

Année 2014-2015
Thomas ROBERT

COMPTE-RENDU DE TP

RECHERCHE D'INFORMATION

Table des matières

1	Scala	3
2	Sauvegarde des résultats et PageLoader	4
3	Crawling	5
4	Indexation	6
4.1	Récupération du texte des pages	6
4.2	Calcul des fréquences des termes dans chaque document (<i>term frequencies</i>) . .	6
4.3	Calcul des fréquences des termes dans les documents (<i>inverse document frequencies</i>)	6
4.4	Calcul des vecteurs tf-idf	7
4.5	Conclusion	7
5	Recherche	8
6	Résultats	9
7	Conclusion	10
	Annexes	11
A	PageLoader.scala	11
B	Crawler.scala	13
C	Stemmer.scala	15
D	Indexer.scala	19
E	Search.scala	21

1 | Scala

Ce TP/mini-projet a été réalisé en Scala, qui est un langage fonctionnel orienté objet (alors que Java est impératif orienté objet). Scala est un langage compilé en bytecode et exécuté en JVM. Il est donc entièrement compatible avec le code écrit en Java : une librairie Scala est utilisable en Java et une librairie Java est utilisable en Scala.

Les deux avantages principaux de ce langage par rapport à Java dans le cadre de ce TP sont la possibilité d'utiliser les opérations fonctionnelles classiques sur les collections et la légèreté de la syntaxe.

En effet, Scala est fourni avec des collections (listes, ensembles, dictionnaires, etc.) qui possède un ensemble très riche d'opérateurs classiques en programmation fonctionnelle (*map*, *reduce*, *fold*, *filter*, *groupBy*, etc.), qui permettent d'appliquer des transformations sur ces collections avec une syntaxe très légère.

Par exemple, la ligne ci-dessous supprime du dictionnaire `urlAndHTMLDatabaseCandidates` (associant url au code source de la page) toutes les pages qui ne sont pas du HTML, puis nettoie le code HTML pour en extraire le texte.

```
1 val urlAndTextDatabase = urlAndHTMLDatabaseCandidates filter { case (_, html) => isHTML(html) } mapValues cleanHTML
```

Le code Java correspondant serait :

```
1 Map<String , String> urlAndTextDatabase = new HashMap<String , String >();
2
3 for(Entry<String , String> pair: urlAndTextDatabase.entrySet()) {
4     String key = pair.getKey();
5     String value = pair.getValue();
6
7     if (isHTML(value)) { // filter
8         value = cleanHTML(value); // mapValues
9         urlAndTextDatabase.put(key, value);
10    }
11 }
```

Soit une syntaxe bien plus lourde, et selon moi bien moins lisible.

L'emploi d'opérateurs fonctionnels courants et de fonctions anonymes permettent donc de transformer des collections de façon très concise.

Scala favorise d'ailleurs globalement la concision, par son mécanisme d'inférence de type qui permet de ne pas avoir à répéter inutilement le type des variables, qui sont déduit de ce que l'on met dans la variable, la légèreté de déclaration des classes, la non-nécessité des points-virgules, etc.¹ Cette concision et puissance du langage permet de passer plus de temps à réfléchir au problème plutôt qu'à taper du code verbeux.

Ce sont donc ces deux avantages majeurs (opérateurs fonctionnels et concision) de Scala qui m'ont le plus servi et « aidé » durant le TP. Mais Scala est un langage complexe et très riche au fonctionnalités bien plus étendues que ce que je viens de présenter.

1. Voir [http://www.wikiwand.com/en/Scala_\(programming_language\)/Syntactic_flexibility](http://www.wikiwand.com/en/Scala_(programming_language)/Syntactic_flexibility)

2 | Sauvegarde des résultats et PageLoader

L'objectif du TP porte sur l'indexation de pages web afin d'y effectuer des recherches. Il est donc nécessaire de télécharger le code source d'un grand nombre de pages web.

Afin de ne pas avoir à recharger plusieurs fois une page web, le module *PageLoader* a été codé. Ce module est appelé pour charger une page, et s'occupe de mettre en cache chaque page chargée. Ainsi, si une page a déjà été chargée (lors d'un lancement précédent), elle sera chargée depuis le cache, permettant de gagner énormément de temps.

Notons que, de manière générale, chaque module sauvegarde son résultat dans un fichier afin que le module suivant puisse être lancé sans avoir besoin de réexécuter les calculs des étapes suivantes.

3 | Crawling

La première étape consiste à crawler le web (suivre les liens de page en page) afin de constituer une « base de données » d'URL.

