

Clustering and Labeling of a V2V Communication Dataset based on CAM and DENM Messages with Malicious Data Injection: Analysis and Implications

Alessandra Blasioli

^aUniversity of Bologna MCs Computer Engineering

Abstract

This study explores vehicular communication systems' potential for enhancing road safety and reducing traffic congestion. By leveraging two initial datasets representing Cooperative Awareness Messages (CAM) and Decentralized Environmental Notification Messages (DENM), the research focuses on data preprocessing, clustering, and the development of a system to identify potential malicious outliers.

Keywords: V2V, clustering, labeling, outliers

1. Introduction

Vehicular communication systems are computer networks where vehicles and roadside units (RSU) serve as communicating nodes, facilitating the exchange of information. This information can include safety warnings and traffic updates. These systems have the potential to significantly reduce accidents and alleviate traffic congestion.

Vehicle-to-vehicle (V2V) communication allows wireless information exchange about the speed and position of nearby vehicles, offering great promise in accident prevention, congestion reduction, and environmental improvement. However, the full benefits of this technology can only be realized when all vehicles can communicate with each other. That's why the National Highway Traffic Safety Administration (NHTSA) has collaborated with the automotive industry and academic institutions for over a decade to make V2V communication's life-saving potential a reality.

In this work, we were provided with two initial datasets containing data from a simulation of a real vehicular network. These datasets represent two different types of messages: Cooperative Awareness Messages (CAM) and Decentralized Environmental Notification Messages (DENM). These are two message standards that we will explain in detail shortly. CAM messages represent vehicle movements, while DENM messages convey information about events such as accidents or traffic incidents that vehicles encounter along their routes.

The objective of this work was to observe the data, perform data preprocessing and clustering, introduce noise into the provided dataset to simulate potential malicious reports, and ultimately develop a system capable of labeling the data based on clustering and identifying potential outliers, which we consider as 'malicious.'

2. Dataset Analysis and Pre-Processing

2.1. Cooperative Awareness Message (CAM)

The Cooperative Awareness Message (CAM), defined by the European Telecommunications Standards Institute (ETSI) in 2011, is a crucial component of cooperative Intelligent Transportation Systems (ITS) networks. It provides a basic awareness service by periodically sending status data to nearby nodes in the network. CAM acts as an application support facility, distributing messages containing information about presence, location, and fundamental status.

CAM messages are broadcasted to nearby ITS stations, allowing them to be aware of each other's presence, positions, movements, and characteristics. The recipient evaluates the message's information to make informed decisions. The ETSI standard offers a predefined set of information for CAM messages, which can be used for various ITS applications. It also allows for the definition of new data elements to cater to specific application needs.

The CAM message format, according to ETSI TS 102 637-2, consists of a header and a body. The header contains message-specific details like version, identifier, and generation time. The CAM message body contains essential information about the sending ITS station, including:

1. A unique identifier for the sending ITS station.
2. The type of ITS station (e.g., mobile, public authority, private).
3. The reference position, which includes latitude, longitude, elevation, and heading.
4. An optional set of CAM parameters, following the standard's recommendations based on the ITS station type.

In addition to specifying the CAM message format, the ETSI standard outlines processes for message handling. It provides guidance on timing requirements for generating CAM messages periodically. The frequency of message generation can be adjusted based on the specific ITS application's needs. The

standard also describes rules for determining when a CAM message should be sent. However, it acknowledges that these rules are general, and architects implementing CAM facilities may have some flexibility in making certain decisions.

2.2. Decentralized Environmental Notification Message (DENM)

DENM, according to the European Telecommunications Standards Institute (ETSI) in 2011, serves as another type of application support facility, primarily providing a notification service regarding road status. While it was designed by ETSI to support active road safety applications, its utility can extend to any ITS application interested in acquiring information about road traffic conditions.

A DENM transmission is initiated by an ITS application that detects a relevant driving environment or traffic event. This application requests the DENM messaging facility to transmit DENM messages to notify others about the event. As per standard specifications, an event is characterized by various parameters:

- Event Type: An identifier associated with the type of detected event (e.g., vehicle breakdown, traffic jam).
- Event Position: Describes the event's location, which can be either a specific point or a geographical area.
- Event Detection Time: Represents when the event is expected to conclude.
- Destination Area: Indicates the geographical area over which the DENM message needs to be distributed among ITS stations.
- Transmission Frequency: Specifies how often DENM messages are issued by the same ITS station.

