

MARITIME TRAFFIC ANOMALY DETECTION FROM AIS
SATELLITE DATA IN NEAR PORT REGIONS

by

Bo Liu

Submitted in partial fulfillment of the
requirements for the degree of
Master of Computer Science

at

Dalhousie University
Halifax, Nova Scotia
August 2015

© Copyright by Bo Liu, 2015

Table of Contents

| | |
|--|-------------|
| List of Tables | iv |
| List of Figures | v |
| Abstract | viii |
| List of Abbreviation and Symbols Used | ix |
| Acknowledgements | x |
| Chapter 1 Introduction | 1 |
| 1.1 Aim and Objectives | 3 |
| 1.2 Research Methodology | 3 |
| 1.3 Scientific Contribution | 6 |
| 1.4 Publications | 6 |
| 1.5 Thesis Outline | 7 |
| Chapter 2 Related Work | 9 |
| 2.1 Trajectory Clustering And Maritime Traffic Clustering | 9 |
| 2.2 Anomaly Detection in Maritime Surveillance | 15 |
| Chapter 3 Normal Traffic Patterns Extraction Model | 19 |
| 3.1 Representing Trajectory Data | 20 |
| 3.2 Normal Moving Trajectories Extraction | 21 |
| 3.2.1 Density-Based Spatial Clustering of Applications with Noise considering Speed and Direction | 21 |
| 3.2.2 Gravity Vector | 28 |
| 3.3 Normal Stopping Areas Extraction | 31 |
| Chapter 4 Anomaly Detection | 36 |
| 4.1 Three Division Distances | 36 |
| 4.2 Anomaly Detection Model | 39 |

| | | |
|---------------------|---|-----------|
| Chapter 5 | Evaluation | 43 |
| 5.1 | Normal Traffic Patterns Extraction Model | 43 |
| 5.1.1 | Juan de Fuca Strait | 44 |
| 5.1.2 | Los Angeles Long Beach | 45 |
| 5.2 | Anomaly Detection Model | 47 |
| 5.2.1 | Experiment On Unlabeled Data Set | 49 |
| 5.2.2 | Experiment On Labeled Data Set | 53 |
| Chapter 6 | Conclusion | 59 |
| 6.1 | Future Work | 60 |
| Bibliography | | 64 |
| Appendix A | Source Code of the Algorithms | 70 |
| A.1 | Source Code of Normal Traffic Patterns Extraction Model | 70 |
| A.1.1 | Code of Clustering Process | 71 |
| A.1.2 | GVs Calculation | 75 |
| A.1.3 | SSPs Calculation | 80 |
| A.2 | Source Code of Anomaly Detection Model | 82 |

List of Tables

| | | |
|-----------|---|----|
| Table 5.1 | Quartile statistics of the three division distances | 52 |
| Table 5.2 | Labels and descriptions by the expert | 56 |
| Table 5.3 | Confusion Matrix I | 58 |
| Table 5.4 | Confusion Matrix II | 58 |

List of Figures

| | | |
|------------|---|----|
| Figure 1.1 | The framework for Maritime Anomaly Detection. This can be roughly regarded as a two-step procedure. First, the normal traffic extraction model extracts the normal traffic patterns from the historical AIS data repository. Then the extracted Gravity Vectors (GVs) and Sampled Stopping Points (SSPs) are employed by the anomaly detection model to decide the abnormality of the new given trajectory. | 4 |
| Figure 3.1 | Three defined relations for the density-based clustering algorithm. | 24 |
| Figure 3.2 | Calculate Gravity Vectors (GV) for one moving cluster. Blue arrows stand for the trajectory points of the cluster, red arrows are the final Gravity Vectors. The length (L) of the GV is the SOG_{avg} of the GV. The width of a grid g is the pre-defined length for partitioning the points. | 29 |
| Figure 3.3 | Border Values of a stopping cluster | 32 |
| Figure 3.4 | An example of a stopping area where simple random sampling cannot work. There are three regions (A,B and C) with higher density than the whole area's average density. | 34 |
| Figure 4.1 | Two Abnormal Cases after considering COG and SOG | 38 |
| Figure 4.2 | One case that the simple combination of the two division distances cannot work. There are two moving clusters in one specific area and the two clusters are with two different average COGs. | 39 |
| Figure 4.3 | The maximum CDD distribution in the area of Juan de Fuca Strait. | 41 |
| Figure 4.4 | The Cumulative Distribution Function of the maximum CDD in the area of Juan de Fuca Strait. | 41 |
| Figure 5.1 | Juan de Fuca Strait and its approaches (west) [25] | 44 |
| Figure 5.2 | The clustering results of JUAN DE FUCA STRAIT area. Points in blue color are the noise. Other different colors stand for different clusters. The cluster in dark green is a stopping cluster. | 45 |
| Figure 5.3 | The Gravity Vectors and Sampled Stopping Points extracted from the clusters in Juan de Fuca Strait area. | 46 |

| | | |
|-------------|--|----|
| Figure 5.4 | Map of Los Angeles Long Beach and the navigation rules defined [25] | 46 |
| Figure 5.5 | Extract Stopping Sampling Points (SSP) from stopping clusters in Los Angeles Port Area. The stopping clusters are first extracted using DBSCAN [15] algorithm shown in (a) and different colors stand for different clusters. Then Algorithm 2 is used for getting the SSPs (shown in figures (b) and (c)) | 48 |
| Figure 5.6 | The clustering results of Los Angeles Long Beach area. This only shows the moving points(in blue) and the corresponding moving clusters (in different colors). | 49 |
| Figure 5.7 | The Gravity Vectors and Sampled Stopping Points extracted from the clusters in Los Angeles Long Beach area. | 50 |
| Figure 5.8 | The Gravity Vectors (open circles) and Sampled Stopping Points (filled circles) extracted from the clusters in JUAN DE FUCA STRAIT area. | 51 |
| Figure 5.9 | The anomaly labeling results of the trajectory data points in JUAN DE FUCA STRAIT area. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD). | 52 |
| Figure 5.10 | First example of the abnormal trajectories detected by our framework. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD). | 53 |
| Figure 5.11 | Second example of the abnormal trajectories detected by our framework. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD). | 54 |
| Figure 5.12 | Third example of the abnormal trajectories detected by our framework. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD). | 54 |
| Figure 5.13 | Fourth example of the abnormal trajectories detected by our framework. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD). | 55 |

| | | |
|-------------|---|----|
| Figure 5.14 | Fifth example of the abnormal trajectories detected by our framework. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD). | 55 |
| Figure 5.15 | Sixth example of the abnormal trajectories detected by our framework. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD). | 56 |
| Figure 6.1 | One situation that our algorithm cannot detect the anomalies. Two normal lanes (Lane <i>A</i> and Lane <i>B</i>) detected by DB-SCANS _D are in red and purple. A new anomaly trajectory <i>C</i> is in blue. | 62 |

Abstract

Maritime traffic monitoring is an important aspect of safety and security in close-to-port operations. While there is a large amount of data with variable quality, decision makers need reliable information about possible situations or threats. In this thesis, we propose a two-component maritime traffic anomaly detection model. First, it extracts normal ship trajectory patterns using, besides ship tracing data, the publicly available IMO Rules. The main result of clustering is a set of generated lanes that can be mapped to those defined in the IMO directives. Then, we show how the second anomaly detection component detects anomalous navigational behaviors based on three specialized division distances with the clusters. It decides for each point if the vessel is anomalous, considering longitude, latitude, direction and speed. This point-based approach is applicable for real-time AIS (Automatic Identification System) surveillance; it is also feasible for analysts to set their own threshold for labeling whole trajectories.

List of Abbreviation and Symbols Used

| | |
|-----------------|---|
| $N_\epsilon(p)$ | Eps-neighborhood of a point p |
| ADD | Absolute Division Distance |
| AIS | Automatic Identification System |
| CDD | Cosine Division Distance |
| COG | Course Over Ground |
| DBSCAN | Density-Based Spatial Clustering of Applications with Noise |
| DBSCANSD | Density-Based Spatial Clustering of Applications with Noise considering Speed and Direction |
| EM | Expectation Maximization |
| GPS | Global Positioning System |
| GV | Gravity Vector |
| IMO | International Maritime Organization |
| MMSI | Maritime Mobile Service Identity |
| RDD | Relative Division Distance |
| SOG | Speed Over Ground |
| SSP | Sampled Stopping Point |
| TSS | Traffic Separation Scheme Boundaries |
| UTC | Coordinated Universal Time |

Acknowledgements

I would like to first express my deepest appreciation to my supervisor, Dr.Stan Matwin, for his support, guidance, encouragement and patience throughout my master program. His valuable suggestions and directions helped to improve the quality of this thesis. He also provides me with precious opportunities to work with multiple projects which can help me open my view.

Besides my supervisor, I would like to thank the rest of my thesis committee: Prof.Raza Abidi, Prof.Ronald Pelot and Prof.Dirk Arnold for their encouragement, insightful comments, and hard questions.

My sincere thanks also goes to Dr.Erico N.de Souza, my project leader. His advices and ideas guided me work on the right paths of the projects. He has also spent a lot time reviewing my papers and this thesis, which is a great help to me.

My parents, Baolong and Changling, who have always been supporting me and have never said no to my any decisions. Without them, this thesis cannot be finished.

I would also like to extend my gratefulness to all my friends and labmates in the Big Data Institute: Dr.Rob Warren, Dr.Xiaoguang Wang, Ahmad, Behrouz, Diana, Habibeh, Xuan, Eman, Baifan, David, Lulu and Vineeth for the stimulating discussions and for all the fun we have had in the past two years.

Last but not least, I would like to acknowledge the generous support of GSTS, Inc., ExactEarth, Inc., Marine Environmental, Observation, Prediction And Response Network (MEOPAR), and the Natural Science and Engineering Research Council of Canada (NSERC).

Chapter 1

Introduction

Today maritime transportation represents 90% of international trade volume [59] and there are more than 50,000 vessels sailing the ocean every day [32]. So the challenges related to safety and security aspects in the maritime domain are of high priority. Modeling ship behavior helps improve security for authorities and one crucial task of maritime surveillance is abnormal vessel behaviours detection. Anomalies in the maritime domain, such as unexpected stops, deviations from regulated routes, inconsistent speed or direction etc., may be related to risks like collisions¹, grounding², sea drunkenness³, smuggling, piracy⁴, etc. To address the issue, different data sources such as data from Automatic Identification System (AIS), synthetic aperture radar, high frequency surface wave radars, infra-red sensors, videos and intelligent reports [27] are used by the maritime authorities. Data generated by AIS are utilized in this thesis.

Automatic Identification System (AIS) is an automatic tracking system to identify and locate vessels by exchanging data with other nearby ships, AIS base stations and satellites. The use of this system has been required by the International Maritime Organization (IMO) since 2004 to enhance safety and efficiency of navigation and improve situational awareness and assessment [21]. IMO was established in Geneva in 1948 with an initial name of Inter-Governmental Maritime Consultative Organization (IMCO). Then in 1982, it was changed to the International Maritime Organization⁵. The IMO's headquarter is in London, United Kingdom and it is a specialised agency of the United Nations with 171 Member States and three Associate Members. IMO is the global standard-setting authority for the safety, security and environmental

¹<http://gcaptain.com/tag/ship-collision/>

²<http://www.wisegeek.com/what-is-ship-grounding.htm>

³<http://www.theglobeandmail.com/news/national/canadian-navy-officer-charged-with-drunkenness-disobeying-orders/article23023966/>

⁴<http://www.maritime-executive.com/article/pirates-attack-two-vessels-in-asia>

⁵https://en.wikipedia.org/wiki/International_Maritime_Organization

performance of international shipping ⁶. According to the regulations defined by IMO, all ships of 300 gross tonnage and upwards engaged on international voyages and cargo ships of 500 gross tonnage and upwards not engaged on international voyages and passenger ships irrespective of size shall be fitted with an automatic identification system and over 400,000 ships worldwide have installed these transponders [6].

AIS messages are automatically broadcasted with a reporting frequency directly proportional to the speed of the vessel [46]. The data can certainly be viewed as big data: a stream of some 80 million messages are generated per day. Two categories of data (static and dynamic information) are transmitted by the AIS transponders. Static information, which must be programmed into the AIS transceiver, includes IMO number, Maritime Mobile Service Identity (MMSI) number, vessel call sign, vessel type, vessel dimensions (length and beam). The dynamic information, which needs to be updated automatically, includes the ship's location (Longitude and Latitude), Course Over Ground (COG), Speed Over Ground (SOG), heading, Rate of Turn (ROT), time in UTC (Coordinated Universal Time), navigational status etc [2]. In this thesis, both static and dynamic information is used to tackle the maritime anomaly detection task.

IMO also controls the navigational lanes that vessels must use when approaching some ports around the world. Specifically, the rules which regulate the general directions in the specific regions for the vessels to navigate are called Traffic Separation Schemes (TSS) [1]. The regulations to establish TSSs mandatory were first proposed by IMO's Maritime Safety Committee meeting in March 1971 and this recommendation was adopted by the IMO later the same year. The first corresponding mandatory traffic scheme is the Dover Straits scheme.⁷ These navigational lanes are definitely good resources for maritime surveillance analysis, but unfortunately, current tools do not offer automatic ways to check whether the vessels are obeying these lanes.

Only checking if the ships are obeying the lanes cannot satisfy the demand of maritime surveillance. TSS regulates the normal geographical lanes for ships to sail, but it does not provide exact speed regulations in specific regions. More specifically, there are no rules to determine the maximum or minimum speed of a particular type of vessels to sail in a given area. Although there are some rules defined by particular

⁶<http://www.imo.org/en/About/Pages/Default.aspx>

⁷<http://www.imo.org/OurWork/Safety/Navigation/Pages/ShipsRouteing.aspx>

ports authorities to restrict the vessels' speed, it is still necessary to apply some methods to extract specific speed information from the real-life dataset. Fortunately, one dynamic type of information, SOG (Speed Over Ground), is included in AIS messages and approaches that are able to utilize this historical speed information for anomaly detection are desirable.

1.1 Aim and Objectives

Based on the above introduction to the background of maritime traffic anomaly detection, the overall research aim of this thesis can be stated as: *Review and analyse existing anomaly detection algorithms for trajectory data, and propose new or modified approaches that are well-suited for the task.*

In order to address the aim, the following objectives are identified:

- 1). Investigate and survey the related approaches proposed for anomaly detection in trajectory data.
- 2). Propose new algorithms fused with IMO rules to tackle anomaly detection task in the maritime domain.
- 3). Demonstrate the effectiveness of the proposed algorithms on real world AIS data sets.

1.2 Research Methodology

To address the aim and objectives stated in Section 1.1, we propose a clustering-based maritime traffic anomaly detection framework. An overview of the entire framework is shown in Figure 1.1. The approach includes two components: the normal traffic patterns extraction model and the anomaly detection model. As can be seen from Figure 1.1, historical AIS data is first sent to the normal patterns extraction model which generates, as output, a set of Gravity Vectors (GV) and Sampled Stopping Points (SSP) via its two sub-components (Normal Moving Trajectories Extraction component and Normal Stopping Areas Extraction component). Afterwards, when new ship trajectory data is to be judged, the next phase model will be applied based on the normal traffic patterns results to decide the new data's abnormality.

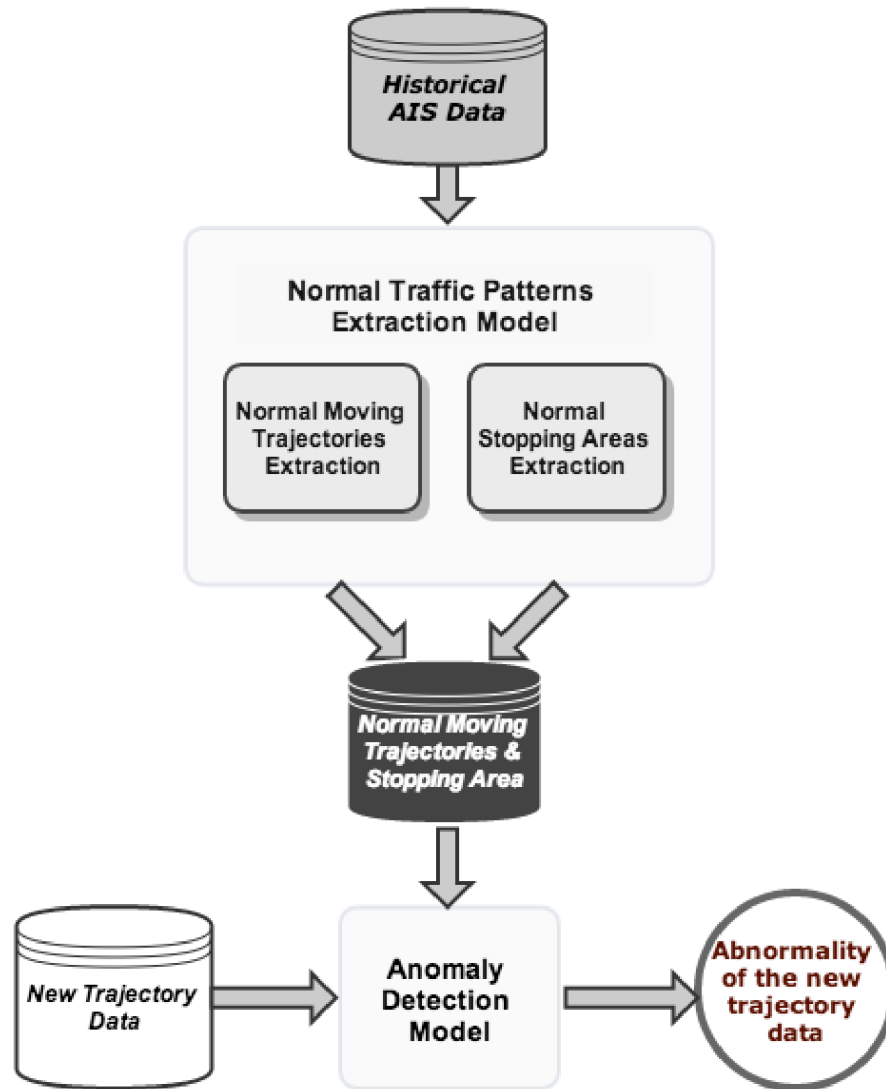


Figure 1.1: The framework for Maritime Anomaly Detection. This can be roughly regarded as a two-step procedure. First, the normal traffic extraction model extracts the normal traffic patterns from the historical AIS data repository. Then the extracted Gravity Vectors (GVs) and Sampled Stopping Points (SSPs) are employed by the anomaly detection model to decide the abnormality of the new given trajectory.

The proposed clustering-based normal traffic patterns extraction model first divides the AIS data into moving and stopping parts, respectively, based on a stopping SOG (Speed Over Ground) threshold of 0.5 knots.

For the case where SOG is not less than 0.5 knots, we propose DBSCANSD (Density-Based Spatial Clustering of Applications with Noise considering Speed and

Direction) as the basic algorithm to detect the main traffic lanes within the data. DBSCANSD is based on DBSCAN (Density-Based Spatial Clustering of Applications with Noise) [15] algorithm, modified to consider that for each point of a cluster, the neighborhood of a given radius has to contain at least a minimum number of points with similar SOGs (Speed Over Ground) and COGs (Course Over Ground). The experiments demonstrate that the clusters can reflect normal patterns of the vessels. However, it is not feasible to employ all the points of the clusters to decide the new arriving trajectory data's abnormality. As a result, a representation called Gravity Vector (GV) is proposed. A GV is a vector composed of 5 features: average COG, average SOG, average Latitude, average Longitude and Median Distance. The Median Distance of the GV is the median of the distances between all the trajectory points and the GV in the cell. The Median Distance is to measure the width of navigational lanes. Thus, the output of the moving pattern extraction model can be a set of GVs which carries information about speed, direction, location and even the relative width of the lanes. More details about the process can be found in Section 3.2.

For the case where the SOG is less than 0.5 knots (stop areas), the original DBSCAN algorithm [15] is executed because speed (SOG) and direction (COG) are not important factors. Another type of vector is created as output, called Sampled Stopping Point (SSP), which is dependent only on the geographic shape of the region. The experiments show that a small number of SSPs can represent a stopping area well in terms of the shape of the area. More details about the process can be found in Section 3.3.

Once the normal patterns are extracted from the historical AIS data, the second phase is the anomaly detection model (presented in Chapter 4), can be applied for the new arriving data. In this process, three specialized distances including Absolute Division Distance (*ADD*), Relative Division Distance (*RDD*) and Cosine Division Distance (*CDD*) are proposed to decide the abnormality of one specific data point. Specifically, *ADD* is employed in the stopping points abnormality detection phase while *RDD* and *CDD* are used for the moving part. Furthermore, one anomaly detection algorithm using the three specialized division distances is proposed in Section 4.2. Although the approach is point-based, which is applicable for real-time AIS surveillance, it is also flexible enough for analysts to set their own threshold for

labeling whole trajectories.

1.3 Scientific Contribution

This thesis is based on previously published work by the author. The published work has also been updated and extended with new theoretical and empirical results in this thesis. The main contributions of this thesis are listed as following:

- Proposal of the algorithm of DBSCANSD (Density-Based Spatial Clustering of Applications with Noise considering Speed and Direction), a novel point-based trajectory clustering algorithm which is applicable for various domains.
- Proposal of a normal maritime traffic extraction model based on *stops – and – moves* approach [53].
- Proposal of a normal maritime traffic extraction model associating with the IMO rules.
- Proposal of two methods to represent the results of clustering, which are the Gravity Vector (GV) for moving clusters and the Sampled Stopping Point (SSP) for stopping clusters.
- Proposal of three specialized division distances measures for maritime anomaly detection.
- Implementation of the proposed maritime anomaly detection framework.
- Evaluation of the proposed maritime anomaly detection framework's effectiveness using both visualization and quantitative approaches.

1.4 Publications

In this section, the publications published during the author's master study are listed. The first two publications are highly relevant to the thesis while the other two are less relevant to the thesis. Publication 2) first introduces the normal traffic extraction model and has been extended with more detailed explanation in Chapter 3. The anomaly detection model proposed in Chapter 4 is extended from the Publication 1).

- 1) Bo Liu, Erico N.de Souza, Casey Hilliard and Stan Matwin, **Ship Movement Anomaly Detection Using Specialized Distance Measures**, IEEE International Conference on Information Fusion (FUSION 2015), 6-9 July 2015, Washington, DC, USA. (to appear)
- 2) Bo Liu, Erico N.de Souza, Stan Matwin and Marcin Sydow, **Knowledge-based Clustering of Ship Trajectories Using Density-based Approach**, Proceedings of IEEE International Conference on Big Data (IEEE BigData 2014), 27-30 October 2014, Washington, DC, USA. pp. 603-608.
- 3) Xiaoguang Wang, Xuan Liu, Bo Liu, Stan Matwin and Erico N.de Souza, **Vessel Route Anomaly Detection with Hadoop MapReduce**, Proceedings of IEEE International Conference on Big Data (IEEE BigData 2014), 27-30 October 2014, Washington, DC, USA. pp. 25-30.
- 4) Robert Warren and Bo Liu, **Language, Cultural Influences and Intelligence in Historical Gazetteers of the Great War**, Proceedings of IEEE International Conference on Big Data (IEEE BigData 2014), 27-30 October 2014, Washington, DC, USA. pp. 70-72.