Ce module de crawling comprend globalement deux parties, une partie qui permet d'extraire et de nettoyer des URL contenues dans une page web, et une partie qui s'occupe de crawler de façon récursive les pages web trouvées.

La première partie n'est pas vraiment intéressante, il s'agit simplement d'être capable reconstruire une URL à partir d'un lien et d'une URL parente, et d'extraire les liens (donc les balises <a>) d'une page web. Tout cela se fait avec des expressions régulières.

La deuxième partie est relativement simple mais plus intéressante.

```
1  def crawl(source: String, limit: Int): Set[String] = {  
2  
3      def crawlRec(set: Set[String], nextHrefs: List[String]): Set[String] =  
4          {  
5              if (set.size > limit)  
6                  set  
7              else {  
8                  // crawl new links to add to set  
9                  val toAdd = getUrls(nextHrefs.head) filter (x => !(set contains x))  
10                 // recursive call with new links minus processed link  
11                 crawlRec(set ++ toAdd, nextHrefs.tail :: toAdd)  
12             }  
13         }  
14     val init = getUrls(source)  
15     crawlRec(init.toSet, init)  
16 }
```

La fonction de crawling récursive prend une base d'URL et une liste d'URL à parcourir. Elle prend le premier lien de la liste, en extrait les URL enfants, filtre les URL qui ne sont pas déjà dans la base, et ajoute ces URL à la base et à la liste, donc on aura extrait le premier élément puisqu'il a été traité.

On remarque qu'encore une fois, Scala nous permet d'écrire cela très « joliment ».

4 | Indexation

La phase d'indexation consiste, à partir de la base de données d'URL générée par la phase de crawling, à générer un modèle vectoriel tf-idf du corpus.

4.1 Récupération du texte des pages

Pour cela, on commence par récupérer le code source de la page pointée par chaque URL et on vérifie que ce code est du HTML. Si oui, on conserve la page dans le corpus et on extrait de la page HTML le texte brut en utilisant la librairie Jsoup. On dispose donc d'une collection `Map[url : String, texte : String]`.

4.2 Calcul des fréquences des termes dans chaque document (*term frequencies*)

Pour chaque page, on calcule alors les *term frequencies*.

On commence donc par découper les textes à chaque espace afin de récupérer une liste de mots (on ne conserve que les mots constitués de lettres et d'au moins 3 caractères). On *stemme* alors chaque mot de la liste en utilisant un *Stemmer de Porter* en Scala².

A partir de la liste de mots *stemmés*, on calcule une *map* du nombre d'occurrence de chaque mot `Map[mot : String, nbOccs : Int]`, et ce pour chaque document.

On traite ensuite chaque `Map[mot, nbOccs]` pour la transformer en une *map* des poids *tf* en calculant pour chaque nombre d'occurrence $\log(\text{nbOccs})$ que l'on divise par le nombre de mots dans le texte.

Au final, on a donc une `Map[url : String, Map[mot : String, tf : Double]]`.

4.3 Calcul des fréquences des termes dans les documents (*inverse document frequencies*)

On veut maintenant calculer le nombre de documents où chaque terme apparaît, quel que soit le nombre d'apparitions de ce terme dans chaque document.

Pour cela, on part de notre base de *tf*, `Map[url : String, Map[mot : String, tf : Double]]`. Pour chaque URL, on extrait la liste des mots contenus dans le document (les clés de la `Map[mot, tf]`). On concatène toutes les listes obtenues et on a donc une liste dans laquelle un mot apparaît à chaque fois qu'il est contenu dans un document.

Il suffit donc de compter le nombre d'occurrences de chaque mot dans cette liste pour savoir dans combien de documents ce mot est présent. On a donc une `Map[mot : String, nbOccsDocs : Int]`.

Enfin, pour calculer les *inverse document frequencies*, on calcule juste pour chaque mot $\log(\text{nbDocs} / \text{nbOccsDocs})$. On obtient donc `Map[mot : String, idf : Double]`.

2. <https://github.com/aztek/porterstemmer>

4.4 Calcul des vecteurs tf-idf

Enfin, une fonction permet de calculer le vecteur tf-idf normalisé d'un document en multipliant chaque valeur tf par la valeur idf correspondante dans le modèle et en divisant chaque élément par la norme du vecteur.

4.5 Conclusion

A la fin de cette étape, on dispose donc de deux collections, une collection de *term frequencies* `Map[url : String, Map[mot : String, tf : Double]]` et une collection de *inverse document frequencies* `Map[mot : String, idf : Double]`.

Ces deux collections sont sauvegardées dans des fichiers JSON.

