

ADTB Project: Semantic Search using PostgreSQL

Konstantinos Giantsios
giantsik@csd.auth.gr

Contents

1 Data	1
2 Populating DB	2
2.1 Creation of semantic vectors	2
2.2 Populating DB and querying	3
3 Web application	4
4 Installation	4
5 Experiments	7
5.1 Results for k=5	7
5.2 Results for k=10	9
5.3 Results for k=20	13

1 Data

The data that were utilized to populate the database are coming from the Semantic Scholar Academic Corpus ¹. In particular a subset of the Semantic Scholar Open Research Corpus (S2ORC) [1] was used. S2ORC is a large corpus of 81.1M English-language academic papers spanning many academic disciplines. Rich metadata, paper abstracts, resolved bibliographic references, as well as structured full text for 8.1M open access papers. For the construction of our subset we have taken only 3k papers that are written in english, have abstract and have a link to open access repository to download the full text.

A slight preprocessing has been applied before I feed forward them to the model to get the semantic vectors and upload them to the database. The text from the title and abstract was lower-cased and cleaned from any control character ('n', 't') and latex remnants. Moreover, I have downloaded the full text for only 1k papers to save space. I uploaded to the database

¹<https://www.semanticscholar.org/product/api>

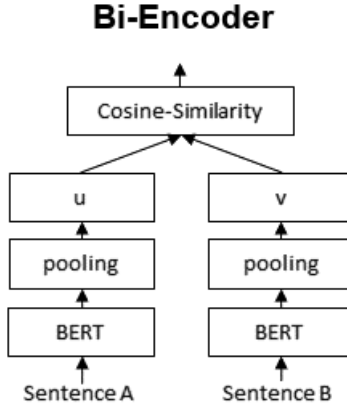


Figure 1: Illustration of a Bi-encoder. [2]

only those papers that have full text. Finally, I have chosen to keep only 4 fields: id, title, abstract, local_link, where local link is the local path to the full text of the paper. The full text of the paper is in pdf format. You can find

2 Populating DB

2.1 Creation of semantic vectors

A well known method to extract semantic vectors from a text is by encoding your text using an model. Bi-Encoders produce for a given text an embedding. In Figure 1 they pass to a BERT independently the sentences A and B, which result in the sentence embeddings u and v. These sentence embedding can then be compared using cosine similarity.

Bi-Encoders are used whenever you need a sentence embedding in a vector space for efficient comparison. Applications are for example Information Retrieval / Semantic Search or Clustering. With a Bi-Encoder, you compute the embedding for each sentence, which takes only 5 seconds [2].

For the creation of our embeddings/vectors we used the all-mpnet-base-v2 model ². This model is an all-round model tuned for many use-cases. Trained on a large and diverse dataset of over 1 billion training pairs. Specifically, the pairs included a lot of with academic context and it ranked second in the semantic search performance benchmark ³ held by sbert framework ⁴, this is the reason why I chose it. Moreover, this model is based on mpnet-base ⁵. MPNet is a novel pre-training method that inherits the advantages of BERT and XLNet and avoids their limitations. MPNet leverages the dependency among predicted tokens through permuted language modeling (vs. MLM in BERT), and takes auxiliary position information as

²<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

³https://www.sbert.net/docs/pretrained_models.html

⁴<https://www.sbert.net/index.html>

⁵<https://huggingface.co/microsoft/mpnet-base>

Model	Factorization	
MLM	$\log P(\text{sentence} \mid \text{the task is [M] [M]})$	$+\log P(\text{classification} \mid \text{the task is [M] [M]})$
PLM	$\log P(\text{sentence} \mid \text{the task is})$	$+\log P(\text{classification} \mid \text{the task is sentence})$
MPNet	$\log P(\text{sentence} \mid \text{the task is [M] [M]})$	$+\log P(\text{classification} \mid \text{the task is sentence [M]})$

Figure 2: An example sentence “the task is sentence classification” to illustrate the conditional information of MLM, PLM and MPNet. [3]

```
[sem_search=# SELECT column_name FROM information_schema.columns WHERE TABLE_NAME = 'Paper';
column_name
-----
embedding
id
title
abstract
local_link
(5 rows)
```

Figure 3: Table Paper Schema

input to make the model see a full sentence and thus reducing the position discrepancy (vs. PLM in XLNet) [3]. In Figure 2 you observe the illustration of conditional information of MLM, PLM and MPNet.

