# Highway Safety Analytics And Modeling

**Dominique Lord**
**Xiao Qin**
**Srinivas R. Geedipally**

# Highway Safety Analytics and Modeling

by
Dominique Lord, Texas A&M University
Xiao Qin, University of Wisconsin-Milwaukee
Srinivas R. Geedipally, Texas A&M Transportation Institute

To be published by Elsevier

The primary purpose of this textbook is to provide information for practitioners, engineers, scientists and researchers who are interested in analyzing safety data in order to make engineering- or policy-based decisions. This book provides the latest tools and methods documented in the literature for analyzing crash data, some of which have in fact been developed or introduced by the authors. The textbook covers all aspects of the decision-making process, from collecting and assembling data to making decisions based on the results of the analyses. Several examples and case studies are provided to help understand models and methods commonly used for analyzing crash data. Where warranted, helpful hints and suggestions are provided by the authors in the text to support the analysis and interpretation of crash data.

(With the exception of Chapter 1, the word counts for all chapters are between 12,000 and 15,000 words)

## CHAPTER 1 – INTRODUCTION

### 1.1 MOTIVATION

### 1.2 IMPORTANT FEATURES OF THIS TEXTBOOK

### 1.3 ORGANIZATION OF TEXTBOOK

#### 1.3.1 Part I: THEORY AND BACKBROUND

#### 1.3.2 Part II: HIGHWAY SAFETY ANALYSES

#### 1.3.3 Part III: ALTERNATIVE SAFETY ANALYSES

#### 1.3.4 Appendices

### 1.4 FUTURE CHALLENGES AND OPPORTUNITIES

### 1.5 REFERENCES

## Part I: THEORY AND BACKBROUND

## CHAPTER 2 – FUNDAMENTALS AND DATA COLLECTION

**2.1 INTRODUCTION**

**2.2 CRASH PROCESS: DRIVERS, ROADWAYS, AND VEHICLES**

**2.3 CRASH PROCESS: ANALYTICAL FRAMEWORK**

**2.4 SOURCES OF DATA AND DATA COLLECTION PROCEDURES**

**2.4.1 Traditional Data**

*2.4.1.1 Crash Data*

*2.4.1.2 Roadway Data*

*2.4.1.3 Traffic Flow Data*

*2.4.1.4 Supplemental Data*

*2.4.1.5 Other Safety-Related Data and Relevant Databases*

**2.4.2 Naturalistic Driving Data**

**2.4.3 Disruptive Technological and Crowdsourcing Data**

**2.4.4 Data Issues**

**2.5 ASSEMBLING DATA**

**2.6 4-STAGE MODELING FRAMEWORK**

**2.6.1 Determine Modeling Objective Matrix**

**2.6.2 Establish Appropriate Process to Develop Models**

**2.6.3 Determine Inferential Goals**

**2.6.4 Select Computation Techniques and Tools**

**2.7 METHODS FOR EVALUATING MODEL PERFORMANCE**

**2.7.1 Likelihood-Based Methods**

**2.7.2 Error-Based Methods**

**2.8 HEURISTIC METHODS FOR MODEL SELECTION**

## CHAPTER 4 – CRASH-SEVERITY MODELING

## Part II: HIGHWAY SAFETY ANALYSES

## CHAPTER 5 – EXPLORATORY ANALYSES OF SAFETY DATA

## CHAPTER 6 – CROSS-SECTIONAL AND PANEL STUDIES IN SAFETY

**CHAPTER 7 – BEFORE–AFTER STUDIES IN HIGHWAY SAFETY**

## CHAPTER 8 – IDENTIFICATION OF HAZARDOUS SITES

# CHAPTER 10 – CAPACITY, MOBILITY, AND SAFETY

## 10.1 INTRODUCTION

## 10.2 MODELING SPACE BETWEEN VEHICLES

## 10.3 SAFETY AS A FUNCTION OF TRAFFIC FLOW

# Part III: ALTERNATIVE SAFETY ANALYSES

# CHAPTER 11 – SURROGATE SAFETY MEASURES

**12.6. NEURAL NETWORK**

        **12.6.1. Multilayer Perceptron (MLP) Neural Network**

        **12.6.2. Convolutional Neural Networks (CNN)**

        **12.6.3. Long Short-Term Memory - Recurrent Neural Networks (LSTM-RNN)**

        **12.6.4  Bayesian Neural Networks (BNN)**

**12.7. SUPPORT VECTOR MACHINES (SVM)**

**12.8. SENSITIVITY ANALYSIS**

**12.9 REFERENCES**

# APPENDICES

**Appendix A: Negative Binomial Regression Models and Estimation Methods**
**Appendix B: Summary of Crash-Frequency and Crash-Severity Models in Highway Safety**
**Appendix C: Computing Codes**
**Appendix D: List of Exercise Datasets**