# Network analysis methods for smart inspection in the transport domain
Bruin, G.J. de

# Network Analysis Methods
# for Smart Inspection
# in the Transport Domain

Gerrit Jan de Bruin

Human Environment and Transport
Inspectorate
*Ministry of Infrastructure
and Water Management*

Universiteit
Leiden

# Network Analysis Methods
# for Smart Inspection
# in the Transport Domain

# Proefschrift

ter verkrijging van
de graad van doctor aan de Universiteit Leiden
op gezag van rector magnificus prof.dr.ir. H. Bijl,
volgens besluit van het college voor promoties
te verdedigen op donderdag 16 november 2023
klokke 13.45 uur

door

Gerrit Jan de Bruin
geboren te Amersfoort
in 1993

# Preface

This thesis is part of an extensive collaboration between the Dutch Ministry of Infrastructure and Water Management (I&W) and Leiden University. My first personal encounter with I&W was during my master's study in Analytical Chemistry. I was determined to work on a subject with societal relevance, leading me to Jasper van Vliet, who worked at the Inspectie Leefomgeving en Transport (ILT), part of I&W. The aim of my master thesis was to conduct efficient compliance monitoring of cargo ship's fuel by using the information from chemical sensors. The thesis made it to a letter to the parliament [88, 89]. Afterward, Jasper invited me to participate in a new Ph.D. project of the Ministry. In this project, the ambition is to arrive at intelligence-led vehicle inspection by risk assessments. Two research directions were launched to explore the risk assessment of vehicles: (1) the application of machine learning techniques and (2) the application of network science techniques.

During my Ph.D. research, many developments occurred related to the topics of the thesis. I would like to mention two specific events that have had an impact on my research. The first one is the introduction of the General Data Protection Regulation (GDPR) in 2018. This law requires the use of transparent models that allow for an explanation of the results achieved. The new law led to increased awareness of the importance of fair data and fair models. The second event is the upset of the Dutch childcare benefits ("De Toeslagenaffaire"). In 2019 it became painfully clear how things can go wrong when authorities are (1) relying on biased data and (2) using models that are not validated fairly. I will address the two points (biases and non-validated models) in my thesis, although I work with non-personal data. They are in particular relevant for the proposed procedure to implement a smart inspection of cargo ships in Chapter 6.

Working both at Leiden University and at the ILT allowed me to interact with the wonderful world of academia and to stay in close contact with a governmental organization that makes a big impact by ensuring safe transportation and reducing the environmental pollution in the Netherlands. This enriching combination has helped me to create this thesis, for which I am grateful.

Gerrit Jan de Bruin, Utrecht, March 8th, 2023

# Contents

# List of Abbreviations

# List of Definitions

# List of Figures

# List of Tables