The ITS application must provide this information to the DENM facility, which is responsible for periodically transmitting DENM messages within the specified destination area and frequency. When the event's status changes, the ITS application notifies the DENM facility to update the information in the DENM message. Once the expiration time is reached, the DENM facility ceases to send DENM messages. Alternatively, events can be explicitly canceled by sending DENM messages to inform about the situation.

When an ITS station receives a DENM message, it assesses the information's relevance and determines the necessary actions (e.g., notifying the driver). It may also forward the message to neighboring ITS stations, whether vehicles or roadside units, to disseminate the information in the destination area specified by the originating ITS station. Unlike CAM messages, DENMs can be forwarded multiple hops away from the sending ITS station, covering longer distances. Roadside ITS stations, in particular, play a crucial role in collecting broadcasted information from vehicle ITS stations, processing it, and forwarding it to a central ITS station, benefiting traffic efficiency and management.

The ETSI TS 102 637-2 standard defines the structure of a DENM message, which consists of a header (similar to the CAM message) and a body containing event information. The event information is categorized into three sections:

- Management: Includes general event information, event source identification, event version, expiration time, sending

- frequency, event reliability, and an indicator of whether the event is false.
- Situation: Provides specific event details, including the event's cause, sub-cause, and severity.
- Location: Contains event location data, including hazard coordinates and extra trace location data.

In addition to message format and data elements, the ETSI standard offers guidance on DENM handling processes, including rules for forwarding DENM messages, detecting outdated messages, and managing situations where multiple ITS stations notify the same event. However, certain aspects, such as the communication path between central and roadside ITS stations for transmitting DENM information, remain undefined in the standard.

2.3. Pre-Processing

In the initial phase of our work, our focus was on preprocessing the data to prepare it for optimal utilization with the subsequent clustering algorithms. Both datasets share a similar structure, with semicolons separating fields within each row. We found it most convenient to manage these datasets using a pandas DataFrame. A pandas DataFrame is a two-dimensional data structure in Python used for data analysis and manipulation. It resembles a table with rows and columns, where each row represents a data point, and each column represents a specific attribute.

By employing a pandas DataFrame, we were able to handle the data effectively. Regarding the clustering process, we encountered the need to convert certain fields in the dataset into float data types. This was necessary because the chosen clustering algorithm, X-Means, which we will explain in detail shortly, requires float values as inputs.

Within the dataset messages, there were some fields represented in hexadecimal values. We addressed this by employing a function called 'hex_to_float' to convert these hexadecimal values into float format. Additionally, we encountered numerous fields with NaN (Not-a-Number) values, which we subsequently dropped from the dataset.

For fields that originally contained letters, we adopted an approach to convert them into binary format and then further into float values. This was done to retain the underlying meaning of the fields while ensuring compatibility with the clustering algorithm.

Through these preprocessing steps, we made it possible to analyze the entire datasets using the chosen clustering algorithm.

3. Dataset Analysis and Clustering based on Latitude and Longitude Parameters

In this section, our primary focus is the analysis of the CAM and DENM datasets using latitude and longitude parameters. The main objective is to identify messages originating from the same source, which show contradictory trajectories between the DENM and CAM datasets. By carefully comparing these trajectories, our goal is to detect any discrepancies that may indicate the presence of malicious messages.

In general, when a vehicle follows a particular trajectory as indicated by CAM messages, any DENM message it sends should originate from a location that aligns with the CAM trajectory. If a report deviates significantly beyond a predetermined threshold, it will be considered as malicious.

3.1. Clustering

The K-Means algorithm is widely used for clustering and requires specifying the desired number of clusters in advance. This can be a challenging task as the correct number of clusters may not be known beforehand. This is where X-Means comes into play.

The objective of the X-Means algorithm is to automatically determine the optimal number of clusters in the data without requiring a predefined specification. It is based on the concept of recursion, introducing an evaluation of clustering quality during execution.

