

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

Investigation of crimes often is dealt with vast amount of data that has to be manually viewed to obtain clues and evidence. An optimal visualization of data is helpful in reducing the manual work and providing more accurate insights and patterns of the crime. This project aims to apply visual analytics into investigation of the kidnapping of several GASTech company employees. Using datasets that capture the newspaper articles, social media data and the information of the employees, a visualization application is created using R Shiny as the platform. The insights derived can help draw patterns and obtain clues on metrics such as relation of employees within the company, the reputation of the company in Kronos and the unfolding of events on the day of incident.

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 (IEEE 2021). 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 (Wang, C. et al. 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 (Yang S. and Jiang, X. 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 (Wang, C. et al. 2014).

One way to overcome this issue of cluttering will be to divide the network graph into sub graphs. When the graph is divided, the density of the visualisation will be reduced, with the readability enhanced.

3.3 Geospatial maps

Static geospatial maps tends to show many different points of interests, which might overload the user with content. Hence, to enhance the use of geospatial maps, we intend to include in the interactivity function so as to allow users to click and explore different points as desired.

4. DESIGN FRAMEWORK

This application makes use of the open-source R language to conduct visual analysis. The application design considerations are as follows:

- Utilise standard R packages to create reproducible text and visual analysis
- Utilise the embedded Shiny Web Application in R to translate the codes into a webpage for users' ease of understanding
- Provides interactivity functions for users to navigate through the app to discover trends and insights

The design of the application will consists of five major tabs for navigation at the top panel.

- i) Introduction - Describes the main purpose of our application
- ii) History of GASTech - Sub tabs of Text Analysis and Network Graph to respectively discover insights from the newspaper articles and employee relationships
- iii) Message Stream Exploration - Explores tweets from different users on the incident day itself to sieve out notable keywords from otherwise spam information
- iv) Risk Level Timeline tab - Further split into Call Centre Reports and Microblog Messages to detect how the public risk levels changed over time on the incident day
- v) Message Stream Geomap - Shows locations where those key messages appeared

Coupled with the interactivity aspect to choose different options from dropdown boxes to sliders for exploration, the combination of these views gives the user the flexibility and

autonomy to navigate through the app to find out more information regarding the Kronos Incident.

To facilitate users in their exploration around the application, we have also provided a user guide for their reference.

4.1 Data Preparation

The data provided are a set of historical news reports, data on GASTech employees and email headers from two weeks of internal GASTech company email. They are collated in the following tables as shown below. Detailed steps on preparation of the data can be found at Data Preparation-Part 1 and Detailed steps on preparation of the data can be found at Data Preparation- Part2

File Name	About
cleanArticles.csv	File containing various newsgroups and their news article
cleanEmail.csv	Email headers from two weeks
cleanEmployee.csv	Employees' data
csv-1700-1830.csv	Message stream - from 1700H to 1830H
csv-1831-2000.csv	Message stream - from 1831H to 2000H
csv-2001-2131.csv	Message stream - from 2001H to 2135H
df_unmap_labeled.csv	Manually geo-labelled data for key call centre reports

Table 1 List of datasets used

These datasets are then loaded into R and further data preparation is conducted to clean these raw data mainly with the tidyverse package.

4.2 Analytic Techniques used in Shiny App

A variety of standard R packages were used to conduct text and visual analysis on the datasets to draw useful insights. The following techniques are used for analysis.

4.2.1 Text Analysis

Our application will build several text analysis outputs to break down large chunks of text data from newspaper articles and tweets into understandable visualisation for users to view:

Comparison Cloud

Comparison clouds are used to visualise the similarity and differences of important words used by different newsgroups. Wordclouds are able to extract out keyword metadata and the frequency of the appearance of a particular word in the articles is determined by the size of the word in the word-cloud (Halvey M.J. and Keane, M.T. 2007). This type of visualisation is useful in the quick pickup of prominent terms to determine the significance of those words. In turn, based on the keywords extracted out, we are also able to get a sense of the sentiments and attitude of these newsgroups regarding their articles about GASTech. Our application will have the option for choosing between viewing keywords from "Content of Articles" or "Titles of Articles." Moreover, users can choose the newsgroups they would have to like to compare from a dropdown list of all the newsgroups available.

