

Driving Behavior Analysis through CAN Bus Data in an Uncontrolled Environment

Umberto Fugiglando*, Emanuele Massaro*, Paolo Santi*[‡], Sebastiano Milardo*[†], Kacem Abida[¶],
Rainer Stahlmann[§], Florian Netter[§], and Carlo Ratti*

*MIT Senseable City Lab, Cambridge, MA, USA

[‡]Istituto di Informatica e Telematica del CNR, Pisa, Italy

[†]University of Palermo, Palermo, Italy

[§]AUDI AG

[¶]VW Group Electronics Research Laboratory

email: {umbertof, emassaro, psanti, milardo, ratti}@mit.edu,

rainer.stahlmann@audi.de, kacem.abida@vw.com

Abstract—Cars can nowadays record several thousands of signals through the CAN bus technology and potentially provide real-time information on the car, the driver and the surrounding environment. This paper proposes a new method for the analysis and classification of driver behavior using a selected subset of CAN bus signals, specifically gas pedal position, brake pedal pressure, steering wheel angle, steering wheel momentum, velocity, RPM, frontal and lateral acceleration. Data has been collected in a completely uncontrolled experiment, where 64 people drove 10 cars for or a total of over 2000 driving trips without any type of pre-determined driving instruction on a wide variety of road scenarios. We propose an unsupervised learning technique that clusters drivers in different groups, and offers a validation method to test the robustness of clustering in a wide range of experimental settings. The minimal amount of data needed to preserve robust driver clustering is also computed. The presented study provides a new methodology for near-real-time classification of driver behavior in uncontrolled environments.

Keywords: Driving behavior, CAN bus, feature extraction, unsupervised learning, drivers segmentation.

I. INTRODUCTION

Modern cars are equipped with several hundreds of sensors and electronic control units (ECUs) [1] that, beyond guaranteeing an optimal functioning of the engine, provide the driver with more safety, control and entertainment. These almost real-time data provide information on the car, the driver and the surrounding environment and can be used to study, analyze, predict and understand a large variety of problems, such as traffic congestion, vehicle energy consumption and emissions, urban mobility and drivers' habits [2].

This huge amount of diverse data has been made available by the CAN bus technology, a serial broadcast bus developed by Robert Bosch in 1986 [3] that allows communication among the electronic control units devices mounted on the car. CAN technology has become *de facto* a standard in car embedded systems providing access to data from an order of several thousands signals, recording at a sub-Hertz frequency information about the car and its surroundings.

With this technology being implemented in modern cars, the amount and variety of collected data increases and all the aforementioned applications can be extended and improved with respect to the state of art of GPS-based technologies. Data availability is not a restrictive aspect anymore as insights from travels can be collected automatically, without the need to modify the car structure or to specifically design an experiment. Moreover, in the present research we leverage a data stream in the order of few gigabytes per hour, which represents just a significative sub-sample of all the information travelling on the CAN bus: this amount of data will only increase with the advent of new autonomous driving cars [4].

A. Driving behavior

The characterization of driving behavior is not only crucial for accident prevention, as most of car accidents are due to human mishandling, but it is also important for designing driving models, which are the core of algorithms that might make the future of self-driving cars possible [5]. Driving behavior characterization is useful also for car insurance companies to quantify accident risk and provide personalized rates. State-of-art technology implements models mostly based on GPS location, traveled distance and coarse grained speed profile [6], [7]. A richer information like the one coming from CAN bus could better characterize human driving behavior and, consequently, accident risk.

In order to be able to use CAN data to characterize drivers in real application scenarios we need to solve two very challenging problems: (1) providing a methodology for consistently identifying driving behavior in a completely uncontrolled environment, and with very limited knowledge of the surrounding conditions; and (2) minimizing the communication and computational load needed to solve (1). This paper introduces and discusses ideas to tackle these challenges and bring CAN bus based driver characterization closer to reality.

