

# Developing a Search Engine and IR Applications using DBLP dataset and PyLucene

Robert Pinsler  
WKWSCI, NTU  
N1509281G@e.ntu.edu.sg

Yiding Liu  
SCE, NTU  
LIUY0130@e.ntu.edu.sg

Yitong Guan  
SCE, NTU  
GUAN0049@e.ntu.edu.sg

Jenn Bing Ong  
SCE, NTU  
ONGJ0063@e.ntu.edu.sg

## ABSTRACT

This study reports the development of a search engine that consists of indexing the dblp computer science bibliography data and process queries on the data fields in terms of keyword and phrase search. The relevance of the query results is evaluated and a user interface for the search engine is built. To showcase possible applications in information retrieval, the most popular research topic as well as the most similar publication venue and year are retrieved.

## CCS Concepts

•Information systems → Document filtering; Information extraction; Relevance assessment;

## Keywords

Search engine; information retrieval; Lucene; DBLP; indexing; query processing; user interface

## 1. INTRODUCTION

A search engine is an Information Retrieval (IR) system designed to help users find information stored on a computer system. In particular, a web-search engine is designed to search for information on the World Wide Web. The web was initially indexed by humans until the 1990s when more and more web servers went online. Human-maintained lists are subjective and expensive to build and maintain, thus cannot scale the web effectively. It was not until a landmark paper by Google [6] that addresses both the issues of search quality and scalability of an automated web-search engine. Google introduced PageRank, used the anchor text, proximity information, etc. to improve search quality dramatically and built a practical and automated web-search engine by fast crawling and query processing technology [6].

The web is a vast collection of uncontrolled heterogeneous documents; web documents have large variations both in

terms of internal and external meta information. Big data refers to datasets that limit the capacity of traditional data processing systems to conduct effective analysis in terms of its data volume, quality and variety in data representation and acquisition speed; hence requiring new processing technologies [5]. All these add to the complexity of IR, which pose enormous challenges but present unprecedented opportunities for new knowledge discovery in real-time that improves decision making.

The objective of this study is to learn and practice the building blocks of a basic search engine using an open-source Application Programming Interface (API). On top of that, two applications from Information Retrieval (IR) are built. The organization of this report is as follows. Section 2 provides information of the different aspects of the search engine built such as indexing, query processing, retrieval evaluation, and user interface. The IR applications are provided in Section 3, followed by the conclusion in Section 4.

## 2. SEARCH ENGINE

In this section, we describe the functionality of our search engine along with illustrative code snippets and evaluate its performance. It is built on top of PyLucene, a Python extension of the open-source software library Lucene<sup>1</sup> that provides text indexing and search capabilities.

The search engine is able to index and search publication records listed in the dblp computer science bibliography (DBLP) [2]. Each publication record contains meta-data about a published document, e.g. title, author and publication year. The actual content of the publication is not provided. All records are combined in a single XML file. The dataset comprises of various kinds of record types, such as articles published in a journal or magazine (*article*), papers published in a conference or workshop proceedings (*inproceedings*), books, and PhD theses. From those, only documents that are classified as *article* or *inproceedings* are considered. To date, this includes over 3.2 million records.

### 2.1 Indexing

Indexing the DBLP dataset is a multi-step process. First, each record has to be extracted from the XML file before it can be further processed. Next, preprocessing is applied to match the expected input format. Finally, the records are processed and indexed through the PyLucene API. Listing 1 gives an overview of the Indexer class that provides above

<sup>1</sup><https://lucene.apache.org>

functionality. It takes as input the file to be indexed, the desired index storage location and an analyzer instance that we will introduce later. By calling its own index function, it triggers the indexing process (Listing 1, line 12). In the following, we explain each step in more detail.

