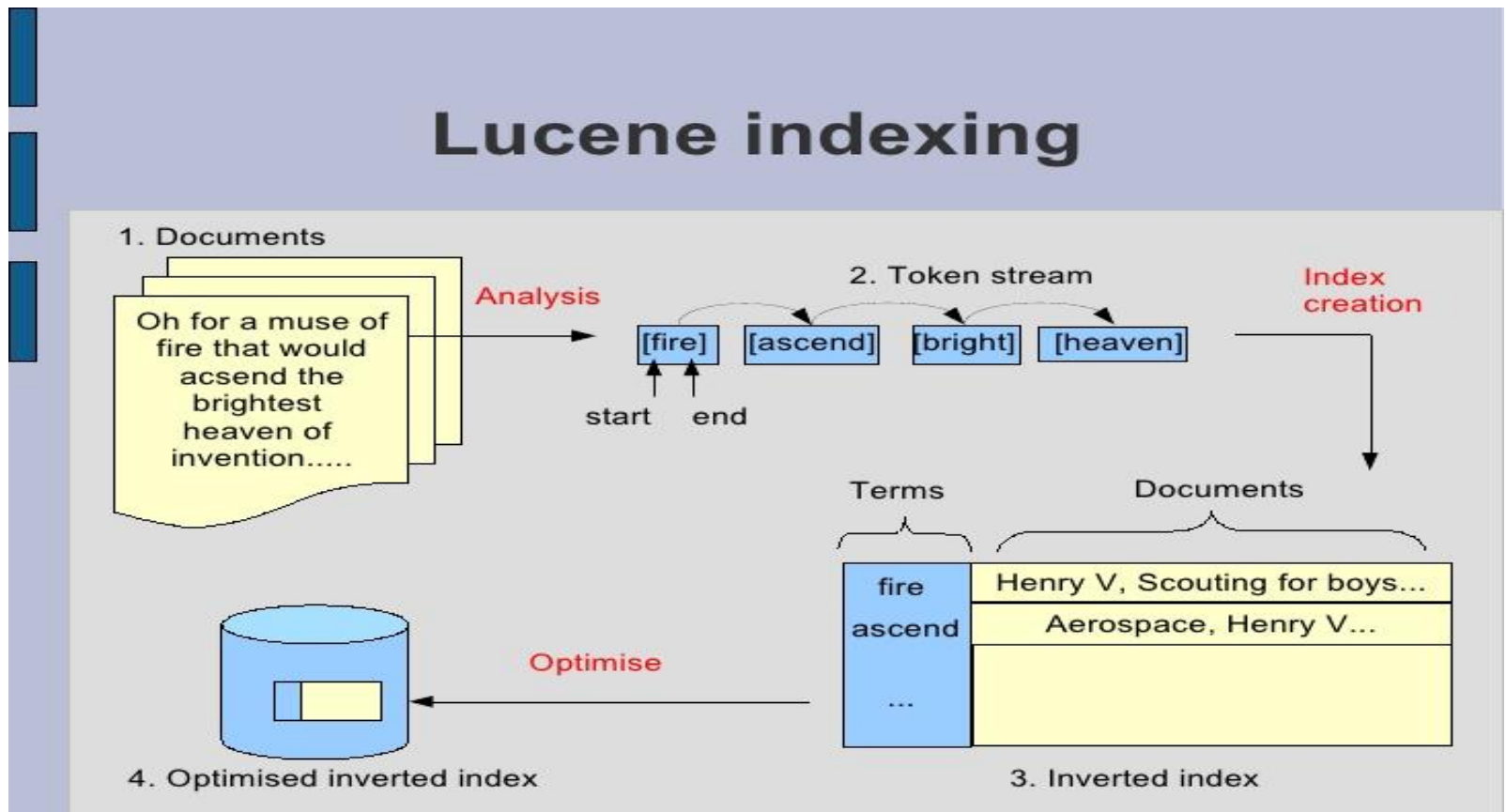


Analyse et traitement du Contenu

Indexation des Documents avec Lucene

- Indexation – construction de l'index – classe IndexWriter



Indexation des Documents avec Lucene

- `Analyzer luceneAnalyzer = new StandardAnalyzer();`
- `IndexWriter indexWriter = new IndexWriter(indexDir, luceneAnalyzer);`
- cette classe peut créer un nouvel index ou ouvrir un index existant et lui ajouter des documents.
- le premier paramètre : ou on stocke les fichiers d'index ;
- le deuxième paramètre : l'analyseur qui sera employé ;
- le dernier paramètre : si vraie, la classe crée un nouvel index ; si faux, il ouvre un index existant

Analyse des Documents avec Lucene

- Analyzer luceneAnalyzer = new **StandardAnalyzer**();

Variantes :

- **WhitespaceAnalyzer** Un analyseur très simple qui sépare juste la marque en
- utilisant l'espace blanc.
- **StopAnalyzer** Enlève les mots anglais communs et de liaison, inutiles pour l'indexation.
- **SnowballAnalyzer** Un analyseur expérimental intéressant qui travaille sur le racines (recherche sur la raining renvoie aussi rained, rain)
- **FrenchAnalyzer** Un analyseur pour le Français

Analyse des Documents avec Lucene

Apache Lucene, ou tout simplement Lucene est un logiciel de recherche et l'indexation des documents API, écrite dans le langage de programmation JAVA. Ce est un logiciel open source de l'Apache Software Foundation sous la Licence Apache.



apache lucene ou tout simplement lucene est un logiciel de recherche et l'indexation des documents api écrite dans le langage de programmation java ce est un logiciel open source de l apache software foundation sous la licence apache



apache lucene ou tout simplement lucene est un logiciel de recherche et l'indexation des documents api écrite dans le langage de programmation java ce est un logiciel open source de l apache software foundation sous a licence apache