More specifically, the goal of the present research is to extract features from CAN bus signals and assess to what extent they are useful for finding similarities among drivers using a clustering algorithm. Given the enormous amount of

data generated by the CAN bus – in the order of a few gigabytes of data per hour – it is not feasible to communicate and process the raw output of the CAN bus in real time to characterize drivers. As such, feasibility of the devised driver characterization methodology is bounded to the definition of a strategy to substantially reduce the amount of data to be processed to perform the driver identification task. Thus, in the second part of the paper we explore different data subsampling methods that allow minimizing data communication between vehicle and infrastructure while guaranteeing robust driver behavior characterization.

The paper is organized as follows. Section II describes the details of the data collection process and the signals considered. Section III is devoted to the clustering of the drivers. Section IV addresses the sampling method question. Finally, section V concludes the paper providing a summary of the future research directions.

B. Related work

In general, research on driving behavior in scientific literature can be classified according two perspectives: (1) the purpose of the research, e.g. driver recognition, maneuver recognition, aggressive or eco-friendly driving detection, *etc.* or (2) the data used for the analyses, i.e. GPS locations, CAN bus data, audio-video data, cellular phone data, car simulator data.

Early studies have been made with the aim of characterizing driving behavior by building a dynamic model to eventually implement a control system that would react like a human, to be used for example in self-driving cars. Models have been proposed to anticipate the driver actions by few seconds [8] or to predict the drivers intended cruising speed up to 20 seconds in advance of reaching that speed [9]. All these works have been validated using data coming from car simulators. Data acquired by a simulator have also been used to quantify the drivers' skills [10].

Some other works, on the other hand, have been conceived to recognize driving maneuvers (e.g. passing, changing lines, turning, starting and stopping) leveraging CAN data: for example, in [11] the drivers were asked by an instructor in the vehicle to perform given maneuvers.

Carmona et al. [12], through a novel hardware tool designed to integrate data from CAN bus, GPS and an Inertial Measurement Unit (IMU), attempt to classify real-time normal and aggressive driver behavior. The classification was performed in an experiment where 10 drivers have been asked to drive the same route twice, in a normal and aggressive way respectively.

CAN sensors have also been coupled with external devices, designed and mounted specifically on the vehicle for the purpose of the experiment, like 3D cameras for eye monitoring or wearable devices used to collect biomedical signals. These experiments are more “human-centric” and are aimed at understanding how drivers' bad habits or distractions are reflected in their way of driving: Choi et al. [13] and, lately, Li et al. [14] detected and classified distraction tasks (e.g. tuning the

radio, interacting with an automatic voice portal) using audio and video data coupled with CAN bus data.

On the other hand, some works focus on the driver recognition problem, which attempts to distinguish different drivers only by looking at the CAN bus data. Wakita et al. [15], using data coming from a car simulator, made a comparison between parametric and nonparametric models, concluding that nonparametric approaches perform better in terms of percentages of drivers correctly recognized. Hallac et al. [1] leveraged the same database used in this work achieving a prediction accuracy of 76.9% for two-driver classification, and 50.1% for five drivers. Miyajima et al. [16], [17] performed driving recognition modelling on pedal operation patterns acquired by CAN bus sensors by means of a cepstral method, both on a car simulator and on real cars involving 276 drivers. However, the exact setting of the experiment, the type of road the drivers used, and how they have been instructed to drive is not specifically mentioned in the paper. Moreover, the vehicle used for data collection (a minivan, [18]), equipped with cameras, microphone, computer rack, power suppliers and amplifiers, suggests that the experimental conditions were far from an everyday context in personal driving.

More recent work uses data coming from mobile phones sensors (accelerometer, gyroscope, magnetometer, GPS, video): in [19], cell phone sensors data have been coupled with CAN bus data as a “ground truth” for isolating acceleration, braking and turning events: the problem of driver recognition was addressed, but the experiment involved only two drivers and reached only 60% of accuracy. Moreover, mobile phones sensors have been used to detect aggressive [20] or drunk [21] drivers.