---

```

1 TAGS = ('article', 'inproceedings')
2 VERSION = Version.LUCENE_CURRENT
3 CREATE = IndexWriterConfig.OpenMode.CREATE
4
5 class Indexer():
6     def __init__(self, data_dir, store_dir, analyzer):
7         store = SimpleFSDirectory(File(store_dir))
8         config = IndexWriterConfig(VERSION, analyzer)
9         config.setOpenMode(CREATE)
10        self.writer = IndexWriter(store, config)
11        self.htmlparser = HTMLParser()
12        self.index(data_dir)
13        self.writer.close()
14
15    def index(self, data_dir):
16        context = etree.iterparse(data_dir, tag=TAGS,
17            events=('end',), dtd_validation=True)
18        for event, elem in context:
19            self.index_document(elem)
20
21    def index_document(self, elem):
22        # indexes extracted element

```

---

**Listing 1: Indexer class<sup>2</sup>**

**Parsing** We use an XML parser to extract the DBLP records from the dataset. Due to the large file size, reading in the whole XML tree at once is impractical. Instead, we use a SAX-like parser from the lxml<sup>3</sup> Python package that sequentially reads the document and emits events when it encounters certain XML tags. This allows to only react to `<article>` and `<inproceedings>` tags, thereby ignoring record types we are not interested in. This makes it very fast and memory-efficient. For each emitted event, we call the `index_document` function. (Listing 1, lines 15-19)

---

```

1 FIELDS =
2     ['title', 'author', 'year', 'journal', 'booktitle']
3 HTMLTAGS = ('i', 'ref', 'sub', 'sup', 'tt')
4
5 def index_document(self, elem):
6     etree.strip_tags(elem, HTMLTAGS)
7     doc = Document()
8     f = StringField('id', elem['key'], Store.YES)
9     doc.add(f)
10    for ch in elem:
11        if ch.tag not in FIELDS:
12            continue
13        if ch.text is None:
14            ch = etree.tostring(ch)
15            ch = self.htmlparser.unescape(ch)
16            ch = etree.fromstring(ch)
17
18        if ch.tag == 'title':
19            f = TextField('title', ch.text, Store.YES)
20        elif ch.tag == 'author':
21            f = TextField('authors', ch.text, Store.YES)
22        elif ch.tag == 'year':
23            f = StringField('year', ch.text, Store.YES)
24        else:
25            f = TextField('venue', ch.text, Store.YES)
26
27        doc.add(f)
28        cf = TextField('content', ch.text, Store.NO)
29        doc.add(cf)
30    self.writer.addDocument(doc)

```

---

**Listing 2: index\_document function**

<sup>2</sup>Listings are slightly modified to improve readability.

<sup>3</sup><http://lxml.de>

**Preprocessing** Given the current element, i.e. an *article* or *inproceedings* record, we can now iterate through its children, i.e. record attributes like title or publication year. This is shown in Listing 2. As some records contain HTML formatting, irrelevant tags are removed and HTML-encoded characters are replaced (Listing 2, lines 5 and 12-15).

**Indexing** Each record is indexed as a Lucene *document*, comprising of the following *fields*: id, title, authors, year and venue. The venue field is derived from the journal or booktitle attribute, depending on the record type. To facilitate the search, an additional *content* field is added that concatenates the values of the above fields except the id. This can be easily achieved by subsequently adding all values to the same field. This method is also used to index multiple authors. (Listing 2, lines 19-20 and 27-28).

All fields are treated as strings<sup>4</sup>. Contents of the original fields (i.e. without the *content* field) are stored in the index, which allows to retrieve them later during search. This is useful to enhance search results with extra details when presented to the user. The values of the fields *id* and *year* are stored as-is, which is a property of the `StringField` class. All other fields are tokenized ( $f_{\text{token}}$ ) and further processed.

The processing applied to each of those `TextFields` is determined by the analyzer that is used, which is known by the index writer (Listing 1, lines 6-10 and Listing 2, line 29). By default, all `TextFields` are processed in the same way, which may include lowercase conversion ( $f_{\text{low}}$ ), removal of stopwords ( $f_{\text{stop}}$ ) and stemming ( $f_{\text{stem}}$ ). In particular, we use the Porter stemmer and the default stopword list of Lucene, which consists of 33 common English words. To assign different analyzers depending on the field, we utilize PyLucene’s `PerFieldAnalyzerWrapper` as shown in Listing 3. In this case, we apply the `StandardAnalyzer` from PyLucene to perform lowercase conversion and stopword removal on each tokenized field except the title, where the `CustomAnalyzer` introduced in Listing 4 is used.

