

Rash Driving Detection Through Data Analytics

Arushi Agrawal^{1*}, Isha Gupta²

^{1,2}Data Analysts, Delhi NCR, India

Corresponding Author: arushiagrawalgwl@gmail.com, Tel.:+91 7753058577

Available online at: www.ijcseonline.org

Abstract—This paper proposes an analytical solution for binary classification of driving behaviour into safe or rash. A methodology is designed and developed which consists of cleaning and scrubbing of raw driving data through identification of best set of parameters, iterative cluster-analysis, creation of target variables by benchmarking against theoretical relations and eventually performing supervised regression using support vector machine on the prepared dataset to classify a fresh driving data point into safe or rash driving. A bad and careless driving is depicted as ‘rash driving’ while a good and efficient driving is depicted as ‘safe driving’. The proposed methodology has the potential of applying to real-world driver profiling system.

Keywords— Classification, rash driving analysis, on board diagnostic, driver profiling, hierarchical clustering, principal component analysis, support vector machine

I. INTRODUCTION

With the ever-rising population, there has been a drastic increase in all over the world since last 10 years. As per statistics [1], more than 73 million passenger cars have been manufactured in the year of 2017. Such rapid increase in the number of vehicles has posed serious problems for entire human race, where most of these problems include rising air pollution, uncontrolled noise pollution and most importantly road accidents. Reckless driving raises nuisance in travel, transport and logistics. The repercussions not only affect the rider and his/her family but also the fellow passengers, drivers, pedestrians and various fleet related businesses such as cab service providers and insurance companies. The losses are of both the life and money.

etc.[2]-[3]. Rash driving is one of the major causes of millions of road accidents and deaths. This rash driving is a serious problem as it affects not only the driver meeting the accident but fellow passengers, fellow drivers, pedestrians and their family. Therefore, it is of high importance to reflect one’s driving to the person himself as well as concerned authorities for relevant action and improvement. According to the reinforcement theory proposed by BF Skinner [15], individual behaviour is a function of its consequences. It is based on “law of effect”, i.e., individual’s behaviour with positive consequences tends to be repeated, but individual’s behaviour with negative consequences tends not to be. Hence, a methodology designed with a positive or negative rewarding road safety behaviour can be effective.

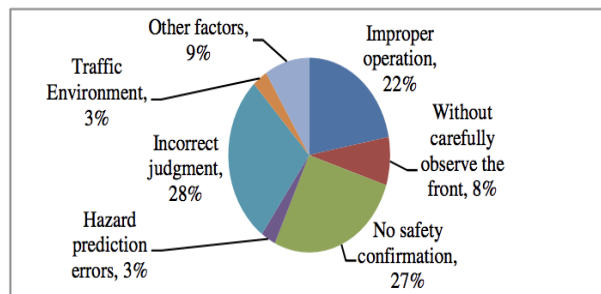


Figure 1: Distribution of traffic accidents statistics [1]

According to the statistics shown in Fig. 1, majority of accidents are caused due to individual poor driving habits. They include improper operation, incorrect judgment, and

On-Board Diagnostic (OBD) Device

OBD (On Board diagnostic) is a standard vehicle diagnostic device which gives the real time inside view of the vehicle and its electronic components. The device could access the vehicle sensors status by getting connected to OBD port. It keeps a check on vehicle engine, fuel level, battery voltage, temperature, etc. Almost all high-end recently manufactured devices already have OBD-II (advanced version on OBD) installed. This device provides detailed driving and vehicle parameters which makes it an important component in rash driving detection and driver profiling solutions. Data captured from the OBD-II port is accessible through several commercially available adopters.

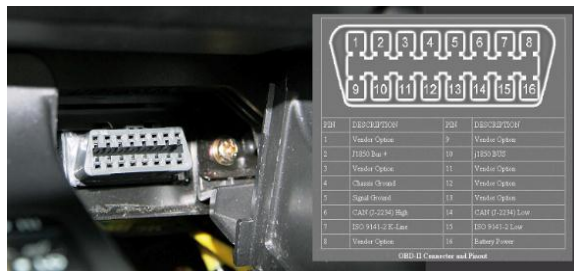


Figure 2 : On board diagnostic device connection in a car [18]

Contribution of the paper

The objective of this paper is to provide a cheap, feasible and continuous monitoring solution to curb rash driving habits. If a driving turns into rash at any moment, it can be both displayed on vehicle dashboard in real time as well as stored on a remote server for penalization. The proposed analytical solution classifies a driving behaviour into safe or rash driving on the basis of the parameters captured in OBD-II device. It is capable of generating a binary score at a constant time interval and hence, providing a real-time live alert in order to inculcate safe driving behavior. If implemented on administrative level, automatic penalization can be implemented by traffic authorities.

