

A Visual Investigation of the Kronos Incident

Connie XIA Yi Jing
Singapore Management University
connie.xia.2020@mitb.smu.edu.sg

Nikitha BANDA
Singapore Management University
nikithab.2020@mitb.smu.edu.sg

TAN Kar Yee
Singapore Management University
karyee.tan.2020@mitb.smu.edu.sg

ABSTRACT

This is the abstract.

It consists of two paragraphs.

1. INTRODUCTION

The fictitious Kronos Incident saw the disappearance of several employees from the Tethys-based GASTech in January 2014 after a successful initial public offering (IPO) of the company. Given that GASTech has not been very environmentally friendly in its operations of a natural gas production site in the island country of Kronos, it was suspected that a Kronos-based organisation (POK) is involved in the disappearance of the employees, as a form of retaliation. In order to have a better idea on what exactly transpired to lead to the vanishing of the GASTech employees, we will be applying visual analytical techniques on the datasets provided.

This study will be handling visualisations on newspaper articles, employee records and emails, call center reports and microblog tweets before structuring them into an interactive web application. Users can then investigate the application and understand more about GASTech's reputation. Furthermore, one can navigate around the app to find out how certain events unfolded on the incident day itself.

2. MOTIVATION AND OBJECTIVES

The motivation behind this study is to look into analytical techniques to visualise large chunks of text data effectively using R Studio. By doing so, we are able to better understand the relationships among people and organisations of importance, as well as see how multiple events of high consequences unfolded in Abila on the incident day.

This interactive Shiny app aims to provide information on:

1. Media portrayal of GASTech over the years

2. Relationships among GASTech, POK, the APA and Government
3. Meaningful event reports during the incident day
4. Risks identified during the incident day and their corresponding locations

3. REVIEW & CRITICS OF PAST WORKS

This study is based on the VAST Challenge 2021, which in turn is adapted from a similar VAST Challenge in 2014. Literature review is conducted on the previous VAST Challenge 2014 submissions to look at the analytical techniques used to solve the challenge back then, even though the exact questions were slightly different. While useful, some of the techniques adopted have certain areas that can be further improved.

3.1 Text Visualisations

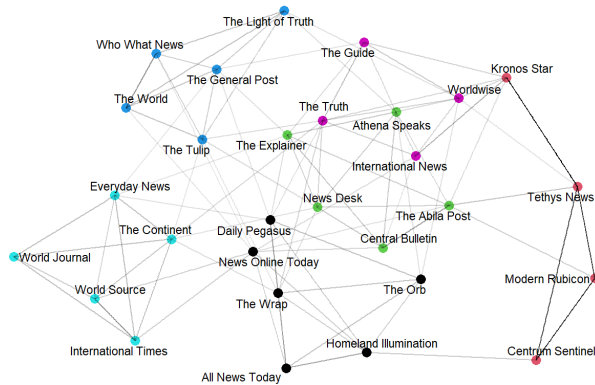
A study conducted by Peking University (2014) on Mini Challenge 1 presented their text analysis in a form of a timeline to showcase different events occurring between January 20 – 21. Articles in the form of text boxes were layered over the timeline for comparison. While it showcased all the news reporting of different events occurring over the two-day time period, it might be difficult for a user to interpret the main concepts of those articles. Hence, a better alternative might be to utilise a word cloud function to pull out key words of the articles for view and interpretation. In addition, interactive comparisons of different newsgroups can also be performed, giving the user flexibility to choose the newsgroups they are interested in to view and evaluate.

While word clouds are generally useful in identifying topic content for a broad overview as shown in the study performed by Tianjing University (2014) on Mini Challenge 3, their results might be less consistent and harder to make sense of due to the presence of spam data. Hence, to be able to distinguish important events from typical chatter, TF-IDF would be a better statistical tool to use.

3.2 Network Graphs

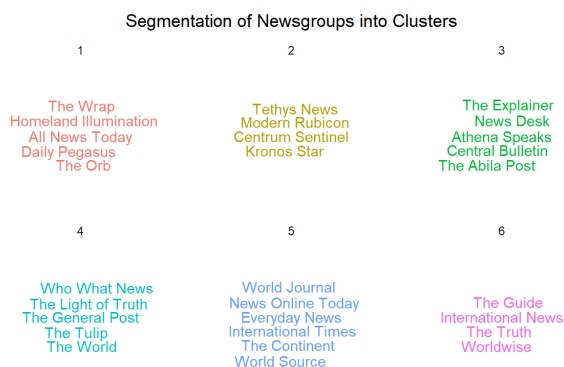
Network graphs are a good visualisation tool to establish the relationships between different parties of interest. By and large, network graphs would be densely populated with nodes and edges if there are numerous parties involved. Yet, this brings about an issue of overcrowding and overlaps of texts, making the entire visualisation looks cluttered, as seen in Fig. x.

tween words. The usage of nodes and links is able to explicit show the relationships between each neighbouring text [cite]. In this case, the first node set is the words found in news articles and the second node set is the newsgroups themselves. That way, a network can be created where newsgroups are connected by their use of the same words, as shown in Fig. X.