In contrast to the present research, in which normal cars have been used, most of the previously cited works used cars developed in specific projects, like the UTDrive project¹ [22] or a specifically designed “vehicle corpora for research” [18], [23]. Finally, uncontrolled experimental settings have been used in the SHRP2 Naturalistic Driving.²

Study, where driving behavior has been analyzed using traditional techniques (thus not through CAN data) and in another large experiment called “EuroFOT” (European large scale Field Operational Test on in-vehicle system)³, where CAN bus data have been used with the only aim of evaluating the impact of 8 different driving assistance systems.

Comprehensive analyses of driving behavior models, tools and experiments can be found in [5], [14], [24]. Summarizing, none of the existing work analyzed usage of CAN bus data for driver classification in a completely uncontrolled and open driving environment. Furthermore, the issue of how to reduce the communication and computational load related to driver classification has, to our best knowledge, never been addressed so far.

¹<http://www.utdallas.edu/research/utdrive/>

²<https://insight.shrp2nds.us/>

³<http://www.eurofot-ip.eu/>

C. Motivations

As it turns out from the previous section, the main novelty of this paper in the field of human driving behavior analysis is the combination of (1) large number of drivers, (2) completely uncontrolled experimental settings and (3) quantity of data recorded.

This sets new limits and possibilities to the present research: limits in terms of the variety of the signals acquired, carrying useful information not supported by “ground truth”, i.e. information we can consider as “true” to which compare the experimental data (for example the “aggressiveness” of the driver, his driving skills or his number of incidents). On the other hand, the framework of the present research opens the way to new CAN-based technologies that could find application in real-life scenarios.

II. DATA COLLECTION

A. Experimental settings

The dataset used in the present research has been collected during an experiment carried out by AUDI AG and Audi Electronics Venture. The data collection experiment took place in the city of Ingolstadt (Germany) and involved 64 different drivers, who have not been instructed in any way on the route they had to drive, on the speed or on the behavior they had to follow during the driving. This gives to the present study its unique characteristic of an experiment under uncontrolled testing conditions. A test fleet of ten Audi A3 vehicles was retrofitted with data loggers. This prototype system enables data acquisition for research purposes.

The data collection phase took place in 2014 with a total of 55 days of experiment. Cars were picked up by the drivers in a central deposit and had to be returned within the same day. Each time a user switched on the car engine, the computer registered a new *session*. A total of 1987 sessions have been recorded, and more than 2135 hours of driving data for each of the 2418 sensors have been acquired. Each user drove an average of 31 sessions, whose average duration was 64 minutes.

CAN bus signals have been recorded on a data logger⁴ and processed in a later phase. The sampling is not uniform due to the particular characteristics of the CAN bus and the signals. Therefore, high frequency signals are constantly sampled at 20 Hz, while low frequency sensors reports their data only when there is a change in their value (e.g. rain sensors, seatbelt sensors, etc.) but for the sake of simplicity all the signals considered in the analysis have been resampled at 4 Hz through linear interpolation.

B. Signals selection

Among the 2418 signals transmitted on the CAN bus, in this work we concentrated the analyses on eight signals:

- Brake pedal pressure (BRK)
- Gas pedal position (GAS)
- Revolutions per minute (R.P.M.)

⁴No personal information on the drivers have been recorded.

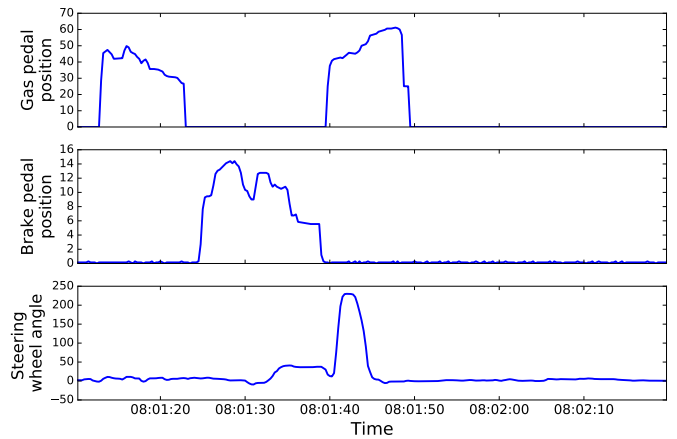


Figure 1: Example of signals acquired by the gas pedal position sensor (top), brake pedal pressure sensor (middle), steering wheel angle sensor (bottom). The three signals have been acquired synchronously.