Analyse des Documents avec Lucene

apache lucene **ou tout** simplement lucene **est un** logiciel **de** recherche **et l**
indexation **des** documents api écrite **dans le** langage **de** programmation java **ce est**
un logiciel open source **de l** apache software foundation **sous la** licence apache

```
graph TD; A["apache lucene ou tout simplement lucene est un logiciel de recherche et l indexation des documents api écrite dans le langage de programmation java ce est un logiciel open source de l apache software foundation sous la licence apache"] --> B((Filtration)); B --> C["apache lucene simplement lucene logiciel recherche indexation documents api écrite langage programmation java logiciel open source apache software foundation licence apache"]
```

Filtration

apache lucene simplement lucene logiciel recherche indexation documents
api écrite langage programmation java logiciel open source apache software
foundation licence apache

Analyse des Documents avec Lucene

Lemmatisation

Sample text: Such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Lovins stemmer: such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Porter stemmer: such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Paice stemmer: such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Analyse des Documents

- **Extraction des phrases**
- l'analyse grammaticale
 - → «syntagme nominal»
 - » → «chien noir » ou « petit enfant ».
 - → «syntagme verbal».
 - » → : "Ilyès aime l'informatique".

Analyse des Documents

- **Extraction des Mots composé**
- Un mot composé est constitué de deux ou plusieurs mots attachés pour former un nouveau mot
- Par exemple
 - ➔ « base de donnée », « *pommes de terre* » « *computer science* »

Analyse des Documents

- **Extraction des entités nommées**
- Jim bought 300 shares of Acme Corp. in 2006.
- ➔ [Jim]_{Person} bought 300 shares of [Acme Corp.]_{Organization} in [2006]_{Time}.