2.2 Populating DB and querying

Furthermore, I concatenated the title and abstract and feed-forward the text that was created to get the corresponding embeddings/vectors for every paper. I populated the PostgreSQL database with the 1k papers using 5 fields: id, title, abstract, local_link and embedding, as you can see in Figure 3. Every embedding has 768 dimensions and is normalized.

For the storing of the embedding the creation of my database I have utilised the pgvector plugin ^{6 7}, which was based in PASE [4]. Using pgvector, I could create an indexing scheme called IVFFlat. IVFFlat uses a clustering algorithm such as k-means to divide vectors in the high-dimensional data space into clusters based on implicit clustering properties. Each cluster has a centroid. IVFFlat traverses the centroids of all clusters to identify the n centroids that are nearest to the vector you want to query. IVFFlat traverses and sorts all vectors in the clusters to which the identified n centroids belong. Then, IVFFlat obtains the nearest k vectors. I have decided to use the docker version of pgvector to be more manageable.

For the querying part, given a query and a k, the app first encodes the query to a vector and then makes a query call to get the k nearest neighbours based on the embedding column. The similarity is calculated using the dot-product. After the retrieval of the k nearest neighbours a re-ranking is applied using a different distance function given by the user. Obviously, if the user selects the dot-product as the scoring function, then no re-ranking is applied. List of

⁶<https://github.com/pgvector/pgvector>

⁷<https://github.com/pgvector/pgvector-python>

re-ranking distance functions:

- Bray-Curtis distance = $\sum |u_i - v_i| / \sum |u_i + v_i|$
- Canberra distance = $\sum_i \frac{|u_i - v_i|}{|u_i| + |v_i|}$
- Chebyshev distance = $\max_i |u_i - v_i|$
- City Block (Manhattan) distance = $\sum_i |u_i - v_i|$
- Cosine distance = $1 - \frac{u \cdot v}{\|u\|_2 \|v\|_2}$
- Euclidean distance = $\|u - v\|_2$

Finally, I have added the capability to return an explanation (the lines from the texts that contain the corresponding terms / keywords that match what the user user) utilising fuzzy matching with a high threshold (0.9) and taking the neighbouring words (span) of the matched word.

3 Web application

To facilitate the user, I have created an API using the FastAPI⁸ python library. To access the User Interface the user have to go to this address: <http://127.0.0.1:8000/docs>, assuming that you run the app locally.

First, the user will see the starting page of the UI, as you can see in Figure 4. Here, the user has 3 options. The first one is for health checking. The second one is for the searching according to a given query. The third one is for opening the full text of a paper, given a local path.

For the searching/querying in UI, shown in Figure 5 the user will need to click on /search, then click the Try Out button and then the user will insert the query text, the k nearest papers that he wants the app to return and the distance metric for re-ranking. The user can read the results if he scrolls down, as shown in Figure 6. The app returns the relevant papers showing the title, the local path to the full text and an explanation as described previously. Finally, if the user wants to read the full text, he can click on /open_file, then click on Try Out button and insert the local path, as shown in Figure 7. The app will open a new tab in the default browser with the full text of the paper.

4 Installation

It is recommended to have a Linux-based OS, Python 3.7+ and docker.

