Anomaly Detection in Maritime Data Based on Geometrical Analysis of Trajectories

Behrouz Haji Soleimani*, Erico N. De Souza*, Casey Hilliard[†] and Stan Matwin*[‡]
*Faculty of Computer Science, Dalhousie University, Halifax, NS, Canada.
[†]Dalhousie University, Halifax, NS, Canada.
[‡]Institute for Computer Science, Polish Academy of Sciences, Warsaw, Poland.

Email: behrouz.hajisoleimani@dal.ca, erico.souza@dal.ca, casey.hilliard@dal.ca, stan@cs.dal.ca

Abstract-Anomaly detection is an important use of the Automatic Identification Systems (AIS), because it offers support to users to evaluate if a vessel is in trouble or causing trouble. For instance, it can be used to detect if a ship is doing something that may cause an accident or if it has changed its route to avoid bad weather condition. In this work, a new method for finding anomalies in the ships' movements is proposed. The method analyzes the trajectory of ships from a geometrical perspective. The trajectory of the ship is compared with a near-optimal path that is generated by a graph search algorithm. The proposed method extracts some scale-invariant features from the real trajectory and also from the optimal movement pattern, and it compares the two sets of features to generate an abnormality score. The method is unsupervised and it does not require training. Instead of labeling the trajectories as normal/abnormal it calculates a score value that denotes the extent of abnormality. The scoring scheme provides a ranking system in which the user can sort the trajectories based on their abnormality score. This is useful when dealing with large number of trajectories and the user wants to picks the most abnormal cases. For the evaluation, the method was run on three months data of North Pacific Ocean and score values were generated. Among the entire dataset, 100 randomly chosen trajectories were labeled by an expert. After applying a threshold on the score value, the proposed method had 94% accuracy.

I. INTRODUCTION

Anomaly detection in oceans is a priority for governmental organizations. The transit of goods occurs over the oceans that cover 2/3's of the planet and yet are inhabited by human beings. To help governments with this task, since 2004, the International Maritime Organization (IMO) requires Automatic Identification System (AIS) transponders to be aboard vessels [1] that are above 300 gross tons. Maritime authorities use AIS to track and monitor ships. Originally, AIS transponders used only VHF radio frequencies to broadcast data about the ship, but with the increasing satellite coverage, these transponders are also transmitting using satellites. One problem with AIS protocol is that it generates a large volume of data, which may overwhelm authorities evaluating the vessel trajectory. Another problem presented by [14] is that current systems are "principally navigational aids that do not provide an adequate solution to new threats and the ingenuity of criminal organisations". Under these conditions, it is necessary to develop automatic tools to improve security

awareness. In this context, this work presents a tool to evaluate if the trajectory of cargo ships or tankers are normal or abnormal in the middle of the ocean. The tool is based on a ranking scheme that reduces the computational requirements, and facilitates visualization of data on-demand for users.

The AIS protocol is based on VHF radio messages that transmits ship's identity, position and heading[3]. Each packet contains messages 256 bits long and fits in one of the 4500 time slots per minute [3]. It also uses Time Division Multiple Access (TDMA) to avoid message collisions. The AIS device installed in the ships is connected to the control bridge and a satellite navigation system (e.g. GPS) [3]. There are 23 AIS messages types, and they can be divided between static and dynamic. Static messages contain information such as vessels' IMO number, Maritime Mobile Service Identity (MMSI), name and length. Usually, static messages are sent with the rate of every six minutes [3], [5]. Dynamic messages contain the position, time (UTC second indicating when the report was generated [5]), speed over ground (SOG), etc. Dynamic messages are sent each 2 seconds to 3 minutes depending on the vessel's speed. There are other messages related to security and safety, also specified in the protocol. This creates a huge volume of data that makes very difficult for an operator to promptly detect these anomalies.

Another related problem deals with the data veracity, and it is well known that vessels may tamper with the AIS devices to inform false types of movements. This false information may, for example, serve to hide fishing activity in protected areas. The potential for automatic detection of these anomalies motivates the use of Machine Learning (ML) techniques.