Analyse des Documents

- Unstructured Information Management Architecture (**UIMA**)
- General Architecture for Text Engineering (**GATE**)
- **Apache OpenNLP library**



The **Apache Software Foundati**
<http://www.apache.org/>

General

- [Home](#)
- [Download](#)
- [Maven Dependency](#)

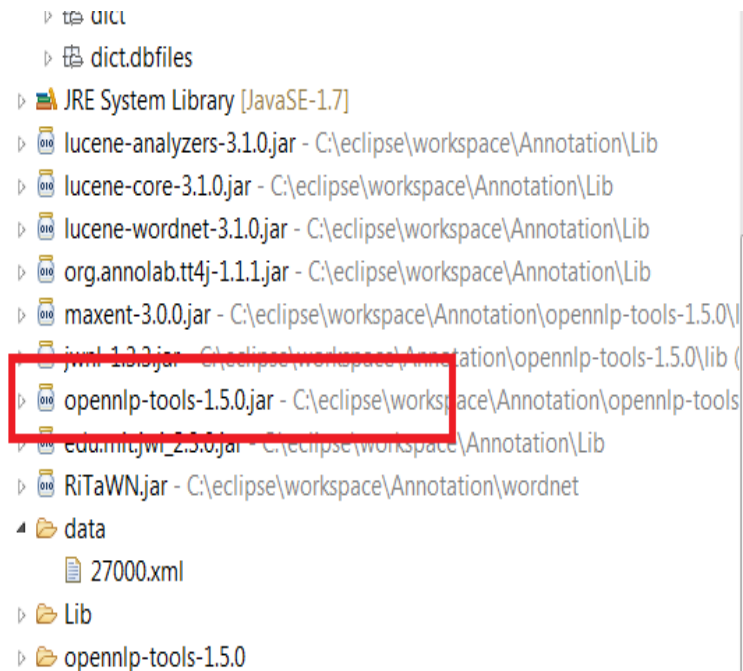
Welcome to Apache OpenNLP

The Apache OpenNLP library is a machine learning based toolkit for the processing

to support the most common NLP tasks such as tokenization, sentence detection,

Analyse des Documents

- **Apache OpenNLP library**



```
195     out2.println(txt2);
196     out2.close();
197     out2.close();
198     return txtlist;
199 }
200 public static void main(String[] args) throws IOException {
201
202     NLP nlp=new NLP();
203     String input=" ";
204     //String content="Pierre Vinken, 61 years old, will join the board as a nonexecut
205     nlp.tag(input);
206     //nlp.sent(content);
207     nlp.parse("los angeles is a state in USA");
208     //nlp.parse("Can anyone help me dig through OpenNLP's horrible documentation? ");
209     //nlp.parse("abdominal external oblique muscle");
210     //nlp.parse("external oblique muscle");
211
212 }
```

Analyse des Documents

- **Apache OpenNLP library**
- **Exemple** : Los Angeles is a state in USA



- (TOP (S (NP (NNS los) (NNS angeles)) (VP (VBZ is) (NP (DT a) (NN state)) (PP (IN in) (NP (NNP USA))))))



- NP los angeles
- NP a state
- NP USA

Analyse des Documents

Part of speech tags¹

CC - Coordinating conjunction	PRP - Personal pronoun
CD - Cardinal number	RB - Adverb
DT - Determiner	RBR - Adverb, comparative
EX - Existential there	RBS - Adverb, superlative
FW - Foreign word	RP - Particle
IN - Preposition or subordinating conjunction	SYM - Symbol
JJ - Adjective	TO - to
JJR - Adjective, comparative	UH - Interjection
JJS - Adjective, superlative	VB - Verb, base form
NN - Noun, singular or mass	VBD - Verb, past tense
NNS - Noun, plural	VBG - Verb, gerund or present participle
NNP - Proper noun, singular	VBN - Verb, past participle
NNPS - Proper noun, plural	VBP - Verb, non-3rd person singular present
PDT - Predeterminer	VBZ - Verb, 3rd person singular present
NP - Noun Phrase.	WDT - Wh-determiner
PP - Prepositional Phrase	WP - Wh-pronoun
VP - Verb Phrase.	WRB - Wh-adverb

¹ <http://bulba.sdsu.edu/jeanette/thesis/PennTags.html>

Analyse des Documents

- **Extraction des concepts**
- Concept se concentre plus sur la sémantique du texte.
- C'est une unité abstraite de connaissances qui explique le sens des mots.
- Il faut utiliser les ressources sémantiques externes pour extraire le concept à partir du texte.

