# COMP90024 Assignment 2
# Big Data Analytics on the Cloud

Group 53
- Parsa Babadi Noroozi (1271605)
- Niket Singla (1288512)
- Jason Phan (1180106)
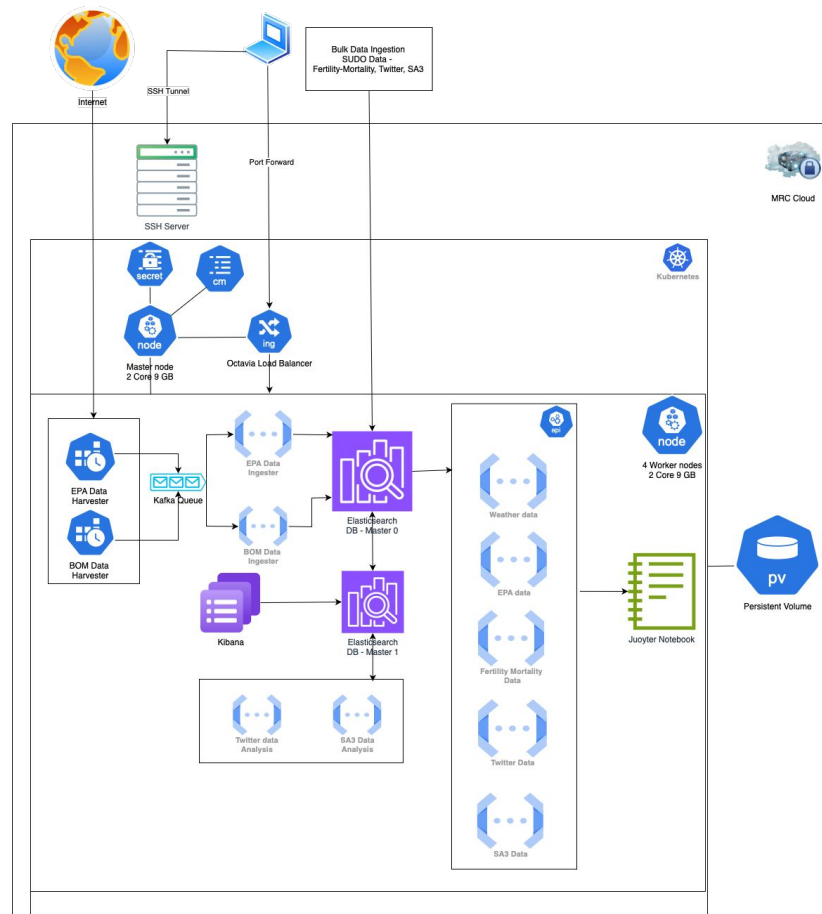- Patipan Rochanapon (1117537)
- Liam Brennan (1269948)

# Team 53 Introduction

| | | | | |
|---|---|---|---|---|
|  |  |  |  |  |
| Niket Singla 1288512 | Liam Brennan 1269948 | Parsa Babadi Noroozi 1271605 | Patipan Rochanapon 1117537 | Jason Phan 1180106 |
| <ul><li>Cluster setup</li><li>Jupyter hub setup on K8 cluster</li><li>Kafka setup</li><li>EPA & BOM data ingestion</li><li>Elasticsearch</li><li>Data analysis API</li><li>Jupyter frontend</li><li>Unit test for all ingestion function and API endpoints</li></ul> | <ul><li>Crash and health risk processing backend</li><li>Joined crash and SA2 analysis API</li><li>Crash and health risk frontend visualisation</li><li>Unit and end-to-end API endpoint testing</li></ul> | <ul><li>Crash and health risk data collection</li><li>Crash and health risk elasticsearch ingestion</li><li>Crash and health risk front end visualisation</li><li>Crash and health risk analysis</li></ul> | <ul><li>Twitter, SA3, SUDO (SA3) Ingestion</li><li>Twitter and SA3 coordinates mapping</li><li>Join Twitter with SUDO (SA3)</li><li>Twitter and SA3 Data Analysis</li></ul> | <ul><li>Twitter API routes</li><li>SA3 data API routes</li><li>Twitter and SA3 Data Analysis</li><li>Twitter sentiment and age/income/education data visualisations</li><li>Elasticsearch queries</li><li>End to End testing for API endpoints</li></ul> |

