

Text - Summarization using Conceptual Statistical Extraction

Vivek Patani

School of Informatics and Computing
Indiana University
Bloomington, Indiana – 47405
Email: vpatani@uemail.iu.edu

Jinsu Kim

School of Informatics and Computing
Indiana University
Bloomington, Indiana – 47405
Email: jk247@uemail.iu.edu

Abstract—There have been increasing number of researches conducted in the domain of text summarization using different types of contexts using different methods. In this paper, we attempt to summarize text with important facts eliminating redundant or irrelevant information using extractive method. In this paper we focus mostly on extracting entities and concepts which occur more commonly and are meaningful. Our approach is a statistical one, wherein we weigh each concept in terms of relations and how frequently they occur.

I. INTRODUCTION

With increasing amount of textual information available, it becomes difficult to find and read important text. Thus, text summarization would be very useful in that it produces important and concise summary and help readers to save time without having to read entire text.

Summary can be generated by **extractive** and **abstractive** methods. Extractive methods create short summary with each words, POS from whole sentences without modification of words. Abstractive summarization is more complex way, which makes new sentences summarize by paraphrasing sentences of source documents.

The paper talks about the selective picking up of entities in the text. We, in this paper weigh **entities** heavily over other concepts. An entity may be a person or any proper noun or moreover any noun. We along with this also collect the verbs associated to them. However, a deep dive into the subject tells us that not all verbs represent an important concept related to the paragraph, hence the term selective.

We divide the approach into 5 sections:

- Sentencifying
- Anaphora Resolution
- Dependency Parsing
- Statistical Inference
- Entity Recognition

We also use a lot of other ideas, but this briefly describes a good overview of the subject we are presenting here. Our focus is on English for different languages hold different semantic syntax, such like languages like Japanese and Chinese have unambiguous sentence ending markers.

A. Relevant Works

These are a few relevant studies & works we came across:

- Several studies summarized text with a machine learning approach. Chuang, W. T., & Yang, J. (2000, July) extracted sentence segments with automatic text summarize of on patent text data. The study used three supervised learning algorithm, C4.5, Bayesian, DistAI, to train and extract important sentence segments, and all of the methods performed well in generating good summary.
- Other than text summarization that have been widely studied, there are some studies explored web-page summarization. Shen et al. (2004, July) studied web-page classification algorithm on LookSmart Web dictionary, web summarization. Nave Bayes and Support Vector Machine were used for classification algorithm to build baseline system. The study showed improvement in using the summarization-based classification algorithm and ensemble classifier than in using text-based algorithm.
- In the Rusu et al. (2009) study, they analyzed text on Reuters newswire article and extracted triplet information in the base of semantic graph. After triplet generation, they used anaphora resolution for entities and semantic normalization and evaluated document with automatic summarization.

II. APPROACH AND METHODOLOGY

A. Sentencifying

To begin with we approach the data with splitting of the set of text into simple sentences. We do this to make relations and dependencies much easier and flexible to detect. The idea here is to break the sentence down at periods with a few reservations such as places like *Dr.*, *Mr.*, *St.*, *etc.* 47% of the articles written in the Wall street journal represent abbreviations^[2]. This problem in NLP is also referred to as Sentence Boundary Disambiguation (SBD), which is the science of detection of the beginning and the ending of a sentence.

B. Anaphora Resolution

The problem of Anaphora Resolution is also termed as pronoun resolution. Let us understand this problem with a brief example, Sentence: Tim walks his dog, he does so every single day. Your mind subconsciously decodes that in this context he is Tim and not the dog. The computer is not so intelligent as you

so it needs a few instruction. There are two approaches to solve this problem, namely Eliminate Constraints & Weighing Preferences. The former uses a set of rules which should be obeyed by both the anaphor and referent and latter method weighs each of the eligible referents and picks the one with the highest score^[3].

C. Dependency Parsing

This is the most important step in summarizing text. Dependency grammar is based on the relationships between words. $A \rightarrow B$, means B is dependent on A or A governs B. Dependency Parser is simply used to realise the predicate, subject and object in a particular context. We can understand a simple selection of dependency by expressing the structure in head – dependent relation, functional categories and some structural categories. Phrases can also be represented by simply the same way. Here is an example of a sentence provided with their dependencies parsed in *fig 1*.