1.5 Thesis Outline

The remainder of the dissertation is organized as follows. Chapter 2 presents the related work. It is divided into two parts. The first part (Section 2.1) is about trajectory clustering in general and ship trajectory clustering in the maritime domain in particular. Section 2.2 is an overview and survey on anomaly detection techniques used for maritime security.

We present our clustering-based maritime anomaly detection framework in Chapter 3 and Chapter 4. In Chapter 3, we propose our maritime traffic patterns extraction model (Section 3.2 and Section 3.3) after giving the definition of a trajectory in the maritime domain (Section 3.1). As we treat the ship trajectories as two different types based on a speed threshold, moving and stopping, which are explained in more details, Section 3.2 and Section 3.3 propose the normal moving patterns extraction

model and normal stopping areas extraction model, respectively. Two different representative methods, Gravity Vector (GV) and Sampled Stopping Point (SSP) that represent the moving clustering results and stopping clustering results, are introduced accordingly in the two sections. Chapter 4 presents the anomaly detection model. As our method is a distance-based model, three division distances are first presented in Section 4.1 and then our abnormal detection algorithm is introduced in Section 4.2.

To evaluate the effectiveness of our framework, experiments are done in Chapter 5. First, in Section 5.1, to evaluate the normal traffic patterns extraction model, regions of Juan de Fuca Strait and Los Angeles Long Beach are selected and the results are presented in Section 5.1.1 and Section 5.1.2. Then in Section 5.2, we conducted another two experiments in the same region of Juan de Fuca Strait. The first one is conducted with the non-labeled data while the second one is done after labelling the data. The results of the first experiment (Section 5.2.1) are shown visually and the second experiment (Section 5.2.2) compares our model's results with the labels by the expert.

Finally, in Chapter 6 we conclude with a summary and discuss our method's limitations and the potential directions for future work.

Chapter 2

Related Work

The problem examined in this thesis can be roughly described as maritime data anomaly detection, a rather specific task. But the increasing demand for detecting and identifying anomalous events from maritime traffic data keeps attracting more and more researchers' attention.

The approach proposed in this thesis is an unsupervised model based on clustering, so in this chapter, a survey on trajectory clustering and other normal traffic extraction techniques are first given in Section 2.1. Since maritime anomaly detection can be considered as a specific type of anomaly detection tasks in maritime domain, we survey previous work in anomaly detection in Section 2.2 and also introduce some representative work done in the maritime domain.

2.1 Trajectory Clustering And Maritime Traffic Clustering

As stated in the previous chapter, our maritime anomaly detection framework is a clustering-based approach, so in this section, we first introduce the problem of trajectory clustering and then provide a survey of various algorithms for solving it. Meanwhile, discussions on the works' limitations for maritime traffic anomaly detection are also conducted while introducing different techniques.

The use of satellites and tracking facilities, such as GPS (Global Positioning System), AIS (Automatic Identification System) and RFID (Radio Frequency Identification Devices) has increased over the past decades, which leads to an increasing number of tracking applications [19]. The data spread over various domains and examples include vehicle position data [18], airplane tracking data [19], ships or vessels tracking data [42][41][59][40][47] and even hurricane or animal movement data [33][17]. One crucial task of the applications is trajectory clustering, which is the process to discover the motion patterns from the tracking data. Exploring different trajectory clustering techniques to extract clusters can help analysts to gain better

insights from the data.

Before investigating different clustering algorithms for trajectories, one point needs to be mentioned. Since the goal of trajectory clustering is to find normal paths by grouping similar trajectories together [32], there are lots of works aiming to find proper measures of similarity between trajectories, such as Euclidean distance [18], DTW (Dynamic Time Warping) [28] and LCSS (Longest Common Subsequence) [61]. A comparison between various similarity methods is conducted in [65] and majority of these measures are alignment based distances, which are not appropriate for the ship trajectory extraction task due to the sparseness of the AIS data. So in this thesis, we focus more on other clustering-based techniques to address the task.

According to [20], clustering is the process of grouping a set of physical or abstract objects into classes of similar objects and the proposed clustering algorithms are usually categorized into four types: partitioning methods (e.g., K-Means [48] and K-Medoids [26]), hierarchical methods like BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies) [64], density-based methods including DBSCAN [15], OPTICS (Ordering Points To Identify the Clustering Structure) [5] and DENCLUE (DENSITY-based CLUstEring) [23]), and grid-based methods such as STING (STatistical INFORMATION Grid-based method) [62] and CLIQUE (CLustering In QUEst) [63]). Among various types of clustering approaches, density-based clustering algorithms can be most suitable for this work because they can find arbitrary shapes of clusters. In the following paragraph, one representative approach is introduced and the reasons why this type of algorithm is suitable is discussed as well.

DBSCAN (Density-based spatial clustering of applications with noise) proposed by Martin *et al.* [15] in 1996 is one of the fundamental density-based clustering techniques. The key idea of this clustering algorithm is that for each point of a cluster the neighborhood of a given radius has to contain at least a minimum number of points. So the algorithm requires two additional parameters, *eps* and *minPts*, as input, where the *eps* is the radius and *minPts* is the minimum number of points required to form a dense region. After applying DBSCAN on a set of points, the points are classified as core points, border points and noise points. Here both core points and border points form clusters while the noise points are filtered out. Compared to other types of clustering algorithm, DBSCAN has two main advantages which make it well-suited

for the maritime trajectory clustering task. Firstly, the algorithm can provide good results in many cases since it is able to discover clusters of arbitrary shapes and it is robust to detect outliers [20]. Secondly, unlike K-Means [48], users are not required to provide the number of expected clusters before running the algorithm. However, one crucial limitation is that the original DBSCAN algorithm can only take into account location data (Longitude and Latitude), which cannot satisfy the requirement of performing this anomaly detection in the maritime domain. Therefore, as part of our work, we extend DBSCAN to adopt other attributes, such as speed and direction, to find normal traffic patterns from the moving trajectory dataset.

To handle specific issues of tracking data, researchers propose various approaches which extend the original DBSCAN in different ways. One method called DB-SMoT (Direction Based Stops and Moves of Trajectories) [47] is a clustering method based on the direction change in a minimal amount of time. For example, if there is a sequence of points with enough direction variance but the duration between the last point and the first point is less than the minimal time threshold, the sequence will not be considered as a final cluster. The main purpose of the algorithm is to find interesting places using single fishing ship trajectory data. The methodology introduced by the authors is simply using DBSCAN to cluster all the data points with a bigger direction variation than the pre-defined threshold. This algorithm is helpful in this specific scenario but cannot be fit for anomaly detection in consideration of the significance of the speed. Different from DB-SMoT, another work done by Palma *et al.* [42] called CB-SMoT (Clustering-Based Stops and Moves of Trajectories) takes speed into account. However, due to the goal of the paper which is also to find interesting places, the method only uses speed as a threshold and then conducts clustering work on the low-speed data set. In this case, the idea of CB-SMoT can be adopted for extracting stopping areas from AIS data. In Section 3.3, we present our approach to cluster and represent the stopping areas in the sea. However, CB-SMoT cannot be applied to extract different clusters with various normal speeds from navigational data set.

One partition-and-group framework, TRACLUS (TRAjectory CLUStering) [33], also uses density-based clustering to detect general lanes. TRACLUS is proposed by Lee *et al.* [33] and is one of the main references in the field of trajectory clustering.

The algorithm is trajectory-based. It tries to solve the clustering problem from the view of line segments of given whole trajectories. The general procedure of TRACCLUS is to first partition each trajectory into line segments and then group close line segments into one particular cluster. The reason why they partition the whole track into short sub-segments before the clustering process is that clustering trajectories as a whole could lose similar portions of the trajectories [33]. Lastly, a representative trajectory is computed for each cluster. One limitation of this strategy is that it cannot consider the speed, which is very important to maritime traffic tasks. Another problem in the adaptation of this algorithm to our task is that TRACCLUS is sensitive to the parameters (*eps* and *MinLns*). Although the paper provides an optimal estimation algorithm for thresholds, the experiments conducted in [33] demonstrate that the generated pair of parameters are not guaranteed to be optimal. When we tried to employ the algorithm for extracting normal ship trajectories, a satisfactory result could not be obtained even after a large number of trials. Fortunately, compared to TRACCLUS, our model is less sensitive to the parameters as long as the analysts have some knowledge of the maritime domain.

One popular scenario of applying trajectory clustering techniques is video surveillance and monitoring systems. After distinguishing the foreground objects from the background [55] and identifying the positions of the moving objects (e.g. cars, pedestrians, etc.) in the foreground image, trajectory clustering approaches can be applied to group similar trajectories of the moving objects. In the work done by Paciarelli *et al.* [44][52], a dynamic model is presented to group similar trajectories into clusters. The model is essentially maintaining a forest of trees to represent the normal patterns. The prefix of one particular node, from the root node to the current node, stands for a specific normal path formed by all the clusters in the branch (including the root and the node itself). When a new trajectory is fed into the system, a prefix matching will be conducted and the best matching cluster will adopt the new arriving trajectory to update its properties, including physical coordinates and variances of the trajectory points inside the cluster. One issue of the approach is that a large number of trees need to be stored, because different trajectories may have different starting points and each start point needs a corresponding root node. Another limitation is that the prefix matching phase requires complete trajectories to be observed, which

is infeasible in maritime traffic tracking. It is common to have gaps in the AIS data stream and the gaps can vary from milliseconds to days. In addition to the above two problems, similar to TRACCLUS [33], it cannot take speed into account, which is a fundamental issue in coastal surveillance. Currently in this thesis, the point-based approach is proposed to address the issue of the inconsistently updated trajectory data and no concern about prefix matching exists.

Inspired by the work done by Paciarelli *et al.* [44][52], Guillaume *et al.* propose a trajectory-based “unsupervised normalcy model” [32] for S-AIS data. This work is also based on stop and move framework [53] and applies corresponding strategies on moving trajectories and stopping areas. Similar to TRACCLUS [33], this technique is also a partition-and-group framework, which first partitions the historical trajectories into sub-segments. To extract normal patterns from moving trajectories, the authors introduce a technique based on the OPTICS algorithm [5]. OPTICS has the advantage of handling situations where clusters exist at different density levels, but it needs more experts’ input knowledge for selecting relevant clusters from different clusters levels, which is more complicated than our non-hierarchical clustering approach. So this is less feasible for applying the algorithm in a large area with complicated traffic routes. Then for stopping areas extraction, the authors present a representation to estimate the locations and the spatial extent of a stopping area. This method is similar to their solution used for moving trajectories, that is, it tries to use one particular node to represent an arbitrary shape of region. This will cause a bias while it comes to coastal stopping areas since the areas might not be in circle shapes. Fortunately, the Sampled Stopping Points proposed in this work can represent arbitrary shapes of regions well.

Ferreira *et al.* [17] propose a strategy called VFKM (Vector Field K-Means) based on vector field fitting. The idea behind their method is to induce a similarity notion on the data set by using streamlines of a single vector field. Two main elements of this algorithm are fitting the vector fields and assignment of trajectories to clusters [17]. Since the approach assigns the trajectories to the vector fields, it is also a trajectory-based method. VFKM is an iterative model similar to original K-Means [48] and two main steps are involved in each iteration. First, the trajectories are partitioned randomly into K clusters (similar to the initialization of original K-Means). Then the

clusters are updated by assigning each trajectory to the vector field that it fits best. These two steps are repeated iteratively until a convergence criterion is met. There are several limitations related to the algorithm. Firstly, as mentioned by the authors [17], the performance of the algorithm is sensitive to the choice of initial clusters. Secondly, although this method can consider other features (e.g. direction and speed) of the trajectory data, it cannot provide a representative trajectory. Instead, the algorithm output is k representative vector fields which cannot be employed in our case. Fortunately, in our thesis, the output of the maritime traffic extraction model is two sets of representative vectors (Gravity Vectors and Sampled Stopping Points), which is feasible for handling the anomaly detection problem.

One work which provides mechanisms to manage traffic in the ocean that is similar to ours is an unsupervised framework called TREAD (Traffic Route Extraction and Anomaly Detection), proposed by Pallotta *et al.* [41][59][40]. It is designed to detect low likelihood behaviors and predict future vessel positions from maritime traffic data. TREAD is a point-based framework and the traffic routes are built by the way-points generated by the clustering process, that is, the route objects are directly formed by the flow vectors of the vessels whose paths connect the derived way-points. The assumption of using waypoints to build the traffic lanes is that the vessels' trajectories are composed of a set of connected straight lines. Another similar methodology proposed by Gariel *et al.* [19] is also a way-point-based clustering framework for handling airspace monitoring. The objective of the method is to identify and group the turning points into way-points and to use the sequence of way-points to represent the aircraft's trajectories. To get the way-points from the historical flight data, two traditional clustering algorithms are employed. When the spatial distribution is sparse, K-MEANS [48] is used, and when the distribution is dense, DBSCAN [15] is used. This point-based idea has a practical advantage of handling trajectories of unequal length or with gaps, which is common due to the low satellite coverage of AIS around the world. The main difference in relation to TREAD [41] relies on the fact that our work tries to directly map the rules defined by IMO into real data available from AIS readings.

Another main advantage of the clustering framework proposed in this thesis is that in the clustering process the IMO rules, particular Traffic Separation Schemes, are

taken into account, which is not the case in previous works. Therefore, our approach maps rules defined by IMO directly to the AIS data, and in that sense it is similar to the Candide [4] system, which is a digital face reconstruction mask. Candide uses a pre-defined digital face mask to be adjusted in a facial picture, and allows the generation of new facial expressions without any extra information. From the rule mapping perspective, this work does the same, with the difference that we use the set of rules defined by IMO and check if vessels are following them.

2.2 Anomaly Detection in Maritime Surveillance

The objective of anomaly detection is to find patterns in data that do not conform to expected behavior and these nonconforming patterns are often referred to as anomalies or outliers [10]. Anomalies or outliers are related to, but distinct from, noise in the data. As defined in [10], noise is some patterns or phenomena in data that are not of interest to the analyst, but act as a hindrance to data analysis. Thus, activities on handling these unwanted noise involve noise removal [57] and noise accommodation [49]. But in anomaly detection tasks, the outliers or the anomalies are the patterns which attract much attention from the analysts in various domains, such as fraud detection for credit cards [54], insurance [60], or health care [35], intrusion detection for cyber-security [34] and military surveillance for enemy activities [10]. In this thesis, our work focuses on detecting the anomaly navigational behaviors of the ships or vessels. Since noise filtering process is not involved in the pre-processing step, the noise caused by AIS devices' issues might also be detected as part of the anomalies.

Maritime Anomaly Detection techniques primarily fall into two categories: statistical modelling [30][46][13] and predictive modelling [41][39][29]. The general idea of statistical techniques for anomaly detection is to fit a statistical model for normal behaviors with the given data set and then apply a statistical inference test to determine if a previously unseen instance belongs to the model [10]. Approaches based on predictive models usually predict future status information (e.g. position, speed and course) of a particular vessel and then compare the real data with the prediction to decide the abnormality. Additionally, there is another type of approach trying to solve the problem based on predefined rules. One representative work was done by Roy *et al.* [50][51] and an expert system for automated anomaly detection was

introduced. As we know, one limitation about these kind of systems is that they are relying on the knowledge of predefined rules, as a result, much effort needs to be done to categorize anomalies. In the papers [50][51], the authors identified a taxonomy of kinematic and geo-spatial concepts in the vessels tracking domain and the taxonomy has to be updated once a new rule is found. Since not much work focuses on rule-based type of techniques, this survey of maritime anomaly detection is more about the other two main categories.

In the maritime domain, the majority of statistical models are built upon the momentary kinematic features (position, course, speed and acceleration rate) of individual vessels.

Laxhammar [30] used a Gaussian Mixture Model (GMM) and a greedy version of Expectation-Maximization (EM) for clustering. A GMM can be regarded as an ensemble model of K multivariate Gaussian distributions (mixture components). GMM is an extension of single Gaussian and it has the advantage of approximating arbitrarily complex distributions in arbitrarily high dimensions. The greedy EM-learning (greedy Expectation-Maximization-learning) algorithm, an extension to the classical EM-algorithm, is employed to determine the optimal number of components (or distributions) and the parameter set for all the distributions. So the greedy EM algorithm is basically an iterative model with mainly two steps in each iteration; the first step is to insert a new component and the second step is to apply classical EM until convergence. Then in order to label new data, the likelihood of the new point will be first calculated based on the probability distribution obtained from the clustering phase. Further, it will be classified as abnormal if the likelihood is below a certain predefined alarm threshold.

According to [31], GMM is not an optimal model because it needs an assumption that the distribution of vessel positions along the major axis of the sea lane segments is uniform. To tackle this issue, in [46], the authors propose to use adaptive Kernel Density Estimator (KDE) for estimating unknown probability densities and modelling arbitrary sea lanes. KDE, also known as the Parzen Window method, is a non-parametric model and the estimated PDF (Probability Density Function) can be determined based merely on the training data, which is superior to GMM. Then in the anomaly detection phase, the anomaly detector is sequentially applied to the

incoming data. The value of new incoming point's density is calculated under the null hypothesis (no anomaly) and this value is then compared with a detector parameter related to false alarms for deciding the new point's abnormality.

A comparison between the two approaches above [30] and [46] is given in [31], demonstrating that the anomaly detection results from both models are not satisfactory as the two models detect the anomalous segments after rather long distances (three kilometers and four kilometers respectively) while an expected effective anomaly detector should detect such behaviors at a shorter distance [31]. Another disadvantage of the two statistical methods is that it is hard to interpret in terms of anomaly types, even if they could find anomaly patterns promptly. Once a point is determined to be abnormal, the likelihood value in [30] and the density value in [46] cannot indicate the types of the abnormality. For example, the analysts cannot know whether the point deviates too far away from the normal lanes or the vessel sails too slow or too fast compared to normal speed.

In [13], Gerben *et al.* propose a technique based on Machine Learning models. In their work, different trajectory alignment kernels (Dynamic Time Warping [28] and Edit Distance [38]) are applied with one-class SVMs (Support Vector Machine) [11] for detecting the outlying trajectories. This trajectory-based method requires that the complete trajectory has been observed before it can be classified as normal or anomalous [45], which is not applicable for sequential anomaly detection in incomplete trajectories (real-time AIS surveillance), unlike our proposed point-based method.

Pallotta *et al.* [41] suggest the use of rule-based and low-likelihood models for anomaly detection. As stated in a previous paragraph, the rule-based techniques generate alerts based on a set of pre-defined rules. So the rules defined in this approach are similar to those in other knowledge-based work, which require maritime domain experts' knowledge. As an example, no specific rules defined for maximum speed in different sea regions can be found in IMO publications [1][25], so the maximum speed pre-defined in a port area can only be accurately estimated by a specialist knowledgeable about the area. Similar to the work done in [30], the low-likelihood anomaly detection aims at detecting deviations from the normal patterns or distributions derived from the training AIS data. For this low-likelihood detection, a Weibull model was employed (a parametric exponential-like model), along with a sliding time

window technique to avoid problems with incomplete and intermittent tracks.

Similar to the work done in [41], Nevell [39] proposes to use a Bayesian approach to predict the future route of a particular vessel for comparison. This methodology is based on a node-sparse network, built from different kinds of coastal nodes. The idea of pre-defining a global maritime traffic network is also adopted in [7] where the network is constructed based on the historical data. But instead of a Bayesian approach, Soleimani *et al.* [7] employ the A^* [14] algorithm to decide the optimal or expected route. In [29], the authors insist that an overall threat is indicated by a sequence of the individual behaviours. Therefore, five specific anomalies (deviation from standard route, unexpected AIS activity, unexpected port arrival, close approach with another ship and entering a zone known as illegal exchanges) are introduced to extend Nevell's work [39] to assess the probability of a higher-level threat based on a constructed Bayesian Network.

One issue with the work described in [39][29] is that it cannot incorporate speed into deciding if a trajectory is anomalous; instead the judgement is based only on position. Another problem is that a pre-defined network may not be applicable in many near-port regions due to the nature of port traffic. Traffic in near-port areas is usually variable and the vessels are not always following straight lanes (optimal routes in [39][29]). Fortunately, both of these problems are handled with our approach.

The anomaly detection model in this thesis uses the results of the clustering framework presented in Chapter 3, which generates normal moving patterns and arbitrary shapes of stopping areas. The proposed anomaly detection method is point based but is capable of handling trajectory tracks. For each track the algorithm will return an anomaly ratio. The ratio is based on three types of specialized distances to take position, direction and speed into consideration.

Chapter 3

Normal Traffic Patterns Extraction Model

Vessels consistently follow different movement patterns in different areas. For instance, cargo ships may travel along straight lines at high speed in the middle of the sea, while they may frequently adjust their directions at low speed in the port or offshore platform areas.

As a consequence, Spaccapietra [53] introduced a model to reason about trajectories, which is called stops and moves. So a trajectory can be treated as a sequence of moves and stops [53], and in the work done by [59][41][32], maritime trajectories analysis is conducted from these two different aspects as well, moving patterns and stopping patterns distinguished by a speed threshold. According to [32], the stop and move model has two main advantages: 1) stopping areas are interesting areas which need to be discovered; 2) stopping points are one important source of noise during trajectory clustering and stopping points do not include any motion information for path modelling. In addition to the above mentioned advantages, another benefit is that it can help to decrease the search space complexity during moving trajectories clustering since the stopping points are filtered out before applying the algorithm.

In this work, we employ a similar method to extract different normal patterns from the historical AIS dataset with a stopping SOG (Speed Over Ground) threshold of 0.5 knots. More specifically, we first divide the historical AIS data into two subsets: the set of moving points (with SOG not less than 0.5 knots) and the set of stopping points (with SOG less than 0.5 knots). The threshold of 0.5 knots is selected based on multiple experiments: we have tried the numbers from 0.1 to 1.0 by increments of 0.05. It shows that there is no significant difference between the clustering results. Since our approach is density-based, the high density stopping areas are mostly contributed by those points with 0 knots. Thus, the approach is not sensitive to the threshold selected by users. Instead of merely clustering the waypoints of the trajectory data [59][41], normal moving trajectories and arbitrary shapes of stopping regions will

be extracted based on the entire moving and stopping dataset in the specific area, respectively.

In this chapter, a definition for a trajectory in the maritime domain (Section 3.1) is first given and then the model for extracting normal patterns from the historical AIS data is proposed. For the case where SOG is not less than 0.5 knots (moving patterns), we propose DBSCANSD (Density-Based Spatial Clustering of Applications with Noise considering Speed and Direction) as the basic algorithm to detect the main traffic lanes within the data. The algorithm's output is a set of Gravity Vectors (GV), which are vectors formed by 5 features: average COG, average SOG, average Latitude, average Longitude and Median Distance. In Section 3.2, for the case where the SOG is less than 0.5 knots (stopping areas), the original DBSCAN [15] algorithm is executed because speed and direction are not important factors. Another type of vector is created as output, labelled Sampled Stopping Point (SSP), which is dependent only on the geographic shape of the region. (Section 3.3)

3.1 Representing Trajectory Data

Definition 1 (*Trajectory*) *A trajectory is a finite sequence $T = ((x_1, t_1), (x_2, t_2), \dots, (x_m, t_m))$. Each data point x_i corresponds to a multi-dimensional feature vector of a moving object at time point t_i , where $t_i < t_{i+1}$ for $i = 1, \dots, m-1$.*

A trajectory can be defined as a data type representing the movement of an object. A trajectory can be represented by a multidimensional time series [3] or a sequence of multi-dimensional points [33]. In this thesis, we modify the definition of its raw form (Definition 1) [45] to adapt it for our work. The new definition of a ship trajectory is as following:

Definition 2 (*Trajectory in Maritime Domain*) *A trajectory is a finite sequence $T = ((x_1, t_1), (x_2, t_2), \dots, (x_m, t_m))$ where x_i is a set of $\langle \text{Latitude, Longitude, COG, SOG} \rangle$ and t_i is the time-stamp.*

Here in Definition 2, each vector x_i is called a **trajectory point**, but in this thesis, we simply use **point** or **data point** where the context is clear.