Instead of using a supervised approach to learn the abnormality ranking, we present an unsupervised ranking system based on how the vessel's trajectory approximates an optimal trajectory calculated from a graph search algorithm. The idea is to build an optimal path, then calculate some geometrical properties of these trajectories (the optimal and real ones). These properties are combined to generate a score that will be used by the user to separate the anomalous tracks from normal ones.

Another key aspect of these features relies on their independence in relation to time component. The user must

order the sequence of points to guarantee that the trajectory keeps the same sequence of points. This also simplifies the database query for the points, since it allows the evaluation of any part of the trajectory as normal or not. This approach also exploits the fact that most vessels in the middle of the ocean are interested only in reducing their costs including fuel consumption and time cost, since ship captains have deadlines to deliver their goods in various ports. This makes them tend to choose a minimal path between two ports.

This work is divided into following sections: Section II presents the related work, then Section III shows our proposed framework. Section IV shows the experimental results, and this work finishes with some conclusions and future work on Section V.

II. RELATED WORK

Various works have presented solutions based on ML algorithms to detect anomalous behaviour in ship trajectories. They can be categorized as clustering methods or classification methods. One of the first methods is based on a clustering method to allow data enrichment and pattern recognition proposed by Pallota et al [12], which presented the Traffic Route Extraction and Anomaly Detection (TREAD) methodology that is an unsupervised and incremental learning approach to the extraction of maritime movement patterns. Their approach also enriches the original tracks with a description of the ship movements (e.g. it adds a label informing the type of movement executed by the ship like 'sailing', or 'stationary'). These labels are added based on the incremental clustering algorithm DBSCAN [4], and are used to group the main routes used by ships. Based on the previous step, their methodology tries to predict the main route that a ship will use.

Another work based on clustering is TRAjectory CLUStering (TRACLUS) [10], which is a partition-and-group framework. The algorithm looks at the whole trajectory, then separates the the trajectory into line segments. In a second stage, it groups these segments looking for a cluster connecting them. Both methods, based on clustering, suffer from a common problem: the dependency on the algorithm parameters. This issue was also observed by [6]. If the user is not careful with the input parameters, the algorithms will generate false clusters.

The proposed method does not suffer from these problems, because it does not rely on input parameters from the user. Instead, the algorithm builds a graph and calculates the optimal paths that ships must use to reduce costs related to fuel. The next step is simply to measure how far away are the real trajectory and the optimal ones.

Other ML methods based on classification can be found in literature, [7] uses self-organizing maps to cluster tracks and then uses Gaussian mixture models for decision making. Another work that uses Gaussian mixture is [8], which proposes a clustering model and finds the parameters of the model by Expectation-Maximization method. [11] clusters movement fragments to extract motifs. These motifs are then used as higher level features in a classification method. [2] points out some limitations of trajectory clustering and proposed a method that tries to model normal trajectories with splines. [13] uses a kernel density estimation method to estimate the density of normal points. They use Parzen windows with adaptive window width using a Gaussian kernel as a density estimator. However, density estimation methods are computationally intensive and may not be practical in real-time systems. [9] compares the performance between Gaussian mixture models [8] and kernel density estimation methods [13].

All methods based on classification suffer from the problem of labeled data. Classification requires the data points to be associated with the labels indicating the patterns the user wishes to detect. Methods based on kernel density (like [13]) are applicable only on specific regions, and they cannot deal with larger geographical areas. This paper's method solves both problems, since it does not require labeled data, and it is computationally robust enough to be used in larger areas.

III. ANOMALY DETECTION FRAMEWORK

The proposed anomaly detection framework analyses the trajectory of each ship individually. The trajectory of a ship can be built by connecting its consecutive GPS coordinates in a time frame. This time frame can be equal to the definition of a trip which begins when the ship departs from a port and ends when it arrives at its desired destination. As in most cases with raw AIS data, the information about the departure and destination points of the ships are not available. To address this limitation, we have elected to use a fixed window and analyze the trajectory in that window. We have used one day as the time window for our analysis. Therefore, each trajectory is divided into multiple tracks based on the one day time window regardless of the departure point and destination point of the ship. In fact, we analyze the behavior of each ship in different days separately.