Notons encore une fois que les collections de Scala et son paradigme fonctionnel permettent d'écrire ces transformations de collection de manière très succincte. Par exemple, voici les fonctions permettant de calculer les TF et IDF d'une base d'URL :

```
1  def getNbOocs(wordsList: Iterable[String]): Map[String, Int] = {
2    wordsList groupBy (x => x) mapValues (_.size)
3  }
4
5  def getFrequencies(wordsList: Iterable[String]): TF = {
6    getNbOocs(wordsList) mapValues (tf => if (tf == 0) 0 else 1 + Math.log10
7      (tf.toDouble) / wordsList.size)
8  }
9
10 def getPageFrequencies(page: String): TF = {
11   getFrequencies(getStems(getWords(page)))
12 }
13
14 def computeTfAndIdf(urlDatabase: Iterable[String]): (TFs, IDF) = {
15   // filter to only keep valid HTMLs
16   val urlAndHTMLDatabaseCandidates = (urlDatabase map (url => url ->
17     PageLoader.getContent(url)) ).toMap
18   val urlAndTextDatabase = urlAndHTMLDatabaseCandidates filter { case (_,
19     html) => isHTML(html) } mapValues cleanHTML
20   val nbDocs = urlAndTextDatabase.size
21
22   // compute termsFrequencies Map[url, Map[term, tf]]
23   val termsFrequencies = urlAndTextDatabase mapValues getPageFrequencies
24
25   // compute documentsFrequencies Map[term, nbDocsWithTerm]
26   val documentsFrequencies = getNbOocs(termsFrequencies.values flatMap (
27     _.keys))
28
29   // compute inverseDocumentsFrequencies Map[term, idf]
30   val inverseDocumentsFrequencies = documentsFrequencies mapValues (
31     nbDocsWithTerm => Math.log10(nbDocs.toDouble / nbDocsWithTerm ))
32
33   (termsFrequencies, inverseDocumentsFrequencies)
```

5 | Recherche

Enfin, on termine par le module de recherche. Pour cela, on demande à l'utilisateur d'entrer une requête. On calcule le vecteur tf-idf de la requête avec les fonctions du module d'indexation vu précédemment.

On dispose donc des tf-idf du corpus et de la requête, il suffit donc de les comparer par produit vectoriel.

On obtient alors une mesure de similarité requête-document pour chaque document sous la forme d'une `Map[url : String, similarité : Double]`. Il suffit alors de transformer cette *map* en séquence de paires (url, similarité) que l'on ordonne par paire.

Voici le code permettant de réaliser tout ça :

```
1  def computeSimilarities(queryTfIdf: Indexer.TFIDF, docsTfIdfs: Indexer.  
    TFIDFs): Map[String, Double] = {  
2      docsTfIdfs mapValues {  
3          //  $\forall \text{ doc, compute } \cos(q, d) = \langle q, d \rangle = \sum q_i * d_i$   
4          dis  $\Rightarrow$  queryTfIdf map {  
5              case (qiTerm, qi)  $\Rightarrow$  qi * dis.getOrElse(qiTerm, 0.0d)  
6          } reduce (_ + _)  
7      }  
8  }  
9  
10 def orderResults(similarities: Map[String, Double]): Seq[(String, Double)]  
    = {  
11     similarities.toSeq.sortBy(_._2)  
12 }  
13  
14 def computeQuery(query: String): Seq[(String, Double)] = {  
15     val queryTf = Indexer.getPageFrequencies(query)  
16     val queryTfIdf = Indexer.computeTfIdf(queryTf, idf)  
17  
18     orderResults(computeSimilarities(queryTfIdf, docsTfIdfs))  
19 }
```


6 | Résultats

On obtient par exemple le résultat suivant pour la recherche « information » :