---

```

1 INDEX_DIR = 'index'
2 DATA_DIR = 'data/dblp.xml'
3
4 config = {'lowercase': True, 'stemming': True,
5         'stopwords': True}
6 title_analyzer = CustomAnalyzer(config)
7 per_field = HashMap()
8 per_field.put("title", title_analyzer)
9 default = StandardAnalyzer(VERSION)
10 analyzer =
11     PerFieldAnalyzerWrapper(default, per_field)
12 Indexer(DATA_DIR, INDEX_DIR, analyzer)

```

---

**Listing 3: Indexing with field-dependent analyzers**

**Evaluation** To evaluate the impact of token processing onto the index, different configurations of lowercase conversion, removal of stopwords and/or stemming are applied to the *title* field. For this, we make use of our `CustomAnalyzer` that gives us full control over those configurations by extending the base analyzer of PyLucene. Note that the class could be easily extended to support other stemmers and stopword lists by further parametrizing the config parameter.

---

```

1 class CustomAnalyzer(PythonAnalyzer):
2     def __init__(self, config):
3         self.lowercase = config['lowercase']

```

---

<sup>4</sup>The *year* attribute may also be stored as integer, allowing for features like range searches. However, this would require additional effort during search that we try to avoid.

```

4     self.stemming = config['stemming']
5     self.stopwords = config['stopwords']
6     PythonAnalyzer.__init__(self)
7
8     def createComponents(self, fieldName, reader):
9         src = StandardTokenizer(VERSION, reader)
10        fltr = StandardFilter(VERSION, src)
11        if self.lowercase:
12            fltr = LowerCaseFilter(VERSION, fltr)
13        if self.stemming:
14            fltr = PorterStemFilter(fltr)
15        if self.stopwords:
16            sw = StopAnalyzer.ENGLISH_STOP_WORDS_SET
17            fltr = StopFilter(VERSION, fltr, sw)
18        return self.TokenStreamComponents(src, fltr)

```

**Listing 4: CustomAnalyzer class**

During our experiments, we measure the speed of indexing (including parsing and analyzing the documents) in seconds,  $t_{\text{ind}}$ , as well as the number of terms in the vocabulary,  $|V_{\text{title}}|$ . The results for the most interesting combinations are shown in Table 1. Note that indexing speed slightly varied over different runs (not shown).

**Table 1: Evaluation of token processing on title field**

$f_{\text{low}}$	$f_{\text{stop}}$	$f_{\text{stem}}$	$t_{\text{ind}}$	$ V_{\text{title}} $
			417.97 (100%)	436 354 (100%)
X			430.57 (+3%)	346 286 (−21%)
	X		414.89 (−1%)	436 321 (±0%)
		X	439.75 (+5%)	365 534 (−16%)
X	X	X	417.94 (±0%)	288 710 (−34%)

Applying lowercase conversion or stemming reduces the vocabulary significantly, whereas the removal of stopwords has nearly no impact. This is expected, since the first two techniques are able to map multiple tokens onto the same term whereas the exclusion of stopwords decreases the vocabulary size only by the number of stopwords. By combining multiple processing steps, the vocabulary can be further reduced, leading to a reduction of to up a third of the original size when no further processing is applied. In terms of indexing speed, we found that the additional processing steps are negligible. In particular, there seems to be no positive correlation between indexing speed and the number of processing steps that are applied. One reason for this is that while processing each document is more time-consuming, the cost for indexing a processed document is actually reduced. Therefore, we conclude that, under the given evaluation criteria, it is recommended to apply all three processing techniques due to the reduction in vocabulary size. However, this also has an effect on other properties of the system, e.g. the quality of search results, which might not be desired.