1. Clone Github repo:

```
git clone https://github.com/CoGian/semantic-search-engine.git
```

⁸<https://fastapi.tiangolo.com/>

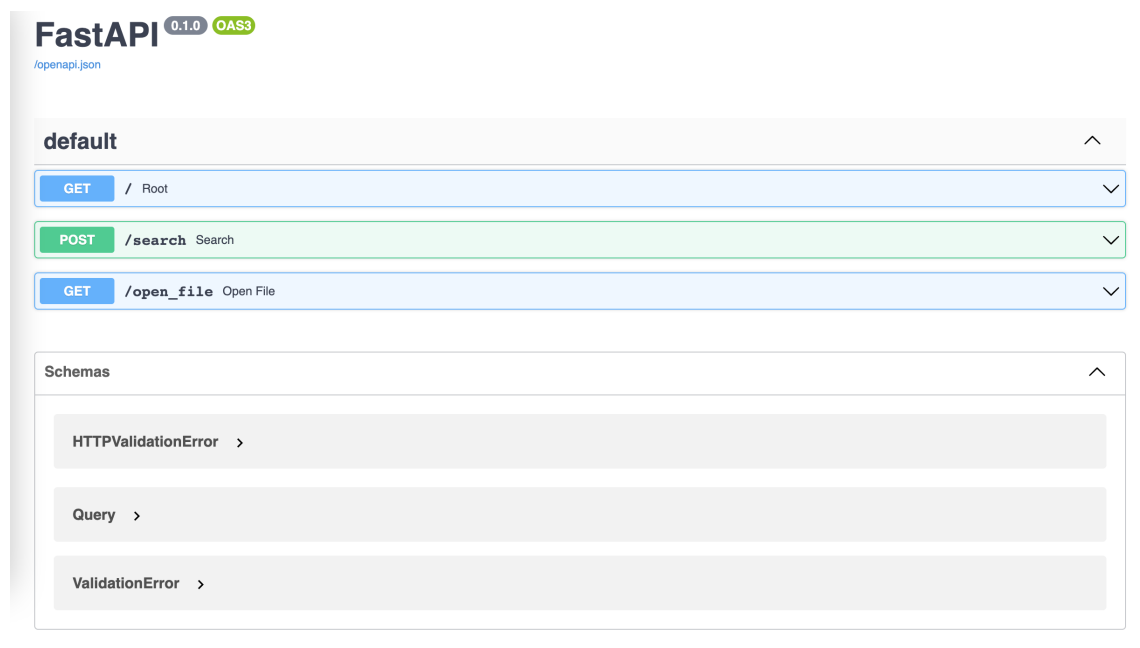


Figure 4: UI starting page



Figure 5: Search/Querying in app

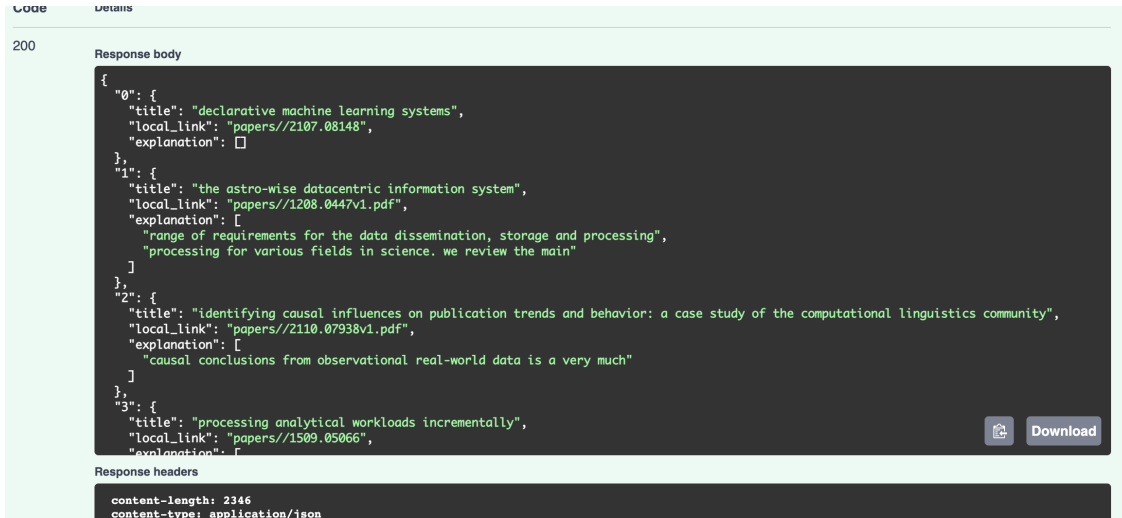


Figure 6: Results from search/querying in app

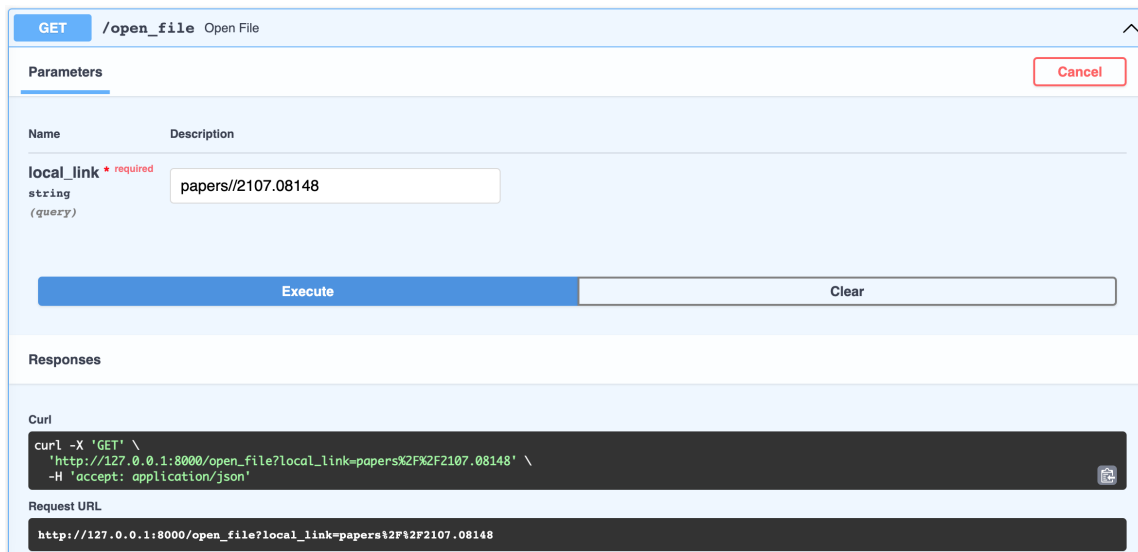


Figure 7: Open file using the app

2. Download & install PostgreSQL

3. Change directory to the repo's directory

```
cd semantic-search-engine
```

4. Download postgresql dump and and install all the dependencies needed :

```
bash install.sh
```

5. Open terminal inside docker container:

```
docker exec -it sem_search_postgres bash
```

6. Connect to postgres:

```
psql -U postgres
```

7. Create DB and close connection:

```
CREATE DATABASE sem_search;
```

8. Load psql dump and exit docker's container:

```
psql -U postgres sem_search < sem_search_export.pgsql
```

9. Run application:

```
source venv/bin/activate
uvicorn src.api:app --reload
```

5 Experiments

Query = "machine learning deep learning natural language processing"

5.1 Results for k=5

- distance=dot_product

```
1 {
2   "0": {
3     "title": "declarative machine learning systems",
4     "local_link": "papers//2107.08148",
5     "explanation": [
6       "declarative machine learning systems the future"
7     ]
8   },
9   "1": {
10    "title": "neural speed reading with structural-jump-lstm",
```

```

11     "local_link": "papers//1904.00761v2.pdf",
12     "explanation": [
13         "neural networks (rnns) can model natural language by sequentially '
14         reading'"
15     ],
16     "2": {
17         "title": "pangu-a: large-scale autoregressive pretrained chinese