- 1 <http://research.microsoft.com/en-us/projects/snssurvey/>,
0.01955816809945687
- 2 <http://www.thefreelibrary.com/Hohhof%2c+Bonnie-a1>, 0.019531222142308673
- 3 <http://www.thefreelibrary.com/Bilal%2c+Dania-a1>, 0.01763201349176205
- 4 <https://zh.wikipedia.org/wiki/%E4%BF%A1%E6%81%AF%E6%AA%A2%E7%B4%A2>,
0.01629096794038907
- 5 <http://trec.nist.gov>, 0.014911142180124952
- 6 https://fa.wikipedia.org/wiki/%D8%A8%D8%A7%D8%B2%DB%8C%D8%A7%D8%A8%DB%8C_%D8%A7%D8%B7%D9%84%D8%A7%D8%B9%D8%A7%D8%AA, 0.014443995373372875
- 7 https://tg.wikipedia.org/wiki/%D2%B6%D1%83%D1%81%D1%82%D1%83%D2%B7%D3%AF%D0%B8_%D0%B8%D0%BD%D1%84%D0%BE%D1%80%D0%BC%D0%B0%D1%82%D1%81%D0%B8%D0%BE%D0%BD%D3%A3, 0.013222401253283245
- 8 <http://www.oxforddictionaries.com/social/cite.html?dictCode=english&entryId=information-retrieval>, 0.013209528647371774
- 9 https://el.wikipedia.org/wiki/%CE%91%CE%BD%CE%AC%CE%BA%CF%84%CE%B7%CF%83%CE%B7_%CF%80%CE%BB%CE%B7%CF%81%CE%BF%CF%86%CE%BF%CF%81%CE%B9%CF%8E%CE%BD,
0.013085771367958876
- 10 https://bg.wikipedia.org/wiki/%D0%98%D0%B7%D0%B2%D0%BB%D0%B8%D1%87%D0%B0%D0%BD%D0%B5_%D0%BD%D0%B0_%D0%B8%D0%BD%D1%84%D0%BE%D1%80%D0%BC%D0%B0%D1%86%D0%B8%D1%8F, 0.012967414816850724
- 11 <http://www.thefreedictionary.com/advertise-with-us.htm>, 0.01287519427530754
- 12 <http://legal-dictionary.thefreedictionary.com/information+resources+management>, 0.012680745177095158
- 13 <http://idioms.thefreedictionary.com/information+resources+management>,
0.012591283318025044
- 14 <http://legal-dictionary.thefreedictionary.com/information+resources>,
0.012495445775518935
- 15 <http://www.britannica.com/topic/1379580/feedback>, 0.012114522486594627
- 16 https://en.wikipedia.org/wiki/Template:Information_science,
0.011964723754558981
- 17 https://en.wikipedia.org/w/index.php?title=Template:Information_science&action=edit, 0.011964723754558981
- 18 <http://www.britannica.com/print/topic/1379580>, 0.011852939474837737
- 19 <http://idioms.thefreedictionary.com/Information+Retrieval>,
0.011791158995515177
- 20 <http://financial-dictionary.thefreedictionary.com/Information+retrieval>,
0.011783448887460997
- 21 <http://medical-dictionary.thefreedictionary.com/Information+retrieval>,
0.011783281512735755
- 22 <http://legal-dictionary.thefreedictionary.com/Information+Retrieval>,
0.0117831818092679
- 23 <http://financial-dictionary.thefreedictionary.com/Information+Retrieval>,
0.011783154653036983
- 24 <http://medical-dictionary.thefreedictionary.com/Information+Retrieval>,
0.011782987907877943
- 25 https://en.wikipedia.org/wiki/Divergence-from-randomness_model,
0.011754402492824782

7 | Conclusion

Ce TP aura donc été l'occasion d'implémenter l'algorithme TF-IDF, sans doute le plus important en recherche d'information et en traitement de documents de manière générale (on peut par exemple l'utiliser en *Machine Learning* par exemple). Cet algorithme a été appliqué à un corpus *crawlé* sur le web.

Cette implémentation pourrait clairement être améliorée. Par exemple, le *stemmer* utilisé est un *stemmer* anglais, mais rien ne nous assure que les pages crawlées le sont, même si nous sommes partis d'une page anglaise. En l'occurrence, dans les résultats ci-dessus, on voit des pages wikipédia dans de nombreuses langues apparaître. Il serait par exemple intéressant de détecter la langue de chaque page afin de ne conserver que des pages en anglais.

Il aurait également été intéressant de mesurer les performances en précision et en rappel de notre « moteur de recherche », mais je n'ai pas eu le temps de le faire.

Annexes

A | PageLoader.scala

```
1 package fr.thomasrobert.inforetrieval
2
3 import java.nio.file.{Paths, Files}
4 import java.io.{File, PrintWriter}
5 import java.net.URL
6
7 import scala.io.Source
8
9 object PageLoader {
10
11   def getSha1(s: String) = {
12     val md = java.security.MessageDigest.getInstance("SHA-1")
13     md.digest(s.getBytes("UTF-8")).map("%02x".format(_)).mkString
14   }
15
16   def saveAsFile(HTML: String, filename: String) = {
17     val writer = new PrintWriter(new File(filename))
18     writer.write(HTML)
19     writer.close()
20   }
21
22   def loadFile(filename: String): String = {
23     try {
24       Source.fromFile(filename).getLines().mkString("\n")
25     }
26     catch {
27       case e: Exception => ""
28     }
29   }
30
31   def loadFileAsSet(filename: String): Set[String] = {
32     try {
33       Source.fromFile(filename).getLines().toSet
34     }
35     catch {
36       case e: Exception => Set()
37     }
38   }
39
40   def getLines(url: String): Iterator[String] = {
41     val connection = new URL(url).openConnection
42     connection.setRequestProperty("User-Agent", "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.0)")
43
44     Source.fromInputStream(connection.getInputStream).getLines()
45   }
46
47   def downloadContent(url: String): String = {
48     try
49       getLines(url).toList.mkString(" ")
50     catch {
51       case e: Exception => ""
52     }
53   }
54 }
```

```
53     }
54
55     def getContent(url: String): String = {
56         val filename = "files/" + getSha1(url) + ".txt"
57         if (Files.exists(Paths.get(filename)))
58             loadFile(filename)
59         else {
60             val HTML = downloadContent(url)
61             saveAsFile(HTML, filename)
62             HTML
63         }
64     }
65
66 }
```

