# Who is suspicious?

Joanna Lian
Singapore Management University
joannalianyz@gmail.com

Ng Yen Ngee
Singapore Management University
ngyenngee@hotmaill.com

## ABSTRACT
Criminal profiling is used as a technique used to identify the perpetrator of a crime by identifying the personality and behavioral characteristics of the offender based upon an analysis of the crime committed. The sensational and dramatic elements of profiling is often portrayed in movies, television shows, and books but how does it fare as an investigative tool against crime in real life? Well, more often than not, there has been doubts within the crime-fighting community about the viability of profiling as an investigative tool. Why? This is because there are views hat a criminal profile only gives a broad indication of the type of person who may have committed the crime. It does not indicate a specific individual who happens to fit the profile. The profiler is therefore unable to say whether it is more probable than not that a specific offender did, in fact, commit the crime. Nevertheless, criminal profiling has proven to be helpful in solving some criminal cases.

Definitely though, more research has to be done before criminal profiling can truly become a useful part of the criminal investigation process. This is a fascinating topic and hence, through our project, we aim to use a range of data visualisation techniques to develop an R Shiny application to maximise the insights obtained from studying the different aspects of an individual's daily life e.g.credit card transactions and email communications in order to successfully profile an individual. The analytics and design choices made in the development of the application and initial findings and future work are discussed.

## 1. INTRODUCTION
The goal of the annual Institute of Electrical and Electronics Engineers (IEEE) Visual Analytics Science and Technology (VAST) Challenge is to advance the field of visual analytics through competition. The 2021 IEEE VAST Challenge brings back a classic challenge from 2014 to see how approaches and techniques have developed since the original release of the challenge. The background of the challenge is as below:

In January, 2014, the leaders of GAStech are celebrating their new-found fortune as a result of the initial public offering of their very successful company. In the midst of this celebration, several employees of GAStech go missing. An organization known as the Protectors of Kronos (POK) is suspected in the disappearance, but things may not be what they seem. It appears that certain employees of GAStech may be involved in the disappearance.

Aligning with our intent to create an application that aims to aid the profiling process, our application will serve as an interactive tool for users to visually investigate and identify which employee profile is indicative of suspicious behaviour. This paper documents our approach to designing and developing the interactive application targeted at crime investigators. This introduction is followed in Section 2 by an explanation of our motivation and objectives. Section 3 details the data used and methodology selected. Section 4 provides a visual overview of the final product. Section 5 concludes the report and offers ideas for further development.

## 2. MOTIVATION & OBJECTIVES
This project was motivated by our findings that there were suspicious activities within GAStech itself which were worth investigating. This project aims to create a data analytics applications to visualize these suspicious activities and relationships for users to judge, who exactly are the suspicious people in GASTech. On a larger scale, this could assist in profiling which can be used for generic crime investigations.

Through our application, users can start to decipher for themselves what defines suspicious behaviour and how can we make use of everyday data to raise red flags on suspicious behaviours. This project aims to understand better the individuals and organizations that are involved in this situation. We do this by exploring the following:

- Conducting exploration data analysis and inferential data analysis on
  - Credit card expenditure data
  - Email headers data
  - Employee categories
- Delivering a R-Shiny app that achieve the following through an interactive user interface design:

– Identifying any anomalous or suspicious behavior.
– Identifying formal (work-related) or informal (non-work related) relationships.

– Discover any associations based on common interest given in the data.

– Discover relationships between CC expenditure, email headers and employee records.
– Decide who are the suspicious GASTech employees * Obtaining a holistic profile on these suspicious employees.

# 3. METHODOLOGY

Our methodology is as below:

- Data preparation using dplyr and other R packages.

- Analysis of VAST 21 data set with background research using some of the following methods:

  – Exploratory Data Analysis (EDA) methods in R.
  – Inferential Analysis methods in R.
  – Network Analysis in R.

- Creating a R shiny dashboard showing our findings/insights and conclusions:

  – R Markdown development for functionality checks
  – R-Shiny app development for user interactivity

## 3.1 Data

We'll be using data sets from Mini-Challenge 1 [MC1] and Mini-Challenge 2 [MC2].

From [MC1], we will be using the email headers and employee records data. From [MC2], we will be extracting insights based on credit card transactions data.

We will also be joining the two data sets based on individuals to analyzing attributes by features across data sets.

## 3.2 Shiny Architecture

Development of the interactive tool was done on Shiny, and R packages were used to build interactive web apps. Shiny is widely adopted in the data analytics industry because:

- It has a framework that makes it user-friendly to collect input values from a web page, with R code written as output values back to the web page.
- The input values can be modified by the user at any time, through interaction with customisable widgets.
- The output values react to changes in input values, with the resulting outputs being reflected immediately.