- Speed (SPD)
- Steering wheel angle (S.W.A.)
- Steering wheel momentum (S.W.M.)
- Frontal acceleration (F. ACC.)
- Lateral acceleration (L. ACC.)

These signals are directly or, in some cases, indirectly related to the interaction between the driver and the vehicle. For instance, pedals and steering wheel signals directly reflect driver’s movements and actions, without any “transfer function” between the input (the driver’s action) and the output (the signal); some other (speed, rpm and accelerations) represent on a phenomenological point of view quantities that a person can “feel” during the driving and could reflect specific driving habits: for example, a driver’s attitude to exceed speed limits. An example of the collected signals is reported in Figure 1.

III. GROUPING DRIVERS’ BEHAVIOR

In this section we propose a methodology that allow us to group in a consistent way the drivers according to common characteristics. This methodology is composed of 4 different steps: A) Features extraction, B) Features normalization, C) Dimensionality reduction and D) Unsupervised Clustering.

A. Feature extraction

Any signal \mathbf{x} in the database can be represented as a set of pairs of the type (x_i, t_i) , where $i \in \mathbb{N}$ and t_i is the timestamp corresponding to the acquisition of the signal value x_i where x_i is a floating point number. From each considered signals we extract the following 7 indicators:

- 1) values of the signal for each sample: x_i .
- 2) difference quotient (discrete first derivative) of the signal between two consecutive samples: $\frac{x_{i+1} - x_i}{t_{i+1} - t_i}$. This measure quantifies the intensity of signal variation over time. Let us now define J as the set of indexes for which the values x_i are singular points (local maxima or minima), i.e. $J = \{i : (x_i - x_{i-1})(x_{i+1} - x_i) < 0\}$, and by $J_{\max} \subset J$ the set of only local maxima. Moreover, let us define on those

Feature	Description
1	Values of the signal for each sample
2	Difference quotient (discrete first derivative)
3	Time interval between two singular points
4	Values of the local maxima
5	Moving mean
6	Moving median
7	Moving standard deviation

Table I: Features definition.

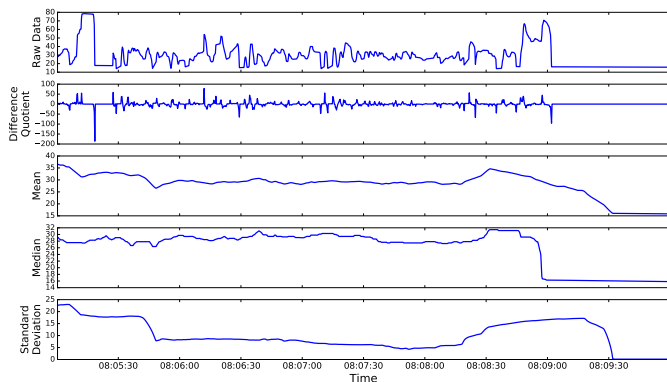


Figure 2: A sample of some of the features extracted from the eight considered signals. In particular, the figure shows the gas pedal angle signal and its difference quotient, mean, median, and standard deviation.

sets a relation \prec , where $j \prec k$ means that j is the largest element of the set that precedes k , i.e. $j = \max\{i \in J : i < k\}$.

- 3) time interval between two singular points: $t_j - t_k$, $j, k \in J$, $j \prec k$. This feature represents the frequency of its peak points, or in other words the rapidity of variation of the signal when it reaches extreme values.
- 4) value of the local maxima: x_j , $j \in J_{\max}$. This feature provides the intensity of the extreme values of the signal. In a temporal window of one minute and remembering the 4 Hz sampling we define the set of indexes $I_i = \{i - 120, \dots, i + 120\}$ and the following.
- 5) moving mean, averaging the values x_i over a temporal window of 1 minute: $\frac{1}{240} \sum_{j \in I_i} x_j$.
- 6) moving median, the median value of the set $\bigcup_{j \in I_i} x_j$.
- 7) moving standard deviation, the variance of the values in the set $\bigcup_{j \in I_i} x_j$.