The packages used to build this visualisation are **tidytext**, **tm** and **wordcloud**. **Tidytext** is used to convert text into a format that is visualisable with the use of 'unnest_tokens' function. **tm** has built in functions such as 'removeStop-words' to help in the removal of unimportant words. Lastly, the comparison cloud will be developed using the **wordcloud**

package.

Textnet and Text Plot

Textnet can be used to help in structuring text data, by showcasing relationships between neighbouring text data of relevance (Trigg, R.H. and Weiser, M. 1986). Hence, textnet can be used for clustering of newsgroups. To represent the data in cluster format, the data has to be converted to two columns. The first column would be the set of words found in news articles and the second is the name of newsgroups themselves. That way, a network can be created where newsgroups are connected by their use of the same words. The thickness of the edge connecting any two nodes represents how similar the two nodes are.

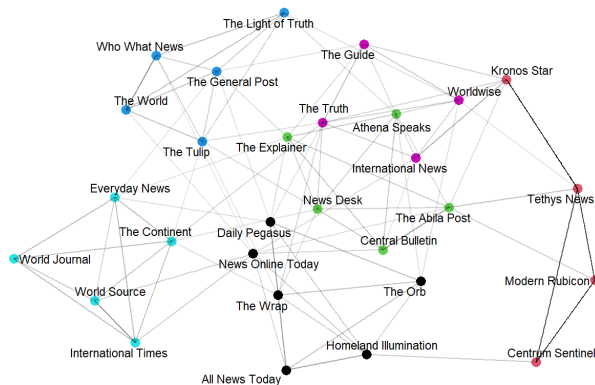


Fig. 1 Relationship between different Newsgroups

The color of the node corresponds to the text communities, with the same color 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 visualize each segmentation as follows:

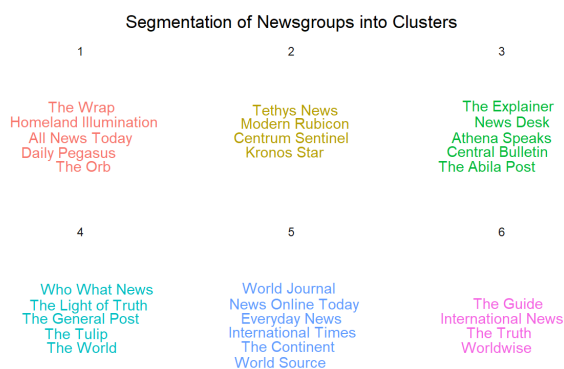


Fig. 2 Segmentation of Newsgroups into Cluster

Upon visualizing the components in each cluster, text plot visualizations are conducted to pull out the word co-occurrences between word-pairs, so as to determine the context and content of each clusters.

Term cooccurrences

showing Cluster 1

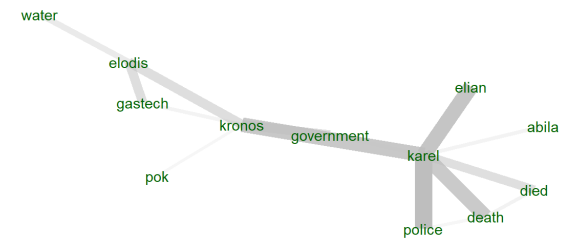


Fig. 3 Text plot showing Word Co-occurrences for Cluster 1

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 (Glen, S.).

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.

Term Frequency-Inverse Document Frequency (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. This approach aims to measure the importance of a word is to a document in a collection/corpus of documents, by including the inverse document frequency (IDF) to balance the term frequency (TF) used in wordclouds (Silge, J. and Robinson, D. 2021). Hence, we will view each hour of tweets as a document and all 5 hours of the dataset as the corpus, to determine the key events that dominate each hour. With that, TF-IDF can be a powerful heuristic to sieve out the most important events occurring during each period from the chatter spread throughout the dataset. Both unigrams and bigrams are plotted to look at both the singular word and word-pairs respectively.

