{avdl}\right) + tf} \cdot w_i^{RSJ}$$

- **BM25F**: PyTerrier uses a java subclass of PerFieldNormWeightingModel setup to implement BM25F as described by [Zaragoza TREC-2004], which is the same set up described in the search engine proposal.

$$\widetilde{tf}_i = \sum_{s=1}^{S} v_s \frac{tf_{si}}{B_s}$$

- The parameters k1 and b are the default for both models and were kept consistent for the experiments.

$$B_s = \left((1-b_s) + b_s\frac{sl_s}{avsl_s}\right), \qquad 0 \le b \le 1$$

$$w_i^{BM25F} = \frac{\widetilde{tf}_i}{k_1 + \widetilde{tf}_i} w_i^{RSJ}$$

References: [pyterrier.readthedocs.io.], [Robertson, S. and Zaragoza, H., 20]

# Retrieval Architecture

## Steps of how retrieval results are computed

1.  The "Application" in the diagram in our case refers to the Google Colab notebook which in the first step issues a query

2.  The query is parsed and an instantiation of a query object takes place.

3.  The query is then handled to the Manager component, which pre-processes the query by applying it to the configured TermPipeline.

4.  The pre-processed query is handled to the Matching component which initializes the Weighting Model (BM25/BM25F) and Document Score Modifiers. Once all components have been instantiated the score computation with respect to the query will take place.

5.  The PostProcessing and PostFiltering takes place and so the score document list is returned to the application.

Reference: [GitHub. (2022). *terrier-org/terrier-core*]

# DEMO

# Results

## Overall results using BM25 vs BM25F

### Overall Results

| Model | MAP | Recall@5 | Recall@10 | Recall@15 | Recall@20 | Recall@30 | Recall@100 | Recall@200 | Recall@500 | Recall@1000 |
|---|---|---|---|---|---|---|---|---|---|---|
| BM25F | 56.24% | 22.37% | 33.87% | 41.14% | 48.23% | 58.09% | 74.40% | 82.65% | 89.55% | 93.28% |
| BM25 | 57.69% | 21.55% | 34.97% | 42.63% | 49.96% | 59.55% | 76.97% | 84.29% | 89.75% | 93.12% |

Recall@k (R@k): The fraction of relevant documents for a query that have been retrieved by rank k.

### Parameters

Parameters k1 and b used were the default values for both models. The field weights for the BM25F model can be seen below:

| Field | Weight |
|---|---|
| Text | 0.5 |
| Preamble | 1.0 |

# Results

## Results per query using BM25 vs BM25F

| Query | Model | MAP | Recall@5 | Recall@10 | Recall@15 | Recall@20 | Recall@30 | Recall@100 | Recall@200 | Recall@500 | Recall@1000 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Commodity derivative | BM25F | 52.43% | 8.77% | 17.54% | 21.05% | 22.81% | 29.82% | 63.16% | 87.72% | 100.00% | 100.00% |
| | BM25 | 53.24% | 8.77% | 15.79% | 21.05% | 21.05% | 26.32% | 71.93% | 92.98% | 100.00% | 100.00% |
| Commodity pool operator | BM25F | 62.95% | 10.81% | 18.92% | 32.43% | 40.54% | 59.46% | 81.08% | 94.59% | 100.00% | 100.00% |
| | BM25 | 70.96% | 10.81% | 21.62% | 35.14% | 45.95% | 62.16% | 97.30% | 100.00% | 100.00% | 100.00% |
| Derivatives clearing organizations | BM25F | 80.99% | 26.32% | 52.63% | 63.16% | 68.42% | 84.21% | 100.00% | 100.00% | 100.00% | 100.00% |
| | BM25 | 86.32% | 26.32% | 52.63% | 63.16% | 78.95% | 89.47% | 100.00% | 100.00% | 100.00% | 100.00% |
| Escheatment | BM25F | 100.00% | 83.33% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% |
| | BM25 | 100.00% | 83.33% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% |
| Liquidity risk | BM25F | 13.84% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 1.89% | 26.42% | 76.42% | 99.06% |
| | BM25 | 14.12% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.94% | 21.70% | 80.19% | 99.06% |
| Major swap participant | BM25F | 81.94% | 16.13% | 32.26% | 48.39% | 61.29% | 74.19% | 87.10% | 93.55% | 96.77% | 100.00% |
| | BM25 | 82.70% | 16.13% | 32.26% | 48.39% | 58.06% | 77.42% | 90.32% | 96.77% | 96.77% | 100.00% |
| National bank | BM25F | 0.50% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 1.42% | 20.28% |
| | BM25 | 0.42% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 0.00% | 18.40% |
| Physical commodity swaps | BM25F | 57.99% | 42.86% | 42.86% | 57.14% | 85.71% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% |
| | BM25 | 44.09% | 28.57% | 42.86% | 57.14% | 85.71% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% |
| Proprietary trading | BM25F | 32.44% | 13.64% | 22.73% | 27.27% | 36.36% | 40.91% | 81.82% | 100.00% | 100.00% | 100.00% |
| | BM25 | 39.52% | 18.18% | 27.27% | 31.82% | 40.91% | 45.45% | 77.27% | 100.00% | 100.00% | 100.00% |
| Swap data repositories | BM25F | 56.14% | 21.05% | 42.11% | 47.37% | 47.37% | 52.63% | 89.47% | 89.47% | 100.00% | 100.00% |
| | BM25 | 60.07% | 21.05% | 42.11% | 47.37% | 52.63% | 57.89% | 94.74% | 100.00% | 100.00% | 100.00% |
| Swap execution facility | BM25F | 66.34% | 14.71% | 23.53% | 35.29% | 47.06% | 55.88% | 88.24% | 100.00% | 100.00% | 100.00% |
| | BM25 | 68.21% | 14.71% | 23.53% | 38.24% | 47.06% | 55.88% | 91.18% | 100.00% | 100.00% | 100.00% |
| Whistleblower | BM25F | 69.28% | 30.77% | 53.85% | 61.54% | 69.23% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% |
| | BM25 | 72.62% | 30.77% | 61.54% | 69.23% | 69.23% | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% |

