

UNIVERSITY OF TECHNOLOGY SYDNEY

DOCTORAL THESIS

---

# Advanced approaches in Traffic Accident modelling

---

*Author:*

Artur Grigorev

*Supervisor:*

Dr. Adriana-Simona Mihaita

Dr. Marian-Andrei RizoIU

Dist. Prof. Fang Chen

*in the*

School of Computer Science

Faculty of Engineering and Information Technology

August 1, 2024



# Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Introduction</b>   | <b>1</b> |
| 1.1      | Research Project Summary  | 2        |
| 1.2      | Stages of the Research  | 3        |
| 1.2.1    | Stage 1. Incident Detection, Incident Duration Modeling, Incident Severity Prediction | 3        |
| 1.2.2    | Stage 2. Incident Impact modelling  | 5        |
| <b>2</b> | <b>Literature review</b>  | <b>9</b> |
| 2.1      | Introduction to Traffic accident analysis   | 10       |
| 2.1.1    | Paper structure   | 11       |
| 2.1.2    | Literature review material and the PRISMA method                                      | 11       |
| 2.1.3    | Incident duration definitions   | 14       |
| 2.2      | Data sets and data availability   | 16       |
| 2.2.1    | Characteristics of traffic incidents  | 17       |
| 2.3      | Incident duration modelling   | 18       |
| 2.3.1    | Traditional accident modelling  | 20       |
| 2.3.2    | Detection and estimation of traffic accident duration from traffic flow data          | 22       |
| 2.4      | Machine Learning in traffic accident analysis   | 22       |
| 2.4.1    | Classification and regression tasks in the incident duration prediction               | 24       |
|          | Classification and regression definitions   | 24       |
|          | Evaluation of prediction accuracy   | 25       |
| 2.4.2    | Feature importance and feature selection  | 25       |
| 2.4.3    | Interpretable models  | 27       |
| 2.4.4    | Imbalanced dataset classification   | 28       |
| 2.4.5    | Boosted models and ensembles  | 29       |
| 2.4.6    | Anomaly/outlier detection   | 29       |
| 2.4.7    | Dimensionality reduction methods  | 31       |
| 2.4.8    | Summary on the use of Machine Learning models   | 31       |
| 2.5      | Deep Learning in traffic accident analysis  | 32       |
| 2.5.1    | Spatial-temporal models for traffic incident modelling                                | 32       |
| 2.5.2    | Textual Accident Description analysis   | 34       |
| 2.6      | Conclusion  | 35       |
| 2.6.1    | Summary of challenges and gaps  | 35       |
| 2.6.2    | Future research directions in incident modelling                                      | 36       |

|          |   |           |
|----------|---|-----------|
| <b>3</b> | <b>Incident duration prediction using a bi-level machine learning framework with outlier removal and intra-extra joint optimisation</b> | <b>43</b> |
| 3.1      | INTRODUCTION  | 44        |
| 3.1.1    | Context   | 44        |
| 3.1.2    | Challenges and contribution   | 44        |
| 3.1.3    | Related works   | 46        |
| 3.2      | METHODOLOGY   | 53        |
| 3.2.1    | Classification and regression definitions   | 53        |
| 3.2.2    | Applicability of knowledge-based incident duration classification guidelines  | 55        |
| 3.2.3    | Selection of baseline machine learning models   | 55        |
| 3.2.4    | Hyper-parameter tuning through randomised search  | 56        |
| 3.2.5    | Model Performance Evaluation  | 57        |
| 3.2.6    | Regression scenarios definition   | 57        |
| 3.2.7    | Outlier removal methods (ORM)   | 58        |
| 3.2.8    | Outliers from ORM point of view   | 59        |
| 3.2.9    | Intra/Extra Joint Optimisation for ML regression prediction (IEO-ML)  | 61        |
| 3.3      | Incident classification results   | 64        |
| 3.3.1    | Binary incident classification results using varying split thresholds   | 64        |
| 3.3.2    | Classification with outlier removal   | 66        |
| 3.3.3    | Multi-class classification  | 67        |
|          | Equally split multi-class classification  | 67        |
|          | Varying multi-class classification via quantile split   | 68        |
| 3.4      | Incident duration prediction using regression: results  | 70        |
| 3.4.1    | Regression scenarios results and comparison   | 70        |
|          | Fusion framework for the incident duration prediction   | 72        |
| 3.4.2    | Outcomes and recommendations  | 75        |
| 3.4.3    | Regression results for proposed IEO-ML model  | 75        |
| 3.4.4    | Bi-level framework implementation   | 77        |
| 3.5      | Feature importance impact and evaluation  | 78        |
| 3.5.1    | Short-term vs long-term incident duration prediction feature importance   | 79        |
|          | Arterial Roads Feature Importance, Sydney Australia.  | 79        |
|          | Motorway Feature Importance, Sydney Australia.  | 79        |
|          | San Francisco Feature Importance, U.S.A.  | 81        |
| 3.6      | CONCLUSIONS   | 81        |
| <b>4</b> | <b>Traffic incident duration prediction via a deep learning framework for text description encoding</b>                                 | <b>85</b> |
| 4.1      | INTRODUCTION  | 86        |
| 4.1.1    | RELATED WORK  | 86        |
| 4.2      | CASE STUDY  | 87        |
| 4.2.1    | Incident description data set and baseline feature set  | 87        |
| 4.2.2    | Traffic flow and speed data   | 87        |



|          |   |            |
|----------|---|------------|
| 4.3      | Methodology . . . . .   | 90         |
| 4.3.1    | LSTM-ANN for the textual incident description encoding . . . . .  | 90         |
|          | The use of MSE versus cross-entropy . . . . .   | 91         |
| 4.3.2    | Artificial Neural Network Encoder for the traffic flow/speed encoding . . . . .   | 91         |
| 4.3.3    | Baseline Machine Learning model selection . . . . .   | 92         |
|          | Model performance evaluation . . . . .  | 93         |
|          | Hyper-parameter tuning for the proposed regression model . . . . .  | 93         |
| 4.3.4    | MAPE versus RMSE comparison and their non-linear relationship . . . . .   | 93         |
| 4.3.5    | Comparison to other baselines . . . . .   | 94         |
| 4.4      | RESULTS . . . . .   | 94         |
| 4.4.1    | Best model selection . . . . .  | 94         |
| 4.4.2    | Parallel coordinates for scenario setup . . . . .   | 96         |
| 4.5      | Conclusion . . . . .  | 97         |
| 4.5.1    | Word importance for severity classification . . . . .   | 97         |
| 4.5.2    | Traffic flow and traffic speed on the day of the incident . . . . .   | 99         |
| <b>5</b> | <b>Spatial-Temporal Traffic Accident Risk Forecasting using Contextual Vision Transformers with Static Map Generation and Coarse-Fine-Coarse Transformers</b> | <b>103</b> |
| 5.1      | Introduction . . . . .  | 104        |
| 5.2      | Methodology . . . . .   | 106        |
| 5.2.1    | Definitions . . . . .   | 106        |
| 5.2.2    | Problem Formulation . . . . .   | 106        |
| 5.2.3    | Contextual Vision Transformer (C-ViT) Model . . . . .   | 107        |
| 5.3      | Experiments and Results . . . . .   | 109        |
| 5.3.1    | Datasets . . . . .  | 109        |
| 5.3.2    | Experiment Setup . . . . .  | 110        |
| 5.3.3    | Evaluation Metrics . . . . .  | 111        |
| 5.3.4    | Baselines . . . . .   | 111        |
| 5.3.5    | Results . . . . .   | 112        |
| 5.4      | Coarse-Fine-Coarse Visual Transformer (CFC-ViT) . . . . .   | 113        |
| 5.5      | Application of the Static Map Generation . . . . .  | 115        |
| 5.5.1    | Pipeline description . . . . .  | 115        |
| 5.5.2    | Description of combination operations . . . . .   | 116        |
| 5.5.3    | S-ARM results . . . . .   | 117        |
| 5.6      | Conclusion . . . . .  | 119        |
| <b>6</b> | <b>Automatic Accident Detection, Segmentation and Duration Prediction using Machine Learning</b>  | <b>123</b> |
| 6.1      | Introduction . . . . .  | 124        |
| 6.2      | Related Works . . . . .   | 126        |
| 6.3      | Methodology . . . . .   | 127        |
| 6.3.1    | Case study . . . . .  | 127        |
|          | CTADS: Accident reports data set . . . . .  | 128        |

|          |  |            |
|----------|--|------------|
|          | PeMS: Traffic speed and flow data set . . . . .  | 128        |
| 6.3.2    | Speed difference estimation definitions . . . . .  | 128        |
| 6.3.3    | Accident duration prediction task definitions . . . . .  | 130        |
| 6.3.4    | Algorithm for vehicle detector station to accident association . . . . .                                       | 131        |
| 6.3.5    | Algorithm for automated disruption segmentation (ADS) . . . . .  | 133        |
| 6.3.6    | Modification of the algorithm for automated real-time early disruption detection                               | 135        |
| 6.4      | Results . . . . .  | 135        |
| 6.4.1    | Data exploration and setup . . . . .   | 135        |
| 6.4.2    | Metric performance comparison . . . . .  | 136        |
| 6.4.3    | Combination of our proposed methodology with modern methods for accident<br>scene segmentation . . . . .       | 139        |
| 6.4.4    | Automated disruption segmentation results . . . . .  | 140        |
| 6.4.5    | Comparison of estimated, reported and manual markup of accident durations                                      | 140        |
| 6.4.6    | Extraction of disruption shapes . . . . .  | 141        |
| 6.4.7    | Accident duration prediction . . . . .   | 142        |
| 6.5      | Ablation study . . . . .   | 143        |
| 6.5.1    | Parameter importance study . . . . .   | 145        |
| 6.6      | Conclusion . . . . .   | 147        |
| 6.7      | Conclusion . . . . .   | 149        |
| <b>7</b> | <b>Discussion, Synthesis and Conclusions</b>   | <b>151</b> |
| 7.1      | Literative review: discussion, synthesis and conclusions . . . . .   | 151        |
| 7.2      | Bi-level framework: discussion, synthesis and conclusions . . . . .  | 152        |
| 7.3      | Data fusion for traffic incident duration prediction: discussion, synthesis and conclusions                    | 153        |
| 7.4      | Visual transformers for traffic accident risk prediction: discussion, synthesis and con-<br>clusions . . . . . | 154        |
| 7.5      | Accident Segmentation: discussion, synthesis and conclusions . . . . .   | 155        |
| 7.6      | Final thesis Conclusion . . . . .  | 156        |
|          | <b>Bibliography</b>  | <b>157</b> |

# List of Figures

|      |  |    |
|------|--|----|
| 1.1  | Figure 1. Visualisation of spatial-temporal incident impact Liu et al., 2017 . . . . .   | 6  |
| 2.1  | Flow diagram for systematic review based on the PRISMA approach . . . . .  | 12 |
| 2.2  | Connected papers graph example . . . . .   | 13 |
| 2.3  | Reviewed publications grouped by year. . . . .   | 14 |
| 2.4  | Estimation of incident impact duration from speed profiles. Source: Haule, H. J., T. Sando, R. Lentz, C.-H. Chuan, and P. Alluri, Evaluating the impact and clearance duration of freeway incidentsHaule et al., 2019a . . . . .   | 15 |
| 2.5  | Machine Learning Pipeline for Traffic Accident Duration prediction . . . . .   | 39 |
| 2.6  | Feature importance for All-to-All regression using XGBoost for San-Francisco, USA Grigorev et al., 2022a . . . . .   | 41 |
| 3.1  | Data profiling for all data sets in our study: Victoria Rd (A) - a) network mapping, d) ecdf - empirical cumulative distribution function g) distribution plot; M7 motorway (M) - b) network mapping, e) ecdf h) distribution plot; San Francisco (SF) - c) network mapping, f) ecdf i) distribution plot. . . . . | 50 |
| 3.2  | The proposed bi-level modelling framework for traffic incident duration prediction. .  | 52 |
| 3.3  | Distribution of incident durations according to MUTCD duration classes: a) Arterial roads, Sydney, Australia b) M7 Motorway, Sydney, Australia c) San-Francisco, USA   | 55 |
| 3.4  | Performance testing of ML models across three different data sets . . . . .  | 56 |
| 3.5  | Data sets with 10% of points with the highest anomaly score removed using Isolation-Forest: a) Arterial roads, Sydney, Australia b) M7 Motorway, Sydney, Australia c) San-Francisco, USA . . . . .   | 60 |
| 3.6  | IEO-ML algorithm with a) Intra joint optimisation schema for the EO-ML algorithm, b) Extra joint optimisation schema for the IO-ML algorithm. Red dot on schema blocks represents output in the form of the best combination of ORM and model hyper-parameters . . . . .   | 60 |
| 3.7  | Incident duration classification using varying thresholds for a) data set AR b) data set M c) data set SF. The red percentage above each set of ML results indicate the percentage split of Subset A and B for that particular $T_c$ . . . . .   | 65 |
| 3.8  | Outlier removal for a) data set AR b) data set M c) data set SF . . . . .  | 66 |
| 3.9  | Multi-class (3-class) classification using quantile splits for a) data set AR b) data set M c) data set SF . . . . .   | 69 |
| 3.10 | Regression using Quantiled Time Folding for a) data set AR b) data set M c) data set SF  | 69 |
| 3.11 | Regression using randomised 10-folds for a) data set AR b) data set M c) data set SF   | 69 |

|      |   |     |
|------|---|-----|
| 3.12 | Pipeline (a) and fusion (b) approaches for the bi-level framework structure . . . . .   | 73  |
| 3.13 | Comparison of the fusion and single model performance for a) data set AR b) data set M c) data set SF. Dashed lines represent the average RMSE score across all folds for each corresponding model . . . . .  | 74  |
| 3.14 | Feature importance for All-to-All regression using XGBoost for a) Arterial roads, Sydney, Australia b) M7 motorway, Sydney, Australia c) San-Francisco, USA . . . . .   | 78  |
| 3.15 | Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of Arterial roads, Sydney, Australia . . . . .  | 80  |
| 3.16 | Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of M7 Motorway, Sydney, Australia . . . . .   | 80  |
| 3.17 | Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of San-Francisco, USA . . . . .   | 81  |
| 3.18 | Binary classification performance using varying incident duration threshold . . . . .   | 83  |
| 3.19 | Performance testing of ML models across three different data sets . . . . .   | 83  |
| 4.1  | a) Traffic speed and b) Traffic flow plots for the VDS associated to incident A-4798 (accident on US-101 Southbound with duration of 31 5-minute iterations - actual reported incident clearance time, without considering the incident recovery time). The red line denotes the start of the accident, and the green line the end of the accident. The blue line denotes the speed evolution in the vicinity of the incident location (drops almost to 20km/h) while the flow is still running at high values due to large numbers of vehicles blocked in traffic. . . . . | 88  |
| 4.2  | The structure of the proposed framework . . . . .   | 89  |
| 4.3  | LSTM sentiment encoder structure. . . . .   | 90  |
| 4.4  | Example of LSTM network training results using 12 units, a ReLU activation function, 10 epochs, 80 hidden units. a) Train-validation score over 20 epochs . . . . .   | 91  |
| 4.5  | The structure of the ANN autoencoder . . . . .  | 92  |
| 4.6  | RMSE vs MAPE results for a) CCS data set, b) CTADS - incident duration c) Random vectors . . . . .  | 94  |
| 4.7  | Regression results for baseline feature set across different ML models. . . . .   | 95  |
| 4.8  | Parallel categories representation for all regression scenarios with GDBT. . . . .  | 97  |
| 4.9  | Word importance estimation using LIME method for incident severity groups . . . . .   | 98  |
| 4.10 | Word importance estimation using LIME method for incident duration groups . . . . .   | 98  |
| 4.11 | Traffic speed and flow during the day of the incident. Part #1 . . . . .  | 100 |
| 4.12 | Traffic speed and flow during the day of the incident. Part #2 . . . . .  | 101 |
| 5.1  | City grid representation for our study. . . . .   | 104 |
| 5.2  | The building blocks of our proposed C-ViT model. . . . .  | 105 |
| 5.3  | Description of the first stage of the historical risk Map encoding. Given a unified single image $X$ , it is then divided into equally-sized image patches that are passed individually to the linear patch embedding layer. . . . .  | 107 |
| 5.4  | Comparison between our proposed C-ViT model and GSNet Wang et al., 2021b, in terms of the number of training parameters. . . . .  | 113 |

|      |   |     |
|------|---|-----|
| 5.5  | Example of GSNet predictions (after training for 2 epochs, when the best performance is observed): a) Actual map of the accident occurrence b) Predicted map of the accident occurrence. . . . .  | 114 |
| 5.6  | Coarse-Fine-Coarse Transformer . . . . .  | 115 |
| 5.7  | The building blocks of our proposed XViT model with Static Map Generation . . . .   | 116 |
| 5.8  | Root Mean Squared Error from performance evaluation of our XViT model for a number of combination operations on NYC data set . . . . .  | 120 |
| 5.9  | Root Mean Squared Error from performance evaluation of our XViT model for a number of combination operations on Chicago data set . . . . .  | 120 |
| 6.1  | Contributions and data-flow schema for association of traffic speed readings with accident reports . . . . .  | 126 |
| 6.2  | CTADS reported accidents for San-Francisco . . . . .  | 128 |
| 6.3  | 1) PeMS data set area coverage for San-Francisco (the map is available at <a href="https://pems.dot.ca.gov/">https://pems.dot.ca.gov/</a> ) 2) Mapping of the Vehicle Detection Stations from PeMS data set. OpenStreetMap excerpt showing San Francisco. Available at: <a href="https://www.openstreetmap.org/#map=12/37.7612/-122.4395">https://www.openstreetmap.org/#map=12/37.7612/-122.4395</a> . . . . . | 129 |
| 6.4  | The application of dilation operation to an image and time series . . . . .   | 134 |
| 6.5  | User input errors located within the CTADS data set . . . . .   | 137 |
| 6.6  | Various metrics applied to difference between recorded speed and speed profile . . .  | 138 |
| 6.7  | Results disruption segmentation algorithm application for accidents a) A-5764, b) A-8119, c) A-9931 . . . . .   | 138 |
| 6.8  | Distribution of accident durations for a) estimated, b) reported accident durations for the area of San Francisco, c) results of manual markup of disruptions observed in traffic speed . . . . .   | 141 |
| 6.9  | Scatter plot for a) estimated and b) reported accident durations for the area of San Francisco, c) results of manual markup of disruptions observed in traffic speed . . . .  | 141 |
| 6.10 | Normalized Wasserstein distance plot for disruption shapes extracted for segmented intervals . . . . .  | 142 |
| 6.11 | Manual markup and algorithm segmentation comparison. Time series segments represented as binary values of 0 and 1. . . . .  | 144 |
| 6.12 | Histogram of F1-score against manual markup for a) reported accident time interval and b) estimated segmentation when algorithm detecting multiple disruption intervals c) estimated segmentation for the single closest interval to reported incident occurrence time . . . . .  | 145 |
| 6.13 | Connection between metric and F1 score . . . . .  | 146 |
| 6.14 | Scatter plots between model parameters and F1 score . . . . .   | 147 |
| 6.15 | Scatter plot between model parameters and F1 score . . . . .  | 148 |



# List of Tables

|     |   |     |
|-----|---|-----|
| 2.1 | Table of temporal features used to describe traffic incident . . . . .  | 18  |
| 2.2 | Table of features used to describe traffic incident across different studies . . . . .  | 19  |
| 2.3 | Metrics used across reviewed papers. . . . .  | 40  |
| 2.4 | Most popular Machine Learning methods used across reviewed papers . . . . .   | 41  |
| 2.5 | Most popular Machine Learning pipeline elements used across reviewed papers. . . .  | 42  |
| 2.6 | Most popular Deep Learning methods used across reviewed papers . . . . .  | 42  |
| 3.1 | Traffic incident features for Sydney Arterial roads (AR) and M7 motorway (M). . . .   | 51  |
| 3.2 | Multi-class classification results for equally-sized 3-class split . . . . .  | 68  |
| 3.3 | MAPE results for all 7 scenarios on data set AR . . . . .   | 71  |
| 3.4 | MAPE results for all 7 scenarios on data set M . . . . .  | 71  |
| 3.5 | MAPE results for all 7 scenarios on data set SF . . . . .   | 71  |
| 3.6 | MAPE results for All-to-All scenario of data set A, using different ORM approaches<br>and incident duration transformation, via the proposed IEO-ML approach. . . . .                           | 76  |
| 3.7 | MAPE results for All-to-All scenario of data set M, using different ORM approaches<br>and incident duration transformation, via the proposed IEO-ML approach. . . . .                           | 76  |
| 3.8 | MAPE results for All-to-All scenario of data set SF, using different approaches for<br>ORM and incident duration transformation, via the proposed IEO-ML approach. . . .                        | 77  |
| 4.1 | Example of the Incident Description values . . . . .  | 90  |
| 4.2 | Top 8 best scenario results for GBDT-enabled framework . . . . .  | 95  |
| 4.3 | Top 8 best results for RF . . . . .   | 95  |
| 4.4 | Top 8 best results for XGBoost . . . . .  | 96  |
| 5.1 | Datasets Statistics . . . . .   | 109 |
| 5.2 | Performance evaluation of our C-ViT model against a number of baseline approaches<br>from the literature over the NYC and Chicago datasets. . . . .   | 110 |
| 5.3 | Performance evaluation of our C-ViT model against a number of baseline approaches<br>from the literature over the high frequency times of accidents in the NYC and Chicago<br>datasets. . . . . | 112 |
| 5.4 | Performance evaluation of our CFC-ViT model on Chicago data set . . . . .   | 114 |
| 5.5 | Performance evaluation of our CFC-ViT model on NYC data set . . . . .   | 114 |
| 5.6 | Performance evaluation of our XViT model for a number of combination operations<br>on NYC data set. Top 20 results. . . . .   | 118 |
| 5.7 | Performance evaluation of our XViT model for a number of combination operations<br>on Chicago data set. Top 20 results. . . . .   | 119 |

|     |  |     |
|-----|--|-----|
| 6.1 | Mean Absolute Error (MAE) Results . . . . .      | 143 |
| 6.2 | Root Mean Squared Error (RMSE) Results . . . . . | 143 |



## Chapter 1

# Introduction

Today, Intelligent Transportation Systems (ITS) are an essential component of transport networks in modern cities. These systems monitor and control transport systems, ensuring safety, increasing efficiency, reducing travel time, and lowering air emissions, which significantly impact the economy and health of city populations.

Traffic congestion is a significant concern for many cities globally. Various factors, such as increased population, workforce concentration in central areas, and lack of efficient public transport modes, contribute to congestion. Two primary forms of congestion occur: a) recurrent traffic congestion during peak hours when traffic demand exceeds road capacity, and b) non-recurrent traffic congestion caused by unplanned events like car accidents, breakdowns, weather, or public demonstrations. Previous studies have shown that nearly 60% of traffic congestion results from non-recurrent incidents with stochastic behavior in space and time Schrank and Lomax, 2002. In Australia, the number of road deaths per year has decreased by 70

Despite ITS systems' efforts to optimize congestion and ease traffic, accidents can occur anywhere and anytime. Transport agencies optimize traffic movements, but accidents can still affect traffic flow and cause congestion, sometimes across multiple adjacent roads. Traffic disruption is an unwanted effect of severe congestion, which can be recurrent (repetitive characteristic of transport networks) or non-recurrent (rarely observed disruptions, traffic incidents).

In Australia, there is currently a lack of advanced incident management and response plan solutions, and most transport management centers rely on staff members' operational experience rather than data-driven approaches. There is a potential to utilize more advanced solutions that utilize extensive information sources published worldwide and modern data-driven techniques like machine learning and deep learning. Additionally, the literature lacks exploration of various modeling capabilities that combine transport modeling and data-driven solutions.

Traffic Incident Management Systems (TIMS) collect data on traffic incidents, including information on different incident duration factors. Accurately predicting the total incident duration shortly after an incident occurred could save operational costs (by providing advice on necessary amounts of equipment and response team size, strategy of incident evaluation depending on its predicted duration) and end-user time (through affecting individual user route planning). Moreover, the clearance time of accidents is highly dependent on the ongoing traffic congestion and several external factors with different degree of importance. Therefore, it is essential to estimate the incident factor importance to improve the accuracy of predictions. Most prior studies related to this topic concentrated on testing

different machine learning models on specific road types like freeways or highways and focuses primarily on different phases of the incident duration such as clearance time, recovery time, and total incident duration Li, Pereira, and Ben-Akiva, 2018a. There is currently a lack of an unified approach that can be applied on all road types, for all accident types, and across various countries with different driving behavior utilizing large amounts of openly available data.

Advances that have the highest impact and which can improve our ability to analyze traffic incidents include:

1. Popularisation of Machine Learning and Deep Learning techniques.
2. Availability of high-performance computing devices.
3. Availability of multiple traffic simulation approaches (meso, macro, hybrid, etc) both for the traffic incident spatial-temporal impact analysis and response strategy evaluation.

Important topics in traffic incident research areas include:

- Traffic incident duration prediction, application of Machine Learning methods and different approaches to data processing.
- The detection of traffic incidents, based on traffic flow data, estimation of traffic incident severity.
- The spatial-temporal incident impact analysis (including impact mapping and estimation of the life-cycle of the traffic incident).
- The response plan modelling and evaluation using different simulation approaches (at micro-scopic, meso-scopic and hybrid levels).

Today, we are open to new approaches and methodologies, which we can use for the incident duration prediction, spatial-temporal incident impact analysis and traffic simulation.

## 1.1 Research Project Summary

**Motivation:** with recent advances in the fields of Machine Learning and Deep Learning we can develop a methodology for traffic incident analysis, which will improve the prediction accuracy of traffic incident duration and spatial-temporal impact estimation.

The primary objective of this research is to construct a comprehensive modeling framework for incident modelling by utilizing advanced machine learning and deep learning methodologies. The goal is to enable accurate prediction of the time duration of reported accidents. The key goals and objectives of the study are outlined below and embody the research blueprint. This research also encapsulates the exploration of data fusion techniques and the development of automated algorithms for accident timeline segmentation. These combined techniques will be instrumental in creating a more precise and detailed understanding of incident duration, leading to more effective and timely incident response strategies.

**Aim 1:** to explore Machine Learning and Deep Learning capabilities for the task of traffic incident duration prediction. In detail, to evaluate the best and state-of-the-art methods of Deep Learning and Machine Learning, different data analysis approaches on the task of traffic incident analysis with the

goal to improve traffic incident duration prediction accuracy. Successful delivery of this aim consists in the development of the methodology for the traffic incident modelling, which allows acceptable accuracy of predicted incident durations.

**Aim 2:** to deliver a modelling framework for the incident impact analysis of the incident impact. In detail, to build a system, which incorporates multiple tasks from the theory of traffic incident analysis: traffic incident detection, incident duration prediction and incident impact analysis (using modelling and simulation). Successful delivery of this aim consists of the development of a system, which allows to detect traffic incident from traffic flow/speed data, estimate and predict their duration and develop a measure to represent their impact.

## 1.2 Stages of the Research

This PhD project will focus on the complex problem of predicting the impact of traffic disruptions in large cities using advanced artificial intelligence algorithms and evaluating the best response plan that traffic authorities can make by synergising traffic simulation modelling of various response scenarios. Important tasks to bring this project to fruition include:

1. The prediction of traffic incident duration employing contemporary Machine Learning techniques and data processing methodologies. This involves utilizing the Machine Learning models to analyze incident data and predict the duration of future incidents accurately.
2. Estimating the impact of an incident, including the evaluation of the temporal lifecycle of the disruption. This phase includes the identification of traffic incidents through the analysis of traffic flow data, which is characterized as time series data obtained from vehicle detection systems.

The solution would improve traffic centres decision making by automatic response plan recommendations. The following sections will detail each stage of the PhD project.

### 1.2.1 Stage 1. Incident Detection, Incident Duration Modeling, Incident Severity Prediction

Most current studies rely on methods for classification and clustering of traffic conditions for doing the incident detection. However, there are very few studies on traffic incidents involving methods for detecting anomalies (such as one-class SVM, isolation forests, etc). Non-recurring traffic incidents are rare and unusual in nature and therefore the detection of a traffic incident can be assessed as the task of detecting anomalies in traffic. By relying on anomaly detection methods, the incident detection system can be adapted to previously unseen situations. Thus, classification and evaluation of the applicable anomaly detection methods in comparison to well-established classification and regression methods will be carried out. Also, road situations detected as anomalous can be extremely valuable for further investigations in the duration of a freshly reported accident.

Incident duration distribution has been modelled as log-normal Sullivan, 1997 and more recently as log-logistics distribution Chung, Walubita, and Choi, 2010, Smith and Smith, 2002. Log-logistic

model has been used more extensively and found to have better goodness-of-fit than log-normal distribution. Also, there are various hazard-based models of traffic incident duration Nam and Mannering, 2000a, Hojati et al., 2013 which employ a hazard function to describe the conditional probability that an incident will end during any particular time interval given that it already lasted until the beginning of the interval. Recent studies also involve multi-component log-logistic models. In Zou et al., 2016a authors describe a g-component log-logistic model and in Li, Pereira, and Ben-Akiva, 2015b describes competing risks mixture model which incorporates multinomial log-logistic model.

An actual distribution estimation can give only approximate information on traffic incident duration. More than that, incident duration distribution is found significantly dependent on incident case parameters (e.g. day/night) Yang et al., n.d. Also, the duration of the incident may be affected by chosen method of incident clearance Li, Pereira, and Ben-Akiva, 2015b. Mao et al., 2019 found that some incident parameters found to be important factors with different contribution to different types of accidents, including weather condition, traffic density, time period, incident location. Also, road factors found to be affecting each one of 4 incident types (rear-end, side wipe, collision with fixtures and rollover) in a different way. These findings draw incident duration distribution estimation as a complex problem dependent on many traffic flow and incident parameters.

Incident duration can be modelled in terms of spatial relations (geometric placement of adjacent lanes, angle of adjacency, different parameters of lanes, including speed limits). Recent studies rely on reported incident parameters Mihaita et al., 2019a, Hamad, Khalil, and Alozi, 2019, but road topology can also play significant role in estimation of the incident probability (e.g. poorly designed junction, wrongly imposed speed limits). According to Curiel, Ramirez, and Bishop, 2018, about 5% of the road junctions are the site of 50% of the accidents in the city of London. Thus, it seems reasonable to analyse incident duration and probability with consideration of the road topology. The task of predicting the duration of an incident usually solved by using Machine Learning methods. Among these methods – tree based methods Ozbay and Kachroo, 1999, fuzzy logic Wang, Chen, and Bell, 2002, Bayesian networks Ozbay and Noyan, 2006a, artificial neural networks Barcellos et al., 2015, Alkheder, Taamneh, and Taamneh, 2017. And recently Ma et al., 2017 studied GBDT as a better performing method for incident duration prediction. Gaussian process regression and artificial neural networks were found to outperform tree methods and SVM in incident duration prediction Hamad, Khalil, and Alozi, 2019.

Also, estimation of incident duration can be reduced to the classification method Mihaita et al., 2019a. To do this, a specific threshold for the duration is set and a prediction is made whether the incident will last longer than a specified time. Artificial neural networks show high average accuracy for prediction of 4 types of incident severity relying on data on the state of the road (lane, condition of the roadway, weather, light, etc.), time and date. Overall accuracy between death, severe, moderate and minor severity accidents was found to be 69-72% Alkheder, Taamneh, and Taamneh, 2017.

Recent studies in machine learning involve interpret-able models. Bayesian networks can produce interpret-able models for incident injury severity prediction O

textasciitilde na, Mujalli, and Calvo, 2011. Bayesian networks also outperform regression models in incident severity prediction (involving three severity indicators: number of fatalities, number of injuries and property damage) Zong, Xu, and Zhang, 2013. Interpret-ability is not specific property of the tree models only and by using knowledge distillation one can extract tree rules from different

prediction models (e.g. Bayesian network Park, Haghani, and Zhang, 2016a). It allows to represent the model as an interpret-able decision tree and estimate feature importance. Methodology: Stage 1 involves the use of anomaly detection, regression and classification methods with the use of reported incident parameters and map data.

**RESEARCH QUESTION 1:** Can we build an universal framework to work with different traffic incident data sets? Traffic incident reports for different countries vary a lot due to the methodology of data collection. How much is the incident duration affected by each of the incident report variables? By using well-established classification and regression methods (GBDT, ANN, etc) and feature importance estimation methods we can select incident report variables that have the highest contribution to the incident duration prediction accuracy. Are these variables common among different data sets?

**RESEARCH QUESTION 2:** How can we utilise anomaly detection methods to improve the traffic incident duration prediction? We will use different anomaly detection methods, including those which can produce “measures of anomaly” for each data point (One-Class SVM, Isolation Forest). It will be used to compare the anomaly detection with regression methods (GBDT, ANN) for the task of incident probability estimation. Similarly, anomaly detection can be used in comparison with recently used classification (GBDT, ANN) and regression (e.g. Random Forest regression) methods, for the task of incident duration estimation. Anomaly detection methods will be used to model different kinds of anomalies in incident reports.

**RESEARCH QUESTION 3:** What are the abilities of different modern Machine Learning methods for the task of incident duration prediction? And how can we model short-term and long-term incidents: what duration split threshold to use and how can we approach a regression task?

### 1.2.2 Stage 2. Incident Impact modelling

By using the spatial-temporal forecast, it is possible to produce an estimation of the road situation development, which can be used in planning of a strategy to eliminate the incident. Also, this kind of forecasting allows to produce incident affect-ability map (which will show how much each road element can be affected and for how long) for every component of the city road network. This kind of data can be used in road planning decisions to reduce incident impacts on road network in long-term and to reveal road elements which are the most sensitive to traffic incidents (e.g. produce wide-spread or long-term congestion). Spatial-temporal incident modelling (e.g. impact forecasting) is prior for the next stage – incident response modelling.

Spatial-temporal impact can be estimated using different approaches to traffic simulation. There are 3 groups of traffic simulation models:

a) **Microscopic** – traffic network is simulated on the level of individual agents (car, pedestrian), relying on rules of movement and interaction (including lane change, acceleration). These models include Car-following models (which relies on real driving behaviour such as keeping a “safe distance” from the leading vehicle Treiber and Kesting, 2013a), such as Gibbs model, Intelligent driver model, etc. All of these methods require a lot of computational resources. Microscopic models are preferable for the Stage 3 of the research since it will include traffic control actions which will affect traffic on micro level (changing of traffic lights, using Visual Message System).

b) **Mesoscopic** - traffic is represented by interactive groups of traffic entities.

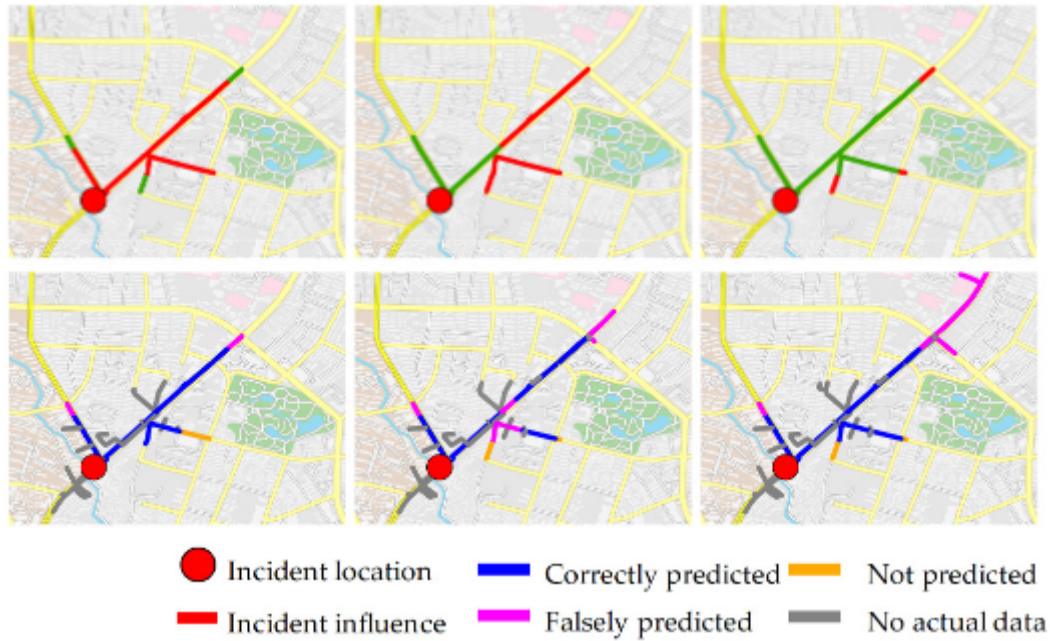


FIGURE 1.1: Figure 1. Visualisation of spatial-temporal incident impact Liu et al., 2017

c) **Macroscopic** – traffic is simulated on the level of traffic flow on road segments. This level of detail considers traffic speed, flow and density) and their relationships. Cases of macro simulation include:

- Binary integer programming(BIP) Chung and Recker, 2012 applied in estimating the temporal and spatial extent of delay caused by freeway accidents.
- Kinematic Shockwave propagation model Liu et al., 2017 - the model leverages a physical traffic shockwave model, analysing different superposition situations of shockwaves. Also, this method has been compared to car-following model and found to be superior in performance (by almost 20 times).
- Temporal Graph-Convolution Network Zhao et al., 2018 (a combination of gated recurrent units and graph-convolution network) applied for real-time traffic forecasting.

Macroscopic modelling of traffic disruptions is preferable for the Stage 2 of the research, since data availability both on traffic flow (on road sections in a close proximity to the source of the disruption) and incident reports. Using macroscopic modelling first will help to research theoretical basis of traffic congestion.

The differential effects of determinants (traffic diversion requirement, crash injury type, number and type of vehicles involved in a crash, day of week and time of day, towing support requirement and damage to the infrastructure) on crash survival probabilities are found to vary considerably across the motorways Tajtehranifard et al., 2016.

Both spatial and temporal aspects of traffic incidents were analysed using real-world spatial-temporal traffic sensor data on road networks Pan et al., 2013. Authors analysed abrupt and long-lasting propagation of the speed changes.



As indicated by many researchers, calculating the consequences of an incident using micro-simulation models is extremely resource-intensive Azevedo et al., 2016-Huang and Pan, 2007-Ozbay and Noyan, 2006a, which made a number of experiments impossible. As known, many transport simulation packages rely solely on the CPU usage (VISSIM, AIMSUN, SUMO, etc).

However, recent advances in the development of parallel computing have made it possible to use GPUs in agent modelling and especially in traffic simulation. Studies of agent models Xiao et al., 2019, and in particular transport micro-agents using GPUs, are already available Heywood, Richmond, and Maddock, 2015, Heywood et al., 2018. New opportunities which transport modelling on the GPUs opens up are significant and fundamental for further research in the area of traffic analysis. The use of GPUs opens a road for a field of incident response experiments involving optimisation based on micro-simulation.

Widely used traffic simulation software AIMSUN has been compared with FLAME GPU Xiao et al., 2019. FLAME GPU shows significant speedup by orders of magnitude (10-100x). The largest simulation executed using both simulators, a one hour simulation containing up to 512,000 vehicles and 1,575,936 detectors, showed a speed-up of 43.8x for the GPU accelerated simulation. While one simulation is being performed on the CPU, up to 50 simulations on the GPU can be done, so finding a solution strategy (by performing multiple simulations) for the incident becomes possible in an acceptable real time. Also, a simulation for the large-scale incident impact becomes possible. Thus, incident impact can be modelled both for small and large-scale environments. This approach has a potential to produce applicable incident response solutions in real-time within meaningful time constraints (during incident response planning). Use of the different tools for the task of traffic control (and not only traffic lights control) draws a complex optimisation task, which in combination with micro-simulation can produce more effective solutions than using only macro-models and simplistic control.

**Methodology:** The research on the current stage will be devoted to the use of micro-simulation modelling of transport network for the task of analysis of spatial-temporal impact during the incident. The main research questions to be solved in this stage are:

**RESEARCH QUESTION 1:** By how much some traffic network nodes are affected more in comparison with others? What are the conclusions from their difference?

**RESEARCH QUESTION 2:** Can we get an incident duration/severity heat-map for the road planning agency (using large-scale micro-simulation of traffic incident impacts)?

**OUTCOMES:** As a result, a map with rated transport nodes for selected traffic network will be produced. Practically it can allow to find areas of interest for transport planning agency. The use of this kind of map can allow to use incident-repelling and incident-attractive route planning, which also will be assessed. A large-scale simulation will be used to evaluate the global impact of incidents on transport nodes on the system as a whole. Incident affect-ability map (which can be defined as a measure of effect on average traffic speed in comparison to a normal traffic condition) can be calculated also for different kinds of incidents. Road elements with highest/lowest produced duration and spatial impact will be studied.

After analysis of duration and severity of impacts using micro-simulation, results on incident duration can be compared to predictions produced by ML methods on stage 1 to assess effectiveness

of micro-simulation for the task of incident duration modelling. The feasible way for this stage of research is to use macro/meso-models or micro-models based on GPU utilisation.

As an additional research, temporal graph-convolution network (which is a macro-model) can be approached to approximate simulation results of micro-model in order to reduce computational resources consumption in further task of incident response modelling.



## **Chapter 2**

# **Literature review**

## 2.1 Introduction to Traffic accident analysis

Today, Artificial Intelligence (AI) is being used to enhance the performance of different industries and businesses, especially the transport industry. AI technologies such as Machine Learning (ML) and Deep Learning (DL) models can be used to address transportation problems such as traffic management, urban mobility and traffic safety. AI models are used to solve traffic prediction, traffic control, road safety planning and traffic flow optimisation problems Abduljabbar et al., 2019; Machin et al., 2018.

Traffic congestion, which in 60% of cases occurs due to unplanned events Schrank and Lomax, 2002, is a significant concern for many cities around the world. Congestion arises due to various factors, including increased population, workforce concentration in central areas, or the lack of efficient public transport modes. Two forms of congestion are typically predominant: a) recurrent traffic congestion during peak hours when traffic demand exceeds the road capacity, and b) non-recurrent traffic congestion caused by stochastic events such as car accidents, breakdowns, weather effects, etc. In Australia, the number of road deaths per year was reduced by 70% since the 1970s. However, the annual economic cost of road crashes was estimated at \$27 billion per annum in 2017 Government, 2017. In Melbourne, Australia more than 640 km of arterial roads are congested during peak hours with 2.9 tons of CO<sub>2</sub> emissions during the years 2014-2015 Linking Melbourne Authority, n.d.

Intelligent Transportation Systems (ITS) are an integral element of transport networks in modern cities. These systems provide monitoring and control of the transport system, ensuring safety, increasing efficiency, reducing travel time, reducing air emissions and thus having a significant impact on the economy and health of the city population. The incorporation of AI techniques into the ITS system has the potential to greatly reduce traffic congestion and its effects on the environment. The main data sources used by Intelligent transportation systems (ITS) are: vehicle detectors (magnetic, infrared, ultrasonic, and microwave), traffic cameras, Global Positioning Systems and Automatic Vehicle Identifiers (e.g. electronic toll collection, access control and speed control) Al-Bordiny, 2014. AI techniques were applied to these kinds of data previously Ma et al., 2020; Benterki et al., 2020. Multiple various measures can be taken by ITS to reduce the impact of incidents (e.g. variable message signs, toll roads, adaptive cruise control, adaptive traffic light control, transport group priority management) Al-Bordiny, 2014.

In Australia, there is currently a lack of incident management and response plans solutions and the majority of transport management centres make decisions based on the operational experience of staff members rather than data-driven lessons learned. There is a true lack in adopting more advanced solutions that can make use of any existing sources of information and modern data-driven techniques such as machine learning /deep learning. Also in the literature, there is a lack of exploring various all modelling capabilities combining transport modelling and data-driven solutions.

Traffic Incident Management Systems (TIMS) collect data on traffic incidents, including information on different incident duration factors. Accurately predicting the total duration shortly after an incident could save operational costs and end-user time (by affecting route planning). Moreover, the clearance time of accidents is highly related to the ongoing traffic congestion and several external factors with different weights of importance. Therefore, it is essential to estimate the incident factor

importance to improve the accuracy of predictions. Most prior studies related to this topic concentrated on testing different machine learning models on specific road types like freeways or highways and focused primarily on different phases of the incident duration such as clearance time, recovery time, and total incident duration Li, Pereira, and Ben-Akiva, 2018b. There is currently a lack of an advanced approach that can be applied on all road types, for all accident types and across various countries with different driving behaviour.

Deep Learning and Machine Learning have become increasingly important tools to improve traffic incident management systems (TIMS). Accurately predicting the total duration of a traffic incident shortly after it occurs is essential to saving operational costs and end-user time, as well as reducing traffic congestion. Understanding the importance of incident factor importance is key to improving the accuracy of predictions. In this paper, we review the literature related to traffic incident duration prediction and spatial-temporal accident modelling. Specifically, we discuss the challenges associated with each modelling step, the complexity of the task, and the most recent advances in this field, with a focus on the potential of deep learning and machine learning for incident duration prediction. Our goal is to provide a comprehensive overview of the most recent advances in this field, and demonstrate the potential of deep learning and machine learning for incident duration prediction and traffic simulation.

### 2.1.1 Paper structure

The paper organisation is detailed as follows:

Section 2.1.2 presents the PRISMA methodology that we have followed for our study, which has revised overall a total of almost 1200 papers on the topic of incident modelling, which have been further filtered and selected down to 75 final papers to provide a comprehensive structured analysis into current gaps and future research directions.

Section 2.2 gives an overview of all the required data sets that one needs to conduct a thorough incident modelling which ranges from accident logs, but also to traffic states such as flow, speed, occupancy, and external related information (weather, events, etc.). We also provide insights into public data sets that have been used for modelling, as many countries restrict access to such data sets due to privacy concerns.

Section 2.3 describes methods of statistical analysis for traffic accident modelling. Section 4 is devoted to the use of Machine Learning in traffic accident analysis including classification and regression tasks, feature selection, imbalanced data set management techniques, anomaly detection, dimensionality reduction, novel machine learning methods and frameworks. Section 5 gives insights into the use of advanced Deep Learning techniques for textual accident report description analysis, accident detection and segmentation from the traffic flow. Finally, in Conclusions we provide a summary of the challenges we have detected as well as future research gaps to be filled.

### 2.1.2 Literature review material and the PRISMA method

The Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) is the review method that has been applied when organising and analysing literature for this paper. The process of reviewing the literature is shown in Figure 2.1. In the first stage, relevant literature has been identified by using publication databases based on keyword search. The list of used keywords:

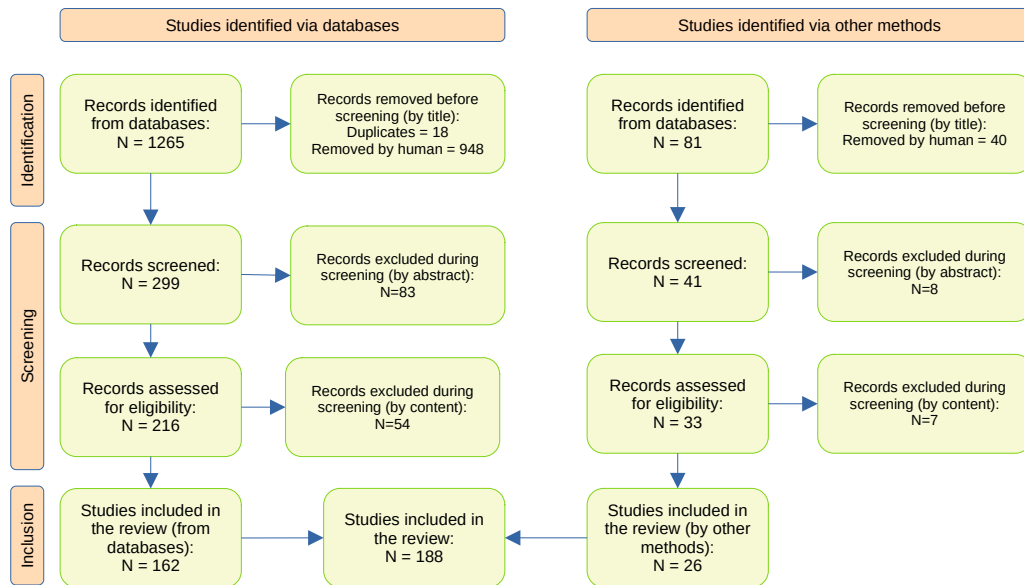


FIGURE 2.1: Flow diagram for systematic review based on the PRISMA approach

- Traffic Incident
- Traffic Incident duration prediction
- Traffic Incident clearance time
- Traffic Incident machine learning
- Traffic Incident Random Forest

Since incidents can be related to different research areas, it is necessary to always specify the 'traffic' as an area. Keywords can include 'clearance time' since this term is very specific to the task of incident duration modelling. Keywords 'machine learning' are the main methods used in incident duration prediction. The use of 'traffic incident random forest' clearly related to the tasks of classification and regression related to the traffic incident duration modelling. By using very specific terms and specifying areas it is possible to locate relevant literature quickly.

The following databases were used for the stage of literature identification:

- ScienceDirect
- Google Scholar
- Research Gate

The alternative source of information on the relevant literature is Connected Papers, which builds a graph of studies based on their semantic similarity. This requires a sample paper to search for similar ones. The search using this approach was performed after the identification of relevant literature using common databases.

The databases were accessed through the University of Technology Sydney, and the publications were limited between 1980 to 2022. In total, 1346 sources were collected, 1265 were found using conventional databases and 81 using Connected Papers.

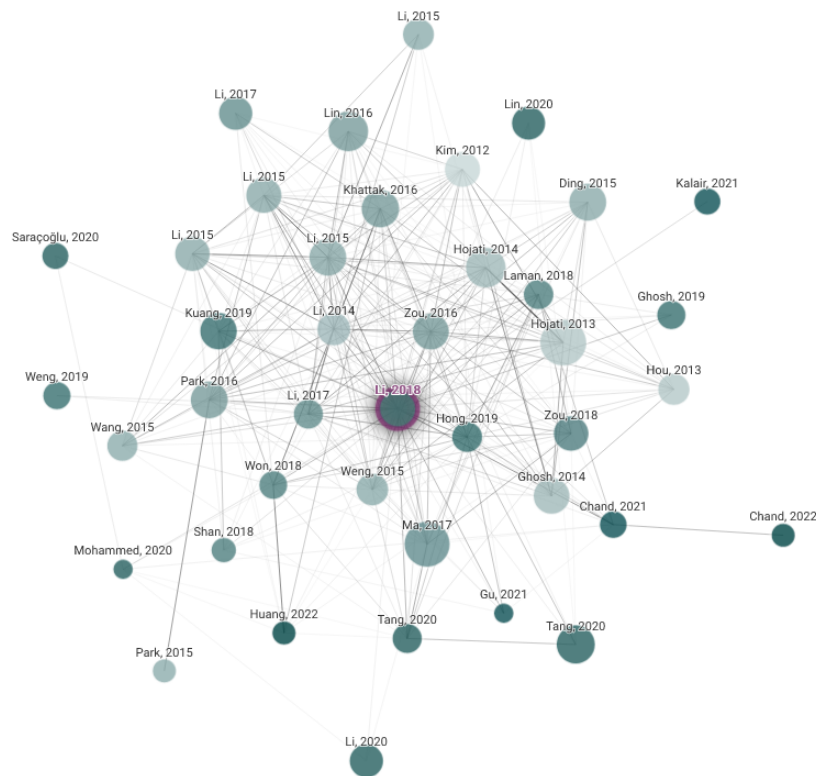


FIGURE 2.2: Connected papers graph example

The PRISMA process for this literature review is detailed in the following:

1. A database of records has been collected using previously described databases and keywords which resulted in 1,265 resources.
2. Resources were screened for duplicates which resulted in 18 resources being removed. Then, records were filtered by paper title, which resulted in the removal of 948 entries mainly due to interference with topics of “internet traffic” (which also relies on the use of machine learning methods and network incident analysis) and “molecular traffic”, presence of many studies related to injury statistics and safety analysis related to traffic accidents.
3. Filtered resources (299 in total) were screened by abstract, which resulted in the removal of 83 records due to the mention of unrelated methodologies and findings.
4. Eligible records (216 in total) were then screened by content (reading of methodology and conclusion sections), which resulted in the removal of 54 records.
5. In total, we obtained 162 records from the database search.
6. The most relevant review (Ruimin Li, F. Pereira, M. Ben-Akiva, Overview of traffic incident duration analysis and prediction) from the previous search have been selected and graph of related papers has been built (see Figure 2.2) and a similar process has been performed for semantic similarity search using the Connected Papers service, which resulted in 26 additional papers selected.

7. In total, we obtain 188 relevant resources for the literature review. The highest amount of papers is dated between 2010 and 2022 with peaks in 2002, 2013, 2016, 2018-2021 (see Figure 2.3). Peaks during this years can be attributed to the introduction of novel Machine Learning methods (e.g. RandomForest in 2002 Liaw, Wiener, et al., 2002, XGBoost in 2016 Chen and Guestrin, 2016)

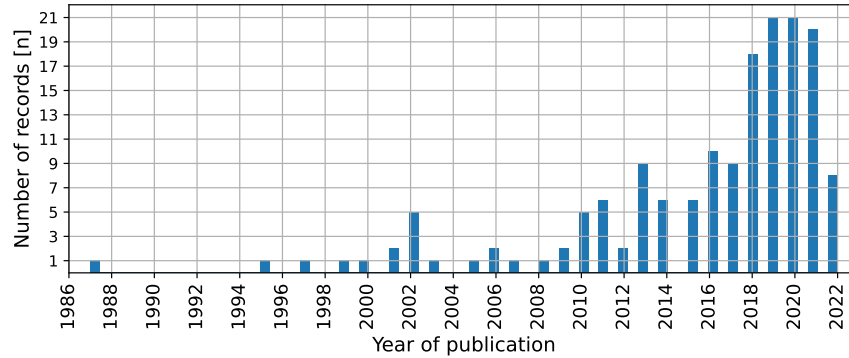


FIGURE 2.3: Reviewed publications grouped by year.

### 2.1.3 Incident duration definitions

Traffic congestion can be recurrent and non-recurrent Adler, Ommeren, and Rietveld, 2013. Non-recurrent traffic congestion is unexpected congestion, caused by random events affecting traffic flow such as traffic incidents, weather phenomena, vehicle breakdowns, hazards, etc. Recurrent traffic congestion is a predictable regularly occurring congestion, which observed in places where traffic flow regularly exceeds road capacity. KIM and CHOI, 2001

The definition of traffic incident duration phases can be found in The Highway Capacity Manual Alkaabi, Dissanayake, and Bird, 2011 and it consists of four phases:

- **Incident Detection:** the time interval between the incident occurrence and its reporting,
- **Incident Response:** the time interval between incident reporting and arrival of the first investigator at the location of the accident,
- **incident Clearance:** the time interval between the arrival of the first investigator and the clearance of the incident,
- **Incident Recovery:** time interval between the clearance of the incident and the return of traffic flow to normal conditions.

Different phases of traffic incident duration (e.g. clearance, recovery time) can be modelled individually, but this type of research is rare because of the complexity of data collection for traffic incidents and the small amount of recorded traffic incidents in real-life data sets Alkaabi, Dissanayake, and Bird, 2011.

Duration of detection, response and clearance phases were modelled separately in the literature so far by using Hazard-based duration modelling Nam and Mannering, 2000b. Researches in Li and

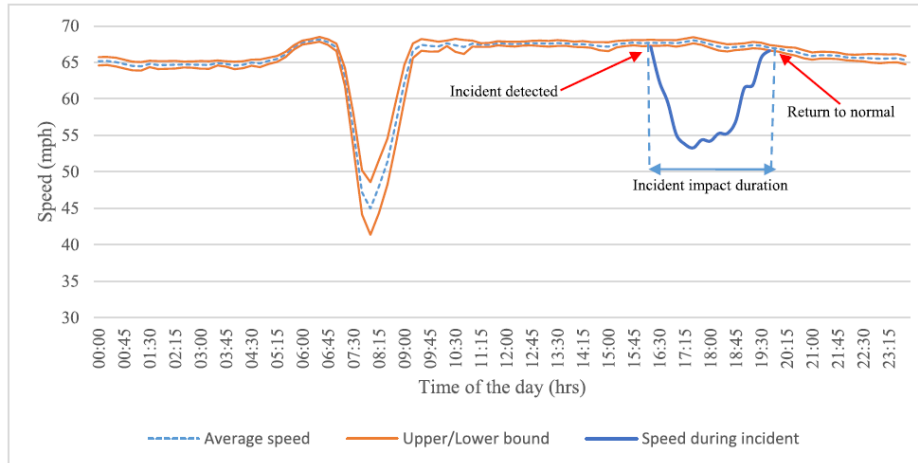


FIGURE 2.4: Estimation of incident impact duration from speed profiles. Source: Haule, H. J., T. Sando, R. Lentz, C.-H. Chuan, and P. Alluri, Evaluating the impact and clearance duration of freeway incidents Haule et al., 2019a

Shang, 2014a focused on the use of multiple types of distributions (Log-normal, Gamma, etc.) for the four-time intervals within the incident duration structure corresponding to: response team preparation time, response team travel time, incident clearance time, total incident duration time. They found the importance of different distributions to approximate different incident duration stages.

Response time (RT) is defined as a time interval comprised of both response team preparation time and travel time to the incident site. RT was modelled in Hou et al., 2013. In another study, recovery time was analysed on freeway segments in the Southeast Queensland (Australia) Hojati et al., 2014 and was derived using historical loop-detector-data and traffic incident characteristics at the time and location of the incident. The event of non-recurrent traffic congestion was detected based on the allowable percentage of speed decrease. The time interval of the incident was determined by forward and backward search in time for time intervals, when traffic speed was unaffected, which appear to be bounding for the traffic incident.

Research published in Zeng and Songchitruksa, 2010 includes the calculation of Recovery Time from the time required for the restoration of travel time during the affected traffic state to travel time during the normal traffic state. That research points to the possibility to use traffic flow data to uncover more precise traffic incident duration instead of relying on the definition given by the response team. As pointed out in Figure 2.4, we can derive the duration of selected phases of traffic incidents from traffic speed data (or possibly traffic flow data).

In conclusion, traffic incident duration consists of multiple phases. Nevertheless, data availability for the duration of such phases is rare to find. Duration of phases and total traffic incident duration (assumable, even more correct than recorded by response teams) can be extracted using traffic flow data assuming the data streams are reliable and free of anomalies or outliers that might affect precision and analysis.

## 2.2 Data sets and data availability

The availability of data related to traffic accidents in recent years has enabled a deeper understanding of the factors that lead to these incidents and their outcomes. This data, composed of information such as the location and time of an accident, the type of vehicle involved, the severity of the crash, casualty statistics and economic cost can provide valuable insights into the causes and outcomes of traffic accidents. By employing machine learning techniques, it is possible to make predictions about the risk of future crashes, classify accidents by severity and predict accident duration to develop strategies for mitigating the risk and severity of accidents. This kind of data analysis can help to inform policymakers and road safety management organisations on how to create safer roads and highways, as well as to make driving safer for everyone.

There are multiple publicly available datasets:

- National Highway Traffic Safety Administration's (NHTSA) Fatality Analysis Reporting System (FARS): This dataset contains detailed information on fatal motor vehicle traffic crashes in the United States occurred since 1975 National Highway Traffic Safety Administration (NHTSA), 2020.
- National Transportation Atlas Database (NTAD) from United States Department of Transportation's Bureau of Transportation Statistics (BTS) U.S. Department of Transportation (USDOT)/Bureau of Transportation Statistics (BTS), 2020: contains detailed information on non-fatal motor vehicle traffic crashes in the United States since 1994.
- European Commission's Road Safety Atlas European Commission, n.d. provides accident statistics for each European country using interactive maps and satellite images.
- UK Road Safety Statistics UK Government, n.d.: This dataset contains detailed information on fatal and non-fatal road traffic accidents in the UK since 1979.
- California Highway Patrol (CHP) Statewide Integrated Traffic Records System (SWITRS) California Highway Patrol (CHP), n.d.: a California-wide data set containing detailed information on motor vehicle collisions reported to California Highway Patrol. Accident report details contain data on the location, severity, road condition and victim data including age and degree of injury. Due to the extensive timeline and precision of reporting, this data set was previously used to analyse the effect of country-scale events on crash severity Waetjen and Shilling, 2021.
- World Health Organization's Global Health Estimates World Health Organization, n.d.: This dataset contains detailed global estimates on road traffic injuries, deaths, and disability-adjusted life years from 1990 to present.
- Australian Road Deaths Database (ARDD) Department of Infrastructure Regional Development and Cities, n.d. provides basic details of road traffic crash fatalities in Australia as reported by the police each month to the State and Territory road safety authorities. The data set includes information on fatal crashes: year, month, day of the week, time, location, crash type and vehicle type involved.



- Countrywide Traffic Accident Dataset (CTADS) - one of the biggest data sets on traffic accidents is , recently released in 2021 Moosavi et al., 2019a Moosavi et al., 2019b, which contains 1.5 million accident reports collected for almost 4.5 years since March 2016, each report containing 49 features obtained from MapQuest and Bing services. This data set was used previously to predict the accident duration Zhao and Deng, 2022; Grigorev et al., 2022a.
- Caltrans Traffic Performance Measurement System (PeMS) Choe, Skabardonis, and Varaiya, 2002 data set contains incident reports with a timeline of events and description of the incident (as a sequence of abbreviations) as it becomes available and status updates from a dispatch unit. Also, the dataset contains 5-minute aggregated traffic speed, traffic flow and traffic occupancy records as well as vehicle detector status. The availability of this data allows analyzing traffic incidents in conjunction to vehicle detector data. A brief description of the incident includes location, area, start time, duration and freeway ID.
- CompassIoT Compass IoT, n.d. is a data set and API of connected vehicle road trips across Australia which aggregates data from 64 different manufacturers across a number of global car brands. It started from 200,000 connected vehicles in 2018 to over 2.2 million trips and billions of data points in 2023. The data set has low data latency, where real-time API can yield data every 5 seconds. Data on braking acceleration and steering can allow identification of dangerous road conditions and risky behaviour characteristic to road accidents.
- TomTom TomTom, n.d. - historical traffic database with information on road speeds, travel times and traffic density. Allows customised queries for route and area analysis providing statistics on travel time, speed.