5.2 Context of Clusters

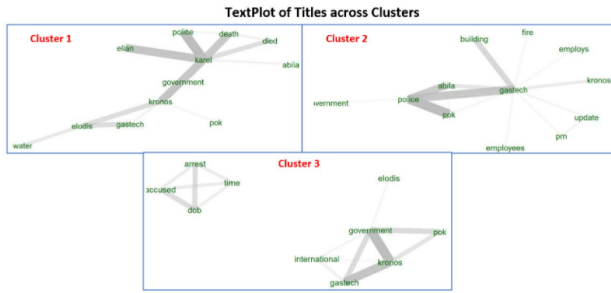


Fig. 6

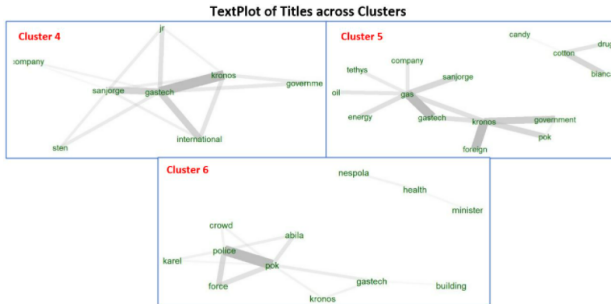


Fig. 7

Clustering separated all of the newsgroups into six different clusters with the following characteristics:

- Cluster 1 - Relationship between police, government and Elian Karel
- Cluster 2 - POK in general and mentions of GASTech employees
- Cluster 3 - Relationship between police, government and POK
- Cluster 4 - CEO of GASTech company and company's investment plans
- Cluster 5 - Nature of GASTech company as well as relationship between government and POK
- Cluster 6 - Tension between police and POK + Health situation in Kronos

5.3 Correlation between Articles

Looking at the highly correlated articles (with $r > 0.9$), we can observe that the highest pair of correlated articles are "The Light of Truth" and "The General Post." These two nodes have the thickest edge connecting between them, suggesting a high similarity between the words used by the two newsgroups. Referencing to Fig. 2 above, we can observe that these two newsgroups also falls in the same cluster - Cluster 4.

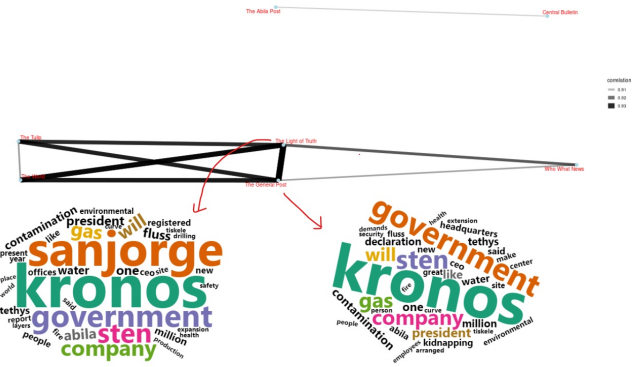


Fig. 8

Fig. 8 also shows the highest frequency of words used in the articles published by these two newsgroups. It appears that both of these published articles mostly contain words like "Kronos," "gas," "government," "Sten" and "company." From this, we can understand that articles published by these 2 newsgroups are mainly about the GASTech company and its operations in Kronos.

