MEMORIAL UNIVERSITY OF NEWFOUNDLAND FACULTY OF SCIENCE DEPARTMENT OF COMPUTER SCIENCE

Text Summarization Computer Science 4750

Jacob House Noah Gallant

Friday, November 30th, 2018

Contents

| 1 | Top | ic | | | | | | | | | | | | | | | | | | | | 3 |
|---|----------------|----------|------------|--|--|--|--|--|--|--|--|--|--|---|--|--|--|--|--|--|--|---|
| | | | Research . | | | | | | | | | | | | | | | | | | | |
| | Our | Implem | entation | | | | | | | | | | | | | | | | | | | 3 |
| | 3.1 | Design | | | | | | | | | | | | | | | | | | | | 3 |
| | 3.2 | Classes | | | | | | | | | | | | | | | | | | | | 3 |
| | | 3.2.1 | Text Class | | | | | | | | | | | | | | | | | | | 3 |
| | | 3.2.2 | Node Class | | | | | | | | | | | | | | | | | | | 4 |
| | | 3.2.3 | Edge Class | | | | | | | | | | | | | | | | | | | 4 |
| 4 | 4 Optimization | | | | | | | | | | | | | 4 | | | | | | | | |
| Α | Den | nonstrat | ion | | | | | | | | | | | | | | | | | | | 6 |

1 Topic

2 Background Research

3 Our Implementation

The implementation chosen mainly followed the graph method. Three primary files were created containing the three main classes of the software. These classes were the node, the edge, and the paragraph class. Combined together these are called by a summarizer Python file which has the role of opening the summary file and ensuring the correct format is used before performing the summarization.

- 3.1 Design
- 3.2 Classes

3.2.1 Text Class

The role of the Text class is to perform all operations on the text. This includes any preprocessing required, the initial object creation, determining further relations, as well as using the determined relationship information to form a reasonably accurate and concise summary.

When first initializing the text, the text to be summarized must be passed into the constructor of the Text class. This text is then saved as an object variable for future use, and proceeds to undergo pre-processing. Pre-processing is done by the preProcessing() method, and is tasked with splitting the text into sentences.

PREPROCESSING() Pre-processing is a crucial step in achieving a high quality summary. Since this is the first modification made to the text, the method in which pre-processing is conducted can have a significant bearing on the quality of the summary. The goal of the pre-processing in this application is primarily to split the text into sentences. This is done by considering all sentence terminating symbols including '.', '!', and '?'. This also considers if a sentence contains a quote, and any other situation in which a terminating symbol should be ignored.

The secondary operation of the pre-processing in our application is to replace special Unicode characters with their plain text equivalent. Many common text editors use Unicode characters as opposed to the plain text characters for many text symbols. Some of these symbols include quotation marks, accents, en dashes and em dashes.

Processing Sentences

Pre-processing removes special characters and splits the text into individual sentence. Once split, each text is processed by creating a node object for each sentence. This node

contains the original sentence, a list of all the words in a sentence which have been passed through a lemmatizer, a list of edges which are connected to the node, as well as the sentence number to keep track of the location of the sentence in the original text. While each sentence Node object is being created, edge objects are being created to connect each node whose sentences are connected in the text. These adjacency edges are the first edges to be created in the graph.

CREATING THE DICTIONARY

After all the sentence nodes have been created, and a list of words in each sentence node has been processed, The Text object will then proceed to loop through each sentence node to create a complete dictionary of all the words found in the text. This dictionary will not only store the available words, but will also store all the sentence node in which the word is contained in. The goal of creating this dictionary is to decrease the complexity in creating the relations between sentences. By creating a dictionary which contains words and the nodes they are contained within, we must simply loop through all the words and create relations when a word is contained within more than 1 node.

CREATING EDGES

Using the now created dictionary of words, we proceeded to loop through each word individually, creating an edge between each of the sentence Node objects which have common words. Each sentence Node object will thus contain both proximity Edge objects, and Edge objects associated with common words.

CREATING A SUMMARY

3.2.2 Node Class

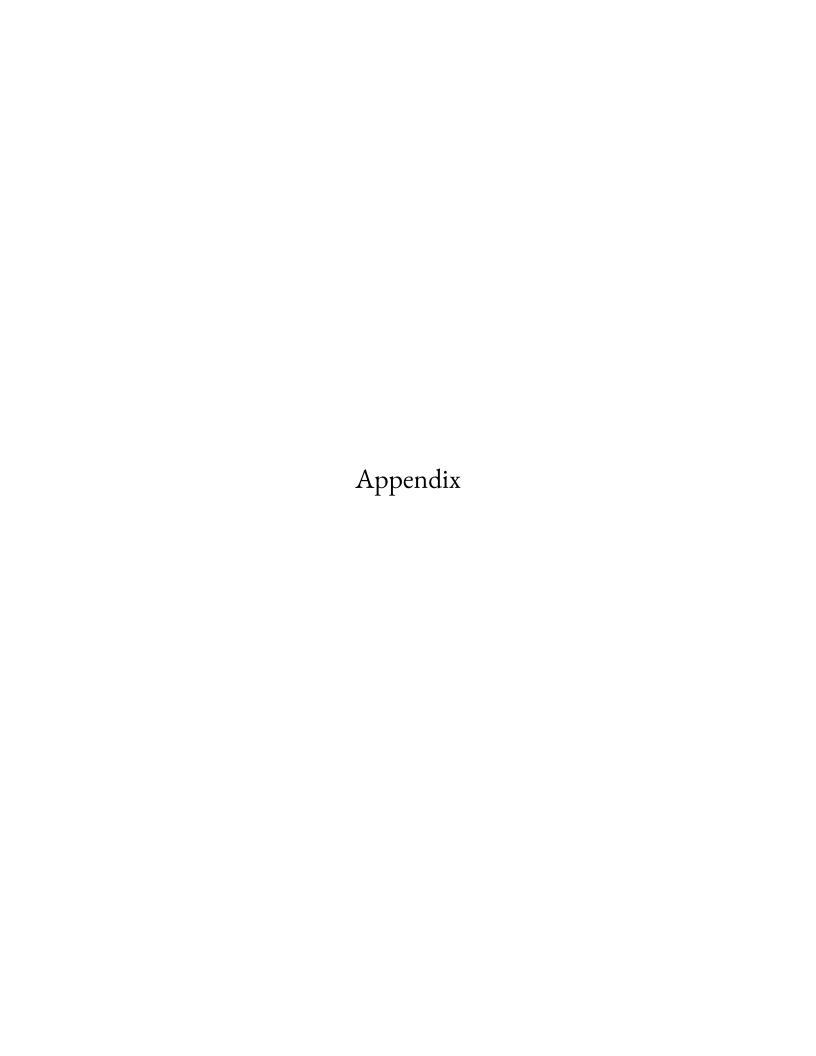
The Node class is the applications representation of a sentence. The Node () constructor takes in a sentence as a string, and will immediately proceed to split the string into words which are saved in a list.

To reduce the complexity of the code, instead of performing lemmatizing in the preprocessing stage as might be expected, it was decided to perform this step at the same time the words in the sentence were being divided. This was decided so as to prevent accessing a word more than was required.

The final task of the preProcessing() method is to use a lemmatizer to remove word endings. When comparing words later in the application to determine relevant sentences, this will improve the results as it will more accurately represent similar words, rather than counting the same words with identical endings as different words and therefore misrepresenting the relations in the text.

3.2.3 Edge Class

4 Optimization



A Demonstration