18             language models with auto-parallel computation",
19         "local_link": "papers//2104.12369",
20         "explanation": [
21             "pangu-a: large-scale autoregressive pretrained chinese language
22             models with auto-parallel computation",
23             "with auto-parallel computation large-scale pretrained language
24             models (plms) have become",
25             "become the new paradigm for natural language processing (nlp). plms"
26             ,
27             "have demonstrated strong performances on natural language
28             understanding and generation",
29             "and generation with \\textit{few-shot in-context} learning. in this
30             work, we",
31             "practice on training large-scale autoregressive language models
32             named pangu-, with"
33         ]
34     },
35     "3": {
36         "title": "parts-of-speech tagger errors do not necessarily degrade
37             accuracy in extracting information from biomedical text",
38         "local_link": "papers//0804.0317",
39         "explanation": [
40             "of muscorian, a generic text processing tool for extracting protein-
41             protein",
42             "comparable performance to biomedical-specific text processing tools.
43             this result was"
44         ]
45     },
46     "4": {
47         "title": "clar: a cross-lingual argument regularizer for semantic role
48             labeling",
49         "local_link": "papers//2011.04732v1.pdf",
50         "explanation": [
51             "a given sentence. although different languages have different
52             argument annotations,",
53             "share common semantic meaning across languages (e.g. adjuncts have