### 2.2.1 Characteristics of traffic incidents

Features used by research studies on traffic incident duration are very diverse. Some researchers didn't perform possible feature manipulations despite data availability (e.g. AM/PM peak hour binary feature or night time). Traffic incident research has the possibility to benefit from comprehensive feature derivation based on time, weather and road network data. By tracking the use of features and assessment of their significance, we can make decisions on the concentration of future efforts in data processing. Also, we can develop strategies for feature extraction considering techniques used by other researchers.

For example, in Li and Shang, 2014a authors use the feature "season" which is represented as a set of four discrete values (summer, winter, spring, autumn). The season can highly affect safety on the road due to weather effects: winter storms, ice on the road, rain showers and other environmental effects (which affect visibility, control and safety) linked to the time of the year. But in some papers Mihaita et al., 2019b authors don't use these features.

Roadway geometry is one of the critical factors which can affect the capability of a road system to withstand the incident impact Al-Bordiny, 2014. Various features are used to define the road structure and use machine learning models: road segment length, road segment centroids, gradient, curvature and general road density surrounding the event area Ziakopoulos, 2021. A comprehensive review of

spatial network (i.e. road network) theory applications provides an in-depth analysis of graph theory indices including betweenness centrality, ringness, route factor, detour index, and alpha index Barthélemy, 2011. All of these indices can be calculated for road networks and incorporated into machine learning pipeline. Accident risk is found to be increasing when traffic speed slows down while traffic density goes up for Yingtian Expressway Liu et al., 2021 which highlights the importance of speed-density diagrams for incident-related traffic state visualisation. In that research, no correlation was observed between traffic flow and crash risk.

Identification of road geometries is important for the analysis of incident occurrence. The impact of variables associated with crash frequency was found to be varying across parts of Tennessee, USA Mohammadnazar et al., 2021. The spatial analysis demonstrated that segment length and median segment width had the highest impact on crash frequency in eastern regions, while commercial land use had the highest connection to crash frequency in southern regions. Multiple studies have indicated that patterns and dependencies in the spatial and temporal dimensions are likely to exist, often represented as clusters or hot spots Al Hamami and Matisziw, 2021.

The following tables 2.1 and 2.2 include incident characteristics used among different research studies, with citing specific ones, unobserved or rarely observed in papers. For the general set of characteristics, readers can refer to Mihaita et al., 2019b; Li, Pereira, and Ben-Akiva, 2018b.

| Feature     | Values  | Reference                |
|-------------|---|--------------------------|
| Peak hour   | {0, 1}  | -                        |
| Weekday     | {0, 1}  | -                        |
| Weekend     | {0, 1}  | Javid and Javid, 2018    |
| Season      | { <i>winter, autumn, summer, spring</i> }                   | Li and Shang, 2014a      |
| Time of day | {0..23}   | Mihaita et al., 2019b    |
| Peak hours  | { <i>Off peak / AM peak(6 – 9 AM) / PM peak(3 – 6 PM)</i> } | Nam and Mannering, 2000b |
| Daytime     | { <i>Evening, night – time</i> }                            | Nam and Mannering, 2000b |

TABLE 2.1: Table of temporal features used to describe traffic incident

## 2.3 Incident duration modelling

Most current studies rely on methods for the classification and clustering of traffic conditions for incident detection. However, there are very few studies on traffic incidents involving methods for detecting anomalies (such as one-class SVM, isolation forests, etc). Non-recurring traffic incidents are rare and unusual in nature and therefore the detection of a traffic incident can be assessed as the task of detecting anomalies in traffic. By relying on anomaly detection methods, the incident detection system can be adapted to previously unseen situations. Thus, classification and evaluation of the applicable anomaly detection methods in comparison to well-established classification and regression methods will be carried out. Also, road situations detected as anomalous can be extremely valuable for further investigations in the duration of a freshly reported accident.

Incident duration can be modelled in terms of spatial relations (geometric placement of adjacent lanes, angle of adjacency, different parameters of lanes, including speed limits). Recent studies rely on reported incident parameters Mihaita et al., 2019a, Hamad, Khalil, and Alozi, 2019, but road topology can also play a significant role in the estimation of the incident probability (e.g. poorly designed

| Feature                               | Values   | Reference                                  |
|---------------------------------------|--|--|
| Incident type                         | { <i>Vehicle fire, out of gas, breakdown, etc.</i> } | KIM and CHOI, 2001                         |
| Number of vehicles involved           | {1..N}   | -  |
| Multiple vehicles involved            | {0, 1}   | Hojati et al., 2014                        |
| Type of vehicle involved #1           | { <i>Motorcycle, Van, Pickup</i> }                   | KIM and CHOI, 2001                         |
| Type of vehicle involved #2           | { <i>Large vehicle</i> }                             | Hou et al., 2013                           |
| Type of vehicle involved #3           | { <i>Truck</i> }                                     | Chung, Chiou, and Lin, 2015                |
| Location of incident on the road      | { <i>For freeways : ramp, left/right shoulder</i> }  | KIM and CHOI, 2001                         |
| Number of lanes                       | {1..N}   | Hojati et al., 2014                        |
| Link capacity                         | {N}  | Hojati et al., 2014                        |
| Average speed at the time of incident | {}   | Hojati et al., 2014                        |
| Number of affected Lanes              | {1..N}   | Mihaita et al., 2019b                      |
| All lines affected                    | {0, 1}   | Javid and Javid, 2018; Hou et al., 2013    |
| Incident Severity                     | {1..N}   | Mihaita et al., 2019b; Haule et al., 2019a |
| Lighting condition                    | { <i>day, night</i> }                                | Haule et al., 2019a                        |
| Secondary crash                       | {0, 1}   | Haule et al., 2019a                        |
| Fire, Injury                          | {0, 1}   | Hou et al., 2013                           |
| Fatality                              | {0, 1}   | Hojati et al., 2014                        |
| Traffic disrupted                     | {0, 1}   | Hojati et al., 2014                        |
| Traffic flow on adjacent lanes        | {N}  | Mihaita et al., 2019b                      |
| Medical required                      | {0, 1}   | Hojati et al., 2014                        |
| Rollover                              | {0, 1}   | Chung, Chiou, and Lin, 2015                |
| Weather #1                            | { <i>Windy, Clear, Rain</i> }                        | Alkaabi, Dissanayake, and Bird, 2011       |
| Weather #2                            | { <i>Sunny, Cloudy, Storm</i> }                      | Chung, Chiou, and Lin, 2015                |
| Weather #3                            | { <i>Rain, Snow, Wind, Fog</i> }                     | Nam and Mannering, 2000b                   |
| Position within road                  | { <i>Inner, Outer, Middle lane</i> }                 | Chung, Chiou, and Lin, 2015                |
| Lane number                           | {1..N}   | Mihaita et al., 2019b                      |
| <b>Other Features</b>                 | <b>Values</b>  | <b>Reference</b>                           |
| Distance from the city center         | { <i>Rkm</i> }                                       | Mihaita et al., 2019b; Hojati et al., 2014 |
| Traffic condition                     | { <i>congested, uncongested</i> }                    | -  |

TABLE 2.2: Table of features used to describe traffic incident across different studies

junction, wrongly imposed speed limits). According to Curiel, Ramirez, and Bishop, 2018, about 5% of the road junctions are the site of 50% of the accidents in the city of London. Thus, it seems reasonable to analyse incident duration and probability with consideration of the road topology. The task of predicting the duration of an incident is usually solved by using Machine Learning methods. Among these methods – tree based methods Ozbay and Kachroo, 1999, fuzzy logic Wang, Chen, and Bell, 2002, Bayesian networks Ozbay and Noyan, 2006a, artificial neural networks Barcellos et al., 2015, Alkheder, Taamneh, and Taamneh, 2017. And recently Ma et al., 2017 studied GBDT as a better performing method for incident duration prediction. Gaussian process regression and artificial neural networks were found to outperform tree methods and SVM in incident duration prediction Hamad, Khalil, and Alozi, 2019.

The majority of prior works has studied the prediction of incident duration on specific types of roads (freeways or motorways), where the data accuracy is higher than on arterial roads; as of 2018, very few applied the prediction strategies on normal arterial roads due to the high modelling complexity and a location mismatching; the majority of traffic incident duration analysis researches focus only on one type of road network (freeways, highways, etc) (Yu and Xia, 2012)-(Chung, Walubita, and Choi, 2011)-(Hojati et al., 2012)-(Zhan, Gan, and Hadi, 2011); this is revealed by a recent state of art published in (Li, Pereira, and Ben-Akiva, 2018b) which emphasises on the difficulty of solving this problem for arterial roads and the lack of studies in this field. Our study proposes a framework capable of predicting the incident duration regardless of the road network or its complexity.

Also, the estimation of incident duration can be reduced to the classification method Mihaita et al., 2019a. To do this, a specific threshold for the duration is set and a prediction is made whether the

incident will last longer than a specified time. Artificial neural networks show high average accuracy for the prediction of 4 types of incident severity relying on data on the state of the road (lane, condition of the roadway, weather, light, etc.), time and date. Overall accuracy between death, severe, moderate and minor severity accidents was found to be 69-72% Alkheder, Taamneh, and Taamneh, 2017.

Bayesian networks can produce interpretable models for incident injury severity prediction O textasciitilde na, Mujalli, and Calvo, 2011. Bayesian networks also outperform regression models in incident severity prediction (involving three severity indicators: number of fatalities, number of injuries and property damage) Zong, Xu, and Zhang, 2013. Interpretability is not a specific property of the tree models only and by using knowledge distillation one can extract tree rules from different prediction models (e.g. Bayesian network Park, Haghani, and Zhang, 2016a). It allows to represent the model as an interpretable decision tree and to estimate the feature importance.

### 2.3.1 Traditional accident modelling

Incident duration distribution has been modelled as log-normal Sullivan, 1997 and more recently as log-logistics distribution Chung, Walubita, and Choi, 2010, Smith and Smith, 2001, Smith and Smith, 2002. Log-logistic model has been used more extensively and found to have better goodness-of-fit than log-normal distribution. Also, there are various hazard-based models of traffic incident duration Nam and Mannering, 2000b, Hojati et al., 2013 which employ a hazard function to describe the conditional probability that an incident will end during any particular time interval given that it already lasted until the beginning of the interval. Recent studies also involve multi-component log-logistic models. In Zou et al., 2016a authors describe a g-component log-logistic model and in Li, Pereira, and Ben-Akiva, 2015c describe a competing risks mixture model which incorporates a multinomial log-logistic model.

An actual distribution estimation can give only approximate information on traffic incident duration. More than that, incident duration distribution is found significantly dependent on incident case parameters (e.g. day/night) Yang et al., n.d. Also, the duration of the incident may be affected by the chosen method of incident clearance Li, Pereira, and Ben-Akiva, 2015c. Mao et al., 2019 found that some incident parameters found to be important factors with different contributions to different types of accidents, including weather conditions, traffic density, time periods and incident location. Also, road factors were found to be affecting each one of the 4 incident types (rear-end, side wipe, collision with fixtures and rollover) in a different way. These findings draw incident duration distribution estimation as a complex problem dependent on many traffic flow and incident parameters.

In a recent study, Haule et al., 2019a, incident clearance time and the total impact duration were modelled using Weibull, log-normal, log-logistic distributions and compared using the Akaike information criterion (AIC) criteria; findings have revealed that log-logistic distribution was outperforming other distributions. As distribution utilisation is highly related to the specificity of each data set, for this study, in which we use three different data sets, we further apply a comparison among several distribution modelling choices by using the AIC criteria.

Also, there are various **hazard-based models** of traffic incident duration Nam and Mannering, 2000b; Hojati] et al., 2014, which employ a hazard function to describe the conditional probability that an incident will end during any particular time interval given that it already lasted until the beginning of the interval. Recent studies also involve multi-component log-logistic models. The authors in

Zou et al., 2016b describe a g-component log-logistic model and Li, Pereira, and Ben-Akiva, 2015c competing risk mixture model which incorporates a multinomial log-logistic model.

Incident duration can be modelled in terms of spatial relations (geometric placement of adjacent lanes, angle of adjacency, different parameters of street lanes, including speed limits). Recent studies rely on reported incident parameters Mihaita et al., 2019b; Hamad, Khalil, and Alozi, 2019, but road topology can also play a significant role in the estimation of the incident probability (e.g. poorly designed junction, wrongly imposed speed limits). According to Curiel, Ramirez, and Bishop, 2018, about 5% of the road junctions represent locations of 50% of the accidents in the city of London. Thus, it seems reasonable to analyse incident duration and probability with consideration of the road topology.

In papers attributed to the early 2000s, authors primarily used one or two distributions to fit traffic incident duration data. But starting from the 2010s, we can observe the use of different distributions within one research for the approximation of different phases of incident duration with comparing them according to the BIC score Li and Shang, 2014a; Alkaabi, Dissanayake, and Bird, 2011 - the Bayesian information criterion (BIC).

In the recent study, clearance duration and impact duration were modelled using Weibul, Log-normal, Log-logistic distribution Haule et al., 2019a:

- Log-logistic model outperformed Weibul and Log-normal models based on the comparison by AIC criteria.
- Incident impact duration was based not on incident durations reported by response teams, but estimated from historical speed data from BlueTOAD device pairs located on road segments.
- Hazard-based modelling approach allowed to estimate the impact of incident parameters on impact and clearance durations. Most of the parameters (night time, severity, EMS involvement, etc) were found to be affecting both durations in the same way but with a different degree of impact. But some characteristics demonstrated the opposite effect on duration (percentage of the lane closure, peak hour, summer/fall season, involvement of towing vehicles).

Incident impact modelled based on traffic flow data has the potential to be more accurate than reported by incident response teams. But it was not compared within this research.

The research underlines diversity in modelling accuracy for different stages of traffic incidents. Because of the significant observed difference in modelling accuracy for impact and clearance duration, authors also proposed to model each incident type (crashes, hazard, etc) separately to see how one model can perform better than another in each case.

The fact that different distributions can be used to approximate different phases of traffic incident Li and Shang, 2014a, implies that we can fit different distributions not only to phases but also split the dataset by specific variables (e.g. peak hour), which can also lead to different estimation in feature importance among resulting datasets (e.g. what is important for peak hour incidents, can be less important than for non-peak hour incidents).

### 2.3.2 Detection and estimation of traffic accident duration from traffic flow data

Main factors contributing to traffic jams are high traffic load (e.g. during peak hours), a bottleneck (a spatial aspect of road geometry) and local disturbances in the flow (e.g. actual traffic incidents which act as a trigger of traffic jam) Treiber and Kesting, 2013b. Incident detection systems can rely on offline or real-time data. The emerging approach is to use computer vision methods to detect traffic incidents using CCTV cameras. The spatial-temporal near-accident condition detection system has been recently proposed leveraging object detection, segmentation and tracking Huang et al., 2020. Advanced Driver Assistance Systems are intended to mitigate or prevent crashes by providing vehicle drivers with the necessary information to avoid collisions. To accomplish this task the truck driver behaviour (speed reduction) encountering vulnerable road users (e.g. cyclists at intersections) have been studied Schindler and Piccinini, 2021. The study of vehicle behaviour using GPS data can provide valuable insight into driver behaviour to further assist general drivers to avoid incidents. Geographically Weighted Poisson Regression (GWPR) models were used to model frequencies of harsh driving behaviour events, which were found to be positively correlated with segment length and presence of traffic lights and negatively with neighbourhood complexity (which is a density area in proximity of the event) Ziakopoulos, 2021. Road features can be used to predict the traffic accident risk since they affect the driver's behaviour. Traffic state identification systems related to traffic jams can rely on traffic speed and density data. The previously used algorithmic approach can allow to identification of interruptions and moving jams Liu et al., 2021.

## 2.4 Machine Learning in traffic accident analysis

The figure 2.5 illustrates the machine learning pipeline used for predicting the duration of traffic accidents. The pipeline consists of 1) data preprocessing (cleaning, data imputation, label encoding, and outlier detection), 2) feature transformation (Principal Component Analysis and Latent Dirichlet Allocation, Log-transformation of target variable), 3) feature selection (e.g. using correlation-based feature selection, univariate feature selection, recursive feature elimination), model training (e.g. using linear regression, support vector machines, k-nearest neighbors, decision trees, random forests, and neural networks), feature importance estimation (using Gini importance, permutation importance, or SHAP), model testing (including cross-validation, confusion matrix, and ROC curve), and results estimation (accuracy, precision, recall, F1-score, RMSE and MAPE). The model is then validated to ensure its accuracy and reliability. The pipeline shows a general way of predicting the duration of stochastic events using machine learning methods and can be used for other similar tasks.

The task of predicting the duration of an incident can be solved by using Machine Learning methods. Among these methods are: tree-based classification methods Smith and Smith, 2001; Ozbay and Kachroo, 1999, fuzzy logic Wang, Chen, and Bell, 2002, Bayesian networks Ozbay and Noyan, 2006a, linear regression analysis (LR) Khattak, Schofer, and Wang, 1995, artificial neural networks (ANN) Wang, Chen, and Bell, 2005; Alkheder, Taamneh, and Taamneh, 2017; Barcellos et al., 2015, support-vector regression (SVR) Wang, Ngan, and Yung, 2018a. Recently GBDT (gradient-boosted decision trees) have been revealed to be a better performing method for incident duration prediction Ma et al., 2017. Gaussian process regression and artificial neural networks were found to outperform tree methods and SVRs in incident duration prediction Hamad, Khalil, and Alozi, 2019.



Extreme Learning Machines (ELM) Huang, Wang, and Lan, 2011 - is a machine learning method, which incorporates a feed-forward neural network initialised with random weights and consequent training step based on produced random feature mapping, designed to avoid overfitting of neural network.

The following methods are commonly used for traffic incident duration modelling: a) gradient boosting decision trees - GBDT (Friedman, 2000) which rely on training a sequence of models, where each model is added consequently to reduce the residuals of prior models; b) extreme gradient decision trees - XGBoost (Chen et al., 2015) which finds the split values by enumerating all the possible splits on all the features (exhaustive search) and contains a regularisation parameter in the objective function; c) random forests - RF (Breiman, 2001) which applies a bootstrap aggregation (bagging, which consists of training models on randomly selected subsets of data) and uses the average (or majority of votes) of multiple decision trees in order to reduce the sensitivity of a single tree model to noise in the data; d) k-nearest neighbours - kNN (Fix and Hodges, 1951) which uses for the prediction on data points the majority of votes or the average from k closest neighbouring data points from the training set (based on a distance metric); e) linear Regressions - LR - a standard predictor using linear equations to model the relation between the features and the regression variable; f) light gradient boosted machines - LightGBM (Ke et al., 2017) which applies gradient boosting to tree-based models; it also uses a Gradient-based One-Side Sampling (GOSS) and excludes data points with small residuals for finding split value. These models can be used for both classification and regression problems (except logistic regression applied to classification only and linear regression to regression problem only).

One of the recent research studies Kuang et al., 2019a presented a two-step approach for traffic incident duration prediction. A cost-sensitive Bayesian network was used to perform binary classification of traffic incidents by choosing a threshold of 30 minutes and then performing regression for each class using the k-nearest neighbours approach. While the approach is functional, one major drawback of the classification problem is to manually choose the class split threshold, as it can lead to severe class imbalance; to overcome this issue, in our study, we perform both a fixed and a varying threshold set-up to find the best class balance for our classification models; even-more, we propose as well a comparison with a multi-class classification approach and debate on the benefits and drawbacks of using classifiers for such problems; we also enhanced more advanced regression models together with outlier removal procedures that would provide a better and more precise prediction of the incident duration precondition in minutes. Overall, the cost sensitivity of incorrect classification can be further extended to the cost-based regression metrics.

In one of the recent researches about the classification of driving state, multiple hyper-optimised ML models were tested, and the entire feature space was visualised using t-SNE(Yi et al., 2019). RandomForest provided the highest accuracy of prediction, but more advanced tree-based models exist that utilise gradient boosting, which we will be using in our research (e.g. gradient-boosted decision trees).

To verify the performance of advanced tree-based methods (such as LightGBM), additional conventional ML models can be used (Chen et al., 2020). LGBM can be compared with conventional ML and non-tree-based models (k-nearest neighbours, Logistic Regression).

The majority of prior works have studied the prediction of incident duration on specific types of roads (freeways or motorways), where the data accuracy is higher than on arterial roads; as of 2018,

very few applied the prediction strategies on normal arterial roads due to the high modelling complexity and a location mismatching; the majority of traffic incident duration analysis researches focus only on one type of road network (freeways, highways, etc) (Yu and Xia, 2012)-(Chung, Walubita, and Choi, 2011)-(Hojati et al., 2012)-(Zhan, Gan, and Hadi, 2011); this is revealed by a recent state of art published in (Li, Pereira, and Ben-Akiva, 2018b) which emphasises on the difficulty of solving this problem for arterial roads and the lack of studies in this field.

### 2.4.1 Classification and regression tasks in the incident duration prediction

The machine learning pipeline follows the general structure including the use of regression/classification models, feature correlation tests, data split according to train-test-validation schema, model calibration and validation, and prediction performance evaluation Gholami et al., 2020. One of the studies on accident hot-spot clustering Al Hamami and Matisziw, 2021 highlights that methods applied for vehicle accident analysis can be applied to various kinds of point-based events like crime, natural hazards, etc. A supervised machine learning model built upon extreme gradient boosting models has been used to predict rail-road accidents (derailments) and to rank, the importance of contributing factors using ANOVA and Gini criteria Bridgelall and Tolliver, 2021. Crash prediction models can be historical and real-time. The comprehensive review on real-time crash prediction approaches has been performed in Hossain et al., 2019.

#### Classification and regression definitions

Using all available data sets and the incident information, we first denote the matrix of traffic incident features as:

$$X = [x_{ij}]_{i=1..N_i}^{j=1..N_f} \quad (2.1)$$

where  $N_i$  is the total number of traffic incident records used in our modelling and  $N_f$  is the total number of features characterising the incident (severity, number of lanes, type, neighbourhood, etc.) according to each specific data set.

The estimation of incident duration can be reduced to the classification task Mihaita et al., 2019b. To do this, a specific threshold for the duration is set and a prediction is made whether the incident will last longer than a specified time. Artificial neural networks show high average accuracy for the prediction of four types of incident severity relying on data obtained from the state of the road (lane, condition of the roadway, weather, light, etc.), time and date. Overall accuracy between deadly, severe, moderate and minor severity accidents was found to be around 69-72% in Alkheder, Taamneh, and Taamneh, 2017.

For the incident duration classification problem the incident duration classification vector is defined as:

$$\begin{cases} Y_c = [y_i^c]_{i \in 1..N} & y_i^c \in \{0, 1\} \\ Y_{mc} = [y_i^{mc}]_{i \in 1..N} & y_i^{mc} \in \{0, 1, 2\} \end{cases} \quad (2.2)$$

where  $N$  is the duration of the traffic incident (in minutes),  $Y_c$  is the vector of binary values taking values in  $\{0, 1\}$ , and  $Y_{mc}$  is the vector of integer values for the multi-class classification problem definition, taking values in  $\{0, 1, 2\}$ . More specifically, a binary classification modelling has the purpose



of identifying short versus long-term incident duration, split by the incident clearance threshold  $T_c$ . Thus the incident duration classification task is to predict  $y_i^c$ , where  $Y_c$  takes one of the binary values:

$$\begin{cases} y_i^c = 0 & \text{if } y_i \leq T_c, \text{ short-term incidents} \\ y_i^c = 1 & \text{if } y_i > T_c, \text{ long-term incidents} \end{cases} \quad (2.3)$$

where  $T_c$  the incident duration threshold.

### Evaluation of prediction accuracy

To evaluate the regression model performance the most commonly used metrics are 1) Mean Absolute Percentage Error (MAPE) and 2) Root Mean Squared Error (RMSE) defined as:

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - F_i}{A_i} \right| \quad (2.4)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (A_i - F_i)^2} \quad (2.5)$$

where  $A_i$  are the actual values and  $F_i$  - the predicted values,  $n$  - number of samples.

The performance of classification models is most commonly evaluated using the Accuracy and F1-score and defined as:

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn}, \quad (2.6)$$

$$F_1 = 2 * \frac{precision * recall}{precision + recall}. \quad (2.7)$$

where  $tn$  represents true negatives,  $fn$  - false negatives,  $tp$  - true positives,  $fp$  - false positives.

For example, we refer to true positives the incidents which have been predicted to be in a specific class (say short-term) and indeed they were short-term upon validation, false positives the incidents which were predicted to be short-term but were not, etc.

F1-score is in general a better performance metric to use when there is an uneven class distribution (especially data sets with fewer incident records). This metric takes into consideration the total number of both false positives and false negatives together with true positives and true negatives.

Various metrics used for the regression tasks, RMSE and MAPE being the most common for the regression task (see Table 2.3).

### 2.4.2 Feature importance and feature selection

There are models of different complexity used to approximate traffic incident duration and duration of its phases. Khattak Khattak, Schofer, and Wang, 1995 used simple linear model to approximate clearance time, which can be defined as a function of incident parameters and coefficients:

$$clearancetime = A_1 * feature_1 + A_2 * feature_2 + \dots + A_N * feature_N \quad (2.8)$$

By using this simple model, we can easily determine the most important features by consequently removing them and estimating the resulting error (e.g. estimate change in mean squared error in relation to the removed feature). This method for determining feature importance is called Recursive Feature Elimination Guyon and Elisseeff, 2003 and is a general method which is applicable to different approximation models.

The Shapley Additive explanation (SHAP) Lundberg and Lee, 2017 provides more advanced approaches for feature importance estimation because it fuses estimation from multiple models trained across many different subsets (which selected both feature-scale and index-scale) of the dataset. These studies motivated the utilisation of the Shap Values for Feature importance ranking and selection and the analysis of such approach on three different data sets, all with different features and information.

The example SHAP plot is provided on Figure 3.14. Each point related to a feature is shown in and represents the SHAP value score (Oy-axis), coloured by its value (from low to high), while the Ox-axis shows the impact of that feature information on the entire prediction output.

It is generally not enough to use all the possible features for the regression analysis of traffic incident duration. Using a high amount of features combined with a small data set size can lead to over-fitting. Some features can be helpful or useless, more or less critical, while others do not impact the prediction results. By performing a feature importance analysis, we can recommend traffic management facilities record the most critical data and omit redundant data related to traffic incidents. Also, we can increase the precision of specific observations (e.g. weather conditions), which were found to play a significant role in some research studies (e.g. during summer and autumn seasons, response team preparation time was higher on freeways in Washington, USA in 2009 Hou et al., 2013, with no noticeable effect on clearance and response team travel time. Peak hours were the most influencing feature on response team preparation delay, which was found to be linked to response procedures (the goal of the response team was to resolve incidents during peak hours as soon as possible). A research study using Beijing traffic incidents data from 2008 Li and Shang, 2014a found the importance of "peak hour" value for the response team travel time and clearance time, but not for the intervention team preparation time. For example, one can use produced decision trees from the tree-ensemble model Chen et al., 2020. A data-driven approach can be used to perform information fusion from different sources Abou El Assad, Mousannif, and Al Moatassime, 2020, which involved the use of Gini-index extracted from Random Forests as a method to estimate feature importance. Nevertheless, the single random model can have a noticeable variance in data mapping when there is a weak connection between features and the target variable by making the feature importance value dependent on the random seed for the model.

A research study using Beijing traffic incidents data from 2008 Li and Shang, 2014a found the importance of "peak hour" value for the response team travel time and clearance time, but not for response team preparation time. Surprisingly, during summer and autumn, response team preparation time was higher, with no noticeable effect on clearance and response team travel time. The type of incident "overturned vehicle" found to have significant effect on response team preparation and incident clearance time with no effect on travel time. The most important factors for the total incident duration time were: 1) bike involvement, which can imply human injury and 2) night shift (10pm-6am), which was linked to higher incident severity and consequently higher clearance time.

One important conclusion is that rather than using the number of affected lanes Mihaita et al.,

2019b (which is only part of information related to road segment), we can also add a resulting value called “all lines affected” Javid and Javid, 2018 or even more informative - ratio of affected lanes (e.g. 50%) indirectly incorporating the number of lanes of a section, where incidents occurred into the model. Also, we can include the number of lanes for the each involved section as a feature.

Modelling of freeway incident response time in Washington State, USA in 2009 Hou et al., 2013 found morning peak hours as the most influencing feature on response team preparation delay, which was found to be linked to response procedures (goal of response team was to resolve incidents during peak hours as soon as possible). Response time was also lower during summer and winter. Because of the response strategy of the transport agency we can observe a clear difference in incident duration due to response priority. Priorities in the elimination of traffic incidents (which consider, for example, number of blocked lanes) can be defined as a part of traffic incident response guidelines. Priority of traffic incident elimination (and corresponding response team actions, like involvement of towing vehicle) can be used as a feature estimated from the traffic incident characteristics (and act as an additional information to the model).

Weekend and nightly incidents were also associated with significantly longer clearance and impact (including recovery) duration of traffic incidents on freeway Haule et al., 2019a. This observation was attributed to a lower number of staff on duty during weekends and nights Haule et al., 2019a; Hou et al., 2013.

The average traffic speed during a 60-min time interval was found to be statistically significant for the traffic incident duration modelling Hojati et al., 2014. One can use aggregated traffic flow data (e.g. represented as a speed of traffic or possibly traffic flow count) within specific time intervals as a feature.

Lighting conditions can be calculated much more precisely than just using binary day/night values Haule et al., 2019a. Using longitude, latitude and time, we can determine the angle between the sun and traffic direction at the time and place of the incident. Also, we can calculate precise lighting conditions based on the elevation of the Sun above the horizon during the time of the incident.

Effects of lighting conditions on driver behaviour were assessed using an interactive driving simulator Hong et al., 2014. According to the research, driver perception was found to be limited during night-time; drivers were also found to be limiting travel speed due to impaired visibility (which includes an incorrect and untimely perception of the road). Installation of road markers (to improve the perception of the road curve) and rumble strips were proposed. The impaired visibility during night-time and the proposed measures point that it is possible to derive visibility of the road structure from the point of traffic incidents and corresponding lighting conditions (distance of driver’s eyesight in relation to road segments).

### 2.4.3 Interpretable models

Interpretable machine learning was defined in Murdoch et al., 2019 as the extraction of relevant knowledge (knowledge which provides insights into the problem, which is then used to guide further actions and discovery) from a machine-learning model concerning relationships either contained in data or learned by the model (e.g. feature importance, dataset visualisation, learned relationships).

Decision-tree-based methods allow us to represent the model as an interpretable decision tree and logic of classification decisions and estimates of feature importance. Interpretability is not strictly

limited to tree models and by using knowledge distillation one can extract tree rules from different prediction models (e.g. Bayesian network Park, Haghani, and Zhang, 2016a ).

Bayesian networks can produce interpretable models for incident injury severity prediction O textasciitilde na, Mujalli, and Calvo, 2011. Bayesian networks also outperform regression models in incident severity prediction (involving three severity indicators: number of fatalities, number of injuries and property damage) Zong, Xu, and Zhang, 2013.

In conclusion, interpretable machine learning is a valuable tool for extracting relevant knowledge from machine learning models which later can be used by decision-makers. Decision tree-based methods and Bayesian networks are particularly useful for representing interpretable models and for estimating the importance of features.

#### 2.4.4 Imbalanced dataset classification

In the case of traffic incident duration prediction, we can have incidents with a duration of more than 60 minutes, which are rare extreme values. We also need to consider that traffic incident duration is modelled using Log-normal, Gamma and Weibull distributions as it has an asymmetric form and long-tailed distribution. That is why methods for the classification of imbalanced data sets can be useful for the classification of these rare incidents.

There are three main approaches for dealing with imbalanced classes:

- Under-sampling - when the size of a bigger subset is reduced to achieve class balance.
- Over-sampling - when the size of a smaller subset is increased (e.g. by synthetically generating additional samples) to achieve class balance.
- Combined approach - one can use both under- and over-sampling to achieve class balance between subsets Prati, Batista, and Monard, 2009.

There are different methods implementing each approach to deal with imbalanced data sets:

- random under-sampling Tahir, Kittler, and Yan, 2012, where the majority class is reduced by randomly choosing samples in order to achieve a higher class balance with the minority class.
- SMOTE - synthetic minority oversampling technique Chawla et al., 2002 which implies the oversampling of minority class by introducing synthetic examples by picking random points on line segments between k neighbours in multidimensional space.
- ADASYN: Adaptive Synthetic Sampling Approach for Imbalanced Learning He et al., 2008 - a method which uses distance distribution among k neighbours to decide the number of generated synthetic samples. The process of sample generation is the same as for SMOTE, with the difference that only the number of samples is adaptive (for more frequent distance among k neighbours, more samples will be generated).
- Utilisation of Support Vector Machines (SVM) to perform under-sampling because in SVM method only few support vectors are important to perform classification (Granular repetitive SVM under-sampling) Tang et al., 2009.

### 2.4.5 Boosted models and ensembles

Different machine learning techniques exist in the literature and each demonstrates different extrapolation performance; AdaBoostSVM and Extreme, for example, are machine learning approaches which has to property to make predictions with reduced overfitting.

AdaBoost is a meta-estimator, a machine learning method which implies a training of a set of weak classifiers with adaptive change to weights samples depending on the correctness of classification (after each boosting operation, weights of samples are adapted so next trained weak estimator gives higher priority to miss-classified samples) Schapire, 2013. Classifiers are then combined into weighted voting (linear) combinations.

AdaBoost and SVM were combined to solve the classification problem of imbalanced datasets Wang, Ngan, and Yung, 2018a using spatial-temporal traffic data for the case of automatic incident classification.

Extreme Learning Machines (ELM) Li et al., 2017 - is a machine learning method, which incorporates a feed-forward neural network initialised with random weights and consequent training step based on produced random feature mapping, designed to avoid overfitting of neural network. The method is two-step: 1) Neural network initialised using random weights (this way we perform feature mapping into ELM feature space) 2) Then Moore-Penrose generalised inverse performed on the hidden layer output matrix and solving feature classification problem using Gaussian estimation. ELM is very fast to train and provides better generalisation than Artificial Neural Network, also demonstrated remarkable efficiency in comparison to SVM, Naive Bayess and ANN on the traffic incident detection task on I-880 freeway in California Li et al., 2017.

### 2.4.6 Anomaly/outlier detection

The majority of studies in the literature have also concentrated on applying state-of-the-art machine learning models mostly for classifying the incident severity (Nguyen, Cai, and Chen, 2017) or their duration Li, Pereira, and Ben-Akiva, 2018b. However, very few have treated the problem of outliers or imbalanced data classes.

Many studies Wang, Ngan, and Yung, 2018a; Barcellos et al., 2015; Li et al., 2017 rely on methods for classification or clustering of traffic conditions for incident detection. However, there are very few studies on traffic incidents involving methods for detecting anomalies (such as one-class SVM, isolation forests, etc). Non-recurring traffic incidents are rare and unusual in nature and therefore the detection of a traffic incident can be assessed as the task of detecting anomalies in traffic. By relying on anomaly detection methods, the incident detection system can be adapted to previously unseen situations. Thus, classification and evaluation of the applicable anomaly detection methods in comparison to well-established classification and regression methods will be carried out. Also, road situations detected as anomalous can be extremely valuable for further investigations in the duration of a freshly reported accident.

Different anomaly detection methods will be used, including those which can produce a measure of an anomaly for each data point (One-Class SVM, Isolation Forest). This will be used to compare anomaly detection with regression methods (GBDT, ANN) for the task of incident probability estimation. Similarly, anomaly detection can be used in comparison with recently used classification

(GBDT, ANN) and regression (e.g. Gaussian process regression) methods, for the task of incident duration estimation. As stated in Ma et al., 2017, tree models perform badly at modelling incident duration with the long-tail distribution. Thus, anomaly detection methods can be used to model such kind of rare duration (which is placed in the long tail) as anomalies of duration.

The IsolationForest (IF) (Liu, Ting, and Zhou, 2008) is an outlier removal method, which uses forests of random split trees. For each tree, the method randomly selects a feature and random feature value. The dataset is divided into two parts in each step until each data point becomes “isolated” (split from the rest of the data). If the data point is an outlier, it will have a small tree depth (i.e. data point gets quickly separated from the rest by selecting values in just a few features). Tree depth is then averaged between all the “isolation” trees and considered an anomaly score (e.g. if the average tree depth for a point is 1.3, the point is easily separable after a small number of splits).

LocalOutlierFactor (LOF) (Breunig et al., 2000) is another outlier removal method, which estimates the anomaly score from the local deviation of density within the  $k$ -nearest neighbourhood. LOF relies on the calculation of local reachability density (LRD), which represents the inverse of the average reachability distance (RD) of neighbouring data points from the selected data point. Reachability distance (RD) represents the distance to the most distant neighbour within a  $k$ -sized neighbourhood ( $k$  is also a hyper-parameter). LOF of data point then represents the relation between LRDs of neighbours and its LRD and can take values above 1 (higher LRD than neighbours), below 1 (lower LRD than neighbours) and 1 (data point has the same density as neighbours). According to the LOF score, we can sort data points and select a specific per cent of data points, which have the highest LOF to be eliminated. LOF method relies on the fact that outliers belong to the area where the density of data points is low, while not outliers belong to the high-density area.

One-class SVM, Covariance estimator, Local outlier factor and Isolation Forest has been applied for outlier detection in machine learning pipeline to predict rail-road accidents Bridgelall and Toliver, 2021. Anomalies detected in historical traffic state data (traffic flow, traffic speed, occupancy) and vehicle trajectory Djenouri et al., 2019, which can be associated with disruptions produced by traffic accidents. Three main approaches for anomaly detection include the use of statistical models, similarity-based models (which rely on data difference measures and neighbourhood estimation methods to find outliers) and pattern-mining methods (which resolve the correlation problem of similarity-based models but are very time-consuming).

Multiple models in many research studies failed to predict extreme values for the traffic incident duration Li and Shang, 2014a; Shen and Huang, 2011. Machine learning method GBDT which demonstrates superior performance on a wide variety of tasks also known as failing at predicting very long incident duration, which is represented as extreme values within part of long-tailed distribution Ma et al., 2017.

In conclusion, anomaly detection can be an effective tool for improving the estimation of incident probability and duration, as it can identify and isolate rare events (or anomalous accident reports) and also can eliminate data records which can contribute to the bias of the prediction model (which can negatively affect model's performance Won, 2020). As such, anomaly detection can be a valuable addition to existing methods when attempting to model complex data sets.



### 2.4.7 Dimensionality reduction methods

Usually, traffic flow data is represented as traffic state readings across multiple (up to hundreds) vehicle detectors in the transport network. Options to use such a high-dimensional data set include: 1) use the closest to the point of interest readings (e.g. traffic flow on 5 closest road segments) Mihaita et al., 2019b 2) use all the data available on traffic flow in the network Mihaita et al., 2019b 3) Perform dimensionality reduction before modelling Wang, Ngan, and Yung, 2018a.

Principal Component Analysis (PCA) is a statistical method that uses an orthogonal transformation to transform multidimensional data into a sequence of linearly uncorrelated variables (components). Each resulting component is a linear combination of input features. The very first component has the greatest variance. The method was proposed in 1901 by Pearson Wold, Esbensen, and Geladi, 1987. PCA was used to reduce 23 dimensions of spatial-temporal signals to 2 dimensions for the task of Automatic Incident Classification proposed in Wang, Ngan, and Yung, 2018a.

Uniform Manifold Approximation and Projection (UMAP) McInnes and Healy, 2018 is a new dimension reduction method (2018), based on the use of Riemann geometry and algebraic topology. UMAP optimises the placement of data in a small dimension space to minimize cross-entropy between two topological representations. The disadvantage of the method is its non-interpretability and non-invertibility (one cannot inverse mapping to feature data). UMAP can be used both for dimensional reduction and data visualisation. To perform clustering and modelling of data points in reduced dimensions and to perform the data visualisation, it is necessary to use different sets of options (e.g. minimal distance between points in reduced dimension needs to be set to zero when preparing data for clustering).

Dimensionality reduction methods can help when we have multidimensional data with a lot of dimensions. A lesser number of dimensions can provide simpler and faster models or allow to process of large data sets within computational resource constraints.

### 2.4.8 Summary on the use of Machine Learning models

The majority of studies rely on Support Vector Machines or Naive Bayes for traffic incident duration prediction (see Table 2.4). The use of more advanced methods like XGBoost or GBDT is rare which is surprising given their effectiveness. This can be explained by a generally slow attribution of both sophisticated models and data-driven approaches to the traffic accident research which we observe in the literature. Also, the complexity of machine learning pipelines has increased in recent years due to the need for the incorporation of more data science knowledge to merge with traditional transport modelling techniques.

Another finding is that various Machine Learning pipeline elements (like dimensionality reduction or feature selection) are rarely used across incident duration prediction studies (see Table 2.5). This reflects a lack of advanced machine learning techniques that can be explored for incident modelling. The lack of popularity among these pipeline approaches leaves room for more innovative ideas for data filtering and feature ranking from the beginning of the incident duration prediction modelling. The use of SHAP and feature selection is found to be lacking in studies, while it may provide a list of entries in the incident reporting form with the highest contribution to the accuracy of the traffic incident duration prediction. An additional description of the most important features may provide

a further increase in prediction accuracy. Bi-level frameworks allow the separation of the task of the incident duration prediction into incident duration classification and regression tasks. Allocating different kinds of models to each task may improve the overall model performance. The use of bi-level frameworks as a technique for prediction performance improvement is also found to be rare. Dimensionality reduction had a low relevance in multiple prior studies due to the small size and low dimensionality of incident reports (see the table of data set sizes used in years prior to 2018 Li, Pereira, and Ben-Akiva, 2018b), but with current advancements in traffic data collection (see Section 2) and significant consequent increase in amount and variety of data collected, we see a rise in the relevance of data pre-processing methods.

## 2.5 Deep Learning in traffic accident analysis

Another example of classification task in accident modelling includes accident detection and accident risk prediction, which can rely on high-resolution traffic data (speed, occupancy, volume) and use Deep Learning methods like Convolutional Neural networks Huang, Wang, and Sharma, 2020. Variational Long Short-Term Memory Encoder has been proposed in Farahani et al., 2020 to perform a short-term traffic flow prediction. It demonstrated better performance than LSTM and Stacked Autoencoder models. Short-term (5 to 30 minutes), mid-term (30 minutes to multiple hours) and long-term prediction (day to multiple days) capabilities of the method were explored. The methodology is in general intended to predict the traffic state using historical data.

Deep learning is becoming increasingly important in the field of traffic accident modelling. By leveraging the power of artificial intelligence and machine learning, deep learning algorithms can be used to detect patterns in historical accident report data that may be indicative of future accidents Nguyen et al., 2018. This methodology can help to develop more effective strategies for preventing accidents, as well as allow traffic management agencies to analyze existing road structures for the probability of accident risk. Deep learning has been used to develop predictive models that can identify accident-prone areas Ren et al., 2018 and predict accident risk based on driver behaviour Shi et al., 2018, as well as to predict the impact of various events on traffic flow Yu et al., 2017.

Overall, the majority of studies implementing Deep Learning techniques for traffic incident duration prediction rely on classic ANN/MLP and rarely use recurrent or convolution networks (see the summarising Table 2.6 of the most popular deep learning techniques used so far). The use of recurrent networks is mostly attributed to the analysis of textual incident reports or messages over social networks. This short summarising of Deep learning shows that there is a significant gap in the current incident modelling to leverage such powerful techniques and this is mostly related to the data availability - traffic flow counts, speed, details traffic signal control, etc. There are different approaches that can be used to improve the incident duration prediction and in the following subsections we detail the most common ones as follows.

### 2.5.1 Spatial-temporal models for traffic incident modelling

In the past years, traffic accident research has seen an increased use of data-driven methods. Different problems were addressed, including: 1) traffic accident duration prediction methods (see methods in Li et al., 2018), 2) accident detection methods (see Parsa et al., 2019), 3) estimation of severity



(reader can refer to ) Singh and Yadav, 2021, and more recently, 4) the development of spatial-temporal modelling methods which have allowed to perform accident risk prediction using high-dimensional spatial, semantic and temporal data sets (see the work of Wang et al., 2021b). The use of such methods has enhanced the automated analysis of traffic data together with the increasing number of publicly available data sets. Traffic accident risk prediction allows to: a) detect high-risk areas within a traffic network, which may facilitate the decision-making inside traffic management authorities, b) to allocate resources and assess the road design to reduce the number of accidents in the future, c) to predict timely high-risk situations on the road and d) to allow implementation of risk-reducing traffic management strategies.

One of the first works on traffic accident risk prediction using Deep Learning has been performed with human mobility data using a Stack Denoise Autoencoder (SDAE) on the Japan traffic network Chen, Chen, and Hsieh, 2016, but traffic flow and time-related matters (including periodicity) were not considered. Another research Ren et al., 2017 relied on the LSTM network to improve the risk prediction in comparison to SAE by considering, in addition, the air quality, traffic flow and weather data, represented as short-term and periodic components. Authors in Zhou et al., 2020a proposed also a Coarse and Fine grained prediction on the target accident risk map. RiskOracle Zhou et al., 2020b relied on a Graph-Convolution network, utilizing hierarchical coarse-to-fine modelling and proposing minute-level predictions in comparison to day-level Yuan et al., 2018 and hour-level Chen, Chen, and Hsieh, 2016. In Yuan et al., 2018 authors have constructed over the ConvLSTM by highlighting the spatial heterogeneity problem and proposing an ensemble of region-specific ConvLSTM models (Hetero-ConvLSTM); they considered weather, the environment and the road condition in Iowa, US for over 8 years of observations, but points of interest (POIs) were not considered. Semantic features, coarse and fine-grained risk maps were considered in Wang et al., 2021a, where also Graph-convolution neural networks and attention-based LSTMs were used. A more recent work in Wang et al., 2021b represents the State-of-Art (SoTA) in the field of accident risk prediction, where the authors proposed a weighted loss function to address the zero-inflated issue (increase in the number of zero-risk grid cells due to the increase in the granularity of predictions) and made an ensemble of models by processing semantic and geo features.

Deep learning algorithms can provide a better understanding of accident risks and impacts in traffic networks in which drivers are operating, allowing traffic management agencies to develop more effective strategies for preventing accidents or mitigating their impact. Additionally, deep learning algorithms have the potential to improve the accuracy and speed of emergency response. However, the use of deep learning in traffic accident prediction is not without drawbacks, such as the need for large amounts of data and the potential for bias in the results due to the algorithm's reliance on existing data Tommasi et al., 2017 (these algorithms may need to be tested against extrapolation between time intervals, different areas and reporting source). Nevertheless, deep learning methodology allows to develop important tools for traffic authorities, as it can provide valuable insights into the causes of accidents, and data patterns which can point to potential risks and impacts and in total, help traffic management authorities to develop better strategies for preventing, responding and reducing the impact of traffic accidents.

### 2.5.2 Textual Accident Description analysis

Traffic accident reports usually contain a textual description of the accident Moosavi et al., 2019a Moosavi et al., 2019b. In recent years, multiple systems were presented to detect traffic accidents using text analysis of social networks content Vallejos et al., 2021; Salas, Georgakis, and Petalas, 2017; Ali et al., 2021. Various methods were also proposed for traffic-related sentiment analysis of social networks: sentiment classification using ontology and latent Dirichlet allocation Ali et al., 2019, the use of gated recurrent unit (GRU) model and generative adversarial networks to estimate the traffic information sentiment Cao, Wang, and Lin, 2018. Overall, sentiment analysis has been performed for various traffic rules including 'yellow light rule' using social network analysis Cao, Wang, and Lin, 2018; Lu et al., 2021.

Also, accident reports can contain a category and subcategory definition of the accident (e.g. types of vehicles involved, multiple or single-vehicle crash, etc). The unique property of text description of an accident is that it can contain information regarding event categories not predefined in the accident reporting form Vallejos et al., 2021.

A typical pipeline for textual description preprocessing includes Ali et al., 2021: 1) Tokenization - text being split into a list of words called tokens, 2) Stop-word removal - the removal of pronouns, prepositions, symbols and articles not providing any valuable information for accident description, 3) Lemmatization and stemming - words are reduced to their base form (e.g. involved -> involve, injuries -> injury, reported -> report) or to their root form (e.g. injuries, injury -> injur) 4) Case-conversion - text is converted into lower case, where the difference between uppercase and lowercase words is not relevant 5) Part-of-speech (POS) tagging - each word gets its type associated to it (e.g. traffic -> noun, stop -> verb), 6) text representation conversion, which relies on a Bag-Of-Words (BOW) representation (each word is represented as a one-hot encoded vector) or on a neural-network-based word embedding method like Word2Vec or FastText, which capture semantic similarities between words. The Word2Vec approach has a significant limitation - the inability to represent a new word which was absent in the training data set with a vector. The FastText resolves this issue by representing each word as a sum of related n-gram vectors.

After the data preparation and representation conversion, various recurrent models are then used to perform tasks related to text analysis.

Incident description features were used in topical text modelling Das, Mohanty, and Bhattacharyya, 2019. Previously, the LSTM architecture has been used for the task of detection of incidents from social media data Zhang, Chen, and Zhu, 2018. LSTM was also successfully used for stock price time-series prediction Sen and Dhar, 2018, making it applicable for the modelling of traffic flow/speed time-series data.

Text analysis of accident reports is vital for understanding the underlying causes of traffic accidents. By providing insight into the information provided to describe accidents, text analysis can help to identify dangerous conditions that lead to accidents. Traffic accident descriptions can contain inaccuracies due to human factor (e.g. inaccurate accident timeline), which highlights the importance of automated accident detection and timeline estimation from traffic state data Taghipour et al., 2022. Ultimately, text analysis is a powerful tool for gaining even deeper insight into the causes of traffic accidents unconstrained by accident reporting forms and developing strategies to reduce them.

## 2.6 Conclusion

Today, we can use several methods for solving the incident duration prediction, among which we name the spatial-temporal incident impact analysis and traffic simulation. However, each of the above topics present several challenging aspects which can range from data collection, cleaning, anomaly detection, up to strategies of using the best artificial intelligence model, or various techniques to model and estimate the impact of traffic incidents across the network and regions. This motivated our current literature review which was organised to provide the reader with a systematic review and understanding of the complexities of each modelling step.

### 2.6.1 Summary of challenges and gaps

Traffic accident analysis has the following challenges which we address in this review:

- Traffic accident reports contain multiple characteristics of a traffic accident. Each characteristic can have various effects on accident modelling performance. To solve the problem of determining which data to collect and what details (features) we need to use feature importance estimation methods.
- The task of predicting traffic accident duration group (e.g. predicting if an accident will be short-term or long-term) can create a problem of imbalanced classification due to uneven number of accidents in duration groups. To solve this issue, multiple approaches for imbalanced classification exist including the use of specific models, metrics and data processing techniques.
- Traffic accident duration distribution usually follows log-normal or log-logistic distribution which is a skewed distribution. Machine learning models show better performance with normally distributed predicted variables. Therefore, the target variable needs to be processed to enhance predictions.
- Accident reports also may contain reporting-specific errors or anomalies in reporting. In general, the outlier removal procedure improves the performance of machine learning models.
- Accident reports can represent very large data sets with high dimensionality. For example, CTADS contains 1.5 million accident reports, 49 features each. When working with high-dimensionality data it may be necessary to use dimensionality reduction to reduce the model training time and memory requirements.
- In particular cases, data availability on accidents can be low. In that case, we need to seek ways to improve the extrapolation performance of our models. The extrapolation ability and noise resistance of ML models can be improved by using model ensembles and bi-level frameworks.
- The specifics of accident report data sets is the presence of textual accident descriptions, which can contain valuable information to enhance the prediction performance of accident modelling. Various natural language processing techniques are of high importance for the task of utilizing accident description.

- Historical traffic speed or traffic flow data can be available for accident reports. In this case, time-series modelling techniques can be utilized to extract useful information to enhance the performance of accident modelling.
- One of the main challenges in incident identification from traffic flow is to provide descriptive statistics for abrupt changes in traffic state Liu et al., 2021.
- There are multiple novel machine learning and deep learning models, which have not been used in traffic accident modelling so far, but many hybrid or advanced frameworks can still be applied to enhance the performance of the incident duration modelling

Traffic jam identification methods (which rely on vehicle detector data) can be further extended from the use of algorithms to the use of machine learning methods intended for time series processing. Further research on the topic of incident-related traffic state identification can be performed using merged data sets of traffic accident reports and traffic states (flow, speed, occupancy) recorded in their proximity Liu et al., 2021.

The use of anomaly detection methods in urban traffic data was found to be seldom Djenouri et al., 2019 which implies that future research can be performed in that direction. The use of Machine Learning methods in transportation was found not being used to its full advantage Behrooz and Hayeri, 2022: 74% of papers were found to be relying on prediction methods like XGBoost, Random Forest, LSTM and Multilayer Perceptron with only minimal use of sophisticated ML methods.

The main limitations of using Deep Learning models for traffic incident duration modelling are: 1) data availability: deep learning models require large volumes of data for training, which can be difficult to access or cannot be provided by the traffic management authorities due to privacy or security concerns, 2) Data Quality: deep learning models are sensitive to data quality, including outliers, missing values, and user-input errors (like incorrect labelling, misreported incident duration), 3) interpretability: deep learning models are often represented as "black box" models (which means that relationships developed between inputs and outputs are hidden inside the model) Olden and Jackson, 2002, making it hard to understand how the model arriving at a particular result, which may limit the model deployment due to possible model bias towards data. The absence of interpretability is particularly critical since black-box models can't be considered reliable in traffic safety applications.

## 2.6.2 Future research directions in incident modelling

The application of traditional clustering methods can provide insights on spatial-temporal patterns and hot spots of traffic accidents Al Hamami and Matisziw, 2021. The evolution of cluster size over time can provide valuable insight into contributing factors, which lead to accident hot spot appearance, disappearance, growth and decline. The procedure of accident hot-spot detection and their evolution prediction can provide important information to traffic management authorities. Runtime performance of the online algorithms to find sub-trajectory outliers (which may help to detect accidents in real-time) was found to be low and may require the development of more efficient methods Djenouri et al., 2019.

The following topics of future research can be addressed:

- Data set integration and fusion models. There are various data sets which exist adjacent in time and space to incident reports like PeMS, which contains data on traffic flow, speed and

occupancy in the proximity of traffic accident reports from the CTADS data set, which may be used to enhance traffic incident duration modelling. The availability of multiple data sets of different types that describe the accident may require the use of data fusion models Wang et al., 2021b and/or feature embedding methods Grigorev et al., 2022c.

- The use of textual data: the textual incident report can contain information unconstrained to reporting form, which may be used to enhance the incident duration prediction accuracy. This approach requires knowledge and techniques from the field of Natural Language Processing.
- The use of sophisticated ML and DL models. Multiple studies on traffic crashes indicate nonlinear relationships and threshold effects between independent variables and dependent variables Parsa et al., 2020; Yang, Chen, and Yuan, 2021. Also, it was highlighted by one of the previous reviews that this study area relies mostly on classical ML models,
- The use of sophisticated ML pipeline elements like anomaly detection, hyperparameter optimization and dimensionality reduction, oversampling and undersampling found to be rarely used and may enhance the incident duration prediction performance
- The effect of combined use of Traffic Incident Duration Prediction with Variable Message Signs (VMC) on traffic flow can be studied further to assess the study the potential of VMS use for traffic incident impact mitigation Ghosh, 2019
- Further study on real-time incident reporting data can be performed using PeMS data set, which includes a timeline of textual incident description availability over time
- The use of feature importance estimation techniques to assess the impact of specific reported values on the incident duration prediction accuracy. In particular, the effect of weather conditions on incident duration can be studied Hamad et al., 2020a
- The requirement of extrapolation tests: the deployment of the incident duration prediction model requires multiple considerations like model bias for data or time/space extrapolation performance Hamad et al., 2020a.
- The rise in relevance of advanced data pre-processing methods. In particular, dimensionality reduction techniques become more relevant due to the significant increase in the amount and variety of data collected across traffic networks.

In conclusion, traffic incident duration prediction is a complex and important task, which may benefit from further research with the use of sophisticated models of artificial intelligence. Intelligent models, such as those based on machine learning, can provide predictions with high accuracy. The use of these models can help traffic management authorities to improve traffic flow and reduce the impact of traffic incidents. Further research is needed to improve the accuracy of these models, such as data set integration, complex and hybrid ML and DL models, the use of textual data, anomaly detection, and hyperparameter optimization. This research has the potential to enhance traffic incident duration prediction performance, ultimately leading to improved traffic flow and reduced impact from traffic incidents.

## **Acknowledgments**

This work has been done as part of the ARC Linkage Project LP180100114. This research is funded by iMOVE CRC (project 5-013) and supported by the Cooperative Research Centres program, an Australian Government initiative."

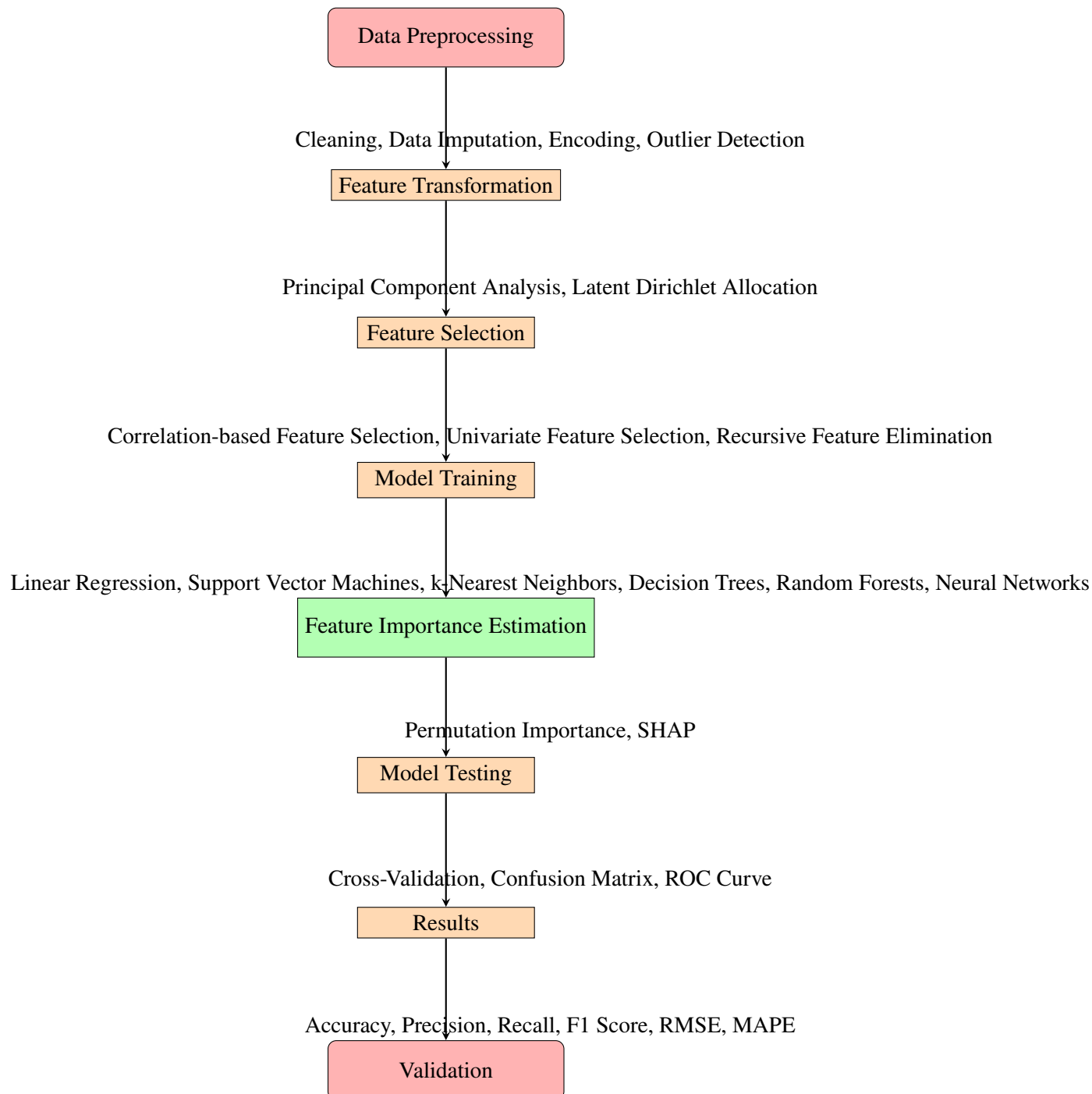


FIGURE 2.5: Machine Learning Pipeline for Traffic Accident Duration prediction

| Metric    | Studies   |
|-----------|---|
| MAE       | Valenti, Lelli, and Cucina, 2010a Yu et al., 2016 Ghosh et al., 2016 Wei and Lee, 2007 Shang, Xie, and Yu, 2022 Tang et al., 2020a Wang, Li, and Guo, 2018 Pereira, Rodrigues, and Ben-Akiva, 2013 Lee and Wei, 2010a Hamad et al., 2020a Park, Haghani, and Zhang, 2016b Li and Shang, 2014b Mohammed, Abdullah, and Al Hussaini, 2021 Al-Najada and Mahgoub, 2017 Hamad, Khalil, and Alozi, 2020 Zou et al., 2021   |
| RMSE      | Valenti, Lelli, and Cucina, 2010a Yu et al., 2016 Ghosh et al., 2016 Wei and Lee, 2007 Shang, Xie, and Yu, 2022 Tang et al., 2020a Wang, Li, and Guo, 2018 Li, 2015 Hamad et al., 2020a Li, Pereira, and Ben-Akiva, 2015c Li and Shang, 2014b Mohammed, Abdullah, and Al Hussaini, 2021 Al-Najada and Mahgoub, 2017 Hamad, Khalil, and Alozi, 2020 Li et al., 2020b Grigorev et al., 2022c Ozen et al., 2019 Lin and Li, 2020 Zou et al., 2021 Ghosh and Dauwels, 2022 Araghi et al., 2014 Grigorev et al., 2022a   |
| MAPE      | Valenti, Lelli, and Cucina, 2010a Yu et al., 2016 Ghosh et al., 2018 Li, Pereira, and Ben-Akiva, 2015a Wei and Lee, 2007 Zou et al., 2018a Kalair and Connaughton, 2021 Haule et al., 2019b Shang, Xie, and Yu, 2022 Tang et al., 2020a Wang, Li, and Guo, 2018 Li, 2015 Reza and Pulugurtha, 2019 Pereira, Rodrigues, and Ben-Akiva, 2013 Lin, Wang, and Sadek, 2016 Lee and Wei, 2010a Mihaita et al., 2019b Li, Pereira, and Ben-Akiva, 2015c Li and Shang, 2014b Mohammed, Abdullah, and Al Hussaini, 2021 Li, Pereira, and Ben-Akiva, 2018b Tang et al., 2020b Grigorev et al., 2022c Ozen et al., 2019 Lin and Li, 2020 Zou et al., 2021 Ghosh, 2019 Ghosh and Dauwels, 2022 Araghi et al., 2014 Grigorev et al., 2022a Chung, 2010 Zou et al., 2016b Kuang et al., 2019b |
| MSE       | Al-Najada and Mahgoub, 2017 Grigorev et al., 2022c  |
| AUC       | Zheng et al., 2021 Motamed et al., 2016 Zhu et al., 2021  |
| Recall    | Saracoglu and Ozen, 2020 Mihaita et al., 2019b Grigorev et al., 2022a   |
| Precision | Saracoglu and Ozen, 2020 Mihaita et al., 2019b  |
| F1        | Saracoglu and Ozen, 2020 Mihaita et al., 2019b Grigorev et al., 2022a   |

TABLE 2.3: Metrics used across reviewed papers.