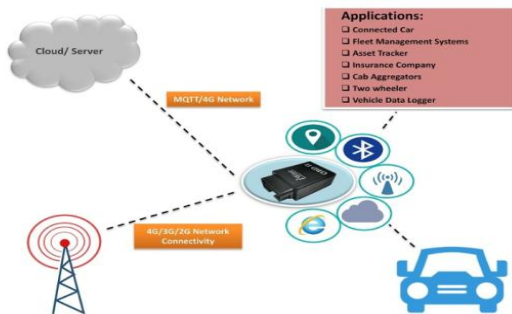


Figure 3: Connected resources with OBD device [17]

Implementation

For evaluating and assigning scores to a driver based on the parameters collected, following technique is adopted: The OBD-II device is installed in the vehicle. When the driver is driving, various parameters are captured by the OBD-II device. An Arduino UNO board is connected with the OBD device that would apply generated model on the collected data to generate a binary score (1 for rash, 0 for safe). This category will be displayed on vehicle dashboard as well as sent to a remote server along with corresponding OBD parameters, on real time basis. Further this data collected on server will be used for:

1. Digital penalization of drivers crossing threshold “rashness” frequency,
2. Re-training and updation of model with latest data in order to produce results at par with the ever changing traffic conditions, road infrastructure and hence permissible

driving limits. Similar to android update, model could be updated in the Arduino boards at negligible cost.

The solution will serve following purposes:

- **Short Term Benefit** - Display of score on dashboard will provide a real time ‘Alert!’ making the driver conscious of his/her driving at critical moments.
- **Long Term Benefit**- After a macroscopic view of rash drivers through central servers, traffic authorities can conduct an easy penalization of all rash drivers.
- **Other Consumer Benefits**- The device can also be integrated with various other organizations. For example,
 - The device can be linked with the driving license number of the vehicle owner. Hence, those who are holding learning license and fail to obey the rules of proper driving can be penalized by not granting permanent license immediately. Permanent licenses of those disobeying traffic laws frequently can be taken away.
 - The insurance companies can use driving behavior as a segregation benchmark for good drivers to provide premium concessions and other offers

All the stated applications of proposed solution will serve as a behavioural deterrent for rash drivers curbing rash driving in longer run.

Organization of the paper

Section I contains the introduction and background of this paper and proposed solution. It depicts the contribution and basic implementation idea of the solution presented. Section II gives the literature review depicting challenges present in the previous works done related to driving analysis and then providing a comparison of the same with this paper. Section III is the detailed description of the proposed idea with the help of flow charts, data information and explanation of applied analytical techniques, generated plots and tables. Section IV discusses the results in the form of confusion matrices and parameters of accuracy, specificity and sensitivity. Section V talks about the scope of work and some limitations. Section VI concludes the research work with future directions.

II. RELATED WORK

For rash driving detection, most previous methods have deployed additionally dedicated sensors like OBD devices inside cars, smartphone based mobile applications or roadside sensors. The technologies include either direct detection through mobile inbuilt sensors, or exporting driving data to a remote server or cloud with further application of various analytical techniques on the basis of certain theoretical thresholds of violation. For instance: algorithm with self-created target variables, which disproves its authenticity.

In paper [3], Shi-Huang Chen et al designed a classification algorithm using Adaboost to determine whether the current driving behavior belongs to safe driving or not. They used 4 vehicle operation parameters - vehicle speed, engine speed, throttle position and engine load, via OBD interference, to generate ratios for classification into safe or rash.

In paper [11], Manish R. Prabhu et al provided an elaborated survey of various methods for analyzing driver behavior. They described the inbuilt smartphone sensors such as accelerometer, gyroscope, global positioning system (GPS), camera, gravity sensor, rotational vector sensor for analyzing driver behavior. The paper also mentions various challenges associated with these smartphone sensors.