5.4 Email Flow and Distribution

Based on the segregation between work related and non-work related emails, we are able to derive the following insights

5.4.1 Work Related Emails

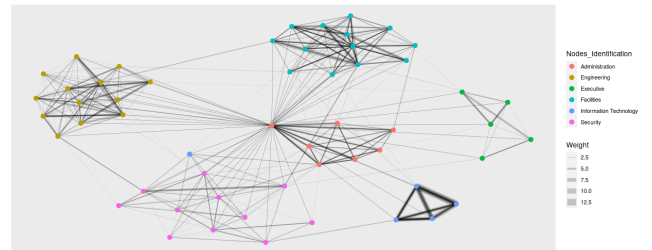


Fig. 9

- Heavy weight assigned to Information Technology, eliminating any suspicious activity occurring in this department
- Employees in Engineering department seem to have more communication among themselves than the other departments
- One particular IT employee who seems to be communicating with the Security and Administration department but not with their own department, which raises suspicion

5.4.2 Non-work Related Emails

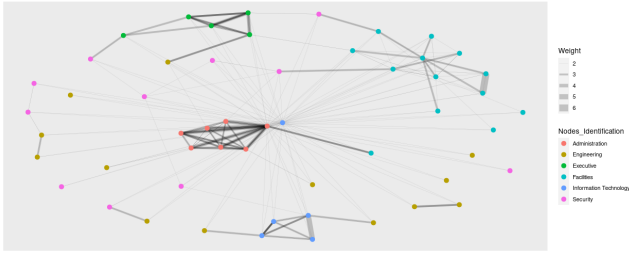


Fig. 10

- Executive Department has more non-work related emails than work-related emails
- One IT employee seems to be the most common receiver and sender of non-work related emails
- Heavy transmission of unofficial emails among Administration Department

5.5 Microblog Tweets Analysis

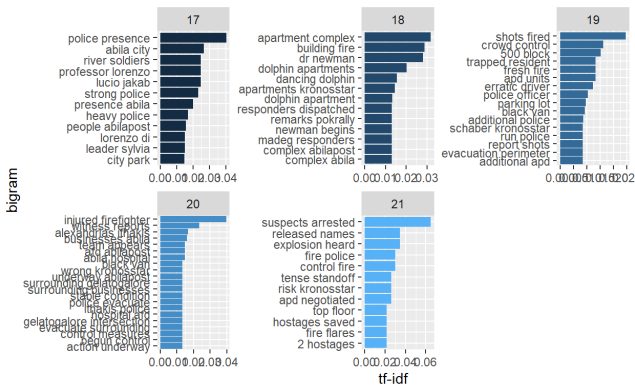


Fig. 11

From this bigram plotted, we are able to extract certain important information as follows:

- 1800H+: The fire was likely related to a building, the apartment complex of Dancing Dolphon
- 1900H+: Shots were fired, with a trapped resident and black van being other points of interest.
- 2000H+: A firefighter was injured. A phrase “alexandrias ithakis” appears (this is the cross-street where the standoff happened). Several businesses appear to be affected as well.
- 2100H+: Besides the arrests already spotted in the wordcloud, an alarming event of explosions heard also appears.

5.6 Risk Levels From Call Center Reports

The basic timeline of major event identified over the course of the evening with its respective risk ratings (on a scale of 1-5, 5 being highest) are as follows:

1. The first big event was a fire detected at about 1830H (exact time: 1842H). An ambulance and fire truck was dispatched to the fire. Crowd control was also requested. The risk appears to have elevated, but appeared under control (risk level 3).

2. The next big event happened from 1915H: what appeared as a vehicle accident (exact time: 1919H) turned out to be a rogue black van running amok in the crowd. It was pursued by police units. Risk level would have gone up to 4, due to the uncontrollable nature of the van.
3. From 1930H, at least one officer was down because of the event (exact time: 1941H). A dire emergency was called, and additional support was requested. Risk level would have been at 5 (given that the situation was termed “dire”).
4. Nearing the end after the black van and fire incident seemed to have died down, from 2045H there were several reports of crime scene investigations, with continued suspicious reports. Risk level would be at 3 (moderate).