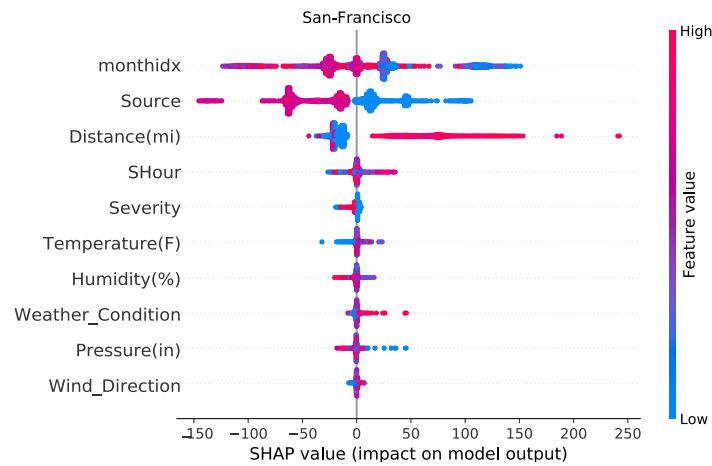


FIGURE 2.6: Feature importance for All-to-All regression using XGBoost for San-Francisco, USA Grigorev et al., 2022a

| Method                          | Studies  |
|---------------------------------|--|
| Random Forest                   | Shang et al., 2019Motamed et al., 2016Mohammed, Abdullah, and Al Hussaini, 2021  |
| XGBoost                         | Mihaita et al., 2019bTang et al., 2020bParsa et al., 2020Grigorev et al., 2022cGrigorev et al., 2022a  |
| Support-vector machines (SVM)   | Yu et al., 2016Zheng et al., 2021Shang et al., 2019Huang et al., 2020Tang et al., 2020aMotamed et al., 2016Park, Haghani, and Zhang, 2016bMohammed, Abdullah, and Al Hussaini, 2021Li, Pereira, and Ben-Akiva, 2018bHamad, Khalil, and Alozi, 2020Lin and Li, 2020Hossain et al., 2019Ghosh et al., 2016Shang, Xie, and Yu, 2022Mohammed, Abdullah, and Al Hussaini, 2021Tang et al., 2020bWu, Chen, and Zheng, 2011Ghosh, 2019Ghosh and Dauwels, 2022 |
| Linear Regression               | Ghasri et al., 2016Pereira, Rodrigues, and Ben-Akiva, 2013Lin, Wang, and Sadek, 2016Ozen et al., 2019Lin, Wang, and Sadek, n.d.Kuang et al., 2019b   |
| Naive Bayes                     | Lu, 2021aOzbay and Noyan, 2006bZheng et al., 2021Shang et al., 2019Tang et al., 2020aMotamed et al., 2016Park, Haghani, and Zhang, 2016bLi, Pereira, and Ben-Akiva, 2018bTang et al., 2020bKim and Chang, 2012Zou et al., 2021Ghosh and Dauwels, 2022Lin, Wang, and Sadek, 2015Wang et al., 2022Hossain et al., 2019Kuang et al., 2019b  |
| Decision Tree                   | Saracoglu and Ozen, 2020Ozbay and Noyan, 2006bWang, Li, and Guo, 2018Hamad, Khalil, and Alozi, 2020  |
| Gradient-boosted Decision Trees | Mihaita et al., 2019bGrigorev et al., 2022cGrigorev et al., 2022a  |
| K-Nearest Neighbours            | Lu, 2021aLin, Wang, and Sadek, 2015  |

TABLE 2.4: Most popular Machine Learning methods used across reviewed papers

| Method                               | Studies  |
|--------------------------------------|--|
| Principal component Analysis (PCA)   | Bridgelall and Tolliver, 2021 Wang, Ngan, and Yung, 2018b  |
| Linear discriminant analysis (LDA)   | Li, Pereira, and Ben-Akiva, 2015a Shang, Xie, and Yu, 2022 Pereira, Rodrigues, and Ben-Akiva, 2013 |
| SHapley Additive exPlanations (SHAP) | Kalair and Connaughton, 2021 Parsa et al., 2020 Grigorev et al., 2022a                             |
| Feature Selection                    | Vlahogianni and Karlaftis, 2013 Shang et al., 2019 Lee and Wei, 2010a                              |
| Clustering                           | Ghosh et al., 2016 Kalair and Connaughton, 2021 Ghosh, 2019 Lin, Wang, and Sadek, 2015             |
| Ensemble                             | Motamed et al., 2016 Hamad, Khalil, and Alozi, 2020  |
| Bi-level frameworks                  | Kuang et al., 2019a Grigorev et al., 2022a   |

TABLE 2.5: Most popular Machine Learning pipeline elements used across reviewed papers.

| Network type                       | Studies  |
|------------------------------------|--|
| Artificial Neural Network (ANN)    | Valenti, Lelli, and Cucina, 2010a Yu et al., 2016 Wei and Lee, 2007 Shang et al., 2019 Motamed et al., 2016 Pereira, Rodrigues, and Ben-Akiva, 2013 Lee and Wei, 2010a Hamad et al., 2020a Mohammed, Abdullah, and Al Hussaini, 2021 Li, Pereira, and Ben-Akiva, 2018b Hamad, Khalil, and Alozi, 2020 Chang and Chang, 2013 Zhu et al., 2021 Grigorev et al., 2022c Lin and Li, 2020 Ghosh, 2019 |
| Multilayer Perceptron (MLP)        | Ghosh et al., 2016 Shang, Xie, and Yu, 2022 Mohammed, Abdullah, and Al Hussaini, 2021 Zhu et al., 2021 Ghosh, 2019   |
| Recurrent Neural Network (RNN)     | Zhu et al., 2021 Wang et al., 2022   |
| Long Short-Term Memory (LSTM)      | Shang, Xie, and Yu, 2022 Zhu et al., 2021 Grigorev et al., 2022c Ghosh, 2019 Wang et al., 2022   |
| Convolutional Neural Network (CNN) | Kalair and Connaughton, 2021   |

TABLE 2.6: Most popular Deep Learning methods used across reviewed papers

## **Chapter 3**

# **Incident duration prediction using a bi-level machine learning framework with outlier removal and intra-extra joint optimisation**

### 3.1 INTRODUCTION

#### 3.1.1 Context

Traffic congestion is a significant concern for many cities around the world. Congestion arises due to various factors, including increased population, workforce concentration in central areas, or the lack of efficient public transport modes. Two forms of congestion are typically predominant: a) recurrent traffic congestion during peak hours when traffic demand exceeds the road capacity, and b) non-recurrent traffic congestion caused by unplanned events such as car accidents, breakdowns, weather, public manifestations etc. Previous studies have shown that almost 60% of traffic congestion is due to non-recurrent incidents with a stochastic behaviour in space and time Schrank and Lomax, 2002. In Australia, the number of road deaths per year has been reduced by 70% since the 1970s. However, the annual economic cost of road crashes was estimated at \$27 billion per annum in 2017 (Government, 2017). Traffic Incident Management Systems (TIMS) collect data on traffic incidents, including information on different incident duration factors. Accurately predicting the total duration shortly after an incident took place could save operational costs and end-user time (through affecting the route planning). Moreover, the clearance time of accidents is highly related to the ongoing traffic congestion and several external factors with different weights of importance. Therefore, it is essential to estimate the incident factor importance to improve the accuracy of predictions. Most prior studies related to this topic concentrated on testing different machine learning models on specific road types like freeways or highways and focused primarily on different phases of the incident duration such as clearance time, recovery time, and the total incident duration Li, Pereira, and Ben-Akiva, 2018a. There is currently a lack of an advanced approach that can be applied on all road types, for all accident types and across various countries with different driving behaviour.

#### 3.1.2 Challenges and contribution

The accuracy of predicting the incident duration is often determined more by the modelling methodology, the feature construction, and the result interpretation rather than by the model in use. In this work, we address several open questions or challenges concerning the prediction of the traffic incident duration.

**The first** challenge is to develop a universal bi-level framework applicable to different incident data sets reported on various road network layouts. The majority of prior works had studied the prediction of incident duration on specific types of roads (freeways or motorways) (Yu and Xia, 2012)-(Chung, Walubita, and Choi, 2011)-(Hojati et al., 2012)-(Zhan, Gan, and Hadi, 2011), where the data accuracy is higher than on arterial roads; as of 2018, very few applied the prediction strategies on normal arterial roads due to the high modelling complexity and a location mismatching; the majority of traffic incident duration analysis studies focus only on one type of road network (freeways, highways, etc.); this is revealed by a recent state-of-the-art paper published in (Li, Pereira, and Ben-Akiva, 2018a) which emphasises the difficulty of solving this problem for arterial roads and the lack of studies in this field. Our study proposes a framework capable of predicting the incident duration regardless of the road network or its complexity.

**Secondly**, the majority of studies in the literature have concentrated on applying state of the art machine learning models mostly for classifying the incident severity (Nguyen, Cai, and Chen, 2017) or their duration (Li, Pereira, and Ben-Akiva, 2018a). However, very few have treated the problem of outliers or imbalanced data classes. Our study addresses both of these issues by proposing a varying threshold procedure that can facilitate binary duration classification threshold selection by considering both class balance and model performance. We also test multi-class classification on data sets split into three equally-sized parts according to incident duration: short, medium or long term. Previous research studies were selecting incident duration thresholds by simple reasoning (e.g. choosing mean, median, percentiles, etc) (Kuang et al., 2019a)-(Zou et al., 2018b)-(Li, Pereira, and Ben-Akiva, 2015b)-(li, 2014). We, on the contrary, test multiple different thresholds for three different data sets. Furthermore, we propose our own optimisation approach which we denote intra-extra joint optimisation (IEO) together with an outlier removal procedure (ORM) and advanced machine learning modelling.

**Thirdly**, we further solve the incident duration regression problem and also perform different regression scenarios to test the extrapolation performance of ML models on various incident data sets. We utilise thresholds selected during the classification threshold evaluation procedure to analyse the extrapolation performance by training ML models and making predictions on several duration subsets. It allows us to find the best ML model and the best extrapolation approach for the regression problem on each duration subset (e.g. short-term incidents) of each data set. For the regression problem, we also detect the most influential factors that affect the incident duration that traffic centres need to prioritise in order to predict incident duration with higher accuracy. Our end goal is to improve the extrapolation ability of machine learning models on the task of incident duration prediction and find the best modelling approach for short-term and long-term incidents.

**Lastly**, the majority of studies are primarily focusing on choosing a single winning algorithm or approach that works for a specific case study. Unfortunately, we show that the performance of ML models is highly affected by the data set and the chosen methodology: data quality, the available features, and the additional parameter tuning and optimisation techniques applied in this work. We try to develop the universal framework for traffic incident duration prediction applicable to different traffic incident data sets. We choose and adapt the best modelling approaches to each data set and show how this can affect the accuracy and performance of the models. This method allows high flexibility that can be applied for classification and regression predictions on various network types and different data sets.

The most similar research to the current work was published in Kuang et al., 2019a and relied only on one data set, one method for classification (Bayesian network), one method for regression (K-nearest neighbours), and authors selected static threshold (30 min) to alleviate the class-imbalance problem. This current paper provides a significant contribution by advancing on multiple aspects from a large pallet of machine learning models to multiple data sets with unique features, up to outlier removal and joint optimisation.

**Overall, our main paper contributions are the following:**

- to the best of our knowledge, this is the first research study proposing a bi-level prediction framework using a large pallet of several machine learning models applied for both incident duration classification and regression, with the scope of predicting the incident duration on different road

types across two different countries (Australia and U.S.A.). Overall, our methodology is agnostic of the location, the network, or the size of the network and can be adapted to any new incident log data set that can be made available.

- we propose a binary versus multi-class classification approach in order to find the best optimal threshold to identify short versus long-term incidents via both quantile analysis and varying threshold data split.
- we propose a novel intra/extra joint optimisation algorithm that integrates baseline ML models with outlier removal and hyper-parameter optimisation techniques across the validation cycle.
- we propose several extrapolation scenarios of analysing the impact of missing logs in the precision of the prediction model and reflect on what type of logs should be best used for tailoring to the prediction problem needs.
- we conduct a feature importance selection using the SHAP method, which allows graphical interpretation of variables impact on the model output, before we conclude on the most important factors affecting traffic incidents.

Overall, this research lays the foundation stone of bi-level predictive methodologies regarding the traffic incident duration and can provide accurate information for both the end-user route choice modelling as well as for the operational centres which need to optimise their operations under non-recurrent traffic congestion. Moreover, this work contributes to our ongoing objective to build a real-time platform for predicting traffic congestion and to evaluate the incident impact during peak hours (see our previous works published in (Mihaita et al., 2019b)-(Shafiei et al., 2020)-(Mao et al., 2021)).

The paper is organised as follows: Section 1 discusses related works, Section 2 presents the data sources available for this study, Section 3 showcases the methodology, Section 4 presents the numerical results for binary and multi-class classification tasks, Section 5 presents the numerical results of the regression part of the framework, Section 6 details on the feature importance evaluation and Section 7 is reserved for conclusions and future perspectives.

### 3.1.3 Related works

**Incident data interpretation:** The definition of traffic incident duration phases is provided in the Highway Capacity Manual Alkaabi, Dissanayake, and Bird, 2011, and it consists of the following time-intervals: 1) **incident detection time** which is the time interval between the incident occurrence and its reporting, 2) **incident response time** standing for the time interval between the incident reporting and the arrival of the first investigator at the location of the accident, 3) **incident clearance time** representing the time interval between the arrival of the first investigator and the clearance of the incident, 4) **incident recovery time** indicating the time interval between the clearance of the incident and the return of traffic flows to normal conditions.

The **total incident duration** is the time interval between the first incident log, and the returning of traffic flows to normal conditions. In our work, we use the term **incident duration** for the time lapse between the detection of an incident and the clearance of the incident, as officially reported in traffic logs provided by local traffic authorities. Therefore we do not include the incident recovery time as

this information is not recorded in the three data sets provided. However, different phases of traffic incident duration (e.g. clearance, recovery time) can be modelled individually upon availability; this type of research is rare because of the complexity of data collection for traffic incidents and small amounts of recorded traffic incidents in real-life datasets Li, Pereira, and Ben-Akiva, 2018a; Alkaabi, Dissanayake, and Bird, 2011.

When it comes to the data interpretation in the literature, the incident duration distribution has been modelled as log-normal Sullivan, 1997 and more recently as log-logistics distribution Chung, Walubita, and Choi, 2010; Smith and Smith, 2001. In a recent study, Haule et al., 2019a, incident clearance time and the total impact duration were modelled using Weibull, log-normal, log-logistic distributions and compared using the Akaike information criterion (AIC) criteria; findings have revealed that log-logistic distribution was outperforming other distributions. As distribution utilisation is highly related to the specificity of each data set, for this study, in which we use three different data sets, we further apply a comparison among several distribution modelling choices by using the AIC criteria.

According to Wali, Khattak, and Liu, n.d., different statistical methods were applied to model traffic incidents: 1) fixed parameter regression 2) random parameter regression 3) quantile regression. In this study, log-transformation of the target variable (incident duration) also applied. Random parameter regression found to give better statistical fit for incident duration models than fixed parameter regression, and therefore provide more accurate predictions of incident durations. It also highlighted that fixed-parameter regression model may give non-accurate incident duration predictions due to over/under- estimation of dependency between variables and incident durations. Also, there were no substantial difference found between fixed parameter regression and quantile regression in the case of 2015 Virginia incident data set. The benefit of quantiled regression is the ability to model the relationship of any quantile (rather than only average incident duration) of the incident duration vector with a set of explanatory variables Khattak et al., 2016. Ordinary Least Squares model can provide the predicted mean of the incident duration. On the contrary, quantile regression provides estimates for every quantile, which represented as a conditional distribution of incident durations, without providing single value as the incident duration prediction. Quantile regression coefficients represent the change in the incident duration in a given quantile category in relation to independent variables. Similar to this approach, variable importance can be estimated within each traffic incident duration group.

**Machine Learning for incident duration prediction:** While several statistical modelling techniques have been applied previously, more recently, new approaches in machine learning (ML) modelling have emerged as a more advanced way of predicting the incident duration due to their capacity to easily account for new data sources, as well as for removing the linearity assumptions between features and the predicted class Hojati et al., 2014. Examples of such approaches are: artificial neural networks (ANNs) Lopes et al., 2013, genetic algorithms Lee and Wei, 2010b, support vector machines (SVMs) Valenti, Lelli, and Cucina, 2010b, k-Nearest-Neighbours (kNNs) Wen et al., 2013 and decision-trees (DTs) He et al., 2013. The recently proposed Gradient-Boosted Decision Trees (GBDTs) have been shown to provide superior prediction performance when compared to Random Forests, SVMs and ANNs Ma et al., 2017. However, it is known that GBDT can easily over-fit when



the prediction target has a long-tail distribution, as is the case of the traffic incident duration distribution Ma et al., 2017. XGBoost Chen and Guestrin, 2016 is another decision-tree enhancement method that has gained popularity recently in the machine learning community due to its tree boosting capability, loss function regularisation and adaptive learning rate. It was employed in several international competitions, winning 17 out of the 29 Kaggle competitions singled out on the 2015 Kaggle blog; it was also employed by every team in the top-10 in the 2015 KDDCup Bekkerman, 2015 for solving various problems such as store sales prediction, web text classification, hazard risk prediction, and product categorisation. XGBoost's popularity is also due to its scalability (it can run on a single machine, as well as on distributed and paralleled clusters), its capacity to handle sparse data and its ability to handle instance weights in approximate tree learning (see the recent paper published by Chen and Guestrin, 2016 where authors proposed an end-to-end tree boosting system with cache-aware and sparsity learning features). While each of these methods has its advantages and disadvantages, building a fast and reliable prediction framework that could be applied for real-time operations represents a true challenge.

One of the recent research studies Kuang et al., 2019a presented a two-step approach for traffic incident duration prediction. A cost-sensitive Bayesian network was used to perform binary classification of traffic incidents by choosing a threshold of 30 minutes and then performing regression for each class using kNN. While the approach is functional, one major drawback for the classification problem is to manually choose the class split threshold, as it can lead to severe class imbalance; to overcome this issue, in our study, we perform both a fixed and a varying threshold set-up to find the best class balance for our classification models; even-more, we propose as well a comparison with a multi-class classification approach and debate on the benefits and drawbacks of using classifiers for such problems; we also enhanced more advanced regression models together with outlier removal procedures that would provide a better and more precise prediction of the incident duration precondition in minutes. Overall, the cost sensitivity of incorrect classification can be further extended to the cost-based regression metrics. We propose our enhanced ML models with a proposed intra and extra joint optimisation technique and outlier removal procedure to have even more precise predictions.

In one of the recent research studies on applying machine learning, which was related to the classification of driving state, multiple hyper-optimised ML models were tested, and entire feature space was visualised using t-SNE for entire feature space visualisation (Yi et al., 2019). RandomForest provided the highest prediction accuracy, but more advanced tree-based models exist that utilise gradient boosting, which we will be using in our research (e.g. gradient boosted decision trees).

To verify the performance of advanced tree-based methods (as LGBM - Light Gradient Boosted Model), additional conventional ML models can be used (Chen et al., 2020). We decided to also include LGBM and compare it to conventional ML models with non-tree based models (KNN and Logistic Regression).

**On the feature selection:** It is generally not enough to use all the possible features for the regression analysis of traffic incident durations. Using a high amount of features combined with a small data set size can lead to over-fitting. Some features can be helpful or useless, more or less critical, while others do not impact the prediction results significantly. By performing a feature importance analysis, we can recommend to traffic management facilities to record the most critical data and omit redundant data



related to traffic incidents. For example, one can increase the prediction accuracy by using as additional features the weather conditions, which were found to play a significant role in some research studies (e.g. during the summer and autumn seasons in Washington – USA in 2009, the preparation time of the rescue team was higher on freeways Hou et al., 2013). In some countries with cold weather, the response times can be much higher, while in regions with sunny weather most of the year, the weather impact on the intervention team can be neglected. Overall, the weather impact on the traffic incident duration prediction needs customised via a data-driven feature important analysis. Peak hours were the most influencing feature on response team preparation delay, which was found to be linked to response procedures (the goal of the response team was to resolve incidents during peak hours as soon as possible). A research study using Beijing traffic incidents data from 2008 Li and Shang, 2014a found the importance of "peak hour" value for the response team travel time and clearance time, but not for the intervention team preparation time. Our study conducts a feature importance ranking based on the best performing ML models we have proposed and provides a detailed overview of their impact. Different approaches to feature importance estimation use tree-based models (e.g. Random Forest, Light Gradient Boosted Machines - LGBM, extreme gradient bosoting models - XGBoost). For example, one can use produced decision trees from the tree-ensemble model Chen et al., 2020. A data-driven approach was used to perform information fusion from different sources Abou El Assad, Mousannif, and Al Moatassime, 2020, which involved the use of Gini-index extracted from Random Forests as a method to estimate feature importance. Nevertheless, the single random model can have a noticeable variance in data mapping when there is a weak connection between features and the target variable by making the feature importance value dependent on the random seed for the model. The Shapley Additive explanation (SHAP) Lundberg and Lee, 2017 provides a more advanced approach for feature importance estimation because it fuses estimation from multiple models trained across many different subsets (which selected both feature-scale and index-scale) of the dataset. These studies motivated the utilisation of the Shap Values for our feature importance ranking across three different data sets, all with different features and incident information.

**On the future application of our research:** In comparison with other work, the research proposed in our paper comes not only with a significant prediction capability for all types of incident data sets with various features, but it can be further extended for solving the route scheduling problem within traffic simulation modelling, which will incorporate the adaptation of agents to occurring traffic incidents. Apart from analysing the effects of traffic control measures Knapen et al., 2014, it is possible to analyse the effect of additional information such as the predicted incident duration, which can be performed both for scheduling and online rescheduling of dynamic agent re-routing. Furthermore, simulation can be performed with and without such information to estimate the possible benefits of the incident duration prediction modelling within the traffic system. Also, using an online rescheduling procedure requires the simulation to be performed at the level of dynamic agents within a micro-simulation model, which could benefit from new re-routing schemes when traffic disruptions occur along the route.

In order to test the efficiency of the proposed bi-level framework, we have used three different data sets from two different countries: Australia and U.S.A. The three data sets represent incident logs from an arterial road suburb in Sydney, a motorway in Sydney, Australia, and a road area from San

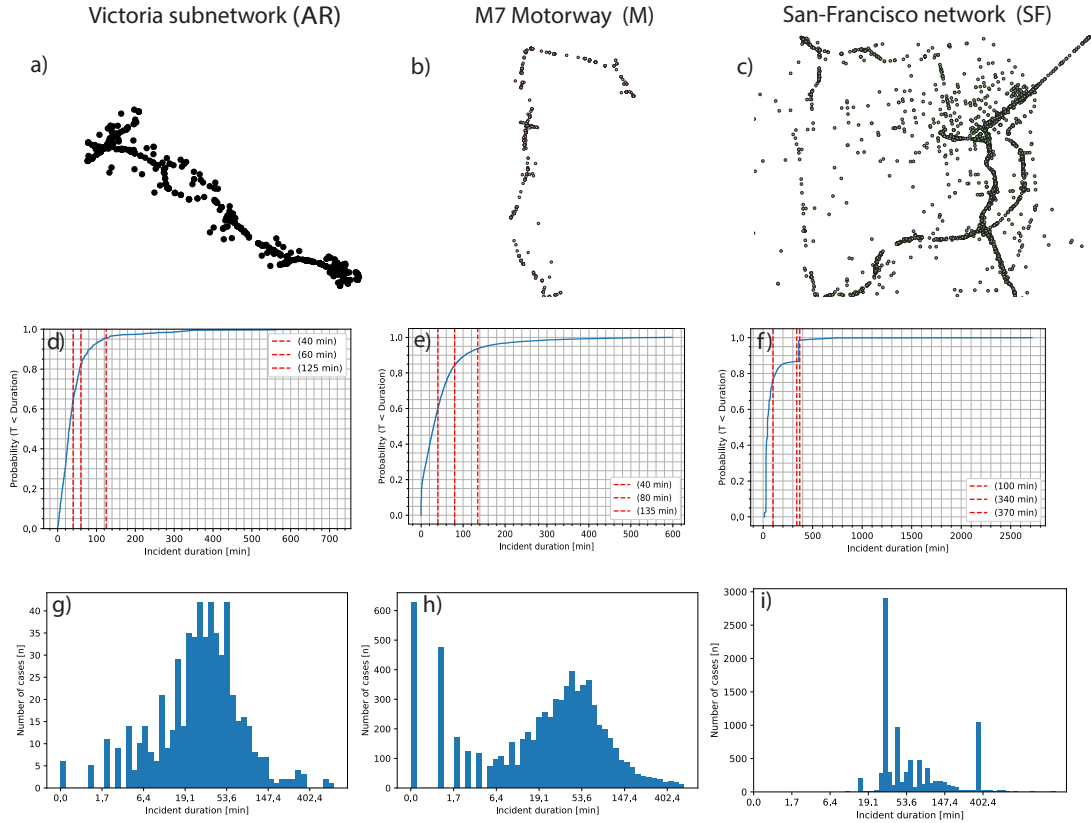


FIGURE 3.1: Data profiling for all data sets in our study: Victoria Rd (A) - a) network mapping, d) ecdf - empirical cumulative distribution function g) distribution plot; M7 motorway (M) - b) network mapping, e) ecdf h) distribution plot; San Francisco (SF) - c) network mapping, f) ecdf i) distribution plot.

Francisco, U.S.A. The data sets are all recorded by different means and allow us to explore the impact of the prediction framework across various types of road networks. The three data sets are represented in Fig. 3.1 and are detailed as follows.

**Victoria Rd - arterial network, Sydney:** The first data set (dataset AR) contains one-year incident logs from the Victoria arterial road from Sydney, Australia (in 2017) (see Table 3.1 for a summary of features, in which the + symbol under each data set column and for each line indicates whether that variable is present or not in the data set - for example, the TZName variable is present in the Arterial Roads data set but not in the M motorway data set). It contains information on 5,134 traffic incidents with different incident types (e.g. hazards, breakdowns, accidents) and subtypes (e.g. work zone, accident with truck). Our current study focuses on 574 “Accidents” since these induce the longest clearance time in the current subnetwork according to the traffic management centre (TMC). Traffic ‘Accidents’ have a mean duration of 44.59 minutes and a maximum of 719 minutes. Weather data represented as average daily temperature (in Celsius) and precipitation rate (in millimetres) are obtained from the Observatory Hill station in Northern Sydney, which is the closest station to the analysis area. Public holiday data represented as boolean values for public and regional holidays in 2017 in New South Wales, Australia. The area geometry features contain the sector ID as defined by TMC, the code of the official area where the accident occurred (as defined by the Bureau of Transport

| Variable            | AR | M | Values  | Description  |
|---------------------|----|---|---|--|
| Location            | +  | + | $\mathbb{N}, \mathbb{N}$                      | $X, Y$ in GDA Lambert coordinates  |
| Hour of day         | +  | + | $\{0, 1, \dots, 23\}$                         | -  |
| Peak Hour           | +  | + | $\{1, 0\}$                                    | Value is 1 if hour belongs to $\{7 \dots 9\}$ or $\{16 \dots 18\}$ hour interval |
| Day of the week     | +  | + | $\{1 \dots 5\}$                               | Weekday numbers from Monday to Friday  |
| Weekend             | +  | + | $\{0, 1\}$                                    | Value is 1 for Saturday and 0 for Sunday   |
| Month of the Year   | +  | + | $\{1, 2, \dots, 12\}$                         | -  |
| Incident Subtype    | +  | + | $\{Bus, car, bicycle, animals, etc.\}$        | Field indicating cause of incident   |
| Affected lanes      | +  | + | $\{1, 2, 3, 4, Alllanes, breakdown, nodata\}$ | Number of lanes affected by the accident   |
| Direction           | +  | + | E, W, N, S, E-W, N-S, One/Both                | Affected traffic direction   |
| Incident Source     | +  | + | $\{ICEMS/ISENTRY, OPERATOR, etc\}$            | Source of the incident report  |
| Unplanned           | +  | + | $\{0, 1\}$                                    | Value is 1 if incident is planned, 0 otherwise                                   |
| Average Temperature | +  | + | $\{11.13C - 32.4C\}$                          | Average temperature for the time of the incident                                 |
| Rainfall            | +  | + | $\{0 - 85mm\}$                                | Rainfall for the time of the incident  |
| Public holidays     | +  | + | $\{0, 1\}$                                    | Value is 1 if days is a public holiday   |
| Sector ID           | +  | + | R+  | Defined by TMC   |
| TZName              | +  | + | R+  | Traffic zone name as Defined by the Bureau of Transport Statistics               |
| Section ID          | +  | + | R+  | Road section on which the incident occurred                                      |
| Section Speed       | +  | + | $R + [Km/h]$                                  | Section speed limit  |
| Section Lanes       | +  | + | $\{1, 2, 3, 4, 5, 6\}$                        | Number of section lanes  |
| Section class       | +  | + | R+  | As defined by TMC  |
| Street ID           | +  | + | R+  | As defined by TMC  |
| Intersection ID     | +  | + | R+  | As defined by TMC  |
| Distance from CBD   | +  | + | R+  | distance between the traffic incident and the city CBD                           |
| Section Capacity    | +  | + | $\{0 \dots 3100 vehicles/hour\}$              | Maximum flow capacity of the section   |

TABLE 3.1: Traffic incident features for Sydney Arterial roads (AR) and M7 motorway (M).

and Statistics), and supplementary information such as section capacity, section speed limit, and the number of lanes. These features are available for all road sections in the Victoria sub-network, and they were extracted from the official traffic simulation model of the Victoria network, developed in Aimsun and previously used by the authors for conducting an incident impact analysis and traffic prediction (Wen et al., 2018).

**M7 motorway, Sydney:** The second data set is a motorway data set (data set M), consisting of 7,194 traffic accidents along the M7 motorway in Sydney, Australia, during the same year 2017. The mean duration of motorway accidents is 47.2 minutes, with a maximum duration of 598 minutes (9.96 hours). This data set also includes weather data (average daily temperature and precipitation). This set of features is similar to the arterial roads data set AR without the geometric features of the lanes (section lanes, section class), intersection ID, distance from the central business district (CBD); this is due to the complexity of mapping of a traffic incident to a correct location along the motorway. We make the observation that for both Data set AR and M, the traffic flow information of the affected road sections was omitted for this study since we found previously no significant improvement to the prediction accuracy (Mihaita et al., 2019b).

**San-Francisco road network:** The last data set is from San-Francisco, U.S.A. (data set SF) and includes information on accidents from all types of roads in the city. It is part of a more considerable initiative entitled "A Countrywide Traffic Accident Dataset", recently released in 2021, which contains 1.5 million accident reports collected for almost 4.5 years since March 2016 (Moosavi et al., 2019b). The SF data set contains 49 features describing the accidents as detailed in (Moosavi et al., 2019b) (due to a large table of feature, we refer the reader to the cited paper and not duplicate this feature information). This study focuses on the "accident" type duration prediction as being the most severe one. We extract and use 8,754 accident records related to the San-Francisco area. As observed from

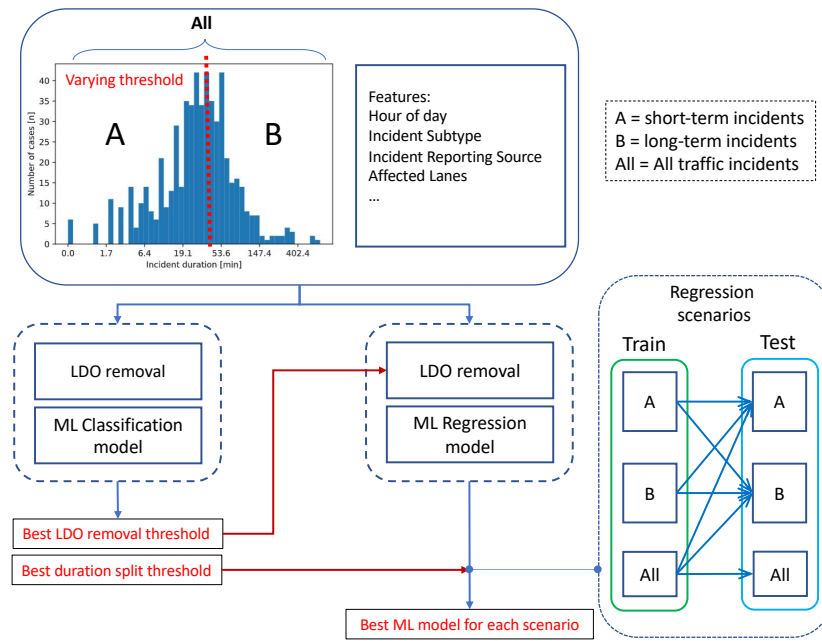


FIGURE 3.2: The proposed bi-level modelling framework for traffic incident duration prediction.

Fig. 3.1 c), a significant part of the accidents occurred along the “US-101” highway and “John F. Foran” Freeway. Accidents have a mean duration of 100 minutes and a max duration of 2,715 minutes.

**Data sets profiling:** Each data set undergoes a profiling procedure by investigating the empirical cumulative distribution functions (ECDF) - as plotted in Fig. 3.1 d), e), f), and their equivalent log-space distribution plots (as represented in Fig. 3.1 g), h), i). The ECDF function presents thresholds of data behaviour (marked in red) across each data set which reveal indicative thresholds of a different behaviour around specific incident duration (see for example Fig. 3.1d) versus Fig. 3.1f) where the first inflection point is around 40min for data set AR versus 100min for data set SF. Findings reveal significant anomalies representative of each data set. For example, data set AR contains a reduced amount of traffic accidents with small incident duration (zero or less than 4 min), data set M contains an increased number of accidents with zero or one-minute duration, while the data set SF despite not presenting any short term incident duration below 17 minutes, it contains a large number of incidents of 29 and 360 minutes which raises the question of either these are outliers in the data set or simply reveal a road network behaviour in terms of incident management in the area; this also might indicate that it will present unique behaviour under the prediction framework and that different processing techniques needs to be applied for this data set. We also observe that the incident duration is long-tail distributed, which is likely to pose difficulties for prediction algorithms due to the presence of extreme values (either small or large).

## 3.2 METHODOLOGY

Clearing accidents in a short time represents a high priority task for traffic management centres (TMC) worldwide. For example, in New South Wales, Australia, the target clearance time for traffic incidents is 45 minutes, but this limit might differ in other countries. Therefore, in the rest of this paper, we will refer to this threshold as “incident clearance threshold ( $T_c$ )” and any incidents cleared before this threshold (e.g.  $< 45$  min) as “short-term”; incidents which lasted more than the clearance threshold (e.g.  $\geq 45$  min) will be referred to as “long-term” traffic incidents. A unique threshold will be derived for each dataset and will be discussed further in this paper. The methodology of this paper has its origins in our previous work applied only for arterial roads (Mihaita et al., 2019b), which we further extend and improve via the joint optimisation and outlier detection enhancements of the prediction framework. The methodology we propose for modelling the incident duration prediction problem is using a bi-level prediction framework combining a classification and regression modelling, as represented in Fig. 3.2. This approach has been constructed by considering the real-time operational goals of TMC and providing short duration prediction into the life-cycle of the incident management.

Based on the initial traffic incident information, the first step is the deployment of a fast classification method which would only predict whether the accident will be either short-term (subset A) or long-term (subset B) - see incoming data set from Fig. 3.2 where the data is split in two parts based on  $T_c$ . Next, we test various duration thresholds and select the optimal  $T_c^o$ , which provides a good class balance and classification performance for each dataset. Once the optimal  $T_c^o$  has been found, a further regression modelling is applied for predicting a more precise duration of future incidents down to the minute level.

Due to the main challenge of this task, we further propose an outlier removal approach (ORM) detailed in Section 3.2.7 and our innovative Intra/Extra Joint Optimisation modelling coupled with several machine learning models trained via a hyper parameter tuning (we denote this approach as IEO-ML and is further detailed in Section 3.2.9).

The boosted regression framework is finally applied under several regression scenarios (see section Section 3.2.6), which are constructed to evaluate the framework capability to predict under all possible situations. For example, when we only have a subset A available (short-term incidents) but the TMC would like to predict long term incident (subset B) we denote this as a Scenario A-to-B (training the models on subset A and making predictions on subset B); all scenarios are constructed based on the assumptions that the framework needs to be robust in order to predict any type of incident durations, under all possible data shortage or lack of information availability. In the following subsection, we further provide the mathematical and theoretical modelling of each of the steps described above.

### 3.2.1 Classification and regression definitions

Using all available data sets and the incident information, we first denote the matrix of traffic incident features as:

$$X = [x_{ij}]_{i=1..N_i}^{j=1..N_f} \quad (3.1)$$

where  $N_i$  is the total number of traffic incident records used in our modelling and  $N_f$  is the total number of features characterising the incident (severity, number of lanes, type, neighbourhood, etc.)

according to each specific data set (see examples provided in Table 3.1). For the incident duration classification problem, we denote the incident duration classification vector as:

$$\begin{cases} Y_c = [y_i^c]_{i \in 1..N} & y_i^c \in \{0, 1\} \\ Y_{mc} = [y_i^{mc}]_{i \in 1..N} & y_i^{mc} \in \{0, 1, 2\} \end{cases} \quad (3.2)$$

where  $N$  is the duration of the traffic incident (in minutes),  $Y_c$  is the vector of binary values taking values in  $\{0, 1\}$ , and  $Y_{mc}$  is the vector of integer values for the multi-class classification problem definition, taking values in  $\{0, 1, 2\}$ . More specifically, in the first stage we create a binary classification modelling with the purpose of identifying short versus long-term incident duration, split by the incident clearance threshold  $T_c$ . Thus our task is to predict  $y_i^c$ , where  $Y_c$  takes one of the binary values:

$$\begin{cases} y_i^c = 0 & \text{if } y_i \leq T_c, \text{ short-term incidents} \\ y_i^c = 1 & \text{if } y_i > T_c, \text{ long-term incidents} \end{cases} \quad (3.3)$$

where the threshold is varied every 5min between  $T_c \in \{20, 25, \dots, 70\}$ . Subsequently, the multi-class method identifies the best two thresholds to separate between short, mid and long-term incident duration. The main purpose of this approach is to test the limits of the class balance which would maintain good model performance, and is expressed as follows:

$$\begin{cases} y_i^{mc} = 0 & \text{if } y_i \leq T_c^1, \text{ short-term incidents} \\ y_i^{mc} = 1 & \text{if } y_i \in [T_c^1, T_c^2], \text{ mid-term incidents} \\ y_i^{mc} = 2 & \text{if } y_i \geq T_c^2, \text{ long-term incidents} \end{cases} \quad (3.4)$$

where  $T_c^1$  and  $T_c^2$  take several values as further detailed in Section 3.3.3. The binary classification approach implemented with a computation time constraint for operational purposes (more details on computation time comparison can be found in ??).

The regression problem is further structured with a more fine-grained incident duration prediction in mind. The main objective motivating the regression modelling consists in more precise information regarding the duration of incidents which can fall into a wide class which contains mostly incident logs with a reported duration between and 0 minutes and 30 minutes (for these cases, the traffic centres require more detailed precision to the minute level as a 5-min accident has different handling procedures than more severe accidents of 30min for example). The incident duration regression vector ( $Y_r$ ) is represented as:

$$Y_r = [y_i^r]_{i \in 1..N}, y_i^r \in \mathbb{N} \quad (3.5)$$

and the regression task is to predict the traffic incident duration  $y_i^r$  based on the traffic incident features  $x_{i,j}$ . The regression models go via an extensive cross-validation procedure with hyper-parameter tuning, with the test of outlier removal using a joint optimisation approach as further detailed in the Section 3.2.4-Section 3.2.7-Section 3.2.9.

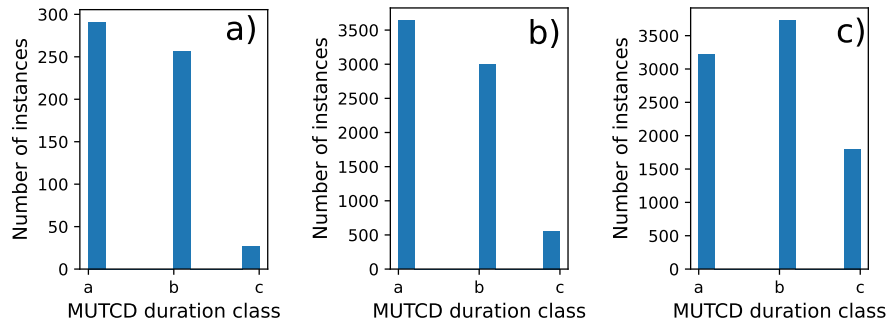


FIGURE 3.3: Distribution of incident durations according to MUTCD duration classes: a) Arterial roads, Sydney, Australia b) M7 Motorway, Sydney, Australia c) San Francisco, USA

### 3.2.2 Applicability of knowledge-based incident duration classification guidelines

According to the The Manual on Uniform Traffic Control Devices (MUTCD) official guidelines Transportation, 2017 Section 6I, traffic incidents divided into three classes: a) Major - with expected duration more than 2 hours b) Intermediate — expected duration of 30 minutes to 2 hours c) Minor — expected duration under 30 minutes.

First, the MUTCD classification seems to be general knowledge-based system and does not consider specifics of each data set / country regulations / specifics of applied incident response guidelines. We approach the classification task from data analysis point of view in relation to application of ML models and try to infer these thresholds from the actual data sets. Shifting to MUTCD classification approach will also make incident duration classes imbalanced (see Figure 3.3). Second, this classification may not be applicable due to road networks heterogeneity Li, 2018 and consequent differences in incident duration distribution. As can be seen from Figure 3.1, all three data sets have different distribution of incident durations and therefore such classification may be biased in each case. Overall, in our study, we aim to provide insights from data analysis point of view.

### 3.2.3 Selection of baseline machine learning models

We have tested and deployed several ML models for both the classification and regression problems for this current work, which have served as baseline models to compare our proposed optimisation approach. These are listed as follows: a) gradient boosting decision trees - GBDT (Friedman, 2000) which rely on training a sequence of models, where each model is added consequently to reduce the residuals of prior models; b) extreme gradient decision trees - XGBoost (Chen et al., 2015) which finds the split values by enumerating over all the possible splits on all the features (exhaustive search) and contains a regularisation parameter in the objective function; c) random forests - RF (Breiman, 2001) which applies a bootstrap aggregation (bagging, which consists of training models on randomly selected subsets of data) and uses the average (or majority of votes) of multiple decision trees in order to reduce the sensitivity of a single tree model to noise in the data; d) k-nearest neighbours - kNN (Fix and Hodges, 1951) which uses for the prediction on data points the majority of votes or the average from k closest neighbouring data points from the training set (based on a distance metric); e) linear Regressions - LR - a standard predictor using linear equations to model the relation



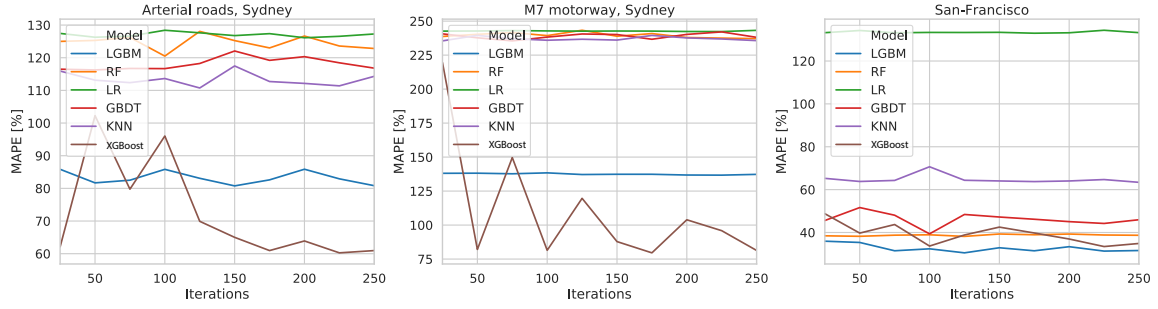


FIGURE 3.4: Performance testing of ML models across three different data sets

between the features and the regression variable; f) light gradient boosted machines - LGBM (Ke et al., 2017) which applies gradient boosting to tree-based models; it also uses a Gradient-based One-Side Sampling (GOSS) and excludes data points with small residuals for finding split value. The models have been used for both classification and regression problems (except logistic regression applied to classification only and linear regression to regression problem only). They are the main base on which we further enhance and develop our outlier and joint optimisation prediction algorithm used in the current bi-level incident duration prediction framework.

### 3.2.4 Hyper-parameter tuning through randomised search

Most machine learning algorithms have a set of hyper-parameters related to the internal design of the algorithm that cannot be fitted from the training data. Both GBDT and XGBoost present dozens of hyper-parameters, out of which the most important ones are max\_depth, learning\_rate, min\_child\_weight, gamma, subsample, colsample\_bytree and scale\_pos\_weight [24]. The hyper-parameters are usually tuned through randomised search and cross-validation. The most extensive search technique is the grid-search, in which several equally spaced points are chosen in the most credible interval for each parameter, and for each point combination, a model is fitted and tested through cross-validation. The grid-search parameter tuning is straightforward; however, the grid-search scales poorly as the number of hyper-parameters increases. In this work, we employ a Randomised-Search (Bergstra and Bengio, 2012) which selects a (small) number of hyper-parameter configurations randomly to use through cross-validation.

To determine the optimal number of iterations for models and data sets, we perform iterative testing. The number of random-search iterations is from 25 to 250 with step 25. For example, on Fig. 3.4, (Arterial roads, Sydney), we see that XGBoost is the best performing model starting from 120 iterations, and it is already close to optimum starting from 175 iterations. The second-best performing model is LGBM, but increasing the number of iterations does not seem to have a significant effect on the model performance which seems to be quite stable without many fluctuations across all evaluation metrics. Other methods perform significantly worse (more than 110% MAPE). For San-Francisco, we see the superior performance of LGBM. The second best is XGBoost. Since there are no metric improvements across iterations for most models, the number of iterations is essential only for XGBoost. According to the results, we decide to search for hyper-parameters for 250 random parameter combinations for each model. We evaluate each combination using a 5-fold cross-validation and then providing results using a 10-fold cross-validation using best combination.



### 3.2.5 Model Performance Evaluation

The performance of classification model is evaluated using the Precision, Recall, Accuracy and F1-score and defined as:

$$Precision = \frac{tp}{tp + fp}, \quad (3.6)$$

$$Recall = \frac{tp}{tp + fn}, \quad (3.7)$$

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn}, \quad (3.8)$$

$$F_1 = 2 * \frac{precision * recall}{precision + recall}. \quad (3.9)$$

where  $tn$  represents true negatives,  $fn$  - false negatives,  $tp$  - true positives,  $fp$  - false positives. For example, We refer to true positives the incidents which have been predicted to be in a specific class (say short-term) and indeed they were short-term upon validation, false positives the incidents which were predicted to be short term but were not, etc.

We use F1-score as a target metric for classification experiments as F1 represents the balance between Precision and Recall, and is in general a better performance metric to use when we are facing an uneven class distribution rather than interpreting the Accuracy results which take into consideration the total number of both false positive, false negative together with the true positives and true negatives; therefore for uneven class balances (especially the ones with fewer incident records), one should rely less on Precision and Accuracy metrics. To evaluate the regression models we use the mean absolute percentage error defined as:

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - F_i}{A_i} \right| \quad (3.10)$$

where  $A_i$  are the actual values and  $F_i$  - the predicted values,  $n$  - number of samples. Other metrics have been calculated but we will keep them concise due to large amount of experiments to show.

### 3.2.6 Regression scenarios definition

The main objective of the bi-level framework is that the regression accuracy can benefit from different setups for different data subsets. For an even better accuracy compared to the classification problems, we are further developing more complex regression models that can provide incident duration prediction at minute-level accuracy. This is the second step of the bi-level prediction framework to be applied when more precision is needed at the minute level regarding the incident duration length. When training such regression models, a crucial step is the size of the data set and the distribution of the target variable (incident duration). Due to the long tail distribution of incident duration and the class imbalance problem previously identified, we need to design and construct various regression models capable of learning from various types of data sets to make accurate predictions. However, with limited information (small data set size), the prediction results can be skewed. This is the primary motivation that led to the construction of several scenarios of model training, validation and prediction that can be applied under both complete or incomplete data sets from traffic centres. By using the classification thresholds identified previously, we split the traffic incident data set into two subsets: subset

A (with duration below threshold  $T_c$ ) and subset B (with duration above threshold  $T_c$ ) as previously defined at the beginning of [Section 3.2](#). We further contract several scenarios of subset combinations for training-validation-testing detailed with the aim of extrapolating the model performance:

**Scenario All-All:** we use the entire data set and apply several regression models using a 10-fold cross-validation approach and different hyper-parameter search methods. This approach will show us the general performance across various methods.

**Scenario A-to-B:** we use subset A (short-term incidents) for training the regression models and evaluate the prediction on subset B (long-term incidents). In this scenario, we will analyse methods to extrapolate to higher values of the target variable.

**Scenario A-to-A:** we use subset A for training the regression models and predict on subset A. In this scenario, we will analyse the prediction ability of methods with long-term incidents excluded (which includes values from the tail of the incident duration distribution).

**Scenario B-to-A:** we use subset B for training the regression models and predict on subset A. In this scenario, we will analyse methods to extrapolate to lower values of the target variable.

**Scenario B-to-B:** we use subset B for training the regression models and predict on subset B. In this scenario, we will analyse the prediction ability on long-term incidents.

**Scenario All-to-A:** we use all the data for training the regression models and predict on each fold within subset A. In this scenario, we will analyse the effect of having access to all types of incident logs in the training phases, both long-term and short-term incidents and how their presence might affect or not, the prediction of short-term incidents duration only. This is to evaluate if using all types of records, including rare events, will help or not to predict better short incidents.

**Scenario All-to-B:** we use all the data for training the regression models and predict on each fold within subset B. In this scenario, we will analyse the effect of having access to all types of incident logs in the training phases, both long-term and short-term incidents and how their presence might affect or not, the prediction of long-term incidents duration only. This is to evaluate if using all types of records, including short-term events, will help or not to predict better long incidents.

Indeed, from an operational perspective the scenario All-to-All is the ideal situation when traffic management centres would have in their data base both long term and short-term incidents. However, from an operational perspective, several records of short incidents for example and not being kept all the time, while long incidents are often being transferred to various other division if they last more than one day, and they become more of a road infrastructure problem rather than an operational problem which requires constant intervention. Therefore, various incident logs can be imbalanced – some containing more short-term incidents, and others more long-term incidents. The main idea is to provide a good deep dive into the effects of data availability on the model training. For example, training any model only on short term incidents as these are the only ones available will most likely not provide good prediction results in case of long-term incidents and vice versa.

### 3.2.7 Outlier removal methods (ORM)

As previously discussed in [Fig. 3.2](#) during the data profiling, we observed that the traffic incident logs contain outliers appearing as either minor incidents, rare traffic incidents with highly long duration and/or as errors in incident reports. Therefore, to reduce the side-effect of outliers on all models, we deploy two commonly used outlier removal methods. The IsolationForest (IF) (Liu, Ting, and Zhou,

2008) is an outlier removal method, which uses forests of randomly split trees. For each tree, the method randomly selects a feature and a random feature value. The data set is divided into two parts in each step until each data point becomes “isolated” (split from the rest of the data). If the data point is an outlier, it will have a small tree depth (e.g. data point gets quickly separated from the rest by selecting values in just a few features). Tree depth is then averaged between all the “isolation” trees and considered an anomaly score (e.g. if the average tree depth for a point is 1.3, the point is easily separable after a small number of splits). LocalOutlierFactor (LOF) (Breunig et al., 2000) is another outlier removal method, which estimates the anomaly score from local deviation of density within k-nearest neighbourhood. LOF relies on the calculation of a local reachability density (LRD), which represents the inverse of the average reachability distance (RD) of neighbouring data points from the selected data point. Reachability distance (RD) represents the distance to the most distant neighbour within a k-sized neighbourhood (k is also hyper-parameter). LOF of data point then represents the relation between LRDs of neighbours and its LRD and can take values: a) above 1 (higher LRD than its neighbours), b) below 1 (lower LRD than neighbours) and c) equal to 1 (data point has the same density as neighbours). According to the LOF score, we can sort data points and select specific per cent of data points, which have the highest LOF to be eliminated. LOF method relies on the fact that outliers belong to the area where the density of data points is low, while regular data points belong to the high-density area. To summarise, the above outlier removal procedures are applied in conjunction with the proposed optimisation framework and regression models and show a significant improvement in prediction accuracy as further detailed in Section 3.4.3.

### 3.2.8 Outliers from ORM point of view

We would like to make the observation that all the incidents have scalar degree of anomaly when applying outlier removal method. herefore, there are no discrete categories of outliers and normal data points from an outlier method point of view. We simply remove a per cent of data points (e.g. 2%) with the highest degree of anomaly. These points are either easily separable using IF method (tree-based) or remain on a low local density area for the LOF method (distance-based).

So does our outlier removal method actually remove long-term incidents failing to distinguish them from outliers? ML methods in our case, find outliers not only by the value of duration but by including all reported variables (e.g. 25 in the case of Arterial roads). Our aim in this work is to remove incident reports which have very rare characteristics overall, which are also known to negatively affect the ML method performance Lu, 2021b.

Fig. 3.5 Showcases Data sets with 10% of points with the highest anomaly score removed using IsolationForest: a) Arterial roads, Sydney, Australia b) M7 Motorway, Sydney, Australia c) San-Francisco, USA. By performing experiments with an outlier removal (isolation forest, 10% of point with the highest anomaly rate removed), we see how many incidents were removed according to each duration interval. An important finding is that outliers do not reside in the area of long-term incidents but rather scattered among the general population of incidents.

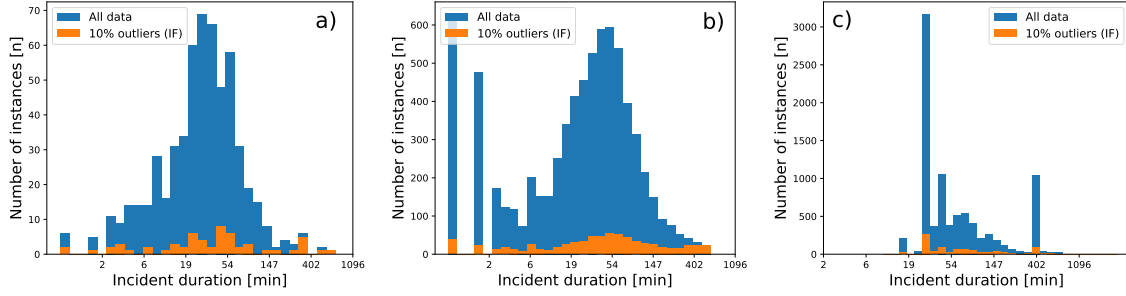


FIGURE 3.5: Data sets with 10% of points with the highest anomaly score removed using IsolationForest: a) Arterial roads, Sydney, Australia b) M7 Motorway, Sydney, Australia c) San-Francisco, USA

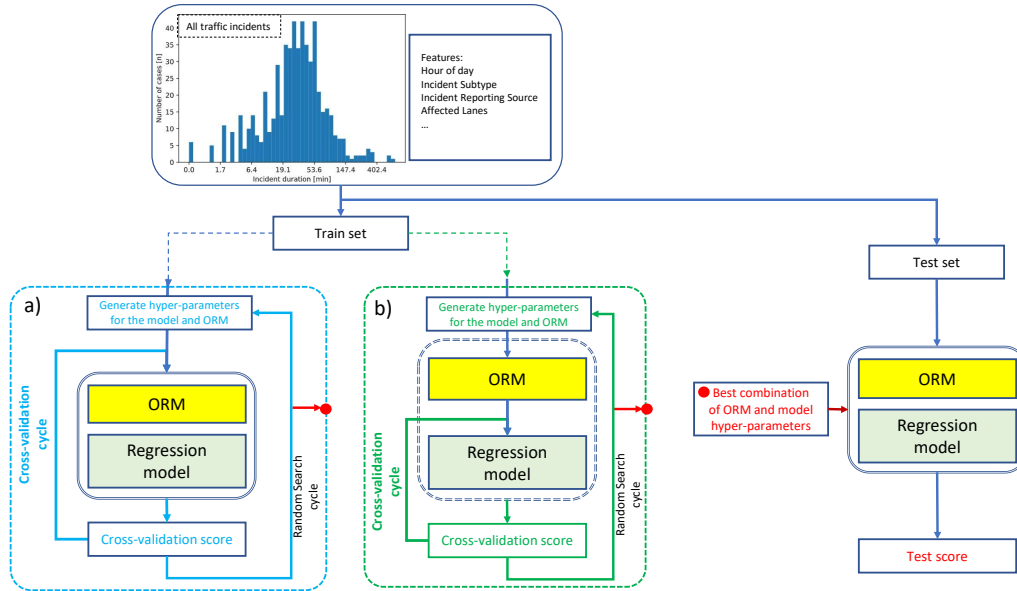


FIGURE 3.6: IEO-ML algorithm with a) Intra joint optimisation schema for the EO-ML algorithm, b) Extra joint optimisation schema for the IO-ML algorithm. Red dot on schema blocks represents output in the form of the best combination of ORM and model hyper-parameters

### 3.2.9 Intra/Extra Joint Optimisation for ML regression prediction (IEO-ML)

This section presents our novel enhancements of ML regression models by constructing an **intra/extra optimisation technique** to jointly optimise the hyper-parameters of the regression models together with previous outlier optimisation methods. In the rest of the paper, we denote this approach as **IEO-ML**, where ML is one of the regression models previously described (GBDT, XGBoost, RF, kNN, LR, LGBM). We introduce this approach for multiple reasons: 1) the traffic incident data is prone to errors during the data collection, which is attributed to human factors (e.g. presence of incidents with 0 and 1-minute durations, for example), 2) an outlier removal performance cannot be assessed on the new dataset with no marking for outliers; thus, we can assess outlier removal performance by looking at model performance with outlier removal applied, use joint outlier removal and modelling to assess the outlier removal performance metrics, 3) both the outlier removal method and models have hyper-parameters forming a single hyper-parameters space, 4) we assume that the outlier removal can be performed either inside (Intra - see Fig. 3.6a)) or outside (Extra - see Fig. 3.6b)) of the cross-validation cycle, and we evaluate the effect of such an approach on the model performance, 5) Intra joint optimisation can provide a more effective outlier removal since common hyper-parameters will be found for different data subsets, which allows ORM to be adapted to different possible combinations of incidents in case of the model deployment and prediction on the newly acquired incident log. Overall we want to compare and observe the impact of each technique on the accuracy of regression models and detect the best combination of Intra/Extra joint optimisation and various ML regression models.

Further, we present our proposed IEO-ML algorithm in conjunction with the two outlier removal methods IF and LOF, and several regressions models. Our approach explores the following combinations of ML models in selected working base (decimal or logarithm) with outlier removal and intra/extra joint optimisation; for example, we denote as *iLOF-LT-MLmodel* a “joint optimisation of any available baseline ML model with LOF in a log-transform base within a cross-validation cycle (an intra optimisation)”. As an observation, ORM has specific hyper-parameters but one parameter in common - the percentage of removed samples, which we assume to be outliers (ORperc). Thus, to solve the ORM problem, we assume that the amount of outliers in each data set (ORperc) can take values up to 5%. EJO is performed only once and before the cross-validation cycle, but IJO is performed within each fold in a number of times which is equal to the number of folds. Thus, ORperc has values in  $\{0, 1 \dots 5\%$  for EJO, in  $\{0, 1/5, \dots, 5/5\}$  for IJO to ensure a comparable amount of removed samples from both approaches. Results for all combinations of the proposed approach inside the incident duration prediction framework are further provided in Section 3.4.3 for eLOF-ML models, iLOF-ML, iIF-ML, eIF-ML (e.g. eIF-ML is a “joint ML optimisation using IF optimised outside (e) of the cross-validation cycle”).



