

Keyword Search Engine (mini-IR)

Zbigniew Dendzik
Institute of Physics, University of Silesia

2025

Introduction

The goal of this exercise set is to introduce a simple but important class of applications of artificial intelligence: information retrieval systems (IR). The task of such a system is to take a textual query (e.g. one or several keywords) and return a list of documents that best match this query.

In practice, similar ideas underlie web search engines, searching in technical documentation, email or notes. Modern production systems use very complex language models (large neural networks) that can understand context and semantic relations between words. In this set, however, we focus on the simplest level: searching a collection of documents using keywords, counting word frequencies and implementing a simple document ranking.

Our “mini search engine” will consist of:

- a class representing a single text document,
- an inverted index: for each word a list of documents where it appears,
- simple one-word and two-word queries and a relevance-based ranking.

Such a simple system is a good starting point for understanding what more advanced IR systems and language models do later on: instead of raw keyword matching they use vector embeddings and semantic similarity.

Example I

Below you will find a skeleton definition of the `Document` class representing a single text document. Each document has an identifier, a title and a body and provides methods for simple text processing.

```
1 class Document
2 {
3     private int id;
4     private String title;
5     private String body;
6
7     public Document(int id, String title, String body)
8     {
9         this.id=id;
```

```

10     this.title=title;
11     this.body=body;
12 }
13
14 public int getId()
15 {
16     return id;
17 }
18
19 public String getTitle()
20 {
21     return title;
22 }
23
24 public String getBody()
25 {
26     return body;
27 }
28
29 //returns an array of "words" after simple normalisation
30 public String[] getWords()
31 {
32     String txt=body.toLowerCase();
33     txt=txt.replaceAll("[^a-z0-9 ]"," ");
34     String[] words=txt.split("\\s+");
35     return words;
36 }
37
38 //returns frequency of a given word in the document
39 public int wordFrequency(String word)
40 {
41     String[] words=getWords();
42     int count=0;
43     for(String s:words)
44     {
45         if(s.equals(word.toLowerCase()))
46             count++;
47     }
48     return count;
49 }
50
51 public java.util.Set<String> wordSet()
52 {
53     java.util.Set<String> set=
54         new java.util.HashSet<String>();
55     for(String s:getWords())
56     {
57         if(s.length()>0)
58             set.add(s);
59     }
60     return set;

```

```

61     }
62
63     public String toString()
64     {
65         return "Document["+id+"]: "+title;
66     }
67 }

```

Listing 1: Document.java

Ex. 1.1

Test the Document class. Write a program that:

- creates several Document objects with different contents,
- prints the id, title and number of words for each document,
- prints the frequency of a few selected words for one chosen document.

```

1 public class TestDocument
2 {
3     public static void main(String[] args)
4     {
5         Document d1=new Document(1,"Movies",
6             "This is a document about movies and series.");
7         Document d2=new Document(2,"Sports",
8             "This document is about football and volleyball.");
9
10        System.out.println(d1);
11        System.out.println("Word count: "+
12            d1.getWords().length);
13        System.out.println("freq 'document': "+
14            d1.wordFrequency("document"));
15
16        System.out.println(d2);
17        System.out.println("Word set: "+
18            d2.wordSet());
19    }
20 }

```

Listing 2: TestDocument.java – skeleton

Example II

The next step is to build a simple inverted index, i.e. a data structure storing information about which documents contain a given word.

```

1 import java.util.*;
2
3 class Index

```

```

4 {
5     private List<Document> documents;
6     private Map<String, List<Integer>> index;
7
8     public Index()
9     {
10         documents=new ArrayList<Document>();
11         index=new HashMap<String, List<Integer>>();
12     }
13
14     public void addDocument(Document d)
15     {
16         documents.add(d);
17         for(String word : d.wordSet())
18         {
19             List<Integer> list=index.get(word);
20             if(list==null)
21             {
22                 list=new ArrayList<Integer>();
23                 index.put(word,list);
24             }
25             if(!list.contains(d.getId()))
26                 list.add(d.getId());
27         }
28     }
29
30     public List<Document> searchWord(String word)
31     {
32         List<Document> result=new ArrayList<Document>();
33         List<Integer> idList=index.get(word.toLowerCase());
34         if(idList!=null)
35         {
36             for(Integer id : idList)
37             {
38                 Document d=findDocumentById(id);
39                 if(d!=null)
40                     result.add(d);
41             }
42         }
43         return result;
44     }
45
46     private Document findDocumentById(int id)
47     {
48         for(Document d : documents)
49         {
50             if(d.getId()==id)
51                 return d;
52         }
53         return null;
54     }

```

Listing 3: Index.java – skeleton

Ex. 1.2

Write a test program TestIndex that:

- creates an Index object,
- adds a dozen documents to it,
- asks the user for a keyword (from keyboard) and prints the titles of documents containing this word.

```
1 import java.util.*;
2
3 public class TestIndex
4 {
5     public static void main(String[] args)
6     {
7         Index index=new Index();
8
9         Document d1=new Document(1,"Movies",
10             "This is a document about movies and series.");
11         Document d2=new Document(2,"Sports",
12             "This document is about football and volleyball.");
13         Document d3=new Document(3,"Music",
14             "A document about film music and concerts.");
15
16         index.addDocument(d1);
17         index.addDocument(d2);
18         index.addDocument(d3);
19
20         Scanner in=new Scanner(System.in);
21         System.out.print("Keyword: ");
22         String word=in.nextLine().toLowerCase();
23
24         List<Document> results=index.searchWord(word);
25         System.out.println("Found documents:");
26         for(Document d : results)
27         {
28             System.out.println(" - "+d.getTitle());
29         }
30     }
31 }
```

Listing 4: TestIndex.java – skeleton

Example III

Now we extend the search engine with multi-word queries and simple relevance ranking.

```
1 import java.util.*;
2
3 class SearchEngine
4 {
5     private Index index;
6     private Set<String> stopWords;
7
8     public SearchEngine(Index index)
9     {
10        this.index=index;
11        stopWords=new HashSet<String>();
12        stopWords.add("and");
13        stopWords.add("the");
14        stopWords.add("or");
15        stopWords.add("a");
16    }
17
18    public List<Document> query(String word)
19    {
20        word=normaliseWord(word);
21        if(stopWords.contains(word))
22            return new ArrayList<Document>();
23        return index.searchWord(word);
24    }
25
26    public List<Document> queryAND(String w1, String w2)
27    {
28        w1=normaliseWord(w1);
29        w2=normaliseWord(w2);
30
31        List<Document> l1=index.searchWord(w1);
32        List<Document> l2=index.searchWord(w2);
33
34        List<Document> result=new ArrayList<Document>();
35        for(Document d1 : l1)
36        {
37            if(l2.contains(d1))
38                result.add(d1);
39        }
40        return result;
41    }
42
43    public List<Document> queryOR(String w1, String w2)
44    {
45        w1=normaliseWord(w1);
46        w2=normaliseWord(w2);
47
48        List<Document> l1=index.searchWord(w1);
```

```

49     List<Document> l2=index.searchWord(w2);
50
51     List<Document> result=new ArrayList<Document>(l1);
52     for(Document d2 : l2)
53     {
54         if(!result.contains(d2))
55             result.add(d2);
56     }
57     return result;
58 }
59
60 public List<Document> sortByRelevance(
61     List<Document> list, String word)
62 {
63     word=normaliseWord(word);
64     Collections.sort(list,
65         new Comparator<Document>()
66         {
67             public int compare(Document d1,
68                 Document d2)
69             {
70                 int f1=d1.wordFrequency(word);
71                 int f2=d2.wordFrequency(word);
72                 return Integer.compare(f2,f1);
73             }
74         });
75     return list;
76 }
77
78 private String normaliseWord(String s)
79 {
80     s=s.toLowerCase();
81     s=s.replaceAll("[^a-z0-9]", "");
82     return s;
83 }
84 }

```

Listing 5: SearchEngine.java – skeleton

Ex. 1.3

Write a test program TestSearchEngine that:

- creates an Index and a SearchEngine,
- adds several documents to the index,
- lets the user choose the query type (one-word, AND, OR),
- prints the list of documents, and for a one-word query: documents sorted by relevance.

```

1 import java.util.*;
2
3 public class TestSearchEngine
4 {
5     public static void main(String[] args)
6     {
7         Index index=new Index();
8         //TODO: add documents...
9
10        SearchEngine se=new SearchEngine(index);
11
12        Scanner in=new Scanner(System.in);
13        System.out.println("1) one-word query");
14        System.out.println("2) AND query");
15        System.out.println("3) OR query");
16        int choice=in.nextInt();
17        in.nextLine();
18
19        if(choice==1)
20        {
21            System.out.print("Word: ");
22            String w=in.nextLine();
23            List<Document> res=se.query(w);
24            res=se.sortByRelevance(res,w);
25            for(Document d : res)
26                System.out.println(d.getTitle());
27        }
28        else if(choice==2)
29        {
30            System.out.print("Word 1: ");
31            String w1=in.nextLine();
32            System.out.print("Word 2: ");
33            String w2=in.nextLine();
34            List<Document> res=se.queryAND(w1,w2);
35            for(Document d : res)
36                System.out.println(d.getTitle());
37        }
38        else if(choice==3)
39        {
40            System.out.print("Word 1: ");
41            String w1=in.nextLine();
42            System.out.print("Word 2: ");
43            String w2=in.nextLine();
44            List<Document> res=se.queryOR(w1,w2);
45            for(Document d : res)
46                System.out.println(d.getTitle());
47        }
48    }
49 }

```

Listing 6: TestSearchEngine.java – skeleton

Example IV (extension)

Modern IR systems and large language models (LLMs) use much more advanced techniques: vector embeddings, semantic similarity, query expansion etc. The simple keyword search engine you build here is a miniature version of the first layer of such systems.

Ex. 1.4*

Extend the search engine with a simple TF-IDF weight:

- compute **TF** (Term Frequency) as the frequency of the word in the document,
- compute **IDF** (Inverse Document Frequency) as $\log(\frac{N}{df})$, where N is the number of documents and df is the number of documents containing the word,
- define the word weight in a document as $TF \cdot IDF$ and use it in document ranking.