# System Architecture

**Following are the major components of our project deployed on MRC**

- Cloud infrastructure, including necessary RAM, storage and processing capacity
- Kubernetes (K8s)
- Kafka
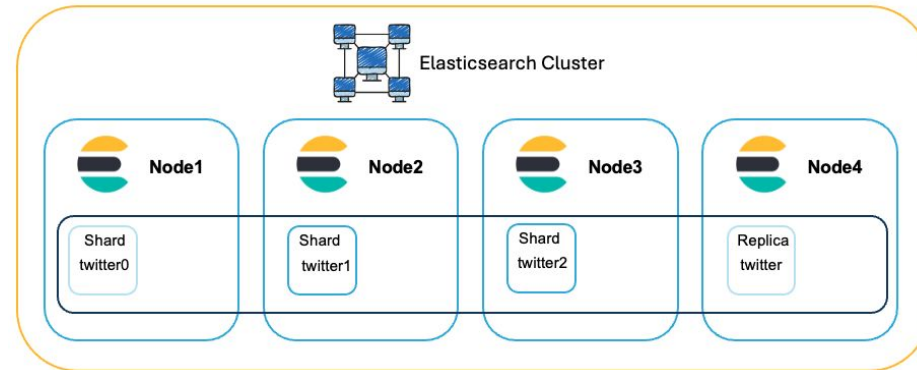- Fission
- Elasticsearch & Kibana
- Jupyter notebook

# Elasticsearch

**Indexing**: Efficient storage and retrieval of data.

**Sharding & Replication**: Distribution and replication of data across nodes for scalability and fault tolerance.

**Ingest Pipelines**: Preprocessing data before indexing using custom pipelines for parsing, enriching, and modifying documents.

**Querying and Analysis**: Retrieving and analyzing data using the Elasticsearch query DSL.

# Data Collection

## Data Collection via API

- Environment Protection Authority Victoria (EPA)
- Bureau of Meteorology (BOM)

## SUDO Dataset

- SA3 Population, Highest Education, and Average Age & Income
- SA2 Health Risk Factors
- Combined SLA11 Premature Mortality & Fertility

## Other External Datasets

- Crash Dataset
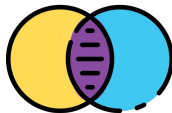- SA2 & SA3 GeoJson Dataset
- Twitter Dataset

# Backend

## Twitter and SUDO (SA3)

- Mapping Twitter coordinates with SA3 polygons.
- Utilize Elasticsearch query with "geo_distance" filter to pinpoint nearest SA3 polygons for each Twitter coordinate.
- Update Twitter index with "sa3_code_2021" to exclude non-Australian coordinates.
- Join Twitter and SUDO index by sa3_code_2021

## Road Crashes and SA2 Health Risks

- Road crash geolocations are joined using intersection queries against SA2 district geometries
- Corresponding SA2 health risk data is incorporated as part of the joining process
- The final joined dataset combines crashes, health risks, and SA2 information together
- This dataset can be analysed using grouping and metric aggregation queries to extract insights

**SA3**       **Join**       **Twitter**

**Crashes**

# Data Analysis - Research Questions

**Weather & Fertility/Mortality Rates**

- Is there a relationship between air quality and temperature?
- Is there a relationship air quality, weather and mortality and fertility?

**Crash & Alcohol Consumption**

- How are car crashes distributed geographically within Victoria?
- How do risk factors such as alcohol consumption affect crash severity?
- Explore the densities of alcohol consumption and car crashes

**Twitter**

- Twitter Use vs. Population Counts
- Is there a relationship between happiness and income levels?
- Is there a relationship between happiness and education?
- Is there a relationship between happiness and age?

# Questions

Thank you

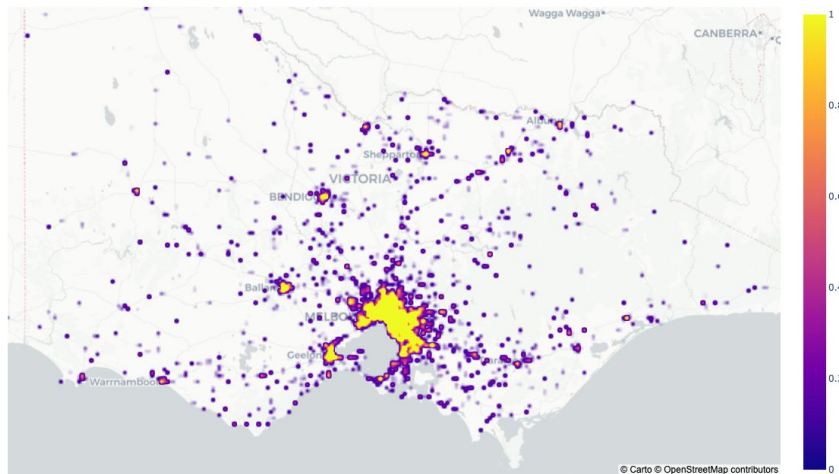# Weather, Air Quality and Fertility-Mortality Data Analysis