Initially, X-Means performs a standard version of the K-Means algorithm with an initial number of clusters. It then calculates a quality measure, such as the Sum of Squared Errors (SSE), to assess the goodness of the obtained clustering.

Subsequently, X-Means checks if dividing a cluster into two sub-clusters could improve the overall quality of clustering. It uses the Akaike Information Criterion (AIC) or another similarity criterion to determine if the division would lead to a significant improvement. If the division is deemed advantageous, the cluster is split into two sub-clusters using the K-Means algorithm.

This process of division and evaluation is iterated for all existing clusters, allowing for the potential creation of new sub-clusters and overall clustering improvement. However, if the division is not advantageous according to the criterion used, the cluster remains intact.

The X-Means algorithm terminates when it is no longer possible to further divide the clusters or when the division does not yield a significant improvement compared to the added complexity. At the end, the final clustering with the automatically determined optimal number of clusters is returned.

The use of the X-Means algorithm has the advantage of avoiding the need to specify the number of clusters in advance, allowing for automatic determination. However, it is important to note that the X-Means algorithm may require more time compared to traditional K-Means, as it involves iterating through multiple divisions and evaluations.

We opted for the utilization of X-Means as we believe it is well-suited for our specific objectives. The X-Means algorithm is part of *pyclustering*, an open-source library for clustering, classification, and data analysis in Python. In our implementation, the first step was to observe the provided data, in order to gain a better understanding of the trajectories followed by the vehicles. These trajectories were recorded using CAM messages, as well as how these vehicles reported simulated events through DENM messages. To this end, we extracted from both datasets only the columns indicating the source and the position coordinates, namely latitude, longitude, and altitude, into two Pandas

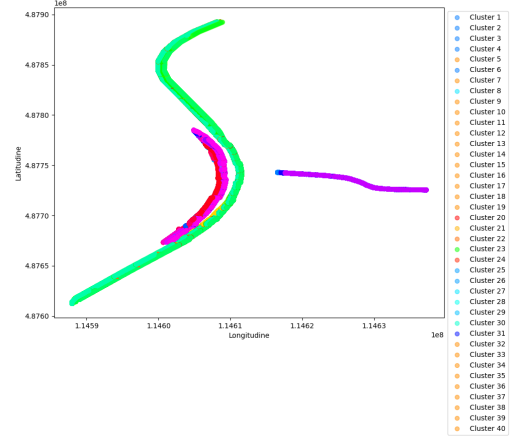


Figure 1: Clustering of the provided CAM dataset

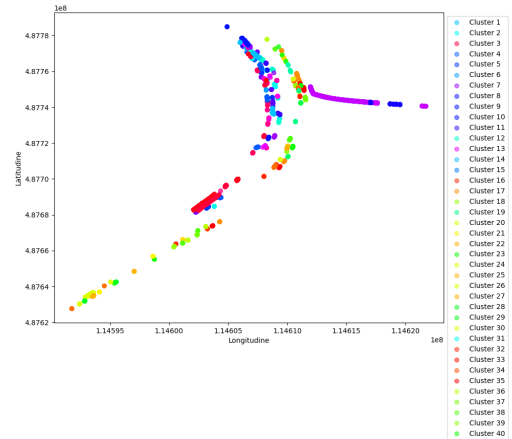


Figure 2: Clustering of the provided DENM dataset

Dataframes. As evident, the two trajectories described by the messages exhibit a high degree of similarity. This is, of course, because the data is clean, and the reports do not contain any potential malicious information.

To prepare the data for clustering, we employed the X-Means algorithm. Firstly, we calculated the "initial.centers" using the "kmeans_plusplus_initializer" method. This method takes the data and an initial number of centers and applies the k-means++ algorithm to determine the initial set of centers. (Note: k-means++ is a method that initializes the centers in a way that improves the efficiency and effectiveness of the k-means algorithm).

Finally, we called the "xmeans" method, which takes the data, initial centers, and a maximum number of centers (40 in our case) as input. This step initiates the X-Means algorithm and performs the clustering process on the provided data.

We followed the same process for both the CAM and DENM datasets.

3.2. "Data Contamination in the DENM Dataset"

The decision was made to introduce noise into this dataset because the primary objective is to simulate a dataset where potential malicious event reports, such as accidents, traffic disruptions, etc., are present. To accurately simulate the possible actions of attackers, it was initially decided to contaminate the dataset in the most realistic manner possible.