3.2 Normal Moving Trajectories Extraction

As stated at the beginning of the chapter, given a data set of AIS historical tracking points, we divide them into stopping and moving parts based on the threshold of 0.5 knots. In this section, we present our moving trajectory extraction model.

As we know, compared to stopping patterns, the moving part plays a more important role in the anomaly detection task since there is a higher collision risk for a vessel sailing at high speed than at low speed in the stopping areas. To extract normal patterns from the moving AIS data, we first propose the new clustering algorithm DBSCANSD (Density-Based Spatial Clustering of Applications with Noise considering Speed and Direction).

3.2.1 Density-Based Spatial Clustering of Applications with Noise considering Speed and Direction

Density-Based Spatial Clustering of Applications with Noise considering Speed and Direction (DBSCANSD), shown in Algorithm 1, is a density-based clustering algorithm based on the algorithm DBSCAN [15]. The main observation behind our approach is that it is common for different types of ships to sail with different velocities in one similar area of the sea. For instance, the cruise speed of a cargo ship can be faster than a fishing ship. Even the same type of vessels can behave differently in relation to direction. Imagine that an oil tanker sails between two different countries. When the vessel is full, its speed is slower than when it is empty.

Definition 3 (*Eps-neighborhood of a point*) Given a database D of moving trajectory points in a specific area, the Eps – neighborhood of a point p , denoted by $N_\epsilon(p)$, is defined by $N_\epsilon(p) = \{q \in D \mid \text{dist}(p, q) < \epsilon\}$

As discussed in Chapter 2, the key idea of DBSCAN [15] is that for each point of a cluster the neighborhood of a given radius has to contain at least a minimum number of points. Here in this thesis, we adopt this idea but also consider two other factors, maximum speed variance ($MaxSpd$) and maximum direction variance ($MaxDir$). The intuition behind this is that the neighbors of a trajectory point should be not only near enough, but also with similar COG (Course Over Ground) and SOG (Speed

Over Ground). Thus, we can modify the definition of Eps-neighborhood (Definition 3) in [15] to Definition 4.

Definition 4 (*Eps-neighborhood of a trajectory point*) Given a database D of moving trajectory points in a specific area, the Eps-neighborhood of a trajectory point p , denoted by $N_\epsilon(p)$, is defined by $N_\epsilon(p) = \{q \in D \mid \text{dist}(p, q) < \epsilon \text{ and } |p.SOG - q.SOG| < \text{MaxSpd} \text{ and } |p.COG - q.COG| < \text{MaxDir}\}$

Note that $\text{dist}(p, q)$ is the Geographical Distance [58] between p and q , instead of Euclidean distance, because it is necessary to take in consideration the Earth's curvature to calculate distances.

Then in the following, we also give the formal definitions of other essential notions for our density-based algorithm. The definitions are changed from the ones originally defined for 2-Dimensional points in the algorithm of DBSCAN [15] to those for trajectory points.

Definition 5 (*Core trajectory point*) Given a database D of moving trajectory points in a specific area, a trajectory point p in D is called a core trajectory point w.r.t. the parameters of ϵ , MinPts , MaxSpd and MaxDir if $|N_\epsilon(p)| \geq \text{MinPts}$

Definition 6 (*Directly density-reachable*) Given a database D of moving trajectory points in a specific area, a trajectory point p in D is directly density-reachable from another trajectory point q in D w.r.t. the parameters of ϵ , MinPts , MaxSpd and MaxDir if $p \in N_\epsilon(q)$ and $|N_\epsilon(q)| \geq \text{MinPts}$

From Definition 6, we can find out whether point q is a core trajectory point according to Definition 5. Because $|N_\epsilon(q)| \geq \text{MinPts}$, q is a core trajectory point. If the point q is also directly density-reachable from the point p , the point q will be also a core trajectory point clearly. But when the point p is not a core point and it is directly density-reachable from core point q , the point p is called **Border Trajectory Point**.

Definition 7 (*Density-reachable*) Given a database D of moving trajectory points in a specific area, a trajectory point p in D is density-reachable from another trajectory point q in D w.r.t. the parameters of ϵ , MinPts , MaxSpd and MaxDir if there is a

chain of trajectory points $p_1, \dots, p_n, p_1 = q, p_n = p$ such that p_{i+1} is directly density-reachable from p_i ($1 \leq i < n$)

From Definition 7, we can find the trajectory point p can be either a border trajectory point or core trajectory point if p is density-reachable from point q . But the point q and other points (except p) in the chain must be core trajectory points. Therefore, this relation is transitive, but it is not symmetric [15].

Definition 8 (*Density-connected*) Given a database D of moving trajectory points in a specific area, a trajectory point p in D is density-connected to another trajectory point q in D w.r.t. the parameters of ϵ , $MinPts$, $MaxSpd$ and $MaxDir$ if there is a trajectory point o such that both p and q is density-reachable from o w.r.t. ϵ , $MinPts$, $MaxSpd$ and $MaxDir$

However, if there exist multiple trajectory points (p_1, \dots, p_n) that are density-reachable from one core trajectory point o , then points p_1, \dots, p_n can be all border trajectory points. In such a cluster, one border trajectory point (p_1) is obviously not density-reachable from any other border points (p_2, \dots, p_n). So it is necessary to define a relationship (Definition 8) for any pairs of two border trajectory points in the same cluster.

With the three density-based relations defined as above, it is easy to see the relationships between them. Figure 3.1 demonstrates these relationships. For example, if there are two trajectory points p and q , and p is directly density-reachable from q , p will be not only density-reachable but also density-connected from q . However, if p is density-reachable from q , it will be also density-connected from q but may be not directly density-reachable from q .

Now the concepts of **Cluster** and **Noise** (Definition 9 and Definition 10 respectively) can be defined based on the above three relations (Directly density-reachable, Density-reachable and Density-connected):

Definition 9 (*cluster*) Given a database D of moving trajectory points in a specific area, a cluster C w.r.t. ϵ and $MinPts$ is a non-empty subset of D satisfying the following two conditions:

1) $\exists p \in C, \forall q \in D$, if q is density-reachable from p w.r.t. the parameters of ϵ , $MinPts$, $MaxSpd$ and $MaxDir$, then $q \in C$

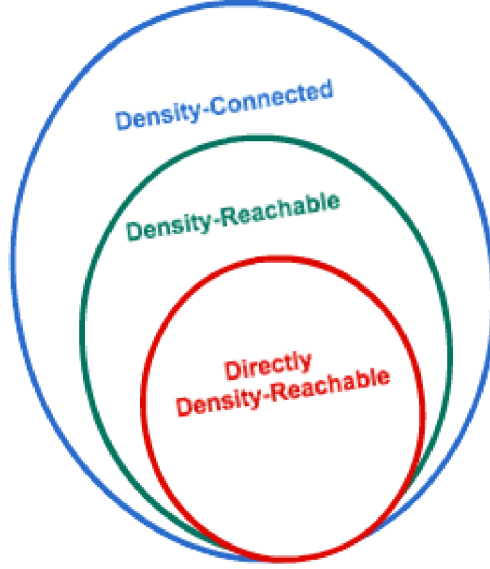


Figure 3.1: Three defined relations for the density-based clustering algorithm.

2) $\exists p \in C, \forall q \in D$, if p is density-connected to q w.r.t. the parameters of ϵ , $MinPts$, $MaxSpd$ and $MaxDir$, then $q \in C$

Definition 10 (noise) Given a database D of moving trajectory points in a specific area, and C_1, \dots, C_k are the clusters of D w.r.t. the parameters of ϵ , $MinPts$, $MaxSpd$ and $MaxDir$, $i = 1, \dots, k$, then the noise is the set of trajectory points in the database D not belonging to any cluster C_i , i.e. $noise = \{ p \in D \mid \forall i: p \notin C_i \}$

Algorithm 1 presents the procedures of DBSCANSD. It requires 5 parameters ($DatasetM$, eps , $MinPts$, $MaxDir$, $MaxSpd$) as input. $DatasetM$ is a list of all the moving points in the trajectories. eps and $MinPts$ are the reachable distance and reachable minimum number of points (see Ester *et al* [15]). The other two, $MaxDir$ and $MaxSpd$, are the two new parameters (maximum direction variance and maximum speed variance). The algorithm starts with a random selected trajectory point and then retrieves all its neighbour points. If its neighbour points' number is exceeding the threshold of $MinPts$, the point is marked as a core trajectory point and all its neighbours are adopted in the cluster. This procedure then iterates until all the points have been visited and the output is a set of clusters.

Algorithm 1 DBSCANS D

```

1: procedure DBSCANS D(DatasetM, eps, MinPts, MaxDir, MaxSpd)
2:   Mark all points in moving dataset DatasetM as unclassified
3:   clusterList  $\leftarrow$  empty list
4:   for each unclassified point P in DatasetM do
5:     Mark P as classified
6:     neighborPts  $\leftarrow$  queryNeighborPoints (DatasetM, P, eps, MinPts,
7:       MaxDir, MaxSpd)
8:     if neighborPts is not NULL then
9:       clusterList.add(neighborPts)
10:    for each cluster C in clusterList do
11:      for each cluster C' in clusterList do
12:        if C and C' are different clusters then
13:          if mergeClusters(C, C') is TRUE then
14:            clusterList.remove(C')
15:  return clusterList
16: procedure QUERYNEIGHBORPOINTS(data, P, eps, MinPts, MaxDir,
17:  MaxSpd)
18:  cluster  $\leftarrow$  empty list
19:  for each point Q in data do
20:    if distance(P,Q) < eps then
21:      if |P.SOG - Q.SOG| < MaxSpd then
22:        if |P.COG - Q.COG| < MaxDir then
23:          cluster.add(Q)
24:  if cluster.size > MinPts then
25:    Mark P as core point
26:    return cluster
27:  return NULL

```

```

26: procedure MERGECLUSTERS(clusterA, clusterB)
27:   merge  $\leftarrow$  FALSE
28:   for each point Q in clusterB do
29:     if point Q is core point and clusterA contains Q then
30:       merge  $\leftarrow$  TRUE
31:       for each point Q' in clusterB do
32:         clusterA.add(Q')
33:       break
34:   return merge

```

Time Complexity Analysis of DBSCANSD algorithm:

procedure QueryNeighbourPoints (lines 15-25):

The dominating operation is the `if` test (lines 18-20), data size: $n = \text{data.size}$. The implementation presented in Algorithm 1 has linear time complexity $O(n)$ where n is the size of the `data` set. However, if some spatial index is used the time complexity can be reduced to $O(\log(n))$.

procedure MergeClusters (lines 26-34):

The dominating operation is the `contains` test operation (line 29), the data size is the number of elements the second cluster $b = \text{clusterB.size}$. In the simple implementation presented in Algorithm 1 the worst time complexity is linear $O(b)$, if we assume the constant time cost of `contains` operation (this can be achieved by keeping a hash-set structure). However, this procedure can be accelerated with more sophisticated data structures, in particular avoiding adding elements one by one.

procedure DBSCANSD (lines 1-14):

Finally, the main procedure DBSCANSD is dominated by the two consecutive loops: point neighbourhoods computation (lines 4-8) and cluster merging (lines 9-13). The data size is the total number of elements $n = \text{DataSetM.size}$. The initialisation (line 3) is linear $O(n)$. Each of the two loops (lines 4-8 and lines 9-13) in the simple implementation presented in Algorithm 1 can be bounded with the quadratic time complexity $O(n^2)$, however, as we mentioned above, it is possible to accelerate it by applying more sophisticated data structures in the sub-procedures. To sum up, the

DBSCANS algorithm has at most quadratic time complexity $O(n^2)$ but there is definitely a room for improvement that will be further explored as the part of the future work.

We use external knowledge based on the Traffic Separation Schemes (TSS) [1] defined by IMO to adjust the parameters of the algorithm; therefore, we are using the lanes defined in IMO publications to fit the clustering results. As we know, once the area is determined, the two parameters *MaxDir* and *MaxSpd* should be fixed. The two parameters are relying on the corresponding port areas. During this thesis work, multiple pairs of values were tested to find the optimal *MaxDir* and *MaxSpd*. Our experiments indicate that 5° and 5 knot are the best parameters for both tested ports (Los Angeles Long Beach area and Juan De Fuca Strait area).

On the other hand, the parameters *eps* and *MinPts* must be defined for each port. Thus, another set of experiments were done to select the best pair for each area. For instance, when a lane generated by the clustering algorithm has any gaps compared to the one in the IMO TSS regulations, the analyst must increase the value of *eps* or decrease that of *MinPts*. In this way, the lanes with lower densities might be adopted as parts of the clusters to fill in the gaps. The two parameters *eps* and *MinPts* selected in this procedure will be directly used in the next stopping area extraction phase.

After applying DBSCANS to the ship trajectory dataset, those geographically close trajectory points with similar direction and speed will be grouped together to form a cluster. Then, we can not only get arbitrary shapes of clusters, but we can also separate those close points with great different normal SOGs and COGs into multiple clusters. Moreover, the algorithm can even treat a ship's acceleration or deceleration as a cluster, as long as the *MaxSpd* is well defined. Similarly, if the *MaxDir* is defined well, a curve shape of navigation will be also treated as one cluster. It should be noted that this point-based clustering approach can force the trajectory points of the same vessel to form multiple clusters if the vessel is not travelling along straight line or the vessel adjusts its speed significantly during navigating. But for one specific trajectory point, it can not be shared by multiple clusters, that is to say, it can only belong to one particular cluster.

3.2.2 Gravity Vector

Although the clustering results can reflect normal patterns of the vessels, it is not feasible to employ all the trajectory points in every cluster to decide the abnormality of the new arriving trajectory data because the time complexity of comparing the new trajectory and the clusters is $O(n * m)$ where n is the total number of the trajectory points in the clustering results and m is the number of points in the new arriving trajectory. As a result, an approach for cluster reduction is necessary.

Some techniques [12][37] that try to map the results to another representative object are proposed for compressing the original clustering results. For example, a spline-based clustering approach is introduced in [12], however, since our approach is a point-based clustering algorithm and the output is multiple sets of points, the idea of using splines is not applicable in this work. Another idea called centroid and envelope has been used for path modeling in vision-based trajectory learning [37]. The centroid can minimally specify the corresponding path and the envelope is used for denoting the path extent. In our case, we adopt a similar idea as centroid and envelope. Here we call it **Gravity Vector (GV)**.

A Gravity Vector is extracted by partitioning a cluster into multiple parts, therefore, each cluster can contain multiple Gravity Vectors. Figure 3.2 presents an example of calculating the Gravity Vectors of one cluster. To partition a cluster, a grid width needs to be first determined. The grid width (g in Figure 3.2) can be decided based on the domain knowledge or multiple experiments' results. Here in this thesis, we choose the length of eps used in Algorithm 1 as that of g . A Gravity Vector is a vector formed by five features: average COG, average SOG, average Latitude, average Longitude and Median Distance. Then one Gravity Vector GV_i can be denoted by:

$$GV_i = \langle COG_{avg}, SOG_{avg}, LAT_{avg}, LON_{avg}, D_{median} \rangle \quad (3.1)$$

As seen in Figure 3.2, all the blue arrows in the figure are belonging to the same cluster and we can see that all the trajectory points in this cluster have similar COGs (directions). Although it is a bit hard to see the speeds from the figure, we know each point in the cluster has minor SOG variation compared with its neighbors' SOGs. In this example, since the cluster is partitioned into 8 grids, the final output for representing the cluster will be a set of 8 Gravity Vectors. The steps of calculating

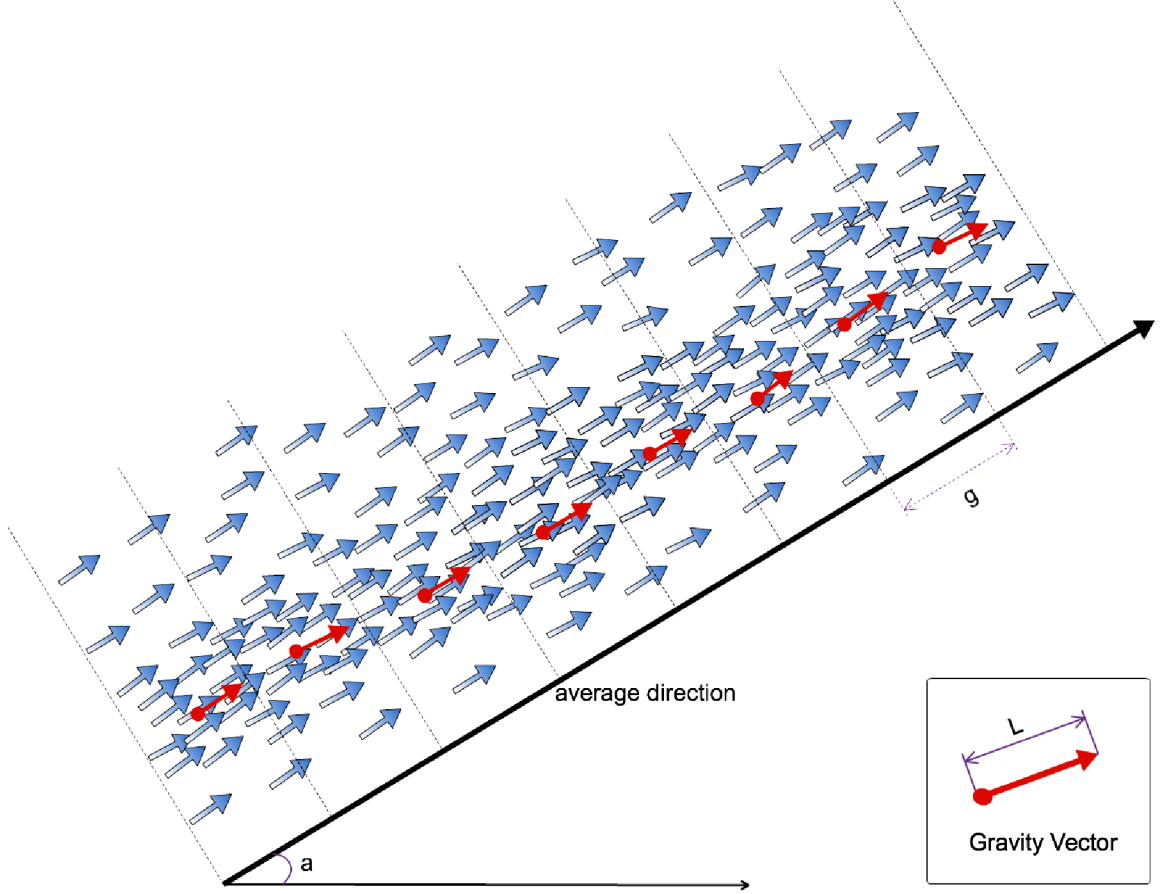


Figure 3.2: Calculate Gravity Vectors (GV) for one moving cluster. Blue arrows stand for the trajectory points of the cluster, red arrows are the final Gravity Vectors. The length (L) of the GV is the SOG_{avg} of the GV. The width of a grid g is the pre-defined length for partitioning the points.

the Gravity Vectors are as following:

Step (1) Calculate the average COG of the whole cluster and let it be COG_w . Note that COG_w is not COG_{avg} in equation 3.1.

Step (2) Partition all the trajectory points in the cluster along the direction of COG_w by a pre-defined grid width (g in Figure 3.2).

Step (3) For each grid of the partitioning results, generate the gravity vector of the grid.

In step (3), assume there are k trajectory points in one grid and TP_i is the i th trajectory point, we can calculate the first 4 features using the following 4 formulas (3.2–3.5).

$$COG_{avg} = \frac{\sum_{i=1}^k TP_i.COG}{k} \quad (3.2)$$

$$SOG_{avg} = \frac{\sum_{i=1}^k TP_i.SOG}{k} \quad (3.3)$$

$$LAT_{avg} = \frac{\sum_{i=1}^k TP_i.Latitude}{k} \quad (3.4)$$

$$LON_{avg} = \frac{\sum_{i=1}^k TP_i.Longitude}{k} \quad (3.5)$$

Then to calculate the last feature D_{median} , we should first calculate the distances between all k points in the grid and the average geographical point (LAT_{avg} , LON_{avg}) generated by formula 3.4 and formula 3.5. After this, we can apply a linear time complexity algorithm for computing median D_{median} based on Hoare's PARTITION algorithm.

In Lemma 1, we present the time complexity of calculating GVs for a cluster.

Lemma 1 *The worst time complexity of calculating a specific cluster's Gravity Vectors is $O(n)$, where n is the total number of points in the cluster.*

Proof. Assume there are n points in one cluster and the cluster is partitioned into m grids. Step (1) takes $O(n)$ time to calculate the average COG of the whole cluster. Step (2) takes $O(n)$ time to map all the points to the axis of average COG and $O(n)$ time to partition them into m grids. In step(3), to calculate one particular GV for the i th grid with k_i points, it takes $O(k_i)$ time to calculate the first 4 features and $O(k_i)$ time to calculate the median distance in linear time using one of the algorithms based on the Hoare's partition algorithm [24], for example the Blum-Floyd-Pratt-Rivest-Tarjan algorithm [8] that has linear worst time.

And this procedure needs to be repeated for m times, the total time for step (3) can be calculated as following:

$$W(n) = \sum_{i=1}^m O(k_i) = O(n)$$

Thus the pessimistic total time complexity for the whole procedure from step (1) to step (3) is $O(n)$. ■

The feature **Median Distance** can be used to measure the width variance of a moving cluster. The work done by Etienne *et al* [16] shows the choice of the statistical decile used to compute the spatio-temporal channel can give a tolerable estimate of this channel's width. However, we use median rather than ninth decile in [16] to provide a more robust width estimation metric. The distance is not the exact width of the channel, instead, it is a relative distance for outlier detection. More detailed information about how to use the GV for anomaly detection is presented in Chapter 4.

3.3 Normal Stopping Areas Extraction

In stopping areas, such as ports and wharfs, there may be different types of ships at anchor with zero-velocity or new arriving vessels entering the area with extremely low speed. In this case, vessels can stop with their prows pointing to any directions and change their headings frequently to arrive at the anchorage berths. Therefore, direction is no longer an essential part for analyzing the risk of collision. The possibility of a low-speed vessel colliding with another ship is very low; therefore this situation was excluded from the clustering process. In other words, only when a ship sails with a higher speed than the stopping threshold in a stopping area can we label the data as outlier. So the task of this stage is to identify the locations and geographical shapes of the stopping areas from the historical data set. Once the stopping areas are identified, the case of grounding (vessels are stopping unexpectedly in an area which is far from stopping areas) is also feasible to be detected.

Based on the analysis of stopping regions, the original DBSCAN [15] algorithm can be employed for stopping points clustering since speed or direction is no longer essential factors. The output is a set of stopping clusters with high density of trajectory points. Intuitively, there is a relatively large number of vessels gathered in each stopping cluster, which indicates the region is a reasonable place for ships to anchor.