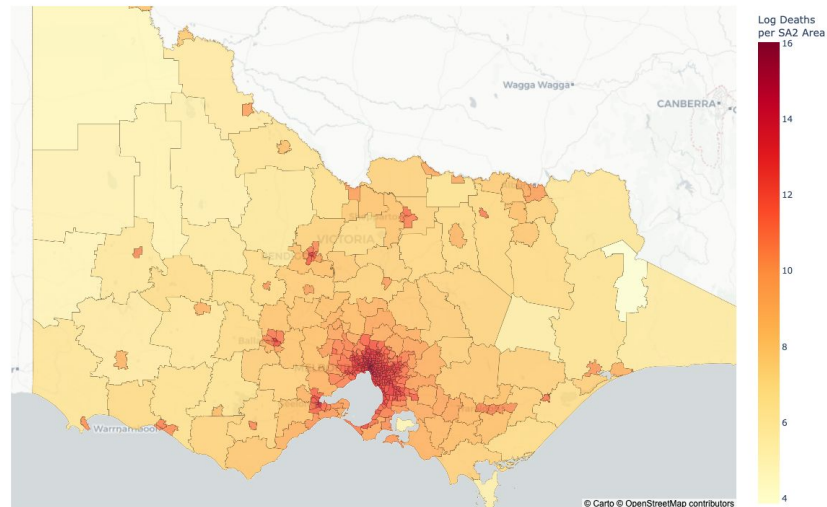# Road Crashes and SA2 Health Risks

**Exploring Crash Incidence**

Swap right map out with new one after demo practice

**Heatmap of Car Crashes in Victoria**



**Distribution of car crash locations**

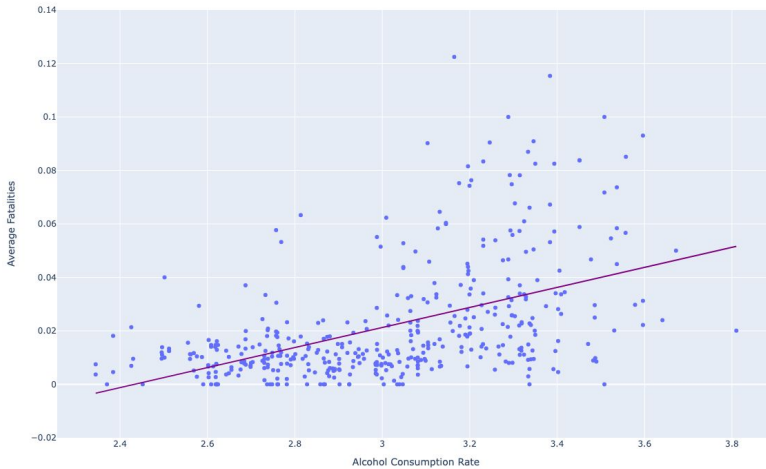**Number of Car Crashes in Victorian SA2 Districts per Region Area**



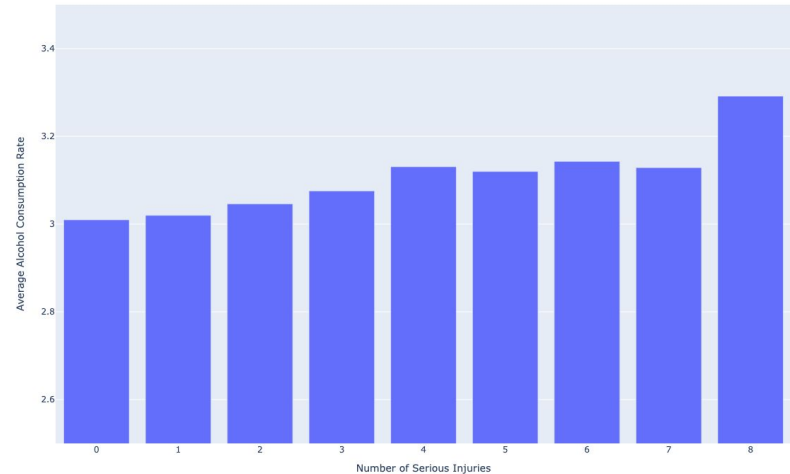**SA2 region crash incidence comparison**

