# University of Huddersfield Repository

Verma, Shrey

Generating Synthetic Automotive Data and Detecting Abnormal Vehicle Behavior Using Unsupervised Machine Learning

**Original Citation**

Verma, Shrey (2022) Generating Synthetic Automotive Data and Detecting Abnormal Vehicle Behavior Using Unsupervised Machine Learning. Masters thesis, University of Huddersfield.

This version is available at http://eprints.hud.ac.uk/id/eprint/35747/

# Generating Synthetic Automotive Data and Detecting Abnormal Vehicle Behavior Using Unsupervised Machine Learning

Shrey Verma

A thesis submitted to the University of Huddersfield in partial fulfilment of the requirements for the degree of Master of Philosophy

February 2022

**Generating synthetic automotive data and detecting abnormal vehicle behavior using unsupervised machine learning**

**Shrey Verma**\*

## Abstract

The amount of data generated, processed, and stored by the modern vehicle is increasing and this is creating the potential to detect abnormal and potentially dangerous situations occurring. The purpose of this thesis is to portray a lack of information in the area of intrusion detection using automotive data and to lay the foundations of research in intrusion detection using unsupervised machine learning. As vehicles continue to become more connected, there is an increased possibility of them being exploitable through a successful cyberattack. An example of a hacked Jeep Cherokee (Miller, Valasek, (2011)) and a remote exploitation strategy using multiple attack vectors (Checkoway et al, (2011)) was the prime exhibition of a situation where the vehicle can be remotely compromised. These examples demonstrate the potential to exploit aspects of the vehicle's communication and control systems, resulting in expected behavior. This thesis is focused on detecting attacks targeting a vehicle by identifying abnormal vehicle behavior, exhibited through control data. To achieve this, synthetic vehicle data containing detectable abnormalities is generated and used for analysis and detection to help detect cyberattacks. Unsupervised machine learning techniques are used as a way to detect abnormal entries in-vehicle data. the synthetic data is generated based on datasets comparable with those generated during normal vehicle operations, before being used to insert manually insert skewness to generate abnormalities, before using and evaluating various unsupervised learning algorithms

## Table of Contents

## Introduction

There has been progressive change in the vehicle industry, whereby innovations are coordinated into the vehicles, changing end-user functionality for reasons such as enhancing well-being, execution, and proficiency. As vehicles turn out to be more connected, they are becoming increasingly susceptible to cyberattacks. Koscher et al. 2010 demonstrated how in situations where it is applicable to control an extensive variety of fundamental capacities: disabling the brakes, specifically targeting the braking system on each singular wheels, halting the engine, and so forth. The user experience while driving has changed at an increasing rate because now there are fewer chances of people getting lost at unknown places, people can operate their phones in the vehicle (without using the phone itself) with limitless connectivity and there are various ways to get entertained for the people sitting inside the vehicle. These changes come with a risk of data and privacy breach, Kyriakidis, (2015) studied that the civilians of developed countries were worried about their data being transmitted from their vehicles. The data, which the vehicle uses to give the user increased feedback in the form of information and entertainment, is at a high risk of being breached. The data can be extracted and can be used for various malicious activities as Enev, (2016) concluded in his paper that a person could be identified with 87% of accuracy using just one sensor and with an accuracy of 99% with the data coming in from five sensors combined.

The necessity to protect the vehicle and its passengers, as well as the vast array of attack vectors, has resulted in a fertile area within the field of cybersecurity in Connected and Autonomous Vehicles (CAVs). Researchers such as Moore et al. (2019) and Han, Kwak, Kim (2018) have been researching this area and processes have been under development for identifying and finding ways that focus on vehicle security and communication.

There is an absence of information detailing how data can be extracted and analyzed, as well as how this can be essential for the security of the vehicle and its occupants. The detection of anomalies in the data extracted from a vehicle's control systems could be a pivotal part of the detection of cyberattack and abnormal behavior of the vehicle. In this thesis, different unsupervised machine learning algorithms have been used to determine their effectiveness in detecting anomalies in the automotive dataset, which provides empirical analysis to determine the more appropriate techniques that can be used for a data generated by the vehicle's control systems. Furthermore, this study will provide the necessary groundwork in creating a framework that can be used to prevent intrusion. Unsupervised machine learning is better suited for our approach as the behavior of a driver may vary from one driver to another which would result in different vehicle data patterns, and through using unsupervised machine learning, the authors were able to overcome this challenge through handling diverse driving patterns. Unsupervised machine learning algorithms is capable of processing vehicle data in real-time, training its model with available data sources, whereas supervised learning only focuses on a given set of labeled data and it can then work out how to detect changes in pattern. (Min-Joo Kang Je-Won Kang, 2016). Unsupervised learning has earlier been used by Gupta, Amruthnath (2018) to detect faults in systems as well as predictive maintenance of systems.

Unsupervised Machine Learning algorithms can be used to detect abnormal behavior in vehicular data. In the future, it is foreseen that this will be an essential and commonplace method as it focuses solely on the live

processing of data simultaneously while the vehicle is in use to detect patterns in the ways of driving which would otherwise be considered as not normal. Previous work by Abdulaziz et al. (2018) demonstrates the detection of intrusion in a vehicle using supervised learning where they already had specific intrusion-related data. There are other related works where unsupervised machine learning is used to detect the behavior of surrounding vehicles (Morris & Trivedi 2009). Unsupervised learning has also been used for driver safety where patterns in the position of a driver while driving a vehicle are observed, trained, and tested to determine safe and unsafe positions/posture of a driver while in control of the vehicle (Veeraraghvan, Atev, Bird, Scharter & Papanikolopoulos, 2005). A challenge still exists to determine what behavior of the vehicle could be considered as abnormal or anomalous, where sensed data makes a sufficiently different and unusual change from other values in a very randomized pattern, or the data changes abruptly from a previously normal profile by a uniform increase or decrease in the values of the data. The unsupervised machine learning should be able to to determine if the vehicle data is normal or abnormal, which could indicate that it has been compromised.

This thesis is structured as follows: first, a literature review of articles concerning related work closer to the field of interest of this thesis is provided. In the next section, the methodology of the exploratory analysis is presented. This leads to the discussion of results that have been achieved by the process and it is followed by a discussion of the results obtained.

## Literature Review

There are limited works published using Unsupervised Machine Learning to detect anomalies in the behavior of a vehicle, as it is currently a developing and fertile research area. Previously researchers have used techniques to prevent security breaches or intrusion, and in this section, a comprehensive analysis and categorization of related research is performed to identify key gaps for further research, At the core of any modern vehicle's internally connected frameworks is the Controller Area Network bus. The CAN transport network is brought together on which most of a vehicle's information movement is communicated. The CAN transport conveys everything from administrator orders, for example the sensor reads, "lower the windows" or "apply the brakes", detailing engine temperature or tire pressure. The approach of CAN transport achieved upgrades in proficiency and a decrease in complexity while likewise reducing wiring costs. Before the advancement of CAN bus innovation, any two vehicular segments expecting to speak with each other would have required a separate physical communication channel between them. (Farsi, Barbosa & Ratcliff, 1999)

As of late, a meaningful step forward in an automotive framework has been made with incorporating various processing components called Electronics Control Units (ECUs). An ECU is used for controlling and checking vehicle subsystems and is of vital importance. The ECU replaces customary mechanical controlling components (Park, Han, Lee, 2015). Automotive networking like Vehicle to Vehicle or Vehicle to Infrastructure requires figuring gadgets to perform communications within the vehicle and among other vehicles (Tuohy et al 2015). The vehicular comparison can be applied to numerous functional traffic frameworks (Lenz et al 1999). Tang et al.

propose to utilize communications to comprehend driving practices, for example, every vehicle's speed and fuel usage (Tang et al, 2014). Jin et al. showed the powerful V2V interchanges relying upon a traffic stream built up a subtle messaging system in the communication. In (Yu, Shi, 2015) productive fuel usage is assessing the paces of the associated vehicles or their separations. It is also true that helpful platooning empowered by the remote interchanges can likewise improve traffic stream. Moore et al., 2017 used a technique to model inter signal arrival delays to detect intrusion in the CAN bus of the vehicles. They focused on the electronic side of the aspect to prevent intrusion in a CAV. They use a simple anomaly detection system that monitors the inter-signal delay of the CAN bus traffic that would provide accurate detection of a regular frequency injection attack. They observed that CAN bus signals for normal settings (50 seconds – 25 seconds with the vehicle on and engine off and the same with plus 25 seconds with vehicle on and engine off), for a fixed ID there has been the regular movement of signals with little noise. The runtime lights sent at a fixed interval with either on or off values is a specific example showing even when not in use the light signals are sent. They observed that the time between signals never differs more than 24% from their mean observations after carrying out calculations.

They attacked three times and the first two attacks were during two basic vehicle operation modes, engine on and engine off, the first attack turned the runtime lights from on which is normal to off which is attacked and the second attack did just the same, these two attacks had to be done repeatedly to inject desired signals into the CAN bus at a rate very high rate to suppress the actual ambient signal which is being sent regularly by the vehicle and these two attacks were designed to be stealth-based. The third attack was less stealthy where a Denial of Service (DoS) attack was launched involving the repeated injections of a signal resulted in shutting down of the vehicle. They also noticed that their hardware implementation allowed Bluetooth connectivity to an Arduino board and hence it could be triggered from outside the vehicle.

The authors have focussed solely on attacks on the vehicle via signals that are being injected through the CAN bus. They demonstrate turning the vehicle off/on at an unwanted time and affecting the functionality of the vehicle with events like disabling brakes/throttle etc. These types of injected attacks do not focus on the data availability in vehicles. While both relate to the security of the vehicle and intrusion detection into the vehicle, there is a major difference with data availability in Connected and autonomous vehicles. Vehicle data logging is a critical area, which has not yet been exploited, and the authors mentioned above solely focus on a signal-based intrusion in a vehicle, which is possible on regular vehicles as well. Connected and autonomous vehicles are also vulnerable to one more threats in addition to the aforementioned.

Abdulaziz et al. (2018) proposed a classification approach for intrusion detection using Supervised Machine Learning. It used two algorithms which were based on K Nearest Neighbour and Support Vector Machines. They used these algorithms to detect Denial of Service and Fuzzy attacks on vehicles. They obtained a dataset from their source which consisted of 300 intrusions of message injections and each intrusion got a timeframe of 3 seconds to 5 seconds to perform malicious activities. They already had a specified dataset that highlighted the

intrusion messages and their approach of Supervised Machine learning with the above techniques paved the way to attain results.

Morris and Trivedi (2009) used unsupervised learning to allow a vehicle to learn patterns in its surroundings automatically and in the process extract natural behavior carried out by a vehicles driven on highways. The focus was on a data-driven approach to learn behavior patterns of a vehicle, driven in an unsupervised process to detect objects in its surroundings. The entire project was conducted because it is easier for a vehicle to learn patterns on when to change lanes, apply brakes, or for a given case even turn directions. The same thing should be possible for a vehicle when processing its surroundings of a foreign object (for example another vehicle) and to act in a way where both the vehicles would know how each other would react to situations that they would be facing on the highway.

Han et al. (2018) used a survival analysis technique to detect intrusion in vehicular networks. The survival analysis model relates to the time taken to detect intrusion through event analysis. This type of analysis focuses on statistics and temporal properties of an event. The method requires information of the measurement object, its survival time, and status. The survival time is a period from the start to finish point and is used for measuring the status, which is used to confirm if the event has survived at that point. Any uncertain data regarding the occurrence of the event is called censored data and they used all this information to conclude the following definition of survival function at the time: (t) = number of events beyond the time (t) / total number of events.

Han et al. (2018) in their paper mention a very different approach to detect intrusions in vehicles. It requires an event to occur which is then checked to a threshold limit, which is repeatedly is updated as the new event occur, and then if the survival probability does not fall within the threshold alert output is triggered.
Diverse communication conventions are created to help with aspects of the event processing (Fan, Dao, Crolla, 2008). Veeraraghavan et al (2005) developed a camera-based system to monitor the activities of an driver with the primary goal to distinguish between safe and unsafe driving positions using unsupervised learning. Unsupervised Clustering was used in this project and was built on the behavior patterns and appearance of a vehicle driver while driving the vehicle, which included aspects such as skin color detection, detection of changes in driver behavior, and differentiation between action models of the driver in regard to his/her general way of being seated in the vehicle. They were successful in presenting two different modes for monitoring safe and unsafe driver activities while driving under challenging circumstances, including proper visibility of the driver at different levels of light exposure. This technique can help to contribute to the driver and automotive safety.

There has also been recent works considering wellbeing issues of a vehicle and intra vehicular communications. Specifically, an Intrusion Detection System (IDS) requires a lot of consideration because of the effectiveness of any implemented techniques to correctly distinguishing the compromises (Muter, Groll, Freiling, 2010). Hoppe et al. proposed an intrusion recognition strategy by utilizing predefined attack sequences, stored in a database.

Larson et al. build up a specification-based methodology, looking at behavioral patterns of the system and how they match to those stored in the database by using a different sensor intended for vehicular compromise situations. Guo et al. (2000) devised a diagnostic system that used external vehicular data segmentation and machine learning algorithms, such as fuzzy systems to determine the data to be good, bad, or unknown. The mechanism was divided into five different fuzzy sets. The system used external vehicle data and was translated it to an internal format, and then used a segmentation layer to segment the signals into sections corresponding to the vehicle's physical state for example acceleration, idle, etc. Then a different layer was used to generate a vector of features for each signal segment obtained in the signal processing layer. The vector elements are numeric inputs, for example wavelet or Fourier coefficient energy that represents a feature of the signal. The next layer generated a super vector of features from the primary signals, and it could also contain features from any of the previous segments to establish a relation with time dependencies. A combination of a fuzzy system and a neural network was used to provide machine learning in the final layer. This was capable of determining the signals that they got in the previous layers and then label them as bad, unknown, or good signals. (Guo, Crossman, Murphy, Coleman, 2000). The system that Guo et al. (2000) introduced was one of the first processes to diagnose vehicular data using machine learning to train a model of signal inputs and use them to compare to live vehicular data and determine them to be good, bad, or unknown. These types of research fueled other vehicular data including the position of the vehicle, speed, and gyroscopic vehicular data to be used to generate and carry out vehicular security in connected and autonomous vehicles where there is a large quantity of data that need to be protected and constantly monitored to detect and prevent intrusion.

Petit and Shladover (2015) focused on cyber-attacks on automated vehicles, formulated various threats, and tabulated them to show priority levels of damage that an attack can cause to the vehicle or the driver. They use the following terminology to define what attacks could mean and what impact they can cause to the security of the vehicle and security and or privacy of the driver: -

- Means – Which describes the attack and on what attack surface it has taken pace.
- Feasibility – This describes the amount of work put into the attack to successfully execute and refers to the technical expertise required to perform an attack.
- Need for physical access to the targeted vehicle – This confirms if the attacker requires access to the vehicle to perform the attack.
- Ease of detection by a driver – This describes if the driver can detect if they are being attacked by feeling that something is wrong in their driving or by seeing a display message on the display panel.
- Ease of detection by the vehicle – This describes if the Intrusion System Detection of the vehicle can detect the attack.
- Probability of the attack success – This is a result of previous criteria which can lead to answers such as if a highly feasible attack which is easily detected is more likely to not succeed.
- The consequence for the vehicle – Describes how the vehicle will perform after all the criteria mentioned above are checked and fulfilled such as entering a Minimal Risk Condition.