```



```

42     more",
    "such linguistic annotation similarity across languages and exploits
      this information",
43     "information to map the target language arguments using a
      transformation",
44     "the space on which source language arguments lie. by doing",
45     "improves srl performance on multiple languages over monolingual and
      polyglot"
46 ]
47 }
48 }

```

- distance=cosine, same as previous
- distance=euclidean, same as previous
- distance=braycurtis, same as previous
- distance=canberra Changed from dot-product results,
position in dot product : position in canberra
0:0, 3:1, 2:2, 4:3, 1:4
- distance=chebyshev Changed from dot-product results,
position in dot product : position in chebyshev
0:0, 3:1, 2:2, 1:3, 4:4
- distance=cityblock, same as dot product results

In my opinion, I believe that the only irrelevant paper that has returned as result is the first one, because it doesn't refer to natural language processing. Therefore, we can agree that we have a 80% precision.

5.2 Results for k=10

- distance=dot_product

```

1 {
2   "0": {
3     "title": "declarative machine learning systems",
4     "local_link": "papers//2107.08148",
5     "explanation": [
6       "declarative machine learning systems the future"
7     ]
8   },
9   "1": {
10    "title": "neural speed reading with structural-jump-lstm",

```

```

11     "local_link": "papers//1904.00761v2.pdf",
12     "explanation": [
13         "neural networks (rnns) can model natural language by sequentially '
14         reading'"
15     ],
16     "2": {
17         "title": "pangu-a: large-scale autoregressive pretrained chinese
18             language models with auto-parallel computation",
19         "local_link": "papers//2104.12369",
20         "explanation": [
21             "pangu-a: large-scale autoregressive pretrained chinese language
22             models with auto-parallel computation",
23             "with auto-parallel computation large-scale pretrained language
24             models (plms) have become",
25             "become the new paradigm for natural language processing (nlp). plms"
26             ,
27             "have demonstrated strong performances on natural language
28             understanding and generation",
29             "and generation with \\textit{few-shot in-context} learning. in this
30             work, we",
31             "practice on training large-scale autoregressive language models
32             named pangu-, with"
33         ]
34     },
35     "3": {
36         "title": "parts-of-speech tagger errors do not necessarily degrade
37             accuracy in extracting information from biomedical text",
38         "local_link": "papers//0804.0317",
39         "explanation": [
40             "of muscorian, a generic text processing tool for extracting protein-
41             protein",
42             "comparable performance to biomedical-specific text processing tools.
43             this result was"
44         ]
45     },
46     "4": {
47         "title": "clar: a cross-lingual argument regularizer for semantic role
48             labeling",
49         "local_link": "papers//2011.04732v1.pdf",
50         "explanation": [
51             "a given sentence. although different languages have different
52             argument annotations,",
53             "share common semantic meaning across languages (e.g. adjuncts have