It should be noted that it is common to have both stopping clusters and moving

clusters in a specific region at the same time, e.g., a shipping lane can pass through a stopping area. So the phenomenon that a vessel is approaching the port (stopping area) with a low speed is not supposed to be considered abnormal. Our approach can handle these kinds of issues because both types of clusters can be generated during the normal traffic patterns extraction step and the abnormality of a point will be decided based on both types clusters.

Although DBSCAN [15] can generate stopping points clusters, these original clustering results are still not adaptable for the next anomaly detection phase because the relatively large number of points in each cluster may pose serious constraints regarding the computational feasibility. Consequently, just like GV, a similar representative point called **Sampled Stopping Point (SSP)** is proposed. Algorithm 2 shows the whole procedure for generating the final SSPs given a set of stopping trajectory points in a specific region.

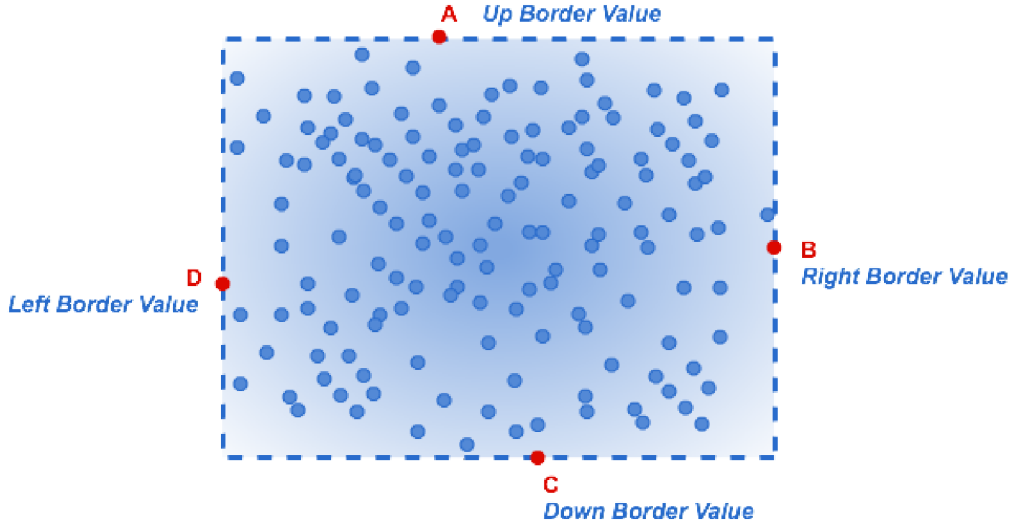


Figure 3.3: Border Values of a stopping cluster

As shown in Algorithm 2, stopping points are first clustered by the DBSCAN [15] algorithm (line 1). The same parameters (eps and $MinPts$) used by DBSCANSD (Section 3.2) are still employed in this stopping areas extraction phase. Then the algorithm starts to generate the SSPs for every cluster. The key idea about SSP is to find the representative points that reflect the shapes and locations of the stopping areas. But before extracting the points from the clustering results, an estimation of

Algorithm 2 Extract SSP From Stopping Dataset

Input: The list of the stopping points of the region, *DatasetS*; Reachable distance, *eps*; Reachable minimum number of points, *MinPts*

Output: The list storing the region's SSP, *resultSSP*

```

1: StoppingPointsClusters  $\leftarrow$  DBSCAN(DatasetS, eps, MinPts)  $\triangleright$  see DBSCAN
   algorithm in [15]
2: resultSSP  $\leftarrow$  empty list
3: for each cluster C in StoppingPointsClusters do
4:   lat1, lat2  $\leftarrow$  minimum and maximum of all the points' Latitude values in
   cluster C
5:   lon1, lon2  $\leftarrow$  minimum and maximum of all the points' Longitude values in
   cluster C
                                      $\triangleright$  estimate the sample size for cluster C
6:   area  $\leftarrow$   $|(\textit{lat1} - \textit{lat2}) * (\textit{lon1} - \textit{lon2})|$ 
7:   if area = 0 then
8:     sample_size  $\leftarrow$  1
9:   else
10:    sample_size  $\leftarrow$  Ceiling(area / ( $\pi * \textit{eps}^2$ ))
                                      $\triangleright$  sample the stopping area points
11:    count  $\leftarrow$  0
12:    while count < sample_size do
13:      randomly select one point P from cluster C
14:      if P is far from all points in resultSSP then
15:        resultSSP.add(P)
16:        count++
17: return resultSSP

```

the number of SSPs needs to be conducted. To address this issue, lines 4–10 first estimate the area of the region and then calculate the sample size. Lines 4 and 5 calculate the four border values of the cluster. Border values are extreme values of one cluster. Figure 3.3 shows an example of border values in a stopping cluster. In Figure 3.3, 4 border points A-D are in red color and Up Border Value is the maximum of all points' Latitude values while Down Border Value is the minimum. Then B and D can be extracted in a similar way as A and C. Any two or three of them (or even all the four points) can be the same point. It is noteworthy that if one cluster crosses the Greenwich meridian, the way to calculate the two border values, $lon1$ and $lon2$, will be changed.

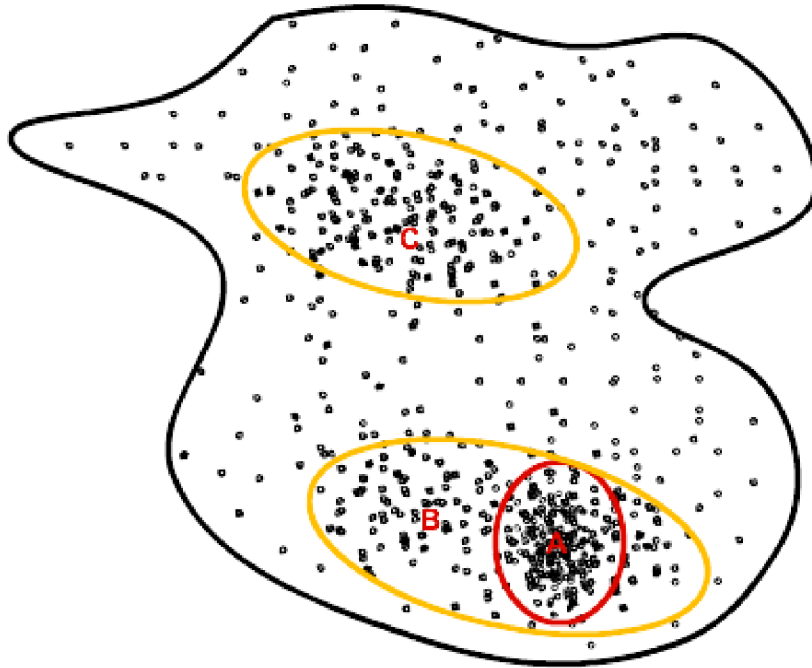


Figure 3.4: An example of a stopping area where simple random sampling cannot work. There are three regions (A,B and C) with higher density than the whole area's average density.

After extracting the 4 border values, the area of the border (blue border lines shown in Figure 3.3) is calculated. With this area result and the eps used for clustering, the sampling number is estimated (lines 6–10). Lastly, we use a random selection strategy to sample the stopping points. However, there still exists some sub-regions with higher density of points compared with other sub-regions in the same cluster and as a result, random selection tends to select points from these higher-density

micro-regions.

An example is shown in Figure 3.4. Figure 3.4 shows a stopping area with three higher density sub-regions and all the black points are the stopping points of the clustering result. We can see that the sub-region A has the highest density while B and C have relatively lower density which are still higher than the whole average. So in this example, if we simply employ a random selection technique, the majority of the final points shall be extracted from region A and then regions B and C . There might be only a limited number of points that are sampled from the rest of the whole stopping area, which is not supposed to happen in terms of representing the whole area's shape. As a consequence, we need to modulate the sampling process to stratified sampling to give each sub-region equal probability to be selected. This can be easily done with lines 11–16; the distance condition in line 14 can be defined by *eps* as well. The algorithm has linear time complexity.

Thus, with the Gravity Vectors and Sampled Stopping Points extracted from the moving and stopping AIS data separately, the anomaly detection task can be conducted, the techniques for which are presented in the next chapter.

Chapter 4

Anomaly Detection

In this chapter, we present our anomaly detection model. Similar to the normal traffic pattern extraction model, we also treat stopping and moving separately. A ship trajectory in a particular area, especially in the near-port areas, can consist of both stopping points and moving points. Hence, given one incoming trajectory data set (one sequence of trajectory points with same identification), every point in this trajectory should first be labeled as stopping or moving according to the SOG threshold (0.5 knots) and then be checked for abnormality. In the following part of the section, three division distances are first presented in Section 4.1 and then our abnormal detection model is introduced in Section 4.2.

4.1 Three Division Distances

In this section, three specialized division distances are proposed to label every data point: Absolute Division Distance (ADD), Relative Division Distance (RDD) and Cosine Division Distance (CDD). *ADD* is employed in the stopping points abnormality detection phase while *RDD* and *CDD* are used for the moving part. The definitions of the three division distances are given in the following. Here we use **target points** to represent the points of a new coming trajectory to be labeled.

The **Absolute Division Distance** between a target point P_t and a Sampled Stopping Point P_s is defined as:

$$D_{absolute} = Distance((P_t.Lat, P_t.Lon), (P_s.Lat, P_s.Lon)) \quad (4.1)$$

As can be seen in Definition 4.1, ADD is actually the Geographical Distance [58] between Latitude and Longitude values of the target point(P_t) and the sampled stopping point(P_s).

The **Relative Division Distance** between a target point P_t and a Gravity Vector GV is defined as:

$$D_{relative} = \frac{Distance((P_t.Lat, P_t.Lon), (GV.Lat, GV.Lon))}{GV.MedianDistance} \quad (4.2)$$

For moving trajectory points, we adopt RDD ($D_{relative}$) rather than ADD ($D_{absolute}$) because the normal moving lanes could have different widths according to their surrounding geographical environment. The routes in a narrow strait can be more cramped than those in the open sea. Relative distance, the ratio of a point's distance from the centroid to the median distance of all the points in one cluster from the centroid, is one efficient metric in clustering-based approaches for detecting outliers [56]. We replace the centroid of each cluster with our Gravity Vectors, which can be regarded as surrogates for a centroid when we only consider Longitude and Latitude. As stated in [36] and previous section, every cluster can have more than one GV, in which case the median distance of one GV can be used to measure the width variance of a moving cluster. Another work that confirms this approach is presented by Etienne *et al.* [16] in which it is shown that the choice of the statistical decile used to compute the spatio-temporal channel can give a tolerable estimate of this channel's width.

The next division distance is CDD, which is employed to involve direction and speed of moving objects. The **Cosine Division Distance** between a target point P_t and a Gravity Vector GV is defined as:

$$D_{cosine} = \cos \alpha \times \frac{\min(P_t.SOG, GV.SOG)}{\max(P_t.SOG, GV.SOG)} \quad (4.3)$$

where:

α is the angle between the two directions, that is, the difference between P_t 's COG and GV 's COG.

The intuition behind CDD is to combine the angle ($\leq 180^\circ$) between the two directions (angle α is defined by COG differences between P_t and GV) and the difference between the two speeds. Figure 4.1 shows two abnormal cases that consider

COG and SOG. In Figure 4.1(a), we can see that the speeds of the two vectors are in the same length L while the angle a is too large. This could be explained as a ship crossing a normal lane in a nearly opposite direction, which may cause a collision. In Figure 4.1(b), the speed of the target point P_t is much lower than GV 's although they have similar COG. In this case, the slowly sailing vessel may also be in a high risk of collision in relation to other ships moving at a normal rate of speed. As a consequence, Cosine Division Distance (CDD) is proposed. It can be easily seen that $\cos \alpha$ accounts for the angles' difference and the ratio between two SOGs can reflect the speeds' difference.

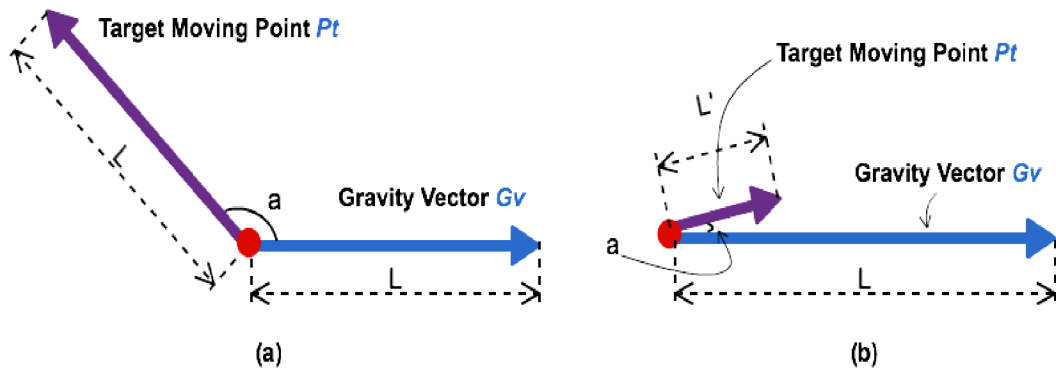


Figure 4.1: Two Abnormal Cases after considering COG and SOG

The Cosine Division Distance proposed in this work combines both COG and SOG in its calculation. A possible alternative is to calculate the distance of each component separately and then check the normality, but this alternative may increase the misclassification error. One example is presented in Figure 4.2. In this scenario, it is assumed that a certain point P_t is in the middle of two different grids belonging to different clusters. If this approach of separately calculating each component is employed, the target point, P_t , will first be considered normal with respect to G_v' on the basis of direction, and also with respect to G_v on the basis of length (value of speed). However, P_t is an abnormal point because it is in the direction of G_v' , but with a much faster speed.

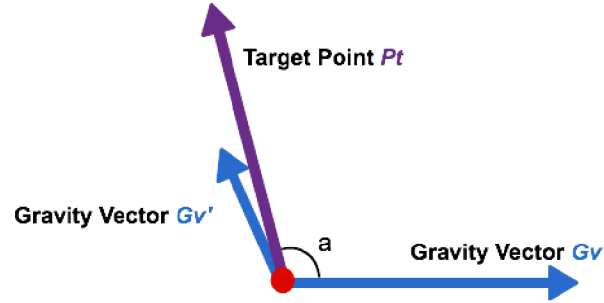


Figure 4.2: One case that the simple combination of the two division distances cannot work. There are two moving clusters in one specific area and the two clusters are with two different average COGs.

4.2 Anomaly Detection Model

After having the three division distances defined, we can present our anomaly detection model.

As shown in Algorithm 3, the anomaly detection process is completed in two steps. Lines 3-6 employ ADD to label the stopping points of the trajectory and Lines 7-14 use RDD and CDD to decide the moving points' labels. Lastly, the ratio of abnormal points to the number of all points is returned as the value of abnormality (Lines 15-18). The abnormality can be interpreted as a confidence ratio in a trajectory point being abnormal and it is beneficial for the end users to choose a level of confidence while retrieving the abnormal results of the whole system.

This algorithm requires three thresholds as inputs, which can be estimated through experiments. In this work, we first take out one month of data in a specific area and calculate all the division distances, then select specific values based on distribution and quartile analysis.

Figure 4.3 shows the distribution of the maximum CDD of moving points in the area of Juan de Fuca Strait. Figure 4.4 illustrates the CDF (Cumulative Distribution Function) of the maximum CDD in the area. The red dotted line in Figure 4.3 is the 5% separation line, which means the percentage of all the CDDs not greater than the separation coordinate is 5%. Similarly, in Figure 4.4, the red dotted line is also a 5% separation line and we can see the corresponding coordinate in y-axis is 0.05. It is obvious that over 90% of the data points' CDD are greater than 0.5. Thus, in this case, we can choose 0.5 or a smaller value as the CDD threshold for this area and

Algorithm 3 Detect abnormality of the target trajectory

Input: (1) The target trajectory dataset, D ; (2) The lists of Sampled Stopping Points and Gravity Vectors from the previous model, SSP and GV ; (3) Three thresholds, $add_threshold$, $rdd_threshold$ and $cdd_threshold$

Output: The abnormality rate, $abnormality$

- 1: Separate D into two sub-datasets based on the speed threshold, moving dataset D_m and stopping dataset D_s
 - 2: Initialize all labels of points in D_m and D_s as Normal
 - ▷ label stopping points of the target trajectory
 - 3: **for** each data point S in D_s **do**
 - 4: $ADD_s \leftarrow \text{minimum}(ADD(S, SSP))$
 - 5: **if** $ADD_s > add_threshold$ **then**
 - 6: $S.label \leftarrow \text{Abnormal}$
 - ▷ label moving points of the target trajectory
 - 7: **for** each data point M in D_m **do**
 - 8: $RDD_m \leftarrow \text{minimum}(RDD(M, GV))$
 - 9: **if** $RDD_m > rdd_threshold$ **then**
 - 10: $M.label \leftarrow \text{Abnormal}$
 - 11: **else**
 - 12: $CDD_m \leftarrow \text{maximum}(CDD(M, GV))$
 - 13: **if** $CDD_m < cdd_threshold$ **then**
 - 14: $M.label \leftarrow \text{Abnormal}$
 - 15: $count_ab \leftarrow$ the number of the abnormal points in D
 - 16: $count_all \leftarrow$ the total number of all points in D
 - 17: $abnormality \leftarrow count_ab / count_all$
 - 18: **return** $abnormality$
-

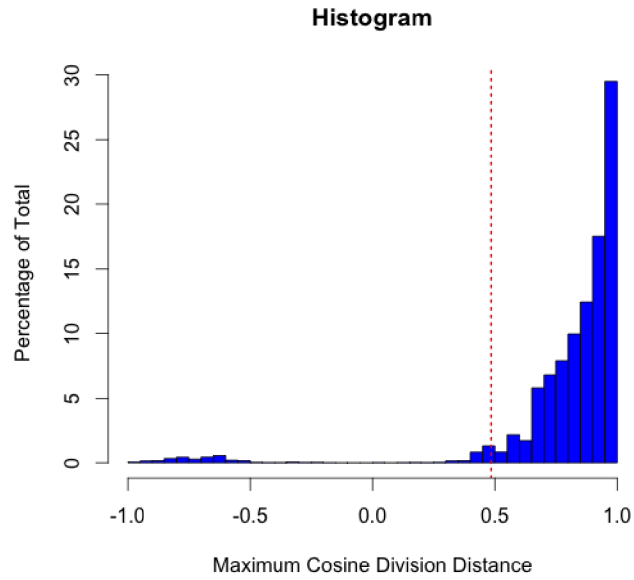


Figure 4.3: The maximum CDD distribution in the area of Juan de Fuca Strait.

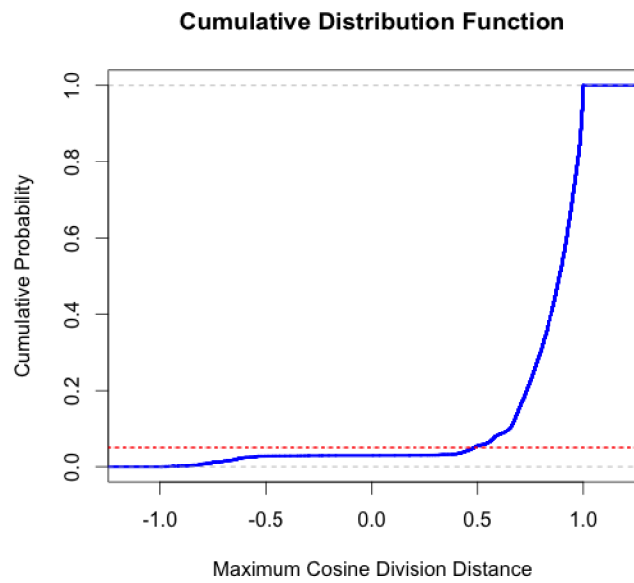


Figure 4.4: The Cumulative Distribution Function of the maximum CDD in the area of Juan de Fuca Strait.

then when a new trajectory dataset is given, we can use this value as the basis for anomaly detection.

Chapter 5

Evaluation

In this chapter, experiments are done to evaluate the effectiveness of the proposed approach. First, to evaluate the normal traffic patterns extraction model, regions of Juan de Fuca Strait (Section 5.1.1) and Los Angeles Long Beach (Section 5.1.2) are selected and the results are presented in Section 5.1. Then in Section 5.2, we conducted two other experiments in the same region of Juan de Fuca Strait. The first one is conducted with the non-labeled data while the second one is done after labelling the data. The results of the first experiment (Section 5.2.1) are shown visually and the second experiment (Section 5.2.2) compares our models results with the labels by the expert.

5.1 Normal Traffic Patterns Extraction Model

In this section, we evaluate the effectiveness of our maritime normal patterns detection model in two regions. The data set contains non-anonymized messages from ships in the region of Juan de Fuca Strait and the region of Los Angeles Long Beach. The data set is in CSV (Comma Separated Values) format. Every row in the data set represents one particular trajectory point of a vessel and every point includes hundreds of attributes (columns). But in this work, we use five attributes: MMSI, Longitude, Latitude, SOG and COG. The movement rules for port areas defined by IMO extracted from the IMO publication [25] are shown in Figure 5.1 and 5.4. Evaluation of the normal patterns detection model is related to the problem of evaluating clusters quality and this is always not a trivial task. Our strategy is to use our approach to generate the lanes and then compare the results with these rules.

Both experiments are conducted on a laptop of MacOS with an Intel Core i5 (1.3GHz) processor and a RAM of 4GB (1600 MHz DDR3). The implementation of the algorithm is a single-threaded, single-process Java program.

5.1.1 Juan de Fuca Strait

The Strait of Juan de Fuca separates the south coast of Vancouver Island from the north coast of State of Washington. The entrance of it lies between Cape Flattery (48°23'N.,124°44'W.) and Carmanah Point (48°37'N.,124°45'W.) [9]. Figure 5.1 shows a map of this area and the predefined rules for it.

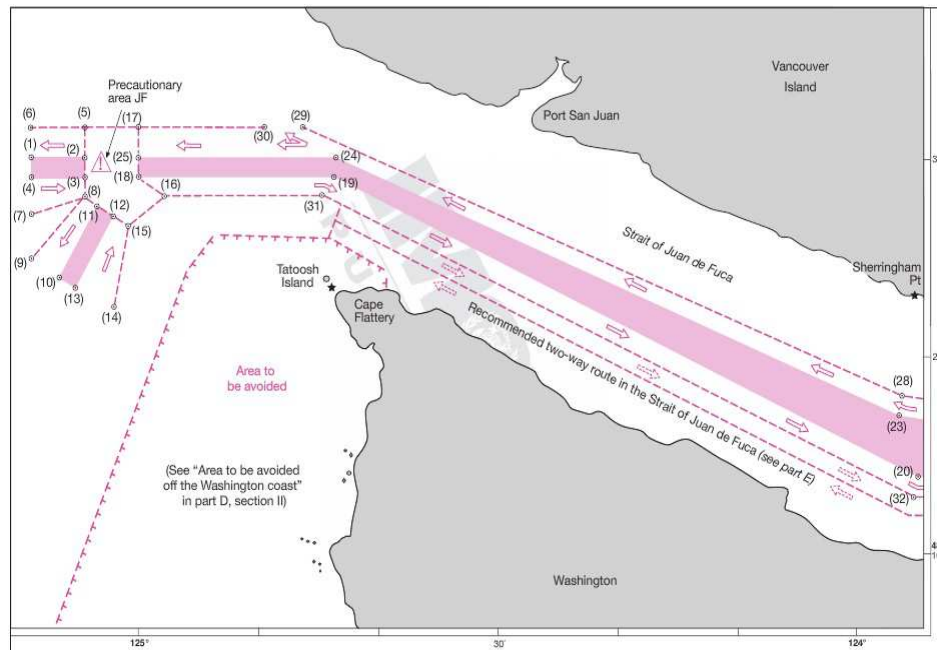


Figure 5.1: Juan de Fuca Strait and its approaches (west) [25]

The data set prepared for this experiment is two-months of trajectory data from November 1 to December 31 in 2012. It consists of 67,850 trajectory points. The whole dataset is not used; instead 46,000 records (40,000 moving points and 6,000 stopping points distinguished by the SOG threshold 0.5 knot) are selected randomly. As can be observed that there are much more moving points than stopping points, which is caused by the geographical environment of the area. Since this is a strait region with only one small stopping area (shown in dark green in Figure 5.2), it is more common for a vessel to approach to or exit from the port in this strait area.