5.7 Map of Key Messages with Geolocation

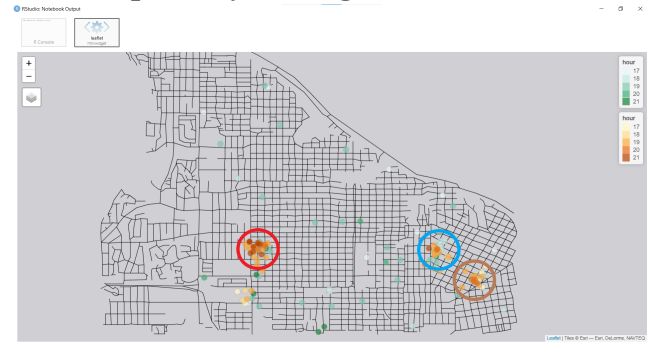


Fig. 12

From Fig. 12, we are able to observe that there are three locations with high density of message origination, suggesting that major events are occurring around those location clusters,

The red cluster appears to be the place with highest risk and potential consequences, due to the high number of negative tweets and several serious call center reports spread over time. Both types of messages mention the standoff and police casualties.

The blue cluster also has many call center reports, but lesser negative tweets. Closer inspection by clicking on the points show that these tweets are mainly about the fire which appeared to have been under control, and the call center reports are mostly about dispatch requests.

The brown cluster seems to be about the explosion that occurred late. Closer inspection shows that many of such messages are sent by a spam account footfingers.

6. CONCLUSION AND FUTURE WORK

In conclusion, this study attempts to draw out insights related to the disappearance of the GASTech employees from a collection of artifacts and emails related to the GASTech company. The analysis in our web application is scoped to investigate how events unfolded on the day of incident and prove the reputation of GASTech since its establishment.

Further developments to the application can be implemented such that users can upload their own datasets, like text corpus into the application to conduct text and visual analysis on new datasets.

7. REFERENCES

references: - id: vast2021 title: "VAST Challenge 2021" author: IEEE type: website publisher: IEEE URL: '<https://vast-challenge.github.io/2021/>' issued: year: 2021 - id: pekingvast2014 title: "VAST Challenge 2014 Mini Challenge 1" author: - family: Wang given: Chenglong - family: Cheng given: Siming - family: Miao given: Zhengjie - family: Yuan given: Xiaoru type: website publisher: Peking University URL: '<https://www.cs.umd.edu/hcil/varepository/VAST%20Challenge%202014/challenges/MC1%20-%20Disappearance%20at%20GASTech/entries/Peking%20University/>' issued: year: 2014 - id: tianvast2014 title: "VAST Challenge 2014 Mini Challenge 3" author: - family: Yang given: Siqi - family: Jiang given: Xinyi type: website publisher: Tianjing University URL: '<https://www.cs.umd.edu/hcil/varepository/VAST%20Challenge%202014/challenges/MC3%20-%20Real-Time,%20Streaming%20Social%20Media/entries/Tianjin%20University/>' issued: year: 2014 - id: halveykeane title: "An Assessment of Tag Presentation Techniques" author: - family: "Halvey" given: "Martin J." - family: "Keane" given: "Mark T." container-title: "WWW '07: Proceedings of the 16th international conference on World Wide Web" URL: '<https://doi.org/10.1145/1242572.1242826>' DOI: 10.1145/1242572.1242826 publisher: Association for Computing Machinery page: 1313-1314 type: article-journal issued: year: 2007 month: 5 - id: triggweiser title: "TEXTNET: A Network-Based Approach to Text Handling" author: - family: "Trigg" given: "Randall H." - family: Weiser given: Mark container-title: "ACM Transactions on Information Systems" Volume: 11 URL: '<https://doi.org/10.1145/5401.5402>' DOI: 10.1145/5401.5402 issue: 1 publisher: Association for Computing Machinery page: 1-23 type: article-journal issued: year: 1986 month: 1 - id: glen title: "Correlation Coefficient: Simple Definition, Formula, Easy Steps" author: - family: Glen given: Stephanie type: website URL: '<https://www.statisticshowto.com/probability-and-statistics/correlation-coefficient-formula/>' publisher: Statistics How To - id: silge title: Analysing word and document frequency: tf-idf author: - family: Silge given: Julia - family: David given: Robinson container-title: "Text Mining with R: A Tidy Approach" URL: '<https://www.tidytextmining.com/tfidf.html>' type: book issued: year: 2021 month: 6 ...