## 2.2 Search

The index can now be used to retrieve a ranked list of documents to a given query, which may be a combination of keywords and phrases. Those search capabilities are bundled in the Searcher class presented in Listing 5. The main search logic is part of the *search* function as outlined in Listing 6, whereas phrase queries are processed using *extract\_pq*.

```

1 class Searcher():
2     def __init__(self, index_dir, analyzer):
3         directory = SimpleFSDirectory(File(index_dir))
4         reader = DirectoryReader.open(directory)
5         self.searcher = IndexSearcher(reader)
6         self.analyzer = analyzer

```

```

7     def search(self, query, adv_query, N=0):
8         # perform search for given query
9
10    def extract_pq(self, q, field, slop=0, boost=5):
11        # extract phrases from given query

```

**Listing 5: Searcher class**

```

1 def search(self, query, adv_query, N):
2     bq = BooleanQuery()
3     if query != '':
4         # handle phrases in standard search
5         pq, query = extract_pq(query, 'content')
6         if pq is not None:
7             bq.add(pq, BooleanClause.Occur.MUST)
8     if query != '':
9         # handle remaining keywords
10        qparser = QueryParser(VERSION, 'content',
11                               self.analyzer)
12        q = qparser.parse(query)
13        bq.add(q, BooleanClause.Occur.SHOULD)
14
15    if adv_query is not None:
16        for field, query in adv_query.iteritems():
17            # handle phrases in advanced search
18            pq, query = extract_pq(query, field)
19            if pq is not None:
20                bq.add(pq, BooleanClause.Occur.MUST)
21            if query != '':
22                # handle remaining keywords
23                qparser = QueryParser(VERSION, field,
24                                       self.analyzer)
25                q = qparser.parse(query)
26                bq.add(q, BooleanClause.Occur.MUST)
27    return self.searcher.search(bq, N).scoreDocs

```

**Listing 6: Query handling**

**Keyword search** Our search allows free text keyword queries on any combination of the attributes title, authors, year and venue. By default, all attributes are considered. This is achieved by utilizing the *content* field of the index. We refer to this as a standard search (Listing 6, lines 8-13). Moreover, it is possible to query a combination of specific fields through an advanced search. Those keywords are directly matched against the respective fields in the index. In this case, documents must contain the provided keywords to be returned in the result set (Listing 6, lines 21-26). Note that keywords are handled after phrase extraction (described below). Standard and advanced search can be combined.

**Phrase queries** Phrase queries are supported using double quotation marks (e.g. “information retrieval”) in both standard and advanced search (Listing 6, lines 3-7 and 17-20) through the *extract\_pq* function outlined in Listing 7. There is no restriction for the number of phrases within a query. Documents containing a particular phrase receive higher weight during scoring. This requires an exact match up to the processing on the respective field (e.g. stemming). If a document does not contain the phrase, it may still be considered when there is a match for at least one of the keywords within the phrase. This considerably increases recall while risking an increase of false positives. All those rules are compactly represented using the query language of Lucene (Listing 7, lines 11-12). Finally, the function returns the constructed sub-query as well as the remaining keywords.

```

1 def extract_pq(self, q, field, slop=0, boost=5):
2     phrases = re.findall(r'"([^\"]*)"', q)
3     if len(phrases) == 0:
4         return None, q
5

```

```

6 q = re.sub(r'"([^\"]*)"','"',q).strip()
7 bq = BooleanQuery()
8 for phrase in phrases:
9     qparser =
10         QueryParser(VERSION, field, self.analyzer)
11     # handle phrase and single keywords
12     pq = qparser.parse('%s "%s"~%d^%.1f' %
13         (phrase, phrase, slop, boost))
14     bq.add(pq, BooleanClause.Occur.MUST)
15 return bq,q

```

**Listing 7: Phrase extraction**

**Query results** The search returns the  $N$  most relevant documents along with their ranks, scores, ids and snippets, where relevance is determined based on Lucene’s internal scoring function.  $N$  is a configurable parameter that can be set in a separate configuration file we provide. In order to show meaningful snippets, we experimented with scraping the abstracts of retrieved results from publisher websites using some heuristics, but faced performance issues and technical restrictions from publishers. Therefore, we only show information from indexed fields. Additionally, we return the search duration.