First, a function was implemented to determine new coordinates. We chose to vary the coordinates using a Gaussian distribution with a mean of 0 and a standard deviation of 1, multiplied by a factor of 100. This approach allows us to stay within a variation ranging from approximately 100 to 900 meters on the map. Since our dataset's coordinates are expressed in degrees, appropriate conversions were applied to enable the calculation of position changes.

Once the coordinates were calculated, it was decided to modify the data every 10,000 entries to maintain a situation that resembles reality. Subsequently, by duplicating the initial DENM dataset and replacing the coordinate values with the newly calculated ones, we saved the contaminated dataset into a CSV file named 'dataset_denm_dirty.csv'.

After the data contamination phase, we naturally applied the clustering algorithm to the contaminated data as well. Upon immediate observation of the generated graph, it is apparent that some points have been shifted from the previously observed trajectory. However, these points are not significantly distant from the event. Therefore, we may consider the contamination to be realistic and successful.

3.3. Event Type Analysis Based on Time Parameters

A crucial parameter to observe within the dataset is the type of events being reported, specifically indicated by the 'situation_eventType' field. Although DENM messages can encompass various event types, our simulation focuses on three specific event types:

- CauseCodeType_dangerousEndOfQueue = 27;

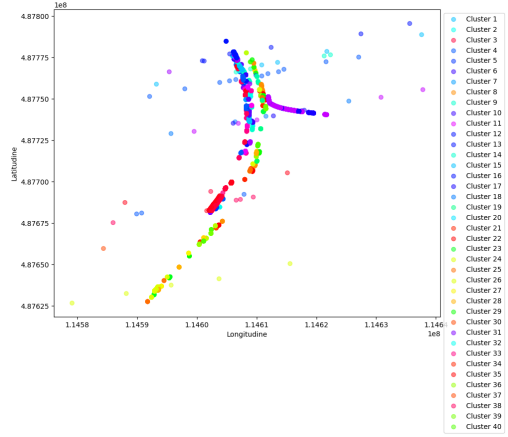


Figure 3: Clustering following data contamination

- CauseCodeType_collisionRisk = 97;
- CauseCodeType_trafficCondition = 1.

To monitor the progression of these reports over time, we leverage two parameters provided by the DENM dataset:

- 'simulation_time,' which represents the simulation time when the DENM message is received by an ITS-S, starting from 0.
- 'detection_time,' indicating the time at which the event is detected by the originating ITS-S.

So, with these parameters, we proceeded to calculate the time difference between the simulation time and the detection time. It's worth noting that the simulation time required an initial conversion for this operation. Unlike the detection time, which starts from 0, we needed to adjust the value to express it in UTC, similar to the detection time.

Once this conversion was applied, we were able to calculate the time differences between when a vehicle perceives an event and when it reports it. By combining this data with the types of events that occurred, we can observe, on average, the clustering of reports into three distinct time periods for the three different events, as we expected.

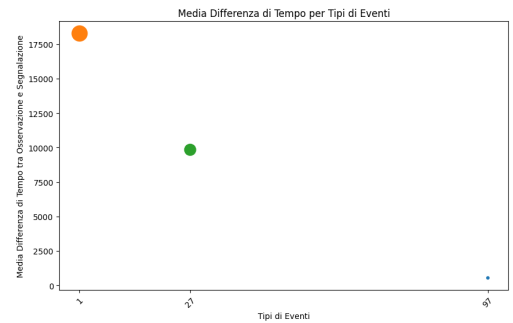


Figure 4: "Average Time Difference by Event Types"

Applying the clustering algorithm to the dataframe, where we have saved the results of the time difference calculations

along with all the previously mentioned data, we observe, as expected, three distinct clusters corresponding to the three events.

4. Labeling

In the final stage of our work, we focused on creating a labeling system capable of identifying potential outliers within the clustering dataset derived from the clustering results. To accomplish this, we opted for the use of an unsupervised learning algorithm known as the **Local Outlier Factor (LOF)**.