In [12], Swati Dixit et al used a combination of bluetooth, smartphone inbuilt accelerometer sensor and cloud based backend for real time driver conduct monitoring. As per a set threshold for acceleration/deceleration, they classified a driving behavior into rash or safe and send the data to the server along with location. Through BAM (Business Activity Monitor) and remote cloud computation, the system overall provides real-time driver monitoring, trip analysis, and vehicle diagnostics.

In paper [13], AmolLakhe et al proposed k-NN classification and haversine algorithm to detect rash driving using smartphone sensors. The accelerometer sensor in mobile application is used to detect rash driving by historical pattern matching. If the driving is rash, GPS coordinates would be send to the nearest police station where further relevant action could be taken.

In [14], Ladly Patel et al proposed an innovative solution using infrared sensors installed on highways. This data collected from the IR sensors is sent to Arduino UNO to calculate the time needed by a specific car for moving from one point to the other on highways. Using this time, it calculates the vehicle speed and sends alerts if an over speed vehicle is detected.

As mentioned above, all the existing solutions use smartphones for data capture, rely on constant internet connectivity and share entire data with common cloud for remote computation. Our solution is devoid of any such dependency. Using smartphone not only distracts driver in navigating or other phone usage while driving but also comes with challenges such as phone-theft, battery discharge and loss of internet connectivity. Constant internet connectivity is an ideal scenario which is difficult to achieve at every location, with every network and every smartphone. Storing all the data from every vehicle on a common server and then applying computation is highly expensive as well as poses serious data privacy issues. Hence to provide a cheaper solution and abide by the data privacy rights, our system

calculates the score on the vehicle itself and shares only the overall score corresponding to each vehicle with only traffic authorities for better regulation. Also as mentioned as a drawback in existing work, GPS is a critical variable in capturing driver's exact location but it breaches privacy which is again violation of privacy rights. With regard to privacy issue, our algorithm doesn't keep track of GPS. The calculation of score doesn't require any constant network connection, rather a GSM module for weekly score transmission as and when it gets the network. Moreover, most of the previously proposed methods have taken into account only few general variables such as speed, acceleration, engine load etc. which limits their scope. We performed research on all the OBD parameters considering their influence on vehicle driving overall and hence, devising a separate initial analysis to choose final variables.

III. METHODOLOGY

The proposed rash driving detection method consists of data collection and acquisition, data pre-processing, principal component analysis, hierarchical clustering, and support vector machine (SVM) classification of driving as either safe or rash. Major constraints for designing the algorithm were the amount and quality of information available in the OBD-II dataset. The algorithm is designed on the data directly captured by OBD-II devices [7], and hence it doesn't contain any information regarding the driver actually being safe or rash. An OBDII dataset for 6 months and 14 cars was trained to build a classification model using Principal Component Analysis (PCA) and two Machine Learning algorithms – hierarchical clustering and Support Vector Machine (SVM) Classification. Clustering was used to classify entire dataset into categories for creating target categorical-variables and SVM classification was used to re-train that supervised data for model creation. The final model is able to classify a new set of the OBDII parameters into – rash or safe driving. This model will be uploaded to multiple Arduino UNO boards which further will be connected to OBDII devices in the vehicles. Thus, a score (1 for rash, 0 for safe) is generated per minute of driving.

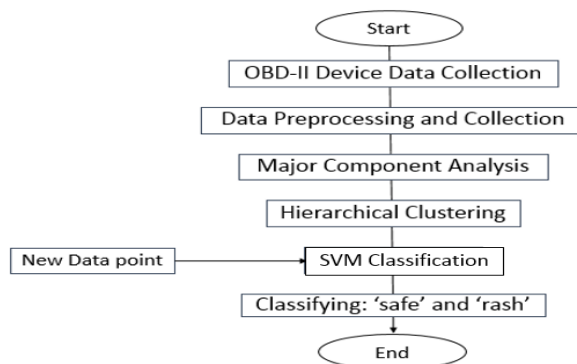


Figure 4: Flowchart of the proposed solution

This category will be displayed on the dashboard, providing a real time alert and making the driver conscious of his/her driving at critical points. With a GSM-module attached, the accumulated score report will be sent to a centralized server on weekly basis, giving a macroscopic view of rash drivers to concerned authorities.

