# Contributions to a system for Open Reproducible Publication Research

## Topic Classification of Publications

Dayoung Lim <D.Lim-2@student.tudelft.nl>
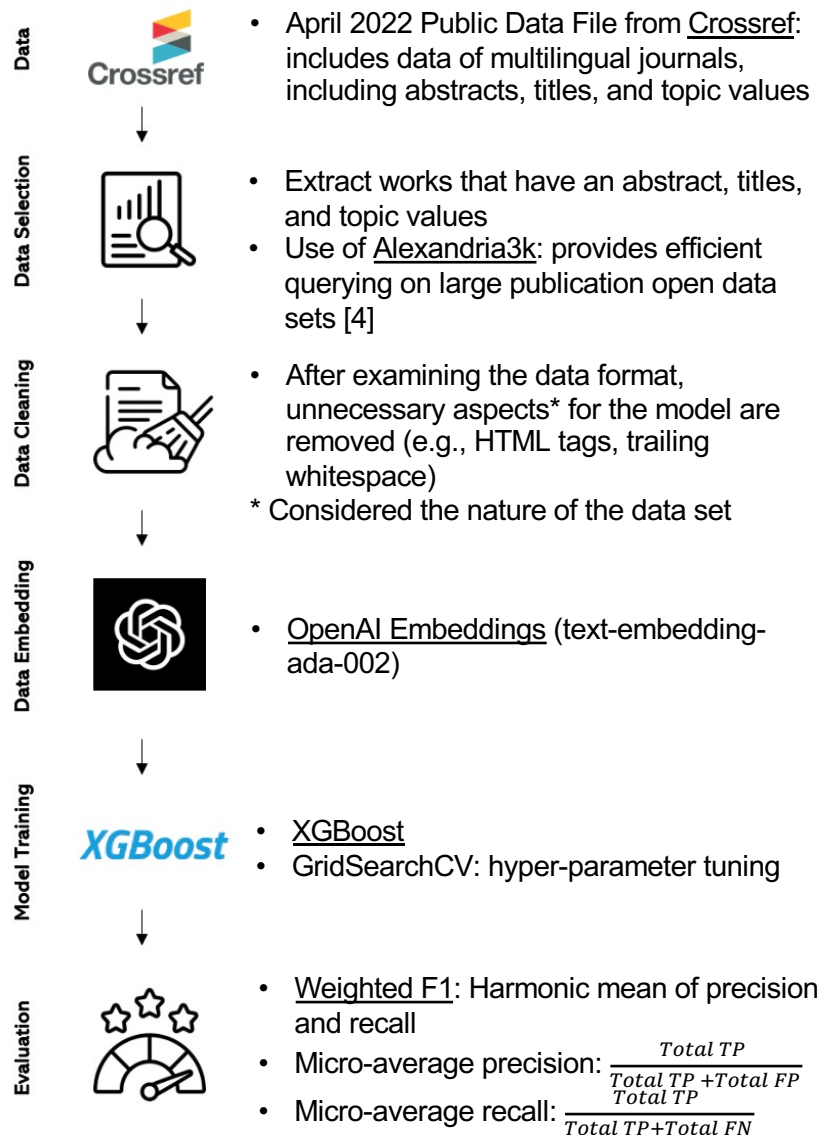
**TU**Delft

## 1. Introduction

**Background**
- Volume of published journals and difficulty in finding journals are increasing [1] → correctly <u>classify</u> publications

**Research Gap**
- Existing classification works with short text are mainly sentence based [2]
- <u>Abstract</u> based classifications are mainly domain specific [3]

**Research Question:** *How can publication topics be identified and matched based on existing journal topic values?*

## 2. Methodology

**Data**
- April 2022 Public Data File from <u>Crossref</u>: includes data of multilingual journals, including abstracts, titles, and topic values

**Data Selection**
- Extract works that have an abstract, titles, and topic values
- Use of <u>Alexandria3k</u>: provides efficient querying on large publication open data sets [4]

**Data Cleaning**
- After examining the data format, unnecessary aspects* for the model are removed (e.g., HTML tags, trailing whitespace)

\* Considered the nature of the data set

**Data Embedding**
- <u>OpenAI Embeddings</u> (text-embedding-ada-002)

**Model Training**
- <u>XGBoost</u>
- GridSearchCV: hyper-parameter tuning

**Evaluation**
- <u>Weighted F1</u>: Harmonic mean of precision and recall
- Micro-average precision: $\frac{Total\ TP}{Total\ TP + Total\ FP}$
- Micro-average recall: $\frac{Total\ TP}{Total\ TP + Total\ FN}$