B | Crawler.scala

```
1 package fr.thomasrobert.inforetrieval
2
3 object Crawler {
4   /**
5    * Return the parts of an URL
6    */
7   def getInfos(url : String): (String, String, String) = {
8     val protocol = ("^([^\:]+)://([^\:]+)/".r findFirstMatchIn url).get.group
9       (1)
10    val domain = ("^([^\:]+)://([^\:]+)/".r findFirstMatchIn url).get.group
11      (2)
12    val directory = ("^(.+/) [^/]*$".r findFirstMatchIn url).get.group(1)
13    (protocol, domain, directory)
14  }
15
16  /**
17   * Get an absolute URL from an HREF and information on the source URL
18   */
19  def getAbsoluteHref(protocol: String, domain: String, directory: String)(
20    href : String): String = {
21    if (matches("^[:]+://.+ ", href))
22      href
23    else if (matches("^//.+ ", href))
24      protocol + ":" + href
25    else if (matches("^/.+ ", href))
26      protocol + "://" + domain + href
27    else
28      directory + href
29  }
30
31  /**
32   * Remove the anchor from an URL
33   */
34  def removeAnchors(s: String): String = {
35    "#.+$".r replaceAllIn(s, "")
36  }
37
38  /**
39   * Utility function that says if "s" matches the regex string "r"
40   */
41  def matches(r : String, s : String): Boolean = r.r.pattern.matcher(s).
42    matches()
43
44  /**
45   * Extract URLs in the page "parentUrl"
46   */
47  def getUrls(parentUrl: String): List[String] = {
48    try {
49      val (protocol, domain, directory) = getInfos(parentUrl)
50
51      val HTML = PageLoader.getContent(parentUrl)
52      val HTMLbody = ("<body [^>]+>(.+)</body>".r findFirstMatchIn HTML).get.
53        group(1)
54
55      val regex = "href=[ '\"]([^\s'\"]+) [ '\"]".r
56
57      val hrefs = (regex findAllIn HTMLbody) map (x => (regex
```

```
        findFirstMatchIn x).get.group(1))

53
54     def getAbsoluteHrefLocal = getAbsoluteHref(protocol, domain,
        directory)(_))
55
56     (hrefs map getAbsoluteHrefLocal map removeAnchors).toSet.toList
57 }
58 catch {
59     case e: Exception => List()
60 }
61 }
62
63 /**
64  * Crawl to look for URLs from a source, stop the crawling when we found
        at least "limit" URLs.
65  */
66 def crawl(source: String, limit: Int): Set[String] = {
67
68     def crawlRec(set: Set[String], nextHrefs: List[String]): Set[String] =
        {
69         if (set.size > limit)
70             set
71         else {
72             // crawl new links to add to set
73             val toAdd = getUrls(nextHrefs.head) filter (x => !(set contains x))
74             // recursive call with new links minus processed link
75             crawlRec(set ++ toAdd, nextHrefs.tail ::: toAdd)
76         }
77     }
78
79     val init = getUrls(source)
80     crawlRec(init.toSet, init)
81 }
82
83 /**
84  * Main
85  */
86 def main(args: Array[String]) {
87     PageLoader.saveAsFile(crawl("https://duckduckgo.com/html/?q=information
        %20retrieval", 3000).mkString("\n"), "files/database.txt")
88     val database = PageLoader.loadFileAsSet("files/database.txt")
89     println(database.size)
90     println(database)
91 }
92 }
```