- The IDS incorporates different modules for collecting and breaking down or analyzing a substantial amount of information. Commonly, the monitoring module detects a signal after it gets extracted. Offline trained and traits were clearly observed in the profiling module. The profiling module refreshes the database for further packets that will arrive in the profiling module if the monitoring module detects another attack. (Min-Joo, Je Won, 2016).
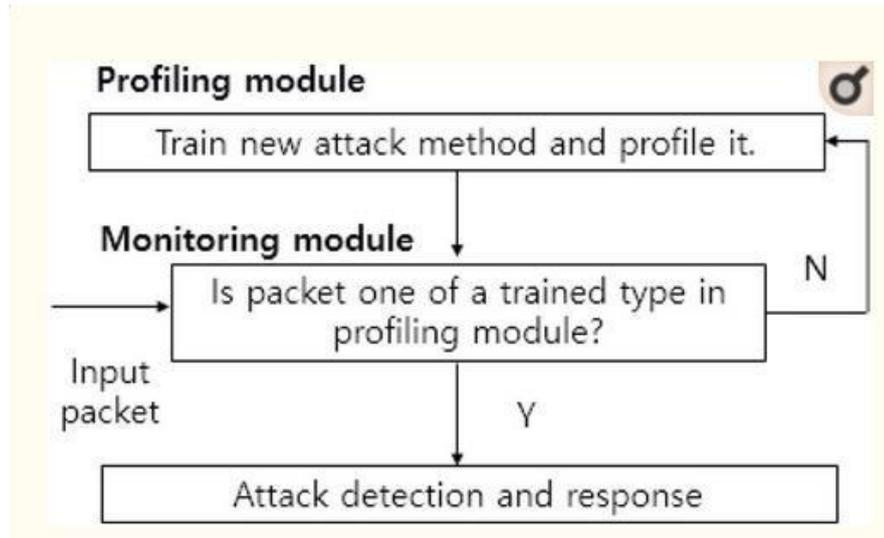


*Figure 1 IDS BASED ON MACHINE LEARNING TECHNIQUES – ARCHITECTURE (Min-Joo, Je Won, 2016)*

Artificial neural systems and support vector machines (SVM) are applied to identification of intrusions, utilizing a statistical model on a piece of packet information (Min-Joo, Je Won, 2016). Chen, Hsu, Shen in 2005 showed recurrence-based encoding strategy that is used for the network packet level and includes the use of an ANN and SVM. The previously mentioned works depend on regulated machine learning techniques, and, along these lines, various marked data are required in the preparation. When contrasted with the methodology, Kayacik et al. use a solo machine learning technique, for example, a self-organizing feature map (SOM) for network-based identification of intrusions. (Min-Joo, Je Won, 2016)

The above-mentioned intrusion detection strategies might be used just for specific attacks that have been already considered in configuration stages (Sun et al, 2015). To adapt to such conditions Machine Learning based IDS methods are used, for the most part, for conventional communication systems. (Deepa, Kavitha, 2012).

Machine Learning was also proposed as an authentication scheme for automobile driver fingerprinting. Machine learning techniques have good generalized capabilities and can be applied    in numerous fields because of their incredible component extraction and introduction capacities (Huang et al, 2018). For example, Xu et al, 2018 proposed a reinforcement learning (RL) based occupation planning calculation by consolidating RL with neural network (NN) to take pursue cost minimization of large information examination on geo-disseminated servers associated with sustainable power sources with unusual limit.

Based on past research, Yijee et al, 2019 used Machine learning algorithms to undertake research on automotive driver fingerprinting. Xun et al. accomplished driver recognition through utilizing a k-NN distance-based approaches. Notwithstanding, there were still a few restrictions. Specifically, the normally recognition accuracy of ten drivers was just 51%, which is only slightly better than random chance. Their approach has show to be poor at distinguishing illegal clients. To begin with, the creators dismissed the potential of the driver's different attributes being firmly related to physical characteristics of the vehicle. These are for instance, the accelerator pedal size, the brake pedal size, and the steering wheel pivot point are facilitated when the vehicle turns. Second, they neglected to watch the way that the attributes relate to the transient length. (Y. Xun, J. Liu, N. Kato, Y. Fang and Y. Zhang,.2020).

Cyber-Physical Systems (CPS) require control, calculation, and communication capacity to manage physical components (Shi et al, 2011). CPS have now been applied to different applications including power stations, human services frameworks, smart grid frameworks, networks, industrial control, and vehicles (Mo, Yilm et al, 2012). As of late, the advancement of the automotive CPS with consistent availability and variety of connection has uncovered basic vulnerabilities that have resulted in new dangers against automotive security. For example, obscure attacks through network infrastructure. Likewise, during driving, there might be peculiarities brought about by the vulnerability of dynamic

situations, and issues may only present themselves after some time. A vulnerability or flaw is an irregular state which may prompt blunders or disappointment of the framework, including lasting, transient, and discontinuous shortcomings (Jo, Minsu et al, 2016). Note that aside from these vulnerabilities, different deficiencies can emerge through a digital attack, which may affect infrastructure in different ways. Vulnerabilities could be discovered and exploited on any sensor in each framework. Uncommonly, the automated CPS is outfitted with a wide range of sensors, for example, (GPS), (IMU), sensor with ultrasonic capabilities, vision, and wheel encoders to improve the wellbeing of travelers in self-sufficient driving. For instance, if an adversary spoofs the GPS of a self-ruling vehicle through the internet, physical harm could occur since it cannot establish definite information as to the position of the vehicle (Koscher, Karl et al, 2010, J. Shin, Y. Baek,Y. Eun and Son, 2017).

A significant issue in digital-physical attacks, identified with hybrid sensors of automotive CPS concerns, is a basic sensor attack. If the in-built sensors are attacked in isolation, it can be difficult to recognize without prior information. The presentation of attack location strategies is constrained to identifying confounded, simultaneous attacks and inconsistencies, or is dependent on known attacks. It is imperative to misuse a lot of notable information that normally include gathering sensor data to recognize these attacks as irregularities in the automotive CPS. Deep Neural Network (DNN) systems, named profound learning, can be an acceptable answer for recognizing abnormality from ordinary circumstances without prior knowledge the attack. This is because those models are equipped for finding confusing relationship among infomation and yield by elevated level reflections without master information even though most of the sensors have the issues. In recent research, researchers have tended to security issues concerning system faults and attacks. Fawzi et al. examined an elevated level framework configuration to improve the security of a general control framework (Fawzi et al, 2012). They have illuminated these issues without thought to the necessities of the focused-on applications since their framework is not intended to target specific applications in explicit spaces (e.g., a self-governing vehicle or modern procedure). Sabaliauskaite et al. examined the strong location framework free of the internet of CPS against digital attacks (Giedre ete al, 2013). They structured the watchdog framework which has savvy sensors to screen and check the physical datasets. As per pre-characterized rules, the framework identifies attacks. It is relevant just to a specific application with characterized basic parameters to be observed (J. Shin ,Y. Baek, ,Y. Eun and Son, 2017).

Although there are a variety of systems in the area of sensor attack recognition advancements, they are centered on quantitative techniques. Among quantitative strategies dependent on explanatory excess as an elective way to deal with process repetition, model-based techniques are recognized as one of the most promising methodologies since the 1990s. They center around producing and assessing an outlier being a reliable indicator to distinguish attacks on sensors. The probability of compromise estimation has been generally used as one of the best  way to deal with identification of an attack a considerable lot of these attack discovery techniques are created dependent on a direct framework. For instance, the

Kalman channel utilizes elements and displaying of a framework and requires extra conditions to anticipate or appraise sensor vulnerabilities (Zhang et al, 2008). It is not reasonable for a nonlinear condition, for example, a self-sufficient vehicle because the Kalman channel for the linear time-invariant framework is accepted (Isermann et al, 1997).

There are Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) as explicit kinds of Recurrent Neural Network (RNNs). Exceptionally, LSTM is equipped for tending to the evaporating slope issue of the RNN, which is connected to the effect of contribution on the concealed states which progressively disappears as the intermittent association of the RNN formation is persistently rehashed (Bengio et al, 1994). The memory cell of LSTM comprises an inward memory cell and some additional entryways (the information, yield, and neglected entryways). The inner memory cell deals with the progression of data and another entryway decides to dispose off the data in the interior memory cell. A data entryway decides to acknowledge new data as information or not and a yield entryway chooses if it needs to use the datasets (Hochreiter et al, 1997). GRU likewise explains the long reliance and has a few entryways. GRU is an altered variant from LSTM by lessening the number of entryways and evacuating the inner memory cell in the memory square (Cho et al, 2014).

### Vulnerable Mode

Physical framework models used by the calculations will enable the controller to recognize the present condition of the physical framework given the observed information. For example, an object identification calculation characterizes the recorded picture from the camera dependent on a pre-prepared AI model. In any case, the aggressor can constrain the framework to get into a hazardous state. For example, the article location calculation can be tricked to group a picture incorrectly with the most elevated likelihood by marginally tweaking the picture (not noticeable to the unaided eye. Then again, the BMS might be attacked to expend more vitality from the battery cell (Ragunathan, Rajkumar,2010).

It is mindful to understand and become familiar with the ordinary function of the physical framework that works along with the controller. The design is prepared by running the CPS (Common Part Sublayer) and observing the system at the desired time to the user, before any attack can occur. Following this, at both sensing and anticipate stage, the computer will analyse the system to distinguish any abnormality or attack and recuperate from them at run time. Additionally, it will likewise be updated and learn new elements, providing the pace of attacks recognized is low in a specific time window (Park, Kim, Chang, 2016).

Commonly, the constant changing behavior of the control loop holding the physical system and the control elements can be designed using mathematical equations and deterministic modeling. One can see complexity in the behavior of physical systems while modeling it through equations. Moreover, there can be various hidden factors posing a threat while modeling. The answer to this issue is data driven statistical modeling (Machine Learning). However, this method of modeling also faces challenges from undiscovered biases and forecasting errors. Therefore, the author proposes a Generative Adversarial Networks (GAN) for such behavior to be captured. This would also enable self-secured control (Parvar, Faruque, 2019).

Generative Adversarial Networks (GAN) inculcates a machine learning framework to produce a stable model. As per GAN, two neural networks exist, discriminator and generator. GAN results in a way that the neural networks can be modeled as a minimax game using game theory. Hence for the generator to succeed, it needs to understand to produce the distribution of the real data for the physical process efficiently, so that the discriminator is unable to differentiate. On to the contrary for the discriminator to succeed, it also needs to understand the distribution of the real data such that the generator cannot deceive it. Hence to beat each other in this game, both would be required to be highly efficient generating and discriminating (Parvar, Faruque, 2019). So at the point of equilibrium, the optimal point for the minimax game, the real data shall get modeled by the generator and the discriminator would fetch a result equaling 0.5 of probability since the real data is same as the output of the generator. (Trippel, Weisse, Xu, 2017).

Determination of GAN engineering for the control configuration would be for the most part because of the accompanying elements:
1) Competition between two systems can give quicker convergence.
2) Two systems are as of now being prepared for purposes required for a secure about control plan.
3) The model will be progressively powerful towards any attacking model particularly ill-disposed models without the requirement for application-explicit information.

There is a requirement for both the neural networks to learn the physical process of the control loop for a minimum number of time steps. Sampling is performed from the signals derived from the sensors and actuators in the loop. For a particular period, the discriminator validates the data and looks for compromises. On the contrary, the generator tries to deceive the discriminator by producing the same data (Parvar, Faruque, 2019). The Conditional Generative Adversarial Network (CGAN) discriminator encapsulates the true dynamic behavior of the control loop for the particular amount of time steps given the conditional data and its time steps. Since its competing with the generator, it discriminates the fake data hence making it more defensive towards adversarial attacks. The physical process of the control loop shall be compromised if there is any attack on the unprotected physical system, sensor, or model. Thus, the discriminator can provide a probability of identifying an anomaly during the period with a particular data and condition. To make a deterministic decision, the conditional anomaly needs to get detected which in turn helps in capturing the physical condition of the system. During the time, when

the probability of a said batch being normal is less than the fake threshold, the control loop is compromised. (Parvar, Faruque, 2019).

### Training and Adjusting

The fundamental challenge of the CGAN design is training two neural systems. The two systems need to converge at an equilibrium, if not the discriminator easily differentiates between the genuine or fake information or the generator consistently creates information which are close to being genuine. Henceforth, at the training stage, the manufacturer is capable to run the controller for a specific timeframe to prepare the CGAN with genuine data. To minimse their loss functions two optimisations will take place during the training phase itself. (Parvar, Faruque, 2019)

After the CGAN has been trained, the two entropies will have small and stable values. In any case, the preparation won't stop after the training stage. The CGAN will be keep on training at the detect and- predict stage given the new clusters of data that was sampled at run-time. In any case, just the consecutive patches of data are used for training that their likelihood of having no anomaly is higher than  characterized trust limit for an enormous timeframe. Then, at the detect and predict stage, the CGAN is used for anomaly locating and detection and recuperating predictions. (Parvar, Faruque, 2019)

## Machine learning technique in automotive system

In creating a vehicle in sight, the car business has been taking two distinct ways. In the more preservationist way, customary automobile organizations start from assistive driving and advance step by step to significant level. Specifically, ADAS frameworks have much evolved and embraced by most significant automakers, for example, BMW, GM, Honda, and Toyota. More dynamically, IT organizations and research foundations legitimately intend to completely promote self-driving, for example, Waymo, Uber, and TuSimple. Observing a modern vehicle framework commonly utilizes various hybrid sensors for observing the environmemnt, for example, LiDARs, radars, GPS, inertial estimation unit (IMU), sonars, and cameras (Lindgren and Chen, 2006). The multi-modular data from these sensors are intertwined and examined to improve framework dependability and safety. Specifically, ML (Machine Learning) procedures are broadly used for handling the information gathered by cameras, and these vision-based ML modules altogether add to the advancement towards computerized and self driving (Lan, Huang, Wang, Liang, Su and Zhu, 2018). Also, perception techniques have additionally been used for drivers and pedestrians. To diminish the interruption from interfaces, hand gesture acknowledgment has been experimented in vehicles. Driver's actions action activity gestures approaches to screen driver behavior and decide to assume control to overtake has been developed. Approaches for recognizing pedestrians have been developed as well. (Endres, Schwartz and Muller, 2011; Zobl, Nieschulz, Geiger, Lang and. Rigoll, 2003).