**Data:** Traffic incident reports (feature vector  $X$ , duration vector  $Y_r$ )

**Input:** HPSm (Hyper-Parameter Space for Model),

ORM: Outlier Removal Method,

HPSor: Hyper Parameter Space for ORM,

Model: ML regression model  $\in \{GBDT, XGBoost, RF, kNN, LR, LGBM\}$ ,

Iters: Number Of Iterations (number of random search steps for hyper-parameter optimisation),

Folds: number of folds for cross-validation,

sample: function for random sampling from the hyper-parameter space,

FoldIndexes: function to get sample indexes for training folds and test fold,

extra: boolean variable stating the use of extra joint optimisation,

intra: boolean variable stating the use of intra joint optimisation,

split: function to split data set into two parts - train/test and validation parts

**Output:** Predicted duration vector  $Y_r$

$x_{tr}, y_{tr}, x_{te}, y_{te} = split(x, y);$

$P = [];$  // temporary cross-validation prediction vector

$results = []$

**for**  $it \leftarrow 1..Iters$  **do**

$HYPm \leftarrow sample(HPSm)$

$HYPor \leftarrow sample(HYPor)$

$idx_{train} = [];$  // indexes of train samples

$idx_{valid} = [];$  // indexes of validation samples

$res = 0;$  // scoring results

**if**  $extra$  **then**

$x = ORM(x, HYPor);$  // if EO then filter the outliers from the feature vector

**for**  $k \leftarrow 1..Folds$  **do**

$idx_{train}, idx_{valid} = FoldIndexes(x, k);$

$x_{train} \leftarrow x_{tr}[idx_0^{train}], \dots, x_{tr}[idx_N^{train}];$  // array of feature vector samples for training

$y_{train} \leftarrow y_{tr}[idx_0^{train}], \dots, y_{tr}[idx_N^{train}];$  // array of duration vector samples for training

$x_{valid} \leftarrow x_{tr}[idx_0^{valid}], \dots, x_{tr}[idx_N^{valid}]$

$y_{valid} \leftarrow y_{tr}[idx_0^{valid}], \dots, y_{tr}[idx_N^{valid}]$

**if**  $intra$  **then**

$x_{filtered}^{train} = ORM(x_{train}, HYPor);$  // if IO then filter outliers

$initialize\_model(Model, HYPm);$  // random hyper-parameter initialisation

$m \leftarrow fit\_model(Model, x_{filtered}^{train}, y_{filtered}^{train});$  // fitting the model to the filtered train set

$y_{pred} \leftarrow predict(m, x_{valid});$  // performing predictions

$P = [P; y_{pred}]$

**end**

$res \leftarrow Metric(y_{tr}, P);$  // scoring the accuracy of predictions using performance metric

$r = [];$  // Initializing hash-array

$r['metric'] = res;$  // populating hash-array with resulting metric

$r['HYPm'] = HYPm$

The algorithm represents the modified cross-validation cycle within the randomised hyper-parameter tuning procedure. We use multiple iterations (in fact, attempts) to find optimal parameters both for the selected model (HYPM) and the outlier removal method (HYPor). On every iteration, we sample hyper-parameter sets from hyper-parameter spaces. Then, if extra joint optimisation selected, an outlier removal procedure performed using all the data before the fold-rotation cycle. Then we perform an n-fold cross-validation procedure, where we split data set into training and testing parts (by preserving ratio between them at F-1:1, where F is the number of folds) according to sequentially generated indexes (e.g. in case of 500 data points, fold 0 will represent indexes from 0 to 100 for the testing set, rest of the folds - indexes from 100 to 500 for the training set, fold 1 - 100-200 for the testing set, rest - 0-100 and 200-500 for the training set, etc). Then, if intra joint optimisation is selected within the cross-validation cycle, we perform outlier removal with sampled hyper-parameters using only the train subset within each train-test split. Hyper-parameters for ORM include the percentage of samples to be removed. After removing outliers, we train a model using a train set and make predictions on the test set.

All arrays with actual and predicted samples collected to be used after the fold-rotation cycle for the model accuracy estimation using specified metric. Since we are selecting test folds in order and making predictions on them, the predicted duration vector will be composed of prediction results composed of these folds. So, first, we collect the resulting metric together with hyper-parameters, actual and predicted labels. To collect data we use hash-array, which is represented as an array, where each element can be addressed by name and not by index as for conventional array. Then we perform the sorting procedure, which will order solutions according to the resulting metric, where we select the best combination of hyper-parameters. Furthermore, finally, we obtain the predicted duration vector by filtering data using the ORM method, training model on the train/test part and making predictions on the validation part.

### 3.3 Incident classification results

This section details the results of the first layer of the bi-level prediction framework related to the classification prediction findings, either via a standard binary classification with varying threshold analysis or via a multi-class classification enhanced by outlier removal procedures.

#### 3.3.1 Binary incident classification results using varying split thresholds

The first classification problem that we address is to predict whether an incident duration will be lower or greater than a selected threshold (we classify short-term versus long-term traffic incidents), which can then be used to supply the initial assessment needs of the traffic management centre (TMC) under fast decision times. For example, an operational clearance threshold for the Sydney TMC has been currently established at 45min based on previous operational field experience; however, choosing a fixed threshold for classification can have a significant impact on the results of any prediction algorithm and is highly dependent on the incident duration distribution chart (as represented in Fig. Fig. 3.1-g, h, i). Fig. 3.2 showcases the data split for the binary classification problem where the threshold  $T_c$  (dashed red line) is varying according to the two set-ups mentioned above: every 5 minutes ( $T_c \in \{20, 25, \dots, 70\}$ ). We name as Subset A all incident duration records which are lower or equal to



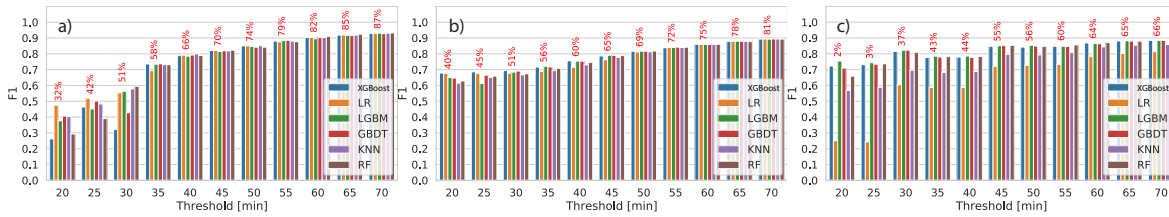


FIGURE 3.7: Incident duration classification using varying thresholds for a) data set AR b) data set M c) data set SF. The red percentage above each set of ML results indicate the percentage split of Subset A and B for that particular  $T_c$ .

$T_c$ , (if  $y_i \leq T_c$ ), and as Subset B all the incident duration records which are higher than  $T_c$  (if  $y_i > T_c$ ). Based on the variation of  $T_c$ , the size of Subsets A and B will have an impact on the prediction algorithms and this impact is further quantified.

The results of the binary classification approach of incident durations using a varying split threshold are detailed in Fig. 3.7 (for a 5-minutes frequency split) across all data sets. More specifically, Fig. 3.7 presents the F1 results obtained for each ML model that we have developed (XGBoost, LR, LGBM, GBDT, kNN, RF); we observe that other performance metrics have been calculated such as Accuracy, Precision and Recall and these are provided in the ??). For example, Fig. 3.7a) showcases the classification results for data set AR in which the blue bar represents the F1-result of the XGBoost classifier ( $F1=0.28$ ) when the data set has been split in Subset A containing incidents with a duration less than 20min (32% of all incident records fall in this subset) and Subset B containing incidents with duration higher than 20min (the rest of 68% of incident records). Therefore, the percentage numbers written in red above each ML result represent the percentage of records lower than the  $T_c$  threshold chosen for this experiment. The split around  $T_c = 20min$  is not ideal given the data imbalance (32% versus 68%) and the low F1 score; therefore further variations have been undertaken which have reported an increased  $F1 = 0.8$  for  $T_c = 45min$ . According to these results, if we use the best performing binary classifier, we need to select a threshold between 35 and 50 minutes because: a) it will reduce the imbalance between classes (and thus reduce the effects of imbalanced classification, which is vital for modelling when using a small data set); b) there is only a tiny improvement in F1-score after  $T_c > 40min$ ; c) it will be a reasonable split for short incidents lower in terms of field operation management. An exciting finding is revealed for  $T_c \in \{20, 25\}min$ : we record an overall lousy performance across all ML models in all data sets (F1-score less than 0.5) while some did not even take effect, such as GBDT; for this reason, we exclude from consideration any thresholds which provide an F1-score of less than 0.5. Furthermore, we set our minimum acceptable F1-score to 0.75, and any model performing lower than this threshold will not be considered for further optimisation. By analysing all sub-figures in Fig. 3.7 which provide both a good F1 score and class balance, we conclude that the optimal thresholds for the binary classification problem are the following: a)  $T_c = 40min$  for the arterial road network in Sydney (Fig. 3.7a:  $F1 = 0.79$  and a class balance of 66% for small incident duration), b)  $T_c = 45min$  for the motorway network in Sydney, (Fig. 3.7b:  $F1 = 0.75$ , class balance = 65%) and c)  $T_c = 45min$  for the San Francisco network (Fig. 3.7c:  $F1 = 0.83$ , class balance=55%).

The other important finding is the cases when  $T_c > 45min$  which present a significant improvement across all models on all performance metrics, with the best result being the one when Subset A incorporates all incidents lower than 70min (which represents the majority of incidents); this is easily

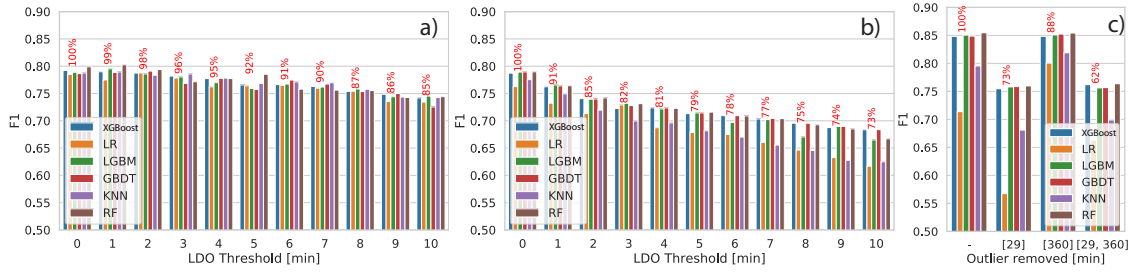


FIGURE 3.8: Outlier removal for a) data set AR b) data set M c) data set SF

explained by the fact that we use almost all the entire data set for training of the models. However, the binary classification can be a rough estimate. If TMCs need a higher prediction precision instead of incidents less than 45min or higher (which can last up to several days), then several regression and multi-class classification models are needed to provide more precise predictions. These will be further detailed in Sections 6 and 7. We will further use the detected optimal thresholds for each data set to perform the split between subset A and B in various scenarios of the incident duration regression problem.

Tree-based models yield similar results. However, in multiple cases (e.g. 35, 45, 50, and 60-minute thresholds for data set AR, 25, 30, 40, 60-minute thresholds for data set M), XGBoost produces a slightly better result than other tree-based models. Thus, we are selecting XGBoost as the best model for the incident duration classification.

### 3.3.2 Classification with outlier removal

After selecting the optimal thresholds for binary classification, we further assess the effect of: a) **low-duration outliers (LDO)** (which we define as reports of incidents with zero or less than a few minutes duration) and b) **high-duration outliers (HDO)** as in the San-Francisco data set, by trying different outlier removal procedures, as depicted in Fig Fig. 3.8.

For example, an LDO Threshold of 1min represents removing outliers below 1 minute (e.g. 0min) and the percentage above each removal test. For example, 99% indicates the number of samples remaining after such removal. Removing these outliers is essential since it represents errors in the incident reporting and may affect the accuracy of prediction. For example, Fig. 3.8a represents the LDO removal from the data set AR, up until 10min reported incident durations; by removing these outliers, we observe that the F1-score does not fall below the acceptable threshold of 0.75 until 5min (this indicates that removing all accidents reported with a duration of 0 or lower than 5min does not reduce the model performance. Therefore, we applied an LDO removal for all traffic incidents for this data set with a duration below 5min. For the data set M, the effect of LDO outlier removal is more significant, as depicted in Fig. 3.8b. This data set contains a lot of incidents with duration of 0 and 1 minute (which represents almost 15% of the entire data set); by removing these, we observe that the highest F1-score drops down to 0.74 across all ML models, which falls below the acceptable threshold for a good prediction accuracy). Therefore, we decide to remove only incidents with duration of 0min or 1min from this data set. Lastly, in the case of the San-Francisco data set, we have a completely different range of outliers since there are no incidents reported with a duration of fewer than 17 minutes (see Fig. 3.8c). There are multiple incidents cleared off at around 29min and 360min (as represented

as well in ??, which can be identified as HDO. However, by removing these HDO data points from the ML model training (representing almost 38% of all incident records), we observe a depreciation of the F1 score from 0.85 to 0.76 for XGBoost, while some models dropped to lower values below 0.7). Therefore, the removal of HDO for the San Francisco data set can not be adopted due to several reasons: 1) we cannot separate “rounded” duration from actually reported duration, 2) the amount of these values is almost half of the data, which becomes property of the data set, 3) these outliers still convey information related to the separation between short-term and long-term traffic accidents and 4) all models perform better when using the entire data set than with outlier removal, which makes the ORM procedure in this case non-necessary. Finally, we observe that the outlier procedure is highly related to the specificity of the data set and the incident area location, not by making default assumptions on either LDO or HDO.

### 3.3.3 Multi-class classification

While binary classification can provide fast insights in the overall incident duration, traffic incidents can have more precise duration definition and can be split (based on the histogram profiling) into short-term, mid-term, long-term. In this case one needs to solve a multi-class classification problem. We have split this problem in two subsection in which we analyse the impact of choosing three equally sized classes, versus quantile varying split thresholds and analyse the best approach.

#### Equally split multi-class classification

Firstly, we analyse the impact of using equally-sized classes (based on duration percentiles of almost 33% from each data set). We use F1-macro to assess the performance of a multi-class classification, defined as the unweighted average of class-wise F1-scores:

$$\text{F1-macro} = \frac{1}{N} \sum_{i=0}^N \text{F1-score}_i \quad (3.11)$$

where  $i$  is the class index and  $N$  is the number of classes. Table 3.2 contains the F1-macro scores across all three data sets for a 3-class prediction problem which can be calculated across each data set independently. For example,  $C1$  for data set AR in Sydney contains incidents between 0 – 24min, while  $C1$  for the SF data set contains incidents between 0 – 30min; similarly, the  $C3$  class for the SF data set contains substantial incidents which can reach up to 2,715min (45h) (this is consistently larger than 710min or 595min in Australia). The F1-macro score is aggregated across all classes, and a low value (below 0.5) indicates that we cannot use a 3-class split for the data set AR (F1-macro=0.35) and M (F1-macro=0.46), but we can do so for the data set SF (F1-macro=0.72). The significant difference between these data sets is the number of records (584 incident records for the data set AR versus 8,754 records for the data set SF), which may affect model performance. The precision of predictions on the data set indicates how many classes we can have to distinguish traffic incidents by duration. However, each data set’s specificity seems to dictate the best classification approach to be done and further justifies the need for a more refined regression prediction approach.

| Dataset     | $[0 - 33\%]_{C_1}$ | $[33 - 66\%]_{C_2}$ | $[66 - 100\%]_{C_3}$ | F1-macro(3-class) | F1 (2-class) |
|-------------|--------------------|---------------------|----------------------|-------------------|--------------|
| Data set AR | 0-24 min           | 25-44 min           | 44-710 min           | 0.35              | 0.79         |
| Data set M  | 0-24 min           | 25-54 min           | 54-598 min           | 0.46              | 0.74         |
| Data set SF | 0-30 min           | 31-71 min           | 72-2,715 min         | 0.72              | 0.85         |

TABLE 3.2: Multi-class classification results for equally-sized 3-class split

### Varying multi-class classification via quantile split

To analyse the effect of splitting data into more varying groups we performed a multi-class classification procedure using quantiles and the F1 results are provided in Figure 3.9 for three data sets: Figure 3.9a) when using the Arterial Roads in Sydney, Australia and b) when using M7 Motorway data set and c) when using the San Francisco data set. The result metric represents an average of F1-scores across classes, where multi-class classification performed as 3 one-vs-all classifications.

The low/high threshold matrix represented in Figure 3.9 indicates a 3-class split performance and allows for the modelling of different size groups separated by quantile thresholds. As an example, the Ox axis in Figure 3.9a) represents the first threshold split ranging from [10% to 80%], while Oy represents the second threshold split percentage ranging from 20% to 90%. The coloured dots represent the F1 scores obtained when splitting the data according to the two thresholds; for example, the combination quantile pair of [10%;20%] gives an F1 score of 0.34, meaning a multi-class split of data logs between the following three classes  $\{C_1 = [0 - 10\%], C_2 = [10\% - 20\%]$  and  $C_3 = [20\% - 100\%]\}$  does not provide good accuracy. Instead, when using the first quantile threshold of 0.3 and the second quantile threshold of 0.6 (meaning  $\{C_1 = [0 - 30\%], C_2 = [30\% - 60\%]$  and  $C_3 = [60\% - 100\%]\}$ ), we obtain the highest F1-macro score,  $F1 = 0.44$ .

In the case of M7 Motorway (see Figure 3.9b), we obtain the best performance for 20% and 60% quantile thresholds (meaning  $\{C_1 = [0 - 20\%], C_2 = [20\% - 60\%]$  and  $C_3 = [60\% - 100\%]\}$ ; 20%, 40%, 40% size grouping. Other options include  $\{20\%, 70\%\}$  and  $\{10\%, 60\%\}$  duration thresholds.

In the case of San-Francisco (see Figure 3.9c), we obtain the best performance for 10% and 90% quantile thresholds (meaning  $\{C_1 = [0 - 10\%], C_2 = [10\% - 90\%]$  and  $C_3 = [90\% - 100\%]\}$ ; this means that the best data split when using quantile thresholds for San Francisco is a  $\{10\%, 80\%, 10\%\}$  size grouping. This is highly explained by the incident distribution plots for the San Francisco area which is different than the rest of data sets.

To further see the impact on error by various incident duration groups we introduce the Quantiled Time-Folding and present the results in Figure 3.10. Incident reports are separated into equally-sized duration groups to perform the procedure of cross-validation (each 9 folds evaluated against 1 excluded fold, repeated 10 times). For all three data sets, incidents with the longest duration have the highest contribution to error, even though they represent only 10% of the data set. Considering this error, we may choose to use the hybrid classification-regression framework, where we perform regression only for intervals with acceptable prediction error. Quantiled Time-Folding can also be useful to see the contribution to error of every duration group and possible extrapolation error towards incidents with unobserved duration groups. Also, the RMSE metric showcased in Figure 3.10 is related to the scale of duration observed in the fold (e.g. high durations can easily translate in high errors), whereas if we

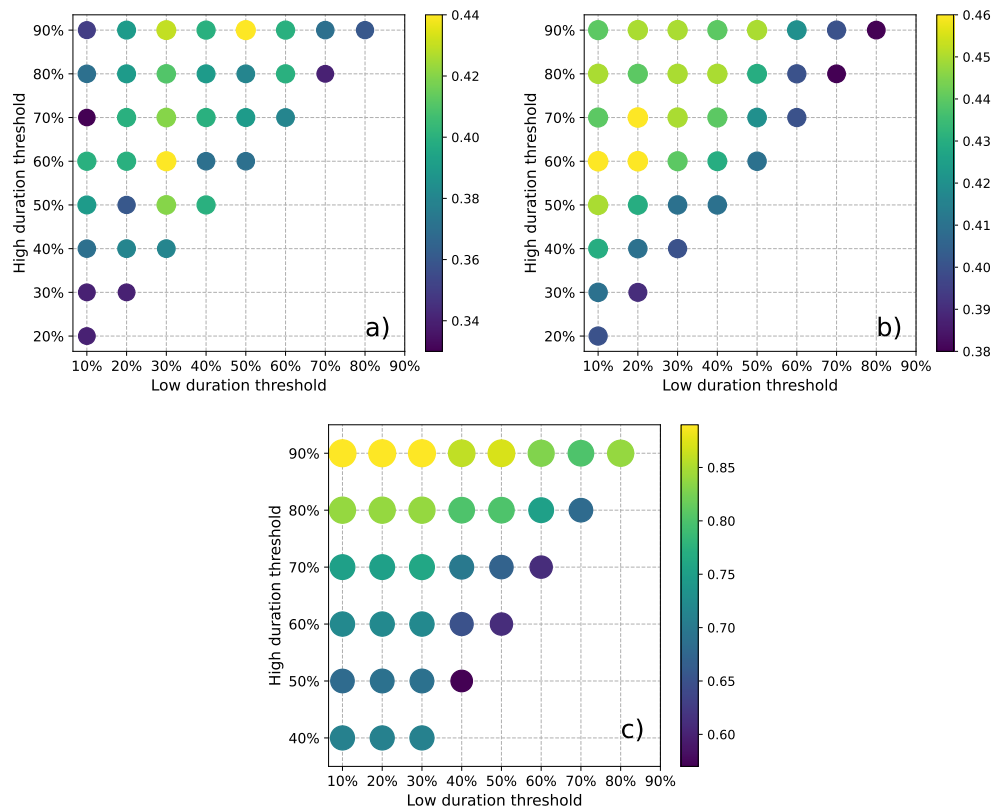


FIGURE 3.9: Multi-class (3-class) classification using quantile splits for a) data set AR  
b) data set M c) data set SF

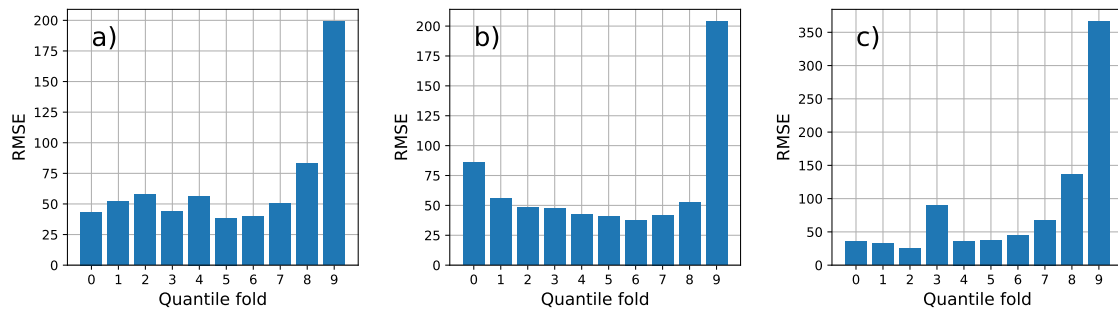


FIGURE 3.10: Regression using Quantiled Time Folding for a) data set AR b) data set M c) data set SF

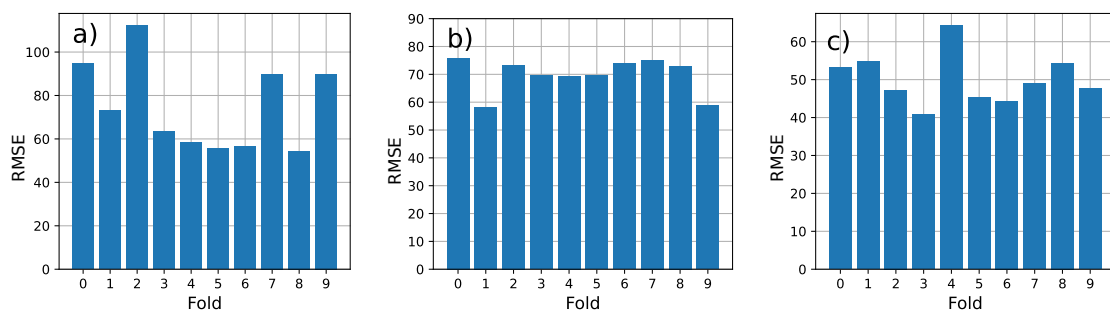


FIGURE 3.11: Regression using randomised 10-folds for a) data set AR b) data set M  
c) data set SF

adopt a regular 10-fold cross validation (see Figure 3.11), the RMSE error remains below 125.0 for most of the folds.

### 3.4 Incident duration prediction using regression: results

The final objective of the bi-level framework is to predict with an accuracy at the minute level the length of a freshly reported incident, regardless of its previous classification as either short, medium or large. Therefore, the second step of the bi-level prediction framework is to develop more advanced regression models that can adjust to each data set independently and over-perform baseline ML models previously used to solve classification problems. When training such regression models, a significant step is the size of the data set and the distribution of the target variable (incident duration). Due to the long tail distribution of incident duration and the class imbalance problem previously identified, we need to design and construct various regression models capable of learning from various types of data sets to make accurate predictions. However, with limited information (small data set size), the prediction results can be skewed (this effect of prediction skewing will be further discussed). This section first presents the regression results obtained across several scenarios of model training, validation and testing, followed by results of our proposed Intra-Extra Optimisation algorithm applied over all baseline ML models.

#### 3.4.1 Regression scenarios results and comparison

In order to find the best set-up that works for traffic incident prediction in TMCs, we test various regression scenarios (detailed previously in Section 3.2.6), which show the extrapolation performance for different ML methods. The outlier removal procedures (LDO, HDO) together with the classification thresholds (which separate short-term and long-term duration of incidents) are selected as described in Section 3.3.1-Section 3.3.2. The primary purpose of this section is to recommend the best scenario set-up for model training and validation when parts of the data set might be hidden. Table 3.3, Table 3.4 and Table 3.5 present the MAPE results for all 7 scenarios (All-to-All, AtoA, AtoB, BtoB, BtoA, AlltoA AlltoB) using all the Baseline ML models across all three data sets (and a dedicated winning regression model across each scenario - last column). Overall, XGBoost seems to be the best regression model in a majority of scenarios across data set AR and M (Table 3.3, Table 3.4): 1) the improvement from using XGBoost shows the lowest MAPE for scenario AtoA of 49.11 and 67.92 correspondingly (predicting short term incidents only using only short term training information), 2) XGBoost also the best performing model for All-to-All regression (59.36% and 85.98% MAPE correspondingly). The main difference between LGBM and XGBoost results is that LGBM struggles with extrapolation to lower values as seen in scenario B-to-A for all data sets: 292.68% vs 77.66% MAPE for data set A, 663.12% vs 180.77% MAPE for data set M, 166.06% vs 32.62% MAPE for data set SF for LGBM and XGBoost correspondingly.

In the SF data set, the LGBM is the best performer reaching a MAPE of 9.34% for the AtoA scenario (which is almost 10 times better than the same scenario for the M data set) and 33.16% MAPE for All-to-All scenario. This is a significant improvement that reveals what model is adapting to what data set, but most importantly, that each data set reacts differently to the seven scenarios.

| Model    | LGBM   | RF           | LR     | GBDT         | KNN    | XGBoost      | Best model |
|----------|--------|--------------|--------|--------------|--------|--------------|------------|
| AlltoAll | 82.76  | 117.28       | 110.99 | 113.41       | 107.79 | <b>59.36</b> | XGBoost    |
| AtoA     | 60.17  | 59.49        | 59.92  | 62.08        | 58.35  | <b>49.11</b> | XGBoost    |
| AtoB     | 64.46  | 64.39        | 64.34  | <b>63.82</b> | 64.68  | 64.39        | GBDT       |
| BtoA     | 292.68 | 381.61       | 367.16 | 348.09       | 349.62 | <b>77.66</b> | XGBoost    |
| BtoB     | 29.52  | <b>25.03</b> | 45.14  | 46.26        | 43.82  | 27.55        | RF         |
| AlltoA   | 117.78 | 121.82       | 175.48 | 176.71       | 120    | <b>51.18</b> | XGBoost    |
| AlltoB   | 34.39  | 37.47        | 32.11  | <b>31.67</b> | 35.57  | 37.46        | GBDT       |

TABLE 3.3: MAPE results for all 7 scenarios on data set AR

| Model    | LGBM   | RF     | LR           | GBDT         | KNN    | XGB           | Best model |
|----------|--------|--------|--------------|--------------|--------|---------------|------------|
| AlltoAll | 135.59 | 226.6  | 229.53       | 229.46       | 229.82 | <b>85.98</b>  | XGBoost    |
| AtoA     | 95.89  | 95.38  | 107.29       | 104.87       | 105.26 | <b>67.92</b>  | XGBoost    |
| AtoB     | 68.78  | 69.01  | 69.49        | <b>68.62</b> | 69.79  | 68.69         | GBDT       |
| BtoA     | 663.12 | 939.59 | 818.08       | 878.47       | 854.81 | <b>180.77</b> | XGBoost    |
| BtoB     | 34.14  | 51.02  | 52.33        | 50.99        | 48.68  | <b>31.18</b>  | XGBoost    |
| AlltoA   | 233.48 | 406.43 | 387.25       | 398.13       | 402.02 | <b>76.71</b>  | XGBoost    |
| AlltoB   | 34.38  | 34.34  | <b>34.21</b> | 34.48        | 36.89  | 34.98         | LR         |

TABLE 3.4: MAPE results for all 7 scenarios on data set M

In the following, we provide a summarised comparison across a selection of few scenarios and their performance.

**Scenario AtoA** uses short-term traffic accidents (below  $T_c$ ) for both training and the prediction. XGBoost shows a significant performance for AR and M data sets compared with other scenarios; more specifically, they outperform by 10% all models in data set AR (MAPE=51.2) and 30% all models in dataset M (MAPE=68.4). For the SF data set, the improvement is even larger (MAPE=12.7), but XGboost loses ground over LGBM, which reaches a MAPE=11.0. The comparison of scenarios AtoA and AlltoA shows that adding incidents with a longer duration can severely affect the prediction performance across all data sets, regardless of the size or location of the incident logs. For the best prediction performance on data sets AR, M and SF, we need to split the data and use separate models for the short-term incidents as predictions become skewed towards longer incident duration. Thus, if we predict short-term incidents using only short-term incidents data logs, we obtain a higher accuracy across all data sets.

| Model    | LGBM         | RF           | LR     | GBDT         | KNN    | XGBoost      | Best model |
|----------|--------------|--------------|--------|--------------|--------|--------------|------------|
| AlltoAll | <b>33.16</b> | 36.88        | 128.42 | 41.85        | 64.24  | 37.03        | LGBM       |
| AtoA     | <b>9.34</b>  | 11.91        | 16.07  | 12.56        | 14.05  | 11.44        | LGBM       |
| AtoB     | 68.08        | 65.77        | 67.21  | <b>65.53</b> | 66.26  | 65.84        | GBDT       |
| BtoA     | 166.06       | 191.55       | 389.07 | 211.61       | 302.46 | <b>32.62</b> | XGBoost    |
| BtoB     | <b>23.69</b> | 28.76        | 70.18  | 31.08        | 37.6   | 27.61        | LGBM       |
| AlltoA   | 45.35        | 50.74        | 218.49 | 60.03        | 99.06  | <b>35.49</b> | XGBoost    |
| AlltoB   | 24.28        | <b>23.97</b> | 45.08  | 25.49        | 30.82  | 24.78        | RF         |

TABLE 3.5: MAPE results for all 7 scenarios on data set SF



**Scenario AtoB** is unique because regression models are trained on Subset A, which contains short-term incident duration logs while they are trying to predict long-term incidents; therefore, the performance is much worse than for AtoA scenario since incidents with long duration are much scarcer and have unique traffic conditions. BtoB scenario shows lower error than AtoB across all three data sets (e.g. BtoB provides 23.69% MAPE and AtoB provides 65.53% MAPE for best models for data set SF). Vice-versa, **Scenario BtoA** shows very high extrapolation errors across all methods to lower values. Adding short-term incidents into the training set of long-term incidents (when we move from BtoA to AlltoA scenario) significantly reduces the error (76.71% MAPE for scenario AlltoA, data set M using XGBoost), but it is still significantly higher than for AtoA scenario (67.92% MAPE for M data set using XGBoost). **Scenario BtoB** shows better performance (e.g. MAPE=31.18% for data set M using XGBoost) than using data addition (such as the case of AlltoB, where MAPE=34.21% using best model) or any extrapolation (as in the case of AtoB, where MAPE=68.62% using best model). By comparing scenarios AtoB and AlltoB we observe a significant performance improvement when adding data for long-term incidents and predicting subset B (from 63.82% to 31.67% MAPE for dataset AR using best model), where error is still higher than for BtoB (25.03%, AR, best model). **Scenario BtoA** shows high prediction errors across all scenarios highlighting a bad extrapolation accuracy when predicting short-term incidents duration using long-term traffic incident data. It means that prediction of the duration of short-term incidents should be performed separately from long-term incidents. Thus, we can't use long-term incidents to predict the duration of short-term incidents and vice versa if we are looking at maximising model performance with limited data set; the second reason lies mainly in different traffic behaviour along with severe accidents that can last for several hours which are harder to clear off - these require similar previous events in order to be predicted for their duration.

### Fusion framework for the incident duration prediction

In comparison to the above proposed framework, we also present a fusion framework approach, which can be applied when the incident duration category is unknown. When an incident occurs, the incident duration category is not known, but we have a historical data on traffic incidents which allows us to predict the incident duration category and apply specialised regression models (oriented towards the prediction on subsets of short-term and long-term incidents). We further propose two possible approaches to this problem:

- **the pipeline approach** (see Fig. 3.12a ): we train a classification model using a historical data available to predict the incident duration category. Then we predict the incident duration category using the available incident reports. This prediction decides which model we need to use for a further regression (either specialised on short-term or long-term incident duration prediction). The prediction result of the specialised model is then considered to be the final prediction. Specialised regression models are trained on their corresponding subsets. In this case, the decision about the incident duration class is made by the classification model only, which becomes the most important part of the model that is highlighted by significantly improved results (see Tables 3.3 to 3.5).
- **the fusion approach** (see Fig. 3.12b ): instead of relying on the classification model to decide on the incident duration subset, we place a decision-making function on the additional "fusion



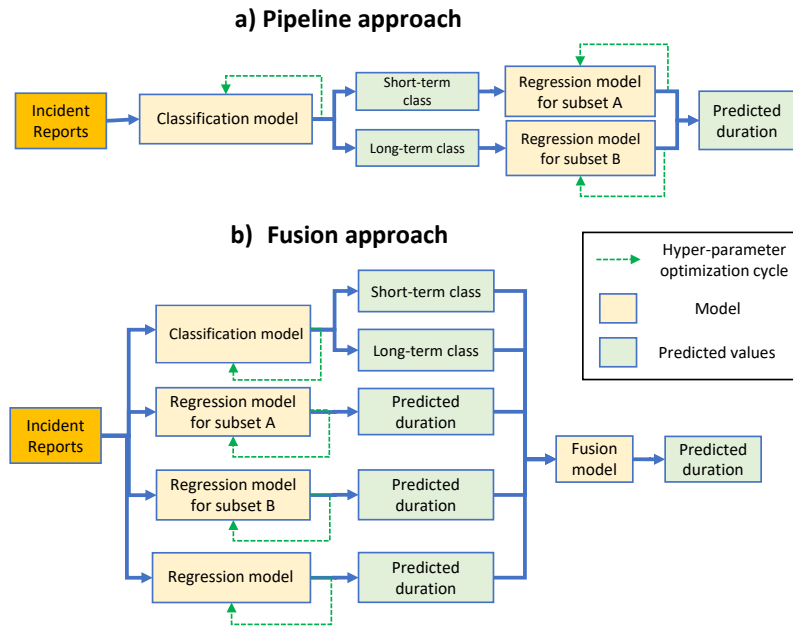


FIGURE 3.12: Pipeline (a) and fusion (b) approaches for the bi-level framework structure

model", which is the global regression model; it now receives the prediction results from the classification model, the regression models specialised on short-term and long-term incidents (subsets A and B) and from the regression model trained on historical data of traffic incidents regardless of the incident duration group. After training all these models on historical data, we perform the incident duration prediction on this historical data. We then use these predictions (such as the predicted incident class, the incident duration predicted by short-term incident duration regression model, the incident duration predicted by the short-term incident duration regression model, and the incident duration predicted by the regression model) in order to train the global fusion model to make a final prediction of the incident duration; we call this the global fusion model and the predicted duration is a result of multiple models fused in a centralised architecture.

The fusion approach can be perceived as the ensemble model, which allows to solve the computational problem of model training. Ensemble models may perform better than single models due to three main reasons: (Dietterich, 2000): a) statistical: without sufficient data, a model can find multiple hypothesis about the data approximation which has the same accuracy. Each of these hypotheses can lean towards its local optima. By averaging hypotheses, we may find a better approximation of the data; b) computational: many machine learning models may get stuck in a local optima (e.g. stochastic gradient descent in the case of neural networks or the greedy split finding in the case of decision trees). An ensemble constructed by models performing local search from many different starting points may provide a better prediction performance than the individual models (Dietterich, 2000; Ballings et al., 2015), c) representational: each model forms an approximation (representation) of the data, which forms a local representation hypothesis. By combining models it is possible to extend the space of representable functions.

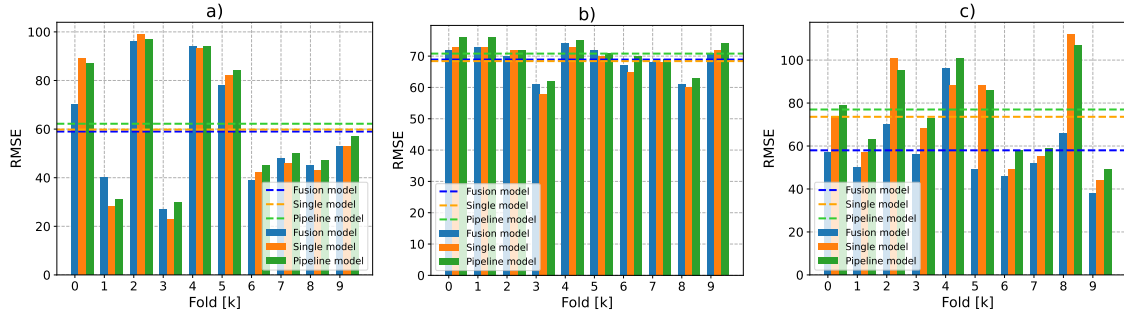


FIGURE 3.13: Comparison of the fusion and single model performance for a) data set AR b) data set M c) data set SF. Dashed lines represent the average RMSE score across all folds for each corresponding model

In the case of a bi-level framework we have statistical, computational and representational reasons to expect a better performance from using an ensemble model rather than a single model, since we use different kinds of models on different subsets (in our case a simple regression model, a classification model, a regression model for subset A, a regression model for subset B). In other words, by splitting the data and by using multiple models we obtain models that are having different local optima (subset A, subset B models) and a different representation of the data (classification and regression models); in this way we can obtain a better prediction performance using model ensemble than using individual models.

Finally, we compare the fusion model, single regression model (e.g. for the data set SF it is the model with the best performance for the task of All-to-All regression) and the pipeline model (where the choice between the regression models depends on the predictions from the classification model) performance on all three data sets in Fig. 3.13. We evaluate all model performance on each fold in a randomised 10 fold cross-validation. We observe that the fusion model performs at least not worse than a single model on all three data sets. We use XGBoost as a fusion model. We also use the corresponding best models for each subset of each data set (see Tables 3.3 to 3.5) with hyperparameter optimisation (e.g. LightGBM as a single model, performing All-to-All regression task for the data set SF, RandomForest as a best classification model for data set SF according to Fig. 3.7). There is a subtle difference in the average RMSE score among the folds for data set A (see Fig. 3.13a) where the average RMSE for the fusion model is 59, for the single regression model is 59.8, for the pipeline model is 62.2). The same is for the data set M (see Fig. 3.13b) where 68.9, 68.4 and 70.8 are the average RMSEs for the fusion, single and pipeline models correspondingly). There is a significant improvement in the average RMSE score for data set SF (see Fig. 3.13c) with an improvement from 73.6 to 58 of the average RMSE when using of the fusion model instead of the single model); the pipeline model didn't show any improvement in the model performance. Overall, results show that data availability (the amount of information available about the incidents, which is high for the data set SF) can significantly affect the performance of the fusion model.

Given the observed performance (from a subtle difference to significant improvement) we recommend to use the fusion approach within bi-level framework for the task of the incident duration prediction.

### 3.4.2 Outcomes and recommendations

From an operational perspective the scenario All-to-All is the ideal situation when traffic management centres would have in their data base both long term and short-term incidents. However, from an operational perspective, several records of short incidents for example and not being kept all the time, while long incidents are often being transferred to various other division if they last more than one day, and they become more of a road infrastructure problem rather than an operational problem which requires constant intervention.

Scenario modelling shows that the baseline ML models are not improving when facing incident duration extrapolation or data addition (e.g. AlltoA versus AlltoA, BtoB versus AlltoB); these two training set-ups badly affect the model performance extrapolating in any direction.

By evaluating regression scenarios, we highlight the importance that incidents from different duration groups need to be modelled separately in order to significantly improve the accuracy of duration predictions (see (see [Table 3.5](#)) for SF data set: if we use all available data, to predict the incident duration, we will have  $MAPE = 33.16\%$  (lower is better), but if we managed to categorise incidents into the short-term group, we could model these incidents with only  $9.34\%$  error, which is a significant improvement. Also, the classification may point us to which data we need to include in the modelling because if we use all data to predict the duration of the short-term incident (scenario AlltoA), we will have a much higher error  $45.35\%$  MAPE than just using the short-term incident for modelling. From the comparison regression of extrapolation scenarios (e.g. scenarios AlltoA versus AtoA), we see how significant can be the impact of having incidents with long duration in the training set when we need to predict the duration of short term incidents, and therefore ML methods become biased towards long-term incidents, which significantly reduces their performance. If we can perform incident duration regression, then we are able to perform incident duration classification as well. We can do this before performing the regression in each group. In other words, our scenario modelling shows modelling advantages of classifying incidents into duration groups.

Therefore, it is essential for the bi-level framework and traffic incident duration prediction to use separate models for short-term and long-term traffic incidents. Moreover, tree-based methods significantly outperforming LR demonstrates that traffic incident regression is a complex non-linear problem that requires more advanced investigations. This aspect was the one that motivated our research to further improve and build a better ML framework for any type of incoming data set, and the results of this novel IEO-ML framework are further detailed in the following section.

### 3.4.3 Regression results for proposed IEO-ML model

In this section, we employ our proposed Intra-extra joint optimisation approach previously presented in [Section 3.2.9](#) and we further present the results of the All-to-All regression scenario, with a log-transformation of incident duration and several outlier removal techniques such as the LocalOutlier-Factor (LOF) and the IsolationForest (IF), previously described in [Section 3.2.7](#). All results across the three data sets are presented in [Table 3.6-Table 3.7-Table 3.8](#).

For the data set A ([Table 3.6](#)), we observe a significant impact of using the log-transformation of the incident duration vector via the resulting MAPE (see Unprocessed versus Log columns). Since the log-transformation provides a significant improvement among majority of ML models, we decide

| $ML_j$      | Log         | Unprocessed | iIF-Log | eIF-Log | eLOF-Log    | iLOF-Log    | Best approach        |
|-------------|-------------|-------------|---------|---------|-------------|-------------|----------------------|
| LGBM        | 80.4        | 81.1        | 79.9    | 82      | <b>78.4</b> | 80.8        | <b>eLOF-Log-LGBM</b> |
| RF          | 80.3        | 121.9       | 79.5    | 80.7    | <b>78.5</b> | 79.1        | <b>eLOF-Log-RF</b>   |
| LR          | <b>80.0</b> | 128.4       | 80.4    | 81.6    | 80.5        | 80.5        | Log-LR               |
| GBDT        | <b>79.4</b> | 128.2       | 82.0    | 81.3    | 81.4        | 83.4        | <b>Log-GBDT</b>      |
| KNN         | 82.9        | 127.4       | 82.3    | 86.2    | 81.7        | <b>81.3</b> | iLOF-Log-kNN         |
| XGBoost     | <b>59.4</b> | 61.1        | 60.8    | 59.8    | 60.9        | 59.9        | <b>Log-XGboost</b>   |
| Best $ML_j$ | XGBoost     | XGBoost     | XGBoost | XGBoost | XGBoost     | XGBoost     |                      |

TABLE 3.6: MAPE results for All-to-All scenario of data set A, using different ORM approaches and incident duration transformation, via the proposed IEO-ML approach.

| $ML_j$      | Log          | Unprocessed | iIF-Log      | eIF-Log      | eLOF-Log | iLOF-Log | Best approach          |
|-------------|--------------|-------------|--------------|--------------|----------|----------|------------------------|
| LGBM        | 124.6        | 138.0       | <b>123.6</b> | 126.8        | 125.1    | 124.1    | <b>iIF-Log-LGBM</b>    |
| RF          | 126.3        | 238.6       | 126.6        | <b>125.7</b> | 127.1    | 126.6    | <b>eIF-Log-RF</b>      |
| LR          | 130.7        | 245.9       | <b>129.8</b> | 129.9        | 131.1    | 131      | <b>iIF-Log-LR</b>      |
| GBDT        | <b>126.7</b> | 240.1       | 126.9        | 126.7        | 127.2    | 126.9    | <b>Log-GBDT</b>        |
| KNN         | 139          | 248.2       | <b>135.1</b> | 137          | 139.4    | 138.2    | <b>iIF-Log-KNN</b>     |
| XGBoost     | 78.6         | 113.2       | <b>77.5</b>  | 80.6         | 78.3     | 79.6     | <b>iIF-Log-XGBoost</b> |
| Best $ML_j$ | XGBoost      | XGBoost     | XGBoost      | XGBoost      | XGBoost  | XGBoost  |                        |

TABLE 3.7: MAPE results for All-to-All scenario of data set M, using different ORM approaches and incident duration transformation, via the proposed IEO-ML approach.

to use it in our outlier removal scenarios. When comparing results across all models, both regular and re-enforced by our IEO approach (column comparison - see Best  $ML_j$  results), we observe that XGBoost is the best performing baseline model for this data set reaching a 59.4 MAPE. Furthermore, when comparing results across regular ML models versus our proposed IEO-ML enhancements (row comparison), then the extra optimisation approaches seem to outperform the intra optimisation approaches (see iIF-Log versus eIF-Log and eLOF-Log versus iLOF-Log columns). The last column indicates the best approach that won across all proposed IEO approaches where for example, eLOF-Log-RF model is read as the extra optimisation method applied together with the Local Outlier Factor and Random Forest over the log scale data transformation; for this data set A results indicate a similar performance between using baseline ML models with log transformation versus enhanced IEO-ML - for example the joint optimisation provides an improvement (eLOF-log-LightBGM, eLOF-log-RF) versus the cases cases when only the baseline ML with the log-transformation was used (e.g. Log-LR, Log-BDT). However, the A data set is very small and has a special behaviour when compared to the others as further results revealed.

For the data set M (Table 3.7), when we use Log-transformation, we observe very high MAPE scores (100% and higher), except for XGBoost, which provides a MAPE of 78.6%. When comparing the models with each other against the IEO enhancements as well (column comparison), using XGboost as a baseline seems to over-perform all the other approaches, with the best results being a MAPE=77.5 for iIF-Log-XGBoost. When comparing against the proposed approaches (row comparison), the Intra joint optimisation using Isolation Forest in log-transform shows the best performance on this data set for four models (iIF-Log-LGBM, iIF-Log-LR, iIF-Log-kNN, iIF-Log-XGBoost), which

| $ML_j$      | Log         | Unprocessed | iIF-Log     | eIF-Log     | eLOF-Log | iLOF-Log    | Best approach           |
|-------------|-------------|-------------|-------------|-------------|----------|-------------|-------------------------|
| LGBM        | 29.9        | 32.6        | 29.7        | <b>29.5</b> | 30.2     | 29.9        | <b>eIF-Log-LGBM</b>     |
| RF          | 28.9        | 38.7        | <b>28.7</b> | 28.9        | 28.8     | 28.9        | <b>iIF-Log-RF</b>       |
| LR          | 72.6        | 140.5       | 72.8        | 73.1        | 73.3     | <b>72.4</b> | <b>iLOF-Log-LR</b>      |
| GBDT        | <b>31.2</b> | 46.3        | 31.5        | 31.4        | 32.4     | 32.2        | <b>Log-GBDT</b>         |
| KNN         | <b>61.5</b> | 108.6       | 61.7        | 62.5        | 62.2     | 61.8        | <b>Log-KNN</b>          |
| XGBoost     | 31.7        | 35.1        | 31.9        | 31.6        | 32.7     | <b>31.0</b> | <b>iLOF-Log-XGBoost</b> |
| Best $ML_j$ | RF          | LGBM        | RF          | RF          | RF       | RF          |                         |

TABLE 3.8: MAPE results for All-to-All scenario of data set SF, using different approaches for ORM and incident duration transformation, via the proposed IEO-ML approach.

can be attributed to data set data structure - outliers can be better analysed using tree-based outlier removal methods rather than distance-based LOF. For the majority of models (4 out of 6), our proposed joint optimisation algorithm obtains the best results for this data set.

For the data set SF (Table 3.8), we observe two competing models - LGBM and Random Forests with a prevalence for Random Forests (column comparison - see Best  $ML_j$  results). Also, we observe a considerably lower MAPE score for the best performing models which reached the lowest threshold of 28.7 across all the data sets used in this study. This reveals the power of more complete and larger data sets which can significantly improve the model performance. When comparing the IEO approaches (row comparison), the intra joint optimisation shows improvement across three models and more specifically for the best performing model on this data set, RF. One consistent finding across all results is the fact that the log-transformation of the incident duration vector should be used at all times for incident duration prediction since it significantly improves predictions accuracy; this is mostly related to the long tail distribution and extreme outliers which can affect the final errors in the model performance evaluation. Overall, the best performing models are considered to be XGBoost and Random Forests.

**To summarise**, every data set has its specifics in the data structure, which make some models and outlier removal methods performing better than others. Thus, it is necessary to deploy different models and outlier removal approaches on every data set. Conventional models (KNN and Linear Regressions) show the highest error which is almost twice in comparison to tree-based models. Thus, tree-based models are preferred options for solving the incident duration prediction together with adapted optimisation and outlier techniques. Overall, we proved that our proposed intra joint optimisation is improving the regression results across multiple data sets (especially data sets M and SF in 7 out of 12 cases). The joint optimisation of the model together with the outlier removal method shows a significant improvement in majority of cases (12 out of 18) across all three data sets.

### 3.4.4 Bi-level framework implementation

The code for the bi-level framework exploring previously described scenarios can be found by the link:

<https://github.com/Future-Mobility-Lab/bi-level-framework>

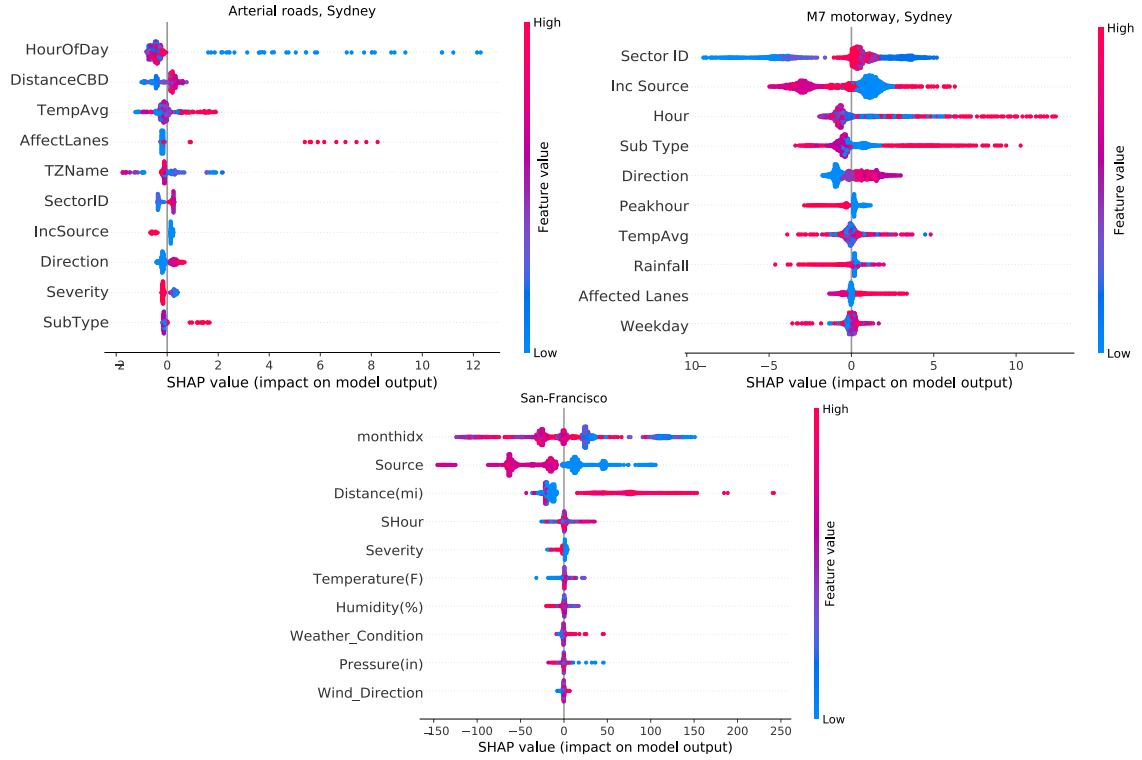


FIGURE 3.14: Feature importance for All-to-All regression using XGBoost for a) Arterial roads, Sydney, Australia b) M7 motorway, Sydney, Australia c) San-Francisco, USA

### 3.5 Feature importance impact and evaluation

Finally, we evaluate the feature importance using a Shapley value calculation in order to estimate the contribution of each feature to the final prediction score. Each point related to a feature is shown in Fig. 3.14 and represents the SHAP value score (Oy-axis), coloured by its value (from low to high), while the Ox-axis shows the impact of that feature information on the entire prediction output. The used models for this feature importance analysis are the winning models of each data set (A, M, or SF) as previously discussed.

The hour-of-the-day when the incident started is among the top 5 features sorted by importance (ranked on the 1<sup>st</sup> place for data set A, 3<sup>rd</sup> for M and 4<sup>th</sup> for SF). For example, Fig. 3.14a) showcases that as the hour of the day increases (getting closer to midnight) the traffic durations are lower as the congestion is lower and rescue teams arrive faster to the accident location; this is the opposite on the motorways as Fig. 3.14b) reflects that rescue teams have a harder time reaching the incident location in the evening, which is mostly explained by the high distance of the motorway from the local incident management centre. The incident reporting source also has a high significance (ranked as 7<sup>th</sup> most important for A, 2<sup>nd</sup> for M, 2<sup>nd</sup> for SF). The Ox-axis on SHAP plots represents the impact on model output (e.g. the effect on the predicted duration value). Even though the average temperature is considered significant, its effect on the regression model output is very small  $[-5min; +5min]$  for data set AR,  $[-5min; +5min]$  for data set M,  $[-25min; +25min]$  for data set SF. The distance from CBD (DistanceCBD) is important in the data set A, as it can point at some problematic areas, therefore causing



a higher incident duration. The number of affected lanes is also an important feature for incident duration prediction on arterial roads in Sydney. The model outputs for the M7 motorway revealed that is highly dependent on the sector ID (similar to the traffic zones in the data set A), which may be linked to the nature of the location or to the distance from incident management agencies. The average daily temperature also affects predictions (3<sup>rd</sup> place in A, 7<sup>th</sup> in M and 6<sup>th</sup> in SF). Weather factors (rainfall) are found to play a significant role in the M and SF data sets (humidity and barometric pressure may be predictors of rainfall). Different incident sub-types in the M data set (e.g. car, motorcycle, truck, multi-vehicle) contribute to the difference in the accident duration. Severity is weakly connected to the incident duration in the A and SF data sets. It is important to note that the SF data set contains 49 features, but 39 are of very low importance for the incident duration prediction. The length of the affected road segment (Distance in SF) may also be an essential feature which is not found in Sydney data sets. Overall, the specificity of each data set is reflected once again not only in the models that may be more successful than others but also in the way that the same model can provide various feature importance due to each country, their unique landscape and different way of dealing with the disruptions.

### 3.5.1 Short-term vs long-term incident duration prediction feature importance

We further perform a comparison of feature importance for the duration prediction of short-term vs long-term traffic incidents, across all data sets.

#### Arterial Roads Feature Importance, Sydney Australia.

Fig. 3.15 showcases the Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of Arterial roads, Sydney, Australia. When analysing long-term incidents, one important observation is the direct influence of the number of affected lanes on the severity and duration of disruptions. However, this feature is found to have low importance for short-term incidents. The farther short-term incidents happen from the CBD, the longer it takes to clear them off. The location of the incident is extremely important for both long and short term incidents, most likely due to the easiness to reach the affected location by the intervention teams. Another important factor affecting the short term incidents in Sydney seems to be the travel patterns for commuting [month of the year, day of week, sectionID, section capacity]. Also, the DayOfWeek (value ranges from 0 to 6), we see that the higher the value (closer to the end of the week), the longer it takes for the incident to clear. Also, some sectors reflected by the SectionID feature demonstrate a lower incident duration, which may highlight that some specific areas of the city are less affected by traffic incidents.

#### Motorway Feature Importance, Sydney Australia.

Fig. 3.16 showcases the Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of M7 Motorway, Sydney, Australia. One immediate observation is the fact that the data has 3 sources of reporting, and this can be seen as three different distributions in the top 1 most important feature ranked in Figure 4a). The source reporting the incidents seems to be the one factor which influence the most the incident duration. When comparing the top features for both short versus long term incidents, these are almost the same in both subsets: average temperature,

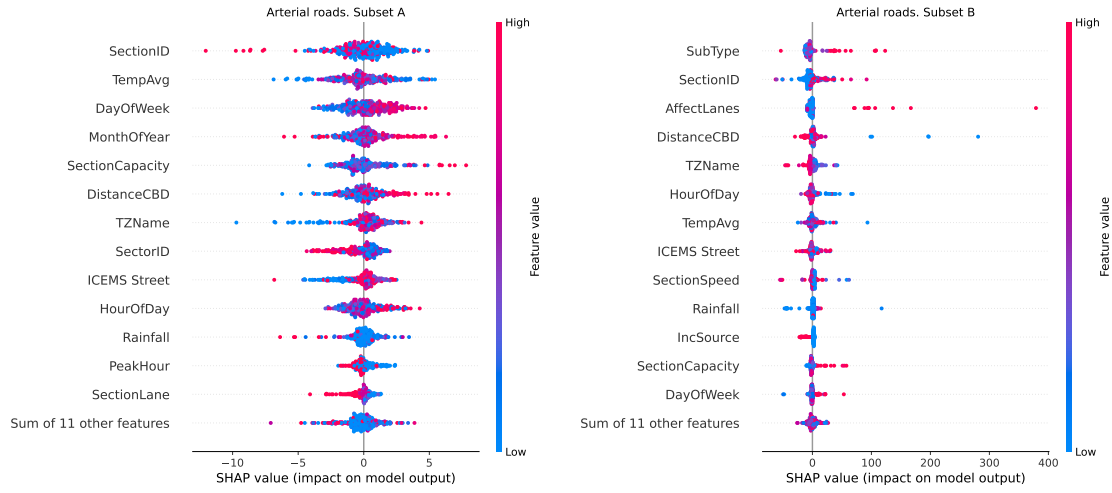


FIGURE 3.15: Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of Arterial roads, Sydney, Australia

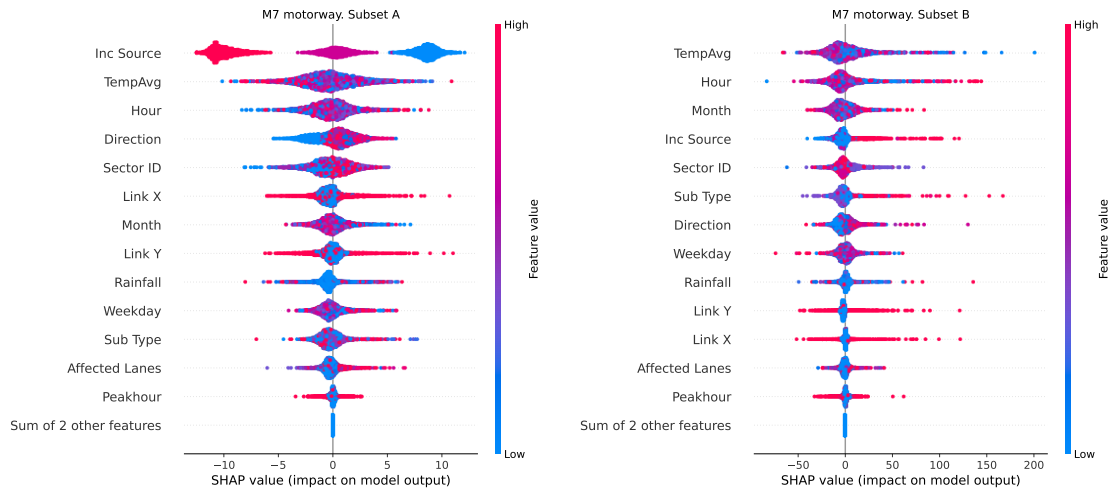


FIGURE 3.16: Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of M7 Motorway, Sydney, Australia



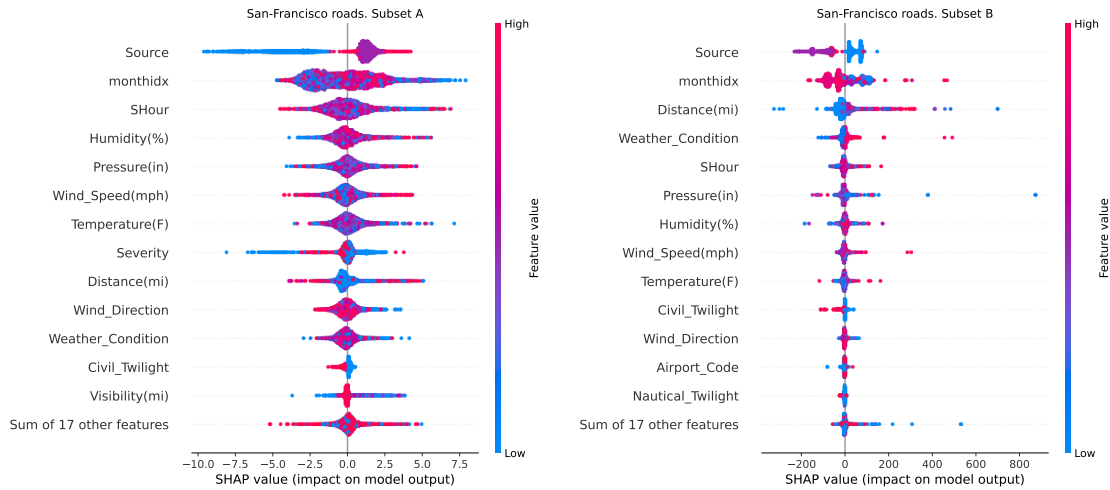


FIGURE 3.17: Feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of San-Francisco, USA

the hour when the incident happened, the Sector ID, the direction of travel and the source of information that reported the incidents. Overall, for this data set, same features can be collected for both types of incidents.

### San Francisco Feature Importance, U.S.A.

Finally, Fig. 3.17 showcases the feature importance for All-to-All regression using XGBoost for a) short-term incidents b) long-term incidents of San-Francisco, USA. This data set is very different than the rest, but as in the case of M7 motorway, the source reporting the incident seems to be most important factor affecting the duration – this is mostly related to the way the information is received to the centre (from road users, from local traffic agents, from video camera surveillance, etc.). We observe that short-term and long-term incident are very different in their nature and incident characteristics found to have different importance in the prediction of the incident duration. For the SF data set, the most important features are Source, monthidx, Shour, regardless of the incident duration. In terms of large accidents however, the distance from the CBD is very important while for small accident the humidity plays an important factor ranking 4th (which might indicate that weather in San Francisco can cause small traffic accidents to happen often). Overall, despite all data sets being different, their specificity and feature important is highly related to their setup, the location of the network and the way the management centre received and handle the disruption. In order to help improve the prevention techniques more effort should be invested in understand which source of incident reporting causes the most errors overall and why.