C | Stemmer.scala

```
1 package fr.thomasrobert.inforetrieval
2
3 /**
4  * Scala implementation of Porter's stemming algorithm.
5  *
6  * See http://snowball.tartarus.org/algorithms/porter/stemmer.html
7  * for description of the algorithm itself.
8  *
9  * @author Evgeny Kotelnikov <evgeny.kotelnikov@gmail.com>
10  */
11 object Stemmer {
12   def stem(word: String): String = {
13     // Deal with plurals and past participles
14     var stem = new Word(word).applyReplaces(
15       "sses" → "ss",
16       "ies" → "i",
17       "ss" → "ss",
18       "s" → "")
19
20     if ((stem matchedBy ((~v~) + "ed")) ||
21         (stem matchedBy ((~v~) + "ing"))) {
22
23       stem = stem.applyReplaces(~v~)("ed" → "", "ing" → "")
24
25       stem = stem.applyReplaces(
26         "at" → "ate",
27         "bl" → "ble",
28         "iz" → "ize",
29         (~d and not(~L or ~S or ~Z)) → singleLetter,
30         (m == 1 and ~o) → "e")
31     } else {
32       stem = stem.applyReplaces(((m > 0) + "eed") → "ee")
33     }
34
35     stem = stem.applyReplaces(((~v~) + "y") → "i")
36
37     // Remove suffixes
38     stem = stem.applyReplaces(m > 0)(
39       "ational" → "ate",
40       "tional" → "tion",
41       "enci" → "ence",
42       "anci" → "ance",
43       "izer" → "ize",
44       "abli" → "able",
45       "alli" → "al",
46       "entli" → "ent",
47       "eli" → "e",
48       "ousli" → "ous",
49       "ization" → "ize",
50       "ation" → "ate",
51       "ator" → "ate",
52       "alism" → "al",
53       "iveness" → "ive",
54       "fulness" → "ful",
55       "ousness" → "ous",
56       "aliti" → "al",
57       "iviti" → "ive",
```

```
58     "biliti" → "ble")
59
60     stem = stem.applyReplaces(m > 0)(
61         "icate" → "ic",
62         "ative" → "",
63         "alize" → "al",
64         "iciti" → "ic",
65         "ical" → "ic",
66         "ful" → "",
67         "ness" → "")
68
69     stem = stem.applyReplaces(m > 1)(
70         "al" → "",
71         "ance" → "",
72         "ence" → "",
73         "er" → "",
74         "ic" → "",
75         "able" → "",
76         "ible" → "",
77         "ant" → "",
78         "ement" → "",
79         "ment" → "",
80         "ent" → "",
81         ((~S or ~T) + "ion") → "",
82         "ou" → "",
83         "ism" → "",
84         "ate" → "",
85         "iti" → "",
86         "ous" → "",
87         "ive" → "",
88         "ize" → "")
89
90     // Tide up a little bit
91     stem = stem.applyReplaces(((m > 1) + "e") → "",
92         (((m == 1) and not(~o)) + "e") → "")
93
94     stem = stem.applyReplaces ((m > 1 and ~d and ~L) → singleLetter)
95
96     stem.toString
97 }
98
99 /**
100  * Pattern that is matched against the word.
101  * Usually, the end of the word is compared to suffix,
102  * and the beginning is checked to satisfy a condition.
103  * @param condition Condition to be checked
104  * @param suffix Expected suffix of the word
105  */
106 private case class Pattern(condition: Condition, suffix: String)
107
108 /**
109  * Condition, that is checked against the beginning of the word
110  * @param predicate Predicate to be applied to the word
111  */
112 private case class Condition(predicate: Word ⇒ Boolean) {
113     def + = new Pattern(this, _: String)
114
115     def unary_~ = this // just syntactic sugar
116
117     def ~ = this
```

```
118
119     def and(condition: Condition) = Condition((word) ⇒ predicate(word) &&
120         condition.predicate(word))
121
122     def or(condition: Condition) = Condition((word) ⇒ predicate(word) ||
123         condition.predicate(word))
124 }
125
126 private def not: Condition ⇒ Condition = {
127     case Condition(predicate) ⇒ Condition(!predicate(_))
128 }
129
130 private val emptyCondition = Condition(_ ⇒ true)
131
132 private object m {
133     def >(measure: Int) = Condition(_.measure > measure)
134
135     def ==(measure: Int) = Condition(_.measure == measure)
136 }
137
138 private val S = Condition(_ endsWith "s")
139 private val Z = Condition(_ endsWith "z")
140 private val L = Condition(_ endsWith "l")
141 private val T = Condition(_ endsWith "t")
142
143 private val d = Condition(_.endsWithCC)
144
145 private val o = Condition(_.endsWithCVC)
146
147 private val v = Condition(_.containsVowels)
148
149 /**
150  * Builder of the stem
151  * @param build Function to be called to build a stem
152  */
153 private case class StemBuilder(build: Word ⇒ Word)
154
155 private def suffixStemBuilder(suffix: String) = StemBuilder(_ + suffix)
156
157 private val singleLetter = StemBuilder(_ trimSuffix 1)
158
159 private class Word(string: String) {
160     val word = string.toLowerCase
161
162     def trimSuffix(suffixLength: Int) = new Word(word substring (0, word.
163         length - suffixLength))
164
165     def endsWith = word endsWith _
166
167     def +(suffix: String) = new Word(word + suffix)
168
169     def satisfies = (_: Condition).predicate(this)
170
171     def hasConsonantAt(position: Int): Boolean =
172         (word.indices contains position) && (word(position) match {
173             case 'a' | 'e' | 'i' | 'o' | 'u' ⇒ false
174             case 'y' if hasConsonantAt(position + 1) ⇒ false
175             case _ ⇒ true
176         })
177 }
```

```
175 def hasVowelAt = !hasConsonantAt(_: Int)
176
177 def containsVowels = word.indices exists hasVowelAt
178
179 def endsWithCC =
180   (word.length > 1) &&
181   (word(word.length - 1) == word(word.length - 2)) &&
182   hasConsonantAt(word.length - 1)
183
184 def endsWithCVC =
185   (word.length > 2) &&
186   hasConsonantAt(word.length - 1) &&
187   hasVowelAt(word.length - 2) &&
188   hasConsonantAt(word.length - 3) &&
189   !(Set('w', 'x', 'y') contains word(word.length - 2))
190
191 /* *
192  * Measure of the word — the number of VCs
193  * @return integer
194  */
195 def measure = word.indices.filter(pos => hasVowelAt(pos) &&
196   hasConsonantAt(pos + 1)).length
197
198 def matchedBy: Pattern => Boolean = {
199   case Pattern(condition, suffix) =>
200     endsWith(suffix) && (trimSuffix(suffix.length) satisfies condition)
201 }
202
203 def applyReplaces(replaces: (Pattern, StemBuilder)*): Word = {
204   for ((pattern, stemBuilder) <- replaces if matchedBy(pattern))
205     return stemBuilder build trimSuffix(pattern.suffix.length)
206   this
207 }
208
209 def applyReplaces(commonCondition: Condition)(replaces: (Pattern,
210   StemBuilder)*): Word =
211   applyReplaces(replaces map {
212     case (Pattern(condition, suffix), stemBuilder) =>
213       (Pattern(commonCondition and condition, suffix), stemBuilder)
214   }: _*)
215
216 override def toString = word
217 }
218
219 private implicit def pimpMyRule[P <% Pattern, SB <% StemBuilder]
220 (rule: (P, SB)): (Pattern, StemBuilder) = (rule._1, rule._2)
221 private implicit def emptyConditionPattern: String => Pattern = Pattern(
222   emptyCondition, _)
223 private implicit def emptySuffixPattern: Condition => Pattern = Pattern(_,
224   "")
225 private implicit def suffixedStemBuilder: String => StemBuilder =
226   suffixStemBuilder
227 }
```