Choice and Control: Based on the discernment and comprehension of the general condition, a ML-based system in modern vehicle frameworks targets creating a sheltered and effective activity plan consistently, with various functions, for example, forecast, lane segmenting and obstruction evasion.

There are various methodologies that uses ML strategies to choose a driving activity legitimately dependent on the tangible info. Beginning from the late 1980s, neural systems have been used to delineate camera picture input linearly at the corners, to keep the vehicle on the right path (Pomerleau, 1989). ML procedures have likewise been used for other control in-vehicle frameworks, for example, engine management. A neural system is used to anticipate explicit fuel utilization and exhaust temperature of a diesel engine for different fuel injection timings (Zhang, Fiddler, and Urtasun, 2016).

V2X Communication: Machine learning techniques are used to identify problems in V2X communication and interrelated vehicular applications. The hybrid centralized strategy utilizing k-means clustering is addressed towards control congestion in vehicular ad hoc networks (VANETs). Ide et al., 2015 has addressed in the paper about LTE connectivity prediction and vehicular traffic prediction along with Poisson dependency network. There are various reinforcement learning models in place for user association, load balancing, vertical handoff in heterogeneous networks, routing for local data storage in vehicular networks and virtual resource allocation on vehicular clouds. (S. Lan, C. Huang, Z. Wang, H. Liang, W. Su, Q. Zhu, 2015)

## Clustering

Mixture of progressive branching methods which consolidate the attributes of various partitional Clustering strategies or partitional and various leveled branching procedures. In this, proficient base up to hybrid leveled branching (BHHC) procedures has been proposed with the end goal of model determination for protein succession arrangement. In the primary stage, a steady partitional branching strategy, for example, pioneer calculation (requested pioneer no update (OLNU) technique) which requires just a single database (DB) filter is used to locate a lot of sub-group agents. In the subsequent stage, either a various leveled agglomerative Clustering (HAC) conspire or a partitional branching calculation—'K-medians' is used on these sub-group agents to acquire a necessary number of groups. In this manner, this mixture conspire is versatile and henceforth would be appropriate for Clustering enormous data and we additionally get a various leveled structure comprising of branches and subbranches and the delegates of which are used for architecture integration. Regardless of whether an increasing number of models are created, order time doesn't expand much as just a piece of the progressive structure is looked at. The test results (Clustering Accuracy (CA) utilizing the models acquired and the calculation time) of the proposed calculations are contrasted and that of the progressive agglomerative plans, K-medians, and nearest neighbor classifier (NNC) techniques.

Clustering is a functioning examination subject in design acknowledgment, data mining, measurements, and AI with various accentuations. We use branching as a device for model choice for design Clustering. It is material for both labelled and unlabeled data as the names are not used while clustering the examples dependent on similarities/dissimilarities measures. The prior Clustering approaches don't satisfactorily consider the way that the informational collection can be excessively enormous and may not fit in the fundamental memory of certain PCs. It is important to inspect the rule of branching to devise productive calculations to limit the I/O tasks and space necessities and to get proper models/reflections to expand the Clustering Accuracy(CA). One such application territory where productive Clustering procedures are required is in bioinformatics

## MiniBatch K-means cluster

Here Small data batches of fixed databases and size are used. The Same iterations for any sample of data are recovered and used in updating in clusters and repeated till needed.

The Method is used in various Machine Learning and data mining approaches due to its simple and efficient technique. It's Straightforward Parallelizability and low computational costs are the key factors in it. (Capó, M., Pérez, A., & Lozano, J. A. (2020) The Term K suggests that it is relatively slow as it requires to contain millions of data into thousands or tens of thousands of data clusters. But these algorithms are improved and designed to perform faster by initiating certain multistage filtering methods. There are certain advantages to it and disadvantages as well. The method is very efficient from the time needed perspective. The Data clusters or sets as mentioned earlier increases gradually in size. For analyzing them the whole main memory of the system faces severe space constraint, hence the algorithm method loses its priority. The Spatial and Temporal cost of the algorithm comes into question in that case. (Hu, Q., Wu, J., Bai, L., Zhang, Y., & Cheng, J, 2017)

The exemplary K-means calculation (Lloyd's calculation) comprises two stages. For a contribution of information purposes of measurements and introductory branch habitats, the task step appoints each point to its nearest group, and the update step refreshes every one of the branch places with the centroid of the focuses allotted to that group. The calculation rehashes until all the group habitats stay unaltered in an iteration. Because of its straightforwardness and general appropriateness, the calculation is one of the most broadly used branching calculations, and is distinguished as one of the top ten data mining calculations. The calculation runs gets delayed because of its less complexity which in turn  have various endeavors attempting to improve its speed (Yufei et al, 2015)

## Birch Clustering

The Term Balanced Iterative Reducing and Clustering using Hierarchies suggests the usage of using it in studies extensively based on its performance in the field of Memory Requirements, run time, quality of clustering, stability and scaling. Zhang, T., Ramakrishnan, R., & Livny, M. (1997)

The Features of this Algorithm are quite simple as it is quite compact there is no need to store individual factors belonging to a cluster. It can cluster a point without having to check against all other data points or clusters. It has served successfully in creating good clusters with a single scan of the data set. It helps to minimize the running time and the resulting quality is adjusted with regards to the available memory. It also helps cancel out outside nonrelevant factors. BIRCH even gives out a warning before the results if certain constraints of the data is not satisfied before clustering the data. Besides, there are no such disadvantages to it. (Lorbeer, B., Kosareva, A.,Deva, B., Softić, D., Ruppel, P., & Küpper, A., 2018)

Clustering calculations are as of late recovering consideration with accessibility of enormous dataentries and the ascent of simultaneous figuring structures. Notwithstanding, most clustering calculations experience the negative effects of two downsides: they don't dwell along with expanding data sizes and regularly asks for legitimate boundary scale which is generally hard to give. A significant model is the group tally, a parameter that by and large is close to difficult to survey. This methodology registers the ideal limit parameter of BIRCH from the information to such an extent that BIRCH does appropriate clustering even without the worldwide clustering stage that is normally the last advance of BIRCH. This is conceivable if the information fulfills certain requirements. If those limitations are not fulfilled, A-BIRCH will give an appropriate admonition before introducing the outcomes. This methodology renders the last worldwide clustering venture of BIRCH pointless much of the time, which brings about two favorable circumstances. In the first place, Pitoli et al, 2017  does not have to know the normal number of branches previously. Second, without the computationally expensive last clustering, the quick BIRCH calculation will turn out to be considerably quicker. For huge informational indexes, Pitoli et al, 2017 presents another variety of BIRCH, which they call MBD-BIRCH, which is of a specific bit of leeway related to A-BIRCH yet is free from it and of general advantage. (Pitoli et al, 2017)

Clustering is an unaided learning strategy that branches a lot of given information focuses on very much isolated subsets. Two unmistakable instances of clustering calculations are k-means, and the desire expansion (EM) calculation, this thesis tends to two issues with clustering: (1) clustering calculations, as a rule, don't run well and (2) most calculations caters the number of groups (branch tally) as information. The primary issue is turning out to be increasingly significant. For applications that need to branch, for instance, a great many records, tremendous picture or video databases, or terabytes of sensor information created by the Internet of Things, adaptability is fundamental. The subsequent issue seriously decreases the relevance of clustering in circumstances where the branch check is exceptionally hard to anticipate, for example, for information investigation, include building, and report clustering.

A significant clustering technique is adjusted iterative diminishing and clustering utilizing chains of importance, or BIRCH, which was presented and is one of the quickest clustering calculations accessible. It beats a large portion of the other clustering calculations by up to two sets of size. In this way, BIRCH as of now illuminates the main issue referenced previously. Be that as it may, to accomplish adequate clustering quality, BIRCH requires the branch considers input, in this manner neglecting to understand the subsequent issue. This thesis depicts a

strategy to utilize BIRCH without giving the branch tally yet protecting group quality and speed. This is accomplished as follows. We first expel the worldwide clustering step that is done toward the finish of BIRCH, since this is moderate and regularly requires additional parameters like for example the group consider input. At that point, by investigating the rest of the piece of the BIRCH calculation, which we call tree-BIRCH, we distinguish three manners by which tree-BIRCH can turn out badly: branch parting, group consolidating, and supercluster parting. This information at that point empowers us to improve tree- BIRCH and register an ideal limit parameter from the information. With the subsequent calculation, the client needs to give neither any parameters to the last clustering step like for example the branch check, since there is no last clustering step any longer, nor the edge parameter, since this is registered consequently.

The edge parameter is registered from two properties of the information, the most extreme branch range and the base distance between groups. For unusual circumstances in which this isn't the situation, we portray one potential technique of acquiring them. Following a thought in Bach and Jordan, we propose to take in these properties from agent information. Above, supercluster parting was referenced as one of the issues with tree-BIRCH. This drove us to devise another expansion of BIRCH, MBD-BIRCH, which is extensively diminishing supercluster parting, however to the detriment of speed. Notwithstanding, MBD-BIRCH, as well, is still a lot quicker than a large portion of the other clustering calculations. Likewise, tree-BIRCH can be used as an online calculation. Full-BIRCH can't be used on the web, since it needs an end at which the worldwide clustering could be run, yet online calculations never end. Be that as it may, tree-BIRCH ordinarily experiences poor clustering quality. In this way, we center on improving the clustering nature of tree-BIRCH.

BIRCH uses three parameters: the stretching factor Br, the limit T, and the branch check k. While the information focuses are gone into BIRCH, a stretch adjusted tree, the branch highlights tree, or CF tree, of progressive groups, is constructed. Every hub speaks to a branch in the group order, halfway hubs are super groups and the leaf hubs are the genuine branches. The fanning factor Br is the greatest number of youngsters a hub can have. This is a worldwide parameter. Each hub contains the most significant data of having a place group, the branch highlights (CF). From those, the group habitats, where are the components of the branch, and the branch radii can be figured for each branch. Each new point begins at the root and continues to stroll down the tree, continually affecting the sub-cluster with the closest focus until the walk closes at a specific leaf hub.

Once showed up at a leaf, the new point is added to this leaf branch if this would not build the sweep of the group past the edge T. In any case, another group is made with the new point as its lone part. Along these lines, the edge parameter controls the size of the branch

## Expanding Factors of Birch

The most extreme number of CF sub-groups in every hub. If another information case shows up to such an extent that the quantity of sub-groups outperforms the spreading factor, then that hub should partition into two hubs with the sub-branches redistributed in each. The parent sub- branch of that hub ought to be withdrawn and two unique

sub-groups are included as guardians of the two split hubs. It sets the default estimation of this parameter to 50 (Kulkarni, 2019).

## Agglomerative Clustering

Agglomerative or hierarchical Clustering has gained importance in the world due to the rapid and exponential rise in data across the world. Very often this information is unlabeled and there is minimal prior area information accessible. One important huddle in taking care of these immense information assortments is the computational cost. There have been plans to improve the effectiveness by presenting a lot of techniques for agglomerative various leveled grouping. Rather than building a data cluster dependent on raw information, some ways or techniques construct a progressive system dependent on the clustering of centroids. These centroids speak to clustering of nearby focuses in the information space,

This algorithm is very easy to understand, because of its relatively straightforward methodology. At the same time, it is embedded with certain disadvantages

- The Best Results are rarely provided
    - o Certain arbitrary decisions are made following this methodology
    - o The Dendrogram is Misinterpreted (Yu, Liu, X., Zhou, X., & Song, (2015))

Clustering is a significant method for information examination in true situations since manual labeling of the information is typically costly. Besides earlier information required to encourage manual labeling is regularly inaccessible or inadequate. Under such conditions clustering is a progressively reasonable choice over administered learning draws near, for example, characterization and relapse. Productive methods for information clustering have been read for a considerable length of time because of its noteworthy ramifications in certifiable applications where the measure of information is rarely enormous, and the collection of information is regularly quickening. A clustering strategy that requires less computational expense can be advantageous by and large information mining and information disclosure, just as in explicit spaces for example bioinformatics, web utilization checking and interpersonal organization examination. Owing to the influence of web applications, cell phones, and system of sensors, the volume of information to be broke down develops a lot quicker than computational force, particularly lately. This surge of information focuses on proficiency in creating clustering techniques.

The proficiency issue of various leveled clustering can be considered to be standard clustering strategies as it is commonly appropriate to most kinds of information. In correlation with partitional clustering calculations, for example, K-means, various leveled approaches have a greater expense, with an intricacy, they do not require any predicated parameter subsequently are progressively reasonable for taking care of genuine information were finding an [appropriate arrangement of parameters can be dubious (Bouguettaya et al., 2021).

Progressive clustering can go the two different ways, accumulating from singular focuses to the most significant level branch or isolating from a top group to nuclear information objects. Our center is the base up approach which is called the agglomerative methodology because computers and Information Technology expense can be

diminished if the base up process begins from someplace in the progressive system and the lower some portion of the chain of importance is worked by a more affordable technique, for example, partitional clustering. This would not function admirably on the top-down methodology known as disruptive various leveled clustering since it is infamous for its significant expense and confirming center level sub-branches by singular information focuses would, in any case, be costly.

It is conceivable to utilize a various leveled way to deal with create center level sub-branches at that point apply partitional calculations on these sub-groups. Anyway, predicated parameters like K despite everything should be resolved. Another conceivable method to improve effectiveness in various leveled clustering is to perform highlight extraction or choice, which may decrease information dimensionality. Anyway, that procedure frequently requires space information on information. It likewise makes the clustering results reliant on the presentation of the component extraction or choice calculations (Bouguettaya et al., 2021).

A progressive calculation yields a Dendrogram speaking to the settled gathering of examples and likeness levels at which groupings change. The clustering procedure is performed by combining the most comparable examples in the group set to frame a greater one. Examined the diverse various leveled clustering calculations. Various leveled clustering approaches produce groups of a higher caliber. Be that as it may, these methodologies experience the ill effects of high time cost. The effectiveness of various leveled calculations can be improved with the help of list structures. BIRCH calculation receives the idea of clustering highlights to catch the data of a group. The branches that are constructed so far by the calculation are sorted out into the CF tree. The leaf hub of the CF tree is a sub-group rather than a solitary information point.

As of late, the transformative calculation has been brought into clustering. As a sort of stochastic pursuit strategy, it can regularly be very powerful in finding ideal arrangements. Anyway, the proficient angle is somewhat an issue as a transformative procedure is tedious. To accelerate a various leveled agglomerative clustering process, GPU can unquestionably be used. This examination doesn't include GPU even though the proposed strategy can incorporate GPU to additionally upgrade the execution time. (Bouguettaya et al., 2021)

## Spectral Co Clustering

The objective of co-Clustering is to cluster or arrange the data in rows and columns from the matrix for data input. It has an advantage as it defeats a few constraints related to Traditional clustering techniques by permitting automated discovery in the similarity of the subsets based on the attributes. That is why there are no sufficient grounds to validate the data. It very often outperforms traditional clustering algorithms like the k-means algorithm (Huang, S., Wang, H., Li, D., Yang, Y., & Li, T. (2015).