```

```

    more",
42     "such linguistic annotation similarity across languages and exploits
        this information",
43     "information to map the target language arguments using a
        transformation",
44     "the space on which source language arguments lie. by doing",
45     "improves srl performance on multiple languages over monolingual and
        polyglot"
46 ]
47 },
48 "5": {
49     "title": "opinion mining of text documents written in macedonian
        language",
50     "local_link": "papers//1411.4472",
51     "explanation": [
52         "text documents written in macedonian language the ability to extract
            ",
53         "and surveys. in this paper machine learning techniques are used",
54         "which are written in macedonian language. the posts are classified"
55     ]
56 },
57 "6": {
58     "title": "training a restricted boltzmann machine for classification by
        labeling model samples",
59     "local_link": "papers//1509.01053v1.pdf",
60     "explanation": [
61         "training a restricted boltzmann machine for classification by
            labeling",
62         "classification performance competitive to semi-supervised learning
            if we first train"
63     ]
64 },
65 "7": {
66     "title": "identifying causal influences on publication trends and
        behavior: a case study of the computational linguistics community",
67     "local_link": "papers//2110.07938v1.pdf",
68     "explanation": [
69         "trending tasks and techniques (e.g., deep learning, embeddings,
            generative, and",
70         "china on propensity of researching languages beyond english, and the
            "
71     ]
72 },
73 "8": {

```

```

74     "title": "determining whether the non-protein-coding dna sequences are
        in a complex interactive relationship by using an artificial
        intelligence method",
75     "local_link": "papers//1708.04019v1.pdf",
76     "explanation": [
77         "ai research have enabled automatically learning representations of
            high dimensional",
78         "data without feature engineering, using deep neural networks.
            therefore, in",
79         "whether a dna sequence is natural or artificial. the trained"
80     ]
81 },
82 "9": {
83     "title": "deep poetry: a chinese classical poetry generation system",
84     "local_link": "papers//1911.08212v1.pdf",
85     "explanation": [
86         "deep poetry: a chinese classical",
87         "classical poetry generation system called deep poetry. existing
            systems for",
88         "multi-modal input. unlike previous systems, deep poetry uses neural
            networks"
89     ]
90 }
91 }

```

- distance=cosine, same as previous
- distance=euclidean, same as previous
- distance=braycurtis Changed from dot-product results,
position in dot product : position in braycurtis
0:0, 1:1, 2:2, 3:3, 4:4, 6:5, 5:6, 8:7, 9:8, 7:9
- distance=canberra Changed from dot-product results,
position in dot product : position in canberra
0:0, 3:1, 2:2, 4:3, 1:4, 6:5, 5:6, 8:7, 9:8, 7:9
- distance=chebyshev
Changed from dot-product results,
position in dot product : position in chebyshev
0:0, 3:1, 2:2, 1:3, 6:4, 9:5, 7:6, 5:7, 8:8, 4:9
- distance=cityblock
Changed from dot-product results,
position in dot product : position in cityblock

0:0, 1:1, 2:2, 3:3, 4:4, 6:5, 5:6, 7:7, 9:8, 8:9

In my opinion, I believe that here we have 3 irrelevant papers that have returned as result, because they don't refer to natural language processing. Therefore, we can agree that we have a 70% precision. Furthermore, we can observe that changing the distance metric, can change the ranking of the returned results.

5.3 Results for k=20

- distance=dot_product

```
1 {
2   "0": {
3     "title": "declarative machine learning systems",
4     "local_link": "papers//2107.08148",
5     "explanation": [
6       "declarative machine learning systems the future"
7     ]
8   },
9   "1": {
10    "title": "neural speed reading with structural-jump-lstm",
11    "local_link": "papers//1904.00761v2.pdf",
12    "explanation": [
13      "neural networks (rnns) can model natural language by sequentially '
14        reading'"
15    ]
16  },
17  "2": {
18    "title": "pangu-a: large-scale autoregressive pretrained chinese
19      language models with auto-parallel computation",
20    "local_link": "papers//2104.12369",
21    "explanation": [
22      "pangu-a: large-scale autoregressive pretrained chinese language
23        models with auto-parallel computation",
24      "with auto-parallel computation large-scale pretrained language
25        models (plms) have become",
26      "become the new paradigm for natural language processing (nlp). plms"
27      ,
28      "have demonstrated strong performances on natural language
29        understanding and generation",
30      "and generation with \\textit{few-shot in-context} learning. in this
31        work, we",
32      "practice on training large-scale autoregressive language models
33        named pangu-, with"
34    ]
35  }
36 }
```

```

27 },
28 "3": {
29     "title": "parts-of-speech tagger errors do not necessarily degrade
        accuracy in extracting information from biomedical text",
30     "local_link": "papers//0804.0317",
31     "explanation": [
32         "of muscorian, a generic text processing tool for extracting protein-
            protein",
33         "comparable performance to biomedical-specific text processing tools.
            this result was"
34     ]
35 },
36 "4": {
37     "title": "clar: a cross-lingual argument regularizer for semantic role
        labeling",
38     "local_link": "papers//2011.04732v1.pdf",
39     "explanation": [
40         "a given sentence. although different languages have different
            argument annotations,",
41         "share common semantic meaning across languages (e.g. adjuncts have
            more",
42         "such linguistic annotation similarity across languages and exploits
            this information",
43         "information to map the target language arguments using a
            transformation",
44         "the space on which source language arguments lie. by doing",
45         "improves srl performance on multiple languages over monolingual and
            polyglot"
46     ]
47 },
48 "5": {
49     "title": "opinion mining of text documents written in macedonian
        language",
50     "local_link": "papers//1411.4472",
51     "explanation": [
52         "text documents written in macedonian language the ability to extract
            ",
53         "and surveys. in this paper machine learning techniques are used",
54         "which are written in macedonian language. the posts are classified"
55     ]
56 },
57 "6": {
58     "title": "training a restricted boltzmann machine for classification by
        labeling model samples",