D | Indexer.scala

```
1 package fr.thomasrobert.inforetrieval
2
3 import org.jsoup.Jsoup
4 import org.json4s.native.Serialization.{read, write}
5 import org.json4s.DefaultFormats
6
7 object Indexer {
8
9   // formats for JSON lib
10  implicit val formats = DefaultFormats
11
12  // types
13  type IDF = Map[String, Double]
14  type TF = Map[String, Double]
15  type TFIDF = Map[String, Double]
16  type TFs = Map[String, TF]
17  type TFIDFs = Map[String, TFIDF]
18
19  // file save
20  def saveAsFile(tfs: TFs, idf: IDF) = {
21    PageLoader.saveAsFile(write(tfs), "files/tfs.json")
22    PageLoader.saveAsFile(write(idf), "files/idf.json")
23  }
24
25  def loadFromFile(): (TFs, IDF) = {
26    val tfs = read[TFs](PageLoader.loadFile("files/tfs.json"))
27    val idf = read[IDF](PageLoader.loadFile("files/idf.json"))
28    (tfs, idf)
29  }
30
31  // process HTML & words
32  def isHTML(page: String): Boolean = page.toLowerCase.replaceAll("[ \t]", "").contains("<body")
33  def cleanHTML(page: String): String = Jsoup.parse(page).text
34
35  def getWords(s: String): Iterable[String] = {
36    val string = "[^a-zA-Z]".r.replaceAllIn(s, " ")
37    string.split " " filter (x => x.length > 3) filter (x => x forall (l => l.isLetter))
38  }
39
40  def getStems(words: Iterable[String]): Iterable[String] = {
41    words map Stemmer.stem
42  }
43
44  // process frequencies (tf & idf)
45  def getNbOcocs(wordsList: Iterable[String]): Map[String, Int] = {
46    wordsList groupBy (x => x) mapValues (_.size)
47  }
48
49  def getFrequencies(wordsList: Iterable[String]): TF = {
50    getNbOcocs(wordsList) mapValues (tf => if (tf == 0) 0 else 1 + Math.log10
51      (tf.toDouble) / wordsList.size)
52  }
53
54  def getPageFrequencies(page: String): TF = {
55    getFrequencies(getStems(getWords(page)))
56  }
```

```
55 }
56
57 def computeTfAndIdf(urlDatabase: Iterable[String]): (TFs, IDF) = {
58
59     // filter to only keep valid HTMLs
60     val urlAndHTMLDatabaseCandidates = (urlDatabase map (url => url =>
61         PageLoader.getContent(url)) ).toMap
62     val urlAndTextDatabase = urlAndHTMLDatabaseCandidates filter { case (_,
63         html) => isHTML(html) } mapValues cleanHTML
64     val nbDocs = urlAndTextDatabase.size
65
66     // compute termsFrequencies Map[url, Map[term, tf]]
67     val termsFrequencies = urlAndTextDatabase mapValues getPageFrequencies
68
69     // compute documentsFrequencies Map[term, nbDocsWithTerm]
70     val documentsFrequencies = getNbOccs(termsFrequencies.values flatMap (_
71         .keys))
72
73     // compute inverseDocumentsFrequencies Map[term, idf]
74     val inverseDocumentsFrequencies = documentsFrequencies mapValues (
75         nbDocsWithTerm => Math.log10(nbDocs.toDouble / nbDocsWithTerm ))
76
77     (termsFrequencies, inverseDocumentsFrequencies)
78 }
79
80 def computeTfIdf(tfs: TF, idf: IDF): TFIDF = {
81     val tfIdfs = tfs map { case (term, tf) => (term, tf * idf.getOrElse(term,
82         0.0d)) }
83     val norm = Math.sqrt( tfIdfs.values.map(x => x * x).fold(0.0d)(_ + _) )
84     tfIdfs mapValues (_ / norm)
85 }
86
87 def computeTfIdfs(tfs: TFs, idf: IDF): TFIDFs = {
88     tfs mapValues (tf => computeTfIdf(tf, idf))
89 }
90
91 // main
92 def main(args: Array[String]) {
93
94     val (termsFrequencies, inverseDocumentsFrequencies) = computeTfAndIdf(
95         PageLoader.loadFileAsSet("files/database.txt"))
96     saveAsFile(termsFrequencies, inverseDocumentsFrequencies)
97     // val (termsFrequencies, inverseDocumentsFrequencies) = loadFromFile
98     ()
99
100     println(termsFrequencies.size + " docs")
101     println(termsFrequencies.head)
102     println(inverseDocumentsFrequencies)
103 }
104 }
```