# Conclusion

## Final thoughts and further research for BM25F and regulation information retrieval

**Conclusions**

- BM25 is still the model to beat!

- BM25F showed potential for the more difficult queries where even BM25 had low performance.
  - The preamble is still an important section of regulation and the highly valuable content it contains continues to prove evasive for search engines to capture in relevance rankings.

**Further research**

- Field length normalization
  - Preambles had varying lengths that could affect the results. Further research into field length normalization

- Improved use of queries for retrieval
  - It would be beneficial for multi-word queries to rank documents that contain 'word 1' AND 'word2' higher

- N-grams
  - Using n-grams could better capture the proximity of terms in the text and queries to improve relevant document retrieval for domain-specific topics that are made up of general terms

# Division of Work

| Task | Team Member | Description | Percentage of work |
|---|---|---|---|
| Data Import and Structure | GD | Bring in and prepare data, both structured and unstructured | 15% |
| Define Queries | GD | Define 15-20 queries and relevant document pairs | 5% |
| Indexing of Data | DD, AT | Using team coding, index data to be scored-Including but not limited to, stop word removal, stemming, tf/idf calculations | 5% |
| Implementation of BM25 | DD, AT | Using team coding, implement BM25 model on the unstructured document data. | 20% |
| Implementation of BM25F | GD, DD, AT | Using team coding, implement BM25F model on structured document data | 20% |
| Optimization of Parameters | DD, AT | Optimize BM25F weight parameters and hyper parameters for both models (potentially using gradient descent) | 15% |
| Results Evaluation | GD | Calculate precision, recall and F-scores to determine quality of retravel results | 2.5% |
| Presentation Write-Up | GD, DD, AT | Write presentation of | 15% |
| Record Presentation/Demo | GD, DD, AT | All team members will record and present search engine and results | 2.5% |

# Sources

Na, S., Kang, I., & Lee, J. (2008). Improving Term Frequency Normalization for Multi-topical Documents and Application to Language Modeling Approaches. *ECIR*.

Singhal, A., Buckley, C., & Mitra, M. (1996). Pivoted document length normalization. *SIGIR Forum, 51*, 176-184.

pyterrier.readthedocs.io. (n.d.). *Installing and Configuring — PyTerrier 0.8.1 documentation*. [online] Available at: https://pyterrier.readthedocs.io/en/latest/installation.html.

terrier.org. (n.d.). *Terrier IR Platform - Homepage*. [online] Available at: http://terrier.org [Accessed 14 Apr. 2022].

GitHub. (2022). *terrier-org/terrier-core*. [online] Available at: https://github.com/terrier-org/terrier-core/blob/5.x/doc/basicArchitecture.md [Accessed 14 Apr. 2022].

[Zaragoza TREC-2004] . H. Zaragoza, N. Craswell, M. Taylor, S. Saria, S. Robertson: Microsoft Cambridge at TREC 13: Web and Hard Tracks. In Proc. of TREC 2004

[Robertson, S. and Zaragoza, H., 2009.] The Probabilistic Relevance Framework: BM25 and Beyond. *Foundations and Trends® in Information Retrieval*, 3(4), pp.333-389.