OBDII parameters as metrics of analysis

Reasons of choosing major critical parameters against which the driver is evaluated is described ahead. Engine RPM, is defined as the rotations per minute of the crankshaft connecting piston rods to pistons heads. Cars should be driven at an optimal RPM relative to other parameters for an efficient driving. Throttle position, engine load and engine speed sensors provide information regarding appropriate usage of clutch, accelerator and brakes. Engine load is the torque output of the engine. This is required since the clutch, accelerator and brake pedals should be used as frugally as possible to lessen overall wear and tear of the system. Thus, for example, if a car is to be stopped, the speed should be gradually reduced - much in advance. The current condition of vehicle engine also majorly influences the safeness of driving. The coolant temperature keeps car from overheating. MAF (Mass air flow) allows car engine to regulate its performance and control air intake in order to keep the engine running in its best condition. Highly unregulated rash driving uses more fuel and creates high variations of fuel pressure in the system. Hence, fuel level and fuel pressure are significant parameters. Other considerable parameters are chosen with reference to environment: ambient air temperature, barometric pressure and air intake temperature. These parameters keep fluctuating while the vehicle is in motion, hence it is important to keep a check on the time elapsed since the time engine started depicting the importance of variable engine runtime

Analysis of variables and their influence on dataset

Initially 11 variables were present in the dataset obtained from the source [7]. Such large datasets have higher extent of multicollinearity among variables. Moreover, visualization and understanding of a multi-dimensional hyperspace is difficult. Due to the expected presence of redundancy and correlation among parameters, Principal Component Analysis (PCA) is used to reduce the original variables into a smaller set of independent, mutually exclusive and exhaustive variables, explaining majority of the variance in the entire dataset.

PCA identifies dimensions (or principal components) along which the variation in the data is maximum. Dimension 1 depicts the maximum variation while dimension 6 depicts the minimum. These dimensions i.e. newly created variables have a linear relation with original ones.

Table 1: Variable contribution in dimensions of PCA

Variables	Variable Contribution In Dimensions					
	Dimension 1	Dimension 2	Dimension 3	Dimension 4	Dimension 5	Dimension 6
BAROMETRIC_PRESSURE	4.41	1.75	0.04	1.66	24.34	59.7
ENGINE_COOLANT_TEMP	0.07	2.97	28.31	12.27	15.74	11.84
FUEL_LEVEL	0	0.72	49.18	9.71	3.34	0.18
ENGINE_LOAD	12.04	9.08	4.61	9.1	5.94	0.41
AMBIENT_AIR_TEMP	1.55	32.4	4.88	9.78	14.93	0
ENGINE_RPM	21.23	0.03	3.22	7.08	5.68	7.04
INTAKE_MANIFOLD_PRESSURE	4.51	17.98	0.41	12.73	5.18	7.08
MAF	24.97	4.12	0.12	1.36	0	0.22
AIR_INTAKE_TEMP	4.87	14.47	6.1	10.31	21.11	4.12
SPEED	19.31	0.04	2.11	13.56	0.05	9.4
ENGINE_RUNTIME	7.04	16.43	1.03	12.44	3.69	0.02

The amount of variance retained by each principal component (PC) is measured by an eigenvalue. The eigenvalues are examined, to determine the number of principal components to be considered. An *eigenvalue* > 1 indicates that PCs account for more variance than accounted by one of the original variables in standardized data.

Table 2: Eigen values and variances for each dimension

	Eigenvalue	Variance Percentage	Cumulative Variance Percent
Dimension 1	3.062	27.839	27.839
Dimension 2	1.471	13.372	41.212
Dimension 3	1.305	11.862	53.074
Dimension 4	1.146	10.419	63.492
Dimension 5	1	9.088	72.581
Dimension 6	0.902	8.196	80.777
Dimension 7	0.688	6.25	87.027
Dimension 8	0.616	5.596	92.623
Dimension 9	0.414	3.761	96.384
Dimension 10	0.237	2.156	98.541
Dimension 11	0.161	1.459	100

The Scree Plot is the plot of eigenvalues in descending order which is studied to decide the optimal number of principal components to be considered. A cut off is decided on the basis of added percentage contribution for each dimension. As per the scree plot below, if top 6 principal components are considered, then 80.8% of variation in entire dataset will be captured. So, instead of taking the original 11 variables, as per the coordination matrix, top 6 dimensions or principal components are considered.