Table I summarizes the features defined above for a quick reference, while Figure 2 shows a plot of a sample signal and some of the features.

B. Features normalization

For any given signal \mathbf{x} of floating point type, we denote by $\mathbf{w}^{k,u}$ the vector of the feature k for user u , obtained by calculating the functions defined above on the vector \mathbf{x} , joining all the sessions of the same user. We then normalize each

vector $\mathbf{w}^{k,u}$ in the following way. Outliers removals has done by keeping only the values between the 2nd and 98th percentile. We consider the vector $\mathbf{w}^{k,u}$ as a set of statistical samples that are used to build frequency histograms.

In order to get for each user histograms with the same bins, we define the set

$$W^k = \bigcup_{u \in \mathcal{U}} \bigcup_i \{w_i^{k,u}\},$$

where \mathcal{U} is the set of users, and partition the interval $[\min W^k, \max W^k]$ into 10 equal intervals⁵ (bins) b_1^k, \dots, b_{10}^k . Then, for each user and for each indicator, the histogram $H^{k,u}$ for the vector $\mathbf{w}^{k,u}$ with bins b_1^k, \dots, b_{10}^k can now be computed, i.e. each bar of the histogram has a value $h_i^{k,u}$ which is the number of items of the vector $\mathbf{w}^{k,u}$ belonging to interval b_i^k . Finally, all the histograms are normalized, obtaining new values $\tilde{h}_1, \dots, \tilde{h}_{10}$ according to the formula

$$\tilde{h}_i^{k,u} = \frac{h_i^{k,u}}{\sum_{j=1}^{10} h_j^{k,u}},$$

so that $\sum_{i=1}^{10} \tilde{h}_i^{k,u} = 1$.

According to our definition, features in form of histograms can be interpreted as a discrete version of the sample distributions of the indicator vectors. This definition, along with its probabilistic interpretation, has two main advantages: it allows to perform analyses on objects which have a probabilistic meaning, while on the other hand it keeps machine learning algorithms relatively simple due to the low dimensionality of the data.

In the following analyses, for data homogeneity we consider users who drove in total at least 10 hours, reducing the number of considered users to 54 from the initial 64.

C. Dimensionality Reduction

In this section we use the K -means clustering algorithm [25] to leverage the features defined in the previous section with the aim of grouping drivers upon common similarities. This is a novel approach in this field and therefore it requires an assessment of the validity of the method in terms of robustness and scalability.

It is worth remarking that the vectors $H^{k,u}$ are 10-dimensional data-points, being them histograms with 10 bins. In order to plot them on bi-dimensional space, therefore, a dimensionality reduction technique has to be performed. In this work we use *Principal Component Analysis* (PCA), a well known statistical procedure that decreases the dimensionality of a space projecting it into another one whose dimensions (principal components) are orthogonal to each other and such that the variance of the projected data-points on the principal components is maximized [25].

⁵The number 10 has been chosen after some preliminary analyses. The rationale for choosing the number of bins was to have a sufficient number of bins to well represent the shape of the probability density distribution, but small enough to keep the computation of the machine learning algorithms feasible.

	Features						
	1	2	3	4	5	6	7
BRK	1.00	0.99	1.00	0.96	1.00	0.66	0.89
GAS	0.90	0.98	0.93	0.85	0.79	0.96	0.78
R.P.M.	0.61	0.95	0.57	0.78	0.70	0.98	0.73
SPD	0.61	0.88	0.54	0.77	0.55	0.91	0.65
S.W.A.	0.92	0.99	0.92	0.97	0.95	0.97	0.80
S.W.M.	0.79	0.96	0.79	0.94	0.89	0.98	0.88
F. ACC.	0.82	0.94	0.76	0.81	0.87	0.97	0.75
L. ACC.	0.99	0.99	0.99	1.00	1.00	0.98	0.99

Table II: Total variance of the original data explained by the first two principal components, for each combination of signal and feature.