Spectral Co clustering and co-clustering are notable procedures in information analysis, and late work has

stretched out absurd clustering to square, symmetric tensors and hyper matrices got from a system. The authors builds up another tensor co-clustering strategy that all the while groups the lines, sections, and cuts of a nonnegative three-mode tensor and sums up to tensors with any number of modes. The calculation depends on another arbitrary walk model which the authors call the super-spacey irregular surfer. The authors show that our technique out-performs cutting-edge co-clustering strategies on a few engineered datasets with ground truth groups and then utilizes the calculation to break down a few genuine world datasets. (Huang, S., Wang, H., Li, D., Yang, Y., & Li, T. (2015).

Clustering is a principal task in AI that intends to dole out firmly related substances to a similar gathering. Conventional techniques streamline some total proportion of the quality of pairwise connections (e.g., similitudes) between things. Ghostly clustering is an especially incredible procedure for processing the branches when the pairwise similitudes are encoded into the nearness lattice of a chart. In any case, many charts like datasets are more normally depicted by higher-request associations among a few elements. For example, multilayer or multiplex systems portray the cooperation between a few charts all the while with hub layer connections. Non-negative tensors are a typical portrayal for a considerable lot of these higher-request datasets. For example, the entry in a third-request tensor may speak to the comparability between things and in the layer. Here the authors build up the General Tensor Spectral Co-clustering (GTSC) system for clustering tensor information. The calculation takes as information a nonnegative tensor, which might be scanty, non-square, and hilter kilter, and yields subsets of files from each measurement (co-groups). Hidden the technique is another stochastic procedure that models higher-request Markov chains, which the authors call a super-spacey irregular walk. This is used to sum up thoughts from phantom clustering dependent on arbitrary walks (Huang, S., Wang, H., Li, D., Yang, Y., & Li, T. (2015). Tao et al, 2016 presented a variation on the notable conductance measure from unearthly chart dividing that the authors call one-sided conductance and portray how this gives a tensor parcel quality measurement; this is likened to Chung's utilization of disseminations to frightfully segment coordinated diagrams.

One-sided conductance is the leave likelihood from a set after our new super-spacey arbitrary walk model (Tao et al, 2016). The objective of co-clustering is to at the same time branch the lines and segments of an information network. It conquers a few impediments related to conventional clustering strategies by permitting programmed revelation of likeness dependent on a subset of characteristics. Be that as it may, diverse co-clustering models normally produce particular outcomes since every calculation has its own inclination because of the advancement of various standards. In the meantime, SCCE is a network deterioration-based methodology that can be defined as a bipartite chart segment issue and unravel it productively with the chose eigenvectors. As far as Hsiao and Chang could know, this is the main work on utilizing ghastly calculation for the co-clustering group. Broad examinations on benchmark datasets exhibit the adequacy of the proposed technique. Our investigation additionally shows that SCCE has some great benefits contrasted and many best-in-class strategies.

Various methodologies have additionally been used as base models to deal with idea float. Contrasted and

customary clustering strategies, co-clustering has a bit of leeway in finding the concealed structure of datasets and foreseeing the missing datasets by utilizing the connection between two elements. (Hsiao, Chang, 2008)

It has been demonstrated that the presentation of co-clustering techniques can be improved by exploiting group learning. SCCE performs outfit errands on push branches (push labeling) and section groups (segment labeling) of an informational collection at the same time, targeting getting an enhanced co-clustering result. All the more explicitly, the baseline branches and segment groups are displayed as two vertices of a bipartite diagram. As needs are, the bipartite chart segment issue can be settled by discovering the least cut vertex segments in the bipartite diagram between push labeling and section labeling the total ongoing years, a few unearthly methodologies are proposed to improve the presentation of clustering. There are likewise a few deals with stretching out otherworldly ways to deal with co-clustering. Otherworldly methodologies change the issue of co-clustering as a parcel issue on a bipartite chart. Contrasted and previous calculations particularly when managing the issue of archive word co-clustering, ghostly methodologies consistently produce better outcomes. In particular, the records and words are displayed as two vertices of a bipartite chart. At that point, the phantom co-clustering calculations are used to limit the edge loads of the vertices in various subgraphs by tackling an eigenvalue framework.

Information Theoretic Co Clustering (ITCC) utilises the information table to create a probability distribution of two random variables before clustering (Huang, Wang, Li, Yan Yang, Li., 2015). ITCC expands the common data between the branched irregular factors and entwines both line and section clustering at all stages. ITCC performs push clustering by surveying the closeness of each column circulation. The segment clustering is performed comparatively, and this procedure is iterated till it meets a nearby least. BCC is a partitional co-clustering definition that is driven by the quest for a decent grid estimate. The investigation of the BCC prompts the base Bregman data standard and is ensured to accomplish nearby optimality. (Huang, Wang, Li, Yan Yang, Li., 2015) In the creators demonstrate that the examination dependent on this standard produces an exquisite meta calculation, uncommon instances of which incorporate most recently realized interchange minimization-based clustering calculations. SCC is defined as a bipartite diagram parcel issue and can be settled by limiting the edge loads of the vertices in various subgraphs with the chose eigenvectors. Every one of these calculations has brought upgrades contrasted and past methodologies. In any case, Shah et al 2008 could see that even though the three calculations can meet to a nearby least individually, they may arrive at various ones. (Shah et al, 2008)

Utilizing clustering techniques to take care of the clustering issue isn't new. Diverse clustering calculations are used to produce base branches. At that point these base branches are joined by agreement work and the conclusive outcome is used to supplant the highlights of unique information. Be that as it may, most customary clustering strategies just work on push labeling.

A Dirichlet Process-based Co-Clustering Ensemble model (DPCCE) was proposed. DPCCE gives a Dirichlet procedure earlier over the information network allotments. In detail, DPCCE loosens up the standard co-clustering

presumption that line branches and segment groups are free, giving an approach to show setting explicit autonomy of line and section branches. The creators determine free Dirichlet process priors for the line and segment branches with the goal that the quantities of the groups are unbounded from the earlier. The genuine quantities of groups can be gained from perceptions.

Subsequently, the co-groups are not limited to the customary lattice parcel, however, structure settled segments with the base co-clustering. Be that as it may, DPCCE can't perform outfit undertakings on push groups and segment branches of a dataset at the same time.

A relational multi-manifold co-clustering ensemble (RMCCE) was proposed. RMCCE is an asymmetric nonnegative network tri-factorization-based approach and can utilize complex troupe figuring out how to improve the presentation of co-clustering. In any case, not at all like the current lattice factorization based co-clustering calculations, there is a complex coefficient vector that must be advanced in RMCCE, which represents a difficult undertaking. Besides, RMCCE can't utilize the various base co-clustering, i.e., RMCCE can't exploit different co-clustering calculations to get preferred prescient execution over could be acquired from any of the calculations.

SCCE makes up the weaknesses of these outfit approaches by performing clustering undertakings on both line labeling and section labeling all the while and in the meantime, in light of the perception on the three co-clustering calculations, an idea about the suspicion that the three calculations might be consolidated together with the goal that the authors can exploit the datasets of every one of them and arrive at a superior least which is nearer to the worldwide one. (Huang, Wang, Li, Yan Yang, Li., 2015)

## DB Scan Clustering

The integration of machine learning techniques highly interesting due to its implications for limitless applications. (García, Moraga, Valenzuela, Crawford, 2019) the DB-scan requires no supervision for learning and is used with the sole intent of using it in the binary process of continuous swarm data analysis algorithms. The DB-scan operator contributes to the binary process is analyzed systematically through the specifically designed random operators.

These Algorithm is also equipped with certain advantages and disadvantages.

### Advantages: -

- It can be used to detect arbitrary clusters
- It Requires two points to function which are completely independent of the order sequence
- It is very robust towards outlier (irrelevant data) detection.
- It can be used to detect clusters which are completely surrounded by other clusters

Disadvantages: -

- For multiprocessor systems, it cannot be partitioned for workload distribution
- It is very sensitive to clustering parameters
- It generally fails to identify clusters and very tricky to handle clusters if the data is too sparse or their densities vary
- Sampling affects density measures. (Dang, Shilpa, 2015).

Branching analysis is an unaided learning technique that isolates the information focuses into a few explicit bundles or clustering, with the end goal that the information focuses on similar clustering have comparative properties and information focuses in various clustering have various properties in some sense.

It includes various techniques dependent on various separation measures. For example, K-Means (separation between focuses), Affinity spread (chart separation), Mean-move (separation between focuses), DBSCAN (separation between closest focuses), Gaussian blends (Mahalanobis separation to focuses), Spectral clustering (diagram separation), and so on.

Midway, all branching techniques utilize a similar methodology for example first the authors ascertained likenesses and afterward they use it to branch the information focuses into clustering or groups. Then they concentrated on the Density-based spatial clustering of uses with commotion (DBSCAN) branching technique.

Thickness Based Clustering alludes to solo learning strategies that distinguish particular clustering/groups in the information, in light of the possibility that a branch in information space is an adjoining locale of high point thickness, isolated from other such groups by touching areas of depressed spot thickness.

Thickness (DBSCAN) is a base calculation for thickness-based clustering. It can find branches of various shapes and sizes from a lot of information, which is containing clamour and anomalies. (M. Pawar, A. Pandey and S. Bhargav , 2018)

Why do we need DB Clustering when we do have K Mean clustering methods?

K-Means clustering may branch inexactly related perceptions together. Each perception turns into a piece of some group in the long run, regardless of whether the perceptions are dissipated far away in the vector space. Since branches rely upon the mean estimation of group components, every datum point assumes a job in framing the branch. A slight change in information focuses may influence the clustering result. This issue is extraordinarily decreased in DBSCAN because of the way branches are framed. This is normally not a major issue except if we go over some odd shape information.

## Algorithmic steps for DB scan clustering

- The calculation continues by self-assertively getting a point in the dataset (until the total of what focuses have been visited).

- If there are at any rate 'minPoint' focuses inside a sweep of 'ε' to the point, then we believe every one of these focuses to be a piece of a similar branch.

- The groups are then extended by recursively rehashing the local figuring for each neighbouring point

The outputs of the above algorithms were recorded with various skewness indices and results were discussed on the basis of their accuracy for different amount of user input anomalies.

## Parameter Estimation for DB scan algorithm

Each datum mining task has the issue of parameters. Each parameter impacts the calculation in explicit manners. For DBSCAN, the parameters ε and minPts are required.

- **minPts**: As a general guideline, a base minPts can be gotten from the quantity of measurements D in the informational index, as $minPts \geq D + 1$. The low worth $minPts = 1$ doesn't bode well, as then every point on its own will as of now be a group. With $minPts \leq 2$, the outcome will be equivalent to of progressive clustering with the single connection metric, with the dendrogram cut at stature ε. Along these lines, minPts must be picked at any rate 3. Nonetheless, bigger datasets are normally better for data with clamour and will yield increasingly noteworthy groups. As a general guideline, $minPts = 2 \cdot dim$ can be used, yet it might be important to pick bigger qualities for enormous information, for loud information or for information that contains numerous copies.

- **ε**: The incentive for ε would then be able to be picked by utilizing a k-separation diagram, plotting the separation to the $k = minPts-1$ closest neighbor requested from the biggest to the littlest worth. Great estimations of ε are the place this plot shows an "elbow": if ε is picked excessively little, an enormous piece of the information won't be grouped; though for a too high estimation of ε, branches will consolidate and most of items will be in a similar group.

Little estimations of ε are best, and as a dependable guideline, just a little portion of focuses ought to be inside this separation of one another.

- Separation work: The decision of separation work is firmly connected to the decision of ε, and majorly affects the results. All in all, it will be important to initially distinguish a sensible proportion of closeness for the informational collection, before the parameter ε can be picked. There is no estimation for this parameter, yet the separation capacities should be picked properly for the informational collection. (Chauhan, 2017)

## Spectral Bi-Clustering

Worldwide investigations of RNA articulation levels are valuable for arranging datasets and generally speaking phenotypes. Frequently these order issues are connected, and one needs to discover "marker datasets" that are differentially communicated specifically sets of "conditions." Kluger et al, 2003r has built up a strategy that all the while groups qualities and conditions, finding unmistakable "checkerboard" designs in grids of quality articulation information, in the event that they exist. In a disease setting, they relate to qualities that are extraordinarily up-or down managed in areas of sudden extrapolation. Kluger et al's, 2003, strategies, biclustering, depends on the perception that box type figures in grids of articulation information are present in eigenvectors relating to trademark articulation designs across qualities or conditions. Also, these eigenvectors can be promptly distinguished by regularly used direct polynomial math draws near, specifically the solitary worth deterioration (SVD), combined with firmly incorporated standardization steps. The author presents various variations of the methodology, contingent upon whether the standardization over qualities and conditions is done autonomously or in a joint type. Kluger et al 2003 at that point apply otherworldly biclustering to a determination of openly accessible malignant growth articulation data and analyse how much the methodology can distinguish checkerboard structures. Moreover, the author thought about the exhibition of our biclustering techniques against various sensible benchmarks (e.g., direct use of SVD or standardized slices to crude information**.** (Kluger et al, 2003)

## Uses of Spectral Biclustering

### Microarray Analysis to Classify Genes and Phenotypes

The World built up a strategy that at the same time groups qualities and conditions. The strategy depends on the accompanying two suppositions:

Two datasets that are co-regulated are relied upon to have connected articulation levels, which may be hard to see because of commotion. Alter et al, 1997 could acquire better gauges of the relationships between quality articulation profiles by averaging over various states of a similar sort.

In like manner, the articulation profiles for each two states of a similar kind are relied upon to be related, and this relationship can be better seen when arrived at the midpoint of over arrangements of datasets of comparable articulation profiles. (Alter et al, 1997)

These presumptions are upheld by straightforward investigations of an assortment of run of the mill microarray sets. For instance, introduced a dataset on five sorts of mind tumors, and afterward used a managed learning methodology to choose datasets that were profoundly connected with class qualification. They put together this work with respect to the outright articulation levels of datasets in 42 examples taken from these five kinds of tumors. Utilizing these information, Hur et al, 2002 estimated the connection be tween's the articulation levels of datasets that are exceptionally communicated in just one kind of tumor and discovered just moderate degrees of relationship. Be that as it may, if Hur et al, 2002 rather normal the articulation levels of every quality over all examples of a similar tumor type (getting vectors with five sections speaking to the midpoints of the five sorts of tumors), the segment of the datasets dependent on relationship be tween's the five-dimensional vectors are progressively clear. (Hur et al, 2002)

Hofman et al, 1999 expected that in the sections each square is consistent. The subsequent factor, signified $\rho i$, speaks to the inclination of quality  to be communicated under every test condition. The last factor, indicated $\chi j$, speaks to the general inclination of datasets to be communicated under condition j. It was accepted that the microarray articulation information to be a cumulative outcome of the result of these three components (Hofman et al, 1999).