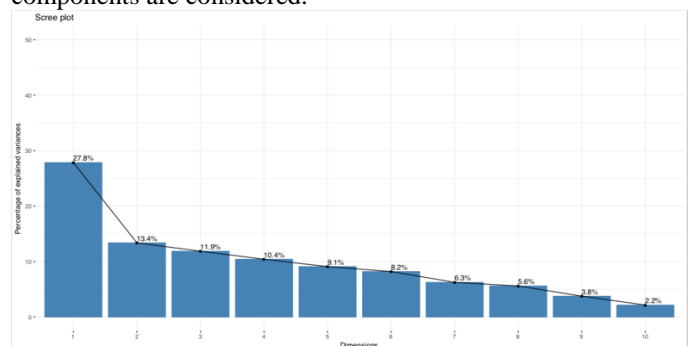


Figure 5: Scree plot - Bar plot of dimensions in descending order of variation contribution in dataset

For getting a better insight regarding the principal components, correlation within variables as well as with the principal components is visualized using variable correlation plot. It is interpreted as:

1. Positively correlated variables are grouped together
2. Negatively correlated variables are positioned on opposite quadrants.
3. The distance between variables and the origin measures the quality of the variables on the factor map. Variables that are away from the origin are well represented on the factor map.

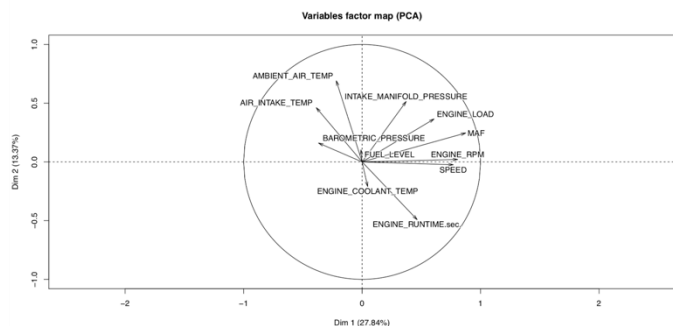


Figure 6: Correlation circle

The quality of representation of the variables on factor map is called \cos^2 . A high \cos^2 (for MAF, Engine RPM, Speed, Ambient air temperature) indicates a good representation of the variable on the principal component. In this case the variable is positioned close to the circumference of the correlation circle and position vector is longer. While a low \cos^2 (engine coolant temperature, Fuel level) indicates that the variable is not perfectly represented by the PCs. In this case the variable is closer to the centre of the circle with a shorter position vector.

After understanding the principal components thoroughly, dataset consisting of above depicted variables was finalised for building algorithm.

Hierarchical Clustering and Algorithm

Hierarchical clustering is used to classify the data set into the 2 groups –safe and rash [16]. The paper use agglomerative (bottom-up) clustering method. The algorithm works as follows:

- Step 1** - Assign each observation to its own cluster.
- Step 2** - Then, compute the similarity between each of the clusters and join the two most similar clusters.
- Step 3** -Finally, repeat above steps until there is only a single cluster left.

After applying hierarchical clustering, two clusters were obtained, and the properties of each cluster were studied and analysed individually to label them as safe or rash driving behaviour. After this, a final dataset was prepared along with a binary target variable named driving type – ‘safe’ or ‘rash’.

SVM Classification

On the prepared dataset after hierarchical clustering, SVM classification was performed to create final model for classifying new data points into safe or rash [10]. The algorithm works as follows:

- Step 1** - Map the input variables to a high dimensional vector space.
 - Step 2** - Generate set of hyperplanes that can be used for classification
 - Step 3** - Choose the hyperplane that has the largest distance to the nearest training data points.
- The closest points that identify this line are known as support vectors. And the region they define around the line is known as the margin.