# Road Crashes and SA2 Health Risks

## Exploring Health Risk Influences



Crash severity (approximated by fatalities) vs
alcohol consumption rate



Alcohol consumption rate vs
crash severity (approximated by number of serious injuries)
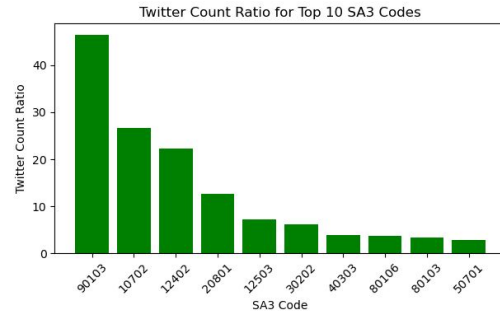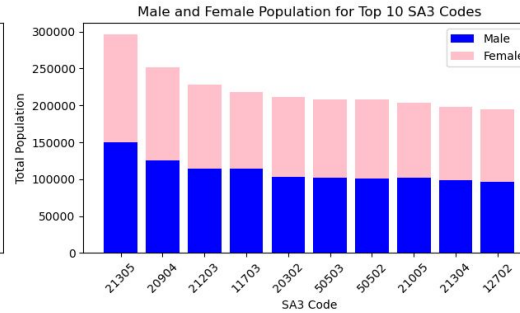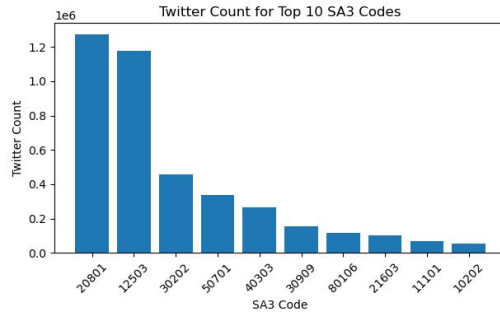
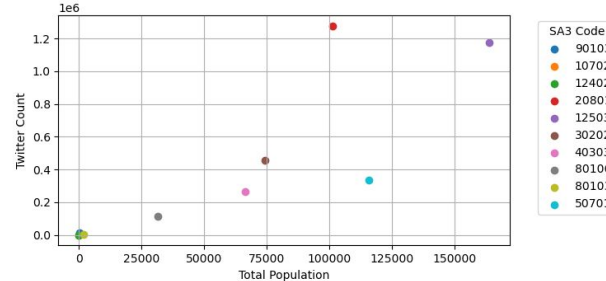# Road Crashes and SA2 Health Risks