After applying the two clustering algorithms, 14 different clusters including 13 moving clusters and 1 stopping cluster are extracted. The result can be seen in Figure 5.2. Points in blue are the original traffic points while others are the cluster points. The size (total number of the points composing the clusters) of moving clustering

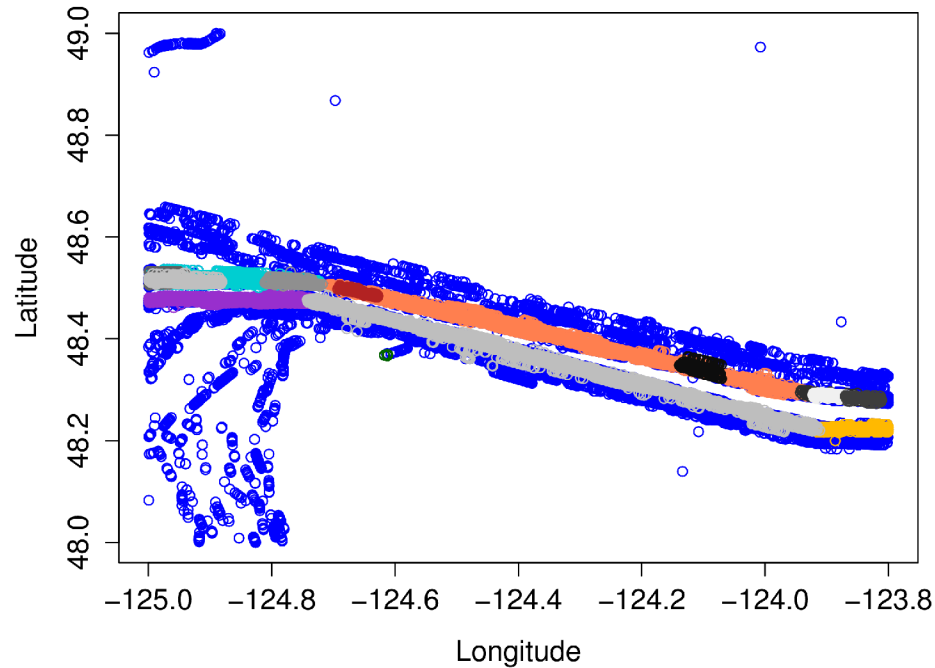


Figure 5.2: The clustering results of JUAN DE FUCA STRAIT area. Points in blue color are the noise. Other different colors stand for different clusters. The cluster in dark green is a stopping cluster.

results is 20,817 while that of the stopping cluster is 4,800.

Figure 5.3 shows the gravity vectors of the moving clusters and the sampled stopping point of the stopping cluster. After this extraction phase, only 302 points are generated which include 301 GVs and only one SSP. From the two figures we can see that the clustering results can reflect the traffic trends in this strait and the GVs as well as the SSP can represent the clustering results too. The total running time of this experiment is 520.7 seconds (403.8 seconds for moving points clustering and 116.9 seconds for stopping points clustering).

5.1.2 Los Angeles Long Beach

The approach has also been tested in the area of Los Angeles Long Beach, where the traffic is more complicated and heavier. Figure 5.4 shows a map of the area with two rules defined. In this experiment, the same period (November 1 to December 31

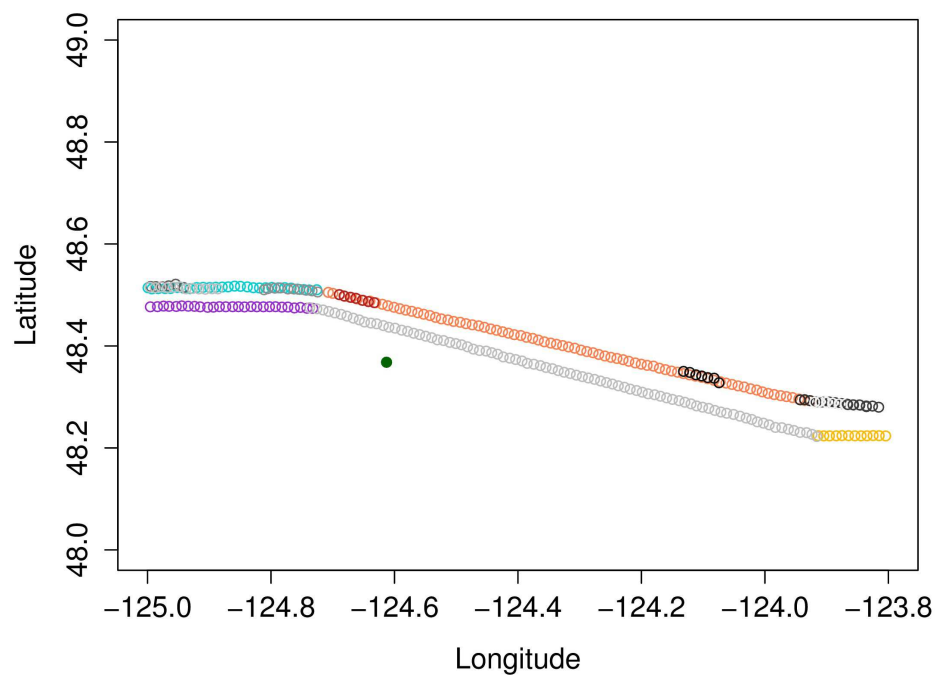


Figure 5.3: The Gravity Vectors and Sampled Stopping Points extracted from the clusters in Juan de Fuca Strait area.

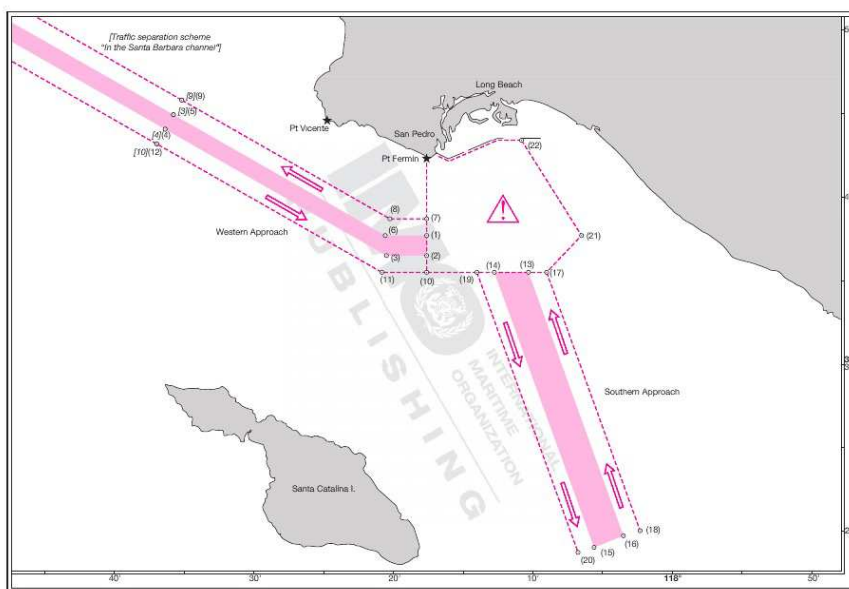


Figure 5.4: Map of Los Angeles Long Beach and the navigation rules defined [25]

in year 2012) dataset has been prepared. There are 327,694 records (99,937 moving points and 227,757 stopping points) in this dataset. We adopt all the moving points in this dataset for moving clusters generation and the first 20,000 stopping points for stopping area generation.

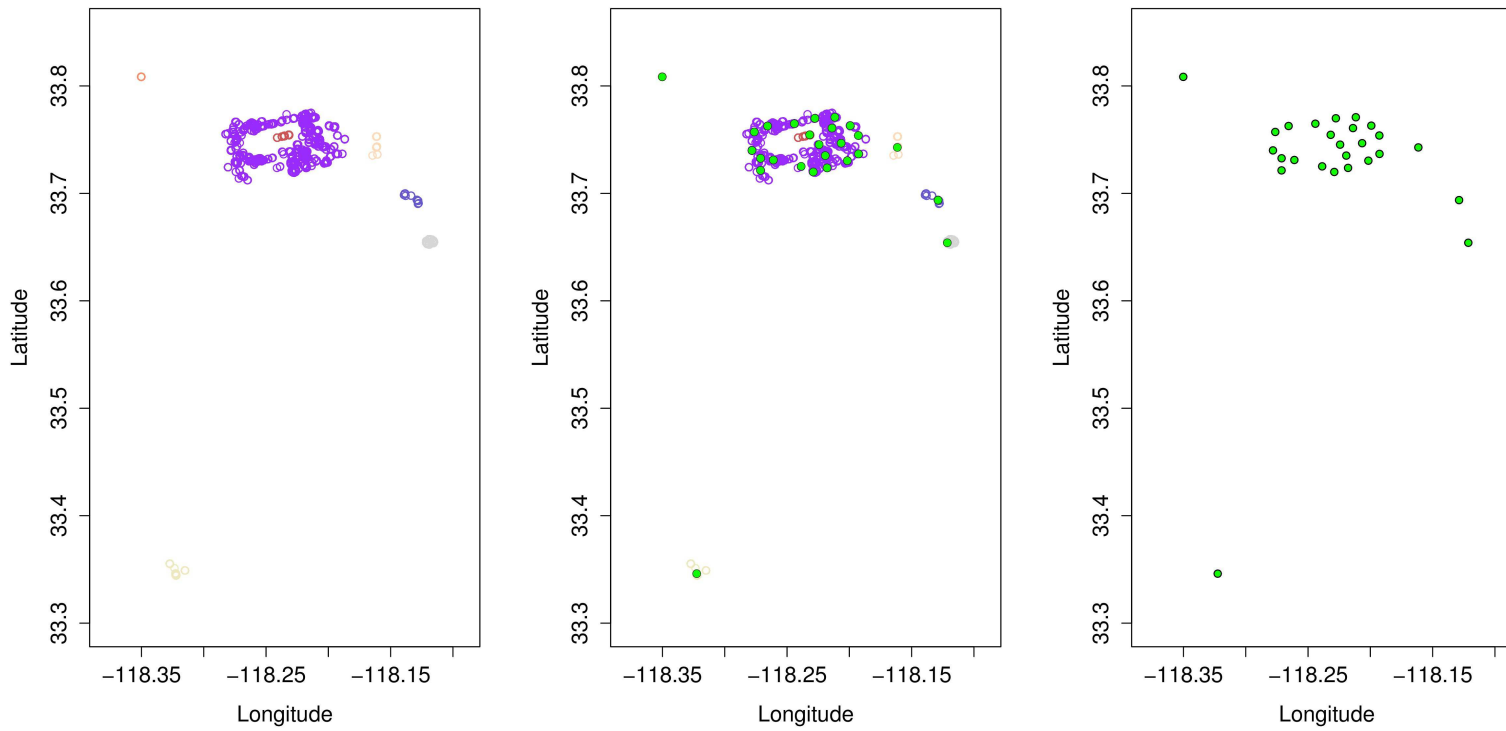
Figure 5.5 shows the result of stopping clusters sampling procedure in the area of Los Angeles Long Beach. There are 7 different stopping clusters in Figure 5.5(a) with 19,089 stopping points generated by DBSCAN algorithm [15]. Then after applying Algorithm 2, only 26 Stopping Sampled Points (shown in green color in Figure 5.5(b) and 5.5(c)) are selected with excellent quality in terms of the representativeness.

Then the algorithm DBSCANS is applied to moving points. Ultimately, 48,404 points are selected to form 51 moving clusters from the original 99,937 points. This result is shown in Figure 5.6. From the figure we can see that in the stopping area, there are multiple moving clusters too. This is because vessels in the area have to change their headings frequently to arrive at the specified anchor location and our algorithm will treat this curve-shape movement as multiple moving clusters with slight COG differences. The last step for this normal traffic pattern extraction is to calculate the Gravity Vectors of the moving clusters. The GVs of the area can be seen in Figure 5.7 and the SSPs (filled circles in dark green color) extracted in the previous step are also shown in the same figure. The total running time of this experiment is 669.8 seconds (290.8 seconds for moving points clustering and 379.0 seconds for stopping points clustering).

5.2 Anomaly Detection Model

In this section, we evaluate the effectiveness of our maritime anomaly detection model in the region of Juan de Fuca Strait. The evaluation work contains two parts, the first one is conducted with the non-labeled data while the second one is done after labelling the data. The results of the first experiment are shown visually and the second experiment compares our model's results with the labels by the expert.

The data set which was prepared for normal traffic patterns extraction phase comprises two months of trajectory data from November 1 to December 31 in 2012 and contains 67,850 trajectory points. The whole data set is not used for extracting normal patterns, instead 46,000 records (40,000 moving points and 6,000 stopping



(a) Stopping clusters by DBSCAN [15] (b) Stopping clusters and Sampled Stopping Points generated (c) Sampled Stopping Points generated

Figure 5.5: Extract Stopping Sampling Points (SSP) from stopping clusters in Los Angeles Port Area. The stopping clusters are first extracted using DBSCAN [15] algorithm shown in (a) and different colors stand for different clusters. Then Algorithm 2 is used for getting the SSPs (shown in figures (b) and (c))

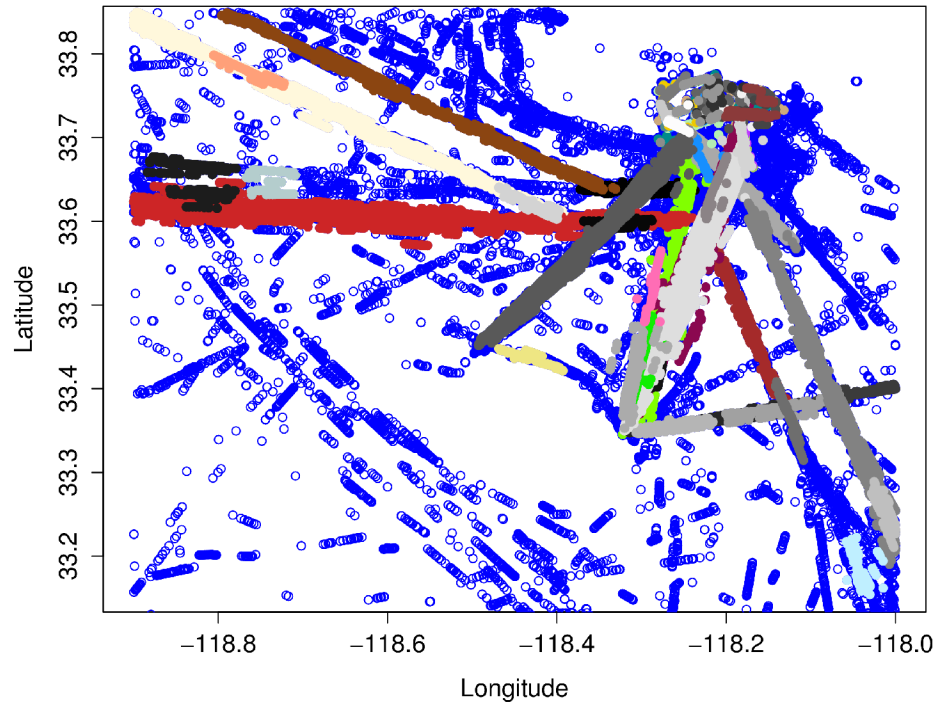


Figure 5.6: The clustering results of Los Angeles Long Beach area. This only shows the moving points (in blue) and the corresponding moving clusters (in different colors).

points distinguished by the SOG threshold 0.5 knots) are selected and then the rest of the data set is used for estimating the three thresholds for the anomaly detection phase.

Afterwards, to evaluate our anomaly detection model in both experiments, we chose the first half of January (January 1st to January 15th) in 2013, as our target trajectory data set. This second dataset consists of 284 different trajectories with 17,431 points.

5.2.1 Experiment On Unlabeled Data Set

After applying the normal traffic extraction model, 16 different clusters, including 15 moving clusters and 1 stopping cluster are identified. Figure 5.8 shows the gravity vectors of the moving clusters and the sampled stopping point of the stopping cluster. After this extraction phase, only 389 points are generated which include 388 GVs and only one SSP.

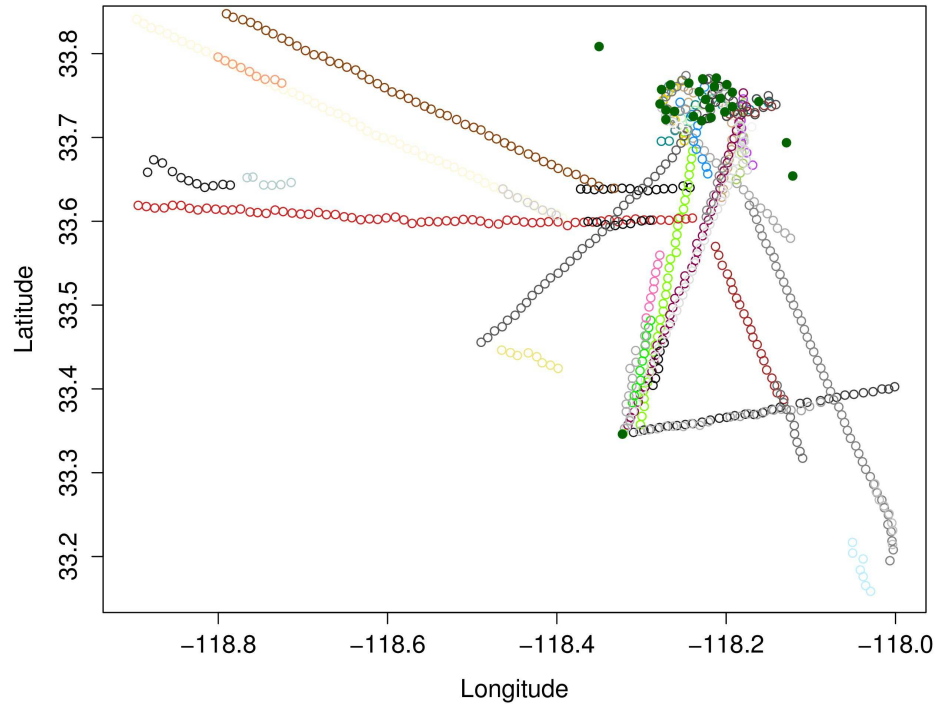


Figure 5.7: The Gravity Vectors and Sampled Stopping Points extracted from the clusters in Los Angeles Long Beach area.

The next step before detecting anomalous trajectories is to estimate the three thresholds ($add_threshold$, $rdd_threshold$, $cdd_threshold$) to be used in Algorithm 3. As stated before, the remaining 21,850 trajectory records are chosen for this phase. There are 10,825 stopping points and 11,025 moving points in this subset.

The quartile values of ADD and RDD of the subset are shown in Table 5.1 (column 1 and column 2). The authors tested different thresholds to be used as the anomaly detector. After various tests, for the area of Strait of Juan de Fuca the best threshold value was to consider 95% of the data as normal in relation to distance, and from this sub-set another 95% of the data to be considered normal in relation to speed and direction. The main objective is to reduce the number of false alarms (vessels considered abnormal, while they are normal), and this filtered data will later be evaluated by a human expert that will give the final decision. The model is flexible to allow changing this threshold value depending on the geographical area under evaluation.

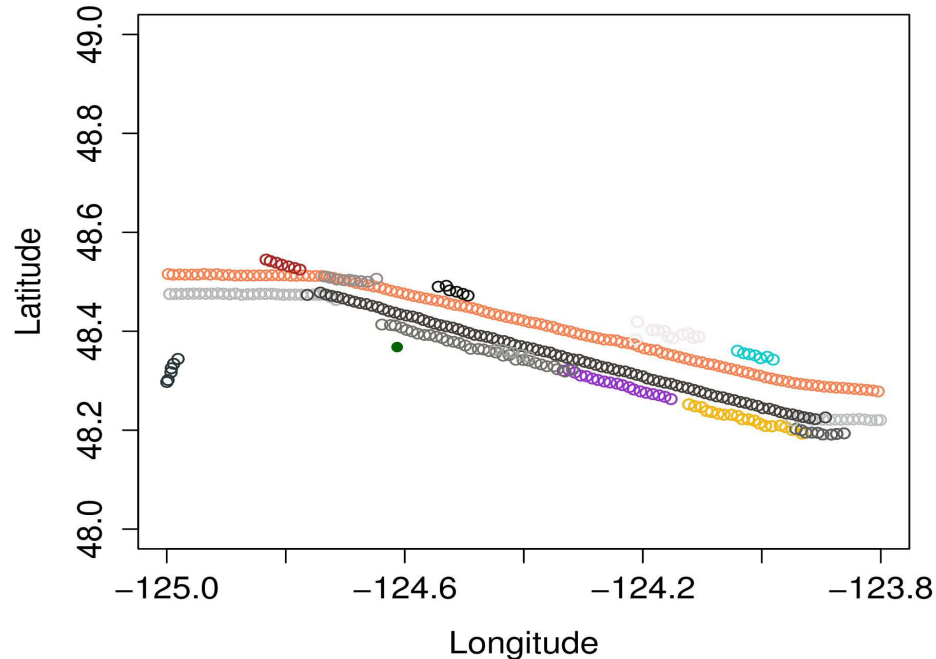


Figure 5.8: The Gravity Vectors (open circles) and Sampled Stopping Points (filled circles) extracted from the clusters in JUAN DE FUCA STRAIT area.

So we choose the sample quantiles of 0.95 for both ADD and RDD. The corresponding thresholds in this case, *add_threshold* and *rdd_threshold* in Algorithm 3, are 97.290 and 5.938. After calculating the RDD threshold, the statistic for CDD is obtained (shown in the 3rd column in Table 5.1). Then we select 0.05 as the possibility to decide the third threshold (0.485) which can be employed as our CDD threshold (*cdd_threshold* in Algorithm 3).

With the extracted normal patterns and the thresholds estimated, we start to evaluate the capacity for detecting abnormal trajectories.

