

Nearest Neighbor Search over Vectorized Lexico-Syntactic Patterns for Relation Extraction from Financial Documents

Pawan Kumar Rajpoot *

UtilizeAI Research

Bangalore, India

pawan.rajpoot2411@gmail.com

Ankur Parikh *

UtilizeAI Research

Bangalore, India

ankur.parikh85@gmail.com

Abstract

Relation extraction (RE) has achieved remarkable progress with the help of pre-trained language models. However, existing RE models are usually incapable of handling two situations: implicit expressions and long-tail relation classes, caused by language complexity and data sparsity. Further, these approaches and models are largely inaccessible to users who don't have direct access to large language models (LLMs) and/or infrastructure for supervised training or fine-tuning. Rule-based systems also struggle with implicit expressions. Apart from this, Real world financial documents such as various 10-X reports (including 10-K, 10-Q, etc.) of publicly traded companies pose another challenge to rule-based systems in terms of longer and complex sentences. In this paper, we introduce a simple approach that consults training relations at test time through a nearest-neighbor search over dense vectors of lexico-syntactic patterns and provides a simple yet effective means to tackle the above issues. We evaluate our approach on REFinD and show that our method achieves state-of-the-art performance. We further show that it can provide a good start for human in the loop setup when a small number of annotations are available and it is also beneficial when domain experts can provide high quality patterns. Our code is available at ¹.

1 Introduction

Relation extraction (RE) from text is a fundamental problem in NLP and information retrieval, which facilitates various tasks like knowledge graph construction, question answering and semantic search. Recent studies (Zhang et al., 2020; Zeng et al., 2020; Lin et al., 2020; Wang and Lu, 2020; Cheng et al., 2020; Zhong and Chen, 2021) in supervised RE take advantage of pre-trained language models

¹https://github.com/pawan2411/PAN-DL_Refind

*Equal Contribution


On December 16, 2016, Uni Line Corp entered into a share purchase agreement (the Purchase Agreement) with Porter Group Limited (PGL) to acquire all issued and outstanding shares of PGL. 

Figure 1: Relation Extraction example, here both organizations are connected with "acquired by" relation.

(PLMs) and achieve SOTA performances by fine-tuning PLMs with a relation classifier. However, (Wan et al., 2022) observes that existing RE models are usually incapable of handling two RE-specific situations: implicit expressions and long-tail relation types.

Implicit expression refers to the situation whereas relation is expressed as the underlying message that is not explicitly stated or shown. In Figure 1, relation "acquired_by(organization, organization)" occurs implicitly. Such underlying messages can easily confuse the relation classifier.

The other problem of long-tail relation classes is caused by data sparsity in training. For example, the REFinD dataset (Kaur et al., 2023) comprises 45.5 % of the no_relation instances. The most frequent class in the dataset - "per:title:title" has 4,468 training examples, while over 14 out of 22 classes have less than 500 examples. The majority class can easily dominate model predictions and lead to low performance on long-tail classes.

Recently, ICL (In-Context Learning) based approach (Wan et al., 2023) is utilized for RE tasks. The approach achieves improvements over not only existing GPT-3 baselines, but also on fully-supervised baselines even with only a limited number of demonstrations provided in the prompt. Specifically, it achieves SOTA performances on the Semeval and SciERC datasets, and competitive performances on the TACRED and (Zhang et al., 2017a) ACE05 datasets. (Rajpoot and Parikh, 2023) utilized the GPT-4 under ICL framework on REFinD and achieved 3rd rank in the shared task.

However, retrieval of examples to demonstrate is a key factor in the overall performance on these

pipelines. Finding efficient demonstrates often relies on learning-based retrieval (Ye et al., 2023; Rubin et al., 2022). These learning-based retrievers use annotated data and a LLM. This type of retrieval strategy comes with the increased cost (API, infrastructure etc.), time as more experiments are required because most LLMs are black box and it also needs special expertise.

Apart from the implicit expression challenge mentioned above, REFinD poses another challenge to rule-based systems in terms of longer and complex sentences. For example, (Kaur et al., 2023) cites that the average sentence length in the REFinD dataset is 53.7 while the average sentence length in the TACRED dataset (Zhang et al., 2017b) is 36.2. Further, As per (Kaur et al., 2023), REFinD includes more complex sentences than TACRED, with an average entity-pair distance of 11, compared to 8 in TACRED. Because of this, writing rules at surface text level is a challenge. Hence, rules at lexico-syntactic level is the need of the hour. However, strict matching of these rules can yield high precision but low recall result due to accuracy of syntactic parsing. Hence, a robust fuzzy pattern matching system is required.