## 2.3 Evaluation

**Table 2: Confusion matrix**

	Relevant	Nonrelevant
Retrieved	true positive (tp)	false positive (fp)
Not Retrieved	false negative (fn)	true negative (tn)

In order to evaluate the quality of search results, we run 10 queries on our search engine and measure performance in terms of precision, recall and the  $F_1$  score against a labeled test collection. Using Table 2, precision (P) and recall (R) are defined as:

$$P = \frac{tp}{tp + fp} \quad R = \frac{tp}{tp + fn}$$

The  $F_1$  score combines both through their harmonic mean:

$$F_1 = \frac{2PR}{P+R}$$

As we are specifically interested in the performance for the top-10 results, the slightly modified metrics precision@K and recall@K are used, where  $K = 10$  in our case. In contrast to the usual recall metric, it is possible that a perfect system cannot achieve a recall@K value of 1. This happens when there are more than  $K$  relevant documents in the test collection. To account for this, we state the overall number of relevant documents for each test query,  $n_{rel}$ , as well as the precision that is achieved after  $n_{rel}$  documents have been retrieved by the system, which is also known as R-precision.

As the search engine is specifically built to retrieve DBLP records, it is reasonable to use the same<sup>5</sup> dataset for evaluation. However, due to the large collection size, we are not able to manually assign relevance judgments to each document for all queries. A naive alternative would be to use a small random subset of the documents. In many cases, this subset will not contain any relevant documents at all though. Large web companies instead employ an approach called *pooling*, where the top-N results from different search

<sup>5</sup>Note that if we wanted to utilize the results to improve performance, special care has to be taken to avoid overfitting.

engines across all test queries are merged into a single collection, which is then labeled by human judges.

A similar approach can be applied by using the DBLP search API, which accepts standard and advanced queries with both keywords and phrases as defined in Section 2.2. Results are returned in XML format and can be post-filtered to remove irrelevant publication types. One issue is that hits are ranked by a coarse relevance score that does not always allow to extract a distinct top 10 for a given query. If there are hits beyond the top 10 that have the same score as the 10-th result, we include them as well. To improve diversity, we additionally augment the test collection with a disjunct set of random documents from the original dataset. In total, the test collection comprises of 300 documents.

We manually added relevance judgments for every document in the test collection with respect to each query. Following TREC conventions, we derive test queries from topics with specific information needs and use the TREC file format to store topic definitions and relevance judgments. The results for each query are shown in Table 3.

As can be seen, precision@10 is constantly very high except for the third query due to the low absolute number of relevant documents (only two out of the three received documents are relevant). In contrast to that, recall@10—and therefore also the  $F_1$  score—significantly drops when  $n_{rel} > 10$ , which is expected. In those cases, R-precision is much more reliable. Overall, the evaluation suggests that the results of our search engine are both precise and complete. However, one issue we encountered during evaluation is that it is very hard, for both human and machine, to infer whether a publication record truly matches the underlying information need, given only the metadata. Therefore, very often the only viable option is to perform simple term matching. This should be taken into consideration when evaluating the results.

## 2.4 User Interface

There are two ways to interact with the search engine: via a text-based command line interface or via a web UI. The latter is built with Flask<sup>6</sup>, a lightweight Python web framework that leverages Jinja2<sup>7</sup> as a templating engine. This allows to easily create HTML documents from within Python. Additionally, we incorporate the web framework Bootstrap<sup>8</sup> to achieve responsive web design. Figure 1 depicts the search results for an example query. As the UI is not our focus, we omit further implementation details; the interested reader is referred to the source code.