IV. RESULTS AND DISCUSSION

The paper validates the accuracy obtained after hierarchical clustering with existing criteria [3] for normal and rash driving behaviours, which defines:

- Normal vehicle driving condition as the relative ratio of the vehicle speed and engine speed being maintained between 0.9 and 1.3 (test in the same gear); the relative ratio of the engine speed and throttle valve being maintained between 0.9 and 1.3; and the engine load being maintained between 20% and 50%.
- Bad vehicle driving condition as the relative ratio of the vehicle speed and engine speed being maintained either above 1.3 or below 0.9; the relative ratio of the engine speed and throttle valve being maintained either above 1.3 or below 0.9; and the engine load being maintained either more than 50% or less than 20%.

For checking the accuracy of hierarchical clustering analysis, a target variable was created based on these ratios. On comparing these theoretical target variables, as per above definitions, with the predicted ones, an accuracy of 80.47% was obtained.

Next to it, on performing supervised learning with SVM classification on this dataset with created target variables, a classification accuracy of 99.45% was achieved

Confusion Matrices

Table 3: Model performance parameters

	Sensitivity	Specificity	Accuracy
Hierarchical Clustering	93.927%	67.014%	80.47%
SVM Classification	99.97%	96.83%	99.45%

Table 4: Confusion matrix for SVM model

		Actual variables (on the basis of hierarchical clustering)		
			Bad Driving (1)	Normal Driving (0)
SVM Classification prediction	Bad Driving (1)		8543	54
	Normal Driving (0)		2	1652

Table 5: Confusion matrix for Hierarchical Clustering model

Hierarchical Clustering prediction	Actual variables (on the basis of theoretical ratios)		
		Bad Driving (1)	Normal Driving (0)
	Bad Driving (1)	5847	1328
Normal Driving (0)	378	2698	

Regular updation of the algorithm would be required. For each new updation, the data accumulated on the server will be used. The data used for re-training of model will be bigger, more recent and hence leading to better classification. Similar to an Android update, this model updation in the Arduino electronic chip can be done online adding no further cost to the entire solution.

V. SCOPE AND LIMITATIONS

- Presence of an OBD port connected to a model-encoded Arduino board is must in the vehicle. Although most of the latest vehicles are mandated to have an OBD port.
- Solution is created on OBDII data of cars collected in a certain region. For a larger scale implementation, a more generalized dataset will be required covering wider areas, more vehicles and vaster driving patterns.
- As the solution is keeping a check on driving behavior, people will be discouraged to purchase and install it and hence, it needs to be made mandatory by regulatory authorities for 100% success.

VI. CONCLUSION AND FUTURE SCOPE

Designing and implementing the clustering algorithm on the dataset, provided the target variables (as rash or safe driving). Further any new datapoint, captured in real time, can be classified into safe or rash using Support Vector Machine (SVM) with a really high accuracy.

Future scope of improvement in this solution:

- Hierarchical clustering can be performed on a better and stronger dataset with more theoretical relations about safe or rash driving.
- Training dataset, if captured, with more precision and accuracy can produce better results in hierarchical clustering.
- This solution can be practically implemented in vehicles with OBDII ports through relevant mechanism and administrative infrastructure.

Cost Estimate

Entire solution hardware and other charges involved are one-time cost. Although, cost depends on the hardware manufacturing company, vehicle compatibility and external factors such as overall maintenance and caring. A rough cost estimation, basis average is as follows:

OBD device - INR 500/-

Arduino UNO board* - INR 1200/-

GSM module cost -INR 1100/-

Soldering and base installation charges - INR 500/-

*Atmega board costs cheaper - around INR 500 but it will incur a model updation cost once every 3 months - which makes it costlier overall than Arduino board which itself has negligible model updation cost