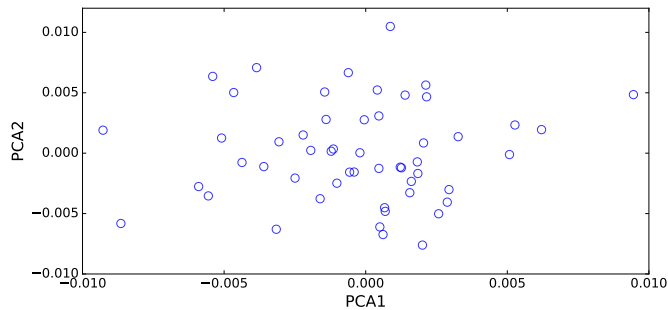


Figure 3: PCA representation for Feature 1 of the gas pedal position signal, where each point represents a different driver.

Table II shows that for most of the combinations of signals and features, the first two principal components explain more than 80% of the total variance of the original high dimensional data. Figure 3, consequently, reports an example of a bidimensional representations of the features (Feature 1 for the gas pedal signal) where each dot corresponds to a driver. It can be noticed that there are no well separated clusters: this can be expected thinking that human behavior typically varies in a range that forms a continuum. For this reason, the word “segmentation” more accurately describes this process than “clustering”: some common behavior can be identified, while some “outliers” slightly deviate from the average.

D. Unsupervised Clustering

Having no previous information about the drivers and their behavior, it is not known *a priori* the number of different attitudes to be detected and whether a driver is correctly classified (as opposed, for example, to [12]). For instance, we cannot tell which of the datapoints represent “aggressive”, “dynamic” or “eco-friendly” drivers, as this information is not accessible to us. This remarks motivate the choice of clustering techniques, being part of the *unsupervised learning* approaches to data analysis, used when no previous knowledge on the data is available. In fact, unlike *supervised learning*, the former is an exploratory analysis that does not rely on a *ground truth*, a concept identifying the *a priori* known information of the data or the information provided by direct observation, as opposed to information provided by inference.

However, a problem arises when the optimal number of clusters has to be chosen and when the overall quality of the clustering has to be evaluated. Some common techniques try to

Algorithm 1: K -means clustering cross-validation algorithm.

```

for each feature  $k = 1 \dots 7$  do
  for number of clusters  $K = 2 \dots 10$  do
    for number of trials  $i = 1 \dots 40$  do
      for each user  $u \in \mathcal{U}$  do
        randomly permute the elements of vector
           $\mathbf{w}^{k,u}$ ;
           $\mathbf{w}_T^{k,u} =$  first 70% elements of  $\mathbf{w}^{k,u}$ ;
           $\mathbf{w}_V^{k,u} =$  last 30% elements of  $\mathbf{w}^{k,u}$ ;
        compute histograms  $\{H_T^{k,u}\}_{u \in \mathcal{U}}$  and
           $\{H_V^{k,u}\}_{u \in \mathcal{U}}$  as in III-A;
         $T = \{H_T^{k,u}\}_{u \in \mathcal{U}}$  (training set);
         $V = \{H_V^{k,u}\}_{u \in \mathcal{U}}$  (validation set);
         $\mathcal{C}_T = K$ -means clustering on  $T$ ;
         $\mathcal{C}_V = K$ -means clustering on  $V$ ;
         $v_i = \text{V-measure}(\mathcal{C}_T, \mathcal{C}_V)$ ;
       $M_{k,K} = \text{mean}(\mathbf{v})$ ;
       $S_{k,K} = \text{standard-deviation}(\mathbf{v})$ ;

```

address this difficulty, for example the plot of SSE (sum of the squared differences between each observation and its group’s mean [25]) or the silhouette index (a measure of how similar an object is to its own cluster compared to other clusters [26]), but as mentioned above in our case clusters are not well separated and those techniques do not provide useful results.

Inspired by the widely used method of cross-validation used in supervised learning, we propose here a new approach for establishing the optimal number of clusters, based on the concept of “robustness” of the clustering to the road sampling. In fact, remembering that the clusters are made up of distributions that come from sampled data, the clusters should be invariant to a subsampling of the original data. In other words, comparing the clusters generated by different subsampling of the original data, those clusters should be similar.