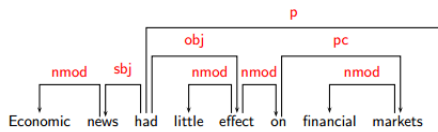


Fig. 1. Dependency Relation

D. Statistical Inference

Once we process the dependencies, we lucidly look for the most important term in all of the text. The term "most important" can be interpreted in many ways. For us, we simply pick the term that has occurred in most number of concepts (A concept is a tuple of (predicate, subject, object)). We pick that term and look for important concepts and then in those concepts we look for relevant concepts. The advantage of this is that we are looking for the best possible term and finding all the things related to it and pick the best one. The downside is that once a *best term* is picked, only relevant concepts are searched, so our scope is limited to a single *best term* and can cause a single point of failure. We also look for the best predicates, subjects and objects to find out more information in regards to the *best term*.

E. Entity Recognition

This is the last part and is optional. Named Entity Recognition (NER) is a very interesting property to know for a certain concept or term. It helps us realise whether if the term or concept is a person and often times also helps us realise the gender. This can be helpful in the selection of templates as we can have a customised one for concepts representing people.

III. PRACTICAL IMPLEMENTATION

We have implemented the above mentioned idea on Github. You can access the readme and follow the details to run the project.

Before you begin, we need to satisfy a few prerequisites, such as:

- Python 3.5+
- Java 1.8+
- MariaDB 10+ (Any MySQL running on default port should do)
- Docker

Let's understand how it works:

- 1) Data to be processed is present in a file and the file is read as input (UTF-8) format.
- 2) Primarily we break down the text into sentences by running it through a sentencer (Stanford NLP provides this).
- 3) We then take each sentence and run it through a coreference relation generator. This is done to resolve any anaphoras that come across.
- 4) Once we complete that we take each anaphora and replace each of those with the actual entity.
- 5) Following that we push it through a Neural Net Dependency Parser (using Stanford NLP). This is the most important step and will cause your quality up or down depending upon how complex sentence structure is.
- 6) Followed by that we save all the data in a MariaDB database. We store it in a triple store form (predicate, subject, object).
- 7) We then perform statistical inferences based on the data stored.
- 8) We find the best term, frequency of each predicate, subject and object and then look for concepts accordingly.
- 9) Best individual concept and best individual term is giving more preference as compared to other things.
- 10) The best term decides the next step, i.e. Templating.
- 11) The best term is processed through a Named Entity Recognition Tool.
- 12) That decides upon what template to pick up and finally that is template we display our best results in.

IV. CONCLUSION

After testing on three different datasets (or inputs) we were successfully able to find the gist of two topics. The topics were:

- A Short story on the Ali Baba and 40 thieves (It would recognise Baba as a human character).
- Microsoft and its progress over the years (It realised Microsoft as an entity).
- Importance of English in our lives (This was not recognised, as the difference between general and special concepts was relatively small).

Our model is really focused on a single idea or concept. We focus on one *best term or concept* and use that to find all relevant relations. This is a good idea if the paragraph or text represents exactly one idea, but that is not always the case. The limitations that this implementation holds is, we only focus on a single idea and do not look for any other ideas on it. Sometimes there maybe multiple ideas as important as the

previous one and we may still miss it completely revealing an incomplete summary of the text. Also another caveat is that the word disambiguation of different senses may mean different in different contexts, but we do not distinguish that while searching for relevant concepts and hence once again may lead to an incomplete summary.

In the future work, we can try text summarization by humans, and evaluate how much our work improved compared to summarization by human. We also can extend our implementation scope to different types of documents such as web page document, which would be very beneficial because more and more text is available in web. Text summarization of different characteristics of multiple documents would also help readers. Another way to improve the approach is by adding a probabilistic model that would make the triples generated independent of the length of the sentence. The problem lies in the way we select the triples and ideas.

ACKNOWLEDGMENT

We would like to thank our Professor, Dr. Damir Cavar and the teaching assistant, Atreyee Mukherjee for giving us this opportunity and bringing out the best in us.

REFERENCES

- [1] Chuang, W. T., & Yang, J. (2000, July). Extracting sentence segments for text summarization: a machine learning approach. In Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval (pp. 152-159). ACM.
- [2] Dali, L., & Fortuna, B. (2008). Triplet extraction from sentences using svm. Proceedings of SiKDD, 2008.
- [3] H. Kopka and P. W. Daly, *A Guide to L^AT_EX*, 3rd ed. Harlow, England: Addison-Wesley, 1999.
- [4] Rusu, D., Fortuna, B., Grobelnik, M., & Mladenic, D. (2009). Semantic Graphs Derived From Triplets with Application in Document Summarization. *Informatica (Slovenia)*, 33(3), 357-362.
- [5] Shen, D., Chen, Z., Yang, Q., Zeng, H. J., Zhang, B., Lu, Y., & Ma, W. Y. (2004, July). Web-page classification through summarization. In Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval (pp. 242-249). ACM.
- [6] Sandra Kubler and Marcus Dickinson
- [7] Stanford NLP Text
- [8] Texting Mine Online