Our anomaly detection framework first tries to find the most normal path, which is the shortest possible path between the start point and the end point of the trajectory. After finding the shortest possible path, it is compared to the actual trajectory that the ship has taken. To compare the real trajectory with its corresponding optimal path, we introduce scale-invariant geometrical features that capture the shape and maneuvering behavior of the trajectories. These features are good descriptors in geometrical terms for analyzing and comparing the shape of the trajectories. Eventually, after extracting these features an anomaly score is calculated and assigned to that trajectory.

A. Finding the shortest possible path

As described before, we first try to find the shortest possible path between the start and end point of the trajectory. A* algorithm is used to find the shortest path, and it is a well known algorithm used in graph theory. A* is guaranteed to find the shortest path if it exists. It requires a connectivity



(b) Density heat-map

Fig. 1: One week data of north pacific ocean.

graph that defines the map. A grid partitioning with a fixed resolution is used to discretize the map of the region of interest and it is used to build the graph. In fact, each cell/node in the grid/graph represents a geographical area and the size of this area depends on the resolution that is used for partitioning.

After constructing the grid, we calculate the number of ships that have passed through each cell. If there exists at least one ship that passed through the cell, we assign it to 1 as passable region and if there is no ship passing the region we assign it to 0 as impassable or obstacle. It should be mentioned that only a small fraction of the dataset is sufficient for constructing the graph and detecting paths and obstacles. In the next step, we run the A* algorithm for each single ship's trajectory, considering its start and end point, to find the shortest possible path for the ship.

Figure 1a represents a binary connectivity graph generated from one week of data in the North Pacific Ocean. Figure 1b illustrates a heat-map generated from the same dataset. The heat-map represents the density of presence of the ships passing from different grids. These kind of heat-maps can be used for detecting the most common paths taken by the vessels. They can also be used as a weighted graph in A* algorithm to find the shortest common path. In this work we used the binary version in the A* algorithm.

B. Feature extraction

As described before, we compare the real trajectory of the ship to the optimal path which is extracted by A*. For this purpose, we introduce four features that capture the geometrical information regarding the shape of the trajectories. These features are extremely useful for comparing the optimal trajectory with the real trajectory in finding anomalies.

1) Length of the trajectory: This is the most intuitive feature that measures the length of the trajectory. This is done by traversing the trajectory and adding up the distance between each sequential pair of GPS coordinates. Since latitude and longitude observations are in a spherical coordinate system, we cannot use Euclidean geometry and Euclidean distance for calculating the distance between two points. For this purpose, the Haversine formula is used to calculate the spherical distance (i.e. great-circle distance) between consecutive points based on their latitude and longitude coordinates. The distance calculations are based on equation (1):

$$h(p_1, p_2) = \sin^2\left(\frac{\Delta\theta}{2}\right) + \cos\theta_1 \cos\theta_2 \sin^2\left(\frac{\Delta\phi}{2}\right) \quad (1a)$$

$$d(p_1, p_2) = 2R \arctan(\frac{\sqrt{h(p_1, p_2)}}{\sqrt{1 - h(p_1, p_2)}})$$
(1b)

where R = 6378137 is the radius of the earth in meters, $p_1 : (\theta_1, \phi_1)$ and $p_2 : (\theta_2, \phi_2)$ are the (latitude, longitude) coordinates of the two observations, $\Delta \theta = \theta_2 - \theta_1$ and $\Delta \phi = \phi_2 - \phi_1$. The length of the trajectory is calculated based on equation (2):

$$L = \sum_{i=1}^{N-1} d(p_i, p_{i+1})$$
(2)

where N is the number of points in the trajectory.

2) Area under the curve: As mentioned before, we have focused on cargo and tanker ships which have more predictable behavior as they tend to use the shortest path and keep their transit costs as low as possible. If we draw a straight line connecting the start and end point of the trajectory and measure the area under the curve of the trajectory it will give us an indication that how much the ship has digressed from that straight line. The more digression from the straight line will result in a bigger area under their curve.