In this step, we first apply our model to the trajectory data points; the labeling results of the data points are shown in Figure 5.9. The red points stand for the GV and the SSP in this area. The green points are normal, while the blue and purple ones are abnormal. More specifically, a blue point means it is too far away from the corresponding GV or SSP, while a purple point represents that its speed or direction is aberrant in the specific location. In this case, 1,534 points (872 in blue and 662 in

Table 5.1: Quartile statistics of the three division distances

| Statistic | ADD | RDD | CDD |
|--------------|----------|----------|---------|
| Min | 0.13 | 0.00537 | -0.9937 |
| 1st Quartile | 3.00 | 0.70300 | 0.7642 |
| Median | 4.46 | 1.04000 | 0.8876 |
| Mean | 36.97 | 1.81500 | 0.8104 |
| 3rd Quartile | 6.89 | 1.58400 | 0.9612 |
| Max | 44250.00 | 52.86000 | 0.9999 |

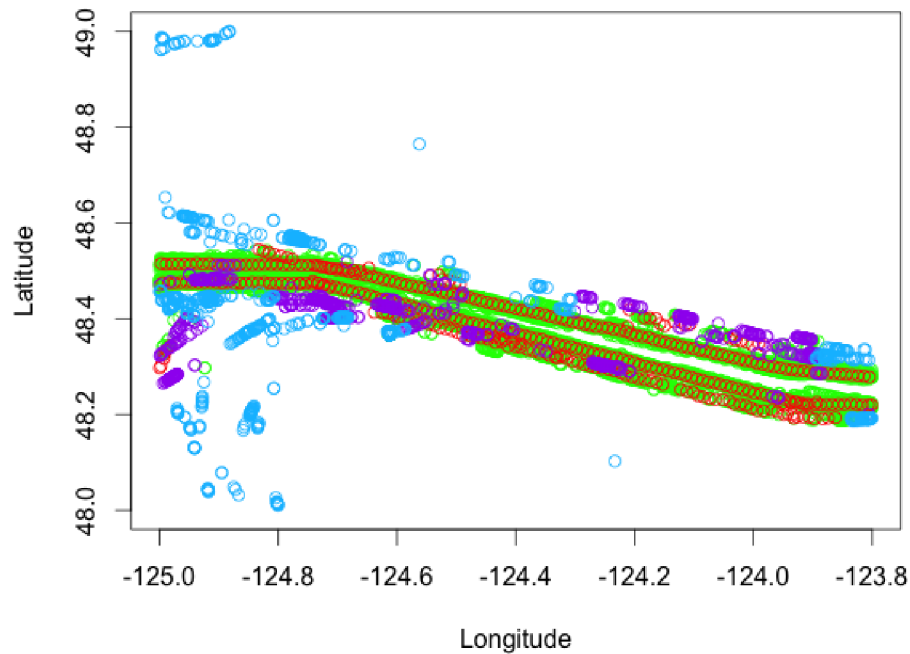


Figure 5.9: The anomaly labeling results of the trajectory data points in JUAN DE FUCA STRAIT area. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD).

purple) of the 17,431 points are finally considered as abnormal.

After finishing the labeling process, we calculate the abnormality ratio of each trajectory. We use a threshold of 0.5 as the minimum confidence rate to extract the abnormal trajectories. Noteworthy, the number 0.5 is adjustable and was chosen

based only on the experiments to reduce the presence of false alarms (normal trajectories considered abnormal by our model). The result is that 22 trajectories are labeled abnormal among the total 284 trajectories. In other words, the abnormality rates of these 22 trajectories are over 0.5. From Figures 5.10 to 5.15, six examples of abnormal cases are illustrated. As can be seen in the figures, both purple points and blue points contribute to the final abnormality of one trajectory.

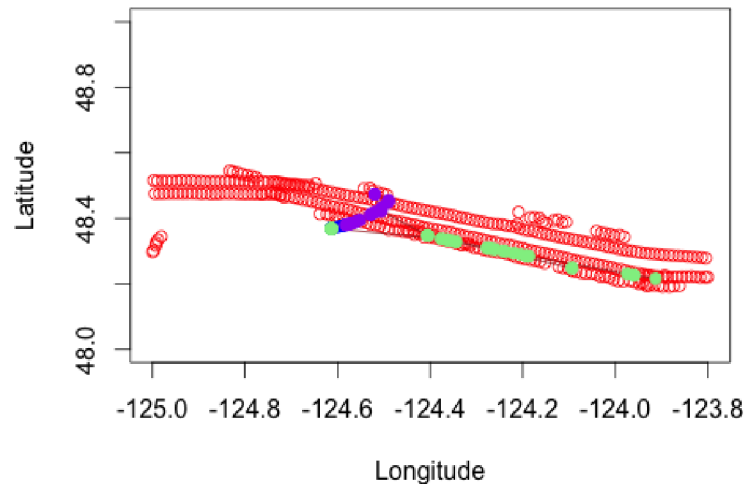


Figure 5.10: First example of the abnormal trajectories detected by our framework. GVs and SSPs are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD).

5.2.2 Experiment On Labeled Data Set

In this experiment, the same data set is labeled by an expert who has multiple years of experience in maritime data analysis. The labelling process is not biased by our model’s results since only the raw AIS data has been provided to the expert.

A track division method is employed by the expert before labelling. More specifically, the AIS data points are divided into distinct tracks on the basis of vessel ID (MMSI), and where temporally sequential points are separated by no more than 4.5 minutes of time. Thus one track, as defined above, may be divided into multiple-sub

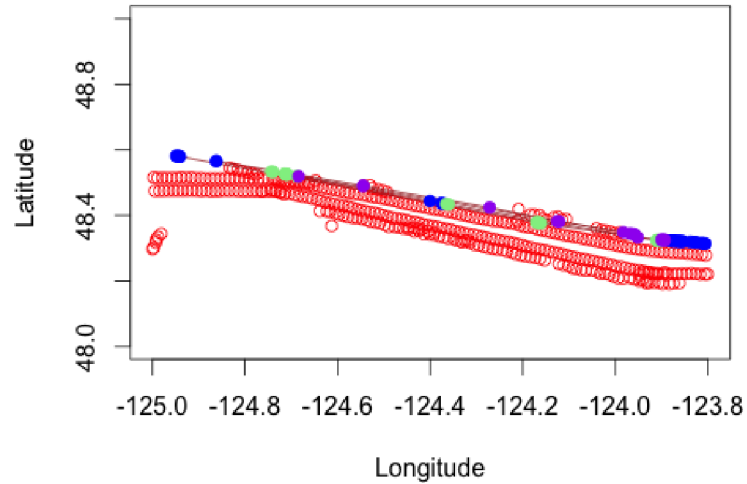


Figure 5.11: Second example of the abnormal trajectories detected by our framework. GV's and SSP's are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD).

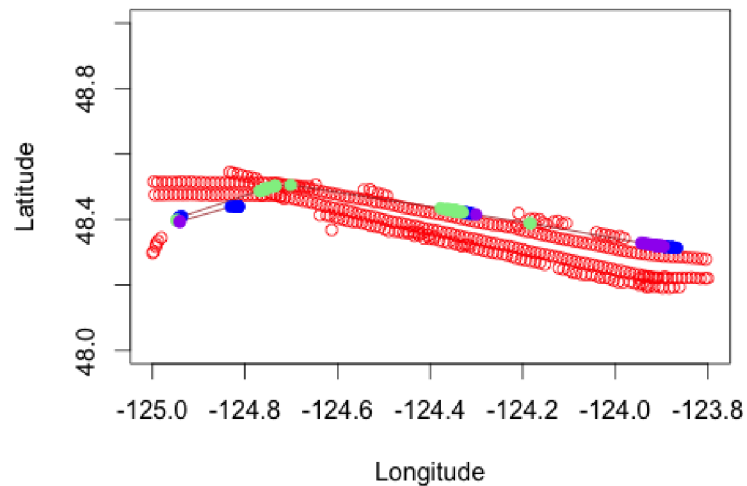


Figure 5.12: Third example of the abnormal trajectories detected by our framework. GV's and SSP's are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD).

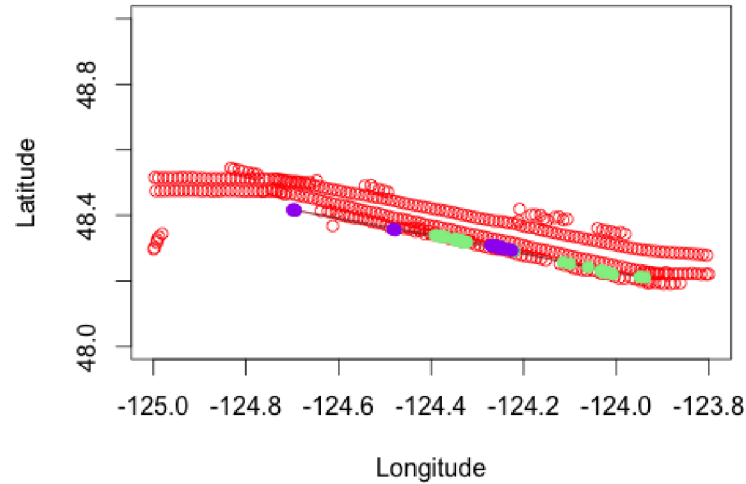


Figure 5.13: Fourth example of the abnormal trajectories detected by our framework. GV's and SSP's are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD).

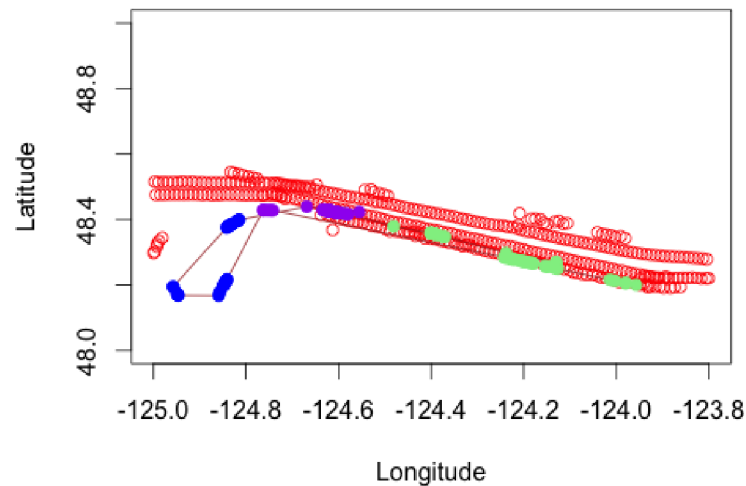


Figure 5.14: Fifth example of the abnormal trajectories detected by our framework. GV's and SSP's are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD).

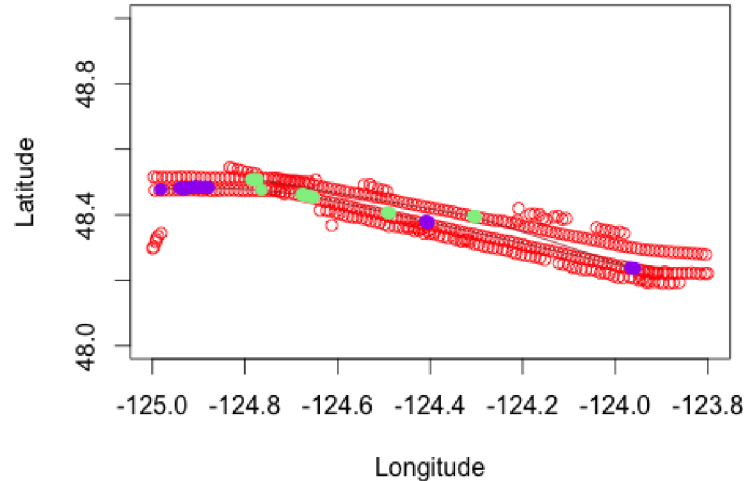


Figure 5.15: Sixth example of the abnormal trajectories detected by our framework. GV's and SSP's are in red and the normal points are in green. The two types of abnormal points are in blue (abnormal in relation to ADD or RDD) and purple (abnormal in relation to CDD).

Table 5.2: Labels and descriptions by the expert

| Label | Description |
|-------------------|--|
| <i>Bad_Pos</i> | Track contains questionable point, far outside track, looks like bad GPS return |
| <i>In_Excl_Zn</i> | Track has significant portion within the exclusionary zone between traffic lanes |
| <i>XING_TSS</i> | Track appears to be crossing lanes of TSS [1] |
| <i>XING_NShor</i> | Track appears to be crossing lanes of near shore two way traffic area |
| <i>Odd_Mvmt</i> | Track shows unusual movement without other explanation |
| <i>Leave_Lane</i> | Vessel was in traffic lane, then veered outside |
| <i>Harbour</i> | Track seems to describe in-harbour navigation or moored vessel |
| <i>Normal</i> | Normal Movement |

tracks on the basis of time. Additionally, this process can result in tracks comprised of single points. For the one-point track case, the length of the track is 0 nautical mile and the track is not assigned with any labels. The final labels and their descriptions

provided by the expert are shown in Table 5.2.

The labels assigned by the expert are not based on the points; instead, they are based on the whole sub-tracks. In other words, as long as one sub-track shows an anomalous pattern, the whole set of points inside the sub-track will be assigned as one kind of abnormal label. Another noteworthy point is that the expert has not taken SOG (speed) into account except for Harbour behavior during his labeling process. That is, whether the speed of the vessel is too fast or too slow near the lane is not considered, but our algorithm can take this into account.

From Table 5.2, we can firstly assume the label of *Normal* as normal patterns based on the description. Then the label *Harbour* can also be considered as a normal case because it is reasonable for a vessel to moor in harbour. Lastly, we observe that all the sub-tracks with the label of *Leave_Lane* only have tiny changes from their route and they still navigate strictly within the normal lanes. So we also classify the label of *Leave_Lane* as normal.

After dividing the 284 tracks (284 different MMSIs), 2,122 sub-tracks are generated. Among them, 680 sub-tracks contain only one point (Length=0 nautical mile). Then 14 sub-tracks are classified as abnormal labels (other than *Leave_Lane*, see Table 5.2) by the expert and the remaining 1,428 tracks are all normal patterns. Thus we can see that the data set is a highly imbalanced data set which can make our work extremely challenging.

At this point we can use our algorithm to label the data set and compare the results with the expert's labels. To compare the results, we first apply the same division method to separate the tracks. We can then use a threshold to decide the whole sub-track's label. In this experiment, we employ 60% as the threshold value. That means that if the portion of abnormal points in one track is greater than 60%, we will label this whole track as abnormal. Using this approach, we find that 131 sub-tracks are classified as abnormal and the remaining 1,311 sub-tracks are normal. Table 5.3 is the confusion matrix for the experiment.

From Table 5.3 we can see that 4 sub-tracks are classified as abnormal by both the expert and our model and 1301 sub-tracks are classified as normal by both too. On the other hand, another 10 sub-tracks are labeled as abnormal by the expert while normal by our model. The remaining 127 sub-tracks are classified as abnormal by our

Table 5.3: Confusion Matrix I

| | Abnormal (Our Model) | Normal (Our Model) |
|---------------------------|-------------------------|-----------------------|
| Abnormal (Expert's Label) | 4 | 10 |
| Normal (Expert's Label) | 127 | 1301 |

Table 5.4: Confusion Matrix II

| | Abnormal (Our Model) | Normal (Our Model) |
|---------------------------|-------------------------|-----------------------|
| Abnormal (Expert's Label) | 4 | 10 |
| Normal (Expert's Label) | 52 | 1376 |

model while normal by the expert. And the overall accuracy of this detection result is 90.49%.

To improve the result, we investigate the 14 sub-tracks designated as abnormal by the expert. Among these, we find that the 4 which are further designated as abnormal by the algorithm are so labeled solely because of their direction. After investigating other tracks, we find that even if the ships deviate far from the lane the expert may still label them as normal. The intuition behind this is straight-forward, as the labelling process is based on Traffic Separation Scheme (TSS) [1] boundaries and the expert cannot affirm that a trajectory point far from the TSS Boundaries is abnormal. Considering this fact, in the following experiment, we ignore the abnormal labels caused by RDD or ADD during the evaluation and we choose a lower threshold for deciding whole sub-tracks' labels. In the previous experiment we use 60% while we choose 10% here. The reason why we use this smaller threshold is because we have much less abnormal labels. If we use 60%, we cannot detect any abnormal tracks. It should be noted that the threshold can be adjusted based on the input from domain experts. In real-time application, it is not necessary to have the threshold while labelling the new incoming points instead of tracks.

Table 5.4 presents the improved results and we can see that the overall accuracy has been increased from 90.49% to 95.70% while keeping the same recall for abnormal cases.

Chapter 6

Conclusion

In this thesis, a maritime traffic pattern extraction model has been first proposed. The approach first separates the trajectory data set into moving and stopping subsets, and then employ different strategies in corresponding subsets. The main advantage of this work is that two attributes, speed and direction, are taken into account during the clustering phase. In this way, geographically close trajectory points with similar direction and speed can be grouped together to form a cluster. We also present two methods to represent the results of clustering; that is, Gravity Vector for moving clusters and Sampled Stopping Point for stopping clusters.

Another contribution of the proposed normal traffic extraction model is that it can be easily extended for other scenarios. Besides the maritime anomaly detection task addressed in this thesis, the moving trajectories clustering algorithm (DBSCANSD) and Gravity Vectors are also applicable for other trajectory clustering tasks (e.g. vehicle position data, animal movement data, hurricane monitoring data). Thus, more experiments can be done with other domains' data sets to verify our method's applicability in the future study.

To show the effectiveness of our approach, two real data sets from different regions (Juan de Fuca Strait area and Los Angeles Long Beach area) are used. We have compared the generated results with the rules defined by IMO and it shows that the results can be successfully mapped to the rules. This demonstrates that our model can effectively mine normal patterns in the two areas, which is another contribution of this work.

One limitation of the normal traffic extraction model is that it is sensitive to parameters. The model requires four parameters (eps , $MinPts$, $MaxDir$, $MaxSpd$) during the moving traffic patterns extraction phase and two parameters (eps and $MinPts$) for stopping area generation. So, if the model is applied in another area, a new set of parameters needs to be given to achieve a good clustering result. Thus,

to achieve satisfactory clustering results, it is necessary to possess a good maritime background. Specifically, when we distinguish moving points from stopping points in the first step, a reasonable SOG threshold is required. Here in this thesis, we adopt 0.5 knots as the threshold for distinguishing stopping and moving. Besides, before applying DBSCANSD to moving trajectory points, we cannot simply assume the tolerated angle (*MaxDir*) between two close points' directions and the tolerated difference (*MaxSpd*) between two close points' speeds.

The Gravity Vectors and Sampled Stopping Points detected by the proposed approach can help authorities update their rules (e.g. re-position the buoy markers). Some other research studies like route planning and vessel position prediction can also be conducted based on our model's results.

Based on the Gravity Vectors and Sampled Stopping Points extracted from the clustering phase, we then propose the anomaly detection model for maritime traffic data. An abnormality detection algorithm is presented based on three division distances (Absolute Division Distance, Relative Division Distance and Cosine Division Distance). This model is a fairly straightforward point-based approach, and is capable of handling complicated maritime traffic situations. One advantage is that the clustering process is associated with TSS Boundaries [1], which can assure a reliable clustering result for the following anomaly detection work. Another critical advantage is that besides position information (Longitude and Latitude), the model can also take speed and direction into account while deciding the abnormality of a single trajectory point. The model is also flexible enough for analysts to set their own thresholds to label whole trajectories.

In order to evaluate the effectiveness of the anomaly model, a highly imbalanced data set from Juan de Fuca Strait area is used. There are 2,122 trajectories while only 14 of them are abnormal (imbalance rate $\approx 0.66\%$). Fortunately, as shown in Table 5.4, our model can detect 28.57% of the abnormal tracks while maintaining a relatively high overall accuracy (95.70%).

6.1 Future Work

One of the possible directions of future work is to improve the efficiency of the presented algorithms by applying more sophisticated data structures (like spatial indexes,

for example). This issue is important in the context of big sizes of the processed data. Moreover, there exist a number of publications to address the scalability issue of DBSCAN using MapReduce, such as [22][43]. Thus, DBSCANSD is also flexible to be implemented in parallel as it is extended from DBSCAN. One idea is to first partition the AIS data to multiple sets based on the Longitude and Latitude and then apply our clustering algorithm on each sets. If the data is partitioned well, the time complexity can be reduced significantly which can make it feasible for a much bigger area's data clustering. In this thesis the focus was not on optimising the efficiency of the algorithms but to present the ideas. We envisage algorithm optimisation in the future work.

One may argue about the use of an unsupervised learning approach instead of supervised learning methods. The reason for this is that there is insufficient labeled data to train an effective model. Once we get enough labeled data we will explore other supervised learning techniques for this anomaly detection task. For example, we can try some classification algorithms designed for handling imbalanced data sets and the proposed specialized division distances can be used as the features of the classification models. Then comparisons between those results and our model's outcomes could be done to illustrate the effectiveness of the proposed division distances.

One limitation of the anomaly detection evaluation model is that the labelling process applied by the expert does not consider speed while our work takes this into account. This leads to another possible future direction, that is, more work should be done by the experts while labelling the data set to consider speed. In this way, the false alert rate can also be reduced while the identification of the anomaly trajectories is improved. Another limitation is that the experiments are only conducted with data from Juan de Fuca Strait area and as a result, more experiments in other regions should be done to better illustrate the effectiveness of our approach.

As discussed before, this is a point-based method and the output is two sets of representative vectors (Gravity Vectors and Sampled Stopping Points). As a consequence, one limitation of our clustering algorithm, compared to trajectory-based ones, is that ours cannot take into account the behavior over time. One example of this limitation is illustrated in Figure 6.1. Assume there are two main lanes in the region and they are Lane *A* and Lane *B* (shown in red color and purple color separately).

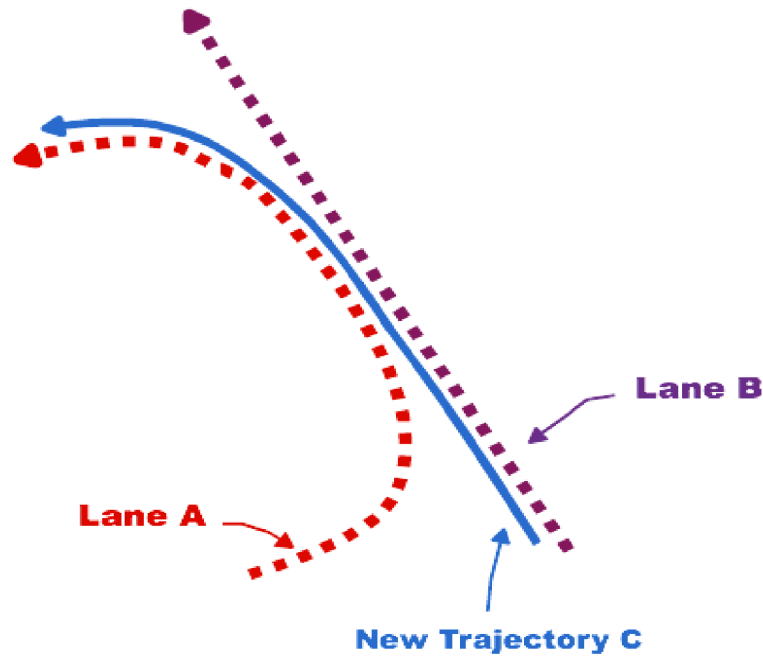


Figure 6.1: One situation that our algorithm cannot detect the anomalies. Two normal lanes (Lane *A* and Lane *B*) detected by DBSCANSD are in red and purple. A new anomaly trajectory *C* is in blue.