```

```

59     "local_link": "papers//1509.01053v1.pdf",
60     "explanation": [
61         "training a restricted boltzmann machine for classification by
           labeling",
62         "classification performance competitive to semi-supervised learning
           if we first train"
63     ]
64 },
65 "7": {
66     "title": "identifying causal influences on publication trends and
           behavior: a case study of the computational linguistics community",
67     "local_link": "papers//2110.07938v1.pdf",
68     "explanation": [
69         "trending tasks and techniques (e.g., deep learning, embeddings,
           generative, and",
70         "china on propensity of researching languages beyond english, and the
           "
71     ]
72 },
73 "8": {
74     "title": "determining whether the non-protein-coding dna sequences are
           in a complex interactive relationship by using an artificial
           intelligence method",
75     "local_link": "papers//1708.04019v1.pdf",
76     "explanation": [
77         "ai research have enabled automatically learning representations of
           high dimensional",
78         "data without feature engineering, using deep neural networks.
           therefore, in",
79         "whether a dna sequence is natural or artificial. the trained"
80     ]
81 },
82 "9": {
83     "title": "deep poetry: a chinese classical poetry generation system",
84     "local_link": "papers//1911.08212v1.pdf",
85     "explanation": [
86         "deep poetry: a chinese classical",
87         "classical poetry generation system called deep poetry. existing
           systems for",
88         "multi-modal input. unlike previous systems, deep poetry uses neural
           networks"
89     ]
90 },
91 "10": {

```

```

92     "title": "attention module is not only a weight: analyzing transformers
          with vector norms",
93     "local_link": "papers//2004.10102",
94     "explanation": [
95         "recently achieved considerable success in natural language
          processing. hence, attention",
96         "bert and a transformer-based neural machine translation system
          include the"
97     ]
98 },
99 "11": {
100     "title": "glass-box program synthesis: a machine learning approach",
101     "local_link": "papers//1709.08669",
102     "explanation": [
103         "glass-box program synthesis: a machine learning approach recently
          proposed",
104         "inspired by number theory, text processing, and algebra. our results
          "
105     ]
106 },
107 "12": {
108     "title": "federated learning for intrusion detection system: concepts,
          challenges and future directions",
109     "local_link": "papers//2106.09527v1.pdf",
110     "explanation": [
111         "federated learning for intrusion detection system:",
112         "and privacy of such devices. machine learning and deep learning",
113         "server. on the contrary, federated learning (fl) fits in
          appropriately",
114         "appropriately as a privacy-preserving decentralized learning
          technique that does not"
115     ]
116 },
117 "13": {
118     "title": "data science with vadalog: bridging machine learning and
          reasoning",
119     "local_link": "papers//1807.08712",
120     "explanation": [
121         "data science with vadalog: bridging machine learning and reasoning
          following",
122         "corporate knowledge and to draw deep insights using machine learning
          ",
123         "data science, typically based on machine learning and statistical
          modelling,"