If we treat the trajectory as a function $f: \phi \rightarrow \theta$, we can integrate the function to calculate the area under the curve of the trajectory. To make this happen, a translation and rotation transformation are needed to change the coordinate system. At first, the origin is transferred to the first point by translating all the points in the trajectory as equation (3).

$$p'_{i} = p_{i} - p_{1}$$
 $i = 1, 2, ..., N$ (3)

Next, the angle between the horizon and the straight line connecting the start and end point is calculated. The rotation is applied to the data afterwards based on the calculated angle ψ in order to make the straight line parallel to the equator as in equation (4):

$$\psi = \arctan(\frac{\theta_N - \theta_1}{\phi_N - \phi_1}) \tag{4a}$$

$$R_{\psi} = \begin{bmatrix} \cos\psi & -\sin\psi\\ \sin\psi & \cos\psi \end{bmatrix}$$
(4b)

$$p_i'' = R_{\psi} p_i'$$
 $i = 1, 2, ..., N$ (4c)

where ψ is the angle between the horizon and the straight line connecting start and end point, R_{ψ} is the rotation matrix and $p''_i : (\phi''_i, \theta''_i) \ i = 1, 2, ..., N$ are the points in the new coordinate system. If we treat the new points p''_i in the new coordinate system as the (input, output) of a function f : $\phi'' \to \theta''$, by integrating the function we can calculate the area under the curve of digression from straight line as in equation (5).

$$A = \int_{\phi_1''}^{\phi_N''} f(\phi^{''}) \mathrm{d}\phi^{''}$$
(5)

Since the trajectory is represented a discrete set of observations, the equation (5) can be numerically calculated by the sum of trapezoidal rule as equation (6).

$$A = \sum_{i=1}^{N-1} |\phi_{i+1} - \phi_i| \frac{|\theta_i + \theta_{i+1}|}{2}$$
(6)

As the trajectories are not ordered based on the longitude and they may go back and forth, we take the absolute value of $\Delta \phi$. Also, the area under the curve below the baseline should not counteract the area above the baseline as they are both digressions from the straight line. Thus, we used the absolute value of $\theta_i + \theta_{i+1}$.

3) Gradient of trajectory with respect to latitude and longitude: The other two features are based on the partial derivatives of the trajectory with respect to latitude and longitude dimension. These features try to measure the total amount of variations in each direction separately. Here the derivations are time independent and they calculate the variations in location. G_{θ} and G_{ϕ} are the gradient features with respect to latitude and longitude respectively ad they are calculated based on equation (7).

$$G_{\theta} = \sum_{i=1}^{N-1} |\theta_{i+1} - \theta_i|$$
(7a)

$$G_{\phi} = \sum_{i=1}^{N-1} |\phi_{i+1} - \phi_i|$$
(7b)

C. Scale-independent measure of abnormality

After extracting the aforementioned features from the real trajectory and the optimal trajectory which is found using the A^* algorithm, they are used as a measure of abnormality in the form of equation (8):

$$f' = \frac{f^{op} - f^{re}}{L^{op}} \qquad f \in \{L, A, G_{\theta}, G_{\phi}\}$$
 (8)

where L^{op} is the length of the optimal trajectory, f^{re} and f^{op} are the extracted features from the real trajectory and its corresponding optimal trajectory, respectively. We normalize the differences between the features of optimal and real trajectories by the length of the optimal trajectory in order to make scale-independent features f'. The f' features express the ratio of digression from the optimal trajectory. These features are scale-independent that means if the trajectories are scaled-up or scaled-down, the value of the features would not change as they calculate the amount of digression in unit length. Negative values of the f^{\prime} features mean that the real trajectory that the ship has taken is longer or it has more variation than the optimal trajectory. Zero value means that they are equal in terms of variations, and positive values means that the real trajectory is even better than the optimal one that we found using A*. For example, if the training data used for building the A* graph is not enough, a ship may take a route which is shorter than any similar route that we used to build the graph. In these cases, we expect a positive value for the features.