Inspired by recent studies (Wan et al., 2022; Khandelwal et al., 2019; Guu et al., 2020; Meng et al., 2021) using k-Nearest Neighbor to retrieve diverse expressions for language generation tasks, we introduce a simple but effective approach that consults training relations at test time through a nearest-neighbor search over dense vectors of lexico-syntactic patterns and provides a simple yet effective means to tackle the above issues. Our method achieves an improvement of 1.18% over baseline (F1-score - 0.7516). We achieved our results using commodity hardware within a day. That’s why our approach is easier to deploy, lightweight and fast. We further show that our approach can provide a good start (F1-score of 0.5122) for human in the loop setup when a small number of annotations (approx. 10% of training data) are available and it is also beneficial (F1-score of 0.6939 with approx. 10% of training data) when domain experts can provide high quality patterns.

2 Preliminary Background

2.1 Task Definition

Let C denote the input context and $e1$ in C , $e2$ in C denote the pair of entity pairs. Given a set of predefined relations classes R , relation extraction

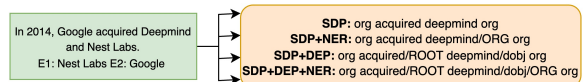


Figure 2: Patterns extracted by our pipeline

aims to predict the relation y in R between the pair of entities ($e1$, $e2$) within the context C , or if there is no predefined relation between them, predict $y="no relation"$.

2.2 Data

The REFinD dataset (Kaur et al., 2023) is the largest relation extraction dataset for financial documents to date. Overall REFinD contains around 29K instances and 22 relations among 8 types of entity pairs. REFinD is created using raw text from various 10-X reports (including 10-K, 10-Q, etc. broadly known as 10-X) of publicly traded companies obtained from US Securities and Exchange Commission.

3 Nearest Neighbor Search over Vectorized Lexico-Syntactic Patterns

3.1 Generating Lexico-Syntactic Patterns

We replaced words representing entities of interest with their entity types given in the dataset.

Instead of conducting nearest neighbor search on a complete sentence, we applied Spacy Dependency Parser² and considered the shortest dependency path (henceforth SDP) between two entities to deal with long and complex sentences with the intuition that considering all sentence words can do more harm in search. SDP is essential for relationship identification in most cases.

We apply Spacy NER on REFinD sentences and replace actual named entities with their types to create higher-level patterns.

We also enriched all SDP words with their Dependency Labels to utilize structure information in our search.

For each sentence, we create 4 patterns: 1. SDP words only (SDP) 2. SDP words with named entities replaced with their types (SDP-NER) 3. SDP words enriched with their Dependency Labels (SDP-DEP) 4. SDP words with named entities replaced with their types and also enriched with their Dependency Labels (SDP-DEP-NER). Example patterns are shown in Figure 2.

²<https://spacy.io/>

3.2 Generating Dense Vectors for Lexico-Syntactic Patterns

We converted all 4 types of Lexico-Syntactic Patterns into Dense Vectors as it performs better than Sparse Vectors. To create a vector, we employed an all-mpnet-base-v2³ sentence encoder. We also created vectors for original sentences using the encoder.

3.3 Creating Class Specific Indices

For each pattern type mentioned above, we created 21 dense vector indices each representing a relationship class except 'no_relation' class. We split 'no_relation' training data instances into 8 splits as per entity-type pairs such as "Person-Organization", "Organization-Organization" etc. and created indices for each split. In this way, there are 29 indices in total for each pattern type. Each element of the index represents a vectorized lexico-syntactic pattern for each training example. For around 11.89% of the training sentences, we faced issues in generating dependency tree and/or SDP. To deal with this, we also created another 29 indices containing dense vectors for original sentences.