Then a vessel enters the region and follows a path such as Trajectory *C*. Since our framework treats the trajectories as sequences of data points, the trajectory will be divided into two sub-sequences and during anomaly detection phase, the first part will be associated with Lane *B* while the second part will be associated with Lane *A*. In this sense, the new trajectory will be assigned as normal. However, the vessel is not supposed to change its navigational lane once it starts its trip on Lane *B* because there are only two shipping routes pre-defined in this region. So in this situation, our model fails to detect the new anomaly trajectory *C*, and this can be improved after considering time, which is one aim of our future work.

The clustering results of the first phase can be used for predicting vessels' routes. Since the GVs extracted carry not only the position information (Longitude and Latitude) but also the kinematic characteristics (course and average speed), it is feasible for researchers to obtain the corresponding expected speed and heading for the specific predicted position. Two potential ways for handling this can be association rules mining techniques or Bayesian approaches. As stated in Chapter 2, anomaly detection approaches based on predictive models usually predict future status information

of a particular vessel and then compare the real data with the prediction to decide the abnormality. Thus, to improve the performance of the model, an ensemble model which incorporates the predictive model can be developed in the future. Specifically, once a trajectory point needs to be labeled, we can consider both its anomalous score (output of our proposed model) and its deviation from the predicted position to get a more confident result.

Bibliography

- [1] Colreg.2/circ.57 new and amended existing traffic separation schemes. http://www.imo.org/blast/blastDataHelper.asp?data_id=14761&filename=57.pdf. Accessed: 2006-05-26.
- [2] *IALA Guidelines on the Universal Automatic Identification System (AIS), Volume 1, Part II-Technical Issues Edition 1.1*. IALA/AISM-20ter rue Schnapper, 78100 Saint Germain en Laye, France, 2002.
- [3] Charu C Aggarwal and Chandan K Reddy. *Data Clustering: Algorithms and Applications*. CRC Press, 2013.
- [4] Jrgen Ahlberg. Candide-3 - an updated parameterised face. Technical report, 2001.
- [5] Mihael Ankerst, Markus M Breunig, Hans-Peter Kriegel, and Jörg Sander. Optics: ordering points to identify the clustering structure. In *ACM Sigmod Record*, volume 28, pages 49–60. ACM, 1999.
- [6] Heather Ball. *Satellite AIS for Dummies*. Wiley, Mississauga, ON, 2013.
- [7] Casey Hilliard Behrouz Haji Soleimani, Erico N. De Souza and Stan Matwin. Anomaly detection in maritime data based on geometrical analysis of trajectories. In *Information Fusion, 2015 18th International Conference on*. IEEE, 2015.
- [8] Manuel Blum, Robert W. Floyd, Vaughan Pratt, Ronald L. Rivest, and Robert E. Tarjan. Time bounds for selection. *J. Comput. Syst. Sci.*, 7(4):448–461, August 1973.
- [9] United States. Defense Mapping Agency. Hydrographic/Topographic Center. *Sailing Directions (enroute) British Columbia*. Pub. (United States. Defense Mapping Agency. Hydrographic/Topographic Center). Defense Mapping Agency, Hydrographic/Topographic Center, 2012.
- [10] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(3):15, 2009.
- [11] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- [12] Anders Dahlbom and Lars Niklasson. Trajectory clustering for coastal surveillance. In *Information Fusion, 2007 10th International Conference on*, pages 1–8. IEEE, 2007.

- [13] Gerben Klaas Dirk De Vries and Maarten Van Someren. Machine learning for vessel trajectories using compression, alignments and domain knowledge. *Expert Systems with Applications*, 39(18):13426–13439, 2012.
- [14] Rina Dechter and Judea Pearl. Generalized best-first search strategies and the optimality of a*. *J. ACM*, 32(3):505–536, July 1985.
- [15] Martin Ester, Hans peter Kriegel, Jrg S, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. pages 226–231. AAAI Press, 1996.
- [16] Laurent Etienne and Thomas Devogele. Spatio-temporal trajectory analysis of mobile objects following the same itinerary. *Advances in Geo-Spatial Information Science*, 17(1):11–34, 2012.
- [17] Nivan Ferreira, James T. Klosowski, Carlos Eduardo Scheidegger, and Cláudio T. Silva. Vector field k-means: Clustering trajectories by fitting multiple vector fields. *CoRR*, abs/1208.5801, 2012.
- [18] Zhouyu Fu, Weiming Hu, and Tieniu Tan. Similarity based vehicle trajectory clustering and anomaly detection. In *Image Processing, 2005. ICIP 2005. IEEE International Conference on*, volume 2, pages II–602. IEEE, 2005.
- [19] Maxime Gariel, Ashok N Srivastava, and Eric Feron. Trajectory clustering and an application to airspace monitoring. *Intelligent Transportation Systems, IEEE Transactions on*, 12(4):1511–1524, 2011.
- [20] Jiawei Han. *Data Mining: Concepts and Techniques*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2005.
- [21] Abbas Harati-Mokhtari, Alan Wall, Philip Brooks, and Jin Wang. Automatic identification system (ais): data reliability and human error implications. *Journal of navigation*, 60(03):373–389, 2007.
- [22] Yaobin He, Haoyu Tan, Wuman Luo, Shengzhong Feng, and Jianping Fan. Mr-dbscan: a scalable mapreduce-based dbscan algorithm for heavily skewed data. *Frontiers of Computer Science*, 8(1):83–99, 2014.
- [23] Alexander Hinneburg and Daniel A Keim. An efficient approach to clustering in large multimedia databases with noise. In *KDD*, volume 98, pages 58–65, 1998.
- [24] C. A. R. Hoare. Algorithm 63: Partition. *Commun. ACM*, 4(7):321–, July 1961.
- [25] IMO. *Ships’ Routing*. International Maritime Organization, London, 2013.
- [26] Leonard Kaufman and Peter Rousseeuw. *Clustering by means of medoids*. North-Holland, 1987.

- [27] Samira Kazemi, Shahrooz Abghari, Niklas Lavesson, Henric Johnson, and Peter Ryman. Open data for anomaly detection in maritime surveillance. *Expert Systems with Applications*, 40(14):5719–5729, 2013.
- [28] Eamonn J Keogh and Michael J Pazzani. Scaling up dynamic time warping for datamining applications. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 285–289. ACM, 2000.
- [29] Richard O Lane, David A Nevell, Steven D Hayward, and Thomas W Beaney. Maritime anomaly detection and threat assessment. In *Information Fusion (FUSION), 2010 13th Conference on*, pages 1–8. IEEE, 2010.
- [30] Rikard Laxhammar. Anomaly detection for sea surveillance. In *Information Fusion, 2008 11th International Conference on*, pages 1–8. IEEE, 2008.
- [31] Rikard Laxhammar, Göran Falkman, and Egils Sviestins. Anomaly detection in sea traffic—a comparison of the gaussian mixture model and the kernel density estimator. In *Information Fusion, 2009. FUSION'09. 12th International Conference on*, pages 756–763. IEEE, 2009.
- [32] Nicolas Le Guillarme and Xavier Lerouvreur. Unsupervised extraction of knowledge from s-ais data for maritime situational awareness. In *Information Fusion (FUSION), 2013 16th International Conference on*, pages 2025–2032. IEEE, 2013.
- [33] Jae-Gil Lee, Jiawei Han, and Kyu-Young Whang. Trajectory clustering: A partition-and-group framework. In *Proceedings of the 2007 ACM SIGMOD Inter. Conf. on Management of Data, SIGMOD '07*, pages 593–604, New York, NY, USA, 2007. ACM.
- [34] Kingsly Leung and Christopher Leckie. Unsupervised anomaly detection in network intrusion detection using clusters. In *Proceedings of the Twenty-eighth Australasian conference on Computer Science-Volume 38*, pages 333–342. Australian Computer Society, Inc., 2005.
- [35] Jing Li, Kuei-Ying Huang, Jionghua Jin, and Jianjun Shi. A survey on statistical methods for health care fraud detection. *Health care management science*, 11(3):275–287, 2008.
- [36] Bo Liu, Erico N de Souza, Stan Matwin, and Marcin Sydow. Knowledge-based clustering of ship trajectories using density-based approach. In *Big Data (Big Data), 2014 IEEE International Conference on*, pages 603–608. IEEE, 2014.
- [37] Brendan Tran Morris and Mohan M. Trivedi. A survey of vision-based trajectory learning and analysis for surveillance. *IEEE Trans. Circuits Syst. Video Techn.*, 18(8):1114–1127, 2008.

- [38] Gonzalo Navarro. A guided tour to approximate string matching. *ACM computing surveys (CSUR)*, 33(1):31–88, 2001.
- [39] David Nevell. Anomaly detection in white shipping. *Mathematics in Defence*, 2009.
- [40] Giuliana Pallotta, Michele Vespe, and Karna Bryan. Traffic knowledge discovery from ais data. In *Information Fusion (FUSION), 2013 16th International Conference on*, pages 1996–2003. IEEE, 2013.
- [41] Giuliana Pallotta, Michele Vespe, and Karna Bryan. Vessel pattern knowledge discovery from ais data: A framework for anomaly detection and route prediction. *Entropy*, 15(6):2218–2245, 2013.
- [42] Andrey Tietbohl Palma, Vania Bogorny, Bart Kuijpers, and Luis Otavio Alvares. A clustering-based approach for discovering interesting places in trajectories. In *Proceedings of the 2008 ACM Symposium on Applied Computing, SAC '08*, pages 863–868, New York, NY, USA, 2008. ACM.
- [43] Md Mostofa Ali Patwary, Diana Palsetia, Ankit Agrawal, Wei-keng Liao, Fredrik Manne, and Alok Choudhary. A new scalable parallel dbscan algorithm using the disjoint-set data structure. In *High Performance Computing, Networking, Storage and Analysis (SC), 2012 International Conference for*, pages 1–11. IEEE, 2012.
- [44] Claudio Piciarelli, Gian Luca Foresti, and Lauro Snidaro. Trajectory clustering and its applications for video surveillance. In *Advanced Video and Signal Based Surveillance, 2005. AVSS 2005. IEEE Conference on*, pages 40–45. IEEE, 2005.
- [45] Laxhammar Rikard. Anomaly detection in trajectory data for surveillance applications. *Studies from the school of science and technology at rebro university 19*, 2011.
- [46] Branko Ristic, Barbara La Scala, Mark Morelande, and Neil Gordon. Statistical analysis of motion patterns in ais data: Anomaly detection and motion prediction. In *Information Fusion, 2008 11th International Conference on*, pages 1–7. IEEE, 2008.
- [47] Jose Antonio MR Rocha, Gabriel Oliveira, Luis O Alvares, Vania Bogorny, and VC Times. Db-smot: a direction-based spatio-temporal clustering method. In *Intelligent systems (IS), 2010 5th IEEE international conference*, pages 114–119. IEEE, 2010.
- [48] Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.

- [49] Peter J Rousseeuw and Annick M Leroy. *Robust regression and outlier detection*, volume 589. John Wiley & Sons, 2005.
- [50] Jean Roy. Rule-based expert system for maritime anomaly detection. In *SPIE Defense, Security, and Sensing*, pages 76662N–76662N. International Society for Optics and Photonics, 2010.
- [51] Jean Roy and Michael Davenport. Categorization of maritime anomalies for notification and alerting purpose. In *NATO workshop on data fusion and anomaly detection for maritime situational awareness, La Spezia, Italy*, pages 15–17, 2009.
- [52] Lauro Snidaro, Claudio Piciarelli, and Gian Luca Foresti. Fusion of trajectory clusters for situation assessment. In *Information Fusion, 2006 9th International Conference on*, pages 1–7. IEEE, 2006.
- [53] Stefano Spaccapietra, Christine Parent, Maria Luisa Damiani, Jose Antonio de Macedo, Fabio Porto, and Christelle Vangenot. A conceptual view on trajectories. *Data Knowl. Eng.*, 65(1):126–146, April 2008.
- [54] Abhinav Srivastava, Amlan Kundu, Shamik Sural, and Arun K Majumdar. Credit card fraud detection using hidden markov model. *Dependable and Secure Computing, IEEE Transactions on*, 5(1):37–48, 2008.
- [55] Chris Stauffer and W Eric L Grimson. Learning patterns of activity using real-time tracking. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(8):747–757, 2000.
- [56] Pang-Ning Tan, Michael Steinbach, and Vipin Kumar. *Introduction to Data Mining, (First Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2005.
- [57] Henry S Teng, Kaihu Chen, and SC Lu. Adaptive real-time anomaly detection using inductively generated sequential patterns. In *Research in Security and Privacy, 1990. Proceedings., 1990 IEEE Computer Society Symposium on*, pages 278–284. IEEE, 1990.
- [58] Chris Veness. Calculate distance, bearing and more between latitude/longitude points. <http://www.movable-type.co.uk/scripts/latlong.html>.
- [59] Michele Vespe, Ingrid Visentini, Karna Bryan, and Paolo Braca. Unsupervised learning of maritime traffic patterns for anomaly detection. In *Data Fusion & Target Tracking Conf. (DF&TT 2012): Algorithms & Applications, 9th IET*, pages 1–5. IET, 2012.
- [60] Stijn Viaene, Richard A Derrig, Bart Baesens, and Guido Dedene. A comparison of state-of-the-art classification techniques for expert automobile insurance claim fraud detection. *Journal of Risk and Insurance*, 69(3):373–421, 2002.

- [61] Michail Vlachos, George Kollios, and Dimitrios Gunopulos. Discovering similar multidimensional trajectories. In *Data Engineering, 2002. Proceedings. 18th International Conference on*, pages 673–684. IEEE, 2002.
- [62] Wei Wang, Jiong Yang, Richard Muntz, et al. Sting: A statistical information grid approach to spatial data mining. In *VLDB*, volume 97, pages 186–195, 1997.
- [63] Mohammed Javeed Zaki, Srinivasan Parthasarathy, Mitsunori Ogihara, Wei Li, et al. New algorithms for fast discovery of association rules. In *KDD*, volume 97, pages 283–286, 1997.
- [64] Tian Zhang, Raghu Ramakrishnan, and Miron Livny. Birch: an efficient data clustering method for very large databases. In *ACM SIGMOD Record*, volume 25, pages 103–114. ACM, 1996.
- [65] Zhang Zhang, Kaiqi Huang, and Tieniu Tan. Comparison of similarity measures for trajectory clustering in outdoor surveillance scenes. In *Pattern Recognition, 2006. ICPR 2006. 18th International Conference on*, volume 3, pages 1135–1138. IEEE, 2006.

Appendix A

Source Code of the Algorithms

In this appendix, some key source code implemented during the thesis is shown. The whole framework is divided into two parts: normal traffic patterns extraction model and anomaly detection model. The first part shown in A.1 has been mostly implemented in Java with a small portion of R code. While the second part shown in A.2 has been implemented totally in R.

A.1 Source Code of Normal Traffic Patterns Extraction Model

As discussed in Chapter 3, the normal traffic extraction model includes two different clustering components (DBSCANSD and DBSCAN) to handle moving trajectory points and stopping areas respectively. Before presenting the code for clustering, it is necessary to first define the classes of **Trajectory Point** and **Cluster**.

```
1 public class TrajectoryPoint {
2     private String mmsi;           //mmsi, the id of the vessel
3     private long timestamp;       //second of UTC time stamp
4     //the following protected variables are to be inherited by GV
5     protected double longitude;   //longitude
6     protected double latitude;    //latitude
7     protected double SOG;         //speed over ground
8     protected double COG;        //course over ground
9     //the following two are useful during the clustering process
10    private boolean isVisited;     //whether it has been visited
11    private boolean isCorePoint;   //whether it is a core point
12    public TrajectoryPoint() {
13        this.isVisited = false;    //initialize the point unvisited
14    }
15    //generate getter and setter functions ...
16 }
```

Listing A.1: Class of TrajectoryPoint

The **cluster** defined in Definition 9 is basically a set of trajectory points. So the Class of Cluster can be written as follows:

```

1 public class Cluster {
2     private ArrayList<TrajectoryPoint> cluster ;
3     private double avgCOG; //the average direction (COG) of the whole
4         cluster
5     public Cluster() {}
6     public double calculateAverageDirection() {
7         double sum = 0;
8         for(int i=0; i<this.cluster.size();i++) {
9             sum = sum+this.cluster.get(i).getCOG();
10        }
11        double avg = sum/(double)(this.cluster.size());
12        return avg;
13    }
14    //generate getter and setter functions ...
}

```

Listing A.2: Class of Cluster

A.1.1 Code of Clustering Process

The clustering process includes two steps which are DBSCANSD for moving points clustering and DBSCAN for stopping points clustering. The two steps have been implemented together in terms of the code's reusability (DBSCANSD can be regarded as an extension of DBSCAN).

```

1 public class DBScanSD {
2     //final clustering results , a global variable
3     private ArrayList<Cluster> resultClusters = new
4         ArrayList<Cluster>();
5     /**
6     * Apply DBSCANSD on the data set. If the data is stoppoing points ,
7         set isStopPoint as true and it will execute the original
8         DBSCAN directly. If the data is moving points , set isStopPoint
9         as false and it will execute DBSCANSD considering speed and
10        direction .
11    */
}

```

```

7 public ArrayList<Cluster> applyDBScanSD(ArrayList<TrajectoryPoint>
  pointsList, double eps, int minPoints, double maxSpd, double
  maxDir, boolean isStopPoint) {
8   for(int index=0;index<pointsList.size();index++) {
9     //we should mark the point as visited, it has no problem
    because here we use a for loop, it can stop
10    ArrayList<TrajectoryPoint> tmpLst = new
      ArrayList<TrajectoryPoint>();
11    TrajectoryPoint p = pointsList.get(index);
12    if(p.isVisited()&&index!=(pointsList.size()-1)&&index%4096!=0)
      continue;
13    tmpLst = isCorePoint(pointsList, p, eps, minPoints, maxSpd,
      maxDir, isStopPoint);
14    if(tmpLst!=null || index==(pointsList.size()-1) || index%4096==0){
15      Cluster c = new Cluster();
16      c.setCluster(tmpLst);
17      if(tmpLst!=null) resultClusters.add(c);
18      int length=resultClusters.size();
19      boolean flag = true;
20      if((index%4096==0) || (index==(pointsList.size()-1))) {
21        while(flag) {
22          flag = false;
23          for(int i=0;i<length;i++){
24            for(int j=0;j<length;j++){
25              if(i!=j){
26                if(i == length) {
27                  flag = true;
28                  continue;
29                }
30                if(mergeClusters(resultClusters.get(i),
                  resultClusters.get(j))) {
31                  resultClusters.remove(j);
32                  j--;
33                  length--;
34                }
35              }
36            }
37          }
38        }

```

```

39         }
40     }
41 }
42     return resultClusters;
43 }
44 /**
45  * Merge two clusters into one cluster
46  */
47 public boolean mergeClusters(Cluster clusterA , Cluster clusterB) {
48     boolean merge = false;
49     if(clusterA.getCluster() == null || clusterB.getCluster() ==
50         null) {
51         return merge;
52     }
53     for(int index = 0; index < clusterB.getCluster().size(); index++)
54     {
55         TrajectoryPoint p = clusterB.getCluster().get(index);
56         if(p.isCorePoint() && clusterA.getCluster().contains(p)) {
57             merge = true;
58             break;
59         }
60     }
61     if(merge) {
62         for(int index=0; index<clusterB.getCluster().size();index++) {
63             if(!clusterA.getCluster().contains(clusterB.getCluster()
64                 .get(index))) {
65                 clusterA.getCluster().add(clusterB.getCluster().get(index));
66             }
67         }
68     }
69     return merge;
70 }
71 /**
72  * Decide if the point p is core point, if yes, return the list
73  * with p and its neighbors, if no, return null.
74  */

```

```

71 public ArrayList<TrajectoryPoint>
    isCorePoint(ArrayList<TrajectoryPoint> lst, TrajectoryPoint p,
    double eps, int minPoints, double maxSpd, double maxDir, boolean
    isStopPoint) {
72     int count = 0;
73     ArrayList<TrajectoryPoint> tmpList = new
        ArrayList<TrajectoryPoint>();
74     for(Iterator<TrajectoryPoint> it = lst.iterator(); it.hasNext();)
        {
75         TrajectoryPoint q = it.next();
76         if(isDensityReachable(p, q, eps, minPoints, maxSpd, maxDir,
            isStopPoint)) {
77             count++;
78             if(!tmpList.contains(q)) {
79                 tmpList.add(q);
80             }
81         }
82     }
83     if(count>=minPoints) {
84         p.setCorePoint(true);
85         p.setVisited(true);
86         return tmpList;
87     }
88     return null;
89 }
90 /**
91  * Decide if the two points are density reachable.
92  */
93 public boolean isDensityReachable(TrajectoryPoint p1,
    TrajectoryPoint p2, double eps, int minPts, double maxSpd,
    double maxDir, boolean isStopPoint) {
94     boolean result = false;
95     if(gpsDistance(p1.getLatitude(), p1.getLongitude(),
        p2.getLatitude(), p2.getLongitude()) <=eps) {
96         //if they are stopping points, we can directly use original
            DBSCAN algorithm without considering speed or direction
97         if(isStopPoint) return true;
98         if(Math.abs(p1.getCOG()-p2.getCOG())<maxDir) {
99             if(Math.abs(p1.getSOG()-p2.getSOG())<maxSpd) {

```

```

100         result = true;
101     }
102 }
103 }
104     return result;
105 }
106 /**
107  * calculalte the gps distance between two trajectory points
108  * considering the curve of the earth.
109 */
109 public static double gpsDistance(double lat1, double lng1, double
110     lat2, double lng2) {
111     double earthRadius = 3958.75;
112     double dLat = Math.toRadians(lat2-lat1);
113     double dLng = Math.toRadians(lng2-lng1);
114     double a = Math.sin(dLat/2) * Math.sin(dLat/2) +
115     Math.cos(Math.toRadians(lat1)) * Math.cos(Math.toRadians(lat2))
116     *
117     Math.sin(dLng/2) * Math.sin(dLng/2);
118     double c = 2 * Math.atan2(Math.sqrt(a), Math.sqrt(1-a));
119     double dist = earthRadius * c;
120     int meterConversion = 1609;
121     return (double) (dist * meterConversion);
122 }

```

Listing A.3: The code of clustering process

A.1.2 GVs Calculation

In this section, the code of the algorithm to calculate a cluster's Gravity Vectors is given. First, as defined in Chapter 3, a Gravity Vector is a vector formed by 5 features: average COG, average SOG, average Latitude, average Longitude and Median Distance, the class of the **Gravity Vector** needs to be first presented. Here the GV is inherited from the Class of Trajectory Point (Listing A.1).

```

1 | public class GravityVector extends TrajectoryPoint{
2 |     private double medianDistance; //the median distance feature

```

```

3 public GravityVector(double longitude, double latitude, double COG,
4     double SOG, double medianDistance) {
5     this.latitude = latitude;
6     this.longitude = longitude;
7     this.COG = COG;
8     this.SOG = SOG;
9     this.medianDistance = medianDistance;
10 }
11 //generate getter and setter functions ...
12 }

```

Listing A.4: Class of Gravity Vector

As demonstrated in Chapter 3, the process of calculating GVs of a cluster needs to map the trajectory points to an axis of the average direction. So here a Class called Mapping Point is also given and it also inherits from the Class of Trajectory Point:

```

1 public class MappingPoint extends TrajectoryPoint{
2     private double mappingtude; //mapping point's coordinate on the
3     exact axis
4     /**
5     * Mapping Point is the point that is mapped onto the axis (average
6     direction of the whole cluster)
7     */
8     public MappingPoint(double longitude, double latitude, double
9     mappingtude, double COG, double SOG) {
10     this.longitude = longitude;
11     this.latitude = latitude;
12     this.mappingtude = mappingtude;
13     this.COG = COG;
14     this.SOG = SOG;
15 }
16 /**
17 * map the trajectory point to the particular axis and return the
18 mapping point
19 */
20 public MappingPoint convertPointToMappingPoint(TrajectoryPoint p,
21     double avgCOG) {
22     double mappingtude = 0;

