CSC3010 Information Retrieval
Tutorial 05

# Lucene 02

Peter CY YAU

# Agenda

- Parsing

- Tokenization

- Core Analysis

# Parsing

- Applications that build their search capabilities upon Lucene may support documents in various formats – HTML, XML, PDF, Word – just to name a few. Lucene does not care about the Parsing of these and other document formats, and it is the responsibility of the application using Lucene to use an appropriate Parser to convert the original format into plain text before passing that plain text to Lucene.

- Plain text passed to Lucene for indexing goes through a process generally called tokenization. Tokenization is the process of breaking input text into small indexing elements – tokens. The way input text is broken into tokens heavily influences how people will then be able to search for that text. For instance, sentences beginnings and endings can be identified to provide for more accurate phrase and proximity searches (though sentence identification is not provided by Lucene).

- In some cases simply breaking the input text into tokens is not enough – a deeper *Analysis* may be needed. Lucene includes both pre- and post-tokenization analysis facilities.

- Pre-tokenization analysis can include (but is not limited to) stripping HTML markup, and transforming or removing text matching arbitrary patterns or sets of fixed strings.

# Tokenization (Cont)

- Stemming – Replacing words with their stems. For instance with English stemming "bikes" is replaced with "bike"; now query "bike" can find both documents containing "bike" and those containing "bikes".

- Stop Words Filtering – Common words like "the", "and" and "a" rarely add any value to a search. Removing them shrinks the index size and increases performance. It may also reduce some "noise" and actually improve search quality.

- Text Normalization – Stripping accents and other character markings can make for better searching.

- Synonym Expansion – Adding in synonyms at the same token position as the current word can mean better matching when users search with words in the synonym set.

- The analysis package provides the mechanism to convert Strings and Readers into tokens that can be indexed by Lucene. There are four main classes in the package from which all analysis processes are derived.

- Analyzer – An Analyzer is responsible for supplying a TokenStream which can be consumed by the indexing and searching processes. See below for more information on implementing your own Analyzer. Most of the time, you can use an anonymous subclass of Analyzer.

- CharFilter – CharFilter extends Reader to transform the text before it is tokenized, while providing corrected character offsets to account for these modifications. This capability allows highlighting to function over the original text when indexed tokens are created from CharFilter-modified text with offsets that are not the same as those in the original text. Tokenizer.setReader(java.io.Reader) accept CharFilters. CharFilters may be chained to perform multiple pre-tokenization modifications.

- Tokenizer – A Tokenizer is a TokenStream and is responsible for breaking up incoming text into tokens. In many cases, an Analyzer will use a Tokenizer as the first step in the analysis process. However, to modify text prior to tokenization, use a CharFilter subclass (see above).

- TokenFilter – A TokenFilter is a TokenStream and is responsible for modifying tokens that have been created by the Tokenizer. Common modifications performed by a TokenFilter are: deletion, stemming, synonym injection, and case folding. Not all Analyzers require TokenFilters.

# Invoking the Analyzer

```java
1  Version matchVersion = Version.LUCENE_XY; // Substitute desired Lucene
   version for XY
2      Analyzer analyzer = new StandardAnalyzer(matchVersion); // or any other
   analyzer
3      TokenStream ts = analyzer.tokenStream("myfield", new StringReader("some
   text goes here"));
4      // The Analyzer class will construct the Tokenizer, TokenFilter(s), and
   CharFilter(s),
5      //   and pass the resulting Reader to the Tokenizer.
6      OffsetAttribute offsetAtt = ts.addAttribute(OffsetAttribute.class);
7
8      try {
9        ts.reset(); // Resets this stream to the beginning. (Required)
10       while (ts.incrementToken()) {
11         // Use AttributeSource.reflectAsString(boolean)
12         // for token stream debugging.
13         System.out.println("token: " + ts.reflectAsString(true));
14
15         System.out.println("token start offset: " +
   offsetAtt.startOffset());
16         System.out.println("  token end offset: " + offsetAtt.endOffset());
17       }
18       ts.end();   // Perform end-of-stream operations, e.g. set the final
   offset.
19     } finally {
20       ts.close(); // Release resources associated with this stream.
21     }
```

# Implementing your own Analyzer and Analysis Components

```java
1 document.add(new Field("f","first ends",...);
2     document.add(new Field("f","starts two",...);
3     indexWriter.addDocument(document);
4
5 ...
6
7 Version matchVersion = Version.LUCENE_XY; // Substitute desired Lucene
  version for XY
8   Analyzer myAnalyzer = new StandardAnalyzer(matchVersion) {
9     public int getPositionIncrementGap(String fieldName) {
10      return 10;
11    }
12  };
13
```

# Token Position Increments

```java
1  public TokenStream tokenStream(final String fieldName, Reader reader) {
2      final TokenStream ts = someAnalyzer.tokenStream(fieldName, reader);
3      TokenStream res = new TokenStream() {
4        CharTermAttribute termAtt = addAttribute(CharTermAttribute.class);
5        PositionIncrementAttribute posIncrAtt =
   addAttribute(PositionIncrementAttribute.class);
6
7        public boolean incrementToken() throws IOException {
8          int extraIncrement = 0;
9          while (true) {
10           boolean hasNext = ts.incrementToken();
11           if (hasNext) {
12             if (stopWords.contains(termAtt.toString())) {
13               extraIncrement += posIncrAtt.getPositionIncrement(); // filter
   this word
14               continue;
15             }
16             if (extraIncrement > 0) {
17
   posIncrAtt.setPositionIncrement(posIncrAtt.getPositionIncrement()+extraIncrem
   ent);
18             }
19           }
20           return hasNext;
21         }
22       }
23     };
24     return res;
25   }
```

# References

Materials in this slides are obtained (included partially modify) from the following resources, including websites, content providers, printed materials, etc.:

- https://lucene.apache.org/

- https://www.lucenetutorial.com/