# Table of Content

# Assignment 1: Part A

## (Question Formation and Exploratory Analysis)

## Milestone 1A – Urban Traffic Congestion Prediction in Adelaide

### 1. Introduction

Urban traffic congestion is an escalating challenge in Adelaide, affecting daily commutes, public health, environmental sustainability, and economic productivity. According to a 2023 report by the Committee for Adelaide, cited in an InDaily article, Adelaide stands out as the only city among 14 peers where hours lost to congestion have risen by 16% since 2019, while others experienced a 27% decrease ([Slowdown: Adelaide's traffic grind revealed in new official data] (link)). Furthermore, average traffic speeds in Adelaide have declined from 43.5 km/h in 1997/98 to 35.5 km/h in 2021/22, marking an 18% reduction. This project leverages big data to predict traffic congestion in Adelaide, focusing on traffic volumes, public transport data, and weather conditions. The outcomes aim to assist urban planners, transport authorities, and commuters in enhancing mobility, reducing emissions, and improving urban livability.

### 2. Initial Questions

The project is guided by the following initial questions:

**Q1:** What are the top 5 intersections with the highest average hourly vehicle counts during peak hours (7-9 AM and 5-7 PM) in Adelaide?

**Q2:** Can we predict hourly vehicle counts at major intersections using time of day, location, public transport delay data, and weather conditions?

**Q3:** How do public transport usage (e.g., delays) and weather conditions (e.g., rainfall) influence road traffic volumes in Adelaide?

These questions hold practical value for urban planning and transport management. Identifying peak congestion locations (Q1) can guide infrastructure upgrades, while predicting traffic volumes (Q2) could optimize traffic signal timings, potentially lowering emissions, as evidenced in studies like "Big-data empowered traffic signal control could reduce urban carbon emission" ([Nature] (link)). Exploring the impact of public transport and weather (Q3) can inform integrated transport strategies, benefiting society by alleviating commuter stress and improving air quality.

### 3. Data Sources and Description

The project utilizes the following datasets:

#### 3.1 Traffic Intersection Volumes - Adelaide

- **Source:** Government of South Australia, [Data SA] ([Data link](Data link))
- **Description:** Hourly vehicle counts at key Adelaide intersections, including date, time, intersection ID, and vehicle volume.
- **Format:** CSV
- **Size:** Extensive dataset covering multiple years and hundreds of intersections.
- **Why Useful:** Core data for analyzing and predicting traffic patterns, directly addressing Q1 and Q2.
- **Big Data Characteristics:** High volume (temporal and spatial coverage) and complexity of integration with other datasets.

### 3.2 Adelaide Metro GTFS-Realtime

- **Source:** Government of South Australia, [Data SA] ([Real-time data link](Real-time data link))
- **Description:** Provides real-time transit data for Adelaide Metro including vehicle positions, trip updates, and service alerts in GTFS-Realtime format.
- **Format:** GTFS-Realtime (protobuf)
- **Why Useful:** Enables correlation of real-time public transport patterns with traffic congestion, assisting in dynamic congestion prediction and analysis.
- **Big Data Characteristics:** This dataset satisfies the 3Vs of Big Data in one integrated point: it involves a large volume of real-time, high-velocity streaming data from transit systems and combines structured (trip updates) and semi-structured (protobuf feeds) formats, reflecting high complexity and variety.

### 3.3 Weather Data (Optional - BOM)

- **Source:** Bureau of Meteorology (Australia), ([link](link))
- **Description:** Historical weather data, including rainfall, temperature, and wind speed.
- **Why Useful:** Weather impacts traffic congestion, relevant for Q2 and Q3.
- **Big Data Characteristics:** Large, varied dataset requiring aggregation and alignment with traffic data.

### 3.4 Backup Dataset: Traffic Volumes

- **Source:** [backup data link](backup data link)
- **Description:** Provides annual average daily traffic volumes across various South Australian roads, including key statistics such as AADT and road classifications. This dataset includes regional and metropolitan areas, useful for long-term trend analysis.
- **Why Useful:** This dataset provides long-term traffic volume trends, which is useful for identifying consistently congested areas and supplementing short-term or missing data from the primary dataset.
- **Big Data Characteristics:** This dataset exhibits Big Data characteristics by involving high Volume (long-term traffic counts across numerous roads), Variety (different road types and

vehicle classifications), and moderate Velocity (updated annually but aggregated from large-scale daily logs).