3.4 Conducting Nearest Neighbor Search

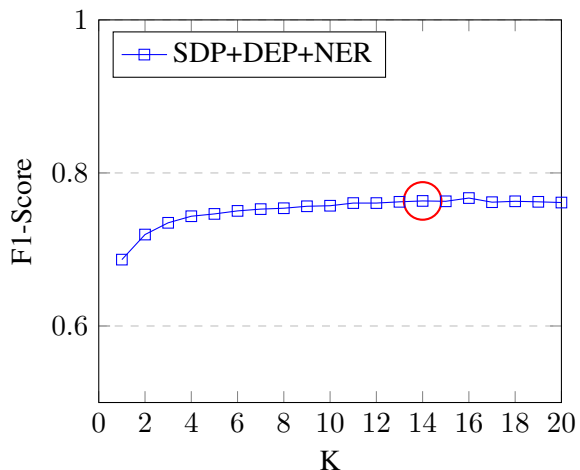
After configuring lexico-syntactic pattern type and value of K, Given a test sentence and an entity-type pair, we first create a vector representing its lexico-syntactic pattern obtained using steps described above. With the entity-type pair, appropriate relation class indices are selected for search. The pattern vector is searched in every appropriate class index using cosine similarity and top K vectors from each class index are obtained. The similarity scores of each of these top K vectors are averaged and the class having the highest similarity score is selected. In the case of syntactic parsing failures, as a fallback strategy, a vector of the original sentence is created and is used against class specific sentence indices in search the same way as mentioned above.

³<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

Pattern	K	F1-score
SDP	14	0.7552
SDP-NER	12	0.7538
SDP-DEP	11	0.7610
SDP-DEP-NER	14	0.7634
Winner on leaderboard (baseline) ⁴	-	0.7516

Table 1: Comparison on performance on REFinD dev data

Figure 3: Sensitivity Analysis



4 Experiment Settings

4.1 Dataset

The REFinD dataset (Kaur et al., 2023) released with the shared task is a part of "Knowledge Discovery from Unstructured Data in Financial Services" (KDF) workshop which is collocated with SIGIR 2023. There are 20070, 4306 and 4300 instances of training data, development data and public test data respectively. The organizers have released training data, development data and public test data with gold labels but haven't released private test data with gold labels. Because of that, we are not able to benchmark our system against the winners of the shared task. Since, leaderboard⁴ and gold labels on development data is available, we have benchmarked our approach against the leaders of development data. We have used training data and public test data to create class specific indices to perform nearest neighbor search for development data sentences.

4.2 Hardware Resources

We have used a laptop with 16GB RAM and Intel® Core™ i7-7500U CPU @ 2.70GHz × 4 CPU to

⁴<https://codalab.lisn.upsaclay.fr/competitions/11770>

produce these results.

4.3 Efforts

Given the dataset, all setup and experiments are conducted within a day.

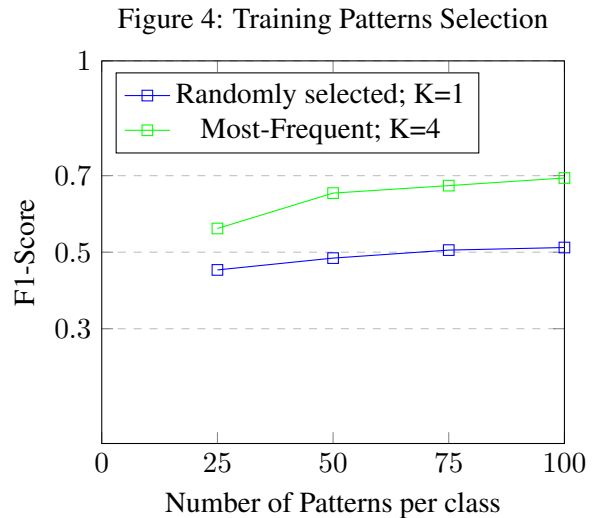
5 Results

We conducted experiments with 4 different pattern types. To find the best value of K , we have created a 10% split from the training data and experimented with different values of K (1 to 20). During evaluation, we faced issues in generating dependency tree and/or SDP for around 8.7% instances and for those instances, indices containing sentence vectors were used as fallback strategy. The results in Table 1 show that our best F1-score is 0.7634 for SDP-DEP-NER pattern and $K=14$ (Top K vectors) and our method shows improvement of 1.18% over baseline. Figure 3 shows how sensitive this approach is with respect to different values of K . This confirms our intuition that there is value in utilizing vectorized lexico-syntactic patterns to deal with long and complex sentences and implicit expressions. Further, splitting instances as per the class and performing lazy classification over these splits can help in dealing with the dataset with long-tail relation classes.

To explore the effectiveness of our approach in human in the loop situation, we conducted a few experiments as shown in Figure 4. We randomly selected N patterns per class from the training data and built indices with those patterns only. We tried different values of N . With $N=100$ and $K=1$ (derived from dev split), we achieved an F1-score of 0.5122 with around 10% of the original training data. It shows that the vectorized lexico-syntactic patterns and the cosine similarity based scoring can be a good start to label similar instances when the number of annotations are less. This method can be used in human in the loop setup to either filter out similar instances (explore) or to find similar instances (exploit) for further human review/annotation.