## 3. Results

**Performance - Initial Run**
- 10,000 data and 50 topic values
- Grid search → max_depth=6, eta=0.5, n_estimators=500
    - max_depth: depth of the tree
    - eta: learning rate
    - n_estimators: number of boosting rounds or trees to build
- Baseline model (<u>BM25 + XGBoost</u>) comparison
    - Difference is in data cleaning stage
    - XGBoost parameters for both experiment are the same
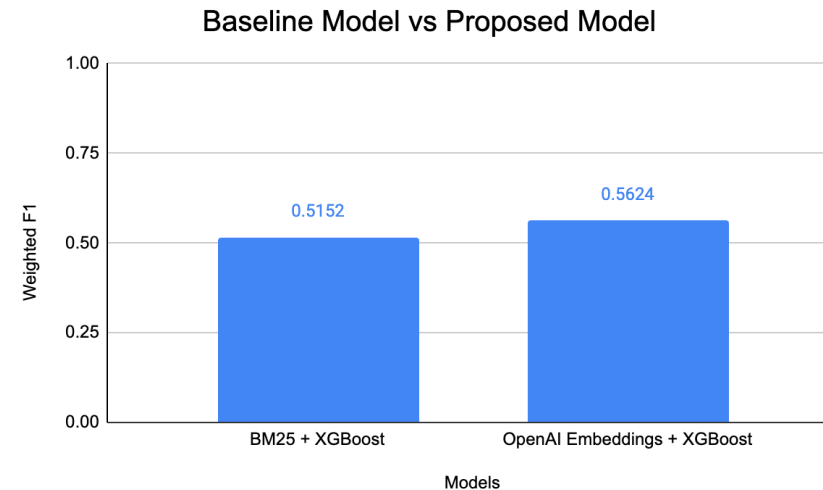- Different Features: Abstract, Abstract + Title, Abstract + Title + Author

Figure 1. Comparison of weighted f1 for baseline and proposed model

**Performance – Final Run**
- 50,000 <u>stratified</u> data
- Abstract + Title
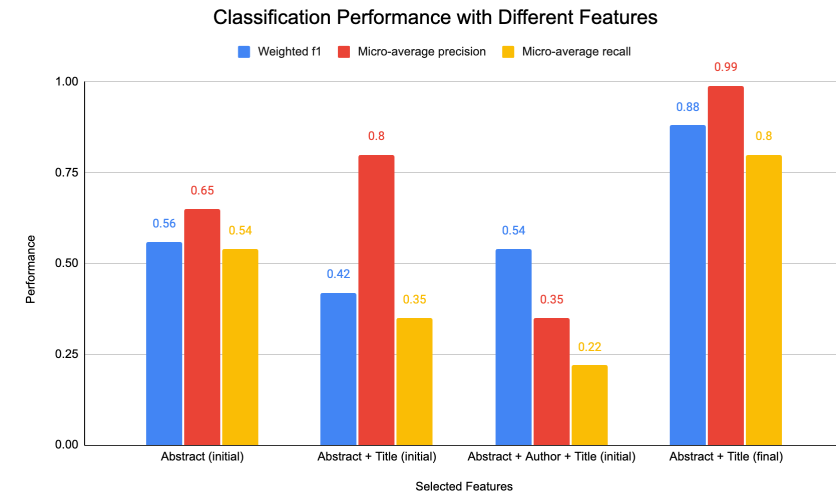- Grid search → max_depth=20, eta=0.8, n_estimators=1000

Figure 2. Performance comparison of initial runs with different features and the final run

**Cost**
- $0.0001/1k token
    - For 10,000 data (50 topic values): ~$0.2
    - For 50,000 stratified data~$2.8

## 4. Conclusion

**Research Question Answer**
- OpenAI Embeddings + XGBoost combination can be used for publication topic classification when the right features are chosen

**Limitations**
- High computational cost → only on a sample of Crossref
    - Classification verified on data with work names
- Correctness of original data has **not** been checked

**Future works**
- Test on works <u>without</u> work names
- Verify its performance on other publication data set
- Usage of <u>newer model</u> for embeddings: text-embedding-3-small/large [5]

## References

[1] AHMED SAJID, N., AHMAD, M., RAHMAN, A., ET AL. 2023. A novel metadata based multi-label document classification technique. *Computer Systems Science and Engineering 46*, 2, 2195–2214.

[2] LEE, K., PALSETIA, D., NARAYANAN, R., PATWARY, MD.M., AGRAWAL, A., AND CHOUDHARY, A. 2011. Twitter trending topic classification. *2011 IEEE 11th International Conference on Data Mining Workshops*.

[3] DEEPIKA, A. AND RADHA, N. 2021. Performance analysis of abstract-based classification of medical journals using Machine Learning Techniques. *Computer Networks and Inventive Communication Technologies*, 613–626.

[4] SPINELLIS, D. alexandria3k: Local relational access to openly-available publication data sets. *GitHub*. https://github.com/dspinellis/alexandria3k.

[5] ZHUANG, J., BRAUNSTEIN, A., NEELAKANTAN, A., ET AL. 2024. *New embedding models and API updates*. https://openai.com/blog/new-embedding-models-and-api-updates.

Responsible Professor: Diomidis Spinellis <D.Spinellis@tudelft.nl>
Supervisor: Georgios Gousios <G.Gousios@tudelft.nl>


Figure 1. Comparison of weighted f1 for baseline and proposed model
Baseline Model vs Proposed Model — BM25 + XGBoost: 0.5152, OpenAI Embeddings + XGBoost: 0.5624


Figure 2. Performance comparison of initial runs with different features and the final run
Classification Performance with Different Features (Weighted f1, Micro-average precision, Micro-average recall)
- Abstract (initial): 0.56, 0.65, 0.54
- Abstract + Title (initial): 0.42, 0.8, 0.35
- Abstract + Author + Title (initial): 0.54, 0.35, 0.22
- Abstract + Title (final): 0.88, 0.99, 0.8