The node color corresponds to the text communities, with the same colour indicating a strong relationship between its components. In this way, clusters are formed to segment the newsgroups with similar characteristics in terms of the types of words used in the news articles.

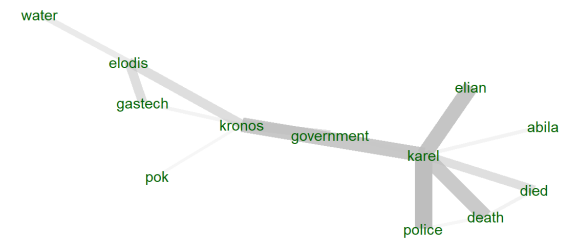
To delve into the cluster segments of different newsgroup, we visualised each segmentation as follows:



Upon visualising the components in each cluster, text plot visualisations are conducted to pull out the word cooccurrences between word-pairs, so as to determine the context and content of each clusters.

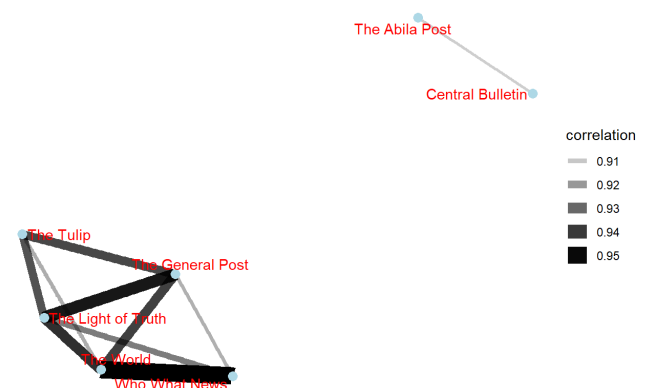
Term cooccurrences

showing Cluster 1



The main R packages used are **textnet**, **ggwordcloud**, **tidytext**, **udpipe** and **textplot**. **Textnet** is currently the only R package available to implement text network techniques in R. To display the word cloud by cluster, **ggwordcloud** was utilised. **Tidytext** helps to convert text into a format that is visualisable with the use of 'unnest_tokens' function. The **udpipe** package provides language-agnostic tokenisation, tagging, lemmatisation and dependency parsing of raw text, which is an essential part in natural language processing (NLP). Lastly, to plot data as a text plot, we will be needing the **textplot** package.

Correlation Graphs Correlation graphs are plotted to determine the correlation between different newsgroups. From this, we are able to determine which newsgroups might be highly related in terms of their reports of certain events over the years. The correlation values are obtained using the widely used Pearson method.



The R packages used are **widyr** and **ggraph**. **Widyr** is able to cast a tidy dataset into a wide matrix, performs an operation such as computing the correlation on it, and then re-tidies the result. 'pairwise_cor' function is found in this package. **ggraph** is then used to plot the relationship between different newsgroups based on their correlation values.

TF-IDF As mentioned in the **Review & Critics of Past Works** section, word cloud alone might not be very useful to visualise a collection of microblog message due to its consistency issue. As such, TF-IDF approach is more helpful when trying to pick out specific key events/topics of relevance.

4.2.2 Network Graphs

Relationships of GASTech Employees

Email Flow

Target Employee

4.2.3 Geospatial mapping

5. APPLICATION INSIGHTS

Nullam semper imperdiet orci, at lacinia est aliquet et. Sed justo nibh, aliquet et velit at, pharetra consequat velit. Nullam nec ligula sagittis, adipiscing nisl sed, varius massa. Mauris quam ante, aliquet a nunc et, faucibus imperdiet libero. Suspendisse odio tortor, bibendum vel semper sit amet, euismod ac ante. Nunc nec dignissim turpis, ac blandit massa. Donec auctor massa ac vestibulum aliquam. Fusce auctor dictum lobortis. Vivamus tortor augue, convallis quis augue sit amet, laoreet tristique quam. Donec id volutpat orci. Suspendisse at mi vel elit accumsan porta ac ut diam. Nulla ut dapibus quam.

6. CONCLUSION AND FUTURE WORK

The study

Further developments to the application can be implemented such that the application is able to take out other text data such as text corpus

References

- [1] Fenner, M. 2012. One-click science marketing. *Nature Materials*. 11, 4 (Mar. 2012), 261–263.
- [2] Meier, R. 2012. *Professional Android 4 Application Development*. John Wiley & Sons, Inc.