## 3.6 CONCLUSIONS

This paper proposed a novel bi-level framework for predicting the incident durations via a unique combination of baseline machine learning models (for both classification and regression), together with an outlier removal procedure and a novel intra-extra joint optimisation technique. The accuracy and

importance of the proposed approach have been proved via three different data sets from 2 countries (Australia and the United States of America) under several scenarios for testing and validation.

**Major contributions:** Firstly, regarding the classification prediction of incidents into short versus long-term: we found that the optimal duration classification thresholds are similar among the three different data sets: 40min for data set AR, 45min for M, 45min for SF. Sydney TIMS also found 45 minutes to be the threshold for incident removal performance evaluation via their on-the-field expertise; this represented a confirmation that our threshold split is in coherence with realistic operational rescue times. Secondly, the best performing and robust models in the classification and regression experiments were the tree-based models (XGBoost, RandomForest, etc.). Thirdly, our extensive regression scenarios demonstrate that the short-term and long-term traffic accidents should be modelled separately. Otherwise, we will observe a drop in performance due to the adverse effect of different scale values in the training set on the model output. Fourthly, our proposed IEO-ML approach outperformed baseline ML models in 12 out of 18 cases (66%), showcasing its strong value to the incident duration prediction problem. Finally, when evaluating the feature importance, we showed that features related to time, location, type of accident, reporting source and weather are among the top 10 critical features in all three data sets. By improving the precision of the most important and removing non-important features from the incident reports, TIMS can significantly improve the quality of data acquisition.

**Limitations of this study:** One of the biggest challenges when studying the problem of incident duration prediction represents data availability. In most cases, the privacy around traffic incidents represents the main reason why data sets are not released publicly. For example, the two data sets from Australia are private and have only been released for the purpose of this study, whereas only the San Francisco data set is made open publicly. Many other countries around the world have not yet fully released their incident logs, and this represents a challenge for this topic. However, if more incident data logs become available, they can represent a good test best for our approach.

Regarding the model performance, we make the observation that the performance of ML methods is highly affected by the data sets and the used methodology. Our approach shows a better performance for 4 of 6 methods in the case of San-Francisco, but if looked more precisely into details, KNN (where there is no improvement) produces an error that is twice as large as the best performing model (GBDT). The same is for data set A when using the LR method. And, with only GBDT left with no improvement may point to the fact that GBDT is robust to outliers and does not need outlier removal (as observed on all three data sets). As can be seen for the data set A, where MAPE is high (80%), there is a very weak connection between features and the labelled data, and thus the performance for all methods is poor. Therefore, there is not much effect from the outlier removal approach on poor data sets or for methods that are weaker by design.

**Future research** can be related to the usage of traffic simulation with information on predicted traffic incident duration included in the decision making process during route planning. For example, the vehicle can consider that a traffic incident is short-term and assume that it will be cleared before arriving at the incident location and therefore reduce its travel time by not planning a route around the incident site. Furthermore, the cost of prediction error and the benefit of traffic accident duration estimation can be estimated from the simulation model, where occasional traffic accidents happen within traffic flow. Also, the benefit of this approach can be estimated for online route planning and

not only at the time of the departure.

Providing additional results for the threshold variation along all data sets such as (Accuracy, Precision and Recall).

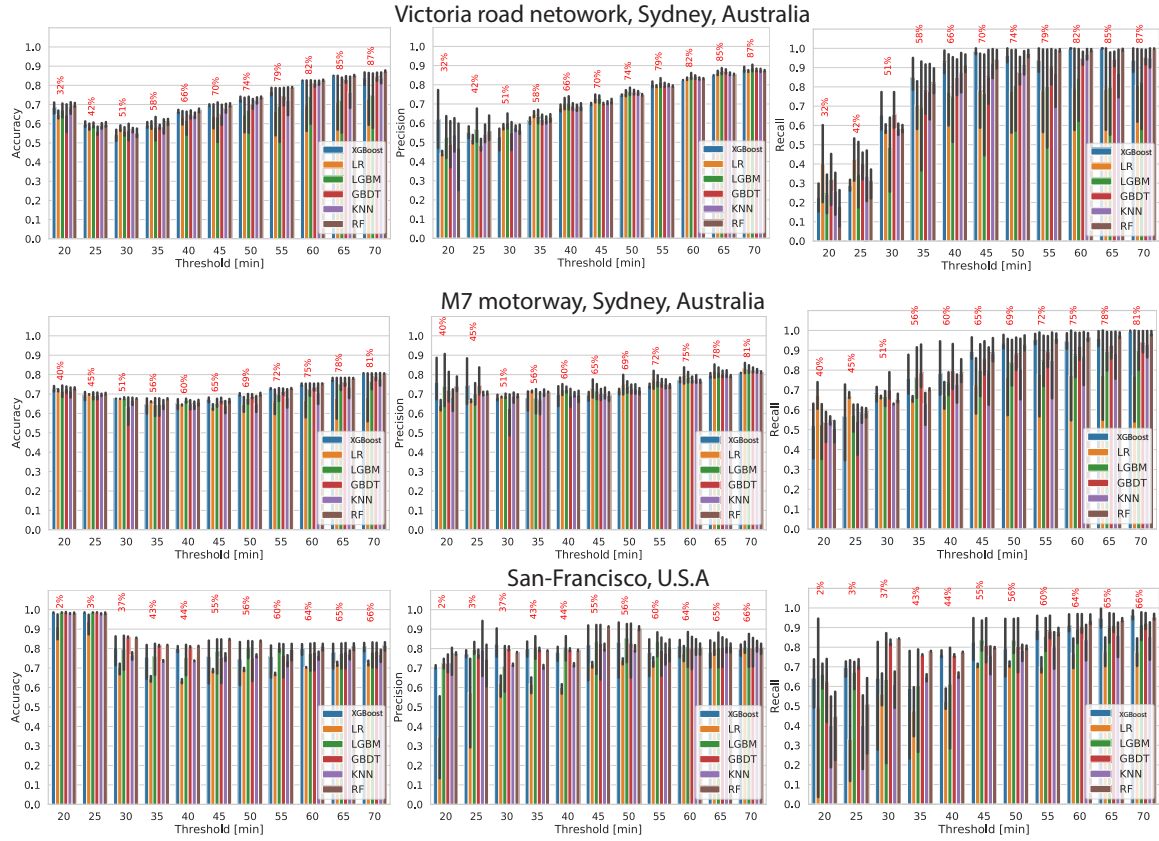


FIGURE 3.18: Binary classification performance using varying incident duration threshold

Providing additional information with regards to the computational time of various baseline ML models across the three data sets. The findings indicate the RF and kNN seem to be the slowest models to train versus LGBM and XGBoost and LR which are faster from a computational time point of view.

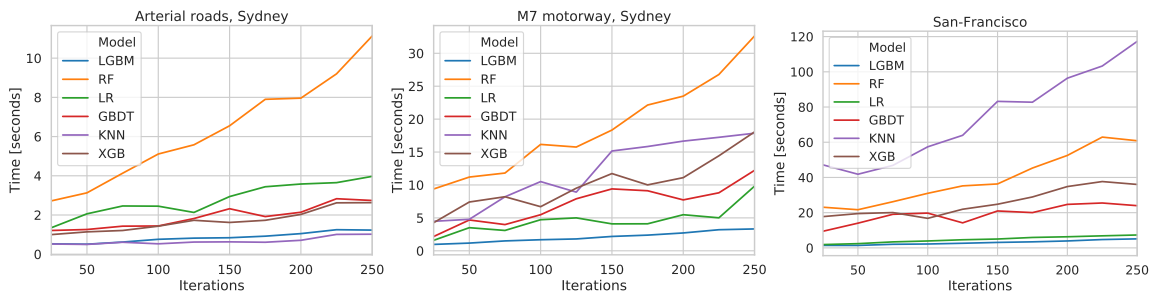


FIGURE 3.19: Performance testing of ML models across three different data sets



## **Chapter 4**

# **Traffic incident duration prediction via a deep learning framework for text description encoding**

## 4.1 INTRODUCTION

When traffic accidents occur, the majority of traffic management centres (TMCs) store a brief textual description and the GPS coordinates of the incident. There is a lot of uncertainty at the beginning of disruptions with regards to how long the traffic incident will last, and most of the time, centres do not have an overview of the length or severity of the disruptions. Therefore, it is extremely insightful for TMCs to be able to utilise the data on historical traffic flows or readily available accident description to predict or improve predictions of the incident duration. In order to improve predictions, we need more information on the factors (both readily-available and historical) which can have an effect on the incident duration prediction accuracy. This paper presents an advanced incident duration prediction framework which makes use of additional incident report variables and past incidents records, merged into a hybrid machine learning (ML) modelling approach with deep learning encoding of additional features (e.g. textual incident description and historical traffic flow in the vicinity of the section). Feature encoding is justified since the traffic incident description and traffic flow/speed measurements have a high dimensionality, which can lead to overfitting when using ML models and it may be worsened by the small size of a typical incident report data set.

This paper is organised as follows: Section 5.1 presents the challenges and reviews the related works; Section 4.2 introduces the data sources we have used as well as our traffic flow mapping algorithm for feature construction; Section 6.3 proposes our modelling framework and explains the ML models we have used, the LSTM sentiment encoder for textual incident descriptions, and the ANN encoder for traffic flow speed; Section 5.3.5 introduces the results before summarising all findings in the Conclusions section.

### 4.1.1 RELATED WORK

There are multiple research papers which use baseline incident reports from TMC with different machine learning models to predict the traffic incident duration Li et al., 2018. The use of traffic flow and incident description features is found to be rare and mostly specific - topical text modelling Das, Mohanty, and Bhattacharyya, 2019 for the task of the incident duration detection, modelling or incident impact prediction by using traffic flows Fukuda et al., 2020. And its scarcity is highlighted since it requires the involvement of additional specific models with a feature fusion approach. In other words, traffic flow data is rarely combined with textual incident description and an actual incident reports since it requires a higher system complexity.

But feature combination can be observed in some specific research studies related to the traffic incident impact prediction, which rely heavily on the historical traffic flow data with and without consideration of features that are describing the incident Fukuda et al., 2020; other works have addressed a similar approach WenTRB2018; Mihaita2019. Also, these works don't focus on the incident duration prediction.

Sometimes, researchers try to apply uniform ML approaches or specific models for all the sub-tasks. Separate RBM models were applied to different kinds of features and feature fusion representing a uniform application of ML method to different data sets Li et al., 2020a. Also, kNN and Bayesian cost-sensitive networks were combined for the task of the incident duration prediction Kuang et al., 2019a. But neither of these research studies investigated a deep dive into their model selection.

Since we have the incident description and incident severity values in our incident reports, we can utilise specific models for the task of sentiment classification. Previously, the LSTM architecture has been compared with Support Vector Machines, Artificial Neural Networks, Deep Belief Networks and Latent Dirichlet Association on the task of detection of incidents from social media data Zhang, Chen, and Zhu, 2018. LSTM was also successfully used for stock price prediction Sen and Dhar, 2018, making it applicable for modelling of traffic flow/speed time-series data. Despite its superior performance, we need to uplift and bring significant modifications to this architecture. Since we are planning to use encoded time series with machine learning methods, we need a controllable size of the feature vector to simultaneously avoid overfitting and provide enough information for ML methods. This is why we propose to use LSTM coupled with ANN, where the ANN feature vector size and the activation function are varied.

## 4.2 CASE STUDY

In this study we assume that textual incident reports as well as historical traffic flows and speed data (including the ones from the moment when an incident happened) are readily available at the moment the incident was reported and sufficient to make the prediction of its duration.

### 4.2.1 Incident description data set and baseline feature set

A Countrywide Traffic Accident Data set (CTADS) has been recently published Moosavi et al., 2019a-Moosavi et al., 2019b, which contains about 1.5 million traffic accident records across 49 states of United States of America from February 2016 to December 2020 (version 4). Each incident report contains 47 features describing the traffic accident. The majority of these traffic accidents were recorded in the state of California. The most notable features include: a) Incident Severity (valued from 1 to 4), b) Start and End Time of the incident (from which the traffic incident duration is derivable), c) The road extent affected by the accident, d) textual Incident Description, d) weather and lighting conditions. For the extended description of features please refer to the original paper describing the data set Moosavi et al., 2019b. This data set allows us to use the textual incident description and, hence, apply a sentiment analysis methodology (based on the incident severity) Alkheder, Taamneh, and Taamneh, 2017. We further refer to these features as a baseline feature set, excluding the textual incident description.

### 4.2.2 Traffic flow and speed data

To collect the data on traffic flows and speed we rely on the Caltrans Performance Measurement System (PeMS) Chen et al., 2001, which provides aggregated 5-minute precision measurements of traffic movements across California. Although there is a lot more data for the Los Angeles area (which may be considered in our future research), we decided to concentrate on the area of the city of San Francisco. We focus on 83 Vehicle Detection Stations (VDS) placed in that area, and we try to manually associate each incident occurred in that area with a VDS in their 500m proximity. VDS in PeMS may have detector failures and incomplete readings, which is common across the data set and should be taken into account. Even though the PeMS data set contains data on reported incidents, we decided to

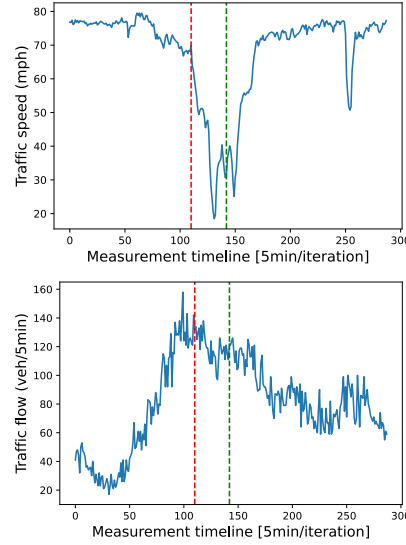


FIGURE 4.1: a) Traffic speed and b) Traffic flow plots for the VDS associated to incident A-4798 (accident on US-101 Southbound with duration of 31 5-minute iterations - actual reported incident clearance time, without considering the incident recovery time). The red line denotes the start of the accident, and the green line the end of the accident. The blue line denotes the speed evolution in the vicinity of the incident location (drops almost to 20km/h) while the flow is still running at high values due to large numbers of vehicles blocked in traffic.

use the descriptions from the Countrywise Traffic Accident Data Set since it provides a high-quality description of each incident (47 features in each incident report) extracted from Bing and MapQuest services.

In total, from 9,275 incidents in the area (extracted from CTADS) we have obtained 1,932 traffic incident reports in a 500m proximity next to VDS stations, which we were able to associate with the correct (no detector faults) and complete traffic flows and speed readings. Incident to VDS association is necessary since both are represented as points and it is not clear which incident is related to which detector since incidents on different separate roads can be in proximity of one detector (also, since we have a different representation of street names in VDS and the incident data sets). The task of VDS-to-incident assignment can be a topic for additional research, but in this paper we summarize our extracted mapping strategy as follows. We extract the following speed and flow readings from each VDS station:

1. Speed – Traffic Speed from the 24h leading to the incident occurrence.
2. Flow – Traffic Flow from the 24h leading to the incident occurrence.
3. Speed7 – Traffic Speed on the same weekday, the week before the incident.
4. Flow7 – Traffic Flow on the same weekday, the week before the incident.
5. SD – the vector difference between the traffic speed on the day of the incident and on the same weekday, the week before the incident.
6. FD – the vector difference between the traffic flow on the day of the incident and on the same weekday, the week before the incident.

Each of these feature vectors contains 288 values, which correspond to 5-minute readings throughout the day. Since each of these vectors have a high dimensionality, we decide to perform dimensionality reduction via an ANN autoencoder.



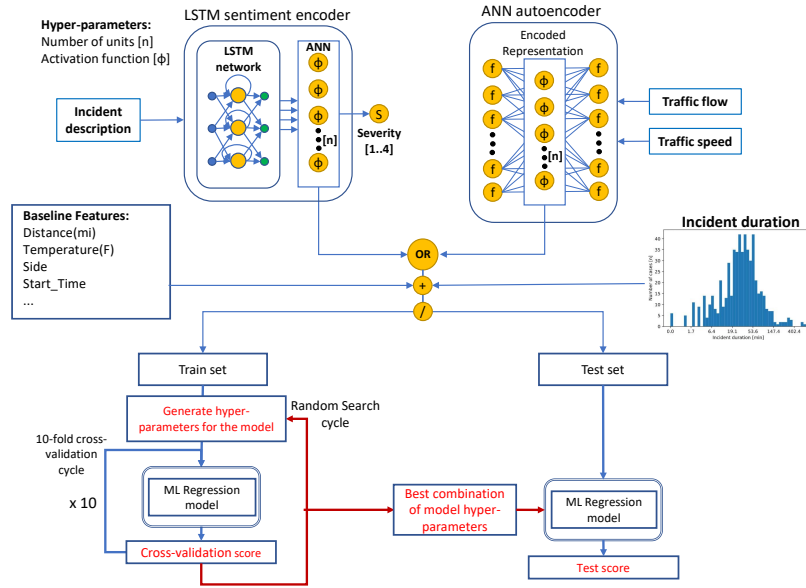


FIGURE 4.2: The structure of the proposed framework

The use of dimensionality reduction is justified since a large number of explanatory variables can cause model overfitting [sawalha2006transferability](#); [sahraei2021](#).

**Regarding the 288 input values on the day of incident:** the traffic data is taken from the time between the incident start and minus 24h before of its occurrence and not during the entire day after the incident has been lodged.

Figure 4.1 shows an example of a traffic speed drop during the incident A-4798. After we have analysed different traffic flows and speed plots we expect that the traffic speed will be the most useful single feature for the task of incident duration prediction as the traffic flow measurement seems to be not affected by the accident (as the majority of vehicles will be waiting for the congestion to clear off the road, and will still be counted as part of the traffic flow). We will also use speed measurements from the weekday, 7 days before the incident in order to obtain the complete picture between what is a regular traffic flow condition versus disrupted traffic condition on the same time and same day of the week. We make the observation that we have also conducted a detailed feature ranking and selection (via SHAP values, forward feature selection, etc.) to several incident data sets which are not presented here due to space limitations.

The point A-4798 point was selected just as an example for a traffic speed drop and its usefulness to the prediction problem; in reality, we have analysed about 100 traffic flow and speed plots before drawing the conclusions (we provide several shapshots of flow and speed reading in the supplementary material **appendix**). As an observation, by adding severity classification probabilities (from the LSTM-ANN model) to the feature vector for the task of incident duration prediction doesn't seem to be useful since we already included Severity, which is a strong feature.

Accident on I-280 Northbound at Exit 57 King St.  
 Right hand shoulder blocked due to accident on I-280 Northbound after Exits 54 54A 54B US-101.  
 Lane blocked due to accident on US-101 Presidio Pkwy Southbound at Exit 438 CA-1.  
 Accident on I-80 Westbound at Exits 1 1C / Bryant St / 8th St.  
 Second lane blocked due to accident on I-80 Eastbound at Exits 2B 2C Harrison St.  
 Lane blocked due to accident on US-101 Golden Gate Brg Southbound at Exit 439 Transit Transfer Facility.  
 Right hand shoulder blocked due to accident on I-280 Northbound at Exit 52 San Jose Ave.  
 Right hand shoulder blocked due to accident on US-101 Southbound at Exits 429B 429C Bay Shore Blvd.  
 Lane blocked on exit ramp due to accident on I-280 Northbound at Exit 55 Cesar Chavez.  
 Right hand shoulder blocked due to accident on I-280 Northbound at Ocean Ave.

TABLE 4.1: Example of the Incident Description values

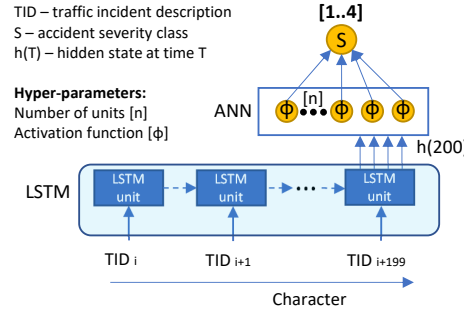


FIGURE 4.3: LSTM sentiment encoder structure.

## 4.3 Methodology

Figure 4.2 shows how we use the data to perform the incident duration prediction. We combine the baseline feature set with either the encoded textual description or the encoded traffic flow/speed values. The encoder parts of both LSTM-ANN network and the ANN autoencoder have hyper-parameters in the form of number of units and used activation functions to ensure an optimal encoding for the specific ML method. After obtaining encoded representations associated with the incident, we search for the optimal hyper-parameters for each ML regression model at each case of the encoded representation. It allows us to adapt the model parameters to work with encoded data and provide the best cross-validation results.

### 4.3.1 LSTM-ANN for the textual incident description encoding

Textual Incident Description in the CTADS data set describes type of disruption caused by the incident and/or location (Table 4.1).

To perform the encoding of the textual description of the incident we use a combination of character-level LSTM and ANN for the sentiment analysis (Figure 4.3). We use the textual incident description from all the available traffic incident reports for the San Francisco area (9,275 incident records). Firstly, we set the target variable for the LSTM classification model as the incident severity (values 1 to 4).

Secondly, we use the encoded representation of the textual description extracted from the LSTM sentiment classification model to use it as additional features for the task of incident duration prediction.

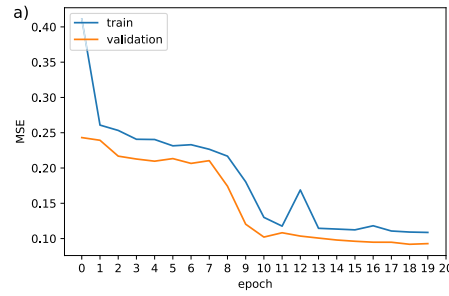


FIGURE 4.4: Example of LSTM network training results using 12 units, a ReLU activation function, 10 epochs, 80 hidden units. a) Train-validation score over 20 epochs

The incident description text is only provided at the beginning of the incident reporting timeline, and no temporal evolution is found across multiple countries for which we analysed the incident logs in our previous work **ag2022**.

Each textual description is formed into repeated strings up to 200 characters in length and each character in that string is then encoded by using one-hot encoding.

In order to showcase the importance of the textual incident description for the tasks of incident duration prediction and incident severity classification, we perform a word importance analysis using the LIME method (provided in the supplementary material **appendix**). We further train an LSTM model with 80-units hidden state vector. We use the encoding of the incident description by using different numbers of neurons and different activation functions. An example of training results for one of the variants is shown on Figure 4.4. Traffic incidents descriptions were used to predict the incident severity. The data set was split into train, validation and test sets by proportion 70:20:10. Training results show that the LSTM sentiment encoder needs at least 15 epochs to converge, so we decided to train each variant of the LSTM sentiment encoder for 15 epochs. We use Root Mean Squared Error (RMSE) as the loss function.

#### The use of MSE versus cross-entropy

MSE is a legitimate metric for the classification when the target feature is represented as an ordered variable **gaudette2009evaluation** in which MSE is preferred instead of the Cross-Entropy (CE) loss in order to reduce the model complexity and the probability of over-fitting. In our research we determined that CE required  $N \times 5$  sized matrix for the intermediate feature vector to the target value classification, while the MSE solution requires only  $N \times 1$  matrix, where  $N$  is size of the intermediate feature vector). MSE loss is also superior to CE loss for class-imbalanced datasets **kato2021mse** and our incident severity feature distribution poses an imbalanced classification problem.

#### 4.3.2 Artificial Neural Network Encoder for the traffic flow/speed encoding

As additional data sources apart from the incident baseline features, we use the general structure of Artificial Neural Network (ANNs) Autoencoder **kramer1991nonlinear** with varying number of neurons and different activation functions in the bottleneck layer to produce the encoded speed/flow data sets. Flow and speed values are normalised to the corresponding maximum observed traffic speed and flow in the data set. To improve the performance of the encoding model we use all the time series

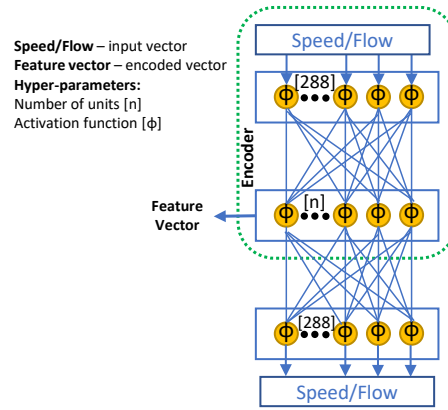


FIGURE 4.5: The structure of the ANN autoencoder

data available for each incident which could be matched to a VDS station. We combine normalised flow and speed data sets to perform the ANN model training which allows the model to grasp actual time series without focusing on speed and flow on an individual level. We do make the observation that while speed and flow could be used as raw features in any ML prediction framework, the benefit of using ANN for auto-encoding is mainly a dimensionality reduction and improved accuracy in case of extreme outliers. Last, we extract the outputs of the ANN autoencoder bottleneck layer and use them as features in the ML models shown in Fig. 4.5.

The following activation units were used in the bottleneck layers of the ANN autoencoder and the LSTM sentiment encoder: a) the Rectified Linear Unit (ReLU) **agarap2018deep** which is a piecewise linear function (output values are  $[0; +\infty]$ ) b) the Exponential Linear Unit (ELU) **trottier2017parametric**, which was developed to reduce bias shift (which leads to weight oscillations) c) the Tanh - a hyperbolic tan function which has the property of equalizing training over layers **kalman1992tanh**; its output can take values in the interval  $(-1; +1)$  d) the Sigmoid activation function which output can take values in the interval  $(0; 1)$ .

### 4.3.3 Baseline Machine Learning model selection

When all encoding has been finalised, we first use the following ML regression models as a baseline to perform the incident duration prediction:

a) gradient boosting decision trees - GBDT **Xia2017TrafficFF** which rely on training a sequence of models, where each model is added consequently to reduce the residuals of prior models; b) extreme gradient decision trees - XGBoost **chen2015xgboost** which rely on an exhaustive search of split values by enumerating over all the possible splits on all the features and contains a regularisation parameter in the objective function; c) random forests - RF **8283291** which applies a bootstrap-aggregation (bagging, which consists of training models on randomly selected subsets of data) and uses the average (or majority of votes) of multiple decision trees in order to reduce the sensitivity of a single tree model to noise in the data d) Support Vector Regression (SVR) machines **drucker1997support** which are characterized by the use of kernels and symmetrical loss function (equal penalization of high and low errors), e) Decision Trees (DT) regression models **breiman1984cart** which rely on the repetitive process of splitting and generates a set of rules which can be used for the value prediction, f) Linear Regression (for which we use standard Ordinary Least Squares optimisation) which represents the

relation between features and the target variable as a linear equation targeting to minimize the residual sum of squares between the actual and the predicted values of the target variable.

### Model performance evaluation

To evaluate the regression models on the task of the incident duration prediction we use the mean absolute percentage error and the root mean squared error defined as:

$$MAPE = \frac{1}{n} \sum_{t=1}^n \left| \frac{A_t - F_t}{A_t} \right| \quad (4.1)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (A_t - F_t)^2} \quad (4.2)$$

where  $A_t$  are the actual values and  $F_t$  - the predicted values,  $n$  - the number of samples. We do make the observation that other performance metrics have been obtained (MAE, SMAPE), but given the current page limitations, we focus on MAPE, RMSE results only.

### Hyper-parameter tuning for the proposed regression model

We use 10-fold cross-validation to overcome the over-fitting problem **geisser1993vol** and to assess the generalization performance of the ML models. In each scenario, the data set is partitioned into 10 folds. The ML regression model is trained on 9 folds to make prediction on the remaining fold. The procedure is then repeated 10 times and the accuracy results are averaged across several repetitions.

#### 4.3.4 MAPE versus RMSE comparison and their non-linear relationship

There is a non-linear relationship between MAPE and RMSE when performing regression, which can be verified by using different regression data sets. We tested this hypothesis on the Concrete Compressive Strength (CCS) Data Set from UCI Machine Learning Repository by using 1000 evaluations of random 9:1 train-test splits using Random Forest evaluated against MAPE and RMSE. **Fig. 4.6a**) presents the MAPE versus RMSE plot in which we observe that, the same MAPE result (e.g. 12%) may be attributed to multiple RMSE results (e.g. from 3.5 to 6.5). A similar situation observed for 45% of MAPE on CTADS using Random Forest (see **Fig. 4.6b**). Therefore, the occurrence of a higher RMSE error when MAPE becomes lower (as in our paper) and vice-versa is a correct result. MAPE vs RMSE compared between 100-units random vectors with 1-10 value interval using 10,000 evaluations (see **Fig. 4.6c**). As can be seen from all three sub-plots, the decrease in MAPE doesn't necessarily mean a decrease in RMSE. For our study we focused on discussing the MAPE metric, which is widely used in the literature on the topic of incident duration prediction since its intuitive meaning (e.g. a 30% MAPE means a 30% deviation of prediction from the actual incident duration) and a less inclination to high errors from outliers such as the case of RMSE. The results are part of the optimal Pareto Front [marked in orange] which showcases that our proposed method can obtain the set of optimal feature combination scenarios rather than only one winning scenario. To conclude, despite an assumption on linear dependence between the RMSE and the MAPE metrics (assumption that both metrics should be reduced in an efficient solution), both in our incident duration case and the CCS data set, we observe

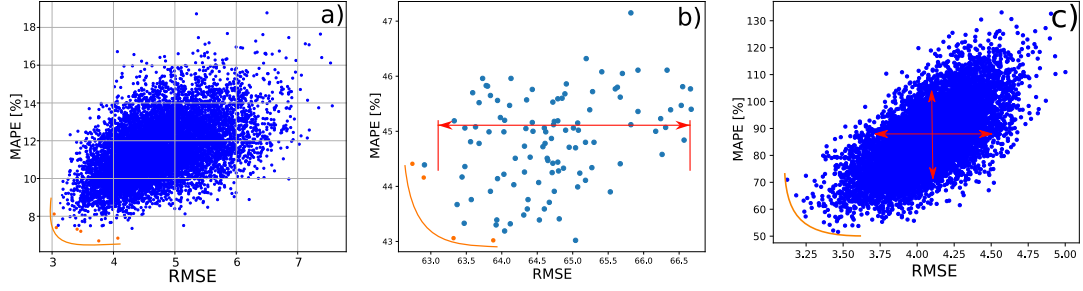


FIGURE 4.6: RMSE vs MAPE results for a) CCS data set, b) CTADS - incident duration c) Random vectors

a Pareto front of efficient solutions (no solution is sufficient in both metrics, making our results stand strong).

### 4.3.5 Comparison to other baselines

It is hard to perform a comparison between different studies on the traffic incident duration prediction since different data sets are used for research purposes Li et al., 2018. Majority of these data sets are also private and rely on different sets of features. CTADS data set appeared only recently (2019) and there is still no uniform convention on which data subset to use as a baseline, since the data set is big (1.5 million records) and heterogeneous (it includes reports from all kinds of traffic networks around United States). Indeed, in our previous work we have compared various ML-DL approaches against logs from Australia and USA, which can be used as extended results.

## 4.4 RESULTS

### 4.4.1 Best model selection

First, we try to find the three best models which show high performance of the baseline feature set consisting of traffic accident reports for which we have available traffic flow counter data. We do so by performing a cross-validation as described in 4.3.3 and a performance evaluation as detailed in 4.3.3. Figure 4.7 shows the average MAPE score for the 10-fold cross-validation obtained across several ML models such as Random Forests (RF), GBDT, XGBoost, kNN, Decision Trees (DTs), Linear Regression (LR) and Support Vector Regression (SVR). Given that the majority of traffic incident duration prediction methods published previously have reported a MAPE score below 50% Li et al., 2018, we select RandomForest, GBDT and XGBoost as the best performing models as their MAPE score falls below 46%. Next, we evaluate these three models against the baseline feature set when we apply our novel modelling approach as previously explained in sections 4.3.1-4.3.2: traffic flow, speed via ANN autoencoding and textual incident description via LSTM sentiment encoding.

There are in total 140 scenarios describing combinations of additional features [7 speed/flow/text features x 5 unit count x 4 activation functions] for each of the top three ML models. Given the restricted space allocation for this article, in Tables 4.2-4.4 we present only the top 8 best scenario results ranked against MAPE for each ML model.

Findings reveal that the encoded textual description is among the top 3 configurations for every regression model as seen from Tables 4.2-4.4. Models also demonstrate a preference for the way

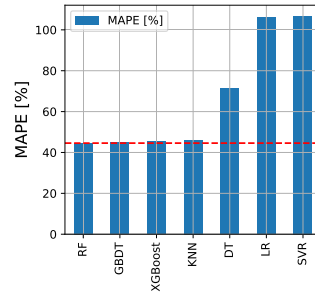


FIGURE 4.7: Regression results for baseline feature set across different ML models.

| AdditionData | units | activation | MAPE  | RMSE  |
|--------------|-------|------------|-------|-------|
| baseline     |       |            | 44.99 | 58.4  |
| LSTM-sent    | 12    | relu       | 41.89 | 65.03 |
| Flow7        | 8     | tanh       | 41.92 | 64.61 |
| LSTM-sent    | 16    | tanh       | 42.05 | 65.04 |
| Speed7       | 16    | tanh       | 42.28 | 63.97 |
| LSTM-sent    | 8     | tanh       | 42.43 | 64.13 |
| LSTM-sent    | 16    | relu       | 42.56 | 65.82 |
| Flow         | 16    | sigmoid    | 42.57 | 64.53 |
| Speed7       | 2     | sigmoid    | 42.59 | 64.76 |

TABLE 4.2: Top 8 best scenario results for GBDT-enabled framework

of encoding: 1) the Tanh activation function forms a majority in the top results for GBDT both for encoding the incident description and flow/speed features (Table 4.2), 2) the ReLU activation function forms a majority in the case of XGBoost (Table 4.4). This observation can point on a preference in the way of encoding features when using specific regression models. The best performing model among the top three finalists, when using all additional features seems to be GBDT: the best results are obtained when encoding the traffic incident description and when using the traffic flow 7 days before the incident with 12 units and the ReLU activation function [ $MAPE = 41.89\%$ , Table 4.2] (therefore including the information on the regular traffic flow profile on the same weekday, together with the incident report proves important for the task of incident duration prediction).

Other models show a higher MAPE or RMSE results for the incident duration prediction (see RF enabled results in Table 4.3 with lowest  $MAPE = 43.2\%$  for a combination of baseline, regular traffic flow, 4 layer units and a tanh activation function); similar findings appear for XGBoost-enabled results

| AdditionData | units | activation | MAPE  | RMSE  |
|--------------|-------|------------|-------|-------|
| baseline     |       |            | 44.58 | 57.6  |
| Flow         | 4     | tanh       | 43.02 | 63.88 |
| LSTM-sent    | 4     | elu        | 43.02 | 65.04 |
| LSTM-sent    | 12    | relu       | 43.06 | 63.32 |
| Flow7        | 16    | sigmoid    | 43.19 | 64.04 |
| Flow         | 16    | elu        | 43.30 | 63.92 |
| FD           | 4     | elu        | 43.32 | 64.12 |
| Flow7        | 16    | tanh       | 43.33 | 63.47 |
| FD           | 4     | sigmoid    | 43.39 | 64.53 |

TABLE 4.3: Top 8 best results for RF



| AdditionData | units | activation | MAPE  | RMSE  |
|--------------|-------|------------|-------|-------|
| baseline     |       |            | 45.44 | 63.41 |
| Flow         | 8     | relu       | 43.44 | 69.93 |
| LSTM-sent    | 4     | tanh       | 43.58 | 71.03 |
| Speed7       | 16    | tanh       | 43.63 | 71.62 |
| SD           | 4     | relu       | 43.73 | 70.58 |
| Speed7       | 16    | relu       | 43.80 | 71.92 |
| LSTM-sent    | 16    | elu        | 43.81 | 70.45 |
| LSTM-sent    | 8     | relu       | 43.82 | 72.19 |
| Flow7        | 2     | relu       | 43.85 | 72.94 |

TABLE 4.4: Top 8 best results for XGBoost

in Table 4.4 with the lowest  $MAPE = 43.44\%$ , when using again the regular flow features, 8 layer units and ReLU activation function. This experiment shows that an accurate incident duration prediction immediately after the event has occurred is possible, leveraging the incident description and the measured traffic flow on the day of accident, which may prove very useful for TMCs to incorporate directly in their incident management platforms. Lower MAPE does not necessarily mean lower RMSE as seen from the baseline and additional data scenarios, but the LSTM sentiment encoding seems to be the approach that obtains the best RMSE score (64.13) when combined indeed with other variations of the activation function and number of hidden units (as shown in Table 4.2).

#### 4.4.2 Parallel coordinates for scenario setup

To supplement the findings, we also provide a parallel categories representation of all the 140 scenarios for the GBDT model in Figure 4.8, which highlights the best combination of activation functions that seem to be working best alongside the character-level LSTM sentiment encoder of traffic flow incident textual description and speed information - mostly from previous daily speed profiling using historical data. The worst results seem to be the ones obtained when using only the speed or flow difference vector alongside the baseline incident features.

Encoding using Sigmoid and Tanh activation units on average performs best, probably because of the limited value range: Tanh and Sigmoid allow encoded representations to take values in ranges  $[-1; +1]$  and  $(0; 1)$  correspondingly, ReLU and ELU can take unlimited positive values. These results indicate which value ranges work best for encoded representation.

Comparison of MSE and CE implementations of LSTM severity classification metric for the purpose of obtaining feature vector representation of Incident Description (see Fig. 4.8) shows that a sentiment classifier with Cross-entropy (lstmSentCE) as a target metric with one-hot encoded severity values is more efficient (left column attributed to lstmSentCE shows more blue rows associated with low metric values than lstmSentMSE - sentiment encoder which predicts severity as a single value). Comparison between the number of units shows preference for 4 units since the presence of the lowest error and absence of the highest error rows. Among the activation units, the Sigmoid is the best performer showing more low error results than other units. This scenario is to show how feature vector representing incident description may be efficiently encoded to be used with conventional GBDT machine learning method: using cross-entropy for the severity classification, using 4 units and the Sigmoid as activation function.



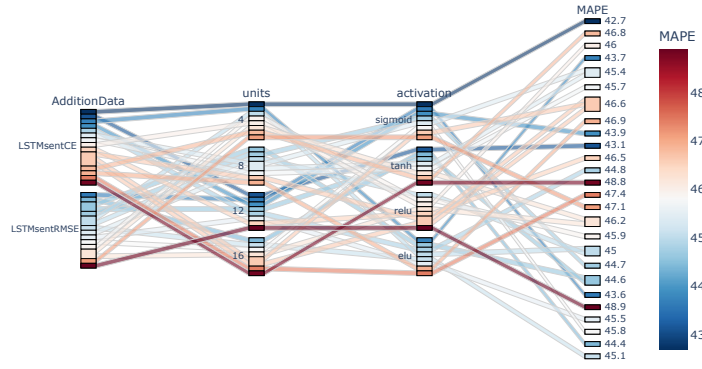


FIGURE 4.8: Parallel categories representation for all regression scenarios with GBDT.

## 4.5 Conclusion

In this paper, we have proposed a novel framework to predict the incident duration using an integration of machine learning with traffic flow and description features encoded via several Deep Learning methods. This approach demonstrates the stable and noticeable improvement across all the performing models. The results give evidence to the importance of using specific deep-learning encoding approaches for all regression models which provide a further boost-up in the model performance from past historical traffic information and the textual incident description. Efficiently encoding incident-related features for the task of incident duration prediction is the first step to model the traffic incident impact on the traffic flow. Further work is currently being focused on exploring the spatial and the temporal dynamic prediction of the incident impact via graph-based modelling approaches. The research has the following limitations: a) we used as study area only San Francisco, but there is a data availability on traffic accidents and traffic flow for the area of California, b) traffic speed and flow were taken into account only before the incident; by collecting traffic count data for longer periods it possible to build traffic speed/flow profiles which may provide more accurate predictions. The societal impact of the research is as follows: the data availability of the predicted incident duration can improve for TMC incident and traffic management (e.g. TMC can announce when an incident is expected to dissipate, how many resources to allocate, etc), which in turn will reduce the time spent by people in the traffic congestion caused by the incident. The code for the paper can be found: [https://github.com/Future-Mobility-Lab/TIDP\\_2022](https://github.com/Future-Mobility-Lab/TIDP_2022).

### 4.5.1 Word importance for severity classification

To estimate word importance in the Incident Description feature, word count matrix has been transformed to a normalized TF-IDF representation (term frequency–inverse document frequency) **TFIDF**. N-gram value range is (1,2). Then linear dimensionality reduction has been performed using truncated singular value decomposition to 50 componenets for 7 iterations. Then we used GBDT classification model to fit incident severity and three quantiled groups (ratio 33%:33%:33% to represent equally sized groups with duration intervals 0-29min, 30-71min and 72-2750min) of the incident duration. Classifier predictions were then analyzed for feature importance using LIME method **LIME**, where every feature represents 1 word or 2 word combination presence in the incident description. One or more combinations of word in the description can contribute to the incident being classified into one of

| Severity group=0    |                   | Severity group=1    |             | Severity group=2    |               | Severity group=3    |                   |
|---------------------|-------------------|---------------------|-------------|---------------------|---------------|---------------------|-------------------|
| Weight <sup>2</sup> | Feature           | Weight <sup>2</sup> | Feature     | Weight <sup>2</sup> | Feature       | Weight <sup>2</sup> | Feature           |
| +0.644              | lanes blocked     | +1.559              | chavez      | +2.805              | 280           | +0.982              | lanes blocked     |
| +0.345              | two lanes blocked | +1.190              | cesar       | +1.909              | 280           | +0.427              | two lanes blocked |
| +0.231              | due to accident   | +0.973              | <BIAS>      | +0.828              | northbound    | +0.365              | due to accident   |
| +0.034              | lanes blocked     | +0.894              | st          | +0.740              | to accident   | +0.174              | on i              |
| +0.007              | due to accident   | +0.475              | i           | +0.736              | accident      | +0.110              | due to            |
| -0.407              | lanes blocked     | +0.467              | to          | +0.721              | i 280         | +0.008              | to accident       |
| -0.620              | blocked           | +0.465              | on          | +0.697              | chavez st     | -0.076              | cesar chavez      |
| -0.689              | i                 | +0.351              | at          | +0.677              | accident on   | -0.127              | lanes blocked     |
| -0.704              | <BIAS>            | +0.309              | northbound  | +0.448              | two lanes     | -0.546              | blocked           |
| -0.748              | to                | +0.307              | cesar       | +0.375              | lanes blocked | -0.621              | i                 |
| -0.760              | st                | +0.289              | chavez      | +0.336              | at cesar      | -0.666              | on                |
| -0.769              | accident          | +0.153              | two         | +0.194              | due to        | -0.672              | <BIAS>            |
| -0.793              | two               | -0.031              | due         | +0.187              | blocked       | -0.678              | to                |
| -0.797              | due on            | -0.101              | blocked     | +0.138              | lanes         | -0.710              | at                |
| -0.818              | at                | -0.125              | due to      | +0.070              | northbound    | -0.762              | two               |
| -0.869              | northbound        | -0.310              | at cesar    | +0.022              | at            | -0.773              | due               |
| -0.924              | 280               | -0.330              | two lanes   | +0.022              | on i          | -0.918              | accident          |
| -0.979              | chavez            | -0.372              | lanes       | -0.160              | northbound    | -0.953              | st                |
| -1.149              | cesar             | -0.466              | lanes       | -0.208              | due           | -0.961              | cesar             |
|                     |                   | -0.647              | blocked     | -0.354              | cesar chavez  | -0.997              | chavez            |
|                     |                   | -0.684              | accident on | -0.358              | at            | -1.116              | 280               |
|                     |                   | -0.692              | i 280       | -0.369              | two           | -1.140              | northbound        |
|                     |                   | -0.711              | chavez st   | -0.498              | on            |                     |                   |
|                     |                   | -0.728              | accident    | -0.509              | to            |                     |                   |
|                     |                   | -0.913              | blocked     | -0.534              | i             |                     |                   |
|                     |                   | -1.993              | 280         | -0.891              | st            |                     |                   |
|                     |                   | -2.728              | northbound  | -0.994              | <BIAS>        |                     |                   |
|                     |                   |                     | 280         | -1.110              | cesar         |                     |                   |
|                     |                   |                     |             | -1.479              | chavez        |                     |                   |

FIGURE 4.9:

| duration group=0    |                | duration group=1    |                | duration group=2    |                |
|---------------------|----------------|---------------------|----------------|---------------------|----------------|
| Weight <sup>2</sup> | Feature        | Weight <sup>2</sup> | Feature        | Weight <sup>2</sup> | Feature        |
| +1.307              | lanes blocked  | +0.548              | 280            | +0.389              | chavez st      |
| +0.653              | two lanes      | +0.444              | northbound     | +0.256              | 280 northbound |
| +0.461              | blocked due    | +0.357              | blocked        | +0.149              | blocked due    |
| +0.422              | lanes          | +0.218              | chavez         | +0.132              | northbound at  |
| +0.326              | to accident    | +0.214              | st             | +0.092              | at cesar       |
| +0.324              | on i           | +0.213              | accident       | +0.075              | cesar chavez   |
| +0.255              | at cesar       | +0.182              | cesar chavez   | +0.068              | to accident    |
| +0.230              | due to         | +0.095              | two lanes      | +0.062              | cesar          |
| +0.216              | northbound at  | +0.091              | cesar          | +0.017              | to             |
| +0.211              | chavez st      | +0.050              | due to         | -0.036              | <BIAS>         |
| +0.177              | accident on    | +0.039              | i 280          | -0.057              | lanes blocked  |
| +0.026              | i 280          | +0.034              | lanes          | -0.080              | due            |
| -0.123              | st             | +0.029              | 280 northbound | -0.088              | at             |
| -0.153              | cesar chavez   | -0.013              | on             | -0.133              | two lanes      |
| -0.232              | blocked        | -0.030              | <BIAS>         | -0.232              | accident       |
| -0.232              | 280 northbound | -0.037              | two            | -0.264              | chavez         |
| -0.275              | i              | -0.069              | to accident    | -0.383              | st             |
| -0.290              | at             | -0.072              | northbound at  | -0.502              | northbound     |
| -0.348              | on             | -0.077              | i              | -0.580              | lanes          |
| -0.405              | chavez         | -0.129              | blocked due    | -0.594              | 280            |
| -0.437              | northbound     | -0.204              | chavez st      | -0.633              | blocked        |
| -0.439              | 280            | -0.655              | lanes blocked  |                     |                |
| -0.440              | to             |                     |                |                     |                |
| -0.449              | due            |                     |                |                     |                |
| -0.485              | accident       |                     |                |                     |                |
| -0.544              | two            |                     |                |                     |                |
| -0.724              | cesar          |                     |                |                     |                |
| -0.918              | <BIAS>         |                     |                |                     |                |

FIGURE 4.10: Word importance estimation using LIME method for incident duration groups

severity groups (Fig. 4.9) - presence of "lanes blocked" and "two lanes blocked" has the highest contribution to the incident being classified into highest (3) or lowest (0) severity group. Severity 1 or 2 is more related to the actual location, which represented as word describing Cesar Chavez St and I-280 Interstate Highway. High positive and opposite high negative contribution of words towards severity group observed for severity groups 1 and 2, where "280" and "chavez" have high opposite contributions, making this groups easily separable. When we perform classification towards equally sized incident duration groups, "lanes blocked" has the highest positive contribution of the incident to be classified into low duration group. If accident happens on Cesar Chavez St, it can be easily classified into low duration group signifying importance of location for the task of incident duration prediction. High negative contribution of "lanes blocked" observed for duration group 1 with the highest contribution of "280" word meaning that incident appears on I-280 Interstate Highway.

#### **4.5.2 Traffic flow and traffic speed on the day of the incident**

The following plots represent recorded traffic speed and flow on the day of the incident and week before in 500m proximity of the incident along the road (see Fig. 4.11 and 4.12). Reports in CTADS data set indicate that the highest impact of traffic incident is attributed to significant decrease in traffic speed, while traffic flow stays the least affected by disruption.

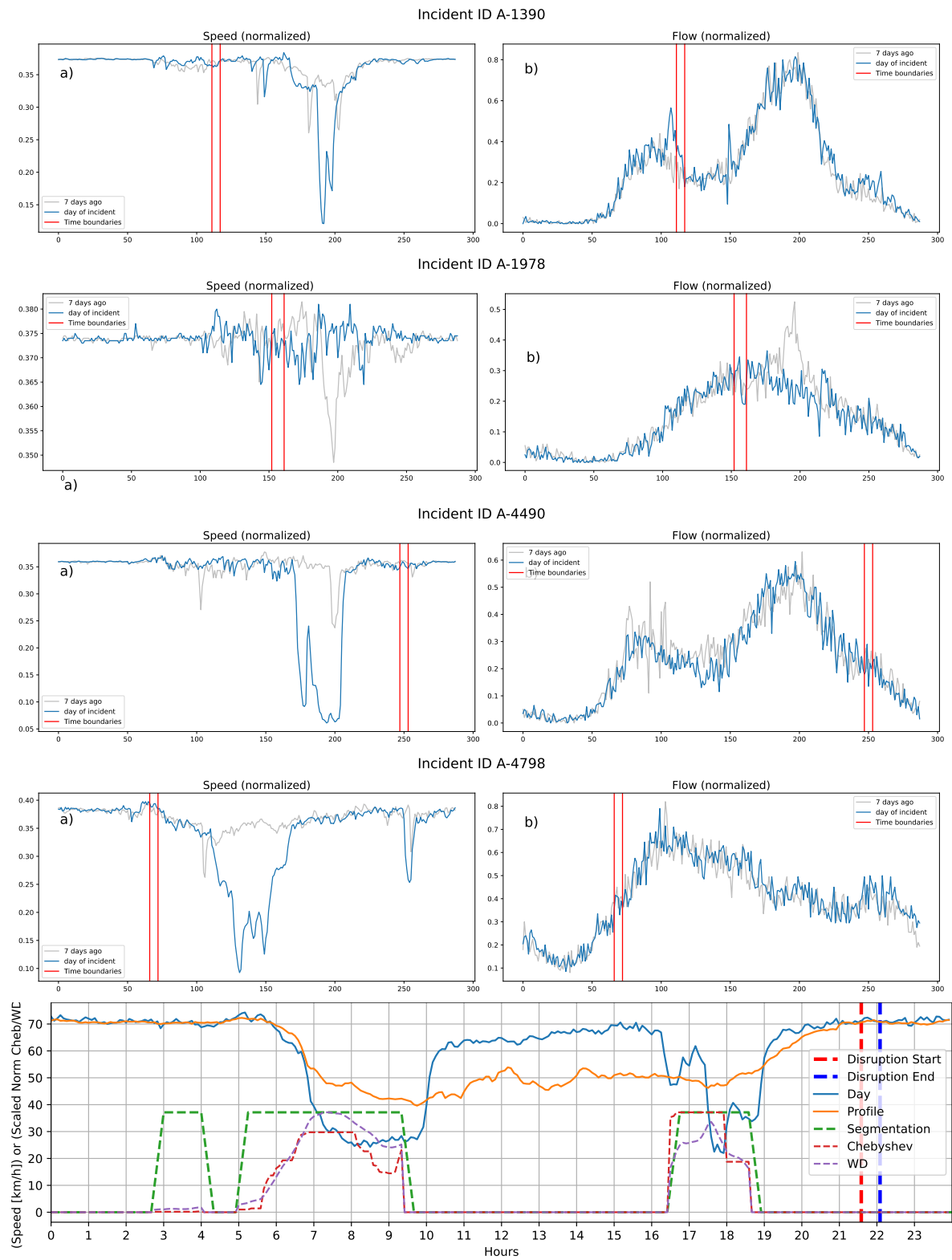


FIGURE 4.11: Traffic speed and flow during the day of the incident. Part #1

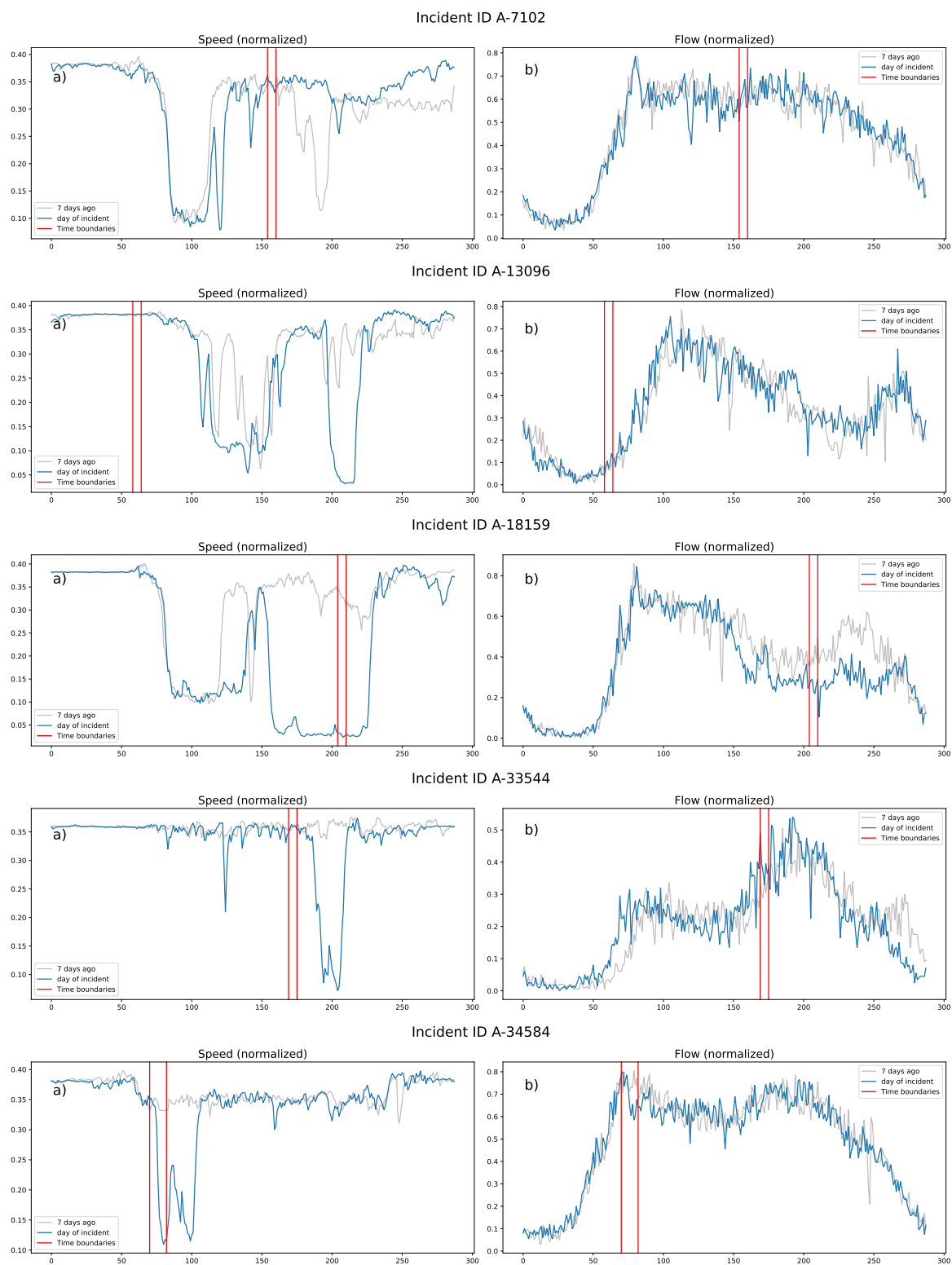


FIGURE 4.12: Traffic speed and flow during the day of the incident. Part #2



## **Chapter 5**

# **Spatial-Temporal Traffic Accident Risk Forecasting using Contextual Vision Transformers with Static Map Generation and Coarse-Fine-Coarse Transformers**



FIGURE 5.1: City grid representation for our study.

## 5.1 Introduction

Traffic accidents represent a major concern for cities around the world due to a significant economical and health impact to their populations. The number of vehicles has been substantially increasing during the past decades, especially in developing countries, which lead to an increase in the number of traffic accidents Organization, 2015. The National Highway Traffic Safety Administration (NHTSA) reports more than 5 million traffic accidents happening in the United States each year Administration, 2013. The World Health Organization also reported 1.35 million fatalities happening worldwide which resulted from traffic accidents in 2016 world2018global.

In the past years, traffic accident research has seen an increased use of computational methods. Different problems were addressed, including: 1) traffic accident duration prediction methods, Li et al., 2018 2) accident detection Parsa et al., 2019, 3) estimation of severity, and more recently, a development of spatial-temporal modelling methods have allowed to perform accident risk prediction using high-dimensional spatial, semantic and temporal data sets Wang et al., 2021b. The use of such methods has enhanced the automated analysis of traffic data together with the increasing number of publicly available data sets. Traffic accident risk prediction allows to: 1) detect high-risk areas within a traffic network, which may facilitate the decision-making inside traffic management authorities, 2) to allocate resources and assess the road design to reduce the number of accidents in the future, 3) to predict timely high-risk situations on the road and 4) to allow an implementation of risk-reducing traffic management strategies.

In the literature, the traffic accident risk forecasting problem is commonly formulated as a time-series forecasting task, where given past historical traffic accidents data for a certain city/region, along with an optional contextual information about those accidents, the objective is to forecast/predict the future traffic accident risk for that city/region. Since the nature of the traffic accident risk problem implicitly involves two types of modelling, e.g. the spatial approach (working on the affected geographic region) and the temporal approach (applied over a period of time), thus this problem is often tackled using at least two different types of model architectures.

One of the first works on traffic accident risk prediction using Deep Learning has been performed with human mobility data using a Stack Denoise Autoencoder (SDAE) on the Japan traffic network Chen, Chen, and Hsieh, 2016, but traffic flow and time-related matters (including periodicity) were



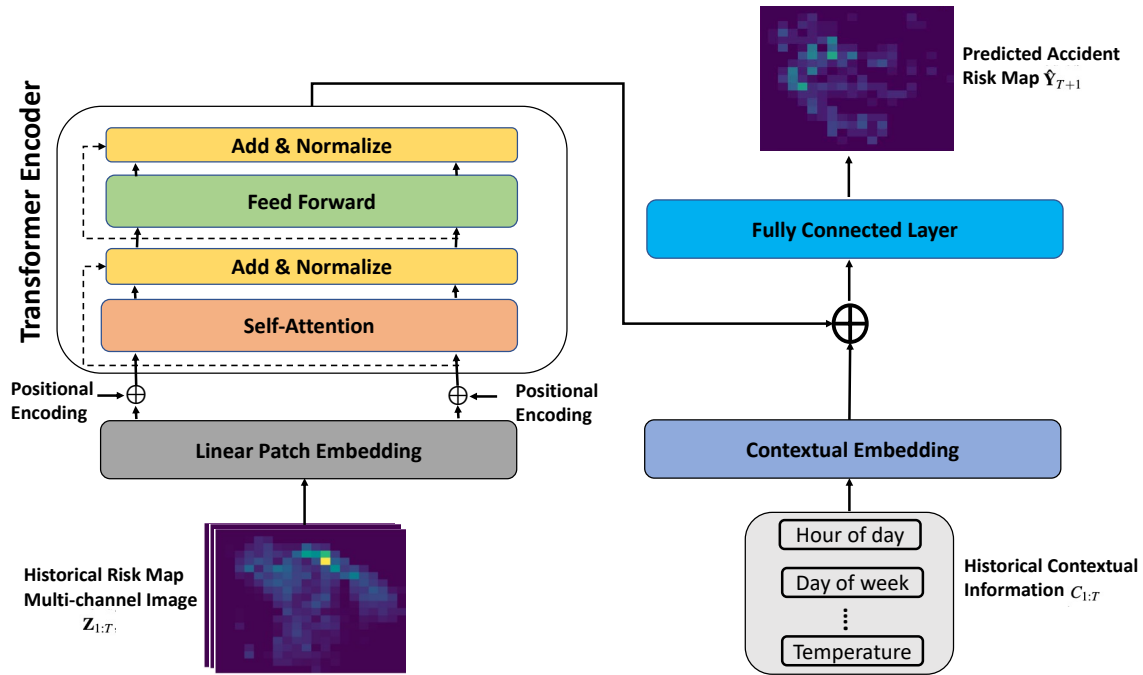


FIGURE 5.2: The building blocks of our proposed C-ViT model.

not considered. Another research Ren et al., 2017 relied on the LSTM network to improve the risk prediction in comparison to SAE by considering in addition the air quality, traffic flow and the weather data, represented as short-term and periodic components. Zhou et al., 2020a proposed also a Coarse and Fine grained prediction on the target accident risk map. RiskOracle Zhou et al., 2020b relied on Graph-Convolution network, utilizing hierarchical coarse-to-fine modelling and proposing minute-level predictions in comparison to day-level Yuan et al., 2018 and hour-level Chen, Chen, and Hsieh, 2016. In Yuan et al., 2018 authors have constructed over the ConvLSTM by highlighting the spatial heterogeneity problem and proposing an ensemble of region-specific ConvLSTM models (Hetero-ConvLSTM); they considered weather, the environment and the road condition in Iowa, US for over 8 years of observations, but POIs were not considered. Semantic features, coarse and fine grained risk maps were considered in Wang et al., 2021a, where also Graph-convolution neural networks and attention-based LSTMs were used. A more recent work in Wang et al., 2021b represents the State-of-Art (SoTA) in the field of accident risk prediction, where the authors propose a weighted loss function to address the zero-inflated issue (increase in the number of zero-risk grid cells due to the increase in the granularity of predictions) and making ensemble of models by processing semantic and geo features.

Current paper relies on the use of original Visual Transformer dosovitskiy2020image; wu2020visual, which has been widely applied to various tasks in areas of Computer Vision. Transformer models have multiple variations including Convolution Neural Network Enhanced Transformer, Hierarchical Transformer, Transformers with Local Attention, Deep Transformer liu2021survey.

