



UNIVERSITÉ D'AVIGNON
ET DES PAYS DE VAUCLUSE

M2 ILSSEN – 2019/20

UE Ingénierie du document et de l'information

UCE Indexation & recherche

Vincent Labatut

TP 1 | Correction

1 Tokénisation basique

Exercice 1

```
1 public int compareTo(Token token)
2 { int result = type.compareTo(token.type);
3   if(result==0)
4     result = docId - token.docId;
5   return result;
6 }
```

Exercice 2

```
1 public boolean equals(Object o)
2 { boolean result = false;
3   if(o instanceof Token)
4   { Token token = (Token) o;
5     result = compareTo(token) == 0;
6   }
7   return result;
8 }
9
10 public String toString()
11 { String result = "(" + type + ", " + docId + ")";
12   return result;
13 }
```

Exercice 3

```
1 public List<String> tokenizeString(String string)
2 { List<String> result = new LinkedList<String>();
3
4   // on segmente la chaîne
5   String temp[] = string.split("[^\\pL\\pN]");
6   // on traite chaque segment
7   for(String tmp: temp)
8   { // si le segment traité n'est pas vide
9     if(!tmp.isEmpty())
10      // on le rajoute à la liste résultat
11      result.add(tmp);
12   }
13
14   return result;
15 }
```

Exercice 4

```
1 public void tokenizeDocument(File document, int docId, List<Token> tokens) throws
2   UnsupportedEncodingException, FileNotFoundException
3 { // on ouvre le flux
4   FileInputStream fis = new FileInputStream(document);
5   InputStreamReader isr = new InputStreamReader(fis, "UTF-8");
```

```

5     Scanner scanner = new Scanner(isr);
6
7     // on traite chaque ligne une par une
8     while(scanner.hasNextLine())
9     { String line = scanner.nextLine();
10      List<String> types = tokenizeString(line);
11      for(String type: types)
12      { // on crée le token approprié
13        Token token = new Token(type, docId);
14        tokens.add(token);
15      }
16    }
17
18    // on ferme le flux
19    scanner.close();
20 }

```

Exercice 5

```

1     public int tokenizeCorpus(List<Token> tokens) throws UnsupportedOperationException,
2         FileNotFoundException
3     { // on récupère les fichiers contenus dans le dossier
4       String folder = FileTools.getCorpusFolder();
5       File fFolder = new File(folder);
6       String fileNames[] = fFolder.list();
7       // note : le tri d'objets File peut varier d'un système à l'autre
8       // mais pas le tri de String.
9       Arrays.sort(fileNames);
10
11      // nombre de documents dans la collection
12      int result = fileNames.length;
13
14      // on les traite un par un
15      int docId = 0;
16      for(String fileName: fileNames)
17      { File file = new File(folder+File.separator+fileName);
18        tokenizeDocument(file, docId, tokens);
19        docId++;
20      }
21
22      return result;
23 }

```

2 Normalisation basique

Exercice 6

```

1     public String normalizeType(String string)
2     { String result = string;
3
4       // on passe en minuscules
5       result = result.toLowerCase();
6
7       // on supprime les signes diacritiques
8       result = java.text.Normalizer.normalize(
9         result, Form.NFD).replaceAll("\\p{InCombiningDiacriticalMarks}+", "");
10
11      // on renvoie null en cas de chaîne vide
12      if(result.isEmpty())
13        result = null;

```

```

14
15     return result;
16 }

```

Exercise 7

```

1  public void normalizeTokens(List<Token> tokens)
2  { // on traite chaque token séparément
3      Iterator<Token> it = tokens.iterator();
4      while(it.hasNext())
5      { // on récupère le token
6          Token token = it.next();
7          // on normalise le type
8          String type = token.getType();
9          String term = normalizeType(type);
10         // si aucun terme n'est renvoyé, on supprime le token
11         if(term==null)
12             it.remove();
13         else
14             token.setType(term);
15     }
16 }

```

3 Fichier inverse simple

Exercise 8

```

1  public int compareTo(Posting posting)
2  { int result = docId - posting.docId;
3      return result;
4  }
5
6  public String toString()
7  { String result = Integer.toString(docId);
8      return result;
9  }
10
11 public boolean equals(Object o)
12 { boolean result = false;
13     if(o instanceof Posting)
14     { Posting posting = (Posting) o;
15         result = compareTo(posting) == 0;
16     }
17     return result;
18 }