D. Calculating the final score

Each of the f' features is a measure of difference from the optimal path. These features can be added together to build up the final anomaly score if they are properly normalized. The length feature L is calculated based on meters, but the other three features, A, G_{θ} and G_{ϕ} are calculated based on degrees. In order to normalize them in the same scale, we multiply the A, G_{θ} and G_{ϕ} features by a constant that is the number of meters in one degree of latitude. Therefore, the final score is calculated based on equation (9):

$$Score = \frac{L^{op} - L^{re}}{L^{op}} + \eta (\frac{A^{op} - A^{re}}{L^{op}} + \frac{G^{op}_{\theta} - G^{re}_{\theta}}{L^{op}} + \frac{G^{op}_{\phi} - G^{re}_{\phi}}{L^{op}})$$
(9)

where $\eta = 111319$ is the number of meters in one degree of latitude, L^{op} , L^{re} , A^{op} , A^{re} , G_{θ}^{op} , G_{θ}^{re} , G_{ϕ}^{op} and G_{ϕ}^{re} are the length of the optimal trajectory, length of the real trajectory, area under the optimal trajectory, area under the real trajectory, latitude variations of optimal trajectory, latitude variations of real trajectory, longitude variations of optimal trajectory and longitude variations of real trajectory, respectively.

If the final score value is greater than or equal to zero it means that the trajectory is completely normal with respect to our algorithm, since it is equal or better than the optimal route that we found using A^* . When the value is negative it is an indication that the ship did not take the shortest possible path. The absolute value of the score gives us the level of abnormality.

If there is a need to label the trajectories as normal or abnormal, a simple threshold value (e.g. -1) on the final scores can be used to label the data. If the score is less than the threshold, the trajectory is considered to be abnormal and if it is greater than the threshold it is normal. In general, the score value is beyond a simple labeling and it can also give us a ranking between the abnormal patterns. One can rank the trajectories based on their score values and pick up the most abnormal ones. This is extremely useful in monitoring applications by interested parties.

IV. EXPERIMENTAL RESULTS

Three months worth of data from June to August 2013 for the North Pacific region were collected from the exactEarth database and used for the analysis. After dividing the trajectories by one day time frame, the dataset contained 39682 tracks. In this work, a track is taken to mean a trajectory for a certain vessel within a specific day. Ship types other than cargo and tanker were removed from the dataset and the proposed anomaly detection method was run on the remaining 21546 tracks to generate score values.

Figure 2 illustrates the distribution of the score values. The 21546 score values are plotted in an ordered way. As it can be seen, there are a few highly negative values, which indicate absolute abnormality by the algorithm, and the rest are either positive or a very small negative number. There are also very few highly positive values that have happened when the A* returned a longer path due to sparseness in the search graph. Lack of training data in some geographical regions leads to such sparseness in the generated graph. Another point that can be observed from the figure is that even by looking at the distribution of the scores, the user can easily pick the right threshold as the bending point of the elbow at the negative side.

For the evaluation of the method, a subset of this dataset was selected for labeling by the expert. 100 tracks from different ships and different geographical regions were chosen for labeling process. The tracks were chosen randomly in such a way that they uniformly cover the entire score space. The reasoning being that we wanted to evaluate the quality of the score function, and our ranking system, and to measure the correlation between the score values to the labels produced by the expert. The maritime expert labeled the data as normal or abnormal. The expert also gave us additional information about the possible reasons for the unusual movement pattern, but we only used the binary labels for the evaluation.

Since the proposed method generates score values instead of labels, we can plot an ROC curve for different score thresholds. Figure 3 illustrates the ROC curve for the proposed score function. As it can be seen, it has very high accuracy and the area under the curve is very good. The reason for



Fig. 2: Distribution of the score values for different tracks. 21546 tracks are ordered by their abnormality score.



Fig. 3: ROC curve for different score thresholds

such a high AUC is that the proposed features can extract the core geometrical properties of a trajectory in a helpful and efficient way.

Table I displays the confusion matrix obtained using -0.75 score threshold. The precision in detecting abnormal cases is 97.62% that is very good. In fact, we only have one false alarm by the method. The recall is 89.13% since the method has missed five abnormal patterns. The F-measure is 93.18% and the Correct Classification Rate (CCR) is 94%. This results show that the proposed score function has a good correlation with the labels that means it can detect anomalies in the vessel movements by analyzing trajectories from the geometrical perspective.