## 3. APPLICATIONS IN IR

We implement two applications based on the DBLP dataset using the nltk<sup>9</sup> Python package. First, we are interested in finding the top-10 most popular research topics of a specific year as shown in Listing 8. The second application retrieves the top-10 most similar publication venues to a given venue and year. Only paper titles are used in both applications.

In this section, related techniques and implementation details of the applications are demonstrated, followed by the experimental results.

<sup>6</sup><http://flask.pocoo.org>

<sup>7</sup><http://jinja.pocoo.org>

<sup>8</sup><http://getbootstrap.com>

<sup>9</sup><http://www.nltk.org>

**Table 3: Evaluation of search results**

Query	$n_{rel}$	Precision@10	Recall@10	$F_1@10$	R-precision
title:user behavior	89	1.00	0.11	0.20	0.93
authors:christopher manning venue:nips	10	1.00	1.00	1.00	1.00
query optimisation 2013	2	0.67	1.00	0.80	0.50
title:"information retrieval"	41	1.00	0.24	0.39	0.98
auto-completion sigir	7	1.00	1.00	1.00	1.00
trec authors:chris buckley	13	1.00	0.77	0.87	1.00
year:2010 title:sentiment analysis	47	1.00	0.21	0.35	1.00
system evaluation	30	0.90	0.30	0.45	0.80
"recommender systems" "data mining"	5	0.83	1.00	0.91	0.80
deep learning venue:icml year:2015	12	1.00	0.67	0.80	1.00

dblp search

"query optimization"



1257 results found (0.28s).

### 1 Dynamic Query Optimization and Query Processing in Multidatabase Systems

Henryk Josinski - 2000 - EDBT PhD Workshop

Relevance score: 3.58

### 2 Object Query Optimization through Detecting Independent Subqueries

Jacek Plodzien, Anna Kraken - 2000 - Inf. Syst.

Relevance score: 3.51

**Figure 1: Web-based search user interface: results for example query "query optimization" (restricted to publications from 2000, advanced query not shown). Only top-2 results are depicted due to space constraints.**

## 3.1 Popular research topics

Before mining popular research topics from paper titles, we need to firstly extract topics from the data. In our application, we identify meaningful topics from titles based on the patterns of words. More precisely, given a query (e.g. "2014"), we have the following process for topic extraction:

1. Retrieve titles of all papers published in 2014 using the index (Listing 8, line 2).
2. Tokenize each title, remove punctuation and stop words (Listing 8, lines 11-13).
3. Use a part-of-speech tagger<sup>10</sup> to assign parts of speech to each word, such as noun, verb, adjective, etc. (Listing 8, line 14).
4. Extract topics from the tagged words based on a given pattern, using regular expression tools (Listing 8, lines 9-10 and 15).

The key point of topic extraction is to define the word pattern of a topic. From the observation of real topic names, we can conclude that:

- A topic name usually consists of adjectives and nouns.
- Adjectives always come before nouns for a topic name.

<sup>10</sup><http://nlp.stanford.edu/software/tagger.shtml>

- Normally, a topic consists of less than four words.

Therefore, in the last step, we define the pattern for topics as `adj+adj+...+noun+...+noun`. Moreover, only unigrams, bigrams and trigrams are considered in this application. We ignore those unigram terms with very high (top 1%) frequency because they may be some widely-used terminology in many different topics. After the topic identification, we count the frequency of each topic that appeared in the titles. Finally, the top-10 most frequent topics will be returned (Listing 8, lines 26-34).

```

1 def get_popular_topics(self, q_year, top_k):
2     titles = self.searcher.search_year(q_year)
3     unigram_dist = {}
4     bigram_dist = {}
5     trigram_dist = {}
6
7     tagset = None
8     tagger = PerceptronTagger()
9     grammar = "NP:
10    {<JJ>*(<NN>|<NNS>)*<NN>(<NN>|<NNS>)*}"
11     cp = nltk.RegexpParser(grammar)
12     for title in titles:
13         title = title.lower()
14         text = word_tokenize(title)
15         sentence =
16             nltk.tag._pos_tag(text, tagset, tagger)
17         result = cp.parse(sentence)
18         for node in list(result):
19             if isinstance(node, nltk.tree.Tree):
20                 entity = zip(*list(node))[0]
21                 if len(entity) == 1:

```



```

20     self.dict_append(entity, unigram_dist)
21     elif len(entity) == 2:
22         self.dict_append(entity, bigram_dist)
23     elif len(entity) == 3:
24         self.dict_append(entity, trigram_dist)
25
26     top_lp = int(len(unigram_dist) * 0.01)
27     ug = Counter(unigram_dist)
28     ug = ug.most_common(top_k + top_lp)[top_lp:]
29     bg = Counter(bigram_dist).most_common(top_k)
30     tg = Counter(trigram_dist).most_common(top_k)
31
32     result = ug + bg + tg
33     return sorted(result, key=lambda k: k[1],
34                   reverse=True)[:top_k]

```

Listing 8: Finding popular topics for a given year

### 3.2 Similar publication venues

For retrieving similar publication venues and years, we additionally make use of the `lda`<sup>11</sup>, `numpy`<sup>12</sup> and `scikit-learn`<sup>13</sup> packages. Firstly we need to measure the similarity between two <venue, year> pairs (e.g. how similar are "TKDE'14" and "CIKM'15") using the paper titles. One simple way is to firstly tokenize every title of each <venue, year> pair, use a bag-of-words model and finally compute the cosine similarity of the word frequency vectors between them, which is defined as:

$$Sim_{i,j} = \frac{\sum_{k \in W} freq_{i,k} \cdot freq_{j,k}}{\sqrt{\sum_{k \in W} freq_{i,k}^2} \sqrt{\sum_{k \in W} freq_{j,k}^2}},$$

where  $i, j$  are two <venue, year> pairs.  $W$  refers to the collection containing all the words which appeared in the corpus.  $freq_{i,k}$  is the frequency of word  $k$  that appeared in the  $i$ -th venue and year.

However, based on our statistical analysis, the dimensionality is too high, which makes it very inefficient to compute similarity. Besides, the data may also be very noisy, containing many words that only appear once in the whole dataset. To address the aforementioned problems, we adopted Latent Dirichlet allocation (LDA)[1] to learn latent topics from the data. Instead of using word frequency vectors, the similarity is defined based on the cosine similarity between the topic distribution of two <venue, year> pairs. Thus, for a given <venue, year> query, we compute topic similarities between it and every other <venue, year> pair, which is given by:

$$Sim_{i,j} = \frac{\sum_{z \in Z} \theta_{i,z} \cdot \theta_{j,z}}{\sqrt{\sum_{z \in Z} \theta_{i,z}^2} \sqrt{\sum_{z \in Z} \theta_{j,z}^2}},$$

where  $Z$  is the topic set and  $\theta_{i,z}$  is the probability of the  $i$ -th <venue, year> pair belonging to topic  $z$ . Then, based on the similarity, the top-10 most similar venues are returned.

### 3.3 Application Experiments

For retrieving popular research topics, we have the results as in Table 4. As we can see, most of the phrases are meaningful topics such as "wireless sensor networks", "cloud computing" and "neural network". Thus, we validate the effectiveness of our proposed method.

<sup>11</sup><https://pythonhosted.org/lda>

<sup>12</sup><http://www.numpy.org/>

<sup>13</sup><http://scikit-learn.org>

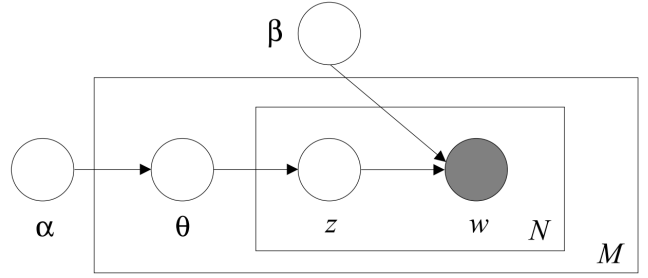


Figure 2: Graph representation of LDA

Table 4: Top-10 most popular research topics of papers published in 2013

Rank	Topic	Frequency
1	wireless sensor networks	1412
2	case study	1182
3	performance analysis	568
4	cognitive radio networks	468
5	special issue	452
6	performance evaluation	404
7	cloud computing	373
8	genetic algorithm	341
9	empirical study	340
10	neural network	340