So far, risk accident prediction relied mostly upon graph-based methods and spatial-temporal modelling. While this approach worked for limited case study applications, we highly believe that in order to scale it up, this approach can benefit from using visual analysis techniques. Thus, in this work we

are re-formulating the problem of traffic accident risk forecasting and we are proposing a novel approach inspired by one of the recent best performing deep learning based architectures for computer vision tasks, the vision transformers **dosovitskiy2020image**. In our proposed model we jointly model and take into account the spatio-temporal nature of the traffic accident risk forecasting problem as well as the influence of contextual information on it using a single unified end-to-end model.

An earlier version of this paper was presented at the IEEE ITSC 2022 Conference and was published in its Proceedings **saleh2022traffic**. The current paper provides a significant expansion (including two new added sections) of our previous work to further improve accident risk prediction results. The current paper expansion includes results on Coarse-Fine-Coarse and Static Map Visual transformer architectures.

In Section 5.2, a detailed description about the proposed methodology will be presented. Then, in Section 5.3, we will introduce the datasets we utilised for training and evaluating the performance of our approach, the experiments setup and the baseline approaches from the literature we compared our approach against. Next, we introduce the Coarse-Fine-Coarse Transformer architecture to improve accident risk prediction results in Section 5.4. Then, we propose an incorporation of static maps into ViT architecture in Section 5.5. Finally, in Section 5.6, we conclude our paper.

The code for the paper can be found by the following link: <https://github.com/Future-Mobility-Lab/ViT-traffic-accident-risk>

## 5.2 Methodology

In this section, we will first start with definitions and the problem formulation for the traffic accident risk forecasting task. Then, we will present and discuss the details of our proposed contextual vision transformer (C-ViT) model (as shown in Fig. 5.2).

### 5.2.1 Definitions

**Grid Representation:** Given a city area bounded by certain latitude and longitude coordinates, we partition it into a grid form with  $I$  rows  $\times$   $J$  columns (as shown in Fig. 5.1), where each cell share the same size.

**Traffic Accident Risk:** At any given time  $t$ , the traffic accident risk  $Y_t^i$  for a grid cell  $i$  is defined by the summation of the different types of traffic accidents occurred at that grid cell. Similar to Wang et al., 2021b, we have three types of traffic accidents and each one has a corresponding value, namely a minor accident has a value of 1, an injured accident a value of 2 and a fatal accident has a value of 3. For instance, if a grid cell incurred three fatal accidents and two minor accidents, the traffic accident risk for it then would be 11.

### 5.2.2 Problem Formulation

In our formulation of the traffic accident prediction problem, we re-cast it as an image regression task instead of the traditional formulation as a time-series prediction task. This new formulation enables us to natively model the spatio-temporal nature of the traffic accident prediction problem in an end-to-end fashion without the need to have a combination of more than one architecture to address it. To that

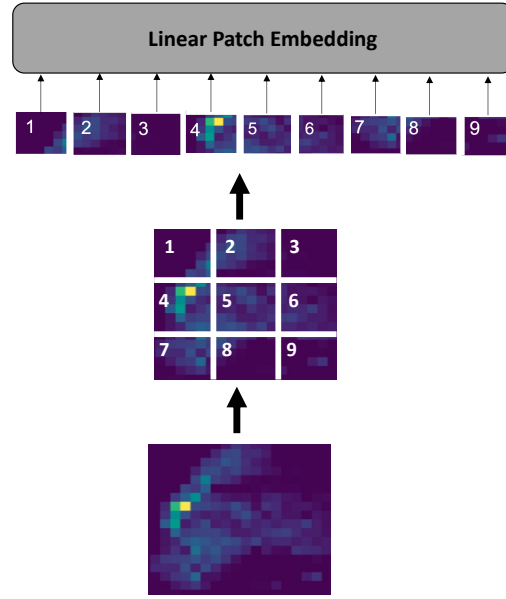


FIGURE 5.3: Description of the first stage of the historical risk Map encoding. Given a unified single image  $X$ , it is then divided into equally-sized image patches that are passed individually to the linear patch embedding layer.

end, given historical observations in the form of a traffic accident risk map  $\mathbf{Z}_{1:T}$ , where  $\mathbf{Z} \in \mathbb{R}^{I \times J}$  over time period  $[1 : T]$ , we represent these observations as an image  $X$  with a resolution of  $I \times J$  and its number of channels to be  $T$ . Then we feed it to our proposed C-ViT model that fuse it together with the historical contextual information  $C_{1:T}$  to predict/regress the future accident risk map in the next hour  $\hat{\mathbf{Y}}_{T+1}$ , where  $\mathbf{Y} \in \mathbb{R}^{I \times J}$ .

### 5.2.3 Contextual Vision Transformer (C-ViT) Model

Given the aforementioned formulation, we compile the traffic accident risk maps  $\mathbf{Z}_{1:T}$  as a unified single image with size  $T \times I \times J$ , where  $T$  is the number of channels,  $I$  is the image's height and  $J$  is the image's width, which we pass as an input to our proposed novel C-ViT model. Our C-ViT model's architecture is inspired by the recently introduced vision transformer network [dosovitskiy2020image](#) that has been achieving competitive results to the convolutional neural network (ConvNet) architecture for image classification tasks [wu2020visual](#); [dosovitskiy2020image](#). The main building blocks of our C-ViT model are three components, namely the historical traffic accident risk map encoding stage, the historical contextual information encoding stage and the transformer encoder stage. In the following we will analyse deeper each component.

**Historical Risk Map Encoding:** Given the historical risk maps as a unified single image  $X$  with size  $T \times I \times J$ , we first encode it into a representation that could be easily digested and learned using our transformer encoder. As it was shown in [vaswani2017attention](#), transformer encoders can work better with input data as a sequence of tokens. Thus, we divide the unified single image into a sequence of equally sub-images  $X_p$  which we refer to it as an image patch sequence. We can think of the image patches as a sub-spatial regions of a number of cells within the city's grid representation that we defined in Section 5.2.2. The rationale behind this patching process is derived by the assumption that

grid cells that are spatially closer to each others will have some geographical and spatial correlations that could potentially be exploited by our model for conducting a better traffic accident risk forecasting.

Here  $X_p$  has a size of  $N \times T \times P \times P$ , where  $P$  is the height/width of the image patch and  $N$  is the total number of sequences of image patches, which is defined by  $N = IJ/P^2$ . The operation of dividing the unified single image into a sequence of image patches  $X_p$  can be shown in Fig. 5.3. The image patches sequence are then individually passed through a linear embedding layer which is essentially a learn-able linear projection operation in order to get a sequence of trainable flattened image patches of size  $D$ , which we refer to as patch embeddings. Additionally, similar to **dosovitskiy2020image**, we have an extra learnable embedding token appended before the sequence of patch embeddings to be passed to the transformer encoder and we refer to this embedding as a “regression token”. The regression token embedding acts as an image representation which its output is transformed inside the transformer encoder into the predicted accident risk map  $\hat{Y}_{T+1}$ .

Since the transformer encoder does not have the notion of order in its input sequence tokens, an additional position embeddings are added to each patch embedding. There are a number of pathways to define position embedding, and in our current model we follow the formulation introduced in **vaswani2017attention**. In this formulation, the position encoding  $PE$  vector is defined by using a wide spectrum of frequencies of sine/cosine functions as follows:

$$\begin{aligned} PE_{(a,2k)} &= \sin(a/10000^{2k/D}) \\ PE_{(a,2k+1)} &= \cos(a/10000^{2k/D}) \end{aligned} \quad (5.1)$$

where  $a$  represents the position, and  $k$  is the dimension. From the above formulation, once can conclude that for each dimension  $k$  of  $PE$  vector, it has a corresponding sinusoid that spans a frequency range from  $2\pi$  to  $10000 \cdot 2\pi$ . In other words, this will allow the model to be mindful of the order in the sequential patch embedding by using unique relative positions. The dimension of the  $PE$  vector is similar to the linear patch embedding layer’s dimension which is  $D$ .

**Historical Contextual Information Encoding** As discussed in Section 5.2.2, besides the historical accident risk maps, our C-ViT model takes into account also the historical contextual information  $C_{1:T}$  for the city grid representation. In our model and similar to Wang et al., 2021b, we took into account the following contextual features: 1) the time period of the day, 2) the day of the week, 3) whether the day is a holiday or not, 4) the weather condition (clear, cloud,..etc), 5) the weather temperature, and 6) traffic condition (inflow and outflow). Given those contextual features, we encode them via a learnable linear embedding layer of dimension  $D$ , whose output is fused together with the output from the transformer encoder via a concatenation operation.

**Transformer Encoder** The main building block of our transformer encoder is the multi-head self-attention module **vaswani2017attention**. In total we have six layers inside our transformer encoder. Internally, each layer is composed of a both self-attention head and feed-forward fully connected sub-layers. Additionally, each sub-layer is followed by two residual connections and a normalisation operation. The multi-head self-attention, or the multi-scaled dot-product attention, works based on the mapping between the so-called ‘query’ vectors and the pair (key, value) vectors. The dimension of the query and key vectors is  $d_k$ , where the values vector dimension is  $d_v$ . The attention operation itself is

| Dataset | Attributes         | Range/Count              |
|---------|--------------------|--------------------------|
| NYC     | Reporting Duration | 1 Jan 2013 - 31 Dec 2013 |
|         | Accidents          | 147K                     |
|         | Taxi Trips         | 173,179K                 |
|         | POIs               | 15,625                   |
|         | Weathers           | 8,760                    |
|         | Road Network       | 103K                     |
| Chicago | Reporting Duration | 1 Feb 2016 - 30 Sep 2016 |
|         | Accidents          | 44K                      |
|         | Taxi Trips         | 1,744K                   |
|         | Weathers           | 5,832                    |
|         | Road Network       | 56K                      |

TABLE 5.1: Datasets Statistics

computed by taking the dot-product between the query and the key vectors divided by the square root of  $d_k$  before finally passing them to the softmax function to get their weights by their values. Since the scaled dot-product attention operation is done multiple times, the queries, keys and values vectors are extended into matrices  $Q, K, V$  respectively. The following formula is the description of how the scaled dot-product attention operation is calculated:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (5.2)$$

## 5.3 Experiments and Results

In this section, we first present the datasets we utilised for training and evaluating the performance of our proposed approach. Then, we provide the details of the setup for our experiments, the evaluation metrics and the compared baseline approaches from the literature. Finally, the quantitative and qualitative results of our proposed approach on real-life datasets are evaluated and discussed.

### 5.3.1 Datasets

In this study we use two publicly available real datasets for the traffic accident risk forecasting problem, namely NYC<sup>1</sup> and Chicago<sup>2</sup>. As it can be seen from Table 5.1, both datasets have historical traffic accidents and historical taxi trips. The historical traffic accident data contains: time, date, location (latitude and longitude), the number of casualties, the weather condition (clear, cloudy, rainy, snowy or mist), the temperature and the road segment data (i.e. road length, width and type). The NYC dataset has an additional Point of interest (POI) data regarding locations (i.e. residence, school, culture facility, recreation, social service, transportation and commercial). The historical taxi trips include the location and times of pick-up and drop-offs and this data is used to calculate the inflow/outflow of the traffic condition in each area.

<sup>1</sup><https://opendata.cityofnewyork.us/>

<sup>2</sup><https://data.cityofchicago.org/>

TABLE 5.2: Performance evaluation of our C-ViT model against a number of baseline approaches from the literature over the NYC and Chicago datasets.

| Dataset                      | NYC           |               |               | Chicago       |               |               |
|------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Model                        | RMSE ↓        | Recall ↑      | MAP ↑         | RMSE ↓        | Recall ↑      | MAP ↑         |
| RNN-GRU chung2014empirical   | 8.3375        | 28.09%        | 0.1228        | 12.6482       | 17.83%        | 0.0664        |
| SDCAE chen2018sdcae          | 7.9774        | 30.81%        | 0.1594        | 11.3382       | 18.78%        | 0.0753        |
| H-ConvLSTM Yuan et al., 2018 | 7.9731        | 30.42%        | 0.1454        | 11.3033       | 18.43%        | 0.0716        |
| GCN wu2019graph              | 7.7358        | 31.78%        | 0.1623        | 11.0835       | 18.95%        | 0.0805        |
| GSNet Wang et al., 2021b     | 7.6151        | 33.16%        | 0.1787        | 11.3726       | 19.92%        | 0.0822        |
| C-ViT (ours)                 | <b>7.0053</b> | <b>33.86%</b> | <b>0.1875</b> | <b>9.4456</b> | <b>20.93%</b> | <b>0.0980</b> |

### 5.3.2 Experiment Setup

Before we train and evaluate our proposed C-ViT model, we first pre-process the two datasets. The first pre-processing stage was to perform a grid representation by dividing each city map of the two datasets (i.e. NYC and Chicago) into equally-sized grid cells each with a dimension of  $(2KM \times 2KM)$ . Secondly, similar to Wang et al., 2021b, we group all the accidents that happened in each grid cell based on their location over the reported duration time for each dataset (for each grid cells with no road segments/accidents, we set its traffic accident risk to zero). Thirdly, we split the data-sets into training, validation and testing. The strategy we followed for the splitting is similar to Wang et al., 2021b, where we use 60% for training, 20% for validation and 20% for testing while making sure that there is no overlapping accidents based on time (i.e. no accident happened in specific grid cell on specific time is shared between the three data splits). It is worth noting that the traffic accidents periodicity according to the two datasets was set to 1 hour. Finally, each data split is standardised by a mean and standard deviation normalisation so that it could help in accelerating the training process.

Regarding the implementation details of our C-ViT model, the size of the historical traffic risk maps  $X$  was set to  $7 \times 20 \times 20$  which corresponds to a total 7 historical traffic accident risks across the city grid with  $I$  rows  $\times$   $J$  columns of size 20. Here we chose 7 historical accident risks specifically to conform with the work done in the literature chen2018sdcae; Wang et al., 2021b for a fair comparison provided later in the paper. For each grid cell, the 7 historical accident risks comes from the most recent accident risks in past 3 hours in addition to the past accident risks in the last 4 weeks. The prediction horizon of the traffic accident risk was set to 1 (i.e next hour) similar to chen2018sdcae; Wang et al., 2021b. The hyper-parameters for our C-ViT model itself were set according to the model performance on the validation split. To that end, the  $D$  dimension for the linear patch embedding, the position embedding layer and the linear embedding layer of the historical contextual encoder was set to 64. The resolution of input patches  $P$  to the patch embedding layer was set to 5. The number of self-attention heads were set to 8 and the final output fully connected layer of our C-ViT model was set to 128. Since we formulated the traffic accident risk prediction task as an image regression task, we have therefore optimised our C-ViT model during the training phase using a weighted mean-squared error (MSE) loss function. The reason for using the weighted MSE loss function instead of using the standard MSE loss function, is to try to combat the unbalanced nature of the traffic risk prediction problem, also known as the zero-inflated problem bao2019spatiotemporal. The procedure for weighting our loss function is motivated by the focal loss introduced in lin2017focal, where we holistically divided the total training

samples into four distinctive classes based on their traffic accident risk values. Those risk values are (0, 1, 2,  $\geq 3$ ). Similar to Wang et al., 2021b, the loss function weights were set to 0.05, 0.2, 0.25 and 0.5 respectively. In total, we have trained our C-ViT model for 200 epochs using the Adam optimiser with a learning rate of 0.003 and the batch size was set to 32.

### 5.3.3 Evaluation Metrics

In order to evaluate the performance of our trained C-ViT model, we utilised the three commonly used metrics for the traffic accident risk prediction task **ma2018point**; Wang et al., 2021b, namely root mean squared error (RMSE), Recall and mean average precision (MAP). The three evaluation metrics are calculated as follows:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{n=1}^N (Y_n - \hat{Y}_n)^2}, \quad (5.3)$$

$$\text{Recall} = \frac{1}{N} \sum_{n=1}^N \frac{|H_n \cap A_n|}{|A_n|}, \quad (5.4)$$

$$\text{MAP} = \frac{1}{N} \sum_{n=1}^N \frac{\sum_{j=1}^{|A_n|} \text{PR}(j) \times \text{REC}(j)}{|A_n|}, \quad (5.5)$$

where  $N$  is the total number of samples to be evaluated,  $Y_n, \hat{Y}_n$  are the ground truth and the predicted risk values for all grid cells of sample  $n$  respectively.  $A_n$  corresponds to the set of grid cells of sample  $n$  that have an actual/true traffic accident risk values.  $H_n$  corresponds to the set of grid cells within  $A_n$  with the highest traffic accident risk values. On the other hand,  $\text{PR}(j)$  corresponds to the precision of the grid cells starting at 1 and ending at grid cell  $j$ . Similarly,  $\text{REC}(j)$  corresponds to the recall value for grid cell  $j$  which is set to 1 in case there was a traffic accident risk at it and set to 0 otherwise.

Based on the definition of these three evaluation metrics, we can deduce that the lower the score of RMSE is, the better is the quality of prediction coming out of the model. On the other hand, the higher the recall and MAP scores are, the better is the accuracy of the model.

### 5.3.4 Baselines

We have compared the performance of our proposed C-ViT model to 5 different baseline approaches from the literature and in the following we will briefly describe each approach:

- **RNN-GRU chung2014empirical**: This model is based on one variant of deep recurrent neural networks (RNN), the gated recurrent unit (GRU) model. This model casts the traffic accident risk forecasting problem as a time-series prediction problem and tries to model the temporal dependency among historical traffic accidents risk.
- **SDCAE chen2018sdcae**: This model is based on the stacked denoised convolution auto-encoder architecture, which focuses mainly on capturing/modelling the spatial features between different cells within a city grid area for a better prediction of the traffic accident risk.



TABLE 5.3: Performance evaluation of our C-ViT model against a number of baseline approaches from the literature over the high frequency times of accidents in the NYC and Chicago datasets.

| Dataset                      | NYC           |               |               | Chicago       |               |               |
|------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|
| Model                        | RMSE ↓        | Recall ↑      | MAP ↑         | RMSE ↓        | Recall ↑      | MAP ↑         |
| RNN-GRU chung2014empirical   | 7.3546        | 30.76%        | 0.1301        | 9.0421        | 18.66%        | 0.0758        |
| SDCAE chen2018sdcae          | 7.2806        | 31.22%        | 0.1536        | 8.7543        | 20.58%        | 0.1002        |
| H-ConvLSTM Yuan et al., 2018 | 7.2750        | 31.43%        | 0.1498        | 8.5437        | 18.93%        | 0.0770        |
| GCN wu2019graph              | 7.0958        | 33.04%        | 0.1647        | 8.4484        | 20.42%        | 0.0933        |
| GSNet Wang et al., 2021b     | 6.7758        | 34.15%        | 0.1769        | 8.6420        | 21.12%        | 0.1052        |
| C-ViT (ours)                 | <b>6.2658</b> | <b>34.46%</b> | <b>0.1802</b> | <b>7.0353</b> | <b>21.95%</b> | <b>0.1247</b> |

- **H-ConvLSTM Yuan et al., 2018:** As the name implies, this model combines both deep convolution layers with RNN-based LSTM layers to extract the spatio-temporal features of the traffic accident risk problem by having a sliding window over the city’s grid cells; this allows to have sub-regions that could potentially capture the heterogeneity among the different types of spatial regions.
- **GCN wu2019graph:** This model is a deep learning model that relies on graph convolution neural network to represent the historical traffic accident data as a graph to capture the long-term spatio-temporal dependency among historical traffic accidents risk data.
- **GSNet Wang et al., 2021b:** A recent model that learns the complex spatial-temporal correlations of traffic accidents risk by using a combination of GCN, LSTM and attention mechanism. To the best of our knowledge, GSNet is currently the SOTA method on the NYC and Chicago data-sets.

### 5.3.5 Results

In Table 5.2, we report the results of our C-ViT model in comparison to the aforementioned baseline approaches from the literature over the total testing splits for both NYC and Chicago data-sets. As it can be noticed, our model has outperformed all the baseline approaches from the literature in terms of RMSE, recall and MAP scores over the two data-sets. It is worth noting from the results, that those models (our C-ViT, GSNet, GCN and H-ConvLSTM) which account for the spatio-temporal property of the traffic accident risk prediction problem, are the top performing approaches on the two data-sets.

The closest competitor baseline approach to our C-ViT model, was the GSNet, which to the best of our knowledge, was the SOTA on the two data-sets before our proposed approach. As it can be seen, our C-ViT model has improved the RMSE, recall and MAP scores in comparison to GSNet especially across the Chicago dataset by a relatively large margin. Furthermore, our C-ViT has more competitive advantage over GSNet in terms of the efficiency. As it can be shown in Fig. 5.4, the number of parameters required by our C-ViT model for training are far lower than those needed for GSNet (saving more than 23x parameters) which makes our approach more suitable for real-time deployment.



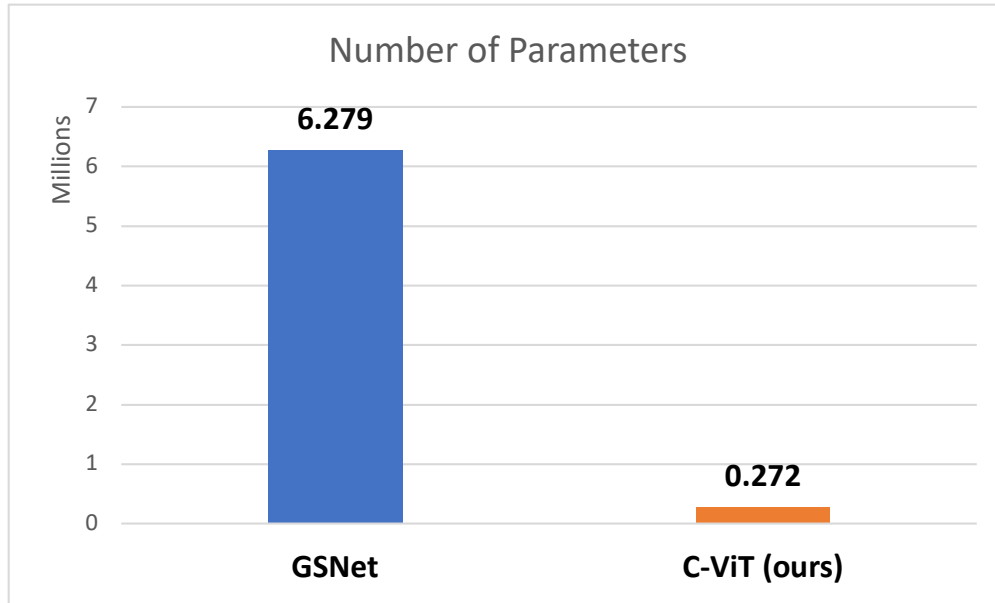


FIGURE 5.4: Comparison between our proposed C-ViT model and GSNet Wang et al., 2021b, in terms of the number of training parameters.

In order to further evaluate the performance of our proposed C-ViT model, in Table 5.3 we report the RMSE, recall and MAP scores of our model when compared to all the other baseline approaches over peak hours of frequent traffic accidents that resulted from the testing split of both the NYC and Chicago data-sets. Those times of high frequency of traffic accidents are essentially during morning/evening rush hours which are within 7:00-9:00 AM and 04:00-07:00 PM. As it can be seen from the reported results, our C-ViT model continues to achieve more robust results than all other compared baseline approaches. This further prove the utility and quality of our proposed approach that it has a consistent performance across different settings.

## 5.4 Coarse-Fine-Coarse Visual Transformer (CFC-ViT)

One of the issues to solve in the topic of accident risk prediction is the zero-inflated issue - the imbalance between the amount of non-zero and zero accident risk cells. This issue can be resolved by using a comparison mask or variations of focal loss Wang et al., 2021b. Another issue, which is usually ignored is the fine granularity of accident risk map. For example, in the grid representation, cells can be separated and of minimal 1x1 cell size (see Fig. 5.5).

Current computer vision methods applied to the task of accident risk prediction produce ‘blurred’ results due to intrinsic limitations of convolution network architectures li2021survey; gu2018recent. To resolve this issue we propose an alternative approach which consists of up-scaling the patches before the embedding to allow fine-grained processing by internal layers of transformer, and then down-scaling embedding to match the original output shape 5.6. Up-scaling may be a necessary step in the of use of asymmetric convolutional networks for segmentation to increase the detailization of results lo2019efficient since there are no upscaling layers at the final part of the asymmetric network. We upscale patches before the embedding to perform fine-grained processing of these patches. The dimensionality of embedding is also increased proportionally to the patch size; processed embedding

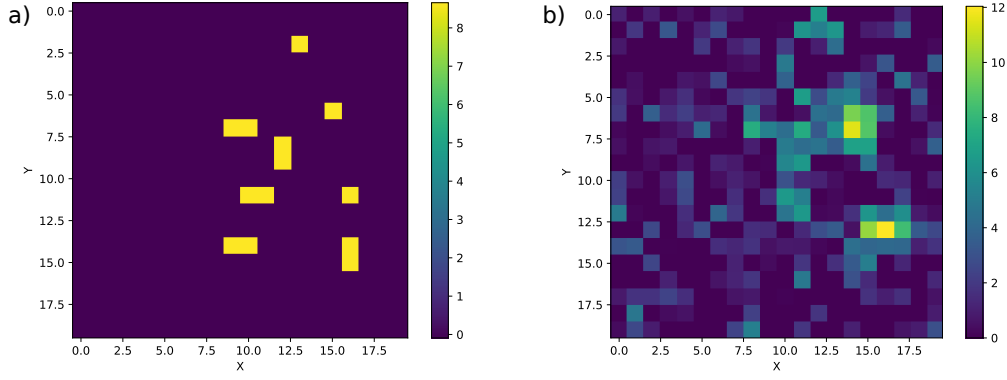


FIGURE 5.5: Example of GSNet predictions (after training for 2 epochs, when the best performance is observed): a) Actual map of the accident occurrence b) Predicted map of the accident occurrence.

TABLE 5.4: Performance evaluation of our CFC-ViT model on Chicago data set

|   | RMSE        | Recall | MAP   | HFT-RMSE    | HFT-Recall | HFT-MAP | Data set | CFC Scale factors                 |
|---|-------------|--------|-------|-------------|------------|---------|----------|-----------------------------------|
| 0 | <b>8.62</b> | 19.38  | 0.08  | <b>6.48</b> | 19.89      | 0.09    | chicago  | 4x, 0.25x                         |
| 1 | 9.24        | 18.84  | 0.06  | 6.85        | 20.85      | 0.07    | chicago  | 2x, 0.5x                          |
| — | —           | —      | —     | —           | —          | —       | —        | —                                 |
| - | 9.45        | 20.93  | 0.098 | 7.035       | 21.95      | 0.125   | chicago  | 1x, <b>baseline (multi-epoch)</b> |

then downscaled by the same rate. This allows the network to form intermediate results of higher dimensionality, which when down-scaled, will produce more fine-grained image.

Results for the Chicago data set show a significant improvement in the RMSE metric results both for 2x and 4x scale factors (see Table 5.4 where the RMSE is 8.62 as compared to 9.45 translating in a 8.78% improvement). There is an inverse dependence observed between the scale factor and the Recall or MAP metrics: the increase in the scale factor lowers RMSE but MAP and recall also decrease. However, given the robustness of the RMSE metric, the improvement is consistent.

Results for the NYC data set show that the prediction performance can increase at a specific scale factor (2x) and decrease at different scale factor (4x) (see Table 5.5. These results suggest that the optimal scale factor for each data set can exist, which leads to a deciated optimization task of finding the optimal scale factor value. Results both for NYC and Chicago data sets show a non-linear dependency between the RMSE, the MAP (mean average precision) or Recall metrics. These metrics are intended for different purposes (RMSE for the regression, MAP and Recall for the classification results) and therefore can produce different results based on the characteristics of the predicted values.

Overall, our new proposed CFC-ViT approach shows an improvement in the RMSE results, but

TABLE 5.5: Performance evaluation of our CFC-ViT model on NYC data set

|   | RMSE        | Recall | MAP    | HFT-RMSE    | HFT-Recall | HFT-MAP | Data set | CFC Scale factors                 |
|---|-------------|--------|--------|-------------|------------|---------|----------|-----------------------------------|
| 0 | 7.09        | 33.17  | 0.1808 | 6.45        | 33.90      | 0.1751  | NYC      | 4x, 0.25x                         |
| 1 | <b>6.81</b> | 32.15  | 0.1838 | <b>6.14</b> | 33.24      | 0.176   | NYC      | 2x, 0.5x                          |
| — | —           | —      | —      | —           | —          | —       | —        | —                                 |
| - | 7.0053      | 33.86  | 0.1875 | 6.2658      | 34.46      | 0.1802  | NYC      | 1x, <b>baseline (multi-epoch)</b> |

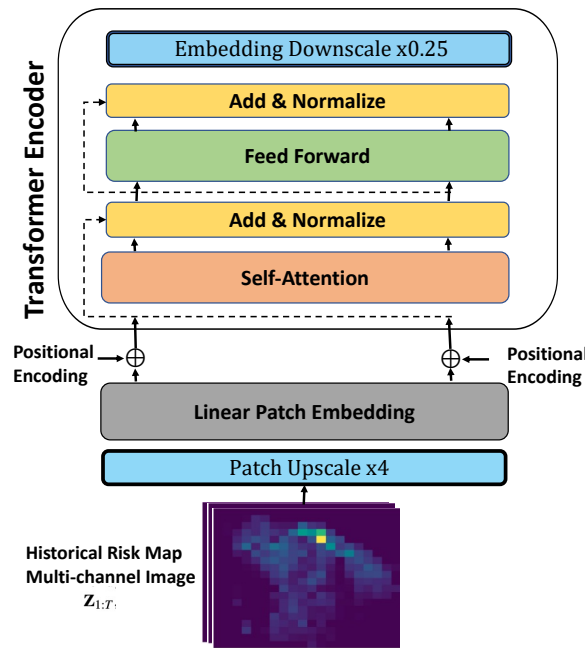


FIGURE 5.6: Coarse-Fine-Coarse Transformer

these results and other metrics depend on the scale factor parameter. The optimal scale factor can vary for each data set and can be found using other optimization techniques.

## 5.5 Application of the Static Map Generation

The use of Attention layers is a computer vision technique which implies an estimation of attention maps from different images. Since each image may have different areas of attention, the attention map is generated for every case of prediction (which we can call the dynamic attention estimation). But in the case of accident risk prediction, we predict on the same area each time. Therefore, we can use the statically generated attention map (static attention estimation). We evaluate multiple scenarios of combining dynamic (DA) and static attention (SA) estimations using varying combination operations. To further utilise the advantage of a non-volatile area, we also try to generate the Static Accident Risk Map (S-ARM) so our network needs to predict the offset of the accident risk (relative accident risk) from the statically generated risk map values instead of predicting the absolute accident risk values. Therefore, another contribution of this work is to further combine the Predicted Offset Accident Risk Map (PO-ARM) with the Static Accident Risk Map (S-ARM) (see Fig. 5.7).

### 5.5.1 Pipeline description

The generalisation performance of the Transformer model can be greatly improved by using one-epoch training [komatsuzaki2019one](#). Therefore we use results obtained from one-epoch training in further scenarios. Other parameters of the setup are the same as in Section 5.3.2.

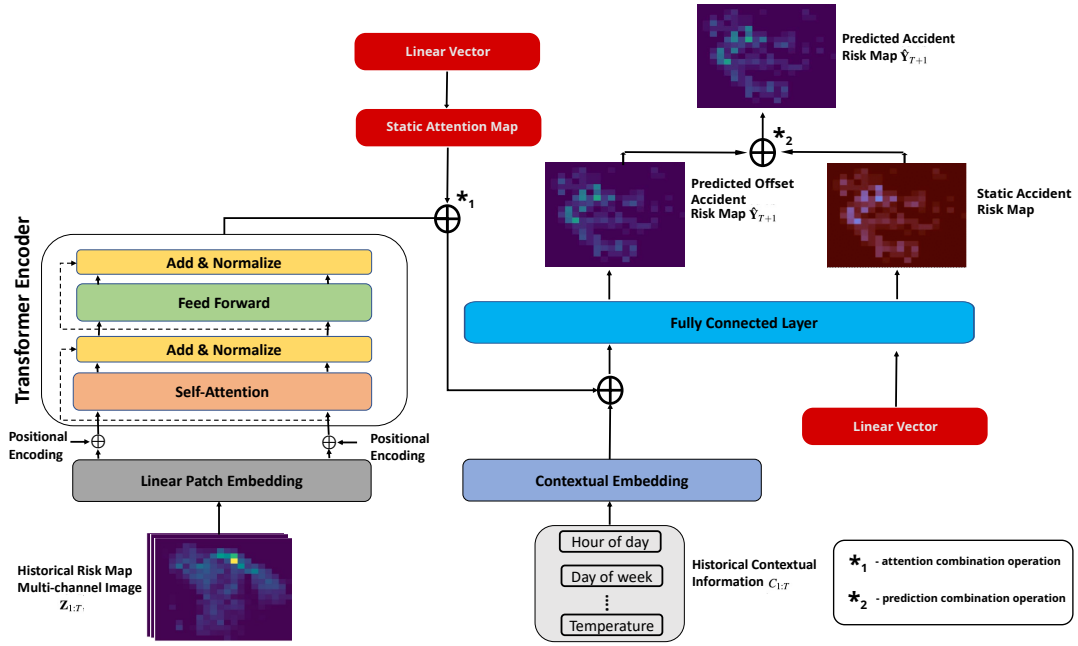


FIGURE 5.7: The building blocks of our proposed XViT model with Static Map Generation

### 5.5.2 Description of combination operations

The use of static map generation at the beginning of the attention layer as well as near the network output can remove the necessity for the network to predict the absolute risk values (static map is assumed to act as a static image and network is required to predict the relative risk from the one in a static map). We test multiple different approaches to achieve the benefit of using the static map generation. Different constraint functions can be used to limit the range of values observed from the static map. Also, the actual static map can be combined differently with the final and the intermediate network values.

Combination operations for the attention layer that we have considered:

1. None - using only the original pipeline structure with no static map,
2.  $\tanh(\text{map})+x$  - the static map is bounded by the tanh function in order to obtain the static map values distribution between  $(-1,+1)$ , combined with layer input values,
3.  $\tanh(\text{map}) * x$  - the same as above, but combined using a multiplication operation,
4. map - we use the static map instead of the attention layer inputs,
5.  $\text{map} + x$  - we combine the static map with the attention layer inputs using the “plus (+)” operation,
6.  $\text{sigmoid}(\text{map})+x$  - static map values are distributed between  $(0,+1)$  and are further combined with the attention layer inputs by using the “plus (+0)” operation (a linear offset combination),
7.  $\text{sigmoid}(\text{map}) * x$  - the same as above, but combined using the “multiplication (\*)” operation.

Combination operations for the network output that we have implemented:

1.  $\tanh(\text{map})+x$  in which the static map values (map) are combined with the intermediate predictions (transformer output -  $x$ ),
2.  $\text{sigmoid}(\text{map})+x$  - same as above, but using sigmoid as a static map constraining function,
3.  $\tanh(\text{map}) * x$  - tanh is used as a static map constraining function, and the static map values (map) are combined with the intermediate predictions by using the “multiplication (\*)” operation,
4.  $\text{head}(\text{map}+\text{non})$  - the static map values combined with the non-risk features and passed through the feed forward neural network,
5.  $\text{head}(\text{map})+\text{head2}(x+\text{non})$  - the static map is passed through a separate feed forward neural network, while other predictions together with the non-risk features are passed through the second network of the same structure,
6.  $\text{head}(\text{map}) * \text{head}(x)$  - the accident features are passed through the same network as the static map values and then combined using the “multiplication (\*)” operation,
7.  $\text{head}(\text{non})+\text{head}(x)$  - the non-risk features and the accident risk features are passed through the same network, and then combined using the plus operation,
8. None ( $\text{head}(x+\text{non})$ )- this is the original ViT implementation.

The constraint functions (tanh and sigmoid) are tested with the assumption that values close to the actual normalised accident risk values will be observed right after the network parameter initialisation. Due to the variation, we name this derivative model an XViT model.

### 5.5.3 S-ARM results

The results for the Chicago and NYC data sets are provided in Tables 5.6-5.7, RMSE results conveniently represented on Figures 5.9-5.8. The results are also provided for the high-frequency hours (HFT) - meaning the RMSE errors obtained only when using the HFT hours when more traffic is normally expected in the city. The use of the static map generation didn't show an improvement on the NYC data set. In fact, the results are a bit worse but closely related to the baseline (7.05 RMSE for the best combination vs 7.00 RMSE using original baseline multi-epoch approach); however we observe that there is an improvement in the recall results for the high-frequency hours (34.84 for the best combination vs 34.25 when using the multi-epoch baseline). But results for the Chicago data set show a very significant improvement across all the metrics (e.g. from 9.45 to 9.01 in RMSE, from 20.93 to 22.24 in Recall, from 21.95 to 23.46 in HFT-Recall). More than that, the 1-epoch training also shows a significant improvement in case of the baseline ViT structure (from 9.45 to 9.25 RMSE, from 20.93 to 21.77 Recall, from 7.035 to 6.93 in HFT-RMSE).

This slight reduction in the model performance in case of the NYC data set and significant improvement in case of the Chicago data set can be interpreted through the concept of local optima and data set size. There may be multiple local optima for the accident risk approximation across historical accident risk records (e.g. multiple average risk maps for different months). This optima can have an

TABLE 5.6: Performance evaluation of our XViT model for a number of combination operations on NYC data set. Top 20 results.

|    | RMSE   | Recall | MAP    | HFT-RMSE | HFT-Recall | HFT-MAP | Data set | S-ARM Combination Operations      |
|----|--------|--------|--------|----------|------------|---------|----------|-----------------------------------|
| 0  | 7.05   | 33.72  | 0.19   | 6.46     | 34.84      | 0.18    | nyc      | <b>tanh(map)*x, tanh(map)+x</b>   |
| 1  | 7.07   | 33.49  | 0.19   | 6.48     | 34.43      | 0.18    | nyc      | sigmoid(map)+x, sigmoid(map)+x    |
| 2  | 7.10   | 33.21  | 0.19   | 6.50     | 33.87      | 0.18    | nyc      | sigmoid(map)+x, head(map)+head(x) |
| 3  | 7.10   | 33.57  | 0.19   | 6.50     | 34.15      | 0.18    | nyc      | sigmoid(map)+x, tanh(map)+x       |
| 4  | 7.11   | 33.26  | 0.19   | 6.49     | 33.76      | 0.18    | nyc      | none, tanh(map)*x                 |
| 5  | 7.11   | 33.38  | 0.19   | 6.52     | 34.53      | 0.19    | nyc      | sigmoid(map)*x, tanh(map)+x       |
| 6  | 7.11   | 33.39  | 0.19   | 6.52     | 34.71      | 0.18    | nyc      | tanh(map)*x, sigmoid(map)+x       |
| 7  | 7.11   | 33.85  | 0.19   | 6.50     | 34.81      | 0.19    | nyc      | sigmoid(map)*x, tanh(map)*x       |
| 8  | 7.12   | 33.21  | 0.19   | 6.52     | 34.50      | 0.18    | nyc      | tanh(map)*x, head(map)+head(x)    |
| 9  | 7.13   | 33.33  | 0.19   | 6.54     | 34.71      | 0.18    | nyc      | tanh(map)+x, sigmoid(map)*x       |
| 10 | 7.13   | 33.36  | 0.19   | 6.52     | 34.64      | 0.19    | nyc      | sigmoid(map)+x, sigmoid(map)*x    |
| 11 | 7.13   | 33.43  | 0.19   | 6.54     | 34.53      | 0.19    | nyc      | tanh(map)+x, tanh(map)+x          |
| 12 | 7.14   | 32.88  | 0.19   | 6.55     | 33.41      | 0.18    | nyc      | map+x, head(map)+head(x)          |
| 13 | 7.14   | 33.09  | 0.19   | 6.55     | 34.18      | 0.19    | nyc      | map+x, sigmoid(map)+x             |
| 14 | 7.14   | 33.35  | 0.19   | 6.55     | 34.25      | 0.19    | nyc      | none, tanh(map)+x                 |
| 15 | 7.14   | 33.36  | 0.18   | 6.52     | 34.67      | 0.18    | nyc      | tanh(map)+x, tanh(map)*x          |
| 16 | 7.14   | 33.55  | 0.19   | 6.54     | 34.81      | 0.19    | nyc      | sigmoid(map)*x, sigmoid(map)+x    |
| 17 | 7.15   | 32.87  | 0.19   | 6.55     | 33.87      | 0.19    | nyc      | sigmoid(map)*x, head(map)+head(x) |
| 18 | 7.15   | 33.09  | 0.19   | 6.56     | 33.87      | 0.18    | nyc      | map+x, tanh(map)+x                |
| 19 | 7.15   | 33.24  | 0.19   | 6.55     | 34.71      | 0.19    | nyc      | map+x, sigmoid(map)*x             |
| 20 | 7.15   | 33.30  | 0.19   | 6.57     | 34.36      | 0.18    | nyc      | none, sigmoid(map)+x              |
| -  | -      | -      | -      | -        | -          | -       | -        | -                                 |
| -  | 7.25   | 33.39  | 0.19   | 6.53     | 34.25      | 0.19    | nyc      | <b>baseline (1-epoch)</b>         |
| -  | 7.0053 | 33.86  | 0.1875 | 6.2658   | 34.46      | 0.1802  | nyc      | <b>baseline (multi-epoch)</b>     |

TABLE 5.7: Performance evaluation of our XViT model for a number of combination operations on Chicago data set. Top 20 results.

|    | RMSE | Recall | MAP   | HFT-RMSE | HFT-Recall | HFT-MAP | Data set | S-ARM Combination Operations           |
|----|------|--------|-------|----------|------------|---------|----------|--|
| 0  | 9.01 | 22.24  | 0.11  | 6.80     | 23.46      | 0.13    | chicago  | <b>tanh(map)*x, tanh(map)+x</b>        |
| 1  | 9.06 | 22.30  | 0.10  | 6.85     | 23.59      | 0.13    | chicago  | tanh(map)*x, head(map)+head(x)         |
| 2  | 9.09 | 22.30  | 0.10  | 6.81     | 23.18      | 0.11    | chicago  | none, head(non)+head(x)                |
| 3  | 9.10 | 21.77  | 0.10  | 6.80     | 22.63      | 0.13    | chicago  | none, head(map)+head2(x+non)           |
| 4  | 9.12 | 21.88  | 0.10  | 6.80     | 23.05      | 0.13    | chicago  | sigmoid(map)*x, head(map)+head2(x+non) |
| 5  | 9.14 | 22.90  | 0.11  | 6.82     | 24.14      | 0.12    | chicago  | sigmoid(map)*x, head(non)+head(x)      |
| 6  | 9.15 | 22.12  | 0.11  | 6.91     | 22.63      | 0.13    | chicago  | none, sigmoid(map)+x                   |
| 7  | 9.15 | 22.18  | 0.11  | 6.94     | 22.91      | 0.13    | chicago  | tanh(map)*x, sigmoid(map)+x            |
| 8  | 9.16 | 21.59  | 0.11  | 6.98     | 21.81      | 0.13    | chicago  | tanh(map)*x, sigmoid(map)*x            |
| 9  | 9.17 | 22.72  | 0.10  | 6.82     | 24.01      | 0.12    | chicago  | tanh(map)*x, head(non)+head(x)         |
| 10 | 9.19 | 21.59  | 0.10  | 6.83     | 22.63      | 0.12    | chicago  | tanh(map)*x, head(map)+head2(x+non)    |
| 11 | 9.22 | 22.00  | 0.11  | 6.88     | 22.63      | 0.13    | chicago  | tanh(map)*x, head(x+non)               |
| 12 | 9.22 | 22.18  | 0.11  | 6.97     | 23.18      | 0.13    | chicago  | none, tanh(map)+x                      |
| 13 | 9.22 | 22.24  | 0.11  | 7.00     | 23.18      | 0.13    | chicago  | sigmoid(map)*x, tanh(map)+x            |
| 14 | 9.22 | 22.24  | 0.11  | 7.00     | 23.32      | 0.13    | chicago  | sigmoid(map)*x, head(map)+head(x)      |
| 15 | 9.24 | 22.00  | 0.11  | 6.99     | 22.77      | 0.13    | chicago  | none, head(map)+head(x)                |
| 16 | 9.25 | 21.77  | 0.10  | 6.93     | 22.36      | 0.12    | chicago  | <b>baseline (1-epoch)</b>              |
| 17 | 9.29 | 22.00  | 0.11  | 6.96     | 22.63      | 0.13    | chicago  | sigmoid(map)*x, head(x+non)            |
| 18 | 9.29 | 22.12  | 0.11  | 6.97     | 22.77      | 0.13    | chicago  | tanh(map)+x, head(map)+head2(x+non)    |
| 19 | 9.31 | 21.88  | 0.10  | 7.08     | 22.50      | 0.13    | chicago  | tanh(map)+x, sigmoid(map)*x            |
| 20 | 9.34 | 22.60  | 0.10  | 6.97     | 23.87      | 0.12    | chicago  | sigmoid(map)+x, head(non)+head(x)      |
| -  | -    | -      | -     | -        | -          | -       | -        | -                                      |
| 16 | 9.25 | 21.77  | 0.10  | 6.93     | 22.36      | 0.12    | chicago  | <b>baseline (1-epoch)</b>              |
| -  | 9.45 | 20.93  | 0.098 | 7.035    | 21.95      | 0.125   | chicago  | <b>baseline (multi-epoch)</b>          |

ability to show a good approximation of the accident risk, but since the road networks and the city structures change over time, different local optima can appear over time as well. So finding just one static accident map may not be optimal for a large data set, but may show benefit in the case of small data set (Chicago has just 44K accident records in comparison to 147K for NYC attributing to 1 full year of records and these are mostly short-time accidents in Chicago - just 8 months). We conclude that there is evidence that the proposed method and the use of multiple static maps can be a topic of the future research which can bring improvement over large data sets.

Another important observation is that the same set of combination operations gives the best results in the case of the ViT network with a generated static map: "tanh(map)\*x" in the attention layer and "tanh(map)+x" near the network output. This not only signifies the use of the constraint function tanh, but also shows where to use each combination operator (addition and multiplication). We also observe that the use of non-risk features is not present among the top 20 results for NYC data set (see Table 5.6, Figure 5.8), while for the Chicago data set it is present in 9 combinations out of 20, which may indicate the difference in quality of these features in both data sets.

## 5.6 Conclusion

In this work, we have presented a novel approach for the task of traffic accident risk forecasting. In our approach we re-formulated the problem as an image regression problem and introduced a unique

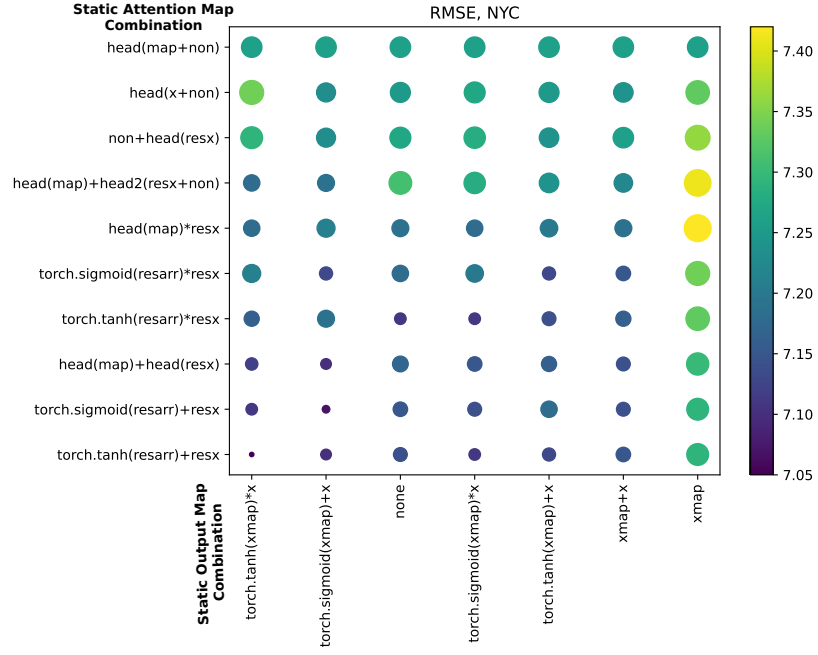


FIGURE 5.8: Root Mean Squared Error from performance evaluation of our XViT model for a number of combination operations on NYC data set

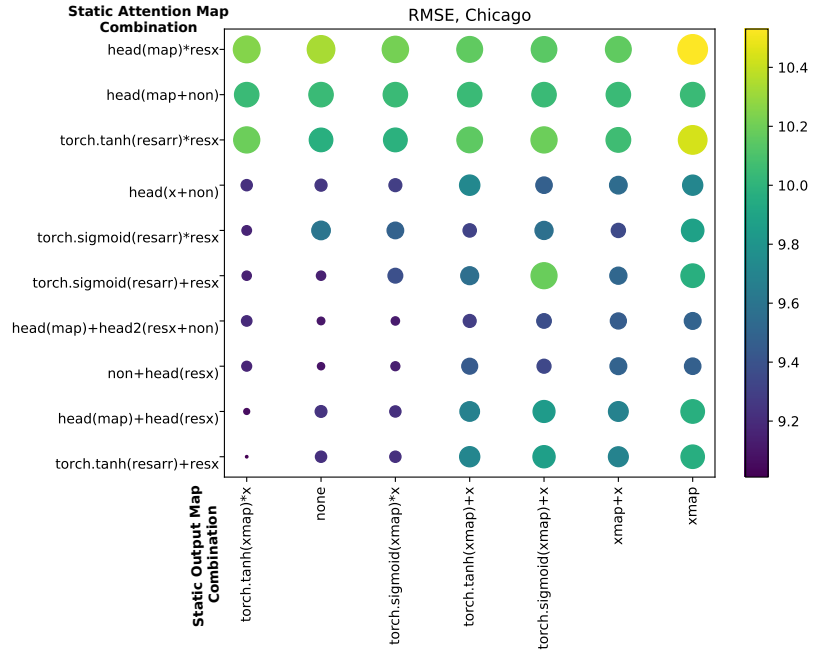


FIGURE 5.9: Root Mean Squared Error from performance evaluation of our XViT model for a number of combination operations on Chicago data set



contextual vision transformer network (C-ViT) that can efficiently model the traffic accident risk forecasting task from both spatial and temporal perspectives. The proposed approach has been evaluated on two publicly available data sets for the traffic accident risk problem. Furthermore, our proposed C-ViT model has been compared against a number of baseline approaches from the literature and it has outperformed them with a large margin while only requiring less than 23 times the number of training parameters.

The combination of static accident risk map with the ViT model (XVit) provides an even more significant improvement over the previous method in case of the New-York data set, thus establishing the new SoTA in the study area. The operation combination method has a potential for improvement (e.g. more different combination methods and constraint functions can be tested). Improvements in results obtained in the current research can also highlight the applicability of vision transformers for non-visual tasks.

The Coarse-Fine-Coarse Visual Transformer (CFC-Vit) architecture allows for fine-grained processing of the accident risk map and introduces an additional scale factor parameter which affects (and may improve) the prediction performance. There is a non-linear dependence between RMSE and the scale factor observed for both data sets, which may suggest that the optimal scale factor for the accident risk map processing may exist and be different for each data set.

Overall, the use of visual transformers and its variations for traffic accident risk prediction outperforms previously used approaches. Further applications of image and video processing methods may provide further improved results and open alternative approaches for the task of accident risk prediction.



## **Chapter 6**

# **Automatic Accident Detection, Segmentation and Duration Prediction using Machine Learning**

## 6.1 Introduction

The number of vehicles has been substantially increasing during the past decades, which currently leads to an increase in the number of traffic accidents Organization, 2015. The National Highway Traffic Safety Administration (NHTSA) reported more than 5 million traffic accidents happening in the United States during year 2013 Administration, 2013. Traffic Managements Agencies usually rely on Traffic Incident Management Systems (TIMS) to collect data on traffic accidents, including information on various accident, traffic state and environmental conditions. Accurately predicting the total duration of an incident shortly after it is being finished, will help in improving the effectiveness of accident response by providing important information to decide the required resources to be allocated (response team size, equipment, traffic control measures) Kim and Chang, 2011. Traffic accident is a rare event with stochastic nature. The effect of the accident can be observed as an anomalous state in the time series of traffic flow Theofilatos et al., 2016.

**Challenges:** The traffic accident analysis may be a challenging task due to incorrect or incomplete accident reports, including the set and the quality of the accident characteristics that have been reported. Accident reports can contain user-input errors related to the accident duration such as: 1) an approximate reporting of accident's start and end time 2) reporting of the accident start time could have been done after the incident finished in reality 3) a 'placeholder' accident duration reporting (filling report with the approximate duration value due to unavailability of data by the moment of reporting). In our previous research Grigorev et al., 2022b; Mihaita et al., 2019b we found that timeline-related errors are present in accident reports across three different data sets from both Australia and the United States of America, which creates the possibility of observing that such errors can occur in other data sets from around the world as well, due to multiple human and technical factors that can arise. To forecast the accident impact it is crucial to have an accurate and correct data regarding the observed disruption timeline. We emphasise that disruptions observed in a recorded traffic state can be automatically segmented and associated with a reported accident at the same time and location as when the accident occurred, which allows to eliminate user-input errors from reports and improve the accident duration prediction performance in many traffic management centres around the world. To help address this issue, in our paper we propose various methods for a correct traffic disruption segmentation, the method for an association between vehicle detector stations and accident reports.

Another important challenge is that many incident data sets around the world are private and not shared for public investigation; for those open data sets, there are several missing information fields, or even worse, incomplete information regarding the traffic conditions in the vicinity of the accidents. Even often publish crash data sets are limited in size as well and contain a very small number of records. This represents a tight constraint when testing one framework over multiple countries with different traffic rules and regulations. For our studies we have oriented our attention towards two big open data sets - CTADS (Countrywise traffic accident data set) which contains 1.5 million accident reports and the Caltrans Performance Measurement System (PeMS) which provides data on traffic flow, traffic occupancy and traffic speed across California. Despite both being extensive data sets, vehicle detector station readings from PeMS are not associated with traffic accident reports from CTADS either by time, location or coverage area. The lack of such association makes it impossible to analyse the relation between accidents and their effects on traffic flow and speed. To address this challenge, in our paper

we introduce the following mapping algorithm which will secure several steps such as :

- an association of Vehicle Detection Stations (VDS) with reported accidents in their proximity,
- a segmentation of traffic speed disruptions from detector readings,
- an association of detector stations with reported accidents (we will further show that this step is necessary due to many detected user-input errors in accident reports).

As a result, we obtain traffic disruptions segmented by the traffic speed associated with reported accidents. This association makes it possible to perform various important tasks of the accident analysis: 1) prediction of the traffic accident impact on the traffic speed based on accident reports, 2) prediction of the traffic accident duration derived directly from the effect of disruption on the traffic speed (impact-based duration), 3) analysis of disruption propagation (each detected disruption can be studied for spatial-temporal impact within the traffic network). Through this work, we will focus on the prediction of the impact-based accident duration and lay the foundation for a further research.

Overall, the main contributions (summarised in Figure 6.1) of our paper are as follows:

1. We conduct a fusion methodology of two large data sets (CTADS and PeMS) for a detailed traffic accident analysis. To the best of our knowledge, this is the first research study proposing the methodology for merging of two large data sets of such nature, which allows an association between observed disruptions in traffic flow and the reported accidents.
2. We propose a novel methodology for the disruption mining using a combination of different metrics (which we further find to have properties important for disruption segmentation): a) the Wasserstein metric, which allows us to measure the disruption severity and b) the Chebyshev metric, which provides a higher selectivity for the disruption mining and a rectangular shape of the disrupted segments, allowing an automated disruption segmentation. We detail all unique properties of both metrics utilized together to allow an accurate disruption segmentation.
3. We perform the estimation of traffic accident disruption duration from traffic speed via the above metrics which allows us to alleviate user-input errors in accident reports.
4. We evaluate multiple machine learning models by comparing both the reported and the estimated accident duration predictions extracted from traffic speed disruptions.
5. We introduce a new modelling approach which focuses on the amount and shape of the the disruption associated with an accident, which allows a further analysis and modelling of accident impact.

Overall, this research forms the foundation for a new early traffic accident disruption detection, traffic disruption speed impact analysis and the use of observed traffic accident durations for correcting errors in user reports. Moreover, this work contributes to our ongoing objective to build a real-time platform for predicting traffic congestion and to evaluate the incident impact (see our previous works published in (Mihaita et al., 2019b)-(Shafiei et al., 2020)-(Mao et al., 2021)).

The paper is further organised as follows: Section 6.2 discusses related works, Section 6.3.1 presents the data sources available for this study, Section 6.3 showcases the methodology, Section 5.3.5 presents the disruption segmentation results, showcases the result of data set fusion and Section 6.7 provides conclusions and future perspectives.

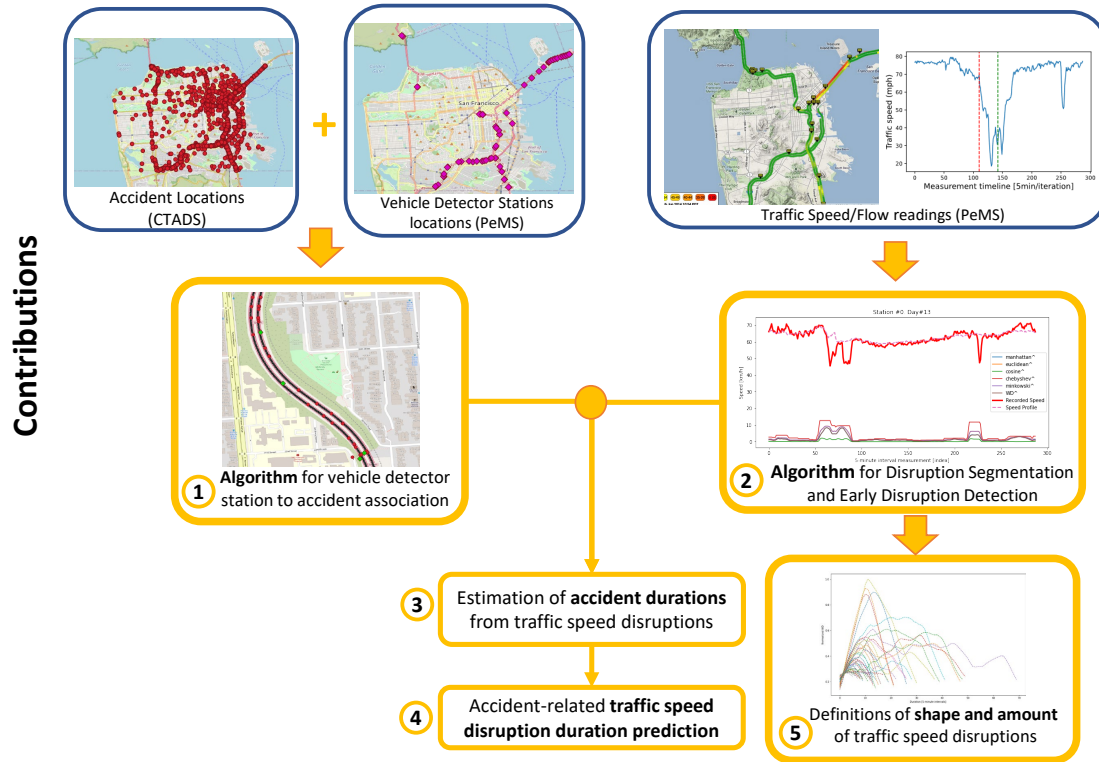


FIGURE 6.1: Contributions and data-flow schema for association of traffic speed readings with accident reports

## 6.2 Related Works

Multiple studies rely on user-input-based incident reports from Traffic Management Centers (TMC) with different machine learning models to predict the traffic incident duration Li, Pereira, and Ben-Akiva, 2018a. The use of traffic flow features is found to be rare and mostly specific - incident detection and incident impact prediction by using traffic flow Fukuda et al., 2020. In other words, traffic flow data is rarely combined with actual incident reports since it requires a higher system complexity and extensive data collection.

There were numerous studies related to accident detection from traffic flow using anomaly detection techniques Parsa et al., 2019. Various methods used for anomaly detection in time series are applicable for the task of traffic disruption detection. The ability to perform the detection of actual disruption, which should give us actual shapes of disruptions and time intervals allows in-depth analysis of usual accident statistics including the effect of the type of accident on the pattern of disruption in traffic flow. By integrating data on traffic state with accident reports we are able to further connect traffic flow disruption patterns to various accident characteristics (hour of the day, weather conditions, crash type, type of vehicle involved - truck/car Eboli, Forciniti, and Mazzulla, 2020, the effect of road pavement types Tsubota et al., 2018, road design and road operation Yannis et al., 2016, etc).

Various machine learning models are used to solve the task of traffic accident duration prediction Li, Pereira, and Ben-Akiva, 2018a including k-nearest neighbours (KNN) and Bayesian networks Kuang et al., 2019a, Recursive Boltzman Machines and Support Vector Machines(SVM) Xiao, 2021 and Random Forests (RF) Hamad et al., 2020b.

The definition of traffic incident duration phases is provided in the Highway Capacity Manual Alkaabi, Dissanayake, and Bird, 2011 and includes the following time-intervals: 1) incident detection - the time interval between the incident occurrence and its reporting, 2) incident response - time between the incident reporting and the arrival of the response team, 3) incident clearance time between the arrival of the response team and the clearance of the incident, 4) incident recovery - the time between the clearance of the incident and the return of traffic state to normal conditions. In this research, we rely on total incident duration - the time between incident occurrence and return of the state to normal conditions. Also, we analyse the subset of traffic incidents - traffic accidents. As we found during the data investigation, traffic accident duration is reported at the time when the incident is cleared by the response team, which doesn't include the duration of the effect that the accident produces on traffic flow. Traffic incident duration prediction studies rely on incident reports without emphasizing on the duration of observed incident effects. In this research we try to solve this issue by proposing the methodology for disruption segmentation from traffic speed.

Analysis of the effect of traffic incidents has been performed previously using Caltrans PeMS data, where the measure of incident impact was represented as a cumulative travel time delay Miller and Gupta, 2012, which is an aggregated value. However, traffic state recovery from disruptions is not necessarily following a single pattern - it may be slowly dissipating, we may observe secondary crashes, it may have a high or low impact, etc. Traffic accident duration prediction methodology relies on reported traffic accidents, but actual reports may contain user-input errors and be misaligned with the actual shape of disruption produced by the accident. Therefore, the approach for disruption segmentation may provide the accident duration estimated from the actual shape of disruption in traffic flow.