```



```

124     "significant step forward towards combining machine learning and
        reasoning in"
125 ]
126 },
127 "14": {
128     "title": "object relational graph with teacher-recommended learning for
        video captioning",
129     "local_link": "papers//2002.11566",
130     "explanation": [
131         "object relational graph with teacher-recommended learning for video
            captioning taking",
132         "information from both vision and language is critical for the",
133         "meanwhile, we design a teacher-recommended learning (trl) method to
            make",
134         "use of the successful external language model (elm) to integrate"
135     ]
136 },
137 "15": {
138     "title": "neural feature selection for learning to rank",
139     "local_link": "papers//2102.11345v1.pdf",
140     "explanation": [
141         "neural feature selection for learning to rank learning to",
142         "of information retrieval (ir) where machine learning models are
            employed"
143     ]
144 },
145 "16": {
146     "title": "3d scene parsing via class-wise adaptation",
147     "local_link": "papers//1812.03622",
148     "explanation": []
149 },
150 "17": {
151     "title": "tomayto, tomahto. beyond token-level answer equivalence for
        question answering evaluation",
152     "local_link": "papers//2202.07654v1.pdf",
153     "explanation": []
154 },
155 "18": {
156     "title": "fast domain adaptation for neural machine translation",
157     "local_link": "papers//1612.06897",
158     "explanation": [
159         "fast domain adaptation for neural machine translation neural machine
            translation",
160         "of text from one human language into another. the basic",

```

```

161     "tasks/benchmarks at least for some language pairs. however, many of"
162     "subjective evaluation metric on two language pairs. with our
      adaptation"
163   ]
164 },
165 "19": {
166   "title": "trepan reloaded: a knowledge-driven approach to explaining
      artificial neural networks",
167   "local_link": "papers//1906.08362v2.pdf",
168   "explanation": []
169 }
170 }

```

- distance=cosine, same as dot product results
- distance=euclidean, same as dot product results
- distance=braycurtis Changed from dot-product results,
position in dot product : position in braycurtis
0:0, 1:1, 2:2, 3:3, 4:4, 6:5, 5:6, 8:7, 9:8, 10:9, 14:10, 7:11, 11:12, 12:13, 17:14, 13:15,
16:16, 15:17, 18:18, 19:19
- distance=canberra Changed from dot-product results,
position in dot product : position in canberra
0:0, 3:1, 2:2, 4:3, 1:4, 6:5, 5:6, 14:7, 8:8, 10:9, 9:10, 12:11, 7:12, 16:13, 17:14, 11:15,
15:16, 13:17, 18:18, 19:19
- distance=chebyshev
Changed from dot-product results,
position in dot product : position in chebyshev
0:0, 3:1, 2:2, 19:3, 15:4, 11:5, 18:6, 13:7, 1:8, 12:9, 6:10, 9:11, 7:12, 16:13, 14:14, 5:15,
8:16, 4:17, 17:18, 10:19
- distance=cityblock
Changed from dot-product results,
position in dot product : position in cityblock
0:0, 1:1, 2:2, 3:3, 4:4, 6:5, 5:6, 10:7, 7:8, 9:9, 8:10, 14:11, 12:12, 17:13, 11:14, 16:15,
13:16, 15:17, 18:18, 19:19

In my opinion, I believe that here we have 8 irrelevant papers that have returned as result, because they don't refer to natural language processing. Therefore, we can agree that we have a 60% precision. Furthermore, we can observe that changing the distance metric, can change even more the ranking of the returned results.

References

- [1] Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Michael Kinney, and Daniel S. Weld. S2orc: The semantic scholar open research corpus. In *ACL*, 2020.
- [2] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks, 2019.
- [3] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pre-training for language understanding. *CoRR*, abs/2004.09297, 2020.
- [4] Wen Yang, Tao Li, Gai Fang, and Hong Wei. Pase: Postgresql ultra-high-dimensional approximate nearest neighbor search extension. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, SIGMOD '20, page 2241–2253, New York, NY, USA, 2020. Association for Computing Machinery.