Figure 4 illustrates some of the abnormal trajectories taken by cargo and tanker ships. Since we are analyzing open



TABLE I: Confusion matrix for -0.75 score threshold

Fig. 4: Some examples of abnormal paths taken by cargo and tanker ships.

ocean area in North Pacific Ocean, any unusual trajectory bending or looping behavior should be detected as abnormal. All the score values in the figure are below -1.5 and they are detected as abnormal, since -0.75 score threshold is used for labeling.

V. CONCLUSION

We propose a new anomaly detection framework which is developed for detecting abnormalities in the cargo and tankers' movements. It is based on comparing the ship's trajectory with a near-optimal trajectory found using A* algorithm. To what extent does the real trajectory differ from the optimal trajectory is the basis for calculation of an abnormality score value. For this purpose, four scaleindependent geometric features were introduced and extracted from the trajectories. Providing a score value instead of labels has an advantage that the user can define a threshold for deciding on the label based on how sensitive they want to be in detecting anomalies. Moreover, it gives a ranking of trajectories from the most abnormal to the most normal. Experiments showed that the generated score values highly correlate with the labels provided by the expert. The area under the ROC curve (i.e. AUC) was very high that shows the effectiveness of the method. Another advantage of the proposed method is that it does not have any limitation on the size of the region of interest and it can be applied on very large geographical areas. Moreover, it does not require any training or parameter tuning which makes it robust and reliable when working with different datasets.

ACKNOWLEDGEMENT

The authors acknowledge the generous support of GSTS, Inc., ExactEarth, Inc. and the Natural Science and Engineering Research Council of Canada for this research.

REFERENCES

- [1] Heather Ball. *Satellite AIS for Dummies*. Wiley, Mississauga, ON, 2013.
- [2] A Dahlbom and L. Niklasson. Trajectory clustering for coastal surveillance. In *Information Fusion*, 2007 10th International Conference on, pages 1–8, July 2007.
- [3] T. Eriksen, AN. Skauen, B. Narheim, O. Helleren, . Olsen, and R.B. Olsen. Tracking ship traffic with space-based ais: Experience gained in first months of operations. In *Waterside Security Conference (WSS)*, 2010 International, pages 1–8, Nov 2010.
- [4] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Michael Wimmer, and Xiaowei Xu. Incremental clustering for mining in a data warehousing environment. In *Proceedings of the 24rd International Conference on Very Large Data Bases*, VLDB '98, pages 323–333, San Francisco, CA, USA, 1998. Morgan Kaufmann Publishers Inc.
- [5] ITU. M.1371 : Technical characteristics for an automatic identification system using time-division multiple access in the vhf maritime mobile band. Technical Report OGC 11-052r4 OGC 11-052r4, ITU, September 2014.
- [6] Chen Jiashun. A new trajectory clustering algorithm based on traclus. In Computer Science and Network Technology (ICCSNT), 2012 2nd International Conference on, pages 783–787, Dec 2012.
- [7] James B. Kraiman, Scott L. Arouh, and Michael L. Webb. Automated anomaly detection processor, 2002.
- [8] R. Laxhammar. Anomaly detection for sea surveillance. In *Informa*tion Fusion, 2008 11th International Conference on, pages 1–8, June 2008.
- [9] R. Laxhammar, G. Falkman, and E. 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, July 2009.
- [10] Jae-Gil Lee, Jiawei Han, and Kyu-Young Whang. Trajectory clustering: A partition-and-group framework. In *Proceedings of the 2007* ACM SIGMOD International Conference on Management of Data, SIGMOD '07, pages 593–604, New York, NY, USA, 2007. ACM.
- [11] Xiaolei Li, Jiawei Han, and Sangkyum Kim. Motion-alert: Automatic anomaly detection in massive moving objects. In Sharad Mehrotra, DanielD. Zeng, Hsinchun Chen, Bhavani Thuraisingham, and Fei-Yue Wang, editors, *Intelligence and Security Informatics*, volume 3975 of *Lecture Notes in Computer Science*, pages 166–177. Springer Berlin Heidelberg, 2006.
- [12] 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, June 2013.
- [13] B. Ristic, B. La Scala, M. Morelande, and N. 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, June 2008.
- [14] A. Vandecasteele and A. Napoli. Spatial ontologies for detecting abnormal maritime behaviour. In OCEANS, 2012 - Yeosu, pages 1–7, 2012.