## 6.3 Methodology

The new proposed framework is represented in Figure 6.1 which we support across some initial definitions for our modeling approach (see next sub-section). First, we associate the road segments with their corresponding Vehicle Detector Stations (VDS) from the Caltrans PeMS data set, as well with the locations of reported accidents (see Algorithms 1 and 2 proposed in sub-section 6.3.4). The main outcome of this algorithm is that traffic accidents will get associated with the traffic flow, speed and occupancy readings from the VDS stations.

Second, we propose a new algorithm for early disruption detection and segmentation, detailed in sub-section 6.3.5. By detecting disruptions that occurred in time-space proximity of reported traffic accidents, we obtain the estimated traffic accident duration. This gives us much more information to include in the model training than just the simple accident duration: 1) the disruption shape in terms of modifications of speed data profiles from the standard patterns 2) the accident duration estimated from the impact on the traffic speed 3) the cumulative accident impact estimation.

### 6.3.1 Case study

Before diving into the methodology, we provide a brief introduction into the data sets in use for showcasing our approach, which helps establishing the modelling base and understanding of the steps taken. We make the observation that the current methodology can be applied on any incident and traffic state

data set which can contain a time component, and is not bounded to the chosen data sets for exemplification.

### CTADS: Accident reports data set

We rely on accident reports from the "Countrywide Traffic Accident Dataset" (CTADS), recently released in 2021 Moosavi et al., 2019a; Moosavi et al., 2019b, which contains 1.5 million accident reports collected for almost 4.5 years since March 2016, each report containing 49 features obtained from MapQuest and Bing services. We select the area of San-Francisco, U.S.A and extract data for 9,275 accidents (see Figure 6.2).

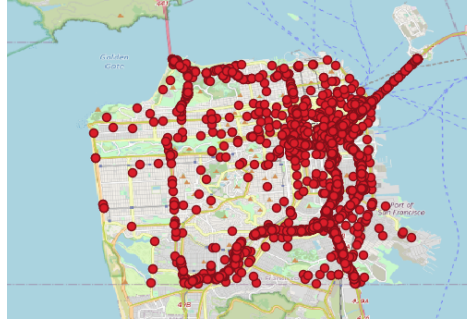


FIGURE 6.2: CTADS reported accidents for San-Francisco

### PeMS: Traffic speed and flow data set

We rely on Caltrans Performance Measurement System (PeMS) Chen et al., 2001 to collect data on traffic flow and speed. This data set provides aggregated 5-minute measurements of traffic flow, speed and occupancy across California. We decided to extract the data for the area of San-Francisco (see Figure 6.3a), which contains 83 Vehicle Detection Stations (VDS) placed in that area (see 6.3b), and we try to associate each traffic accident occurred with each of San-Francisco VDS in their 500m proximity using the algorithm detailed in the following section. In total, from 9,275 accidents in the area (extracted from CTADS) we have obtained 1,932 traffic incident reports which we were able to associate with the correct and complete traffic flow and speed readings from a VDS.

## 6.3.2 Speed difference estimation definitions

In the current study we compare the performance of multiple difference metrics that will help us to correctly estimate the impact of an accident and the deviation from the historical speed patterns. These metrics are defined as follows:

a) The Chebyshev difference is a measure of the maximum difference between corresponding elements of two one-dimensional vectors  $u$  and  $v$  and is expressed as:

$$D_{\text{Cheb}}(u, v) := \max_i |u_i - v_i| \quad (6.1)$$



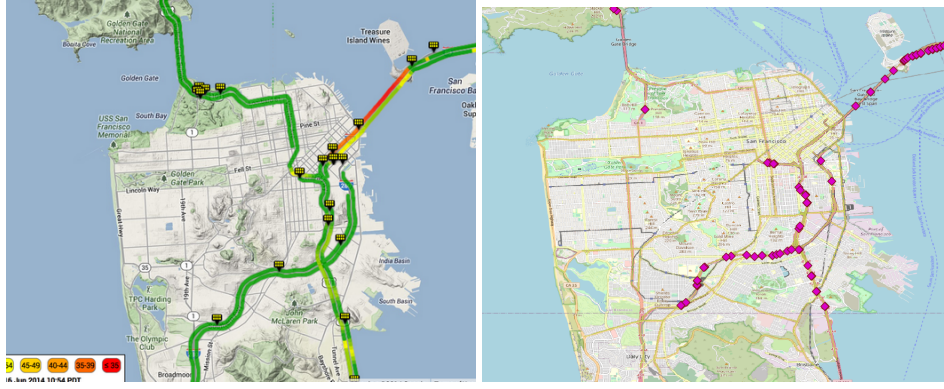


FIGURE 6.3: 1) PeMS data set area coverage for San-Francisco (the map is available at <https://pems.dot.ca.gov/>) 2) Mapping of the Vehicle Detection Stations from PeMS data set. OpenStreetMap excerpt showing San Francisco. Available at: <https://www.openstreetmap.org/#map=12/37.7612/-122.4395>

The metric that best captures the concept of similarity in relation to the application at hand should be considered optimal, regardless of the model used for classification or regression. This hypothesis has been tested and validated across nine datasets and five prediction models François, Wertz, and Verleysen, 2011. This Chebyshev metric, commonly used in data analysis, has been found to outperform other metrics on a variety of tasks.

b) The Wasserstein difference, also known as the earth mover's distance, is a measure of the minimum "work" required to transform one probability distribution  $u$  into another  $v$ . It is expressed as:

$$D_{WD}(u, v) = \inf_{\pi \in \Gamma(u, v)} \int_{\mathbb{R} \times \mathbb{R}} |x - y| d\pi(x, y) \quad (6.2)$$

This metric was introduced by Leonid Kantorovich in 1942 **kantorovich1942translocation** and has found applications in fields such as computer vision, image processing, and natural language processing.

c) The cosine difference, also known as the cosine similarity, is a measure of the similarity between two one-dimensional vectors  $u$  and  $v$ . It is expressed as:

$$D_C(u, v) = \frac{u \cdot v}{|u|_2 |v|_2}. \quad (6.3)$$

This metric is commonly used in information retrieval and has also found applications in recommender systems and document clustering **salton1988term**.

d) The Euclidean difference is a measure of the distance between two one-dimensional arrays  $u$  and  $v$  in a Euclidean space. It is expressed as:

$$D_E(u, v) = \left( \sum (w_i |(u_i - v_i)|^2) \right)^{1/2} \quad (6.4)$$

This metric is commonly used in fields such as machine learning, computer vision, and signal processing.

e) The Minkowski difference is a generalization of the Euclidean difference and is a measure of the distance between two one-dimensional arrays  $u$  and  $v$  in a Minkowski space. It is expressed as:

$$D_M(u, v) = \left( \sum |u_i - v_i|^p \right)^{1/p}. \quad (6.5)$$

This metric is a generalization of other distance metrics, such as the Manhattan distance (when  $p = 1$ ) and the Euclidean distance (when  $p = 2$ ), and is commonly used in fields such as physics, engineering, and data science **minkowski1909** **geometrie**.

f) The Bray-Curtis difference metric **bray1957** **ordination** between two vectors  $\mathbf{u}$  and  $\mathbf{v}$  is given by:

$$D_{BC}(\mathbf{u}, \mathbf{v}) = \frac{\sum_{i=1}^n |u_i - v_i|}{\sum_{i=1}^n (u_i + v_i)}, \quad (6.6)$$

where  $n$  is the number of dimensions in the vectors.

g) The Canberra difference metric **ivankovic2006** **comparison** between two vectors  $\mathbf{u}$  and  $\mathbf{v}$  is given by:

$$D_{Can}(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^n \frac{|u_i - v_i|}{|u_i| + |v_i|}, \quad (6.7)$$

where  $n$  is the number of dimensions in the vectors.

### 6.3.3 Accident duration prediction task definitions

Using all available data sets and the incident information, we first denote the matrix of traffic incident features as:

$$X = [x_{ij}]_{i=1..N_i}^{j=1..N_f} \quad (6.8)$$

where  $N_i$  is the total number of traffic incident records used in our modelling and  $N_f$  is the total number of features characterising the incident (accident severity, vehicles involved, number of lanes, etc) according to the accident report data set.

Traffic Speed represented as a vector with 5-minute averaged readings from Vehicle Detector Stations:

$$S = [s_i]_{i=1..N} \quad (6.9)$$

where  $N$  is the total amount of traffic speed readings.

Within this research we assess the performance of Machine Learning models on tasks of predicting reported and estimated accident duration. We define the task of accident duration prediction as a regression problem.

The incident duration regression vector ( $Y_r$ ) is represented as:

$$Y_r = [y_i^r]_{i \in 1..N}, y_i^r \in \mathbb{N} \quad (6.10)$$

and the regression task is to predict the traffic accident duration  $y_i^r$  based on the traffic incident features  $x_{i,j}$ . The regression models go via an 10-fold cross-validation procedure with hyper-parameter tuning.

To estimate the accident duration prediction performance we use the root mean squared error (RMSE):

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (A_i - F_i)^2} \quad (6.11)$$

where  $A_i$  - actual value,  $F_i$  - predicted value.

#### 6.3.4 Algorithm for vehicle detector station to accident association

In order to match correctly what traffic conditions reflect best the effects of each incident, we further define the association procedure between traffic accidents and VDS (Accident-to-VDS), for the San Francisco area. We observe that only a few traffic accidents have VDS stations in their proximity to allow a good traffic speed and flow extraction, as shown in Figures 6.2 and 6.3.

In order to find the traffic incidents for which we can have traffic flow and speed data, we develop a mapping algorithm (Accident-to-VDS) which consists of two parts (see Algorithm 2-3), defined by the following steps:

1. We extract primary and secondary road lines from Open Street Map.
2. Road segments are then transformed into points at 2-meters equal distance.
3. Each VDS station and accident are mapped to the closest road point (up to 10m distance).
4. From this step we use the following algorithm to process the point-based representation of VDS, accidents and road segments (see Algorithm 2). The `vdsPoints` array contains tuple of form (VDS ID, x and y coordinates), each point in `accidentPoints` contains an array `visitedBy` (initialized to be empty) to maintain a list of stations in proximity of the accident and `assignedVDS` as a resulting nearest VDS station to the accident along the road.

The algorithm relies on a recursive function to implement the process of visiting road points (see Algorithm 3). The association part of the algorithm works as follows:

1. We select the current VDS station.
2. We move (jump by points) in all possible directions available from the starting and forthcoming points in a 3m radius. This radius allows us to move along the road jumping between road points.

---

**Algorithm 2:** Accident-to-VDS: Accident to VDS mapping algorithm

---

```

Input: point
Output: None
Access global arrays: roadPoints, accidentPoints, vdsPoints
Function visitNearestPoints(VDSID, point, currentHops)
    accidents := findNearestAccidents(point, accidentPoints, 10m)
    for i = 0 to length(accidents) do
        a := accidents[i]
        a.visitedBy.append([VDSID, currentHops]);           //Recording visits from stations to
        internal accident list
    end
    if currentHops < 500/2 then
        ;                                                    //Limiting the travel distance from VDS
        roadpoints:=findNearestRoadPoints(point, roadPoints, 3m) for i = 0 to length(roadpoints) do
            rp := roadpoints[i]
            if VDSID not in rp.visitedBy then
                ;                                           //Preventing the infinite recursion
                rp.visitedBy.append(VDSID) visitNearestPoints(point, currentHops + 1)
            end
        end
    else
        Return
    end

```

---

Movement in all possible directions allows to grasp the propagation of the traffic congestion associated with the accident. The maximum available distance is set to 500m (250 jumps) and allows to limit the observable impact distance.

3. By moving across points we collect traffic incidents in the 5m proximity of each point and associate them with the current VDS station.

---

**Algorithm 3:** The recursive function for traveling across road points

---

```

Input: roadPoints, accidentPoints, vdsPoints
Output: assignedAccidents
for i := 0 to length(vdsPoints) do
    vds := vdsPoints[i]
    visitNearestPoints(vds, 0)
end
assignedAccidents = []
for i := 0 to length(accidentPoints) do
    accident := accidentPoints[i]
    if length(accident.visitedBy) > 0 then
        accident.assignedVDS = sort(accident.visitedBy, sortvalue = hops)[0]; //Choosing
        closes VDS station
        assignedAccidents.append(accident)
    end
end
return assignedAccidents

```

---

The algorithm is recursive and relies on the list of visited points for each VDS. At the end of the algorithm, we have a subset of traffic accidents with their associated VDS which allows us to extract the traffic flow and speed in the vicinity of the accident. Ideally, all traffic accidents should

have associated traffic flow but given their unavailability (due to detector coverage), we select accident reports which have associated traffic flow information currently available from the PeMS data set.

### 6.3.5 Algorithm for automated disruption segmentation (ADS)

Once the accidents have been mapped and associated to their VDS stations which allows us to select the flow/speed that match the day of the incident, etc, we are using the extracted traffic state parameters to propose a new automated disruption segmentation (ADS) method. The algorithm for the segmentation of disruptions via traffic speed works as follows:

1. A time series pre-processing step prepares all the data for segmentation (see Alg. 4):
  - (a) Calculate the average monthly profile for daily traffic speed measurements;
  - (b) Iterate over the traffic speed time series using a moving window of 1-hour time interval (in total there are twelve measurements of 5-minute each)
  - (c) On each iteration perform a comparison of a 12-unit window between the monthly profile and the current day of measurements. The resulting single value is added to the resulting time series sequence.
  - (d) Calculated the time series differences (TS) choosing the above defined metrics will be then adjusted by selectivity (using the power function, which will keep values closer to one for the least affected by the function and minor values the most suppressed) and normalized to produce nTS and pTS arrays respectively.
2. The time-series segmentation step (see Algorithm 5):
  - (a) A first order derivative (dTTS) is calculated for the resulting time series of the previous stage (nTS), which returns positive peaks when entering the disruption and negative peaks when exiting the disruption state.
  - (b) We iteration over resulting derivative time series to record the opening and closing of each disruption in each time series. If two consecutive positive peaks (opening times) are observed then we choose the largest one between the two (we will further debate on this aspect in our future work plans). We repeat the same for consecutive negative peaks.
  - (c) We then associate the detected disruptions with the accident reports: for each accident report, we extract the traffic speed time series on the day of the accident and if both opening and closing times are recorded, we perform an association of the accident with these times and extract the actual time series sequence for further analysis.

**Enhancing selectivity:** We use the convolution with the kernel (1,1,1), which attributes to the morphological dilation operation, to facilitate the work of the segmentation algorithm. By applying this convolution we make multiple consequent differences to be accumulated ; for example, assuming we have a sequence of 0.3, 0.1, 0.1, 0.2 and 0.2 as differences for each 5-minute step, therefore a total of 0.9 change over 4 iterations. The convolution (1, 1, 1) will produce the values of 0.5, 0.4, and 0.5 by making a sequence of high values from the sequence of small changes (see Figure 6.4). The dilation

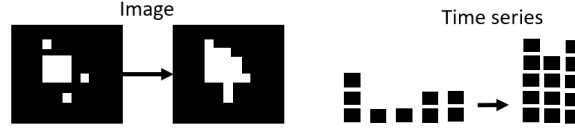


FIGURE 6.4: The application of dilation operation to an image and time series

operation is primarily used in computer vision tasks to make connected groups from closely placed scattered points to facilitate a further image analysis.

To obtain the monthly profile, the traffic speed measurement sequence was obtained for a duration of 1 month from the VDS before the accident occurred, and was done separately for each accident. This sequence then gets reshaped into a matrix of the form  $[number\_of\_days; 288]$ , where columns contain the total number of measurements across an entire day ( $24 \times 12 = 288$ ). The monthly average was then calculated across axis 1 (number of days) to obtain a vector with 288 values of measurements. This vector gets recalculated for a number of days of observations from each detector to be comparable with the VDS daily measurements.

---

**Algorithm 4:** Algorithm for automated disruption segmentation. Part 1

---

```

Input: monthlyProfile, speedReadings, selectivity, shift
Output: cTS
; //Accidents array contains a day number, starting and ending index for segmented
  traffic disruptions
step := 1
windowSize := 12
i := windowSize
lastDiff = 0
DS = []
while i < length(speedReadings) do
  A := speedReadings[i - windowSize : i]; //Look-back window of readings
  B := monthlyProfile[i - windowSize : i]
  diff := metric(A, B)
  DS.append(diff)
  lastDiff = diff
end
for i = 0 to windowSize do
  ; //Padding array with the latest observed value to obtain full-day readings
  DS.append(lastDiff)
end
pTS = power(TS, selectivity); //The use of power function to improve selectivity of
  significant disruptions
nTS = cyclicShift(shift); //The use of cyclic shift operation
nTS = normalize(pTS)
dTS = derivative(nTS); //First order derivative allows to decompose metric results
  into positive and negative change to the disruption amount
cTS = convolution(dTS, [1, 1, 1])
return cTS

```

---

As an observation, the constants *pThreshold* and *nThreshold* represent thresholds for change that observed in the time series of the metric derivative; they allow us to define a positive and negative change of the difference metric, the selectivity defines power function coefficient to suppress the non-significant and filter the most significant disruptions.

**Algorithm 5:** Algorithm for automated disruption segmentation. Part 2

---

**Input:**  $cTS, pThreshold, nThreshold, selectivity$   
**Output:** *Accidents*

```

; //Accidents array contains a day number, starting and ending index for segmented
  traffic disruptions
state := 0
Accidents = []
for i := 0 to length(cTS) do
  if cTS[i] > pThreshold then
    ; //Significant positive peak identifies the start of disruption
    if state <> +1 then
      state = +1
      enteridx = i
    else
      if cTS[i] > cTS[enteridx] then enteridx = i;
      ; //Choosing the largest change from previously observed
    end
  end
  if cTS[i] < nThreshold then
    ; //Significant negative peak identifies the end of disruption
    if state <> -1 then
      state = -1 exitidx = i
    else
      if cTS[i] < cTS[enteridx] then exitidx = i;
    end
  end
  if i mod 288 == 0 and i > 0 then
    ; //Reset segmentation procedure at the end of each day
    state = 0
    Accident.append([i div 288, enteridx, exitidx])
  end
end
return Accidents

```

---

**6.3.6 Modification of the algorithm for automated real-time early disruption detection**

Since our proposed algorithm doesn't look into the future and calculates different metrics based on the currently observed traffic speed and a few measurements in the past (11 units in the current study, equivalent to 55 minutes from the past), we can perform an early accident detection which will consist in calculating and comparing the first-order differential (FOD) of Chebyshev metric based on the monthly profile. The detection of significant positive peaks (e.g. 0.3-0.5 of normalized difference metric) can identify the amount of disruption in real-time. The end of the disruption can be detected using the same approach in real-time as well by observing a significant negative peak.

**6.4 Results****6.4.1 Data exploration and setup**

CTADS data set contains traffic accident reports, which after an initial data mining investigation, we found to contain several user-input errors; for example, a lot of traffic accident durations have been rounded to 30 or 360 minutes (see Fig. 6.5d)); or the incident start time which was reported is unrelated

to any disruptions observed by the vehicle detector stations in the proximity - see Figure 6.5 in which we have provided two different examples of speed recorded during two different accidents A-5198 and A-4490; the red lines indicate the official reported start and end time of the accidents, while in reality the accidents have had a long lag in spreading across the network - see Fig. 6.5a) or were reported much later than the official speed drop was recorded - see Fig. 6.5b).

At this step we observed a significant amount of user-input errors in accident reports, which affect the accident duration/impact analysis: 1) accidents can be reported earlier or later than its occurrence (observable disruption misalignment in time) 2) a report can be filled with "placeholder" duration values not representing the actual accident duration 3) there may be no observable disruption in traffic speed despite the accident report (due to placement and management of the accident) (false positive) 4) there may be accident-related traffic disruptions not grasped by accident reports (false negative). Therefore, incorrect accident start time, duration and end time, unreported presence or absence of disruption make it necessary to estimate accident duration characteristics from traffic state data instead of relying on user reports. In this paper our proposed methodology is really meant to solve the user-reporting issues related to traffic accidents and to be applied automatically on any data set, regardless of its nature or geo-location.

The use of PeMS data set allows to estimate the impact of accidents on the traffic states (flow, speed). For our scenarios, we choose the area of San-Francisco with accidents recorded from 2016 to 2020 in the CTADS data set. We then obtain Vehicle Detector Station locations from PeMS, the road network shape from OpenStreetMap and we perform an association of CTADS accident reports with VDS stations along the road within 500m proximity. We then try to segment the disruption time interval occurred on the day of an accident. Further, we associate observed disruptions in the traffic speed series with actual accident reports. The purpose of this step is to reduce user-input errors in accident reports and to enhance the modelling of traffic disruptions with an analysis of traffic speed.

#### 6.4.2 Metric performance comparison

We apply the difference metrics detailed earlier in Section 6.3 to a monthly traffic speed/flow profile (monthly readings averaged to one day) and reading on the day of the traffic accident. There are two approaches to applying the difference calculation: 1) a global difference - when we try to find the difference between the monthly profile and traffic flow/speed readings on the day of the accident; the global approach is too broad and will not allow the actual comparison between disruptions localized in time (metric results can be very similar between the very long subtle disruption and abrupt but impactful one). We measure the amount of difference that occurred within a moving time window (we choose twelve 5-minute time intervals equivalent to one hour). Traffic speed/flow readings from the moving window are taken right before the currently observed value to ensure that the difference estimation algorithm is not looking into the future.

To compare the metric performances we provide an example of speed readings from one of the detector stations. Each difference metric demonstrates its specifics as represented in Figure 6.6: 1) the Chebyshev metric, which we define as the maximum difference between the monthly profile and the observed readings, produces a noticeably rectangular shape and demonstrates a higher selectivity towards major disruptions than other metrics; the Chebyshev metric will be further used for the automated accident segmentation; 2) the use of Cosine metric allows to detect the change in the traffic state



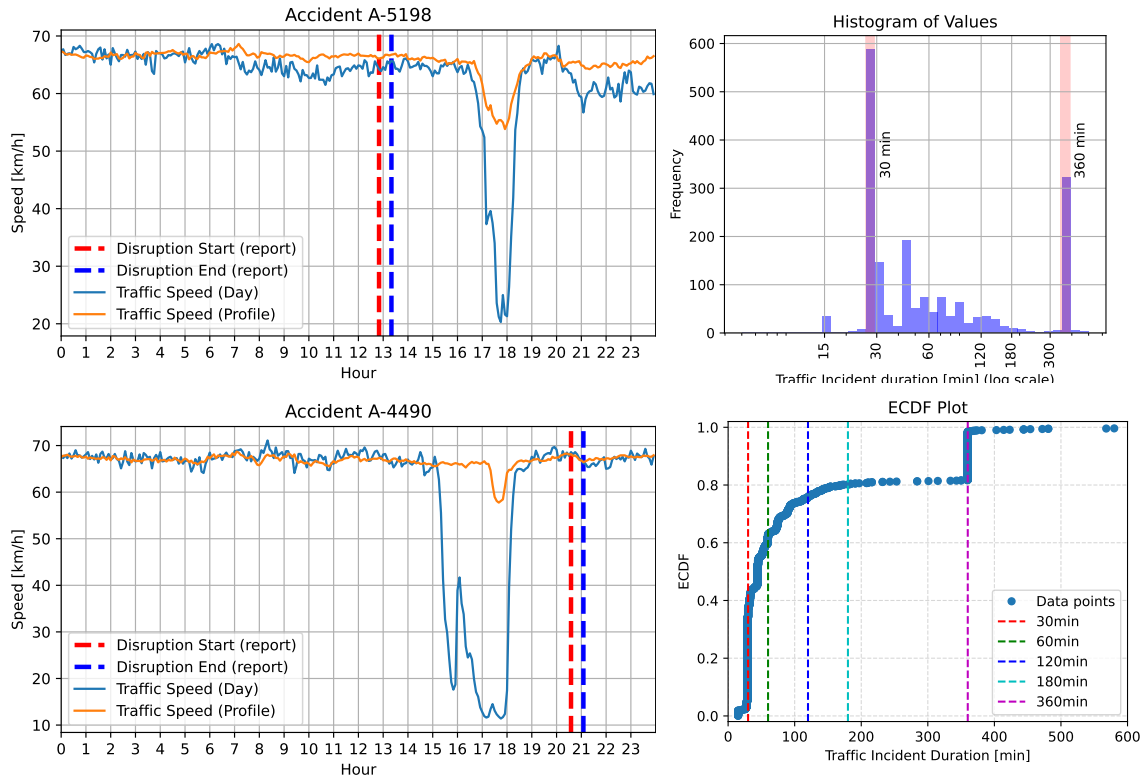


FIGURE 6.5: User input errors located within the CTADS data set

- speed decrease and increase both represented as positive peak values, 3) the Wasserstein difference allows for smooth representation of the amount of disruption (conceptually, it measures the amount of work necessary to change one shape into another, which we can rephrase as the amount of work produced by an accident to deviate the traffic state from the normal operation), 4) the Minkowsky, Euclidean and Manhattan difference metrics show little to no difference to the Wasserstein distance; we choose to use the Wasserstein difference since its connection to physical interpretation.

Examples of applying our proposed algorithm are presented on Figure 6.7. The 'Disruption Start' and 'Disruption End', which are represented as dashed blue and red vertical lines correspondingly show a reported accident timeline. The 'Day' (blue line) represents the traffic speed on the day of the incident and 'Profile' shows the average speed for every 5-minute interval across 14 days of measurements. Application of the 'Wasserstein distance (WD)' shows a gradual measurement of the observed disruption, while the 'Chebyshev' metric shows the nearly rectangular outline of a time interval where disruption is observed. This 'rectangular' result of the 'Chebyshev' metric was the main consideration for the development of the presented algorithm. One of the main observations from the figures as well as from the procedure of manual markup was that accidents were primarily reported 1-2 hours after the return of traffic state to normal conditions (which we define as the end of disruption). Other observation is that accident timeline is often misreported as a 'rounded' value of either 30 or 360 minutes. Application of both metrics shows a clear outline of disruption shape observed in traffic speed.

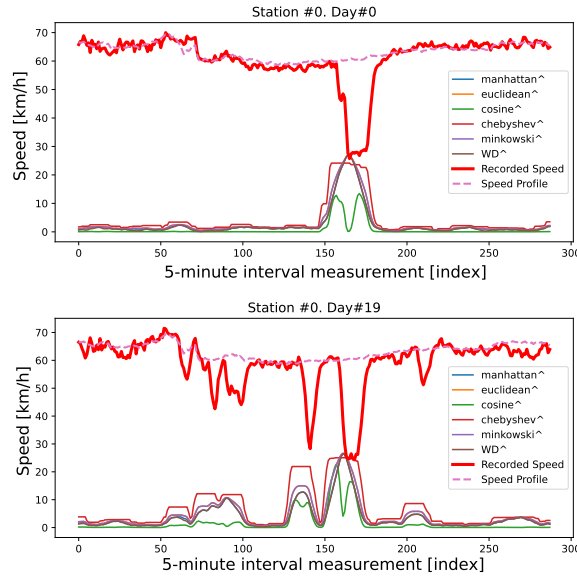


FIGURE 6.6: Various metrics applied to difference between recorded speed and speed profile

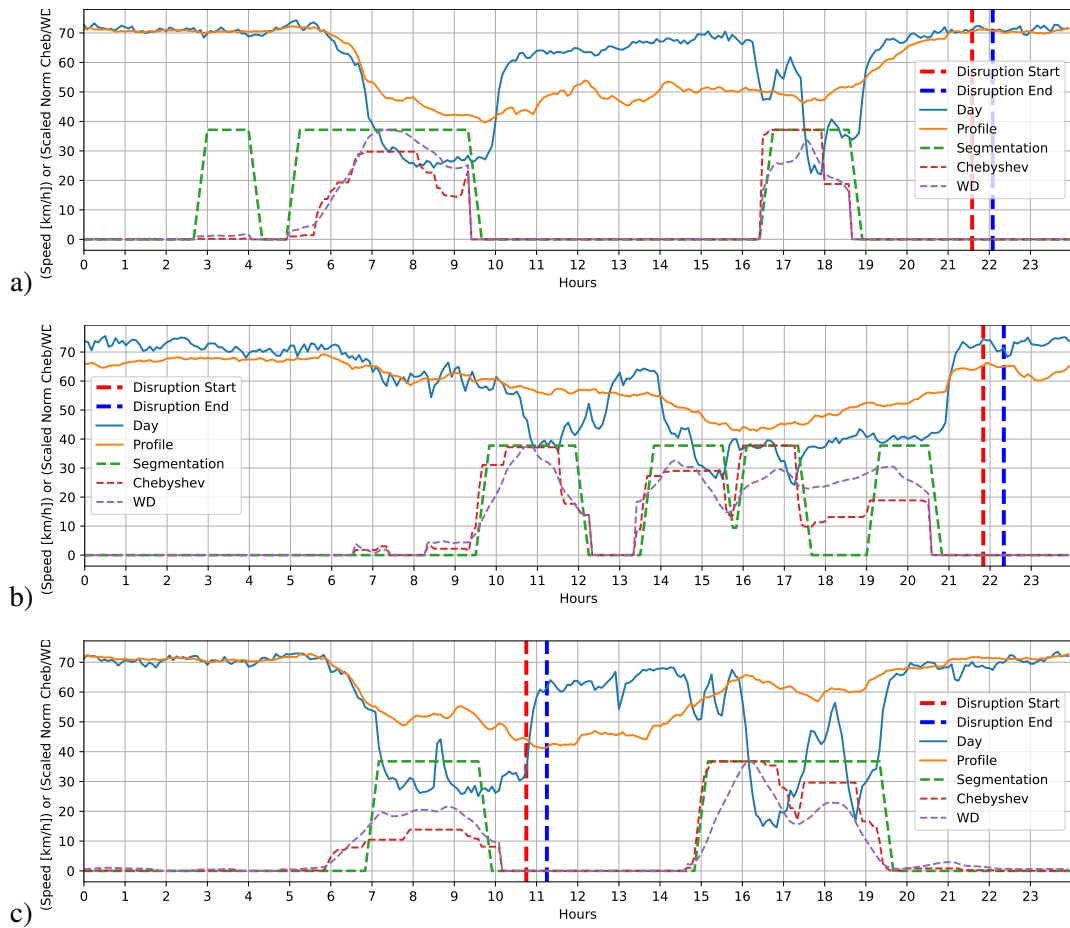


FIGURE 6.7: Results disruption segmentation algorithm application for accidents a) A-5764, b) A-8119, c) A-9931

### 6.4.3 Combination of our proposed methodology with modern methods for accident scene segmentation

Image segmentation methods can be utilized to output a degree to which an accident is observed in an image [zhang2021exploring](#), ultimately helping to create an accident timeline. By leveraging the power of semantic segmentation, the method can quantify the extent of the accident by assigning scores or probabilities to different elements within the scene. Here's how this can be done: 1) Accident-related object detection (Spatial analysis): Semantic segmentation can identify accident-related objects in the scene, such as damaged vehicles, debris, or injured pedestrians; by calculating the proportion of these objects within the segmented image, it is possible to assign a degree or score that represents the severity or extent of the accident at that specific moment, 2) Temporal analysis: by analyzing the segmented images over time, we can track changes in the accident scene, such as the motion of vehicles or the appearance of new accident-related elements. This enables the creation of a timeline that reflects the progression of the accident and the associated changes in the severity or extent of the event, 3) Probability-based analysis: advanced segmentation methods can output probability maps that indicate the likelihood of each pixel belonging to a specific class or label; by analyzing these probability maps, it is possible to compute a score that represents the degree of accident occurrence within the scene over the timeline, 4) Accident phase classification: The degree to which an accident is observed can also be used to classify distinct phases of the accident, such as pre-collision, impact, and post-collision. By evaluating the changes in accident-related object proportions or scores over time, the segmentation method can identify critical moments or transitions between different accident phases. This information can be used to construct a detailed accident timeline that highlights the key events and their corresponding degrees of severity.

There is potential to connect our proposed methodology with accident scene segmentation research approaches to create a more comprehensive and accurate framework for analyzing traffic accidents and predicting disruption durations. Here's how the two research approaches can be integrated: 1) Improved incident duration prediction: The segmentation output from the first research can be used as input for the early detection and disruption segmentation algorithm in the second research. This would allow for a more accurate identification of critical moments in the accident timeline and better prediction of incident durations, 2) Integration of mathematical metrics: the Wasserstein and Chebyshev metrics proposed in the second research can be used to refine the segmentation results obtained from the accident scene segmentation over timeline. This would help to improve the performance of the accident scene segmentation and contribute to a more accurate incident duration prediction, 3) Joint machine learning model: The speed disruption segmentation from our research can be combined with the semantic segmentation methods to create a joint machine learning model. This integrated model could leverage both the event-driven dynamic context and the mathematical metrics for segmentation to improve its predictions for incident duration and severity. By connecting these two research approaches, a more comprehensive framework for analyzing traffic accidents and predicting disruption durations can be developed. This integrated approach would benefit from the strengths of both methods, enabling more accurate and reliable predictions for incident durations. Ultimately, this could lead to improvements in road safety, emergency response, and traffic management.

In conclusion, image segmentation methods can be employed to not only segment the accident scene but also to quantify the degree to which an accident is observed in an image. This information

can be used to create an accident timeline that reflects the progression of the accident, the severity of the event, and the critical moments when interventions or safety measures could have been taken. This approach in combination with our proposed disruption segmentation method can potentially contribute to better accident analysis, road safety improvements, and more effective emergency response strategies.

#### 6.4.4 Automated disruption segmentation results

Figure 6.7 presents the results obtained from our algorithm for the automated disruption segmentation. The segmentation line (dotted blue) represents the estimated disruption intervals represented as 0 and 1 to perform our visualisation investigation better. Figure 6.7a) shows that there may be multiple observed disruptions in a  $300 \times 5 = 1500$  time interval. Due to errors in accident reports regarding the starting time and the duration of the accident, it is non-trivial to determine which disruption is associated with the accident. The situation may be easier in the case when only one disruption is observed during the day. According to our algorithm, we select the largest disruption on the day the accident was reported. Figures 6.7b) and 6.7c) highlight additional specific situations which need to be considered: 1) higher traffic speed at the end of the day than observed from the monthly profile, 2) unstable traffic speed approaching normal traffic conditions with high frequency, 3) slight misalignment of disruption intervals with the visually observed disruption intervals. All these problems can be addressed by using manual segmentation with deployment of Deep Learning models since there are advanced computer vision methods proposed in recent years (e.g. autoencoders for segmentation).

#### 6.4.5 Comparison of estimated, reported and manual markup of accident durations

There is a significant difference between the estimated and the reported accident durations that we would like to highlight: 1) the reported accident durations contain a large amount of 30 and 360 minutes duration values (nearly 40% of data - see Figure 6.8a)) while the estimated accident durations using our approach have an average duration of 58 minutes, while the reported is 108 minutes (which is by assumption skewed due to 360 placeholder values), 3) the estimated accident durations are distributed between 90 and 355 minutes (0.10 and 0.90 quantiles correspondingly) (see 6.9b)), while the reported durations are distributed between 29 and 360 minutes (see 6.9a) and manually detected disruptions distributed between 75 and 440 minutes), which highlights that disruptions observed from traffic speed are much shorter than reported in the original data set, 4) There is no noticeable correlation between observed and reported durations with high amount of horizontal anomalies in reported accident durations (see Figure 6.9). Traffic accident duration is most common to follow log-normal or log-logistic distribution Li, Pereira, and Ben-Akiva, 2018a and on resulting plots, we see that accident reports are found to represent log-normal distribution to less extent than manual markup or estimated accident duration.

To perform the ablation study, we perform a manual markup of disruptions observed in traffic speed for 800 accidents, which will be discussed in the corresponding section.

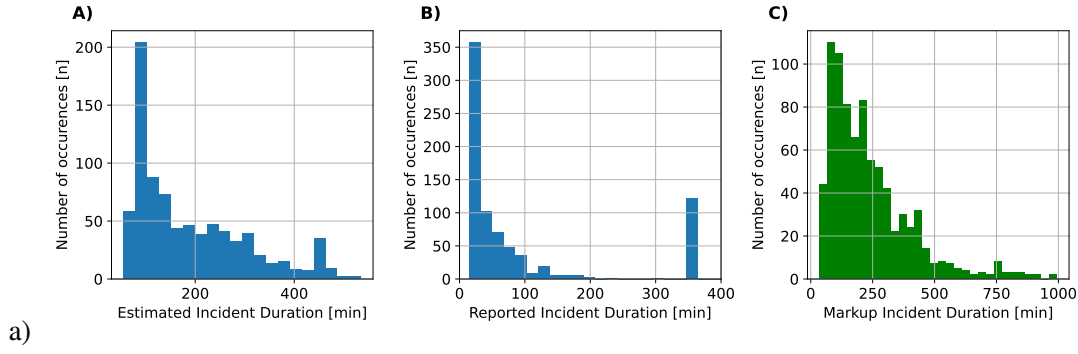


FIGURE 6.8: Distribution of accident durations for a) estimated, b) reported accident durations for the area of San Francisco, c) results of manual markup of disruptions observed in traffic speed

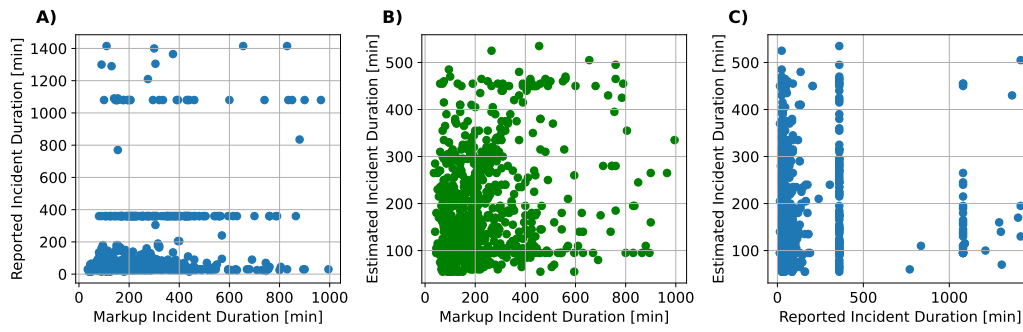


FIGURE 6.9: Scatter plot for a) estimated and b) reported accident durations for the area of San Francisco, c) results of manual markup of disruptions observed in traffic speed

#### 6.4.6 Extraction of disruption shapes

In previous subsections we applied a Chebyshev metric to perform segmentation of disruptions. To analyse the disruption impact we apply the Wasserstein difference between monthly speed profile and daily traffic speeds and extract the corresponding disruption intervals. Wasserstein difference, originally named an Earth Mover distance, has an intuitive physical interpretation - the minimum "cost" of altering one pile of earth into the other, which is assumed to be the amount of earth that needs to be moved times the mean distance it has to be moved. In application to traffic state, it is the minimum amount of work necessary to alter the traffic state to disrupted condition, or in other words - the amount of disruption. We compare normalized metric values since every at every vehicle detector station there is a different average traffic speed. As in our proposed algorithm, we use a 12-units moving window (one hour) to estimate the Wasserstein difference between traffic speed measurements and provide the plot for the first 40 segmented disruptions, which allows for shape analysis of traffic disruption amount (see Figure 6.10): 1) We observe the similarity between multiple disruptions - they have a 'hill' shape, 2) there are secondary (double 'hill') and long-lasting disruptions. The observed shapes can be defined through the parametric equation to perform the classification of disruption effects and facilitate the prediction of disruption impact timeline since we observe that high-peak fast-ascending disruptions have a probability to end sooner than slowly ascending ones. The analysis of the speed of ascendance has potential to perform the early classification of disruptions, which is planned for further

research.

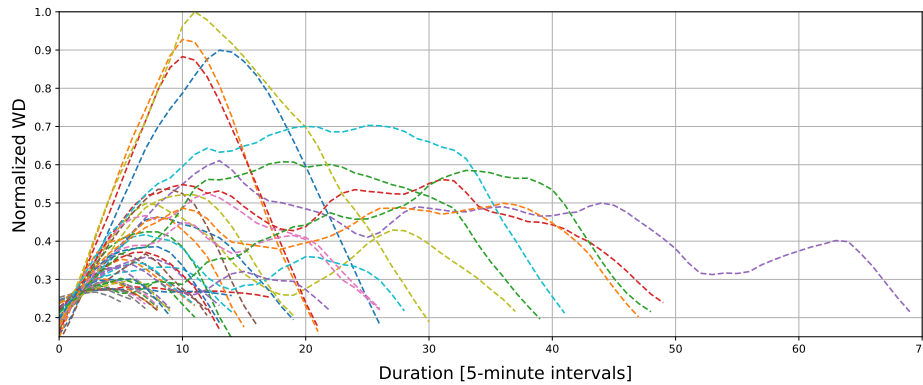


FIGURE 6.10: Normalized Wasserstein distance plot for disruption shapes extracted for segmented intervals

### 6.4.7 Accident duration prediction

We further compare a regression model prediction performance on the CTADS data set by using on the training data set both our estimated versus the reported accident durations. We report results of a 10-fold cross-validation over 820 accident reports for which we performed a Vehicle Detector Station association and manual markup of traffic disruptions from traffic speed for ablation study. Firstly, we need to consider that the performance using reported durations from CTADS can be affected because of the presence of user-input errors in the form of placeholder values. Secondly, the nature of estimated accident durations is different since accident response teams usually report the end of the accident at the moment they finished the accident clearance, without estimating the time for the traffic to return to a normal condition, which would require additional presence, calculations and access to measurements.

We have further extended the current results by adding newtables with several machine learning models on the task of predicting a target variable.

Table 6.1 shows the Mean Absolute Error (MAE) results. The model with the lowest MAE is the CatBoost model, with an estimated MAE of 17.55, followed by the Ridge Regression with an estimated MAE of 17.87. The highest MAE is reported by the Linear Regression model (76.76). The CatBoost model outperforms all the other models by a significant margin, with the next best model (Ridge Regression) having an estimated MAE that is only slightly lower.

Table 6.2 shows the Root Mean Squared Error (RMSE) results. Here, the CatBoost model also has the lowest RMSE, with an estimated value of 22.55. The next best model is the Ridge Regression with an estimated RMSE of 22.21. The highest RMSE is reported by the SVM model, with an estimated value of 208.29. As with the MAE results, the CatBoost model outperforms all the other models by a significant margin. All the methods use default parameters as they are presented in Scikit-learn **scikit-learn** and corresponding modules.

When we are using accident reports to predict the estimated accident duration, we obtain a better performance using the RMSE metric across all the regression models, which may be connected to the lower amount of long accident durations than reported.

Overall, the CatBoost model consistently outperforms all the other models across all metrics.

TABLE 6.1: Mean Absolute Error (MAE) Results

| Model                               | Reported   | Manual   | Estimated   |
|-------------------------------------|------------|----------|-------------|
| KNN Kuang et al., 2019a             | 44.11      | 26.22    | 19.73       |
| RandomForest Hamad et al., 2020b    | 26.52      | 21.89    | 17.21       |
| XGBoost Chen and Guestrin, 2016     | 24.22      | 23.06    | 18.29       |
| LinearRegression                    | 76.76      | 24.12    | 17.82       |
| LightGBM Ke et al., 2017            | 36.57      | 22.43    | 18.26       |
| SVM Xiao, 2021                      | 84.82      | 23.70    | 17.55       |
| GBDT ye2009stochastic               | 26.50      | 22.37    | 17.46       |
| CatBoost dorogush2018catboost       | 23.96      | 21.58    | 17.55       |
| NeuralNetwork gallant1990perceptron | 55.34      | 24.33    | 19.27       |
| RidgeRegression mcdonald2009ridge   | 84.72      | 24.26    | 17.87       |
| Target                              | (Reported) | (Manual) | (Estimated) |

TABLE 6.2: Root Mean Squared Error (RMSE) Results

| Model                                | Reported   | Manual   | Estimated   |
|--------------------------------------|------------|----------|-------------|
| KNN Kuang et al., 2019a              | 142.97     | 35.27    | 24.45       |
| Random_Forest Hamad et al., 2020b    | 93.73      | 29.94    | 21.79       |
| XGBoost Chen and Guestrin, 2016      | 82.67      | 31.75    | 23.61       |
| Linear_Regression                    | 117.53     | 32.54    | 22.35       |
| LightGBM Ke et al., 2017             | 99.77      | 30.55    | 23.58       |
| SVM Xiao, 2021                       | 208.29     | 34.23    | 23.66       |
| GBDT ye2009stochastic                | 73.14      | 30.46    | 22.11       |
| CatBoost dorogush2018catboost        | 73.05      | 29.64    | 22.55       |
| Neural_Network gallant1990perceptron | 124.21     | 33.38    | 23.18       |
| Ridge_Regression mcdonald2009ridge   | 134.71     | 32.48    | 22.21       |
| Target                               | (Reported) | (Manual) | (Estimated) |

## 6.5 Ablation study

In this paper, we propose using the F1 score to estimate the quality of time interval segmentation in binary time series (see Figure 6.11) in which we provide two different examples of different stations with both manual markups of the incidents - red markups- and our segmentation algorithms - blue markups- that is more efficient at detecting multiple incidents throughout the 24h time period and not only one single isolated event. The value on Y-axis shows a positive 1.0 value if the interval contains the disruption. Examples are provided for Accidents with ID A-1024015 and A-1034382 from CTADS data set.

Given a ground truth dataset with original reported accident duration, we perform a manual labelling of segments and obtain a set of predicted segments obtained from our automated segmentation algorithm, we compute the precision and the recall of the algorithm, and then combine them into a single F1 score.

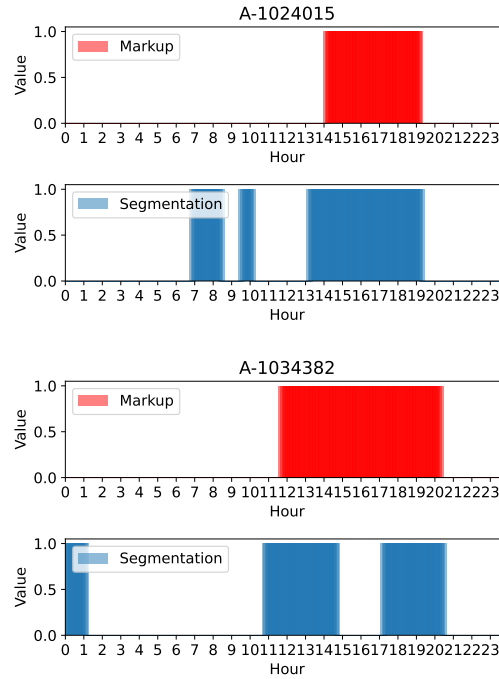


FIGURE 6.11: Manual markup and algorithm segmentation comparison. Time series segments represented as binary values of 0 and 1.

F1-score is a popular metric used to evaluate the quality of binary classification models defined as follows:

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

where true positives are the number of correctly classified positive instances, false positives are the number of negative instances classified as positive, and false negatives are the number of positive instances classified as negative.

F1-score is defined as the harmonic mean of precision and recall, given by:

$$\text{F1-score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

F1-score ranges from 0 to 1, with higher values indicating a better classification performance.

In the case where a time series is represented as a series of points with values of 1 for segmented intervals and 0 for intervals with no segments, F1-score can be applied to estimate the quality of the time interval segmentation.

To apply the F1-score, we need a ground truth dataset with manually labelled segments (and we obtain this manual markup for 820 accidents), and a set of predicted segments obtained from our automated segmentation algorithm. We can use these two sets to compute the precision and recall of the segmentation algorithm, and then combine them into a single F1-score.



Precision measures the proportion of true positives among all the predicted positives. In the context of time interval segmentation, the precision measures the accuracy of the algorithm in detecting the true segments. The Recall measures the proportion of true positives among all the actual positives. In the context of time interval segmentation, the recall measures the completeness of the algorithm in detecting all the true segments.

To apply the F1-score to estimate the quality of time interval segmentation, we can compute the precision and recall for each segment, and then compute the overall F1-score as the weighted average of precision and recall, weighted by the number of segments. This provides a single metric that reflects the quality of the time interval segmentation.

As a result (see Figure 6.12), the official reported incident segmentation is found to be very off (with a mean F1-score of 0.29 - Figure 6.12a)); next, the segmentation done by the algorithm while selecting only the interval closest to the reported timeline yields the highest average F1-score of 0.51 - Figure 6.12c)) with a peak at 0.3; lastly, when considering multiple segmented incident intervals detected from our algorithm, it produced a slightly lower F1 score of 0.47 - Figure 6.12b)), but more evenly distributed. Overall, the algorithm performance that we propose in this paper yields a higher precision in detecting disruptions from time series of traffic speeds than from the reported accident timeline. The use of multiple segments produced by the algorithm can highlight multiple disruptions while producing just a slight decrease in the quality of results. The error for multiple intervals segmentation increases because more additional intervals are considered in the evaluation of the metric, which may lay outside of originally marked intervals (see Figures 6.11 and 6.7).

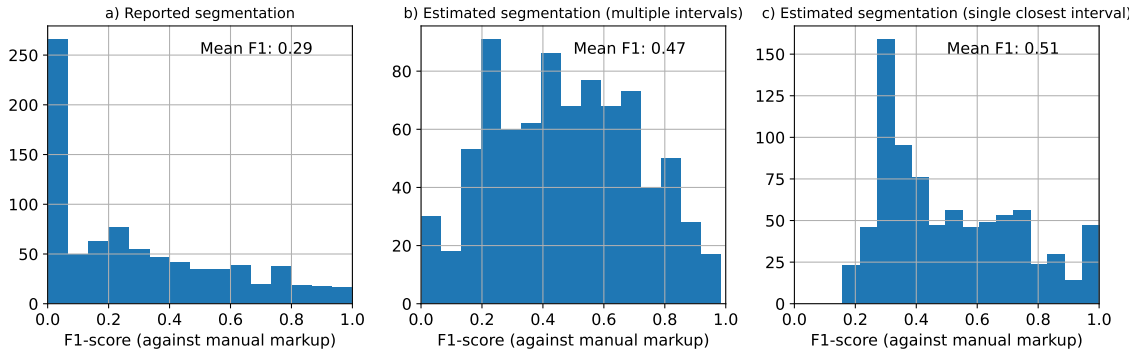


FIGURE 6.12: Histogram of F1-score against manual markup for a) reported accident time interval and b) estimated segmentation when algorithm detecting multiple disruption intervals c) estimated segmentation for the single closest interval to reported incident occurrence time

### 6.5.1 Parameter importance study

For our model, we have the following variables and their intervals of variation:

- **gran:** Granularity, an integer value controlling the level of detail (moving window size) in the metric estimation function. In the provided search space, the range of gran is [2, 40] with a step of 1. Default value is 12.

- **kernel\_size**: A list of float values used as weights in the dilation convolution operation. The search space for the kernel is the size of the convolution [1, ... 4], float values primarily intended to implement pre-processing operation for the day time-series. Default value is 3.
- **selectivity**: A float value between 0.01 and 4.0 that determines the power coefficient in post-processing difference estimations Default value is 2.0.
- **shift**: An integer value between -32 and +32 that represents a cyclic shift of the resulting time series to attribute to a shift in convolution operation and facilitate to overall adaptation to the target segmentation. Default value is 0.
- **threshold**: A float value that serves as a threshold in the interval processing function, which is used to perform the binarization of the normalized output array by disruption degree. In the provided search space, the range of the search space for the threshold is [0.01, 0.99]. Default value is 0.15.

At the begininning we perform a hyper-parameter search across all the mentioned parameters but also include a search among metric list (Bray-Curtis, Canberra, Chebyshev, Manhattan, Correlation, Cosine, Euclidean, Minkowski difference metrics) to determine the best performing difference metric for our algorithm. By performing search across 3,000 iterations we then estimate the avrage f1 score obtained when using each metric (see Figure 6.13). The Chebyshev metric yields higher f1 score than other metrics, possibly due to the structure and interpretation of the metric: high difference between maximum and minimum traffic speed measurements within a time window can indicate the presence of the disruption.

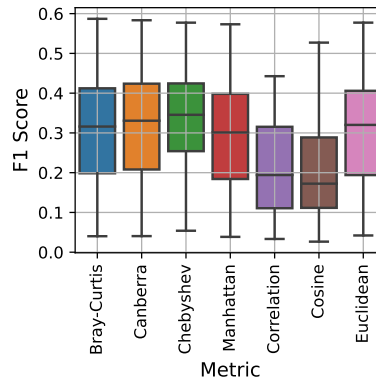


FIGURE 6.13: Connection between metric and F1 score

Our next step is to perform a hyper-parameter search for the Chebyshev difference metric only for 1,000 iterations. We obtained a significant improvement in the average f1 score for multiple interval comparison - 0.62 (a significant improvement from 0.52). As can be seen from scatter plots (see Figure 6.14), there are noticeable positive (kernel size vs f1 score), negative (binarization threshold vs f1 score) and peaking trends (shift vs f1 score) observed in results. Optimal values for the binarization threshold are located at lower values (between 0.01 and 0.4). Overall, the algorithm requires a positive shift in the post-processing function, which contributes to a substantial increase from 0.52 to 0.62 in f1 score when considering the positive shift of the resulting array.

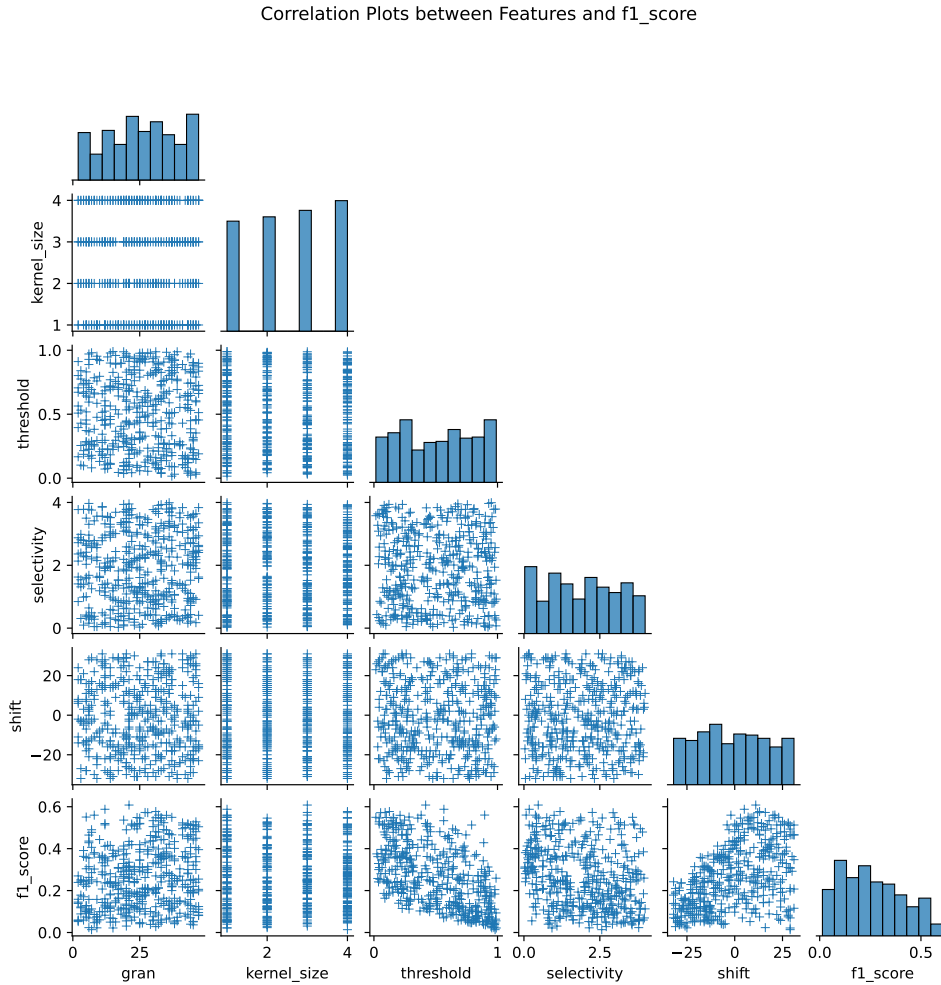


FIGURE 6.14: Scatter plots between model parameters and F1 score

We further provide a Correlation heatmap between algorithm parameters (see Figure 6.15) and the resulting f1 score: 1) The highest Pearson correlation values are with variables threshold (-0.55), shift (0.49), followed by selectivity (-0.32). There are no significant correlations between model parameters themselves.

In conclusion, the hyper-parameter search led to the selection of the Chebyshev metric, which demonstrated the highest average F1 score. Fine-tuning the disruption segmentation algorithm hyper-parameters significantly improved the average F1 score. Trends and optimal parameter values were identified, and the correlation heatmap showed that threshold, shift, and selectivity had the highest Pearson correlation with the F1 score.

## 6.6 Conclusion

Our methodology aims to automatically detect, segment, and extract traffic disruptions and accidents using distance metrics. This approach improves incident prediction accuracy across multiple machine learning models and provides better fit to manual markup of observed traffic speed disruptions. By obtaining the intervals and shapes of traffic disruptions, we can model the impact of accidents with greater precision, using traffic state measurements rather than just reported parameters (duration, start

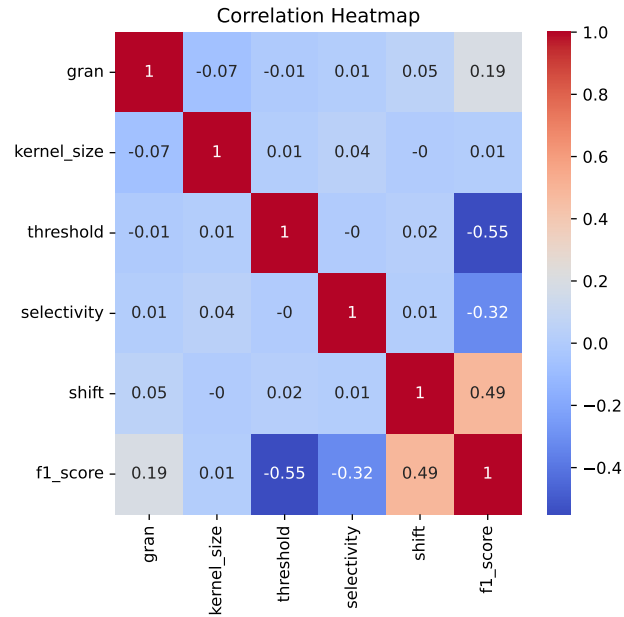


FIGURE 6.15: Scatter plot between model parameters and F1 score

time, etc). This approach provides more data on the accident and allows us to study accident impacts in greater detail.

**Relevance of this work can be summarized in following points:** 1) Enhancement of Traffic Management Systems: Integrate the proposed early detection and disruption segmentation algorithm into existing traffic management systems to improve and automate incident detection and corresponding data collection. This will help to minimize congestion and the overall impact of incidents on traffic flow, 2) highlight of reporting errors to standardize data reporting: Establish standardized guidelines and protocols for reporting traffic incidents, including the accurate reporting of the location, start and end times, number of lanes affected, and other relevant details; this will ensure that data-driven models can accurately predict incident severity and disruption length, 3) highlight the necessity of creating of data standards policies across countries for collecting necessary traffic accident information, 4) development of Incident Response Strategies by utilizing the improved incident prediction models to develop data-driven incident response strategies, including the dynamic traffic rerouting and real-time traffic guidance; this will help to mitigate the impact of traffic incidents on road users and reduce the risk of secondary incidents; 5) Data Fusion for a better traffic accident analysis: due to observed improvement in the quality of prediction arising from data fusion, traffic Authorities can consider integrating data sets from private companies for jointly analysing traffic datasets of various types to improve traffic safety by improving accuracy of traffic incident duration prediction.

**Future research in this area:** 1) Algorithm's complexity can be expanded by incorporating custom kernels, which can be found using hyper-parameter search, 2) Disruption measurements obtained over time can enable the prediction of traffic incident impact propagation with greater accuracy than relying solely on reported values, 3) The proposed methodology can be extended to include disruptions beyond accidents, such as construction or road closures, which can improve the accuracy of impact prediction, 4) Further improvement can also be achieved by performing data fusion and incorporating external data sources, such as weather and events, into the incident impact prediction models. We

are currently modelling the cascading effect on traffic disruptions and how these can be automatically identified based on multiple incoming traffic state streams; the main challenge of detecting subsequent incidents lie in the time-span duration of the first incident which is normally stochastic in nature.

**Limitations of this work:** The current modelling approach has been applied to a San Francisco data set due to its public availability and easiness to access. However, we would like to test the approach on multiple other countries and incident databases across the globe; the main challenge is the lack of both traffic states and traffic accidents logs to be released with synchronised timelines.

## 6.7 Conclusion

The proposed methodology in this paper aims at detecting, segmenting and extracting the observed disruptions in the traffic speed which was modelled together with reported traffic accidents by traffic management centers. The approach is innovative in its distance metric approach for an automatic incident detection coupled with an incident segmentation which has shown to improve the incident prediction by almost 41% in RMSE across multiple machine learning models. By obtaining shapes of disruptions we lay the foundation for accident impact modelling. Many studies still rely on the modelling of reported accident durations and pre-defined parameters, while they can be estimated from traffic state measurements, which gives us more data than just aggregated variables (duration, start time, etc). By having information on how each accident affect the traffic flow, we can study the accident impact with precision.

**Limitations of this work:** The current modelling approach has been applied to a San Francisco data set due to its public availability and easiness to access. However, we would like to test the approach on multiple other countries and incident databases across the globe; the main challenge is the lack of both traffic states and traffic accidents logs to be released with synchronised timelines. **Future works:** We are currently modelling the cascading effect on traffic disruptions and how these can be automatically identified based on multiple incoming traffic state streams; the main challenge of detecting subsequent incidents lie in the time-span duration of the first incident which is normally stochastic in nature.



## Chapter 7

# Discussion, Synthesis and Conclusions

### 7.1 Literative review: discussion, synthesis and conclusions

**Summary of the main findings:** This research identified a significant gap in both the Australian transport management sector and the academic literature with regard to the implementation of advanced methods in incident management and response plan solutions. Current practices largely rely on operational experience rather than data-driven decision-making. The research also highlights the importance of using deep learning and machine learning in traffic incident management systems (TIMS) to accurately predict incident durations and identify the most critical factors influencing incident clearance time. This approach has the potential to save operational costs, reduce end-user time, and decrease traffic congestion. This review also identified several challenges and gaps in the field of traffic accident analysis including imbalanced classification, skewed distribution of accident duration data, reporting errors and anomalies, high-dimensionality of the data, low data availability, incorporation of textual accident descriptions, utilization of historical traffic data, the use of novel machine learning and deep learning models, and incident-related traffic state identification.