The method proposed is described in Algorithm 1 and can be synthesized as follows. For each user u and for each feature k , the vector $\mathbf{w}^{k,u}$ is divided into two different vectors: 70% of its components, taken randomly, form the vector $\mathbf{w}_T^{k,u}$ (*training vectors*), while the other 30% form the vector $\mathbf{w}_V^{k,u}$ (*validation vectors*). After having computed the histograms for the two sets of vectors, a K -means cluster algorithm is performed separately on both the training set and the validation set, producing two different clusterings of the same set of drivers. These two clusterings are then compared using a metric called “V-measure” [27], a score ranging from 0 to 1 and evaluating the similarity of the clusterings: if the clusterings are exactly the same (except for permutations on the labels of each cluster) the score is 1, while the score is closer to 0 as the clusterings are more dissimilar. This operations are repeated for a number of clusters K ranging from 2 to 10. Moreover, being the

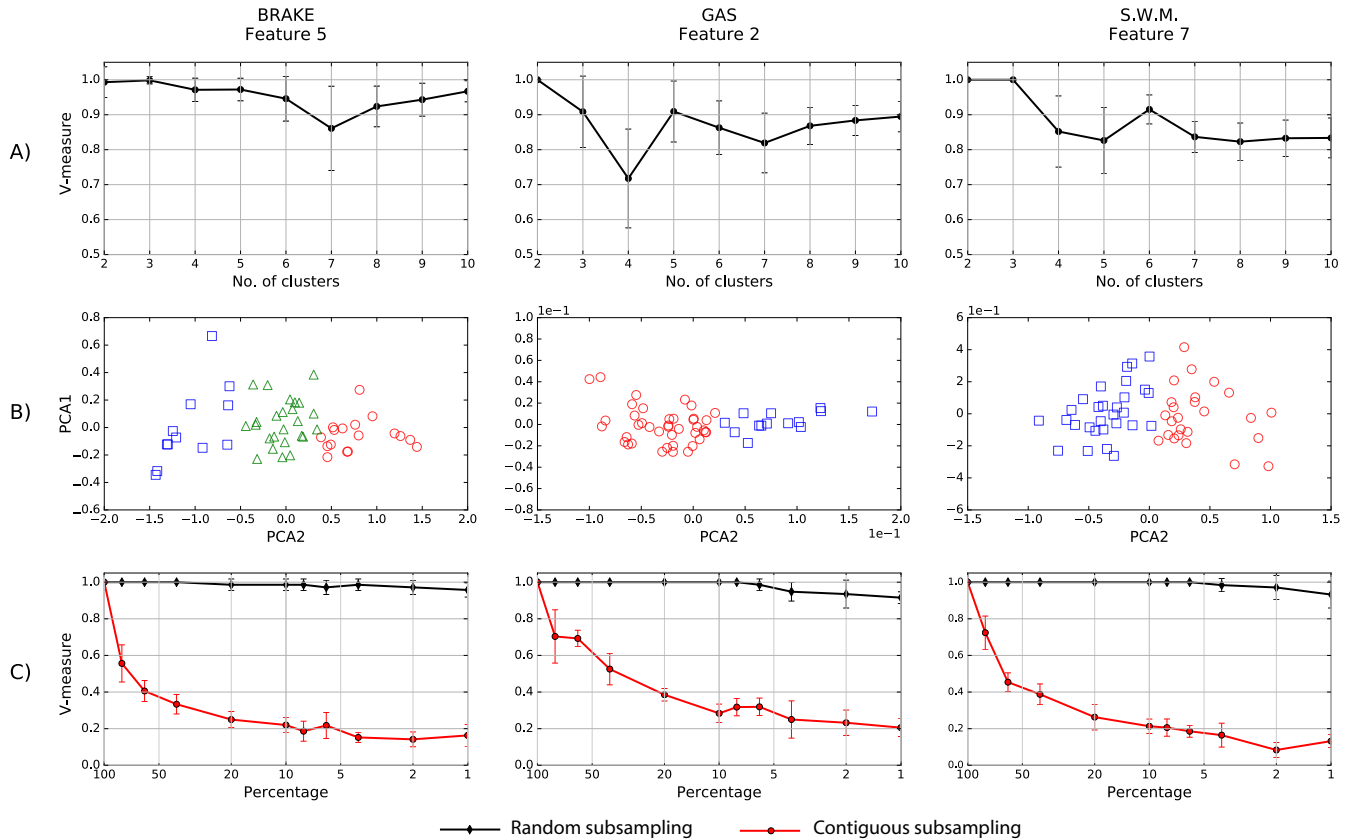


Figure 4: Plot of analyses for selected combinations of signals and features: (A) Output of Algorithm 1, plotting the V-measure for different values of K ; (B) Drivers clusterings for different signals and features. The K -means algorithm has been run on all data in the database and for the optimal values of K as in Table III; (C) Subsampling methods: the graphs show the V-measures of the comparisons of the K -means clusters generated using all the data in the database, with the clusters generated by a subset of the data (validation set), for different sizes of the validation set (100%, 50%, 20%, 10%, 5%, 2%, 1% of the original data). The clusterings use the optimal values of K as in Table III.

subsampling random, for each value of K the algorithm is repeated 40 times: averages and standard deviations of the scores for each value of K are calculated and lead to plots like the ones in Figure 4A.

The optimal number of K that provides a “robust” clusterization is thus defined as the value of K that maximizes the corresponding V-measure in Algorithm 1. Table III provides, for each combination of feature and signal, the optimal values together with mean and variance of their corresponding V-measures. In case of ties of the V-measure, the lowest value of K has been considered as the optimal one.