```

Exercise 9

```

1  public IndexEntry(String term)
2  { this.term = term;
3      this.postings = new ArrayList<Posting>();
4      this.frequency = 0;
5  }
6
7  public int compareTo(IndexEntry entry)
8  { int result = term.compareTo(entry.term);
9      return result;
10 }
11
12 public String toString()
13 { String result = "<" + term + " [" + frequency + "]" + " (";

```

```

14     for(Posting posting: postings)
15         result = result + " " + posting;
16     result = result + " >";
17     return result;
18 }
19
20 public boolean equals(Object o)
21 { boolean result = false;
22   if(o instanceof IndexEntry)
23   { IndexEntry entry = (IndexEntry) o;
24     result = compareTo(entry) == 0;
25   }
26   return result;
27 }

```

Exercise 10

Constructeur de la classe-fille ArrayIndex :

```

1 public class ArrayIndex extends AbstractIndex
2 { public ArrayIndex(int size)
3   { data = new IndexEntry[size];
4   }
5 }

```

Constructeur de la classe-fille HashIndex :

```

1 public class HashIndex extends AbstractIndex
2 { public HashIndex(int size)
3   { data = new HashMap<String,IndexEntry>(size);
4   }
5 }

```

Constructeur de la classe-fille TreeIndex :

```

1 public class TreeIndex extends AbstractIndex
2 { public TreeIndex()
3   { data = new TreeMap<String,IndexEntry>();
4   }
5 }

```

Exercise 11

Dans la classe-fille ArrayIndex :

```

1 public void print()
2 { for(int i=0;i<data.length;i++)
3   System.out.println(i+"\t"+data[i]);
4 }

```

Dans les classes-filles HashIndex et TreeIndex :

```

1 public void print()
2 { int i = 0;
3   for(IndexEntry indexEntry: data.values())
4   { System.out.println(i+"\t"+indexEntry);
5     i++;
6   }
7 }

```

Exercise 12

Dans la classe-fille ArrayIndex :

```

1 public void addEntry(IndexEntry indexEntry, int rank)
2 { data[rank] = indexEntry;
3 }

```

Dans les classes-filles HashIndex et TreeIndex :

```

1 public void addEntry(IndexEntry indexEntry, int rank)
2 { String key = indexEntry.getTerm();
3   data.put(key, indexEntry);
4 }

```

Exercice 13

Dans la classe-fille ArrayIndex, une version avec la recherche dichotomique implémentée à la main :

```

1 public IndexEntry getEntry(String term)
2 { IndexEntry result = null;
3
4   if(data.length>0)
5   { // on effectue une recherche dichotomique classique
6     int s=0,e=data.length-1;
7     do
8     { int m = (s+e)/2;
9       String t = data[m].term;
10      int comp = term.compareTo(t);
11      // term == t
12      if(comp==0)
13        result = data[m];
14      // term > t
15      else if(comp>0)
16        s = m + 1;
17      // term < t
18      else if(comp<0)
19        e = m - 1;
20    }
21    while(result==null && s<=e);
22  }
23
24  return result;
25 }

```

Alternativement, une version utilisant l'implémentation de Arrays :

```

1 public IndexEntry getEntry(String term)
2 { IndexEntry result = null;
3
4   // on a besoin d'un objet IndexEntry bidon
5   IndexEntry entry = new IndexEntry(term);
6   // on applique la méthode du JDK
7   int pos = Arrays.binarySearch(data, entry);
8   // on récupère l'objet adéquat
9   if(pos>=0)
10    result = data[pos];
11   return result;
12 }

```

Dans les classes-filles HashIndex et TreeIndex :

```

1 public IndexEntry getEntry(String term)
2 { IndexEntry result = data.get(term);
3   return result;
4 }

```

Exercice 14

Dans la classe-fille ArrayIndex :

```

1 public int getSize()

```

```
2 { int result = data.length;
3   return result;
4 }
```

Dans les classes-filles `HashIndex` et `TreeIndex` :

```
1 public int getSize()
2 { int result = data.size();
3   return result;
4 }
```