The noun **computer science** has 1 sense

1. computer science, computing -- (the branch of engineering science that studies (with the aid of computers) computable processes and structures)

Analyse des Documents

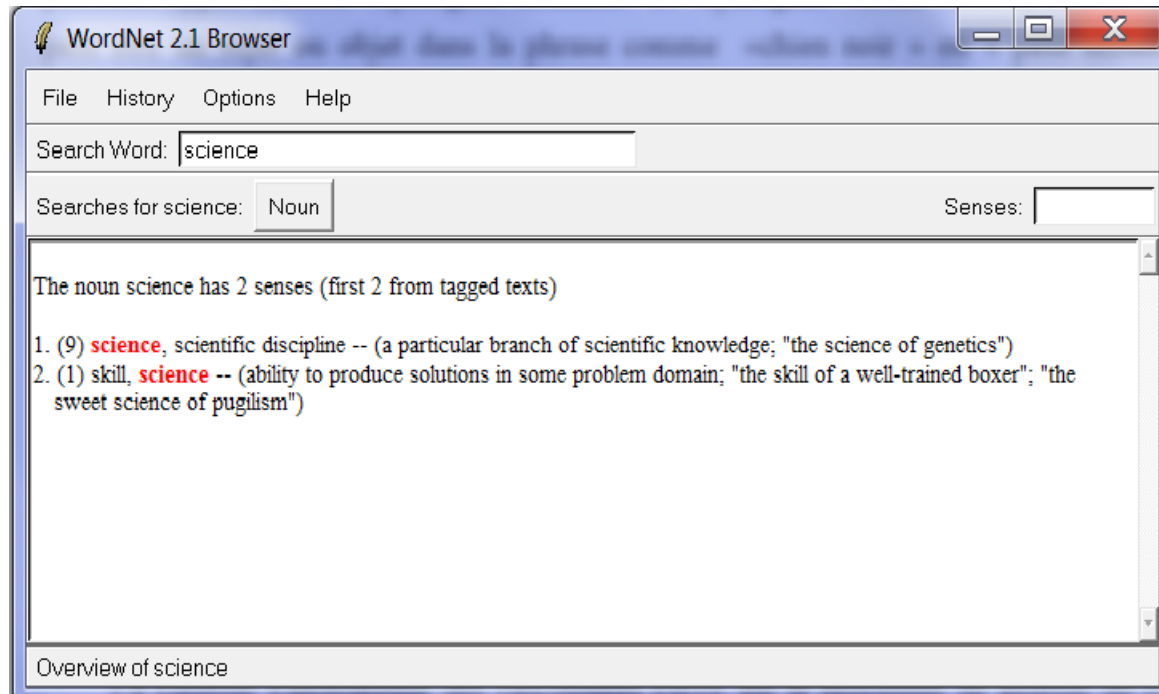
- **Extraction des concepts**

The noun **computer** has 2 senses (first 1 from tagged texts)

1. (6) computer, computing machine, computing device, data processor, electronic computer, information processing system -- (a machine for performing calculations automatically)
2. calculator, reckoner, figurer, estimator, computer -- (an expert at calculation (or at operating calculating machines))

Analyse des Documents

- Extraction des concepts



Indexation des Documents avec Lucene

- `Document doc = new Document();`
-
- `Field contentField1=new Field("contenu", contenu, Field.Store.YES, Field.Index.ANALYZED);`
-
- `doc.add(contentField1);`
-
- `doc.add(new Field("Titre", Titre_principale, Field.Store.YES, Field.Index.ANALYZED));`
-
- `indexWriter.addDocument(doc);`
-
- `indexWriter.optimize();`
- `indexWriter.close();`