**Discussion of the implications:** The findings of this study have several implications for the field of traffic management and incident response. First, they emphasize the need for a shift from experience-based decision-making to data-driven approaches in transport management centres. This will enable more efficient incident management and the allocation of resources. Second, the application of deep learning and machine learning in TIMS can lead to more accurate predictions of incident durations across different road types, accident types, and countries with varying driving behaviour. This enhanced prediction capability will contribute to existing knowledge and inform the development of future research and traffic management strategies.

**Limitations and future research:** The research has some limitations, such as the focus on Australian traffic management centres and the lack of exploration of all possible modelling capabilities. Future research should investigate a more comprehensive approach to incident duration prediction, encompassing various road types and accident types across different countries. Additionally, further studies should aim to improve the integration of transport modelling and data-driven solutions, utilizing the full potential of deep learning and machine learning in TIMS.

Based on the identified challenges and gaps, several potential future research directions have been proposed:

- **Data set integration and fusion models:** Combining data sets such as traffic flow, speed, and occupancy with traffic accident reports can enhance incident duration modelling. This may

require the use of data fusion models or feature embedding methods.

- Utilization of textual data: Incorporating natural language processing techniques to analyze textual accident reports can provide valuable insights and improve prediction accuracy.
- Application of advanced machine learning and deep learning models: Exploring more sophisticated ML and DL models can help identify nonlinear relationships and threshold effects in traffic incident duration prediction.
- Integration of advanced ML pipeline elements: Anomaly detection, hyperparameter optimization, dimensionality reduction, and sampling techniques can enhance prediction performance.
- Real-time incident reporting text analysis: Examining the timeline of textual incident descriptions using available data sets like PeMS.
- Standardization of accident reporting form based on data-driven approaches utilizing the feature importance estimation techniques: assessing the impact of specific factors, such as weather conditions, on incident duration prediction accuracy, detailing and choosing features that have the highest contribution to the prediction accuracy.
- Road Type extrapolation tests and model bias considerations: Ensuring the model's performance remains reliable when extrapolating data to other road networks (e.g. cross-network test) or when applied to different time and space contexts (e.g. time-based cross-validation).
- Advanced data pre-processing methods: Implementing dimensionality reduction and feature extraction techniques (e.g. using autoencoder) to manage the growing volume and variety of data collected in traffic networks.

In conclusion, traffic incident duration prediction is a complex and important task that can benefit from further research involving sophisticated artificial intelligence models. By addressing the identified challenges and gaps, future research has the potential to significantly improve traffic incident duration prediction performance, ultimately leading to enhanced traffic flow and reduced impact from traffic incidents. The study highlights the need for more advanced incident management solutions that leverage deep learning and machine learning techniques to accurately predict traffic incident durations and identify the most important factors affecting incident clearance time. Implementing such data-driven approaches will result in better resource allocation and improved traffic management, ultimately benefiting end-users and society as a whole. Further research should explore broader applications of these techniques in TIMS and investigate their potential in various road types, accident types, and countries.

## 7.2 Bi-level framework: discussion, synthesis and conclusions

### Summary of the main Findings:

The work around a bi-level framework for traffic incident duration prediction presented a universal bi-level framework that addresses several challenges for different road network layouts. The study proposes a framework capable of predicting incident duration regardless of the road network or its



complexity. It addresses the issues of outliers and imbalanced data classes by proposing a varying threshold procedure, optimizes both classification and regression problems, and highlights the most influential factors that affect the incident duration. The research demonstrates that the performance of machine learning models is highly affected by the dataset and the chosen methodology, emphasizing the need for a flexible and adaptable approach. This chapter lays the groundwork for bi-level predictive methodologies regarding traffic incident duration, ultimately aiding incident modelling.

#### **Discussion of the implications:**

The proposed bi-level framework for traffic incident duration prediction has significant implications for the field of study. It fills a gap in the literature by providing a universal framework that can be applied to different traffic incident datasets and various road network types. This approach offers a more comprehensive solution to the problem of incident duration prediction, considering varying incident duration threshold analysis, joint hyper-parameter optimization algorithms, and feature importance selection. By identifying the most influential factors that affect incident duration across three different types of road networks, this research can help traffic authorities prioritize their efforts and improve their decision-making processes. The framework's adaptability also allows for more accurate predictions and better-informed decisions.

#### **Limitations and Future Research:**

While the work in this chapter addressed several challenges in predicting traffic incident duration, it acknowledges some limitations. The performance of machine learning models is highly dependent on the quality and size of the dataset and the chosen methodology. Future research could focus on exploring more advanced machine learning techniques to improve model performance further. Additionally, more extensive and diverse datasets could be used to test the framework and further validate its applicability to various road network types and incident scenarios. Incorporating real-time traffic data and dynamic road conditions could also enhance the prediction accuracy and better reflect real-world complexities.

#### **Conclusions:**

Overall the work around the bi-level framework modeling contributes to the ongoing development of a real-time platform for predicting traffic congestion and evaluating the incident impact during peak hours. The proposed bi-level framework for traffic incident duration prediction is a significant advancement in the field, offering a flexible, adaptable, and comprehensive solution. By addressing the challenges of predicting incident duration on different road network layouts and accounting for various influential factors.

### **7.3 Data fusion for traffic incident duration prediction: discussion, synthesis and conclusions**

**Summary of the main findings:** This chapter proposed a novel framework for predicting incident duration by integrating machine learning prediction methods with traffic flow and textual incident description features encoded via several Deep Learning methods. The approach showed stable and significant improvement across all models. Our research also highlighted the importance of using specific deep-learning encoding approaches for regression models to further enhance performance. The study revealed that encoding incident-related features efficiently is crucial for predicting traffic

incident impacts. We also investigated the importance of words in incident descriptions using the LIME method, which showed that certain word combinations contribute to the classification of incidents into specific severity or duration groups.

**Discussion of the implications:** This research contributes to the ongoing objective of building a real-time platform for predicting traffic congestion and evaluating incident impacts during peak hours. The proposed framework can provide accurate information for both end-user route choice modeling and operational centers looking to optimize their operations under non-recurrent traffic congestion. Furthermore, the study lays the foundation for bi-level predictive methodologies concerning traffic incident duration, which can be beneficial for Traffic Management Centers (TMCs) in improving incident and traffic management.

**Limitations and future research:** The study has some limitations, including focusing only on the San Francisco area and considering traffic speed and flow only one week before the incident. Future research could incorporate more extensive geographical areas, like California, and longer periods of traffic count data to build traffic speed/flow profiles for more accurate predictions. Also, additional methods of time series encoding may be utilized. Further work could also explore the spatial and temporal dynamic prediction of incident impact using graph-based modelling approaches. The availability of the predicted incident duration data integrated with data on traffic flow and textual incident description can improve the TMC incident and traffic management, reducing the time that people spend in traffic congestion caused by incidents.

## 7.4 Visual transformers for traffic accident risk prediction: discussion, synthesis and conclusions

**Summary of the main findings:** This research introduces a novel approach to traffic accident risk forecasting by reformulating the problem as an image regression task and proposing a unique Contextual Vision Transformer network (C-ViT) that efficiently models traffic accident risk from both spatial and temporal perspectives. The proposed approach outperforms existing methods, requiring significantly fewer training parameters. Additionally, incorporating a static accident risk map with the ViT model (XViT) further improves performance, establishing a new state-of-the-art. The Coarse-Fine-Coarse Visual Transformer (CFC-ViT) architecture allows for fine-grained processing of the accident risk map and introduces an additional scale factor parameter, which can enhance prediction performance.

**Discussion of the implications:** The findings of this study demonstrate the potential of visual transformers and their variations for traffic accident risk prediction, surpassing previous approaches. This research highlights the applicability of vision transformers for non-visual tasks and suggests that further applications of image and video processing methods may yield even better results and open alternative approaches for accident risk prediction. The proposed methods can contribute to more accurate and efficient traffic management, accident prevention, and policy-making in urban environments.

**Limitations and future research:** This study has several limitations, such as the potential for improvement in operation combination methods and constraint functions. There may also be a non-linear dependence between RMSE and the scale factor observed for different datasets, suggesting that

an optimal scale factor for accident risk map processing may exist and vary between datasets. Future research can explore alternative combination methods and constraint functions, investigate the optimal scale factor for various datasets, and apply image and video processing methods to further improve accident risk prediction.

**Conclusions:** This research presents a novel approach for traffic accident risk forecasting using visual transformers and their variations, outperforming existing methods. By incorporating static accident risk maps and Coarse-Fine-Coarse Visual Transformer architectures, the proposed methods show significant improvement in prediction performance. These findings can contribute to better traffic management, accident prevention, and policy-making, while also opening new avenues for applying vision transformers to non-visual tasks and exploring image and video processing methods for accident risk prediction.

## 7.5 Accident Segmentation: discussion, synthesis and conclusions

**Summary of the main findings:** This research focused on addressing the challenges in traffic accident analysis due to incorrect or incomplete accident reports and the need for accurate traffic disruption segmentation. We proposed novel methods for traffic disruption segmentation and association between vehicle detector stations and accident reports. We also developed a fusion methodology for combining two large datasets, CTADS and PeMS, to analyze the relationship between traffic accidents and their effects on traffic flow and speed. Through the evaluation of multiple machine learning models, we introduced a new modeling approach that focuses on the amount and shape of the disruption associated with an accident. This research lays the foundation for early traffic accident disruption detection, traffic disruption speed impact analysis, and the use of observed traffic accident durations for correcting errors in user reports.

**Discussion of the implications:** Our findings have significant implications for the field of traffic accident analysis and prediction. By accurately segmenting and analyzing traffic disruptions, we can better understand the impact of accidents on traffic flow and speed. This can lead to improved traffic incident management and more effective allocation of resources in response to accidents. Furthermore, the fusion of two large datasets enables the investigation of traffic accidents across different countries and traffic conditions, contributing to a more comprehensive understanding of the factors influencing accident duration and impact. Our research also provides a foundation for the development of real-time platforms for predicting traffic congestion and evaluating incident impact.

**Limitations and future research:** This research has some limitations, such as the reliance on two large datasets which may not be representative of all traffic conditions worldwide. Additionally, the fusion methodology may not be applicable to all data sources, and the machine learning models tested may not be optimal for all situations. Future research should focus on expanding the analysis to more diverse datasets and investigating other machine learning models for improved prediction performance. Moreover, further research could explore the spatial-temporal impact of disruptions within the traffic network and analyze the influence of various accident characteristics on traffic flow patterns.

**Conclusions:** In conclusion, this research contributes significantly to the field of traffic accident analysis by addressing challenges related to data quality and segmentation, proposing novel methods

for traffic disruption analysis, and evaluating the performance of multiple machine learning models. The findings of this research have important implications for traffic incident management, resource allocation, and the development of real-time platforms for predicting traffic congestion and incident impact. Further research is needed to refine the methodologies and expand their applicability to a broader range of traffic conditions and datasets.

## 7.6 Final thesis Conclusion

In conclusion, this thesis demonstrates the potential for leveraging advanced machine learning, deep learning, and artificial intelligence techniques to improve traffic incident duration prediction, accident risk forecasting, and overall traffic management. By utilizing novel approaches such as bi-level frameworks, contextual vision transformers, estimation of observed incident duration via time series segmentation and integrating deep learning methods with traffic flow and description features, these studies contribute to the development of more accurate and efficient traffic management systems.

These advancements have important societal implications, as improved incident prediction and management can lead to reduced congestion, more efficient allocation of resources, and better-informed policy-making. While there are limitations, there are more areas for future research, such as exploring broader geographical areas, optimizing methodologies, and further investigating the applications of vision transformers for non-visual tasks, these studies lay the groundwork for the continued development and application of cutting-edge techniques in the field of traffic management and incident response.

# Bibliography

- Abduljabbar, Rusul et al. (2019). “Applications of artificial intelligence in transport: An overview”. In: *Sustainability* 11.1, p. 189.
- Abou El Assad, Zouhair Elamrani, Hajar Mousannif, and Hassan Al Moatassime (2020). “A real-time crash prediction fusion framework: An imbalance-aware strategy for collision avoidance systems”. In: *Transportation research part C: emerging technologies* 118, p. 102708.
- Adler, Martin W., Jos van Ommeren, and Piet Rietveld (2013). “Road congestion and incident duration”. In: *Economics of Transportation* 2.4, pp. 109–118. ISSN: 2212-0122. DOI: <https://doi.org/10.1016/j.ecotra.2013.12.003>. URL: <http://www.sciencedirect.com/science/article/pii/S2212012213000269>.
- Administration, National Highway Traffic Safety (2013). *Traffic safety facts 2013*. U.S. department of transportation.
- Al-Bordiny, Mohamed Abdullah Mohamed Ahmed (2014). “Function, Requirements and Applications of Intelligent Transportation Systems (ITS)”. PhD thesis. Faculty of Engineering, Tanta University.
- Al Hamami, M and TC Matisziw (2021). “Measuring the spatiotemporal evolution of accident hot spots”. In: *Accident Analysis and Prevention* 157, p. 106133.
- Al-Najada, Hamzah and Imad Mahgoub (2017). “Real-time incident clearance time prediction using traffic data from internet of mobility sensors”. In: *2017 IEEE 15th Intl Conf on Dependable, Autonomic and Secure Computing, 15th Intl Conf on Pervasive Intelligence and Computing, 3rd Intl Conf on Big Data Intelligence and Computing and Cyber Science and Technology Congress (DASC/PiCom/DataCom/CyberSciTech)*. IEEE, pp. 728–735.
- Ali, Farman et al. (2019). “Transportation sentiment analysis using word embedding and ontology-based topic modeling”. In: *Knowledge-Based Systems* 174, pp. 27–42. ISSN: 0950-7051. DOI: <https://doi.org/10.1016/j.knosys.2019.02.033>. URL: <https://www.sciencedirect.com/science/article/pii/S0950705119300942>.
- Ali, Farman et al. (2021). “Traffic accident detection and condition analysis based on social networking data”. In: *Accident Analysis and Prevention* 151, p. 105973. ISSN: 0001-4575. DOI: <https://doi.org/10.1016/j.aap.2021.105973>. URL: <https://www.sciencedirect.com/science/article/pii/S000145752100004X>.
- Alkaabi, Abdulla Mohammed Saeed, Dilum Dissanayake, and Roger Bird (2011). “Analyzing clearance time of urban traffic accidents in Abu Dhabi, United Arab Emirates, with hazard-based duration modeling method”. In: *Transportation Research Record* 2229.1, pp. 46–54.
- Alkheder, Sharaf, Madhar Taamneh, and Salah Taamneh (2017). “Severity prediction of traffic accident using an artificial neural network”. In: *Journal of Forecasting* 36.1, pp. 100–108.

- Araghi, Bahar N et al. (2014). "A comparative study of k-NN and hazard-based models for incident duration prediction". In: *17th International IEEE Conference on Intelligent Transportation Systems (ITSC)*. IEEE, pp. 1608–1613.
- Azevedo, Tiago et al. (2016). "JADE, TraSMAP and SUMO: A tool-chain for simulating traffic light control". In: *arXiv preprint arXiv:1601.08154*.
- Ballings, Michel et al. (2015). "Evaluating multiple classifiers for stock price direction prediction". In: *Expert Systems with Applications* 42.20, pp. 7046–7056. ISSN: 0957-4174. DOI: <https://doi.org/10.1016/j.eswa.2015.05.013>. URL: <https://www.sciencedirect.com/science/article/pii/S0957417415003334>.
- Barcellos, Pablo et al. (2015). "A novel video based system for detecting and counting vehicles at user-defined virtual loops". In: *Expert Systems with Applications* 42.4, pp. 1845–1856.
- Barthélemy, Marc (2011). "Spatial networks". In: *Physics reports* 499.1-3, pp. 1–101.
- Behrooz, Hojat and Yeganeh M Hayeri (2022). "Machine Learning Applications in Surface Transportation Systems: A Literature Review". In: *Applied Sciences* 12.18, p. 9156.
- Bekkerman, R (2015). *The present and the future of the KDD cup competition: an outsider's perspective*.
- Benterki, Abdelmoudjib et al. (2020). "Artificial intelligence for vehicle behavior anticipation: Hybrid approach based on maneuver classification and trajectory prediction". In: *IEEE Access* 8, pp. 56992–57002.
- Bergstra, James and Y. Bengio (Mar. 2012). "Random Search for Hyper-Parameter Optimization". In: *The Journal of Machine Learning Research* 13, pp. 281–305.
- Breiman, Leo (2001). "Random Forests". In: *Mach. Learn.* 45.1, 5–32. ISSN: 0885-6125. DOI: [10.1023/A:1010933404324](https://doi.org/10.1023/A:1010933404324). URL: <https://doi.org/10.1023/A:1010933404324>.
- Breunig, Markus et al. (June 2000). "LOF: Identifying Density-Based Local Outliers." In: vol. 29, pp. 93–104. DOI: [10.1145/342009.335388](https://doi.org/10.1145/342009.335388).
- Bridgelall, Raj and Denver D Tolliver (2021). "Railroad accident analysis using extreme gradient boosting". In: *Accident Analysis and Prevention* 156, p. 106126.
- California Highway Patrol (CHP) (n.d.). *Statewide Integrated Traffic Records System (SWITRS)*. URL: <https://www.chp.ca.gov/switrs/> (visited on 02/02/2021).
- Cao, Donglin, Shuru Wang, and Dazhen Lin (2018). "Chinese Microblog Users' Sentiment-Based Traffic Condition Analysis". In: *Soft Comput.* 22.21, 7005–7014. ISSN: 1432-7643. DOI: [10.1007/s00500-018-3293-8](https://doi.org/10.1007/s00500-018-3293-8). URL: <https://doi.org/10.1007/s00500-018-3293-8>.
- Chang, Hsin-li and Tse-pin Chang (2013). "Prediction of freeway incident duration based on classification tree analysis". In: *Journal of the Eastern Asia Society for Transportation Studies* 10, pp. 1964–1977.
- Chawla, Nitesh V et al. (2002). "SMOTE: synthetic minority over-sampling technique". In: *Journal of artificial intelligence research* 16, pp. 321–357.
- Chen, Chao et al. (2001). "Freeway performance measurement system: mining loop detector data". In: *Transportation Research Record* 1748.1, pp. 96–102.
- Chen, Shun-Yuan, Yung-Chuan Chen, and Cheng-Hao Hsieh (2016). "Learning Traffic Accident Risk Prediction using Deep Learning". In: *IEEE Transactions on Intelligent Transportation Systems* 17.12, pp. 3409–3418.

- Chen, Tianqi and Carlos Guestrin (2016). “Xgboost: A scalable tree boosting system”. In: *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794.
- Chen, Tianqi et al. (2015). “Xgboost: extreme gradient boosting”. In: *R package version 0.4-2 1.4*.
- Chen, Tianyi et al. (2020). “Predicting lane-changing risk level based on vehicles’ space-series features: A pre-emptive learning approach”. In: *Transportation research part C: emerging technologies* 116, p. 102646.
- Choe, Tom, Alexander Skabardonis, and Pravin Varaiya (2002). “Freeway performance measurement system: operational analysis tool”. In: *Transportation research record* 1811.1, pp. 67–75.
- Chung, Yi-Shih, Yu-Chiun Chiou, and Chia-Hua Lin (2015). “Simultaneous equation modeling of freeway accident duration and lanes blocked”. In: *Analytic Methods in Accident Research* 7, pp. 16–28. ISSN: 2213-6657. DOI: <https://doi.org/10.1016/j.amar.2015.04.003>. URL: <http://www.sciencedirect.com/science/article/pii/S2213665715000275>.
- Chung, Younshik (2010). “Development of an accident duration prediction model on the Korean Freeway Systems”. In: *Accident Analysis and Prevention* 42.1, pp. 282–289.
- Chung, Younshik and Wilfred W Recker (2012). “A methodological approach for estimating temporal and spatial extent of delays caused by freeway accidents”. In: *IEEE Transactions on Intelligent Transportation Systems* 13.3, pp. 1454–1461.
- Chung, Younshik, Lubinda Walubita, and Keechoo Choi (Jan. 2011). “Modeling Accident Duration and Its Mitigation Strategies on South Korean Freeway Systems”. In: *Transportation Research Record Journal of the Transportation Research Board* 2178, pp. 49–57. DOI: [10.3141/2178-06](https://doi.org/10.3141/2178-06).
- Chung, Younshik, Lubinda F Walubita, and Keechoo Choi (2010). “Modeling accident duration and its mitigation strategies on South Korean freeway systems”. In: *Transportation research record* 2178.1, pp. 49–57.
- Compass IoT (n.d.). *Compass IoT*. URL: <https://www.compassiot.com.au/> (visited on 02/02/2021).
- Curiel, Rafael Prieto, Humberto Gonzalez Ramirez, and Steven Richard Bishop (2018). “A novel rare event approach to measure the randomness and concentration of road accidents”. In: *PloS one* 13.8.
- Das, Soumitra, Sachin Mohanty, and Pushpak Bhattacharyya (2019). “Topical text modeling for incident detection on twitter”. In: *Expert Systems with Applications* 122, pp. 182–195.
- Department of Infrastructure Regional Development and Cities (n.d.). *Australian Road Deaths Database - ARDD*. URL: [https://www.bitre.gov.au/statistics/safety/fatal\\\_road\\\_crash\\\_database](https://www.bitre.gov.au/statistics/safety/fatal\_road\_crash\_database) (visited on 02/02/2021).
- Dietterich, Thomas G. (2000). “Ensemble Methods in Machine Learning”. In: *Multiple Classifier Systems*. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 1–15. ISBN: 978-3-540-45014-6.
- Djenouri, Youcef et al. (2019). “A survey on urban traffic anomalies detection algorithms”. In: *IEEE Access* 7, pp. 12192–12205.
- Eboli, Laura, Carmen Forciniti, and Gabriella Mazzulla (2020). “Factors influencing accident severity: an analysis by road accident type”. In: *Transportation Research Procedia* 47. 22nd EURO Working Group on Transportation Meeting, EWGT 2019, 18th – 20th September 2019, Barcelona, Spain, pp. 449–456. ISSN: 2352-1465. DOI: <https://doi.org/10.1016/j.trpro.2020.03.120>. URL: <https://www.sciencedirect.com/science/article/pii/S2352146520303197>.

- European Commission (n.d.). *Road Safety Atlas*. URL: [https://ec.europa.eu/transport/road/\\_safety/specialist/library/atlas/index\\_en.htm](https://ec.europa.eu/transport/road/_safety/specialist/library/atlas/index_en.htm) (visited on 02/02/2021).
- Farahani, Mehrdad et al. (2020). “Short-term traffic flow prediction using variational LSTM networks”. In: *arXiv preprint arXiv:2002.07922*.
- Fix, Evelyn and JL Hodges (1951). “Discriminatory analysis, nonparametric discrimination”. In: François, Damien, Vincent Wertz, and Michel Verleysen (2011). “Choosing the metric: a simple model approach”. In: *Meta-Learning in Computational Intelligence*, pp. 97–115.
- Friedman, Jerome (Nov. 2000). “Greedy Function Approximation: A Gradient Boosting Machine”. In: *The Annals of Statistics* 29. DOI: [10.1214/aos/1013203451](https://doi.org/10.1214/aos/1013203451).
- Fukuda, Shota et al. (2020). “Short-term prediction of traffic flow under incident conditions using graph convolutional recurrent neural network and traffic simulation”. In: *IET Intelligent Transport Systems* 14.8, pp. 936–946.
- Ghasri, Milad et al. (2016). “Hazard-based model for concrete pouring duration using construction site and supply chain parameters”. In: *Automation in Construction* 71, pp. 283–293.
- Gholami, Hamid et al. (2020). “Mapping wind erosion hazard with regression-based machine learning algorithms”. In: *Scientific Reports* 10.1, p. 20494.
- Ghosh, Banishree (2019). “Predicting the duration and impact of the non-recurring road incidents on the transportation network”. PhD thesis.
- Ghosh, Banishree and Justin Dauwels (2022). “Comparison of different Bayesian methods for estimating error bars with incident duration prediction”. In: *Journal of Intelligent Transportation Systems* 26.4, pp. 420–431.
- Ghosh, Banishree et al. (2016). “Predicting the duration of non-recurring road incidents by cluster-specific models”. In: *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, pp. 1522–1527.
- Ghosh, Banishree et al. (2018). “Dynamic prediction of the incident duration using adaptive feature set”. In: *IEEE Transactions on Intelligent Transportation Systems* 20.11, pp. 4019–4031.
- Government, Australian (2017). *Road Safety*. URL: <https://infrastructure.gov.au/roads/safety/>.
- Grigorev, Artur et al. (2022a). “Incident duration prediction using a bi-level machine learning framework with outlier removal and intra-extra joint optimisation”. In: *Transportation research part C: emerging technologies* 141, p. 103721.
- (2022b). “Incident duration prediction using a bi-level machine learning framework with outlier removal and intra-extra joint optimisation”. In: *Transportation Research Part C: Emerging Technologies* 141, p. 103721. ISSN: 0968-090X. DOI: <https://doi.org/10.1016/j.trc.2022.103721>. URL: <https://www.sciencedirect.com/science/article/pii/S0968090X22001589>.
- Grigorev, Artur et al. (2022c). “Traffic incident duration prediction via a deep learning framework for text description encoding”. In: *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, pp. 1770–1777.
- Guyon, Isabelle and André Elisseeff (2003). “An introduction to variable and feature selection”. In: *Journal of machine learning research* 3.Mar, pp. 1157–1182.



- Hamad, Khaled, Mohamad Ali Khalil, and Abdul Razak Alozi (2019). "Predicting Freeway Incident Duration Using Machine Learning". In: *International Journal of Intelligent Transportation Systems Research*, pp. 1–14.
- (2020). "Predicting freeway incident duration using machine learning". In: *International Journal of Intelligent Transportation Systems Research* 18, pp. 367–380.
- Hamad, Khaled et al. (2020a). "Predicting incident duration using random forests". In: *Transportmetrica A: transport science* 16.3, pp. 1269–1293.
- (2020b). "Predicting incident duration using random forests". In: *Transportmetrica A: transport science* 16.3, pp. 1269–1293.
- Haule, Henrick J. et al. (2019a). "Evaluating the impact and clearance duration of freeway incidents". In: *International Journal of Transportation Science and Technology* 8.1, pp. 13–24. ISSN: 2046-0430. DOI: <https://doi.org/10.1016/j.ijtst.2018.06.005>. URL: <http://www.sciencedirect.com/science/article/pii/S2046043018300522>.
- Haule, Henrick J et al. (2019b). "Evaluating the impact and clearance duration of freeway incidents". In: *International journal of transportation science and technology* 8.1, pp. 13–24.
- He, Haibo et al. (July 2008). "ADASYN: Adaptive Synthetic Sampling Approach for Imbalanced Learning". In: pp. 1322–1328. DOI: [10.1109/IJCNN.2008.4633969](https://doi.org/10.1109/IJCNN.2008.4633969).
- He, Qing et al. (2013). "Incident duration prediction with hybrid tree-based quantile regression". In: *Advances in dynamic network modeling in complex transportation systems*. Springer, pp. 287–305.
- Heywood, Peter, Paul Richmond, and Steve Maddock (2015). "Road network simulation using FLAME GPU". In: *European Conference on Parallel Processing*. Springer, pp. 430–441.
- Heywood, Peter et al. (2018). "Data-parallel agent-based microscopic road network simulation using graphics processing units". In: *Simulation Modelling Practice and Theory* 83, pp. 188–200.
- Hojati, Ahmad et al. (June 2012). "Analysing freeway traffic incident duration using an Australian data set". In: *Road and Transport Research* 21, pp. 16–28.
- Hojati, Ahmad Tavassoli et al. (2013). "Hazard based models for freeway traffic incident duration". In: *Accident Analysis and Prevention* 52, pp. 171–181.
- Hojati, Ahmad [Tavassoli et al. (2014). "Modelling total duration of traffic incidents including incident detection and recovery time". In: *Accident Analysis and Prevention* 71, pp. 296–305. ISSN: 0001-4575. DOI: <https://doi.org/10.1016/j.aap.2014.06.006>. URL: <http://www.sciencedirect.com/science/article/pii/S0001457514001791>.
- Hong, Sungmin et al. (Jan. 2014). "The Effect of Road Environment Factors on Freeway Traffic Crash Frequency during Daylight, Twilight, and Night Conditions". In.
- Hossain, Moinul et al. (2019). "Real-time crash prediction models: State-of-the-art, design pathways and ubiquitous requirements". In: *Accident Analysis and Prevention* 124, pp. 66–84.
- Hou, Lin et al. (2013). "Modeling freeway incident response time: A mechanism-based approach". In: *Transportation Research Part C: Emerging Technologies* 28. Euro Transportation: selected paper from the EWGT Meeting, Padova, September 2009, pp. 87–100. ISSN: 0968-090X. DOI: <https://doi.org/10.1016/j.trc.2012.12.005>. URL: <http://www.sciencedirect.com/science/article/pii/S0968090X12001519>.
- Huang, Bo and Xiaohong Pan (2007). "GIS coupled with traffic simulation and optimization for incident response". In: *Computers, environment and urban systems* 31.2, pp. 116–132.

- Huang, Guang-Bin, Dian Hui Wang, and Yuan Lan (2011). "Extreme learning machines: a survey". In: *International journal of machine learning and cybernetics* 2.2, pp. 107–122.
- Huang, Tingting, Shuo Wang, and Anuj Sharma (2020). "Highway crash detection and risk estimation using deep learning". In: *Accident Analysis and Prevention* 135, p. 105392.
- Huang, Xiaohui et al. (2020). "Intelligent intersection: Two-stream convolutional networks for real-time near-accident detection in traffic video". In: *ACM Transactions on Spatial Algorithms and Systems (TSAS)* 6.2, pp. 1–28.
- Javid, Roxana J and Ramina Jahanbakhsh Javid (2018). "A framework for travel time variability analysis using urban traffic incident data". In: *IATSS research* 42.1, pp. 30–38.
- Kalair, Kieran and Colm Connaughton (2021). "Dynamic and interpretable hazard-based models of traffic incident durations". In: *Frontiers in future transportation* 2, p. 669015.
- Ke, Guolin et al. (2017). "Lightgbm: A highly efficient gradient boosting decision tree". In: *Advances in neural information processing systems* 30, pp. 3146–3154.
- Khattak, Asad J, Joseph L Schofer, and Mu-Han Wang (1995). "A simple time sequential procedure for predicting freeway incident duration". In: *Journal of Intelligent Transportation Systems* 2.2, pp. 113–138.
- Khattak, Asad J et al. (2016). "Modeling traffic incident duration using quantile regression". In: *Transportation Research Record* 2554.1, pp. 139–148.
- KIM, Hyung Jin and Hoi-Kyun CHOI (2001). "A Comparative Analysis of Incident Service Time on Urban Freeways". In: *IATSS Research* 25.1, pp. 62–72. ISSN: 0386-1112. DOI: [https://doi.org/10.1016/S0386-1112\(14\)60007-8](https://doi.org/10.1016/S0386-1112(14)60007-8). URL: <http://www.sciencedirect.com/science/article/pii/S0386111214600078>.
- Kim, Woon and Gang-Len Chang (Jan. 2011). "Development of a Hybrid Prediction Model for Freeway Incident Duration: A Case Study in Maryland". In: *International Journal of Intelligent Transportation Systems Research* 10. DOI: [10.1007/s13177-011-0039-8](https://doi.org/10.1007/s13177-011-0039-8).
- (2012). "Development of a hybrid prediction model for freeway incident duration: a case study in Maryland". In: *International journal of intelligent transportation systems research* 10, pp. 22–33.
- Knapen, Luk et al. (2014). "Within day rescheduling microsimulation combined with macro-simulated traffic". In: *Transportation Research Part C: Emerging Technologies* 45, pp. 99–118.
- Kuang, Li et al. (2019a). "Predicting duration of traffic accidents based on cost-sensitive Bayesian network and weighted K-nearest neighbor". In: *Journal of Intelligent Transportation Systems* 23.2, pp. 161–174.
- (2019b). "Predicting duration of traffic accidents based on cost-sensitive Bayesian network and weighted K-nearest neighbor". In: *Journal of Intelligent Transportation Systems* 23.2, pp. 161–174.
- Lee, Ying and Chien-Hung Wei (2010a). "A computerized feature selection method using genetic algorithms to forecast freeway accident duration times". In: *Computer-Aided Civil and Infrastructure Engineering* 25.2, pp. 132–148.
- (2010b). "A computerized feature selection method using genetic algorithms to forecast freeway accident duration times". In: *Computer-Aided Civil and Infrastructure Engineering* 25.2, pp. 132–148.

- Li, Jie et al. (2018). "An overview of graph embedding: Problems, techniques and applications". In: *ACM Transactions on Knowledge Discovery from Data (TKDD)* 12.5, pp. 1–36.
- Li, Linchao et al. (2017). "Traffic incident detection based on extreme machine learning". In: *Journal of Applied Science and Engineering* 20.4, pp. 409–416.
- Li, Linchao et al. (2020a). "A deep fusion model based on restricted Boltzmann machines for traffic accident duration prediction". In: *Engineering Applications of Artificial Intelligence* 93, p. 103686.
- li, Ruimin (Jan. 2014). "Traffic incident duration analysis and prediction models based on the survival analysis approach". In: *IET Intelligent Transport Systems* 9. DOI: [10.1049/iet-its.2014.0036](https://doi.org/10.1049/iet-its.2014.0036).
- Li, Ruimin (2015). "Traffic incident duration analysis and prediction models based on the survival analysis approach". In: *IET intelligent transport systems* 9.4, pp. 351–358.
- Li, Ruimin, Francisco C Pereira, and Moshe E Ben-Akiva (2015a). "Competing risk mixture model and text analysis for sequential incident duration prediction". In: *Transportation Research Part C: Emerging Technologies* 54, pp. 74–85.
- (2015b). "Competing risks mixture model for traffic incident duration prediction". In: *Accident Analysis and Prevention* 75, pp. 192–201.
- (2015c). "Competing risks mixture model for traffic incident duration prediction". In: *Accident Analysis and Prevention* 75, pp. 192–201.
- (2018a). "Overview of traffic incident duration analysis and prediction". In: *European transport research review* 10.2, p. 22.
- (2018b). "Overview of traffic incident duration analysis and prediction". In: *European transport research review* 10.2, pp. 1–13.
- Li, Ruimin and Pan Shang (Nov. 2014a). "Incident Duration Modeling Using Flexible Parametric Hazard-Based Models". In: *Computational intelligence and neuroscience* 2014, p. 723427. DOI: [10.1155/2014/723427](https://doi.org/10.1155/2014/723427). URL: <http://dx.doi.org/10.1155/2014/723427>.
- (2014b). "Incident duration modeling using flexible parametric hazard-based models". In: *Computational intelligence and neuroscience* 2014, pp. 33–33.
- Li, Xiaobing (2018). "Analysis of Large-Scale Traffic Incidents and En Route Diversions Due to Congestion on Freeways". In: ().
- Li, Xiaobing et al. (2020b). "Sequential prediction for large-scale traffic incident duration: Application and comparison of survival models". In: *Transportation research record* 2674.1, pp. 79–93.
- Liaw, Andy, Matthew Wiener, et al. (2002). "Classification and regression by randomForest". In: *R news* 2.3, pp. 18–22.
- Lin, Lei, Qian Wang, and Adel W Sadek (n.d.). "Duration prediction of urban freeway traffic accidents based on the M5P tree and hazard-based duration model". In: ().
- (2015). "A novel variable selection method based on frequent pattern tree for real-time traffic accident risk prediction". In: *Transportation Research Part C: Emerging Technologies* 55, pp. 444–459.
- (2016). "A combined M5P tree and hazard-based duration model for predicting urban freeway traffic accident durations". In: *Accident Analysis and Prevention* 91, pp. 114–126.
- Lin, Yunduan and Ruimin Li (2020). "Real-time traffic accidents post-impact prediction: Based on crowdsourcing data". In: *Accident Analysis and Prevention* 145, p. 105696.

- Linking Melbourne Authority (n.d.). *Department of Economic Development, Jobs, Transport and Resources Annual Report for 2014-2015*. URL: [https://www.parliament.vic.gov.au/file/uploads/Linking\\_Melbourne\\_Authority\\_Annual\\_Report\\_2014-2015\\_CwBGv8WN.pdf](https://www.parliament.vic.gov.au/file/uploads/Linking_Melbourne_Authority_Annual_Report_2014-2015_CwBGv8WN.pdf) (visited on 02/02/2021).
- Liu, Chun et al. (2017). “A dynamic spatiotemporal analysis model for traffic incident influence prediction on urban road networks”. In: *ISPRS International Journal of Geo-Information* 6.11, p. 362.
- Liu, Fei Tony, Kai Ming Ting, and Zhi-Hua Zhou (2008). “Isolation forest”. In: *2008 eighth ieee international conference on data mining*. IEEE, pp. 413–422.
- Liu, Tong et al. (2021). “Using empirical traffic trajectory data for crash risk evaluation under three-phase traffic theory framework”. In: *Accident Analysis and Prevention* 157, p. 106191.
- Lopes, Jorge et al. (2013). *Dynamic forecast of incident clearance time using adaptive artificial neural network models*. Tech. rep.
- Lu, Hao et al. (2021). “Social signal-driven knowledge automation: a focus on social transportation”. In: *IEEE Transactions on Computational Social Systems* 8.3, pp. 737–753.
- Lu, Yi-Chi (2021a). “Detecting outliers for improving the quality of incident duration prediction”. PhD thesis. University of Maryland, College Park.
- (2021b). “Detecting Outliers for Improving the Quality of Incident Duration Prediction”. In.
- Lundberg, Scott and Su-In Lee (2017). “A unified approach to interpreting model predictions”. In: *arXiv preprint arXiv:1705.07874*.
- Ma, Xiaolei et al. (2017). “Prioritizing influential factors for freeway incident clearance time prediction using the gradient boosting decision trees method”. In: *IEEE Transactions on Intelligent Transportation Systems* 18.9, pp. 2303–2310.
- Ma, Yifang et al. (2020). “Artificial intelligence applications in the development of autonomous vehicles: A survey”. In: *IEEE/CAA Journal of Automatica Sinica* 7.2, pp. 315–329.
- Machin, Mirialys et al. (Apr. 2018). “On the use of artificial intelligence techniques in intelligent transportation systems”. In: pp. 332–337. DOI: [10.1109/WCNCW.2018.8369029](https://doi.org/10.1109/WCNCW.2018.8369029).
- Mao, Tuo et al. (2021). “Boosted Genetic Algorithm using Machine Learning for traffic control optimization”. In: *IEEE Transactions on Intelligent Transportation Systems*.
- Mao, Xinhua et al. (2019). “Risk factors affecting traffic accidents at urban weaving sections: Evidence from China”. In: *International journal of environmental research and public health* 16.9, p. 1542.
- McInnes, Leland and John Healy (2018). “UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction”. In: *ArXiv* abs/1802.03426.
- Mihaita, Adriana-Simona et al. (2019a). “Arterial incident duration prediction using a bi-level framework of extreme gradient-tree boosting”. In: *arXiv preprint arXiv:1905.12254*.
- Mihaita, Adriana Simona et al. (2019b). “Arterial incident duration prediction using a bi-level framework of extreme gradient-tree boosting”. In: *CoRR* abs/1905.12254. arXiv: [1905.12254](https://arxiv.org/abs/1905.12254). URL: <http://arxiv.org/abs/1905.12254>.
- Miller, Mahalia and Chetan Gupta (2012). “Mining traffic incidents to forecast impact”. In: *Proceedings of the ACM SIGKDD international workshop on urban computing*, pp. 33–40.
- Mohammadnazar, Amin et al. (2021). “Understanding how relationships between crash frequency and correlates vary for multilane rural highways: Estimating geographically and temporally weighted regression models”. In: *Accident Analysis and Prevention* 157, p. 106146.

- Mohammed, Zainab Ali, Mohammed Najm Abdullah, and Imad Husain Al Hussaini (2021). "Predicting incident duration based on machine learning methods". In: *Iraqi Journal of Computers, Communications, Control and Systems Engineering* 21.1, pp. 1–15.
- Moosavi, Seyed Mohammad et al. (2019a). "Accident Type Detection Using Text Mining of Traffic Accident Reports". In: *IEEE Access* 7, pp. 109960–109976.
- Moosavi, Sobhan et al. (2019b). "A countrywide traffic accident dataset". In: *arXiv preprint arXiv:1906.05409*.
- Motamed, Moggan et al. (2016). "Developing a real-time freeway incident detection model using machine learning techniques". PhD thesis.
- Murdoch, W. James et al. (2019). "Definitions, methods, and applications in interpretable machine learning". In: *Proceedings of the National Academy of Sciences* 116.44, pp. 22071–22080. ISSN: 0027-8424. DOI: [10.1073/pnas.1900654116](https://doi.org/10.1073/pnas.1900654116). eprint: <https://www.pnas.org/content/116/44/22071.full.pdf>. URL: <https://www.pnas.org/content/116/44/22071>.
- Nam, Doohee and Fred Mannering (2000a). "An exploratory hazard-based analysis of highway incident duration". In: *Transportation Research Part A: Policy and Practice* 34.2, pp. 85–102.
- (2000b). "An exploratory hazard-based analysis of highway incident duration". In: *Transportation Research Part A: Policy and Practice* 34.2, pp. 85–102. ISSN: 0965-8564. DOI: [https://doi.org/10.1016/S0965-8564\(98\)00065-2](https://doi.org/10.1016/S0965-8564(98)00065-2). URL: <http://www.sciencedirect.com/science/article/pii/S0965856498000652>.
- National Highway Traffic Safety Administration (NHTSA) (2020). *Fatality Analysis Reporting System (FARS)*. URL: <https://www.nhtsa.gov/research-data/fatality-analysis-reporting-system-fars>.
- Nguyen, Hoang, Chen Cai, and Fang Chen (2017). "Automatic classification of traffic incident's severity using machine learning approaches". In: *IET Intelligent Transport Systems* 11.10, pp. 615–623.
- Nguyen, Hoang et al. (2018). "Deep learning methods in transportation domain: a review". In: *IET Intelligent Transport Systems* 12.9, pp. 998–1004.
- Olden, Julian D and Donald A Jackson (2002). "Illuminating the "black box": a randomization approach for understanding variable contributions in artificial neural networks". In: *Ecological Modelling* 154.1, pp. 135–150. ISSN: 0304-3800. DOI: [https://doi.org/10.1016/S0304-3800\(02\)00064-9](https://doi.org/10.1016/S0304-3800(02)00064-9). URL: <https://www.sciencedirect.com/science/article/pii/S0304380002000649>.
- O
- textasciitilde na, Juan de, Randa Oqab Mujalli, and Francisco J Calvo (2011). "Analysis of traffic accident injury severity on Spanish rural highways using Bayesian networks". In: *Accident Analysis and Prevention* 43.1, pp. 402–411.
- Organization, World Health (2015). *Global status report on road safety 2015*. World Health Organization.
- Ozbay, Kaan and Pushkin Kachroo (1999). "Incident management in intelligent transportation systems". In.
- Ozbay, Kaan and Nebahat Noyan (2006a). "Estimation of incident clearance times using Bayesian Networks approach". In: *Accident Analysis and Prevention* 38.3, pp. 542–555.
- (2006b). "Estimation of incident clearance times using Bayesian Networks approach". In: *Accident Analysis and Prevention* 38.3, pp. 542–555.



- Ozen, Halit et al. (2019). "Multi-step approach to improving accuracy of incident duration estimation: case study of Istanbul". In: *Tehnički vjesnik* 26.6, pp. 1777–1783.
- Pan, Bei et al. (2013). "Forecasting spatiotemporal impact of traffic incidents on road networks". In: *2013 IEEE 13th International Conference on Data Mining*. IEEE, pp. 587–596.
- Park, Hyoshin, Ali Haghani, and Xin Zhang (2016a). "Interpretation of Bayesian neural networks for predicting the duration of detected incidents". In: *Journal of Intelligent Transportation Systems* 20.4, pp. 385–400.
- (2016b). "Interpretation of Bayesian neural networks for predicting the duration of detected incidents". In: *Journal of Intelligent Transportation Systems* 20.4, pp. 385–400.
- Parsa, Amir Bahador et al. (2020). "Toward safer highways, application of XGBoost and SHAP for real-time accident detection and feature analysis". In: *Accident Analysis and Prevention* 136, p. 105405.
- Parsa, Siavash et al. (2019). "Real-Time Detection and Classification of Objects in Traffic Scenes". In: *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, pp. 1095–1100.
- Pereira, Francisco C, Filipe Rodrigues, and Moshe Ben-Akiva (2013). "Text analysis in incident duration prediction". In: *Transportation Research Part C: Emerging Technologies* 37, pp. 177–192.
- Prati, Ronaldo, Gustavo Batista, and Maria-Carolina Monard (Jan. 2009). "Data mining with imbalanced class distributions: Concepts and methods". In: pp. 359–376.
- Ren, Honglei et al. (2018). "A deep learning approach to the citywide traffic accident risk prediction". In: *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, pp. 3346–3351.
- Ren, Zhenhua et al. (2017). "Deep Spatial-Temporal Residual Networks for Citywide Crowd Flows Prediction". In: *arXiv preprint arXiv:1710.00025*.
- Reza, RM Zahid and Srinivas S Pulugurtha (2019). "Forecasting short-term relative changes in travel time on a freeway". In: *Case Studies on Transport Policy* 7.2, pp. 205–217.
- Salas, Angelica, Panagiotis Georgakis, and Yannis Petalas (2017). "Incident detection using data from social media". In: *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 751–755. DOI: [10.1109/ITSC.2017.8317967](https://doi.org/10.1109/ITSC.2017.8317967).
- Saracoglu, Abdulsamet and Halit Ozen (2020). "Estimation of traffic incident duration: A comparative study of decision tree models". In: *Arabian Journal for Science and Engineering* 45.10, pp. 8099–8110.
- Schapiro, Robert E (2013). "Explaining adaboost". In: *Empirical inference*. Springer, pp. 37–52.
- Schindler, Ron and Giulio Bianchi Piccinini (2021). "Truck drivers' behavior in encounters with vulnerable road users at intersections: Results from a test-track experiment". In: *Accident Analysis and Prevention* 159, p. 106289.
- Schrank, David and Tim Lomax (2002). "The 2002 Urban Mobility Report (College Station, TX: Texas Transportation Institute, Texas AandM University, June)". In.
- Sen, Pranab Kumar and Sujit Kumar Dhar (2018). "Stock price prediction using LSTM, RNN and CNN-sliding window model". In: *International journal of scientific and engineering research* 9.11, pp. 2046–2054.

- Shafiei, S et al. (2020). "Short-Term Traffic Prediction Under Non-Recurrent Incident Conditions Integrating Data-Driven Models and Traffic Simulation". In: *Transportation Research Board 99th Annual Meeting*.
- Shang, Qiang, Tian Xie, and Yang Yu (2022). "Prediction of duration of traffic incidents by hybrid deep learning based on multi-source incomplete data". In: *International journal of environmental research and public health* 19.17, p. 10903.
- Shang, Qiang et al. (2019). "A hybrid method for traffic incident duration prediction using BOA-optimized random forest combined with neighborhood components analysis". In: *Journal of Advanced Transportation* 2019.
- Shen, Luou and Min Huang (2011). "Data mining method for incident duration prediction". In: *International Conference on Applied Informatics and Communication*. Springer, pp. 484–492.
- Shi, Xiupeng et al. (2018). *Accident risk prediction based on driving behavior feature learning using CART and XGBoost*. Tech. rep.
- Singh, Swarnima and Vikash Yadav (2021). "An Improved Particle Swarm Optimization for Prediction of Accident Severity". In: *IJEER* 9.3, pp. 42–47.
- Smith, K and B Smith (2001). *Forecasting the Clearance Time of Freeway Accidents Final report of ITS Center project: Incident Duration Forecasting*. Tech. rep. Smart Travel Lab Report No. STL-2001-01.
- Smith, Kevin and Brian L Smith (2002). "Forecasting the clearance time of freeway accidents". In: Sullivan, Edward C (1997). "New model for predicting freeway incidents and incident delays". In: *Journal of Transportation Engineering* 123.4, pp. 267–275.
- Taghipour, Homa et al. (2022). "A novel deep ensemble based approach to detect crashes using sequential traffic data". In: *IATSS Research* 46.1, pp. 122–129. ISSN: 0386-1112. DOI: <https://doi.org/10.1016/j.iatssr.2021.10.004>. URL: <https://www.sciencedirect.com/science/article/pii/S0386111221000455>.
- Tahir, Muhammad Atif, Josef Kittler, and Fei Yan (2012). "Inverse random under sampling for class imbalance problem and its application to multi-label classification". In: *Pattern Recognition* 45.10, pp. 3738–3750.
- Tajtehranifard, Hasti et al. (2016). "Motorway crash duration and its determinants: do durations vary across motorways?" In: *Journal of Advanced Transportation* 50.5, pp. 717–735.
- Tang, Jinjun et al. (2020a). "Statistical and machine-learning methods for clearance time prediction of road incidents: A methodology review". In: *Analytic Methods in Accident Research* 27, p. 100123.
- Tang, Jinjun et al. (2020b). "Traffic incident clearance time prediction and influencing factor analysis using extreme gradient boosting model". In: *Journal of Advanced Transportation* 2020.
- Tang, Yuchun et al. (Mar. 2009). "SVMs Modeling for Highly Imbalanced Classification". In: *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on* 39, pp. 281–288. DOI: [10.1109/TSMCB.2008.2002909](https://doi.org/10.1109/TSMCB.2008.2002909).
- Theofilatos, Athanasios et al. (2016). "Predicting Road Accidents: A Rare-events Modeling Approach". In: *Transportation Research Procedia* 14. Transport Research Arena TRA2016, pp. 3399–3405. ISSN: 2352-1465. DOI: <https://doi.org/10.1016/j.trpro.2016.05.293>. URL: <https://www.sciencedirect.com/science/article/pii/S235214651630299X>.

- Tommasi, Tatiana et al. (2017). "A deeper look at dataset bias". In: *Domain adaptation in computer vision applications*, pp. 37–55.
- TomTom (n.d.). *TomTom*. URL: <https://www.tomtom.com/> (visited on 02/02/2021).
- Transportation, Federal Highway Administration of the United States Department of (2017). *The Manual on Uniform Traffic Control Devices*. URL: <https://mutcd.fhwa.dot.gov/hdm/2009/part6/part6i.htm>.
- Treiber, Martin and Arne Kesting (2013a). "Traffic flow dynamics". In: *Traffic Flow Dynamics: Data, Models and Simulation*, Springer-Verlag Berlin Heidelberg.
- (2013b). "Traffic flow dynamics". In: *Traffic Flow Dynamics: Data, Models and Simulation*, Springer-Verlag Berlin Heidelberg, pp. 983–1000.
- Tsubota, Takahiro et al. (2018). "Effect of Road Pavement Types and Ages on Traffic Accident Risks". In: *Transportation Research Procedia* 34. International Symposium of Transport Simulation (ISTS'18) and the International Workshop on Traffic Data Collection and its Standardization (IWTDCS'18) Emerging Transport Technologies for Next Generation Mobility, pp. 211–218. ISSN: 2352-1465. DOI: <https://doi.org/10.1016/j.trpro.2018.11.034>. URL: <https://www.sciencedirect.com/science/article/pii/S2352146518303235>.
- UK Government (n.d.). *UK Road Safety Statistics*. URL: <https://www.data.gov.uk/dataset/cb7ae6f0-4be6-4935-9277-47e5ce24a11f/road-safety-data> (visited on 02/02/2021).
- U.S. Department of Transportation (USDOT)/Bureau of Transportation Statistics (BTS) (2020). *National Transportation Atlas Database (NTAD)*. URL: <https://www.bts.gov/ntad>.
- Valenti, Gaetano, Maria Lelli, and Domenico Cucina (2010a). "A comparative study of models for the incident duration prediction". In: *European Transport Research Review* 2.2, pp. 103–111.
- (2010b). "A comparative study of models for the incident duration prediction". In: *European Transport Research Review* 2.2, pp. 103–111.
- Vallejos, Sebastian et al. (Feb. 2021). "Mining Social Networks to Detect Traffic Incidents". In: *Information Systems Frontiers* 23. DOI: [10.1007/s10796-020-09994-3](https://doi.org/10.1007/s10796-020-09994-3).
- Vlahogianni, Eleni I and Matthew G Karlaftis (2013). "Fuzzy-entropy neural network freeway incident duration modeling with single and competing uncertainties". In: *Computer-Aided Civil and Infrastructure Engineering* 28.6, pp. 420–433.
- Waetjen, David and Fraser Shilling (2021). "Rapid Reporting of Vehicle Crash Data in California to Understand Impacts from COVID-19 Pandemic on Traffic and Incidents". In.
- Wali, Behram, Asad J Khattak, and Jun Liu (n.d.). "Heterogeneity Assessment in Incident Duration Modelling: Implications for Development of Practical Strategies for Small and Large Scale Incidents". In: ().
- Wang, Li-Li, Henry Y.T. Ngan, and Nelson H.C. Yung (2018a). "Automatic incident classification for large-scale traffic data by adaptive boosting SVM". In: *Information Sciences* 467, pp. 59–73. ISSN: 0020-0255. DOI: <https://doi.org/10.1016/j.ins.2018.07.044>. URL: <http://www.sciencedirect.com/science/article/pii/S0020025518305681>.
- Wang, Li-Li, Henry YT Ngan, and Nelson HC Yung (2018b). "Automatic incident classification for large-scale traffic data by adaptive boosting SVM". In: *Information Sciences* 467, pp. 59–73.
- Wang, Shi, Ruimin Li, and Min Guo (2018). "Application of nonparametric regression in predicting traffic incident duration". In: *Transport* 33.1, pp. 22–31.



- Wang, Shuangcheng et al. (2022). “Research on a dynamic full Bayesian classifier for time-series data with insufficient information”. In: *Applied Intelligence*, pp. 1–17.
- Wang, Shujie et al. (2021a). “Traffic Accident Risk Prediction with Attention-based Multi-Modal Fusion Networks”. In: *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pp. 2031–2040.
- Wang, Wenqun, Haibo Chen, and Margaret Bell (2002). “A study of the characteristics of traffic incident duration on motorways”. In: *Traffic And Transportation Studies (2002)*, pp. 1101–1108.
- Wang, WenQun, Haibo Chen, and MARGARET C Bell (2005). “Vehicle breakdown duration modeling”. In: *Journal of Transportation and Statistics* 8.1, p. 75.
- Wang, Yuxiang et al. (2021b). “GSNET: Graph-Structured Network for Traffic Scene Understanding”. In: *IEEE transactions on pattern analysis and machine intelligence*.
- Wei, Chien-Hung and Ying Lee (2007). “Sequential forecast of incident duration using artificial neural network models”. In: *Accident Analysis and Prevention* 39.5, pp. 944–954.
- Wen, Tao et al. (2018). “Integrated incident decision-support using traffic simulation and data-driven models”. In: *Transportation research record* 2672.42, pp. 247–256.
- Wen, Yuan et al. (2013). “Traffic Incident duration prediction based on k-nearest neighbor”. In: *Applied Mechanics and Materials*. Vol. 253. Trans Tech Publ, pp. 1675–1681.
- Wold, Svante, Kim Esbensen, and Paul Geladi (1987). “Principal component analysis”. In: *Chemometrics and Intelligent Laboratory Systems* 2.1. Proceedings of the Multivariate Statistical Workshop for Geologists and Geochemists, pp. 37 –52. ISSN: 0169-7439. DOI: [https://doi.org/10.1016/0169-7439\(87\)80084-9](https://doi.org/10.1016/0169-7439(87)80084-9). URL: <http://www.sciencedirect.com/science/article/pii/0169743987800849>.
- Won, Minsu (2020). “Outlier analysis to improve the performance of an incident duration estimation and incident management system”. In: *Transportation research record* 2674.5, pp. 486–497.
- World Health Organization (n.d.). *World Health Organization’s Global Health Estimates*. URL: <https://www.who.int/data/gho/data/indicators/indicator-details/GHO/estimated-number-of-road-traffic-deaths> (visited on 02/02/2021).
- Wu, Wei-wei, Shu-yan Chen, and Chang-jiang Zheng (2011). “Traffic incident duration prediction based on support vector regression”. In: *ICCTP 2011: Towards Sustainable Transportation Systems*, pp. 2412–2421.
- Xiao, Jiajian et al. (2019). “A Survey on Agent-Based Simulation using Hardware Accelerators”. In: *ACM Computing Surveys (CSUR)* 51.6, p. 131.
- Xiao, Siyao (2021). “Traffic accident duration prediction based on natural language processing and a hybrid neural network architecture”. In: *2021 International Conference on Neural Networks, Information and Communication Engineering*. Vol. 11933. SPIE, pp. 194–202.
- Yang, Chao, Mingyang Chen, and Quan Yuan (2021). “The application of XGBoost and SHAP to examining the factors in freight truck-related crashes: An exploratory analysis”. In: *Accident Analysis and Prevention* 158, p. 106153.
- Yang, Hong et al. (n.d.). “Development of an automated approach for quantifying spatio-temporal impact of traffic incidents”. In.
- Yannis, George et al. (2016). “Use of Accident Prediction Models in Road Safety Management – An International Inquiry”. In: *Transportation Research Procedia* 14. Transport Research Arena

- TRA2016, pp. 4257–4266. ISSN: 2352-1465. DOI: <https://doi.org/10.1016/j.trpro.2016.05.397>. URL: <https://www.sciencedirect.com/science/article/pii/S2352146516304033>.
- Yi, Dewei et al. (2019). “A machine learning based personalized system for driving state recognition”. In: *Transportation Research Part C: Emerging Technologies* 105, pp. 241–261.
- Yu, B et al. (2016). “A comparison of the performance of ANN and SVM for the prediction of traffic accident duration”. In: *Neural Network World* 26.3, p. 271.
- Yu, Bin and Zhengfeng Xia (July 2012). “A Methodology for Freeway Incident Duration Prediction Using Computerized Historical Database”. In: pp. 3463–3474. ISBN: 978-0-7844-1244-2. DOI: [10.1061/9780784412442.351](https://doi.org/10.1061/9780784412442.351).
- Yu, Rose et al. (2017). “Deep learning: A generic approach for extreme condition traffic forecasting”. In: *Proceedings of the 2017 SIAM international Conference on Data Mining*. SIAM, pp. 777–785.
- Yuan, Jing et al. (2018). “Hetero-ConvLSTM: A Hierarchical Spatiotemporal Network for Accident Risk Prediction”. In: *IEEE Transactions on Intelligent Transportation Systems* 20.3, pp. 764–776.
- Zeng, Xiaosi and Praput Songchitruksa (2010). “Empirical Method for Estimating Traffic Incident Recovery Time”. In: *Transportation Research Record* 2178.1, pp. 119–127. DOI: [10.3141/2178-13](https://doi.org/10.3141/2178-13). eprint: <https://doi.org/10.3141/2178-13>. URL: <https://doi.org/10.3141/2178-13>.
- Zhan, Chengjun, Albert Gan, and Mohammed Hadi (Dec. 2011). “Prediction of Lane Clearance Time of Freeway Incidents Using the M5P Tree Algorithm”. In: *IEEE Transactions on Intelligent Transportation Systems* 12, pp. 1549–1557. DOI: [10.1109/TITS.2011.2161634](https://doi.org/10.1109/TITS.2011.2161634).
- Zhang, Zhiyuan, Zhiwei Chen, and Xianpei Zhu (2018). “Deep learning based incident detection from social media data”. In: *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, pp. 1059–1064.
- Zhao, Ling et al. (2018). “Temporal graph convolutional network for urban traffic flow prediction method”. In: *arXiv preprint arXiv:1811.05320*.
- Zhao, Yuexu and Wei Deng (2022). “Prediction in traffic accident duration based on heterogeneous ensemble learning”. In: *Applied Artificial Intelligence* 36.1, p. 2018643.
- Zheng, Qikang et al. (2021). “Investigating the predictability of crashes on different freeway segments using the real-time crash risk models”. In: *Accident Analysis and Prevention* 159, p. 106213.
- Zhou, Ruiming et al. (2020a). “Foresee: A Multi-Scale Multi-Task Network for Accident Risk Prediction”. In: *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pp. 2021–2030.
- Zhou, Ruiming et al. (2020b). “RiskOracle: A Graph-Convolution Network for Fine-Grained Accident Risk Prediction”. In: *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pp. 1437–1446.
- Zhu, Weiwei et al. (2021). “Dynamic prediction of traffic incident duration on urban expressways: A deep learning approach based on LSTM and MLP”. In: *Journal of intelligent and connected vehicles* 4.2, pp. 80–91.
- Ziakopoulos, Apostolos (2021). “Spatial analysis of harsh driving behavior events in urban networks using high-resolution smartphone and geometric data”. In: *Accident Analysis and Prevention* 157, p. 106189.

- Zong, Fang, Hongguo Xu, and Huiyong Zhang (2013). “Prediction for traffic accident severity: comparing the Bayesian network and regression models”. In: *Mathematical Problems in Engineering* 2013.
- Zou, Yajie et al. (2016a). “Application of finite mixture models for analysing freeway incident clearance time”. In: *Transportmetrica A: Transport Science* 12.2, pp. 99–115.
- (2016b). “Application of finite mixture models for analysing freeway incident clearance time”. In: *Transportmetrica A: Transport Science* 12.2, pp. 99–115.
- Zou, Yajie et al. (2018a). “Jointly analyzing freeway traffic incident clearance and response time using a copula-based approach”. In: *Transportation research part C: emerging technologies* 86, pp. 171–182.
- (2018b). “Jointly analyzing freeway traffic incident clearance and response time using a copula-based approach”. In: *Transportation Research Part C: Emerging Technologies* 86, pp. 171–182. ISSN: 0968-090X. DOI: <https://doi.org/10.1016/j.trc.2017.11.004>. URL: <https://www.sciencedirect.com/science/article/pii/S0968090X17303108>.
- Zou, Yajie et al. (2021). “Application of the bayesian model averaging in analyzing freeway traffic incident clearance time for emergency management”. In: *Journal of advanced transportation* 2021, pp. 1–9.