In order to find similar publication venues and years, we firstly learn the topics as shown in Table 5. These can then be used to find similar venues and years. An example is shown in Table 6, where the most similar venues and years for the query "IEEE Trans. Knowl. Data Eng. 2014" are listed. According to that, "IEEE Trans. Knowl. Data Eng. 2015" and "IEEE Trans. Knowl. Data Eng. 2012" are the most similar venues and years to it, which is quite intuitive, because they are the same journal with different years. Moreover, we also have "CIKM" and "WAIM" included since they are about data engineering as well. Besides, we can observe that the year of the top-ranked ones are more close to 2013, which also suggests the correctness of our results. This shows that by adopting the LDA model combined with cosine similarity, we are able to efficiently and effectively mine similar <venue, year> pairs.

## 4. CONCLUSIONS

In this study, a search engine and two example IR applications have been built using PyLucene and the DBLP dataset. Various aspects of building a search engine such as indexing, query processing, retrieval evaluation and user interface are covered. Moreover, the top-10 most popular topics for each year and most similar <venue, year> pairs are retrieved. The key computational functions, code snippets and algorithms are included to elucidate the process. Evaluations of both the search engine and the applications have underlined the correctness of our approach.

## 5. ACKNOWLEDGEMENTS

The authors would like to thank Assoc. Prof. Aixin Sun for the comprehensive and organized teaching of the various topics in information retrieval and analysis.

**Table 5: Representative words for topics learnt by the LDA algorithm**

Topic	Words
Geographic Information System	data using based sar spatial sensing analysis remote images radar model land image satellite surface detection gis classification resolution hyperspectral
Robotics	robot control based using robots mobile motion robotic planning human autonomous design navigation vision dynamic sensor multi force tracking visual
Wireless Sensor Networks	networks wireless based network mobile routing ad hoc sensor protocol performance control using traffic efficient scheme qos algorithm 802 multi
Software Technology	software testing analysis based using test engineering development study code case model systems java quality source empirical tool approach program
Social Media	social online media networks community network behavior use internet effects study analysis communities understanding games self role game impact influence

**Table 6: Top-10 most similar venues & years to “IEEE Trans. Knowl. Data Eng. 2014”**

Rank	Venue	Year	Similarity
1	IEEE Trans. Knowl. Data Eng.	2015	0.9898
2	IEEE Trans. Knowl. Data Eng.	2012	0.9888
3	DaWaK	2015	0.9887
4	IEEE Trans. Knowl. Data Eng.	2013	0.9856
5	LD4IE@ISWC	2014	0.9850
6	CIKM	2006	0.9845
7	WAIM	2013	0.9830
8	CIKM	2007	0.9823
9	FQAS	2009	0.9815
10	CIKM	2004	0.9803

## 6. REFERENCES

- [1] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003.
- [2] M. Ley et al. dblp computer science bibliography dataset. <http://dblp.uni-trier.de/xml>, 2016. last accessed: 18-March-2016.
- [3] K. W. T. Leung, D. L. Lee and W. C. Lee. PMSE: A personalized mobile search engine. *IEEE Transactions on Knowledge and Data Engineering*, 4:820–834, 2013.
- [4] M. Moran and B. Hunt. Search Engine Marketing, inc.: Driving Search Traffic to Your Company’s Web Site. 2015.
- [5] M. Chen, S. Mao and Y. Liu. Big data: A survey *Mobile Networks and Applications*, 2:171–209, 2014.
- [6] S. Brin and L. Page. Reprint of: The anatomy of a large-scale hypertextual web search engine. *Computer Networks*, 18:3825–3833, 2012.