```

```

18 | double angle = (avgCOG/(double)180) * Math.PI;
19 | if((avgCOG>=0&&avgCOG<90)) {
20 |     mappingtude = (p.getLongitude() + (1.0/Math.tan(angle)) *
21 |         p.getLatitude()) * Math.sin(angle);
22 | } else if((avgCOG >= 270&&avgCOG<360)) {
23 |     mappingtude = (p.getLatitude() - (Math.tan(Math.PI*2-angle))*
24 |         p.getLongitude()) * Math.cos(Math.PI*2 - angle);
25 | } else if(avgCOG >= 90&&avgCOG<180) {
26 |     mappingtude = ((Math.tan(Math.PI-angle))*p.getLongitude() -
27 |         p.getLatitude()) * Math.cos(Math.PI - angle);
28 | } else if(avgCOG >= 180&&avgCOG<270) {
29 |     mappingtude = -(((double)1/Math.tan(angle - Math.PI)) *
30 |         p.getLatitude() + p.getLongitude()) * Math.sin(angle -
31 |         Math.PI);
32 | }
    MappingPoint mp = new MappingPoint(p.getLongitude(),
    p.getLatitude(), mappingtude, p.getCOG(), p.getSOG());
    return mp;
}
// generate setter and getter functions ...
}

```

Listing A.5: Class of Mapping Point

After having the Gravity Vector and Mapping Point classes, we can start to write the code for extracting the GVs from a given cluster.

```

1 | public class GravityVectorExtraction {
2 |     /**
3 |      * extract the GVs from the moving cluster
4 |      * @param cluster: input a cluster for extracting the GVs
5 |      * @return an arraylist of GVs
6 |      */
7 |     public ArrayList<GravityVector> extractGravityVector(Cluster
8 |         cluster) {
9 |         //Step(1): Calculate the average COG of the whole cluster
10 |         double avgCOG = cluster.calculateAverageDirection();
11 |         //Step(2) Map all the points to the axis and partition the points
            based on 0.01 grid width
            ArrayList<MappingPoint> mpLst = new ArrayList<MappingPoint>();

```

```

12   for(int i=0; i<cluster.getCluster().size(); i++) {
13       MappingPoint mp =
           MappingPoint.convertPointToMappingPoint(cluster.getCluster()
           .get(i), avgCOG);
14       mpLst.add(mp);
15   }
16   insertionSort(mpLst);
17   ArrayList<GravityVector> ppL = new ArrayList<GravityVector>();
18   int count = 0;
19   int k = 0;
20   double sum_x = 0;
21   double sum_y = 0;
22   double sum_SOG = 0;
23   double sum_COG = 0;
24   ArrayList<MappingPoint> traPointsTMP = new
           ArrayList<MappingPoint>();
25   double medianDistance = 0;
26   //partition the points and calculate each GV for each cell
27   while(count<=mpLst.size()) {
28       if(count < mpLst.size() && (mpLst.get(count).getMappingtude()
           - mpLst.get(k).getMappingtude() < 0.01)) { //0.01 as the
           pre-defined grid width
29           sum_x = sum_x+mpLst.get(count).getLongitude();
30           sum_y = sum_y+mpLst.get(count).getLatitude();
31           sum_SOG = sum_SOG+mpLst.get(count).getSOG();
32           sum_COG = sum_COG+mpLst.get(count).getCOG();
33           traPointsTMP.add(mpLst.get(count));
34           count++;
35       } else {
36           double x = 0;
37           double y = 0;
38           double sog = 0;
39           double cog = 0;
40           x = sum_x/(double)(count-k);
41           y = sum_y/(double)(count-k);
42           sog = sum_SOG/(double)(count-k);
43           cog = sum_COG/(double)(count-k);
44           //insert median distance calculation
45           double [] distances = new double[traPointsTMP.size()];

```



```

46         for(int i=0; i<traPointsTMP.size();i++) {
47             double lon=traPointsTMP.get(i).getLongitude();
48             double lat=traPointsTMP.get(i).getLatitude();
49             double dist = gpsDistance(lat , lon , y , x);
50             distances [i]=dist;
51         }
52         //medianDistance
53         medianDistance = quartile(distances , 50);
54         //for each cell of the grid, calculate its GV
55         GravityVector gv = new
56             GravityVector(x,y,cog ,sog ,medianDistance);
57         ppL.add(gv);
58         sum_x = 0;
59         sum_y = 0;
60         sum_COG = 0;
61         sum_SOG = 0;
62         k = count;
63         traPointsTMP.clear();
64         if(count==mpLst.size()) break;
65     }
66     return ppL;
67 }
68 /**
69  * insertion sorting
70  */
71 public void insertionSort(ArrayList<MappingPoint> mpl) {
72     for(int i=1; i<mpl.size(); i++) {
73         int k = i;
74         MappingPoint mp = mpl.get(i);
75         boolean insertAlready = false;
76         while(mpl.get(i).getMappingtude()<mpl.get(k-1).getMappingtude())
77             {
78                 if(k==1) {
79                     mpl.remove(i);
80                     mpl.add(0, mp);
81                     insertAlready = true;
82                     break;
83             }

```

```

83         k--;
84     }
85     if (!insertAlready) {
86         mpl.remove(i);
87         mpl.add(k, mp);
88     }
89 }
90 }
91 /**
92  * Retrieve the quartile value from an array, used for generating
93  * relative distance.
94  * @param values The array of data
95  * @param lowerPercent The percent cut off. For the lower quartile
96  * use 25, for the upper-quartile use 75
97  * @return the quartile value
98  */
99 public double quartile(double[] values, double lowerPercent) {
100     if (values == null || values.length == 0) {
101         throw new IllegalArgumentException("The data array either
102             is null or does not contain any data.");
103     }
104     // order the values
105     double[] v = new double[values.length];
106     Arrays.sort(v);
107     int n=0;
108     if(v.length==1) {
109         n = 0;
110     }
111     else n = (int) Math.round(v.length * lowerPercent / 100);
112     return v[n];
113 }

```

Listing A.6: Extract GVs from a cluster

A.1.3 SSPs Calculation

SSP (Sampled Stopping Point) is a type of points to represent the geo-spatial shape of a cluster and it does not need to consider the factors of speed or direction (See Section

3.3). Here the algorithm of extracting SSPs from stopping clusters (Algorithm 2) has been implemented in R.

```

1  ## The main function of SSP extraction process.
2  ## The parameter of stoppingClusters is the stopping clustering results
   generated by DBSCAN;
3  ## The parameter of radius here is the eps defined during the
   clustering process
4  stopPointsSampling<-function(stoppingClusters , radius) {
5    clusterIDs<-unique(stoppingClusters [,"clusterindex"]);
6    end=length(clusterIDs);
7    result<- data.frame(Longitude=numeric(),Latitude=numeric());
8    for(i in 1:end) {
9      id=clusterIDs[i];
10     clusterData<-stoppingClusters [stoppingClusters$clusterindex==id,];
11     sample_num<-sampleNumberEstimate(clusterData , radius);
12     sampledPoints<-samplePointsFromData(clusterData , sample_num, radius);
13     result<-rbind(result , sampledPoints);
14   }
15   return(result);
16 }
17 ## Sample points from stopping clusters such that every sampled point
   is far enough from other sampled points
18 samplePointsFromData<-function(data , sample_num, radius) {
19   result<- sampledPoint<-data [sample(nrow(data) ,1) ,2:3];
20   count=1;
21   while(count<sample_num) {
22     sampledPoint<-data [sample(nrow(data) ,1) ,2:3];
23     end=nrow(result);
24     nearFlag=FALSE;
25     for(i in 1:end) {
26       if(gpsDistance(result [i,] , sampledPoint)<radius) {
27         nearFlag=TRUE;
28         break;
29       }
30     }
31     if(nearFlag) {
32       next;
33     } else {

```

```

34     count=count+1;
35     result<-rbind(result ,sampledPoint);
36   }
37 }
38 return(result);
39 }
40 ## Estimate the number of points to be sampled
41 sampleNumberEstimate<-function(clusterPoints , radius) {
42   min_x<-min(clusterPoints [, "Latitude"]); #downmost point
43   max_x<-max(clusterPoints [, "Latitude"]); #upmost point
44   min_y<-min(clusterPoints [, "Longitude"]); #rightmost point
45   max_y<-max(clusterPoints [, "Longitude"]); #leftmost point
46   areaEstimate<-(max_x-min_x)*(max_y-min_y);
47   numberEstimate<-floor (areaEstimate/(pi*radius ^2))+1;
48   if (areaEstimate==0) {
49     numberEstimate=max((max_x-min_x) ,(max_y-min_y))/radius+1;
50   }
51   return (numberEstimate);
52 }

```

Listing A.7: Extract SSPs from stopping clusters

A.2 Source Code of Anomaly Detection Model

This section is focused on the anomaly detection labeling process. The algorithm uses as input the clusters described in previous sections, and uses them to generate a label that informs what type of anomaly was found. The algorithm associates labels to each AIS data point according to the following:

- **0** stands for distance abnormal
- **1** is normal
- **-1** is direction or speed abnormal
- **9** is data incomplete

The R function code, called `labelAnomalyPointsWithStopMoveclusters`, contains the function that is responsible for the labeling detection. Observe that this function uses other helper functions that are later described. The function inputs are the following:

- The `inputData` parameter is a data frame of the following format: (“MMSI”, “SOG”, “Longitude”, “Latitude”, “COG”);
- The `normalPointsMove` parameter is a data frame of the following format: (“clusterindex”, “Longitude”, “Latitude”, “SOG”, “COG”, “QuartileDistance”)
- The `normalPointsStop` is a data frame with the following format: (“Longitude”, “Latitude”)

The algorithm first executes a check to guarantee that data is in correct format, by verifying if all required fields in `inputData` are complete. If they are not, the algorithm adds a label 9, indicating that data is incomplete. This step guarantees that all AIS coordinates that have missing values are going to be labeled as incomplete.

Next, the algorithm executes a function called `generatefeaturesRelativeDistance`, which generates a set of extra distance descriptors for the `inputData` parameter. This function is presented in next section. These distances will be used to rank which AIS points are more abnormal than others. The function `getIndexOffarPointsRELATIVE` calculates which points are abnormal in relation to distance from the clusters. This function uses two extra parameters as threshold to indicate the anomaly. These values were empirically found during clustering steps.

If the points are still normal in relation to distance from the clusters, it does not mean that they are normal in relation to speed and direction. Then a second set of tests is required to check abnormality. This is started in function `getIndexOfNearPointsRELATIVE`, which simple gets a list of points that are close to the clusters. After getting the list of close points, the algorithm generates another set of distances that are used to estimate direction and speed anomalies. This is done in function `generatefeatures4` that calculates the cosine distance, and speed similarities in relation to the cluster. Next section presents the other sub-routines.

```

1 ## label the input data with normal moving clusters(GVs) and normal
   stopping clusters(SSPs)
2 labelAnomalyPointsWithStopMoveclusters<-function(inputData ,
   normalPointsMove , normalPointsStop) {
3   labeledData<-inputData;
4   labeledData[, "label"]<-labeledData[,1];
5   labeledData[, "label"]<-1;
6   #calulate features for the 1st step
7   #extract those incomplete data

```

```

8   end=nrow(labeledData);
9   for(i in 1:end) {
10      if(!complete.cases(labeledData[i,])) {
11         labeledData[i,"label"]<-9;
12      }
13   }
14   print("Incomplete records labeled finished")
15   featuresOne <- generatefeaturesRelativeDistance(inputData ,
16      normalpointsMove = normalPointsMove , normalpointsStop =
17      normalPointsStop);
16   #below thresholds are for Juan De Fuca strait in paper of
17      "Knowledge-based clustering of ship trajectories using
18      density-based approach"
17   rel_dis_threshold = 5.765;
18   abs_dis_threshold = 94.096
19   #get indices for far points
20   index_far_relative <- getIndexOfFarPointsRELATIVE(featuresOne ,
21      rel_dis_threshold , abs_dis_threshold = abs_dis_threshold);
21   if(nrow(index_far_relative)>0) {
22      for(m in 1:nrow(index_far_relative)) {
23         if(labeledData[index_far_relative[m,1] ,"label"] !=9) {
24            labeledData[index_far_relative[m,1] ,"label"] <-0;
25         }
26      }
27   }
28   print("First step finished.")
29   #2nd step
30   index_near_relative <- getIndexOfNearPointsRELATIVE(featuresOne ,
31      rel_dis_threshold , abs_dis_threshold);
31   if(nrow(index_near_relative)>0) {
32      featureSecond<-generatefeatures4(inputData [index_near_relative
33         [,1],], normalPoints , rel_dis_threshold);
33      cosineThreshold=0.588
34      end2 = nrow(featureSecond);
35      for(j in 1:end2) {
36         if(featureSecond[j,]$cosSimilarity<cosineThreshold) {
37            labeledData[index_near_relative[featureSecond[j,]$INDEX,1] ,
38               "label"] <- (-1);

```

```

39     }
40   }
41   print("Data label process finished!")
42   return(labeledData);
43 }
44 ## generate features considering relative distance
45 ## realpoints: the point to be labeled
46 ## normalpointsMove: moving clusters (GVs)
47 ## normalpointsStop: stopping clusters (SSPs)
48 generatefeaturesRelativeDistance <- function(realpoints ,
      normalpointsMove , normalpointsStop) {
49   end = nrow(realpoints);
50   features = data.frame (distance=numeric() ,
      relative_distance=numeric() , SOGRatio=numeric() , COG=numeric() ,
      isStopPoint=integer());
51   for(i in 1:end) {
52     features =
53       rbind(features , calculateDistancesRelative(realpoints[i ,] ,
      normalpointsMove , normalpointsStop));
54   }
55   colnames(features) <- c("Absolute_Distance" , "Relative_Distance" ,
      "SOGRatio" , "COG" , "isStopPoint");
56   return(features);
57 }
58 ## calculate the relative distances between the point and the clusters
59 ## realpoints: the point to be labeled
60 ## normalpointsMove: moving clusters (GVs)
61 ## normalpointsStop: stopping clusters (SSPs)
62 calculateDistancesRelative<- function(realpoint , normalpoints ,
      normalpointsStop) {
63   index = -1;
64   min_d = 10000000000000000;
65   relative_distance = 10000000000000000;
66   end = nrow(normalpoints);
67   end2=nrow(normalpointsStop);
68   isStopPoint = 0;
69   #use 0.5 as stopping point condition , we can change it to another
      speed threshold
70   if(realpoint[1 , "SOG"] <= 0.5 | !complete.cases(realpoint[1 ,])) {

```

```

70     isStopPoint = 1;
71     for( i in 1:end2) {
72         #actually here I should use absolute distance instead of
           relative_distance
73         #I am too lazy to do this and this will not influence the
           result , just a name
74         distance_tmp = gpsdist(realpoint[1,"Latitude"],
           realpoint[1,"Longitude"],
           normalpointsStop[i,"Latitude"],
           normalpointsStop[i,"Longitude"]);
75         if(distance_tmp < min_d) {
76             relative_distance<-999999999;#stands for NA
77             min_d <-distance_tmp;
78             index <- i;
79         }
80     }
81 } else {
82     for(i in 1:end) {
83         distance = gpsdist(realpoint[1, "Latitude"], realpoint[1,
           "Longitude"], normalpoints[i, "Latitude"],
           normalpoints[i, "Longitude"]);
84         quartileDistance=normalpoints[i, "QuartileDistance"];
85         if(quartileDistance==0) {
86             quartileDistance=1; #to avoid the x/0 error
87         }
88         relative_dis_tmp = distance/quartileDistance;
89         if(relative_dis_tmp < relative_distance) {
90             relative_distance<-relative_dis_tmp
91             min_d <-distance;
92             index <- i;
93         }
94     }
95 }
96 speedratio <- (abs(realpoint[1,"SOG"] -
           normalpoints[index,"SOG"]))/ normalpoints[index,"SOG"];
97 direction_dif<- abs(realpoint[1,"COG"] - normalpoints[index,"COG"]);
98
99 return (c(min_d, relative_distance ,speedratio , direction_dif ,
           isStopPoint ));

```



```

100 }
101 ## used for calculating the geographical distance of two positions
102 ## lat1: latitude of point 1
103 ## lon1: longitude of point 1
104 ## lat2: latitude of point 2
105 ## lon2: longitude of point 2
106 gpsdist<-function(lat1,lon1,lat2,lon2) {
107     R = 6367000;
108     dlat = (lat2-lat1)*pi/180;
109     dlon = (lon2-lon1)*pi/180;
110     radlat1 = lat1*pi/180;
111     radlat2 = lat2*pi/180;
112     a = sin(dlat/2) * sin(dlat/2) + (sin(dlon/2) * sin(dlon/2) *
113         cos(radlat1) * cos(radlat2));
114     c = 2 * atan2(sqrt(a), sqrt(1-a));
115     d = R * c;
116 }
117 ## get the indices of those points which are too far away
118 ## test_set: the data set to be labeled
119 ## dis is the threshold for the distance
120 getIndexOfFarPoints <- function(test_set, dis) {
121     end = nrow(test_set)
122     index = data.frame(index=numeric());
123     for(i in 1:end) {
124         if(test_set[i,]$distance>dis) {
125             index <- rbind(index,i);
126         }
127     }
128     return(index);
129 }
130 ## get the indices of the points that are near the clusters
131 ## they can be used for training or predicting the speed&direction
132     abnormal thing
133 getIndexOfNearPoints <- function(test_set, dis) {
134     end = nrow(test_set)
135     index = data.frame(index=numeric());
136     for(i in 1:end) {
137         if(test_set[i,]$distance<=dis) {

```

```

137         index <- rbind(index,i);
138     }
139 }
140 return(index);
141 }
142 ## get the indices of the points that are near the clusters relatively
143 ## they can be used for training or predicting the speed&direction
144     abnormal thing
145 #this is based on relative distance
146 getIndexOfNearPointsRELATIVE <- function(test_set ,
147     relative_dis_threshold , abs_dis_threshold) {
148     end = nrow(test_set)
149     index = data.frame(index=numeric());
150     for(i in 1:end) {
151         if(test_set[i,"isStopPoint"]==1) {
152             } else {
153                 if(test_set[i,]$Relative_Distance <=
154                     relative_dis_threshold) {
155                     index <- rbind(index,i);
156                 }
157             }
158         }
159     return(index);
160 }
161 ## points far away from normal points
162 ## this is based on relative distance
163 ## relative_dis_threshold: the threshold for moving points
164 ## abs_distance: the threshold for stopping points
165 getIndexOfFarPointsRELATIVE <- function(test_set ,
166     relative_dis_threshold , abs_dis_threshold) {
167     end = nrow(test_set)
168     index = data.frame(index=numeric());
169     for(i in 1:end) {
170         if(test_set[i,"isStopPoint"]==1) {
171             if(test_set[i,]$Absolute_Distance > abs_dis_threshold) {
172                 index<-rbind(index,i);
173             }
174         } else {

```

```

171         if(test_set[i,]$Relative_Distance > relative_dis_threshold)
172             {
173                 index <- rbind(index,i);
174             }
175     }
176     return(index);
177 }
178 ## generate the features for the second labelling step: cosine division
179     distance
180 generatefeatures4 <- function(realpoints, normalpoints,
181     threshold_relative) {
182     end = nrow(realpoints);
183     features = data.frame(relativeDistance=numeric(),
184         distanceSD=numeric(),
185         distance=numeric(), cosSimilarity=numeric(),
186         distanceSimilarity=numeric(), SOG=numeric(), COG=numeric(),
187         INDEX=numeric());
188     print("Start to label cosine division distances")
189     for(i in 1:end) {
190         index = i;
191         if(complete.cases(realpoints[i,])) {
192             tmp<-calculateDistances4(realpoints[i,], index,
193                 normalpoints, threshold_relative);
194             features <- rbind(features,tmp);
195         }
196     }
197     colnames(features) <- c("relativeDistance", "distanceSD",
198         "distance", "cosSimilarity", "distanceSimilarity", "SOGdif",
199         "COGdif", "INDEX");
200     #distanceSD minimum distance considering speed and direction
201     return(features);
202 }
203 ## subroutine in function of "generatefeatures4"
204 calculateDistances4<- function(realpoint, index, normalpoints,
205     threshold) {
206     min_relativeDistance = 1000000000; # used for storing the relative
207         distance considering speed&direction

```

```

197 min_d = 100000000000;      #used for storing the distance to the
    nearest points considering speed&direction
198 cosSimilarity = -1000000;  # cosineSimilarity
199 direction_dif<-100000000; # direction difference
200 speed <- 100000000;       # speed difference
201 min_distance <- 100000000000;#used for storing the nearest points
    without considering speed&direction
202 distanceSimilarity <- 0;
203 end = nrow(normalpoints);
204 for(i in 1:end) {
205     distance = gpsdist(realpoint[1, "Latitude"], realpoint[1,
        "Longitude"], normalpoints[i, "Latitude"], normalpoints[i,
        "Longitude"]);
206     relative_distance = distance/(normalpoints[i,
        "QuartileDistance"] + 0.00000001);
207     # below is to record the nearest point in the space without
        considering speed & direction
208     if(distance<min_distance) {
209         min_distance = distance;
210     }
211     if(relative_distance < threshold) {
212         min_relativeDistance=relative_distance;
213         # consider those points with small speed in the stop area
        as normal, that is, no need to consider direction
        parameter(set its difference to 0)
214         if(realpoint[1,"SOG"]<=0.3& normalpoints[i,"SOG"]<=0.3) {
215             if(distance<min_d) {
216                 min_d <- distance;
217             }
218             speed <- abs(realpoint[1,"SOG"]-normalpoints[i,"SOG"]);
219             direction_dif <- 0;
220
221             cosSimilarity = 1;
222         }
223         else {
224             if(realpoint[1,"SOG"]>normalpoints[i,"SOG"]) {
225                 speedRatio <-
                    normalpoints[i,"SOG"]/realpoint[1,"SOG"];

```

```

226         cosSimilarityTmp <- cos(((abs(realpoint[1,"COG"] -
227             normalpoints[i,"COG"]))/180)*pi) * speedRatio;
228     } else {
229         speedRatio <-
230             realpoint[1,"SOG"]/normalpoints[i,"SOG"];
231         cosSimilarityTmp <- cos(((abs(realpoint[1,"COG"] -
232             normalpoints[i,"COG"]))/180)*pi) * speedRatio;
233     }
234     if(cosSimilarityTmp>cosSimilarity) {
235         min_d <-distance;
236         direction_dif<-
237             abs(realpoint[1,"COG"]-normalpoints[i,"COG"]);
238         speed<-abs(normalpoints[i,"SOG"]-realpoint[1,"SOG"]);
239         cosSimilarity<- cosSimilarityTmp;
240     }
241     }
242     distanceSimilarity <- (min_distance+1)/(min_d+1)
243 }
244 return (c(min_relativeDistance, min_d, min_distance, cosSimilarity,
245     distanceSimilarity, speed, direction_dif, index));
246 }

```

Listing A.8: Anomaly Detection