Microarray exhibit tests for all the while estimating RNA articulation levels of thousands of datasets are getting generally used in genomic explore. They have tremendous guarantee in such territories as uncovering capacity of datasets in different cell populaces, tumor arrangement, medicate target recognizable proof, understanding cell pathways, and forecast of result to treatment. A significant utilization of small-scale exhibit innovation is quality articulation profiling to foresee result in various tumor types. In a bioinformatics setting, Getz et al, 2000 could apply different information mining strategies to malignancy datasets so as to distinguish class differentiation datasets and to order tumors. A halfway rundown of techniques incorporates: (1) information pre preparing (foundation disposal, distinguishing proof of differentially communicated datasets and standardization); (2) unaided clustering and representation strategies (various levelled, SOM, k-means, and SVD); managed AI strategies for order dependent on earlier information (discriminant investigation, support-vector machines, choice trees, neural systems, and k-closest neighbors); and progressively aggressive hereditary system models (requiring a lot of information) that are intended to find natural pathways utilizing such methodologies as pairwise collaborations, ceaseless or Boolean systems (in view of an

arrangement of coupled differential conditions), and probabilistic chart displaying dependent on Bayesian systems .

Our emphasis here is on solo clustering techniques. Unaided procedures are helpful when marks are inaccessible. Models incorporate endeavors to distinguish sub classes of tumors, or work on recognizing branches of datasets that are co controlled or share a similar capacity. Solo strategies have been effective in isolating specific kinds of tumors related with various sorts of leukemia and lymphoma. Be that as it may, solo (and even directed) strategies have had less achievement in apportioning the examples as indicated by tumor type or result in illnesses with different sub orders. Likewise, the strategies Getz et al, 2000 proposed here are identified with a technique for co clustering of words and archives. (Getz et al, 2000)

### Previous application of Spectral biclustering

The possibility of synchronous clustering of lines and sections of a network returns. Techniques for concurrent clustering of datasets and conditions were all the more as of late proposed. The objective was to discover homogeneous submatrices or stable branches that are important for organic procedures. These strategies apply eager iterative inquiry to discover fascinating examples with regards to the lattices, a methodology that is additionally regular in uneven clustering. Interestingly, our methodology is progressively "worldwide," discovering bi-clusters utilizing all segments and lines.

## Partitional Clustering

As opposed to the progressive clustering, a partitional clustering calculation acquires a level parcel of the dataset which advances a predefined model capacity. The most generally used partition clustering calculation is K-means clustering, which over and over allots each item to its nearest group focus and registers the new branch habitats as needs be until the predefined model is met. In view of how the separation between information focuses is figured, different partitional clustering calculations have been created and delegate ones incorporate unearthly clustering, diagram dividing based non-negative framework factorization based methodologies. Contrasting and K-means clustering, these calculations ordinarily produce branches of better quality. (Lee et al., 2007)

Nonetheless, these calculations are all the more computationally included, requiring performing eigen disintegration or dull framework augmentation, making them not adaptable to exceptionally enormous datasets. Blend model or other thickness based clustering calculations yield delicate branch participations, permitting every datum point to be related with various groups with various probabilities. Contrasted and the proposed approach and progressive clustering when all is said in done, partitional clustering calculations endure two significant impediments. To begin with, their exhibition intensely depends on pre-characterized parameters, particularly the quantity of branches, so the nature of information groups cannot be ensured. Second, the resultant branches have a level structure rather than various levelled structure that caught a lot more extravagant relationship among information focuses. A various levelled structure offers a progressively characteristic approach to sort out some genuine articles (e.g., records and website pages) and encourage human clients to peruse the information (Bouguettaya, 1996).

## Hybrid Clustering

Hybrid information clustering consolidates the progressive and partitional strategies to acquire the great nature of the previous and the effectiveness of the last mentioned. Distinctive mixture information clustering calculations have been proposed. A hybrid clustering calculation called CURE was proposed to viably recognize the subjectively molded branches. Given an enormous dataset, CURE draws a lot of information tests from the entire dataset by arbitrary inspecting. The information tests are gathered as a few parcels and those in each segment are mostly branched. The exceptions are then expelled from the dataset. The last branches are gotten by further clustering over the fractional groups delivered in the past advance. Fix is adaptable to huge datasets with a straight time unpredictability. (Lin et al., 2005)

Be that as it may, not the same as the proposed approach, it despite everything requires the client indicated parameter esteems including the quantity of groups and the contracting factor, which may

influence the nature of branches. A Cohesion-based Self-Merging (CSM) clustering calculation is proposed. CSM receives another closeness measure, alluded to as union, to compute the separation between groups. Attachment is characterized dependent on the consolidating tendency of two groups as indicated by the presence of a mutual information point. Since the combining tendency ought not to be dictated by just a couple of focuses, union all things considered considers all the information focuses in the two groups to be blended. This makes the union measure powerful to the presence of anomalies. By utilizing attachment, CSM adequately consolidates the highlights of partition and progressive clustering techniques. In the main stage, it segments the first information space into little groups utilizing k-means. At that point the got little groups are combined utilizing the union similitude measure in a progressive way. Trial examines show that CSM exceeds expectations at both clustering precision and execution time. Nonetheless, since CSM expects clients to determine the quantity of sub-branches expected in the information segment stage and the quantity of conclusive groups, it infers that reasonable parameter esteems should be provided by the space information on a specific database. (Xu and Wunsch, 2010). Interestingly, the proposed approach is area autonomous. The Self-Partition and Self-Merging (SPSM) calculation additionally attempted to lessen the impact of client determined parameters by utilizing a recursive information segment handling. The sub- branches are created by recursively isolating the dataset into four segments. Be that as it may, it despite everything experiences the negative impact of isolating a solitary branch into various parts or clustering information pointers of various classifications into one group. This is brought about by the pre-characterized number (i.e., four) of sub-groups for each apportioning step. Researchers proposed to join a troublesome procedure and an agglomerative system. Both the disruptive and agglomerative segments utilize a sort of various levelled calculations, in this manner this technique doesn't exploit the advantages from partition clustering as in the proposed approach. (Luxburg, 2007)

## Important Applications of Data Clustering

The primary motivation behind utilizing information clustering methods is to improve the exhibition of information access by summing up the information objects into clusterings. Regularly a gathering of clustering techniques or a mix of clustering with different strategies functions admirably. it is proposed to incorporate clustering with a social DBMS for permitting K-means calculation to branch huge data inside a social database the executives framework. Dissimilar to the standard K-means approach that deals with the information and clustering brings about memory, all the information are put away in plate. The presentation of clustering is improved by measurements-based introduction of centroids and accomplished quick intermingling. A parcel and-clustering structure was proposed to find basic sub-directions from a direction database utilizing direction based clustering which is joined with a proposed locale based clustering. This district-based clustering finds the areas having directions of one significant sort. The direction-based clustering misuses the move

examples of directions dependent on their low-level highlights, while the area-based clustering uses progressively broad highlights without considering specific move designs. Both the effectiveness and precision could be improved because of the cooperation between these two diverse clustering strategies. Group avoiding rearranged record is proposed dependent on the partition clustering for effective recovery of archives. Other than the general record data, the branch participation and centroid data are put away in the altered document also. Clustering methods have been applied to examination of miniaturized scale exhibit informational indexes. For instance, Pan et al. proposed to utilize inspecting based network deterioration for quick co-clustering of miniaturized scale exhibit information. It was proposed to distinguish co-controlled quality groups by another tree-based clustering calculation. Subspace based clustering has been proposed to defeat the dimensionality revile of high dimensional information. A separation-based clustering model called n Cluster was proposed to recognize the noteworthy groups by an adaptable measurement parcel approach which permits the covering between various receptacles of a characteristic.

Bouguettaya et al. (2015) proposed an optimal visual dictionary which utilizes subspace-based clustering to retrieve videos efficiently. The procedure included obtaining high quality clusters by determining an optimal subspace combination. This optimal subspace combination should have the maximal discrimination power and only then recursive k-means algorithms are performed over each such dominant subspace. This enables the preservation of high accuracy of the video.

## Applications and Challenges:

There are yet huge difficulties for completely using AI and acknowledging self-driving in vehicles. Specifically, the significant difficulties originate from the security basic and time-basic nature of vehicle frameworks. In the first place, while ML methods, for example, profound learning may give great execution much of the time, it is difficult to reason about their most pessimistic scenario function and strength with less affirmation of the vigor and strength of Machine Learning methods, the primary security of vehicles cannot be ensured. Difficult machine learning techniques are tedious and take up larger runtimes. Current self-driven vehicle models uses costly processors and ALU, together with costly sensors. The commercialization of automated vehicles will require substantially more affordable and cheaper but efficient technology (Salahuddin, Fuqaha and Guizani, 2016).

## Motivation of Approach

Generating synthetic vehicular data and analyzing the same to determine any abnormal behavior in the regular functioning of the vehicle using various unsupervised learning algorithms to further address presence of security concerns in automobiles. The datasets with acceleration, gyroscope and positions values are deemed to be important as mentioned by Miller and Valasek (2014) where an intrusion in a JEEP vehicle caused severe complications, the intrusion was powerful enough to cause crippling of brakes, could take over the control of throttle and had the access to even disable transmission. In such cases of intrusion, the vehicle is supposed to act abnormally, and the behavioral pattern of the vehicle will give us accelerations, gyroscopic and positional values which cannot be considered normal unless the intrusion / attacker chooses to just intrude and not interfere with the normal functioning of the vehicle which in most cases is highly unlikely. All the mentioned intrusion detection strategies might be used just for specific attacks that have been already considered in configuration stages (Sun et al, 2015). To adapt to such condition machine learning based ids methods were used, for the most part, for conventional communication systems. (Deepa, Kavitha, 2012). The motive is to catch fundamental factual and statistical features of information and use them to distinguish any harmful intrusion (Tsaia et al, 2009). Intrusion Detection System utilises artificial neural network (Golovko, Kochurko, 2005) and support vector machine (Hu et al, 2003) are produced for characterizing attack types. The propelled machine learning techniques are seldom used for a vehicular network because the computing intensity and power of the regular ECU is restricted to process the further complex procedure. Be that as it may, the computing intensity of ECU has expanded in recent times to process colossal real time calculations in the modern vehicular framework. (Aurngren, Nielsen, 2005)

## Aims and Objectives

**Aim**: Generating synthetic vehicular data to represent the impact of a cyberattack, before analyzing for the purpose of identifying any abnormal behavior in the regular functioning of the vehicle using various unsupervised learning.

**Objectives**: -

1. Perform an in-depth literature review to determine the state-of-the-art research conducted in the related field.

2. Define and follow a methodology to generating synthetic data from the attained sample data.

3. Use different Unsupervised Machine Learning techniques to identify anomalies in the generated dataset.

4. Observe and analyze results for the various algorithms.

## Methodology



COLLECTION OF DATA

(from AEGIS , mainly consists of accelerations, positions and gyroscopic data)

GENERATION OF SYNTHETIC DATA
(From the sample data obtained earlier)

PRE-PROCESSING THE GENERATED DATA

(Using timestamp as the primary key, combining three data sets into one with only the resultant of each of the three sample datasets)

NORMALISATION OF THE RESULTANT DATASET

USING A GENERATED SCRIPT TO PROVIDE THE DATASET WITH LABELS

(0 For non-anomaly and 1 for anomaly)

USE UNSUPERVISED MACHINE LEARNING ALGORITHMS TO DETERMINE THE MOST COMPATIBLE ALGORITHM FOR THE DATASET BASED ON THE ACCURACY PERCENTAGE
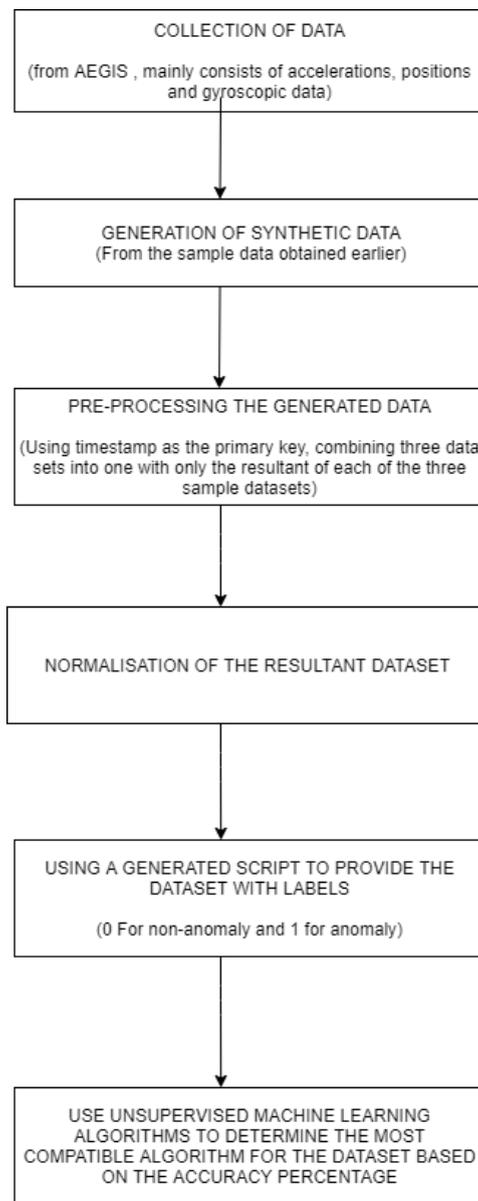
*Figure 2 Flowchart explaining the proceedings of Methodology*

Figure 2 illustrates the process undertaken in this research: data was first collected from a sample data source, (Stocker, Kaiser, Festl, 2017), which was followed by synthetic data generation where the script could take random skewness as input from the range for 1% to 100% as required to further test the accuracy of the algorithms, same script was used to generate data of the same kind but in a larger scale, this was then followed by the pre-processing which was applied on the generated dataset and a resultant dataset was obtained and it was normalized. A script was used to add labels to this dataset (0 for anomaly and 1 for non-anomaly) followed using Unsupervised Machine Learning algorithms to determine the most compatible algorithm in our case or observe the results and come

down to a ground conclusion about the usage of Unsupervised algorithms in detecting anomalies in automotive data. This script was created to take values from the previous data as inputs and linearly increase values at the end to maintain continuity of the model data, this would enable us to have sufficient data to promote smooth functioning and have sufficient result to back up the research.

Skewness taken was random because we wanted the unsupervised model to be trained in a way that any sudden jump or depreciation in values at any point of time could be taken into consideration at once and the training of the model could be carried in such environment.