E | Search.scala

```
1 package fr.thomasrobert.inforetrieval
2
3 object Search {
4
5     val (tfs, idf) = Indexer.loadFromFile()
6     val docsTfIdfs = Indexer.computeTfIdfs(tfs, idf)
7
8     def computeSimilarities(queryTfIdf: Indexer.TFIDF, docsTfIdfs: Indexer.
9         TFIDFs): Map[String, Double] = {
10         docsTfIdfs mapValues {
11             //  $\forall \text{ doc, compute } \cos(q, d) = \langle q, d \rangle = \sum q_i * d_i$ 
12             dis  $\Rightarrow$  queryTfIdf map {
13                 case (qiTerm, qi)  $\Rightarrow$  qi * dis.getOrElse(qiTerm, 0.0d)
14             } reduce (_ + _)
15         }
16
17     def orderResults(similarities: Map[String, Double]): Seq[(String, Double)] = {
18         similarities.toSeq.sortBy(_._2)
19     }
20
21     def computeQuery(query: String): Seq[(String, Double)] = {
22         val queryTf = Indexer.getPageFrequencies(query)
23         val queryTfIdf = Indexer.computeTfIdf(queryTf, idf)
24
25         orderResults(computeSimilarities(queryTfIdf, docsTfIdfs))
26     }
27
28     // main
29     def main(args: Array[String]): Unit = {
30         var continue = true
31         while (continue) {
32             val ln = scala.io.StdIn.readLine()
33             if (ln != null)
34                 computeQuery(ln) take 25 map println
35             else
36                 continue = false
37         }
38     }
39
40 }
```