

# Store Pre-processing Text Result Methods

Text Pre-processing has become a very important task in modern Computer Science fields such as Artificial Intelligent and Machine Learning. For example, in 2013, E. Haddi et. al have used multiples Text Pre-processing methods to analyze the sentiment of online movie reviews [1]. However, it is a very expensive process, and it becomes extremely expensive when the dataset becomes bigger. For example, In one of the research on extracting Named Entities from text, the data scientists have to process millions of online news articles from multiple online publishers including the Los Angeles Times, Reuters, and the New York Times [2].

There are three main components of Pre-processing including *Tokenization*, *Normalization*, and *Substitution* [3] and their output results can be reused many times. However, those tasks require very high computational resources. In order to minimize the computing cost, we can save the output results to file and then we can reload them while needed. So that saving results to a file becomes a mandatory feature. There are some methods to address this problem and the methods' technique can be variable which depends on the data structure of the results. The methods can be put into two main categories which are *Human-readable* (text-based) and *Non-human-readable* (binary).

Since all the results of Text Pre-processing in this assignment is stored in Objects so the problem is narrowed down to "How to store an Java Object to file" and it can be solved using *Object Persistence Methods* [4], they are basically a process that transforms an Object to a format that can be stored on disk or database, however, we are only interested in storing in file due to the scope of this assignment.

- The first approach is to use the Build-in Java Serialization which makes use of object stream classes including *ObjectInputStream* and *ObjectOutputStream* [5] to convert and write an object to a binary file.
- The second approach is to transform the Objects to a human-readable string using the *Data Markup Language* such as XML or *Data Serialization Language* such as JSON and YAML and then those strings to a text file.
- The third approach is to convert the preprocessed data into *Comma-Separated Value* (CSV) format and then write to a text file.

I will use the third approach for my implementation, although it may not be as efficient as using Java Serialization in term of run time and space, however, it is a multi-platforms method that can be used with multiple programming languages and it is also human-readable that is easier for error checking. For each section, I will convert the Map of sanitized lines into a CSV file named `<section>.data.csv`, another index CSV named `<section>.index.csv` will store the Trie of words as well as the index tables of occurrences.

By loading data directly from the document, the *loadDocument()* method runs with the time complexity of  $O(c)$  with  $c$  is the number of characters in the documents. However, by loading the data from the stored CSV files, it only runs with the time-complexity of  $O(w)$  with  $w$  is the number of words in the document. That is a significant improvement because  $c \gg w$ . The detailed time complexity analysis is described in the following pseudo-code fragments.

## Time Complexity Analysis

Load directly from Document Pseudo Code	#Operations
<pre> <b>Indexes</b> &lt;- GetIndexes() <b>SectionMap</b> &lt;- Empty Map of Sections <b>LineNo</b> &lt;- 0 For each <b>Entry</b> in <b>Indexes</b>:     <b>Section</b> &lt;- New Section     <b>Trie</b> &lt;- New Trie     <b>Lines</b> &lt;- Empty Map of Lines     <b>Files</b> &lt;- New File("document.txt")     While <b>LineNumber</b> &lt; <b>Entry.LastLine()</b>:         <b>Line</b> &lt;- <b>Files.ReadLine()</b>         <b>Text</b> &lt;- Sanitize(<b>Line</b>)         <b>Lines.Add(LineNo, Text)</b>         <b>Words</b> &lt;- Tokenize(<b>Text</b>)         For each <b>Word, Pos</b> in <b>Words</b>:             <b>Trie.Insert(Word)</b>             <b>Table</b> &lt;- <b>Trie.Get(Word)</b>             If <b>Table</b> == Null:                 <b>Table</b> &lt;- New IndexTable                 <b>Trie.Set(Word, Table)</b>             <b>Table.Add(Pos)</b>         <b>Section.Trie</b> &lt;- <b>Trie</b>         <b>Section.Lines</b> &lt;- <b>Lines</b>         <b>SectionMap.Add(Entry.Title(), Section)</b> </pre>	<pre> #sections (s) 1 1 #sections (s) 1 1 1 1 1 #lines in sect. (n) 1 #chars in line (c) 1 (expected) 1 (expected) #chars in line (m) #chars in word (c) #chars in word (c) 1 1 #chars in word (c) 1 (expected) 1 1 1 1 (expected) </pre>
Accumulated: $s + s(n(m + m(3c))) \approx snmc$ . Since $c$ can be considered as constant, then the time complexity can be considered as $O(snm)$ which is actually the number of character in the document	