Sample data was obtained from AEGIS which stands for "Advanced Big Data Value Chain for Public Safety and Personal Security" (Stocker, Kaiser, Festl, 2017). The datasets obtained were comprised of data related to:

- Acceleration – This dataset mainly had information about the vehicle in the Cartesian coordinate system along with their timestamps which are all necessary for finding out the acceleration of a moving object. Table 1 is a part of the data which shows the acceleration of the vehicle – in the given timestamp

- Gyroscope – Gyroscopic data is similar to that of the data contained in the acceleration dataset. This one basically measures the angular acceleration of a moving object in the 3-D plane. Table 2 is a part of the data which shows the gyroscope data of the vehicle – in the given timestamp

- Position – This datasheet helps us with the information about the position of an object in relation to 3D space, viz. Latitude and Longitude. Table 3 is a part of the data which shows the latitude and the longitude values of the vehicle – in the given timestamp

| acceleration_id | trip_id | x | y | z | timestamp |
|---|---|---|---|---|---|
| 2 | 3 | -0.91448 | -0.85635 | 0.00024 | 24/03/2019 - 00:00:00 |
| 3 | 3 | -0.51445 | -0.64235 | -0.55463 | 24/03/2019 - 00:00:01 |
| 4 | 3 | 0.33757 | 0.01638 | 0.95508 | 24/3/2019 - 00:00:02 |
| 5 | 3 | 0.34412 | 0.00678 | 0.96455 | 24/03/2019 - 00:00:03 |
| 6 | 3 | 0.35018 | 0.01166 | 0.96184 | 24/03/2019 - 00:00:04 |
| 7 | 3 | 0.34503 | 0.0161 | 0.95344 | 24/03/2019 - 00:00:05 |

*Table 1 - Example of a part of the Generated data of accelerations from the obtained data from AEGIS*

| gyroscope_id | trip_id | x_value | y_value | z_value | timestamp |
|---|---|---|---|---|---|
| 2 | 3 | 1.2575 | 1.50867 | -0.54458 | 24/03/2019 - 00:00:00 |
| 3 | 3 | 1.3519 | 1.39751 | -1.21998 | 24/03/2019 - 00:00:01 |
| 4 | 3 | 1.16632 | 2.25566 | -0.88404 | 24/03/2019 - 00:00:02 |
| 5 | 3 | 1.40622 | 2.1774 | -0.73994 | 24/03/2019 - 00:00:03 |
| 6 | 3 | 1.55011 | 1.47389 | -1.229 | 24/03/2019 - 00:00:04 |
| 7 | 3 | 1.08121 | 1.45494 | -1.00734 | 24/03/2019 - 00:00:05 |

*Table 3 Example of a part of the Generated data of the gyroscopic sensors obtained data from AEGIS*

| pos_id | trip_id | latitude | longitude | altitude | timestamp |
|---|---|---|---|---|---|
| 2 | 3 | 4703.787 | 1527.474 | 363.8 | 24/03/2019 - 00:00:00 |
| 3 | 3 | 4703.783 | 1527.476 | 365 | 24/03/2019 - 00:00:01 |
| 4 | 3 | 4703.785 | 1527.48 | 360 | 24/03/2019 - 00:00:02 |
| 5 | 3 | 4703.785 | 1527.473 | 360.1 | 24/03/2019 - 00:00:03 |
| 6 | 3 | 4703.783 | 1527.473 | 363.6 | 24/03/2019 - 00:00:04 |
| 7 | 3 | 4703.783 | 1527.473 | 365.2 | 24/03/2019 - 00:00:05 |

*Table 4 Example of a part of the Generated data of the positional values of the vehicle obtained data from AEGIS*

The values obtained from AEGIS because they ran simulation of test of a vehicle on track under normal conditions and the values obtained were kept for further simulations and provided as sample data to promote further similar research. The datasets were obtained separately and were merged based on timestamp values. Furthermore, a script was created which would generate data of the similar kind with varied values of skewed data. This depended on the user input, for example, if the user input was .1 then the script would generate data similar to the sample but with 10% skewness. In the following section, the details of our data generation process are presented. This script was created taking values from the previous data as inputs and linearly increasing values at the end to maintain the continuity of the model data.

## Synthetic data generation

The synthetic data generation process is where data is generated through artificial means to produce a dataset that is comparable and unidentifiably different from a real dataset. The process of synthetic data generation used in this thesis case is based on the sample dataset that was obtained from AEGIS (Advanced Big Data Value Chain for Public Safety and Personal Security). A script was created to generate similar data (similar to the sample AEGIS data), and at the same time the script could introduce anomalies as per the user input. We created a Python script which took the sample data, which was in a Comma Separated Value (CSV) file format. An input and used a custom range to add more

data to the sample data and give us another .csv file as outputs and we could use this script to generate numerous amounts of data by just tweaking the custom range in our script. Later, we introduced a section in the script that adds randomness to the generated data, which was named skewness. The script would take a skewness percentage as user input and generate a .csv file with the amount of skewness desired, giving us a long range of anomalies which we would in later. Synthetic data generation was a vital part of the research because of the lack of availability of sufficient real-world data. Generating similar data would save time in hopping from place to place to get similar data and would certainly help get past the hassles of getting permissions and recommendations to use real world value if obtainable. Synthetic data, if not the same as generated in the real world, can be really close to it and certainly depends on how one desires it to be. In our case the data that we used were acceleration, gyroscope and altitude values, once we had the sample data if was easier for us to replicate and expand the available data rather than search for the exact values and data that we need in various real-world places.

## Pre Processing

Pre-Processing was carried out on the three sample datasets. First argument being dataset of position, second argument being dataset of acceleration and the third argument being dataset of gyroscope. Keeping Timestamp as the primary key in the CSV file, a resultant dataset was generated which contained the combined result of all the three (position, acceleration and gyroscope) sample datasets.

Normalization was carried out to the above dataset was processed as Liao, Carneiro, (2016) showed the importance of the technique used in machine learning. It refers to getting data, which would be measured on different scales normally and get them to a common scale to make it easier for further calculation and processing using the min-max technique, as described in the following equation:

$$ Z_i \ = \ \frac{x_i - \min(x)}{\max(x) - min(x)} $$

Where $x_i$ is an individual value within the dataset, $x$, $z_i$ is the new normalised value, $\min(x)$ is the mimum value in the datset $x$ , and finally, $\max(x)$ is the mimum value in the datset $x$.

This equation is used for normalising. The MinMaxScaler function was used for normalisation of data after going through an analysis of various normalisation techniques (Bolstad, Irizarry, Astrand, Speed, 2003). The data is comprised of accelerations, gyroscopic and positional values of a vehicle which has to deal with values relating to the minimum and maximum ranges in which a vehicle could show

changes in behaviour hence the appropriate approach to normalise the data in hand was to go with the above-mentioned technique. After normalising the dataset, it was saved, and a model was created using pickle library on python so that it can be used in future to get predictions from the model at any time.

## Unsupervised Machine Learning

Unsupervised Machine Learning is a type of algorithm that refers to the datasets consisting of input data and does not follow the path of already labelled values or response. In terms of clustering-based algorithms, the clusters are further modelled keeping in mind that the data is similar, and it is defined by a system of metrics called Euclidian distance. The datasets close to this distance will be similar and the model learns to call it similar. Any values henceforth not complying with the Euclidian distance for this case will be categorized as abnormal (Jain, 2008).

This thesis has used this technique to further our detailed analysis in detecting anomalies. The datasets obtained from the data generation process shown above, were combined into the same dataset using a field timestamp as the primary key which would give us the acceleration sensor value, gyroscope sensor value and the position of the vehicle (latitude and longitude) at that given timestamp and so on. The entire pre-processing of the data, before including it to the model and further applying the learning techniques, required two libraries on Python, namely Math and Pandas.

K-Means Cluster model is an unsupervised model from the family of unsupervised machine learning. It makes the cluster by figuring out the feature and attribute of an object and makes a relation to determine its belonging to a category in the cluster. As we need a reference dataset with predefined anomaly for testing the accuracy of the algorithms, K Means Cluster algorithm was carried out to obtain a dataset with labelled values (0 for non-anomaly and 1 for anomaly).

Once we had the labelled data we used this dataset as an input for the various algorithms comprising of: -

- MiniBatch K-means cluster
- Agglomerative Clustering
- Birch Clustering
- SpectralCo Clustering
- SpectralBi Clustering
- DB Scan Clustering

The methods were subjected to skewed data ranging from 5% - 95% and subsequent test runs were carried out to determine the accuracy of each unsupervised algorithm for various percentages of skewness.

## Results

The results presented in Table 4 demonstrate the accuracy values based on the different percentage of skewness. Various percentage of skewness were taken, and algorithms were run on the skewed models. The six different unsupervised machine learning algorithms were tested on various levels of skewness ranging from 5% to 50 %.

| ALGORITHMS | SKEWNESS | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 5% | 10% | 15% | 20% | 30% | 35% | 40% | 50% |
| | ACCURACY OBTAINED FOR THE ABOVE SKEWNESS | | | | | | | |
| MiniBatch K-means | 99 % | 99 % | 99 % | 97 % | 97 % | 96 % | 95 % | 94 % |
| Agglomerative Clustering | 98 % | 95 % | 84 % | 85 % | 84 % | 87 % | 82 % | 80 % |
| Birch Clustering | 80 % | 83 % | 14 % | 12 % | 39 % | 37 % | 21 % | 18 % |
| Spectral CoClustering | 91 % | 91 % | 88 % | 85 % | 84 % | 81 % | 82 % | 81 % |
| Spectral BiClustering | 81 % | 17 % | 18 % | 22 % | 77 % | 23 % | 80 % | 81 % |
| DB Scan Clustering | 78 % | 63 % | 54 % | 55 % | 48 % | 47 % | 51 % | 51 % |

*Table 4 Accuracy results of all algorithms with skewness ranging from 5% to 50%*

First a test was run with skewness inputs that ranged from 5% to 50% with the sole understanding that algorithms will perform well on lower accuracies and the accuracy percentages will slowly decrease as we move up with the skewness percentages. The results we got were mostly similar to what was comprehended but there were a few algorithms that outperformed the others and others underperformed significantly at certain values of skewness.

Minibatch K-Means showed the best accuracy of all algorithms and its accuracy never once dropped below 94 % in the given skewness range. Agglomerative cluster showed normal and as expected behavior for this range and its accuracy dipped with the increase of skewness. Birch Cluster one the other hand had very low and inconsistent accuracy percentages after the first two trials of 5% and 10

% and showed no real pattern in the detection of anomalies within the range of the given skewness. Spectral Co Cluster algorithm showed consistently higher values of accuracy percentages and its output was not greatly affected by the gradual increase of skewness whereas Spectral Bi Cluster showed higher accuracy percentages towards the higher end of the skewness values but failed to show decent accuracies for skewness values of within the range if 10 % - 35 %. DB Scan cluster showed a visible decrease in accuracy within the range.

This made it necessary to use higher skewness to see if the algorithms that outperformed in the range of 5 % – 50 %, would do the same in the larger skewness range and to check if the algorithms with lower accuracy in this range would show some change in the larger range. Various percentage of skewness were taken and algorithms were run on the skewed models. The six different unsupervised machine learning algorithms were tested on various levels of skewness ranging from 60% to 95 %. Table 5 shows the results of skewness ranging from 60% - 95 %

| SKEWNESS | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 60 % | 65 % | 70 % | 75 % | 80% | 85% | 90% | 95% |
| **ALGORITHMS** | | | | | | | | |
| | ACCURACY OBTAINED FOR THE ABOVE SKEWNESS | | | | | | | |
| MiniBatch K-means | 99 % | 96 % | 93 % | 92 % | 94 % | 94 % | 98 % | 96 % |
| Agglomerative Clustering | 90 % | 87 % | 77 % | 77 % | 87 % | 86 % | 89 % | 20 % |
| Birch Clustering | 57 % | 43 % | 45 % | 45 % | 73 % | 44 % | 61 % | 14 % |
| Spectral CoClustering | 82 % | 81 % | 81 % | 81 % | 82 % | 82 % | 82 % | 82 % |
| Spectral Biclustering | 83 % | 83 % | 68 % | 15 % | 13 % | 13 % | 12 % | 13 % |
| DB Scan Clustering | 51 % | 52 % | 51 % | 51 % | 51 % | 51 % | 51 % | 49 % |

*Table 5 Accuracy results of all algorithms with skewness ranging from 60% to 95%*

The process of adding skewness extended till 95%. The sole reason being few algorithms were constantly providing higher accuracies throughout the process and it was better if the process was carried out till the highest skewness addition to monitor the behavior of those algorithms for higher

values of skewness.

Minibatch K –Means continued to outperform with significantly higher values of accuracy than all others and never dropping below 90 % in this region of skewness (60 % - 95 %) as well. Agglomerative clustering showed higher accuracy in this region till the final skewness

value of 95 %, where it showed a huge dip to 20 % from 90 % in the previous simulation. Birch Clustering started to show decrease in accuracy and the accuracy completely decreased to a very nominal value when it reached the final simulation of 95 % skewness. Spectral Co Cluster maintained its higher accuracy percentage and showed no real signs of dropping below 80 % in the entire process where as Spectral Bi Cluster showed a sudden drop from skewness range of 75 % - 95 % with a very nominal value of accuracy. DB Scan Cluster showed no real sign of reduced accuracy from the final simulations of Table 4 and it remained constant at 50 % in the entire final part of the simulation.

## Discussion

The accuracy and the ROC curves obtained provided evidence that different algorithms performed differently for varying levels of skewness. Minibatch K- means cluster algorithm outperformed all the other algorithms with the highest accuracy percentage for all given skewness inputs. It constantly maintained accuracy percentage of more than 90 %.

Spectral Co Cluster constantly maintained accuracy percentages of over 80 % making it another reliable algorithm to predict the dataset we have concerning automotive data. Algorithms like Agglomerative Cluster and Spectral Bi Cluster performed great at certain skewness range and their accuracy decreased at higher anomaly levels. DB scan cluster algorithm had higher accuracy at the lower anomaly levels and slowly decreased before being stable at 50% for the higher part of anomaly percentages.
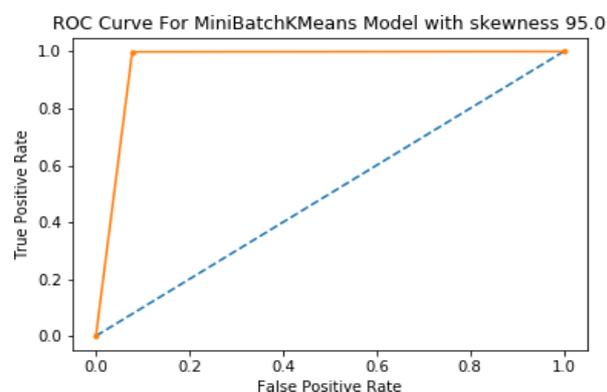


*Figure 5 ROC Curve for 95% Skewness of Minibatch K-Means algorithm*

Minibatch K- means cluster could be the best possible algorithm for our case of automotive data concerning acceleration, position and gyroscopic value. Higher values of skewness did not affect the accuracy levels for this algorithm.
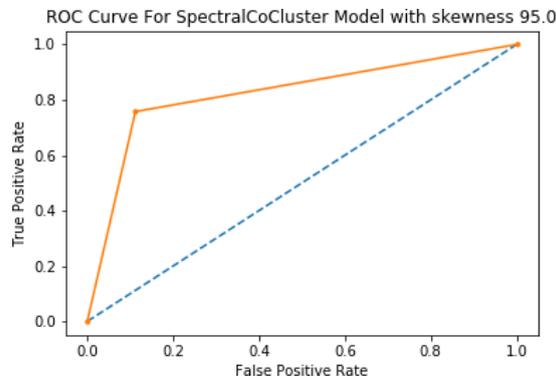


*Figure 6 ROC Curve for 95% Skewness of Spectral Co Cluster algorithm*

Spectral Co Cluster showed high values of accuracy throughout the process and maintained value of over 80 % without showing any signs of fall in accuracy and that algorithm could also be of use in our further research.
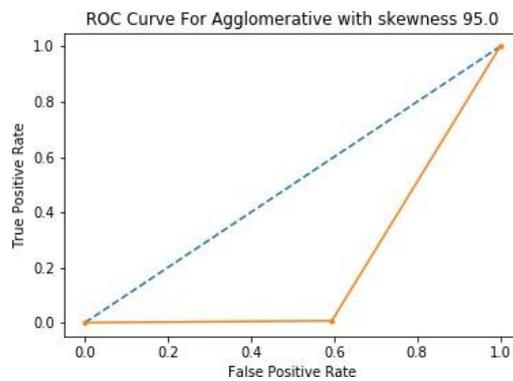


*Figure 7 ROC Curve for 95% Skewness of Agglomerative Cluster algorithm*

Other algorithms like Agglomerative Cluster showed an unusual pattern towards the end of the simulation process with a sudden dip in accuracy percentages, all other algorithms showed a decrease in accuracy either gradually or suddenly throughout the simulation process.