**Exploring Health Risk Influences**



SA2 alcohol consumption counts vs car crash distribution

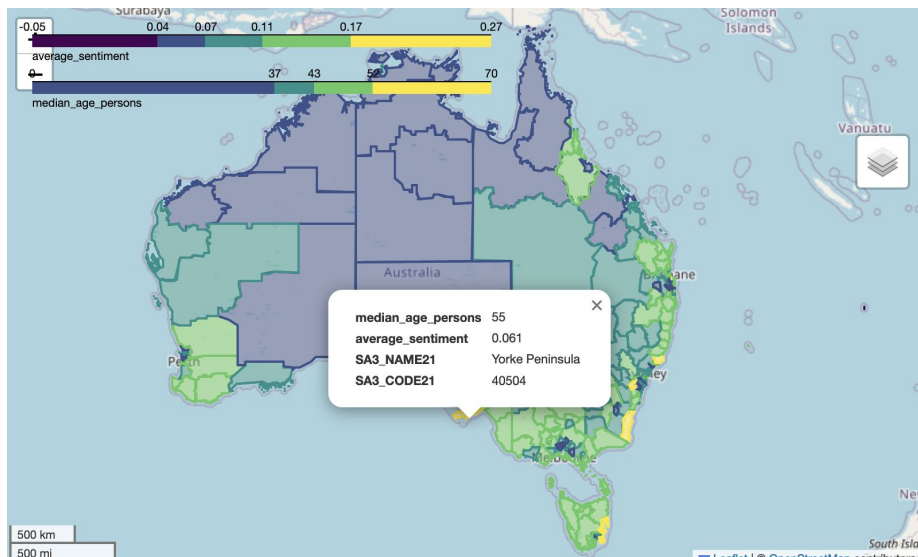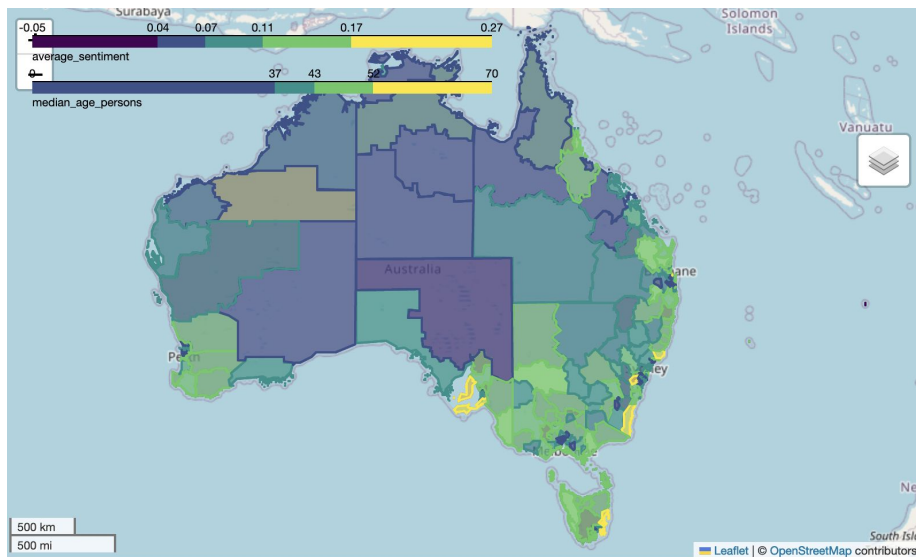# Twitter Sentiment Analysis
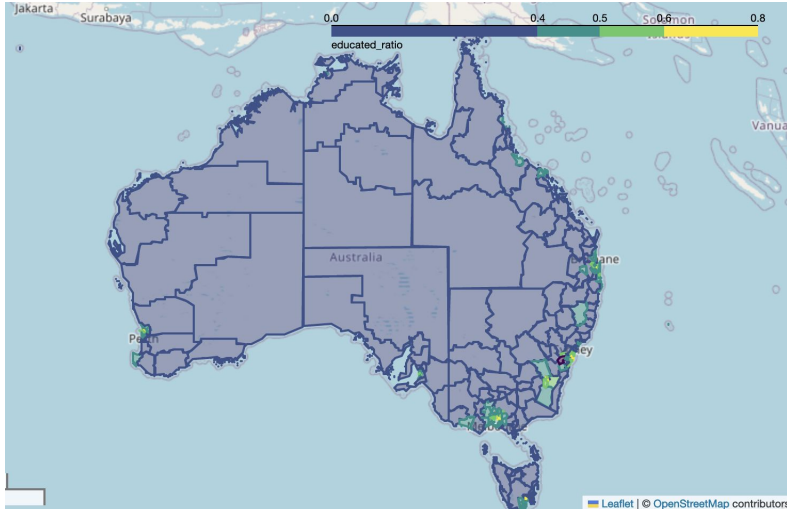
## Twitter Count VS Population

# Twitter Sentiment Analysis

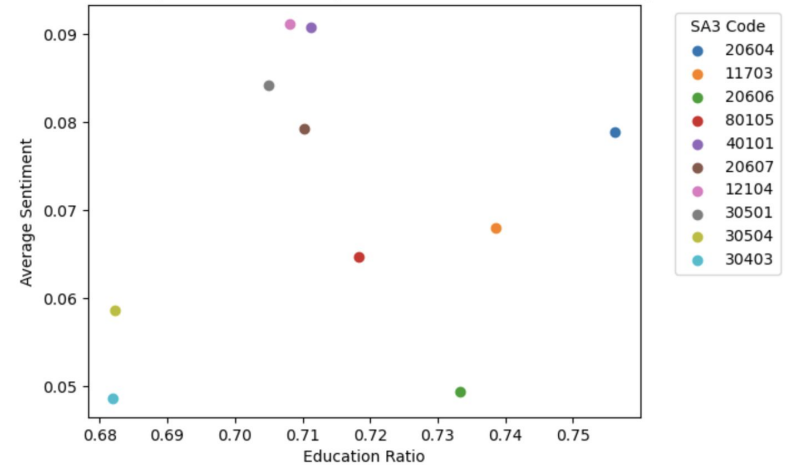**Is there a relationship between happiness and age?**

# Twitter Sentiment Analysis

**Is there a relationship between happiness and education?**





Education Ratio vs Average Sentiment for Top 10 most educated SA3 Codes by Education Ratio

# Twitter Sentiment Analysis

**Is there a relationship between happiness and income levels?**