Indexation des BDs avec Lucene

- `Void indexDocs(IndexWriter writer, Connection conn) throws Exception {`
- `String sql = "select id, name from employe";`
- `Statement stmt = conn.createStatement();`
- `ResultSet rs = stmt.executeQuery(sql);`
- `while (rs.next()) {`
- `Document d = new Document();`
- `d.add(new Field("id", rs.getString("id"), Field.Store.YES,Field.Index.NO));`
- `d.add(new Field("name", rs.getString("name"),`
`Field.Store.NO,Field.Index.TOKENIZED));`
- `writer.addDocument(d);`
- `}`
- `}`

Indexation des Documents avec Lucene

```
11 import java.util.logging.Logger;
12 import org.apache.lucene.analysis.standard.StandardAnalyzer;
13 import org.apache.lucene.document.Document;
14 import org.apache.lucene.document.Field;
15 import org.apache.lucene.index.IndexWriter;
16 import org.apache.lucene.store.FSDirectory;
17 import org.apache.lucene.util.Version;
18
19
20 public class Index {
21
22     public void indexFile(String contenu, String Titre_principale) throws IOException {
23
24         File indexDir=new File("C:/Users/pc/Documents/NetBeansProjects/Lucene/Index");
25
26         IndexWriter indexWriter = new IndexWriter(FSDirectory.open(indexDir),new StandardAnalyzer(Version.LUCENE_30),false
27
28         Document doc = new Document();
29
30
31         Field contentField1=new Field("contenu", contenu, Field.Store.YES, Field.Index.ANALYZED);
32         doc.add(contentField1);
33
34         doc.add(new Field("Titre", Titre_principale, Field.Store.YES, Field.Index.ANALYZED));
35         indexWriter.addDocument(doc);
36         indexWriter.optimize();
37         indexWriter.close();
38
39     }
40
41 }
```

Indexation des Documents avec Lucene

```
public static void main(String[] args) {
    // TODO Auto-generated method stub
    String contenu="nous allons ´etudier l'utilisation d'un index de recherche textuel OpenSource : Lucene";
    String titre="Lucene in Action";

    String contenu2="Pour aller plus loin, nous vous proposons de cr´eer un index pour un site web. Pour cela, v
"un objet qui va r´ecup´eer la page principale d'un site web d´etermin´e par son domaine. Ensuite pour chaque\n" +
"ancree (lien href) pr´esent sur la page, vous parcourrez le lien s'il est sur le m´eme site web et le parcourir\n" +
"r´ecursivement3";
    String titre2="Recherche d'information";

    Index ndx =new Index();
    try {
        ndx.indexFile(contenu, titre);
```

Indexation des Documents avec Lucene

```
public class Recherche {  
  
    public void search(File indexDir, String queryStr) throws IOException, ParseException {  
        int i;  
        Directory directory = FSDirectory.open(indexDir);  
        IndexReader reader = IndexReader.open(directory);  
        IndexSearcher searcher = new IndexSearcher(reader);  
        QueryParser parser = new QueryParser(Version.LUCENE_31, "contenu", new StandardAnalyzer(Version.LUCENE_31));  
        Query query = parser.parse(queryStr);  
  
        TopDocs topDocs = searcher.search(query, 10);  
        ScoreDoc[] hits = topDocs.scoreDocs;  
        for ( i = 0; i < hits.length; i++) {  
  
            int docId = hits[i].doc;  
            Document d = searcher.doc(docId);  
  
            System.out.println("-----");  
            System.out.println(i+1+"- " +d.get("Titre"));  
        }  
    }  
  
    public static void main(String[] args) throws IOException, ParseException {  
  
        File indexDir = new File("C:/Users/pc/Documents/NetBeansProjects/Lucene/Index");  
        Recherche R =new Recherche();  
        R.search(indexDir, "site web");  
    }  
}
```