The design of our Shiny web app is as follows:

1. Main layout - This will consist of two pain parts:

![IMG/1.PNG) + A top navigation bar that will help the user to navigate to the various modules that will be present in the blue area. + An input bar where user can decide for themselves who in GASTech is suspicious.
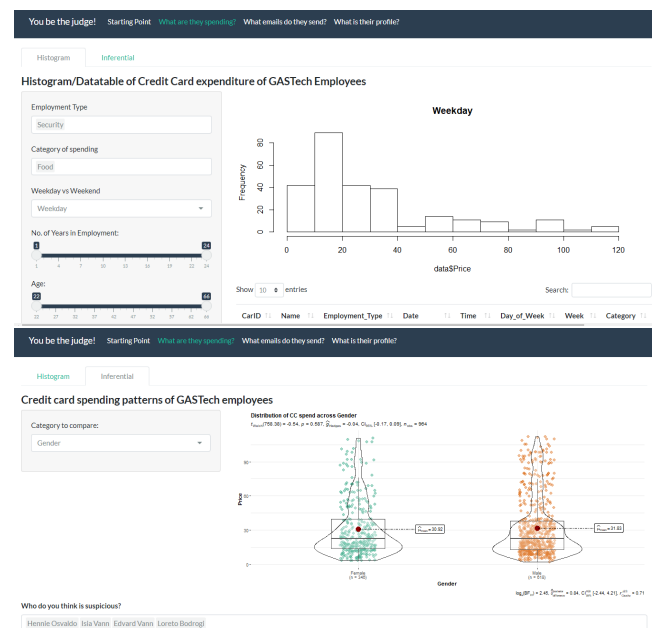
The objective of this layout is when the user switch from one tab to ano

2. First module - What are the employees spending on? We will conduct EDA and inferential analysis on the credit card transaction data. On the side panel, we will have 2 tabs - 1 for EDA (a histogram and a data table) and 1 for inferential analysis.

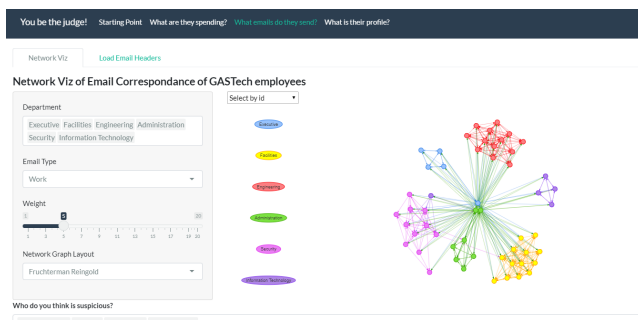   For each of the tabs, we will include options and filters for user interaction.

   The filters applicable for interactions for the histogram include:

   - Employment type (Security, Executive, Administration. . . )
   - Expenditure type (Food, Retail, Company. . . )
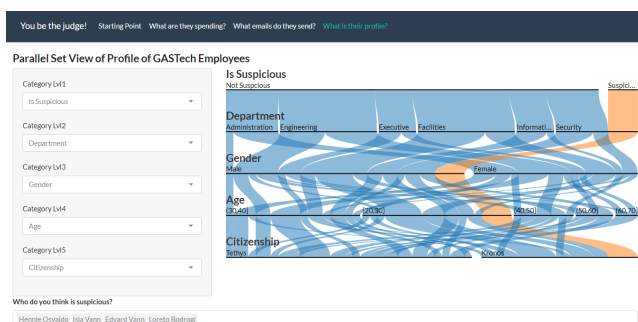   - Time category: Weekend vs Weekday
   - Years of employment
   - Age



Depending on our selection of the various filters, the visualisatins in the main panel will be automatically filtered to represent the data corresponding to the filters.

2. Second module – What emails do they send? We utilise a Network Graph for this visualisation.There will be a filter for the department that the employee belongs to, with the addition of Email Category (work related vs non-work related).

3. Third module - What is their profile? We utilise a Parallel Plot to profile the employees with 5 degrees of interaction as shown below..



To facilitate exploration by users, we have created a user guide as well...

## 3.3 Analysis Techniques

### 3.3.1 Histogram
A histogram is a plot that illustrates the underlying frequency distribution of a set of continuous data. This allows the inspection of the data for its underlying distribution (e.g., normal distribution), outliers, skewness, etc

### 3.3.2 Violin Plot
A violin plot depicts distributions of numeric data for one or more groups using density curves. The width of each curve corresponds with the approximate frequency of data points in each region. Densities are frequently accompanied by an overlaid chart type, such as box plot, to provide additional information.

### 3.3.3 Network Graph
Network graphs show interconnections between a set of entities. Each entity is represented by a Node (or vertice). Connections between nodes are represented through links (or edges). In the context of our shiny app, employees represent the nodes while the links represent the email correspondences.