The higher percentage of anomaly detection and constantly getting higher accuracy percentages as visible in the Minibatch K – Means Cluster and Spectral Co Cluster algorithms can help our cause of detection of misbehavior of an automobile based the process we have followed as mentioned in this thesis.

Various aspects of unsupervised machine learning made it clear about using of certain algorithms over

others, the process of unsupervised learning is preferred because the no two drivers can be the same while driving a vehicle, so the behavioral pattern of the drivers need to be tested and recorded and actions need to performed real time. Using supervised learning will only allow us to test the previously determined data which can be very irrational in our case where a vehicle is compromised and intruded remotely which may result to a lot of anomalies in the behavioral pattern of the vehicle. Unsupervised learning can help determining the anomalies real time and determine the intensity of the intrusion real time as well. We tried using supervised algorithms and faced the problems mentioned above. Supervised Learning techniques can be best used when we have a set case and we can use them train the model and it can show very accurate results if any kind of anomaly is detected but for the Unsupervised machine learning part to detect any kind of anomaly, the model has to be trained real time and then it comes down to the percentage of accuracy for the best algorithm to suit the case. The other challenging part here was to run all the unsupervised learning techniques with fluency, DB Scan clustering algorithm and Agglomerative Clustering algorithm required more usage of RAM in completing the training of the entire data that was to be processed. We had to increase the RAM of the computer to finish training the model and detecting accuracy and use the accuracy for our research purpose.

Machine Learning is a sophisticated way to learn about data intrusion in an automotive and can be implemented in any modern-day vehicle with less complications. The program needs to constantly monitor for minor fluctuations in data and compare it with previously trained data to get to know the intensity of fluctuations. In our case if there is a sudden fall/rise of acceleration values and the values are far away from any normal conditions, it is of utmost possibility that the vehicle could be compromised or there has been a tampering with the system which in both cases the driver should know about and take proper safety measures at that point. We have also taken into consideration the gyroscopic value and constant extreme fluctuations of this value will also fall in such case of intrusion or system tampering/malfunction. The final value we considered is the location of the vehicle which is based on the GPS provided location data as a feed to keep in check the normal behavior of the vehicle system, if there is a sudden jump in the location of the vehicle from a place to another which cannot be possible in real life driving conditions or if there are constant fluctuations in a position of the vehicle the possibility of intrusion in the vehicle is of a high probability. The detection of such cases of anomalies is carried out with unsupervised machine learning algorithms and the Mini Batch K Means cluster is found out to be the best case scenario for detecting anomalies with the highest accuracy percentage values.

## Limitations

This entire process of using unsupervised machine learning to detect anomalies in vehicle data has had various limitations. The skewed data introduction in the dataset was very random and it was difficult to keep track of the skewed data without running actual tests. A small modification in the script that shows which rows and columns has skewed data would be great for further improvements that will be made to this script.

We tried using supervised algorithms and faced the problems mentioned above. Supervised Learning techniques can be best used when we have a set case and we can use them train the model and it can show very accurate results if any kind of anomaly is detected but for the Unsupervised machine learning part to detect any kind of anomaly, the model has to be trained real time and then it comes down to the percentage of accuracy for the best algorithm to suit the case.

The use of DB Scan Clustering was a bit challenging as we constantly faced a lot of errors during the simulation of large data sets.  However, this was mitigated when we partitioned the dataset into smaller sections, but it would have been best if we were could run simulations on the entire dataset at once. A larger RAM in the computer would mitigate this issue which was faced during simulation.

## Practical Significance

This technique can be used by small scale to medium scaled automotive organizations to attain safety from potential cyber-attacks. Since the technique is feasible as does not require a lot of investment upfront to get it started and running, this could be a go to method for vehicles on the cheaper end of the scale. This can further implement automation and can reduce the price of services required for Connected and Automated Vehicles and in turn reducing the overall cost of the vehicle.

## Summary of key findings

Using unsupervised learning enabled the diversity in approach and catered to a wider phenomenon of different types of driving. We found out that the vehicle if being intruded or compromised at different levels in altitude at short span of time after identifying a driving pattern, can be detected based on varied values we get as in the form of vehicular data. Similar for the case of sudden depreciation and jump in acceleration and gyroscopic values of the vehicle, which will result to the conclusion that a vehicle is either compromised or there is a possible tamper of the vehicular wellbeing.

Detecting anomalies in vehicular data using unsupervised machine learning gave us insights on which algorithm was best suited to detect anomalies depending on the accuracy rate of each algorithm. We

were successful in finding that Mini Batch K Means cluster was the best method to suit our research. High accuracy values consistently with a different user input skewness, determined the best possible algorithm and based on the table 5 and table 4 we can conclude the same.

## Conclusion

The thesis was aimed at generating synthetic vehicle data and detecting for anomalies in the same using unsupervised machine learning, this was achieved by creating a script and introducing skewness to it which we would need in the detecting phase. Substantial amount of literature was studied to find out that there are little to no literature that used unsupervised learning as a technique to detect anomalies in vehicle data, this was important for us to initiate our work that has not been done before. We selected six unsupervised machine learning techniques to run the sample and skewed data on and attained various sets of results. The results obtained were further used to determine the best unsupervised machine learning technique for our data set and our case as a whole.

Machine Learning is a sophisticated way to learn about data intrusion in an automotive and can be implemented in any modern-day vehicle with less complications. The program needs to constantly monitor for minor fluctuations in data and compare it with previously trained data to get to know the intensity of fluctuations. In our case if there is a sudden fall/rise of acceleration values and the values are far away from any normal conditions, it is of utmost possibility that the vehicle could be compromised or there has been a tampering with the system which in both cases the driver should know about and take proper safety measures at that point. We have also taken into consideration the gyroscopic value and constant extreme fluctuations of this value will also fall in such case of intrusion or system tampering/malfunction. The final value we considered is the location of the vehicle which is based on the GPS provided location data as a feed to keep in check the normal behavior of the vehicle system, if there is a sudden jump in the location of the vehicle from a place to another which cannot be possible in real life driving conditions or if there are constant fluctuations in a position of the vehicle the possibility of intrusion in the vehicle is of a high probability. The detection of such cases of anomalies is carried out with unsupervised machine learning algorithms and the Mini Batch K Means cluster is found out to be the best-case scenario for detecting anomalies with the highest accuracy percentage values.

## Future Work

This thesis initiates the usage of unsupervised machine learning as a technique to detect anomalies in vehicle data. The potential for this is huge in today's market where the world is driven on connected and automated vehicles (CAV). Slowly but steadily people CAVs will be the forefront of the automotive industries. Most vehicles today are smart vehicles with a variety of sensors and gadgets installed in them, though there are a lot of vehicles which are still not digitized and automated, but sooner rather than later this will be a necessity for the automotive industry. Detecting attacks and maintain safety of the people inside the vehicle will be a major task at hand. This being said this thesis is focused on

detecting cyber-attacks on vehicles by recognition of anomalies in vehicle data, this can be integrated to further possibilities which include using the driver behavioral pattern and implementing it to different vehicles and automating 'auto pilot' techniques on low budget vehicles. The location information attained from the data can help in better monitoring of vehicle location at the remotest of place even without a proper internet connectivity. The entire idea of this thesis is to minimize the cost of intrusion detection in vehicles so that the cheapest of vehicles can come equipped with this feature and they have a layer of cyber security service installed in the car which would prevent data theft and health damage if or when compromised.

# References

1. Guo, H., Crossman, J., Murphey, Y. also, Coleman, M. (2019). Car signal diagnostics utilizing wavelets and AI - IEEE Journals and Magazine. [online] Ieeexplore.ieee.org. Accessible at: https://ieeexplore.ieee.org/theoretical/report/892549 [Accessed 6 Oct. 2019].

2. Miller, C., and Valasek, C. (2014). Undertakings in Automotive Networks and Control Units. Recovered November 11, 2015,

3. Koscher, K., Czeskis, A., Roesner, F., Patel, S., Kohno, T., Checkoway, S. et al. (2010). Exploratory Security Analysis of a Modern Automobile

4. Moore, M., Bridges, R., Combs, F., Starr, M. also, Prowell, S. (2019). Displaying between signal appearance times for exact recognition of CAN transport signal infusion assaults: an information driven way to deal with in-vehicle interruption identification. CISRC '17 Proceedings of the twelfth Annual Conference on Cyber and Information Security Research, Article No. 11.

5. Alshammari, Abdulaziz and Zohdy, Mohamed and Debnath, Debatosh and Corser, George. (2018). Order Approach for Intrusion Detection in Vehicle Systems. Remote Architecture and Technology. 09. 79-94. 10.4236/wet.2018.94007.

6. Veeraraghvan, H., Atev, S., Bird, N., Scharter, P., and Papanikolopoulos, N. (2005). Driver Activity Monitoring through Supervised and Unsupervised Learning. IEEE. doi: 10.1109/ITSC.2005.1520169

7. Morris, B., and Trivedi, M. (2009). Solo Learning of Motion Patterns of Rear Surrounding Vehicles. IEEE. doi: 10.1109/ICVES.2009.5400238

8. Greenberg, A. (2015a). Programmers Remotely Kill a Jeep on the Highway—With Me in It. Recovered November 15, 2015,

9. Checkoway, S., McCoy, D., Kantor, B., Anderson, D., Shacham, H., Savage, S., Koscher, K., Czekis, A., Roesner, F., and Kohno, T. (2011). Exhaustive trial examinations of car assault surfaces. Procedures of USENIX Security 2011

10. Jain A.K. (2008) Data Clustering: 50 Years Beyond K-means. In: Daelemans W., Goethals B.,

Morik K. (eds) Machine Learning and Knowledge Discovery in Databases. ECML PKDD 2008. Talk Notes in Computer Science, vol 5211. Springer, Berlin, Heidelberg

11. Bolstad, B. M., Irizarry, R. A., Åstrand, M., and Speed, T. P. (2003). A correlation of standardization strategies for high thickness oligonucleotide cluster information dependent on difference and inclination. Bioinformatics, 19(2), 185-193. https://doi.org/10.1093/bioinformatics/19.2.185.

12. J. Shin, Y. Baek, Y. Eun and S. H. Son, "Intelligent sensor attack detection and identification for automotive cyber-physical systems," *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, 2017, pp. 1-8, doi: 10.1109/SSCI.2017.8280915.

13. Han, M., Kwak, B., and Kim, H. (2018). Abnormality interruption recognition technique for vehicular organizations dependent on endurance examination. Vehicular Communications, 14, 52-63. doi: 10.1016/j.vehcom.2018.09.004

    a. Thompson JA, Tan J, Greene CS. 2016. Cross-stage standardization of microarray and RNA-seq information for AI

14. Kang, Min-Joo, and Je-Won Kang. "Intrusion Detection System Using Deep Neural Network for In-Vehicle Network Security", PLoS ONE, 2016.

15. Y. Chen, Z. Lin, X. Zhao, G. Wang and Y. Gu, "Deep Learning-Based Classification of Hyperspectral Data," in *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 7, no. 6, pp. 2094-2107, June 2014, doi: 10.1109/JSTARS.2014.2329330.

16. applications. PeerJ 4:e1621 https://doi.org/10.7717/peerj.1621

17. Min-Ju Kang, Je-Won Kang. "A Novel Intrusion 73

18. Detection Method Using Deep Neural Network for In-Vehicle Network Security", 2016 IEEE 83rd Vehicular Technology Conference (VTC Spring), 2016

19. Gupta, T., and Amruthnath, N. (2018). An exploration concentrate on unaided AI calculations for early deficiency identification in prescient upkeep. IEEE. doi: 10.1109/IEA.2018.8387124

20. Farsi, M.; Ratcliff, K.; Barbosa, M.: 'A review of Controller Area Network', Computing &amp; Control Architecture Journal, 1999, 10, (3), p. 113-120, DOI: 10.1049/cce:19990304 IET Digital Library, https://advanced library.theiet.org/content/diaries/10.1049/cce_19990304

21. Stocker, Alexander, Kaiser, Christian, and Festl, Andreas. (2017). Car Sensor Data. An Example Dataset from the AEGIS Big Data Project [Data set]. Zenodo. http://doi.org/10.5281/zenodo.820576

    a.   Lee, J., Han, J., and Whang, K. (2007). Direction branching: A parcel and-gathering structure. In SIGMOD (pp. 593–604).

22. Eisen M., Spellman, P.T., Brown, P.O., and Botstein, D. 1998. Group examination and show of genome-wide articulation designs. Proc. Natl. Acad. Sci. 95: 14863-14868

23. Kemmerer RA and Vigna G. Interruption identification: A concise history and diagram. PC 2002. 10.1109/MC.2002.1012428

24. Y. Xun, J. Liu, N. Kato, Y. Fang and Y. Zhang, "Automobile Driver Fingerprinting: A New Machine Learning Based Authentication Scheme," in *IEEE Transactions on Industrial Informatics*, vol. 16, no. 2, pp. 1417-1426, Feb. 2020, doi: 10.1109/TII.2019.2946626.

25. Getz G., Levine, E., and Domany, E. 2000. Coupled two-way branching investigation of quality microarray information. Proc. Natl. Acad. Sci. 97: 12079-12084.