Load from Stored Result Pseudo Code	#Operations
<pre> <b>Indexes</b> &lt;- GetIndexes() <b>Sections</b> &lt;- Empty Map of Sections For each <b>Index</b> in <b>Indexes</b>:     <b>Section</b> &lt;- New Section     <b>Trie</b> &lt;- New Trie     <b>Lines</b> &lt;- Empty Map of Lines     <b>DFile</b> &lt;- New File(<b>Index.Title()</b> + ".data.csv")     While !<b>DFile.EOF</b>:         <b>LineNo, Text</b> &lt;- <b>DFile.ReadLine()</b>         <b>Lines.Add(LineNo, Text)</b>     <b>IFile</b> &lt;- New File(<b>Index.Title()</b> + ".index.csv")     <b>Word</b> &lt;- Null;     <b>Table</b> &lt;- Null;     While !<b>IFile.EOF</b>:         <b>NextWord, Pos</b> &lt;- <b>IFile.ReadLine()</b>         If <b>NextWord</b> != <b>Word</b>:             <b>Trie.Insert(Word)</b>             <b>Table</b> &lt;- New IndexTable             <b>Trie.Set(Word, Table)</b>         If <b>Table</b> != Null:             <b>Table.Add(Pos)</b>         <b>Section.Trie</b> &lt;- <b>Trie</b>         <b>Section.Lines</b> &lt;- <b>Lines</b>         <b>SectionMap.Add(Entry.Title(), Section)</b> </pre>	<pre> #sections (s) 1 #sections (s) 1 1 1 1 1 #lines in sect. (n) 2 1 (expected) 1 1 1 #words in sect. (w) 2 1 #chars in word (c) 1 #chars in word (c) 1 1 (expected) 1 1 1 1 </pre>
Accumulated: $s + s(n + w(2c)) \approx swc$ . Since $c$ can be considered as constant, then the time complexity can be considered as $O(sw)$ which is actually the number of words in the document	

## References

- [1] E. Haddi, X. Liu and Y. Shi, "The Role of Text Pre-processing in Sentiment Analysis", *Procedia Computer Science*, vol. 17, pp. 26-32, 2013.
- [2] Y. Shinyama and S. Sekine, "Named entity discovery using comparable news articles," in *Proceedings of the 20th international conference on Computational Linguistics*, 2004, p. 848: Association for Computational Linguistics.
- [3] M. Mayo, "A General Approach to Preprocessing Text Data", *Kdnuggets.com*, 2018. [Online]. Available: <https://www.kdnuggets.com/2017/12/general-approach-preprocessing-text-data.html>. [Accessed: 12- Oct- 2018].
- [4] C. JMTauro, R. Kumar Sahai and S. Rani A., "Object Persistence Techniques - A Study of Approaches, Benefits, Limits and Challenges", *International Journal of Computer Applications*, vol. 85, no. 5, pp. 19-27, 2014.
- [5] "The Java™ Tutorials - Object Streams", *docs.oracle.com*, 2018. [Online]. Available: <https://docs.oracle.com/javase/tutorial/essential/io/objectstreams.html>. [Accessed: 12- Oct- 2018].
- [6] S. man and A. Utama Siahaan, "Huffman Text Compression Technique", *International Journal of Computer Science and Engineering*, vol. 3, no. 8, pp. 103-108, 2016.