Indexation des Documents avec Lucene

File Tools Settings Help

Overview Documents Search Files Plugins

Index name: **C:\eclipse\workspace\PDF\indexe**

Number of fields: **5**

Number of documents: **2250**

Number of terms: **69715**

Has deletions? / Optimized?: **No / Yes**

Last modified: **Thu Feb 04 19:51:48 GMT+01:00 2016**

Index version: 152a99a6b0f

Index format: -11 (Lucene 1.3 or prior)

Index functionality: unknown

TermInfos index divisor: N/A

Directory implementation: org.apache.lucene.store.MMapDirectory

Currently opened commit point: segments_3rq (Thu Feb 04 19:51:48 GMT+01:00 2016)

Current commit user data: --

Select fields from the list below, and press button to view top terms in these fields. No selection means all fields.

Available fields and term counts per field:

Name	Term count	%	Decoder
Titre	188	0,27 %	string utf8
contents	65 041	93,3 %	string utf8
page	0	0,00 %	string utf8
sous_titre	2 401	3,44 %	string utf8
sous_titre_di	2 085	2,99 %	string utf8

Select a field and set its value decoder: string utf8 Set

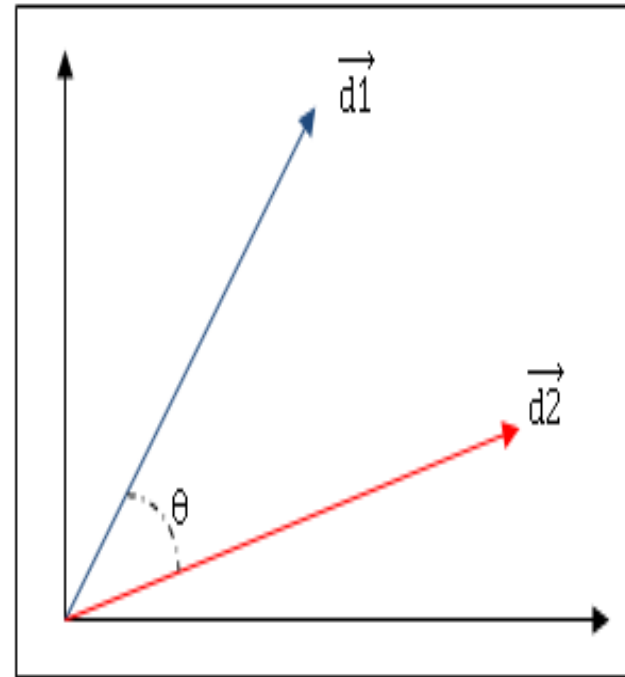
Top ranking terms. (Right-click for more options)

No	Rank	Field	Text
1	1919	contents	de
2	1702	contents	et
3	1696	contents	a
4	1682	contents	le
5	1613	contents	la
6	1566	contents	d
7	1539	contents	un
8	1523	contents	est
9	1515	contents	une
10	1506	contents	l

Index name: **C:\eclipse\workspace\PDF\indexe**

Indexation des Documents avec Lucene

- Requête / Document
- La similarité entre deux vecteurs de documents (d_1 : requête) et (d_2)



- $\text{Cos}(\theta) = \frac{\vec{d}_1 \cdot \vec{d}_2}{|\vec{d}_1| |\vec{d}_2|}$
- Documents et requêtes sont représentés
- comme des vecteurs
- $D1 = \{w_{11}, w_{21}, \dots, w_{N, 1}\}$,
- $q = \{w_{1,q}, w_{2,q}, \dots, w_{N, q}\}$
- $\vec{d}_1 \cdot \vec{q} = \sum w_{d1} * w_{q}$
- Norme $|\vec{d}_1| = \sqrt{w_1^2 + w_2^2 + w_3^2 \dots}$