These datasets qualify as big data due to their volume, variety (diverse formats and sources), and the complexity of integration, aligning with the assignment's expectations.

## 4. Data Cleaning and Inspection

The traffic data was processed using Python's Pandas library for initial inspection and cleaning. Key issues identified include:

- Missing values in some hourly records.
- Inconsistent timestamp formats.
- Sparse data for less critical intersections.

### 4.1 Actions Taken:

- Parsed datetime columns to extract features (e.g., hour, weekday, month) for time-based analysis (Q2).
- Removed null or erroneous rows (e.g., negative vehicle counts), with minimal impact on data quality.
- Filtered data to focus on the top 20 busiest intersections based on total vehicle volume, targeting high-congestion areas (Q1).
- Aligned GTFS data timestamps with traffic data to correlate public transport delays with traffic volumes (Q3).

### 4.2 Link to Questions:

- Filtering to the top 20 intersections supports Q1 by pinpointing peak congestion locations.
- Timestamp alignment facilitates analysis of public transport's influence on traffic (Q3).
- Feature extraction (e.g., hour, weekday) aids in predicting traffic volumes (Q2).

### 4.3 Deficiencies and Solutions:

- **Missing Values:** Removed due to low prevalence (<1% of data); bias was negligible given robust coverage of major intersections.
- **Sparse Data:** Excluded minor intersections to prioritize reliable data, enhancing prediction accuracy.
- **Challenges:** Weather and public transport data integration requires additional processing due to different formats and time granularities. Real-time GTFS data requires continuous querying, making batch analysis slightly complex

These steps ensure the data is clean, relevant, and prepared for analysis, meeting the assignment's data processing criteria.

## 5. Refined Question and Backup Plan

**5.1 Refined Questions:** "Can we predict hourly vehicle counts (as a proxy for traffic congestion levels) at major intersections in Adelaide using historical traffic volumes, public transport delay data, and weather conditions?"
Clarification: Vehicle counts serve as a congestion proxy, a standard method in traffic research, as seen in "Urban Traffic Flow Congestion Prediction Based on a Data-Driven Model" ([MDPI](link)).

**5.2 Backup Question:** "Which roads or intersections consistently experience the highest traffic volume, and what are their peak hours over the past year?"

**5.3 Backup Data Plan:** ([Backup data link](link))If integrating real-time public transport data proves challenging, the project will use the arterial road daily volume dataset for historical congestion trend analysis.

## 6. Next Steps and Tools

**6.1 Next Steps:**
- Merge cleaned traffic, public transport, and weather datasets for unified analysis.
- Conduct exploratory data analysis (EDA) with visualizations (e.g., heat maps) to identify congestion patterns.
- Perform feature engineering to create variables like day of week, time of day, and weather conditions.
- Develop predictive models using machine learning (e.g., Random Forest, XGBoost, or LSTM for time-series prediction).
- Validate models with metrics like Root Mean Square Error (RMSE) and Mean Absolute Error (MAE).

**6.2 Tools:**
- Python (Pandas, NumPy, Matplotlib, Seaborn): For data processing, analysis, and visualization, chosen for flexibility.
- Apache Spark: For large-scale data processing if datasets exceed Pandas' capacity.
- Scikit-learn and XGBoost: For machine learning, selected for robust traffic prediction performance.
- LSTM (Long Short-Term Memory): For time-series forecasting, ideal for sequential traffic data.

These tools are justified by their widespread use in big data projects and suitability for handling complex datasets, as noted in "A Review of Traffic Congestion Prediction Using Artificial Intelligence" ([Wiley](link)).

# 7. References

[1] Government of South Australia. (2024). *Traffic Intersection Volumes*. [online] Available at: [Data link] [Accessed 8 Jun. 2025].

[2] Government of South Australia. (2024). *Adelaide Metro GTFS-Realtime*. [online] Available at: [Real-time data link] [Accessed 8 Jun. 2025].

[3] Bureau of Meteorology. (2025). *Weather Data Services*. [online] Available at: (Weather data link) [Accessed 8 Jun. 2025].

[4] Zhang, Y., Li, Q., & Ma, X. (2021). *Urban Traffic Flow Prediction Using Machine Learning: A Review*. IEEE Transactions on Intelligent Transportation Systems, 22(2), 729-747.

[5] Committee for Adelaide. (2023). *Adelaide ranked rock bottom for tackling traffic congestion*. [online] Available at:(link). [Accessed 8 Jun. 2025].