To explore the effectiveness of our approach when domain experts are available and can provide high quality patterns specially for the task like this which is restricted to a particular domain, types of documents, types of entities and a handful of relations, we conducted a few experiments as shown in Figure 4. To approximate this experiment, we selected N training patterns from each class which

occurs frequently in Top K search when development patterns are classified correctly. We call these patterns the Most-Frequent Patterns. We built indices with those patterns only. As shown in Figure 4, with $N=100$ and $K=4$ (derived from dev split), we achieved an F1-score of 0.6939 (6.95% less than the best result) with around 10% of the original training data. It shows that it can bridge gaps quickly with a small amount of high quality patterns.



6 Conclusion

Our approach consults training relations at test time through a nearest-neighbor search over dense vectors of lexico-syntactic patterns. We evaluated our approach on REFinD and show that our method achieves state-of-the-art performance without any direct access to large language models (LLMs) or supervised training or fine-tuning or any hand-crafted rules. We achieved our results using a commodity hardware within a day. That’s why our approach is easier to deploy, lightweight and fast. We further explore that our approach can provide a good start for human in the loop setup when a small number of annotations are available and it can be also beneficial when domain experts can provide high quality patterns.

7 Limitations

Since our method is based on nearest neighbors search, it’s sensitive to the value of K . Furthermore, our method is also very sensitive to syntactic parsing and NER. Our vectors are not optimal representations because our syntactic patterns are not a natural fit for the sentence encoder.

References

- Fei Cheng, Masayuki Asahara, Ichiro Kobayashi, and Sadao Kurohashi. 2020. [Dynamically updating event representations for temporal relation classification with multi-category learning](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1352–1357, Online. Association for Computational Linguistics.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.
- Simerjot Kaur, Charese Smiley, Akshat Gupta, Joy Sain, Dongsheng Wang, Suchetha Siddagangappa, Toyin Aguda, and Sameena Shah. 2023. Refind: Relation extraction financial dataset. *arXiv preprint arXiv:2305.18322*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. Generalization through memorization: Nearest neighbor language models. *arXiv preprint arXiv:1911.00172*.
- Ying Lin, Heng Ji, Fei Huang, and Lingfei Wu. 2020. [A joint neural model for information extraction with global features](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7999–8009, Online. Association for Computational Linguistics.
- Yuxian Meng, Shi Zong, Xiaoya Li, Xiaofei Sun, Tianwei Zhang, Fei Wu, and Jiwei Li. 2021. Gnn-lm: Language modeling based on global contexts via gnn. *arXiv preprint arXiv:2110.08743*.
- Pawan Kumar Rajpoot and Ankur Parikh. 2023. [Gpt-fnre: In-context learning for financial relation extraction using large language models](#).
- Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. [Learning to retrieve prompts for in-context learning](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2671, Seattle, United States. Association for Computational Linguistics.
- Zhen Wan, Fei Cheng, Zhuoyuan Mao, Qianying Liu, Haiyue Song, Jiwei Li, and Sadao Kurohashi. 2023. [Gpt-re: In-context learning for relation extraction using large language models](#).
- Zhen Wan, Qianying Liu, Zhuoyuan Mao, Fei Cheng, Sadao Kurohashi, and Jiwei Li. 2022. [Rescue implicit and long-tail cases: Nearest neighbor relation extraction](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1731–1738, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jue Wang and Wei Lu. 2020. [Two are better than one: Joint entity and relation extraction with table-sequence encoders](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1706–1721, Online. Association for Computational Linguistics.
- Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023. [Compositional exemplars for in-context learning](#). *arXiv preprint arXiv:2302.05698*.
- Daojian Zeng, Haoran Zhang, and Qianying Liu. 2020. Copymtl: Copy mechanism for joint extraction of entities and relations with multi-task learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 9507–9514.
- Ranran Haoran Zhang, Qianying Liu, Aysa Xuemo Fan, Heng Ji, Daojian Zeng, Fei Cheng, Daisuke Kawahara, and Sadao Kurohashi. 2020. [Minimize exposure bias of Seq2Seq models in joint entity and relation extraction](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 236–246, Online. Association for Computational Linguistics.
- Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017a. [Position-aware attention and supervised data improve slot filling](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)*, pages 35–45.
- Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017b. [Position-aware attention and supervised data improve slot filling](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 35–45, Copenhagen, Denmark. Association for Computational Linguistics.
- Zexuan Zhong and Danqi Chen. 2021. [A frustratingly easy approach for entity and relation extraction](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 50–61, Online. Association for Computational Linguistics.