26. Hofmann T. furthermore, Puzicha, J., 1999. Idle class models for cooperative sifting. In Proceedings of the International Joint Conference in Artificial Intelligence. IJCAI 1999

27. Ben-Hur, A., Elisseeff, A., and Guyon, I. 2002. A strength based technique for finding structure in branched information. Pac. Symp. Biocomput

28. Alter O., Brown, P.O., and Botstein, D. 2000. Particular worth deterioration for genome-wide articulation information preparing and displaying. Proc. Natl. Acad. Sci. 97: 10101-10106

29. Park T, Han C, Lee S. Improvement of the electronic control unit for the rack-impelling cow by-wire utilizing the equipment on top of it reproduction framework. Mechatronics 2015; 15: 899–918. 10.1016/j.mechatronics.2005.05.002

30. Tuohy S, Glavin M, Hughes C, Jones E, Trivedi M, Kilmartin L. Intra-Vehicle Networks: an audit. IEEE Trans. on Intelligent Transportation Systems 2015; 2: 534–545. 10.1109/TITS.2014.2320605

31. . Biswas S, Tatchikou R, Dion F. Vehicle-to-vehicle remote correspondence conventions for upgrading roadway traffic security. IEEE Signal Processing Magazine 2006; 44: 82–97.

32. Fan Y, Dao L, Crolla DA. Incorporated vehicle elements control best in class survey. Vehicle Power and Propulsion Conference 2008.

33. Tsugawa S. Between vehicle correspondences and their applications to wise vehicles: a review. IEEE Intell. Veh. Symp. 2002.

34. Lenz H, Wagner CK, Sollacher R. Multi-anticipative vehicle following model, Eur. Phys. J.B 1999; 7 10.1007/s100510050618

35. M. Pawar, A. Pandey and S. Bhargav, "Non Convex Clustering on Datasets with missing values and noise using EMBDBSCAN & EMBOPTICS," *2018 International Conference on Computing, Power and Communication Technologies (GUCON)*, 2018, pp. 280-288, doi: 10.1109/GUCON.2018.8674914.

36. K. Vatanparvar and M. A. Al Faruque, "Self-Secured Control with Anomaly Detection and Recovery in Automotive Cyber-Physical Systems," *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2019, pp. 788-793, doi: 10.23919/DATE.2019.8714833

37. Tang T, Shi W, Shang H, Wang Y, another vehicle following model with thought of between vehicle correspondence. Nonlinear Dynamics 2014; 76: 2017–2023. 10.1007/s11071-014-1265-9

38. 8. Tang T, Shi W, Shang H, Wang Y. An all-inclusive vehicle following model with thought of the unwavering quality of between vehicle correspondence. Estimation 2014; 58: 286–293. 10.1016/j.measurement.2014.08.051

39. .Jin WL, Recker WW. Immediate data proliferation in a rush hour gridlock stream through between vehicle correspondence. Transp. Res. B 2006; 3.

40. . Kesting A, Treiber M, Helbing D. Network Statistics of Store-and-forward Intervehicle Communication. IEEE Transactions on Intelligent Transportation System 2010. 10.1109/TITS.2009.2037924

41. Yu S, Shi Z. Fuel utilizations and fumes emanations instigated by the helpful versatile voyage control methodology. Worldwide Journal of Modern Physics B 2015; 29(14)

42. Yu S, Shi Z. The impacts of vehicular hole changes with memory on traffic stream in helpful versatile journey control technique. Physica A 2015; 428(15): 206–223. 10.1016/j.physa.2015.01.064

43. Yu S, Shi Z. Elements of associated journey control frameworks considering speed changes with memory criticism. Estimation 2015; 64: 34–48. 10.1016/j.measurement.2014.12.036

44. Yu S, Shi Z. An all-inclusive vehicle following model at signalized crossing points. Physica A 2014; 407(1): 152–159. 10.1016/j.physa.2014.03.081

45. Shudong Huang, Hongjun Wang, Dingcheng Li, 5

46. Yan Yang, Tianrui Li. "Spectral co-clustering ensemble", Knowledge-Based Systems, 2015

47. . Nunen E, Kwakkernaat R, Ploeg J, Netten B. Helpful Competition for Future Mobility. IEEE Transactions on Intelligent Transportation System 2012; 13(3): 1018–1025.

48. Geiger A, Lauer M, Moosmann F, Ranft B, Rapp H, Stiller C, et al. Group Annie WAY's Entry to the 2011 Grand Cooperative Driving Challenge. IEEE Transactions on Intelligent Transportation System 2012; 13(3): 1018–1025. [

49. Lidstrom K, Sjoberg K, Holmberg U, Andersson J, Bergh F, Bjade M, et al. A Modular CACC System Integration and Design. IEEE Transactions on Intelligent Transportation System 2012; 13(3): 1008–1017 10.1109/TITS.2012.2204877

50. Farsi M, Ratcliff K, Barbosa M. A review of Controller Area Network. Processing and Control Engineering Journal 1999; 10.

51. Johansson KH, Aurngren M, Nielsen L. Vehicle uses of regulator region network Handbook of Networked and Embedded Control Systems 2005.

52. Koscher K, Czeskis A, Roesner F, Patel S, Kohno T, Checkoway S, et al. Trial security examination of an advanced car. IEEE Symposium on Security and Privacy, 2010.

53. Charlie M, Chris V. Experiences in Automotive Networks and Control Units. 2013.

54. Checkoway S, McCoy D, Kantor B, Anderson D, Shacham H, Savage S, et al. Exhaustive Experimental Analyses of Automotive Attack Surfaces Proceedings of USENIX Security 2011.

55. kleberger P, Olovsson T, Jonsson E. Security parts of the in-vehicle network in the associated vehicle. Canny Vehicles Symposium (IV) 2011.

56. Hoppe T, Kiltz S, Dittmann J. Security Threats to Automotive CAN Networks—Practical Examples and Selected Short-Term Countermeasures. Procedures of the 27th International Conference SAFECOMP 2008.

57. Larson E, Nilsson, Dennis K, Jonsson E. A way to deal with particular based assault location for in-vehicle organizations. IEEE Intelligent Vehicles Symposium 2008.

58. Muter M, Groll A, Freiling FC. Organized way to deal with oddity discovery for in-vehicle organizations. sixth International Conference on Information Assurance and Security 2010.

59. .Patsakisa C, Delliosb K, Bourochea M. Towards a conveyed secure in-vehicle correspondence engineering for current vehicles. PCs and Security 2014. 10.1016/j.cose.2013.11.003

60. Woo S, Jo HJ, Lee DH. A Practical Wireless Attack on the Connected Car and Security Protocol for In-Vehicle CAN IEEE Trans. on Intelligent Transportation Systems 2015.

61. Sun X, Yan B, Zhang X, Rong C. An Integrated Intrusion Detection Model of Cluster-Based Wireless Sensor Network. PLoS ONE 2015; 10(10) 10.1371/journal.pone.

62. Deepaa AJ, Kavitha V. A Comprehensive Survey on Approaches to Intrusion Detection System. Procedia Engineering 2012. 10.1016/j.proeng.2012.06.248

63. Tsaia C, Hsub Y, Linc C, Lin W. Interruption identification by AI: An audit. Master Systems with Applications 2009. 10.1016/j.eswa.2009.05.029

64. Golovko V, Kochurko P. Interruption Recognition Using Neural Networks. Smart Data Acquisition and Advanced Computing Systems: Technology and Applications 2005

65. Zhang Z, Li J, Manikopoulos C, Jorgenson J, Ucles JP. HIDE: a Hierarchical Network Intrusion Detection System Using Statistical Preprocessing and Neural Network Classification, IEEE Workshop on Information Assurance and Security 2001.

66. Hu W, Liao Y, Vemuri V R. Robust Anomaly Detection Using Support Vector Machines. International Conference on Machine Learning 2003.

67. Hinton G, Deng L, Yu D, Dahl G, Mohamed A, Jaitly N, et al. Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups,Communications Magazine 2012.

68. Bengio Y. Learning deep architectures for AI. Foundat. and Trends Mach. Learn. 2009.

69. Krizhevsky A, Sutskever I, Hinton GE. ImageNet Classification with Deep Convolutional Neural Networks, Advances in Neural Information Processing Systems 25 (NIPS 2012) 2012.

70. Wu J, Peng D, Li Z, Zhao L, Ling H. Network Intrusion Detection Based on a General Regression Neural Network Optimized by an Improved Artificial Immune Algorithm. PLoS ONE 2015; 10(3) 10.1371/journal.pone.0120976

71. Lv Y, Duan Y, Kang W, Li Z, Wang F. Traffic Flow Prediction With Big Data: A Deep Learning Approach. IEEE Trans. on Intelligent Transportation Systems 2015.

72. Zhang J, Wang F, Wang K, Lin W, Xu X, Chen C. Data-driven intelligent transportation systems: A survey. IEEE Trans. on Intelligent Transportation Systems 2011. 10.1109/TITS.2011.2158001

73. Dan C, Ueli M, Jonathan M, Jürgen S, Multi-column deep neural network for traffic sign classification. Neural Networks 2012.

74. Ma X, Yu H, Wang Y, Wang Y. Large-Scale Transportation Network Congestion Evolution Prediction Using Deep Learning Theory. PLoS ONE 2015; 10(3). 10.1371/journal.pone.0119044

75. Hinton GE, Osindero S, Teh Y, A fast learning algorithm for deep belief nets. Neural Computation 2006. 10.1162/neco.2006.18.7.1527

76. Chen WH, Hsu SH, Shen HP, Application of SVM and ANN for intrusion detection, Computers and Operations Research 2005.

77. Deng L, An Overview of Deep-Structured Learning for Information Processing. APSIPA 2011.

78. 48. Bengio Y, Simard P, Frasconi P, Learning long-term dependencies with gradient descent is difficult. IEEE Trans. on Neural Networks 1994. 10.1109/72.279181

79. H. Huang, W. Xia, J. Xiong, J. Yang, G. Zheng, and X. Zhu, "Unsupervisedlearning-based fast beamforming design for downlink MIMO,"IEEEAccess, vol. 7, pp. 7599–7605, 2018.

80. C. Xu, K. Wang, P. Li, R. Xia, S. Guo, and M. Guo, "Renewableenergy-aware big data analytics in geo-distributed data centers with re-inforcement learning,"IEEE Trans. Netw. Sci. Eng., to be published, doi:10.1109/TNSE.2018.2813333.

81. Y. Wang, K. Wang, H. Huang, T. Miyazaki, and S. Guo, "Traffic andcomputation co-offloading with reinforcement learning in fog computingfor industrial applications,"IEEE Trans. Ind. Informat., vol. 15, no. 2,pp. 976–986, Feb. 2019.

82. F. Tang, B. Mao, Z. M. Fadlullah, and N. Kato, "On a novel deep-learning-based intelligent partially overlapping channel assignment in SDN-IoT,"IEEE Commun. Mag., vol. 56, no. 9, pp. 80–86, Sep. 2018

83. W.-F. Hsiao, T.-M. Chang**An incremental cluster-based approach to spam filtering**

84. Expert Sys. Appl., 34 (3) (2008)

85. H. Shan, A. Banerjee, Bayesian co-clustering, in: IEEE International Conference on Data Mining, 2008

86. T. Zhang, R. Ramakrishnan, M. Livny **Birch: an efficient data clustering method for very large databases.** SIGMOD Record, 25 (2) (1996)

87. Pitolli, G., Aniello, L., Laurenza, G., Querzoni, L., & Baldoni, R. (2017). Malware family identification with BIRCH clustering. Paper presented at the 1-6. doi:10.1109/CCST.2017.8167802

88. R. Xu, D. Wunsch**Clustering algorithms in biomedical research: A review** IEEE Reviews in Biomedical Engineering, 3 (2010)

89. C.-R. Lin, K.-H. Liu, M.-S. Chen**Dual clustering: Integrating data clustering over optimization and constraint domains**

90. IEEE Transactions on Knowledge and Data Engineering, 17 (5) (2005), pp. 628-637

91. Marco Casale-Rossi, Pietro Palella, Mario Anton et al., "The World Is Going… Analog & Mixed-Signal! What about EDA?", *Proceedings of the Conference on Design Automation & Test in Europe*, pp. 37, 2014.

92. Korosh Vatanparvar and Mohammad Abdullah Al Faruque, "ACQUA: Adaptive and Cooperative Quality-Aware Control for Automotive Cyber-Physical Systems", *International Conference on Computer-Aided Design (ICCAD)*, pp. 1-8, 2017

93. amza Fawzi, Paulo Tabuada and Suhas Diggavi, "Secure Estimation and Control for Cyber-Physical Systems Under Adversarial Attacks", *IEEE Transactions on Automatic Control*, vol. 59, no. 6, pp. 1454-1467, 2014

94. Korosh Vatanparvar, Sina Faezi et al., "Extended Range Electric Vehicle with Driving Behavior Estimation in Energy Management", *IEEE Transactions on Smart Grid*, 2018.

95. Maral Amir and Tony Givargis, "Priority Neuron: A Resource-Aware Neural Network for Cyber-Physical Systems", *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 11, pp. 2732-2742, 2018.

96. Ragunathan Raj Rajkumar, Insup Lee et al., "Cyber-Physical Systems: The Next Computing Revolution", *47th Design Automation Conference (DAC)*, pp. 731-736, 2010.

97. Sangyoung Park, Younghyun Kim and Naehyuck Chang, "Hybrid Energy Storage Systems and Battery Management for Electric Vehicles", *50th Design Automation Conference (DAC)*, pp. 1-6, 2013.

98. Timothy Trippel, Ofir Weisse, Wenyuan Xu et al., "Walnut: Waging doubt on the integrity of mems accelerometers with acoustic injection attacks", *IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 3-18, 2017.

99. Armin Wasicek, Patricia Derler et al., "Aspect-oriented Modeling of Attacks in Automotive Cyber-Physical Systems", *51st Design Automation Conference (DAC)*, pp. 1-6, 2014

100. Hong Guo, J. A. Crossman, Y. L. Murphey and M. Coleman, "Automotive signal diagnostics using wavelets and machine learning," in *IEEE Transactions on Vehicular Technology*, vol. 49, no. 5, pp. 1650-1662, Sept. 2000, doi: 10.1109/25.892549.

101. Athman Bouguettaya, Qi Yu, Xumin Liu, Xiangmin Zhou, Andy Song, Efficient agglomerative hierarchical clustering, Expert Systems with Applications, Volume 42, Issue 5, 2015, Pages 2785-2797, ISSN 0957-4174,https://doi.org/10.1016/j.eswa.2014.09.054

102. K. Vatanparvar and M. A. Al Faruque, "Self-Secured Control with Anomaly Detection and Recovery in Automotive Cyber-Physical Systems," *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2019, pp. 788-793, doi: 10.23919/DATE.2019.8714833.

103.     C. Ide, F. Hadiji, L. Habel, A. Molina, T. Zaksek, M. Schreckenberg, et al., "Lte connectivity and vehicular traffic prediction based on machine learning approaches", IEEE 82nd Vehicular Technology Conference (VTC Fall), pp. 1-5, 2015.

104.     S. Lan, C. Huang, Z. Wang, H. Liang, W. Su and Q. Zhu, "Design Automation for Intelligent Automotive Systems," *2018 IEEE International Test Conference (ITC)*, 2018, pp. 1-10, doi: 10.1109/TEST.2018.8624723.