REFERENCES

- [1]. Miyaji, M., Danno, M., Oguri, K., "Analysis of driver behaviour based on traffic incidents for driver monitor systems", In Intelligent Vehicles Symposium, IEEE(2008) pp. 930-935.
- [2]. Liang, J., Cheng, X., Chen, X., "The research of car rear-end warning model based on mas and behavior", In: Power Electronics and Intelligent Transportation System, IEEE (2008), pp. 305-309.
- [3]. Chen, S., Pan, J. and Lu, K. (2018), "Driving Behaviour Analysis Based on Vehicle OBD Information and AdaBoost Algorithms", [online] Iaeng.org.
- [4]. "E/E Diagnostic Test Modes". Vehicle E E System Diagnostic Standards Committee. J1979. SAE International. 2017-02-16. doi:10.4271/J1979_201702.
- [5]. "Digital Annex of E/E Diagnostic Test Modes". Vehicle E E System Diagnostic Standards Committee. J1979-DA. SAE International. 2017-02-16. doi: 10.4271/J1979DA_201702.
- [6]. Enrique, et al. "PCA - Principal Component Analysis Essentials." Correlation Matrix : A Quick Start Guide to Analyze, Format and Visualize a Correlation Matrix Using R Software - Easy Guides - Wiki - STHDA, 23 Sept. 2017.
- [7]. Barreto, Cephas. "OBD-II Datasets." RSNA Pneumonia Detection Challenge | Kaggle, 21 Aug. 2018, www.kaggle.com/cephasax/obdii-ds3
- [8]. "Hierarchical Clustering", Saedsayad.com, 2018. [Online]. Available:https://www.saedsayad.com/clustering_hierarchical.htm
- [9]. Khedkar, Sujata et al, "driver evaluation system using mobile phone and OBD-II system", In International Journal of Computer Science and Information Technologies (IJCSIT) 2018, Vol. 6 (3), 2015, 2738-2745
- [10]. Cortes, C., & Vapnik, V. (1995). Support-vector networks. Machine learning, Vol. 20(3), 273-297.
- [11]. Prabhu, M., Deorukhakar, D., Amberkar, S., & Ambre, K. (2016), "A Review on Rash Driving Detection and Alert System", In IOSR Journal of Computer Engineering (IOSR-JCE), Conf. 16051, e-ISSN: 2278-0661, p-ISSN: 2278-8727 pp. 47-49
- [12]. Dixit, S., Thomas, P., & Agarwal, C. (2016). "A Cloud-based Driver Monitoring for Inefficient Driving Behavior using OBD2 Telematics". In International Journal of Computer Applications (0975 – 8887) Vol. 153, Issue. 9, November 2016
- [13]. Lakhe, A., Patel, C., Raut, H., & Shinde, S. (2017), "Smart Phone based Rash Driving Detection". In International Journal of Innovations & Advancement in Computer Science (IJACS) ISSN 2347 – 8616, Vol. 6, Issue. 11, November 2017
- [14]. Patel, L., Gaurav, K., & V, D. (2018), "Detection of rash driving on highways", In International Journal of Engineering Development and Research (IJEDR), Vol. 6, Issue. 3, ISSN: 2321-9939
- [15]. Juneja, P. (2019), Reinforcement Theory of Motivation, Retrieved from https://www.managementstudyguide.com/reinforcement-theory-motivation.htm

- [16]. Zhao, Y., Karypis, G., & Fayyad, U. (2005). Hierarchical clustering algorithms for document datasets. *Data mining and knowledge discovery*, 10(2), 141-168.
- [17]. Iwavesystems.com. (2019). *OBD II Connected Car Device | iWave Systems*. [online] Available at: <https://www.iwavesystems.com/obd-ii-connected-car-device.html>
- [18]. Hareendran (2017). *Hands-on review: When hacking an OBD-II adapter, choose carefully*. [online] Electronic Products. Available at: https://www.electronicproducts.com/Interconnections/Wire_and_Cable/Hands_on_review_When_hacking_an_OBD_II_adapter_choose_carefully.aspx [Accessed 24 Jan. 2019].

Authors Profile

Ms. Isha Gupta pursued Bachelor of Technology from Indian Institute of Technology (IIT), Roorkee in 2017. Right from college she was fascinated by optimization, improvisation and advancement of process. She is passionate about machine learning and along with working as a data analyst, she is pursuing independent research in data science field to create feasible and convenient solutions for day-to-day problems.



Ms. Arushi Agrawal pursued Bachelor of Science from Indian Institute of Technology (IIT) Kanpur in 2017. Coming from Mathematics & statistics background and currently working as a data scientist, she is deeply passionate about Machine Learning and Artificial Intelligence. She strives to make general lifestyle easier by performing research and implement data-driven solutions in everyday life.