### 3.3.4 Parallel Set Plot
Parallel Set Plots (ParSet) are a visualization method for the exploration of categorical data. They focus on the data frequencies instead of individual data points. The technique is built on the axis layout of parallel coordinates, with boxes representing the categories of data (e.g. gender, department,

spending patterns etc.) and parallelograms between the axes showing the relations between categories. Our app enables the user to interactively remap the data to up to five levels of categorization.
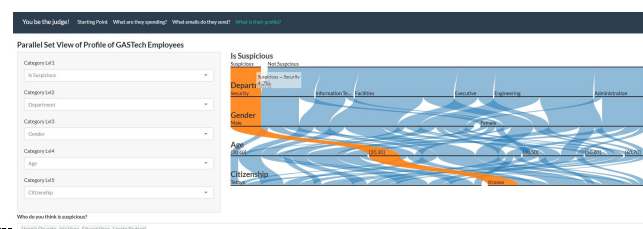
## 4. DESCRIPTION OF PRODUCT & FINDINGS

### 4.1 First Look
From our previous analysis and literature review, we know that these four GAStech employees have familial relationships with members of POK and hence are the most suspicious people in GAStech.

- Hennie Osvaldo
- Isia Vann
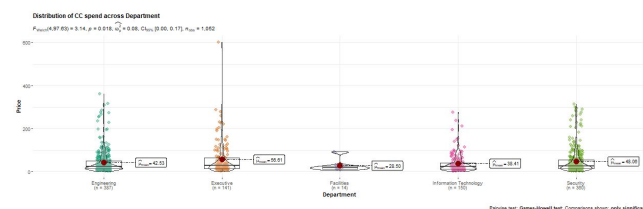- Edvard Vann
- Loreto Bodrogi

Hence, this is our starting point. We use the ParSet view to see their profile:



From here we can observe that these four people have many similarities. They are all males in their twenties and all from the security department, probably with a lot of drive and impulse to possibly plot an internal kidnapping. In our subsequent analysis, we will scrutinize closely the data with these characteristics.

### 4.2 Credit Card Expenditure
Now looking at credit card transactions, we added all of the options so that the data is not filtered and then we sort the data table by price. We can immediately see that there are a few outliers. To ensure that the outliers are worth investigating, we select the inferential tab to visualize statistical analysis.



When we exclude the outliers, the analysis tells us that there is insufficient statistical evidence to reject the null hypothesis that mean of the transactions between the departments are the same. In fact, the means of the transactions hover between \$25-60 which suggests the outliers that we spotted is worth investigating into.

Histogram for Weekday:



Histogram for Weekend:



We make the following observations: * Nils Calixto + Made a purchase of $10,000 at Frydos Autosupply n' More just two days before the kidnapping. + Other transactions are charged to the decimal place, but this particular transaction is very neat at $10,000. + if we compare to the other purchases within the same cateogory, most transactions fall below $500 + This single transaction make Nils Calixto very suspicious. * Axel Calzas + Made a purchase of $1239.41 at Albert's Fine Clothing just 3 days before the kidnapping. + Most purchases made are around $300. It could be that Axel helped pick up for 3-4 person. + We consider Axel Calzas as suspicious first but we should look into his relationships in more detail. * Sten Sanjorge Jr + Made a purchase of $600 at the Chostus Hotel just 3 nights before the kidnapping incident. + Sanjorge Jr was initially counted as one of the kidnapped victims but it turns out that he was just returning to Tethys. (add ref) + The combination of this news and this transaction makes the CEO a suspect as well.

## 4.3   Email Correspondence

We focused straight into the correspondence within the security team and compare the differences between Work and Non Work below:

Network Visualization for Security Department - Work Correspondence:



Network Visualization for Security Department - Non-Work Correspondence:



If we observe the network of Work Correspondence, we observe that the emails sent are sent rather consistently throughout the team, however if we look at the other visualization we can see quite obviously the imbalance. Isia Vann and Hennie Osvaldo seems to be in the center of all these correspondence. There are some members who only receive emails while some members where there are exchange of emails.

Let us scrutinize the non work email correspondence of Isia Vann with the rest of the organization:



We found another member closely related to our current suspicious candidates - Inga Ferro who is in correspondance with at least 3 of the original suspicious candidates which makes this person suspicious as well.

When we select Axel Calzas, we realize that this person do not have non-work correspondence, and any work related emails are department of company based. This eliminates the suspicious off and we removed the name.

## 4.4   Parset Final

Here is the final profile of our judgement of suspicious employees in GASTech:



## 5.   CONCLUSION & FURTHER IDEAS

This paper set out the development of a web application targeted at exploring how we could use visually motivated

tools to facilitate the profiling process. The project was motivated by our realization through the VAST 2021 project that there were suspicious behaviors demonstrated by some individuals which were worth investigating.

The application was developed using the Shiny architecture on R, supported with a range of statistical packages to provide users with a whole range of techniques to derive insights from the data. Definitely though, the range of techniques and ideas executed in this project is non-exhaustive and We suggest that the below can be explored in the future development of the app:

## 6. REFERENCES