The Local Outlier Factor (LOF) is an unsupervised machine learning algorithm designed for outlier detection. It works by evaluating the local density of data points in a dataset. LOF calculates the density of data points by considering their proximity to neighbors. Dense regions have high local density, while sparse regions have low local density. LOF compares the local density of each data point to that of its neighbors. Data points with significantly lower local density compared to neighbors are flagged as potential outliers. Each data point is assigned an LOF score, quantifying its deviation from local density. High LOF scores indicate outliers, while low scores suggest normal data points. LOF requires a contamination factor, specifying the expected proportion of outliers. It influences the outlier detection threshold. A threshold is applied to LOF scores to classify data points as outliers or inliers. It is used to label data points as "outliers" or "inliers" based on their LOF scores.

LOF is a valuable tool for identifying anomalies in various domains, including fraud detection, security, and quality control. This algorithm assesses the local density deviation of data points to identify anomalies within a dataset.

After selecting a contamination factor, in our case, 0.0001%, we employed the 'fit_predict' function to predict the number of outliers within our clustering. To flag these outliers as such, we added a 'malicious' column to our dataset. In this column, we marked entries as 'True' if they were identified as outliers and 'False' otherwise, using the straightforward command '*datadirty['malicious'] = outlier_predictions == -1*'. As we can see, the results are quite satisfactory; the algorithm is able to detect a significant portion of anomalies. By utilizing this system, we can automatically identify malicious messages.

References

- [1] Zhaojun Lu, Qian Wang, Gang Qu, and Zhenglin Liu; BARS: a Blockchain-based Anonymous Reputation System for Trust Management in VANETs
- [2] Canhuang Dai, Xingyu Xiao, Yuzhen Ding, Liang Xiao, Yuliang Tang, Sheng Zhou; Learning Based Security for VANET with Blockchain
- [6] Abid Ali, Muhammad Munwar Iqbal, Sohail Jabbar, Muhammad Nabeel Asghar, Umar Raza, Fadi Al-Turjman; VABLOCK: A blockchain-based secure communication in V2V network using icn network support technology
- [4] anagiotis Papadimitratos, EPFL Levente Buttyan and Tamás Holczer, Budapest University of Technology and Economics, Elmar Schoch, Ulm University, Julien Freudiger and Maxim Raya, EPFL, Zhendong Ma and Frank Kargl, Ulm University, Antonio Kung, Trialog, Jean-Pierre Hubaux, EPFL; Secure Vehicular Communication Systems: Design and Architecture

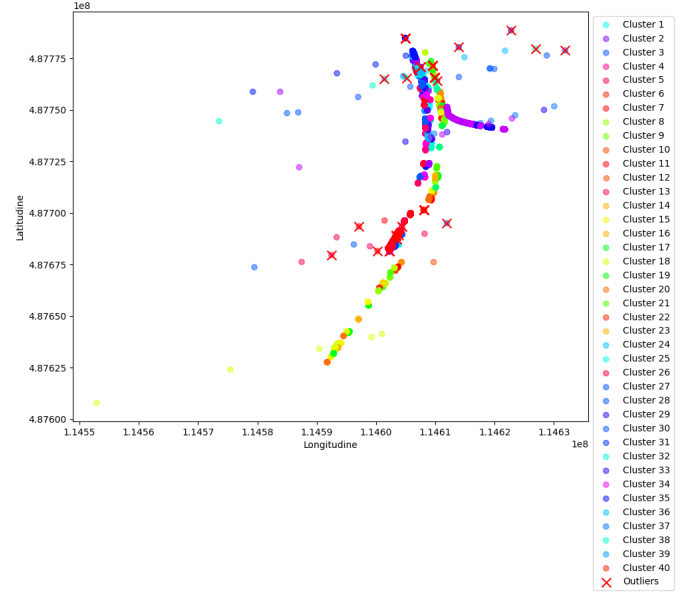


Figure 5: Results of LOF prediction

- [5] Xuanxia Yao, Xinlei Zhang, Huansheng Ning, Pengjian Li; Ad Hoc Networks: Using trust model to ensure reliable data acquisition in VANETs
- [6] Hichem Sedjelmaci, Sidi Mohammed Senouci; An accurate and efficient collaborative intrusion detection framework to secure vehicular networks