Results clearly show that there are some numbers of clusters that separate users in a better way in terms of “robustness”. For example, feature 2 for the gas pedal position separates drivers in two different groups, which keep exactly the same in all the 40 repetitions of the cross-validation algorithm, whilst it is not the same for $K = 4$.

Overall, some features and some signals perform better than other: the brake pressure signal is the one with most promising results, followed by the gas pedal position and the steering wheel. This is a first important result, as it confirms what has been already found in the literature with data from an unstructured experiment [16].

Finally, Figure 4B reports the results of the K -means clustering for a selection of signals (see Figure 5 in the Appendix for a comprehensive chart), with values of K as in Table III.

IV. DATASET REDUCTION

Once we have verified that a consistent, robust clustering of drivers is possible also in completely uncontrolled, open traffic conditions, we tackle the second fundamental aspect for real-life application: the best sampling method and the minimum amount of data required to provide consistent results. In fact, state-of-art technology in car communication uses mobile connectivity to stream data from the car to the server where they are processed, and given the massive volume of the sampled data it is crucial to investigate a lower-bound for this data communication. We compare two methods that involve different spatiotemporal sampling of the data and we study the quality of the clustering with different quantities of analyzed data.

The subsampling of the vectors $\mathbf{w}^{k,u}$ presented in Section III-D is completely random and does not consider any spatial or temporal dimension: in other words, it is an *independent subsampling*. We compare it with a different subsampling

	Features						
	1	2	3	4	5	6	7
BRAKE	2 (0.95, 0.11)	4 (0.99, 0.01)	2 (1.00, 0.00)	5 (1.00, 0.01)	3 (1.00, 0.01)	3 (0.95, 0.05)	2 (0.92, 0.07)
GAS	2 (0.96, 0.06)	2 (1.00, 0.00)	2 (0.93, 0.06)	4 (0.98, 0.03)	2 (1.00, 0.00)	2 (0.99, 0.03)	2 (0.99, 0.03)
R.P.M.	3 (0.99, 0.02)	2 (0.98, 0.05)	2 (0.85, 0.06)	2 (1.00, 0.00)	2 (1.00, 0.00)	6 (0.71, 0.06)	2 (0.92, 0.08)
SPEED	2 (1.00, 0.00)	2 (1.00, 0.02)	3 (0.81, 0.12)	2 (0.98, 0.05)	2 (0.93, 0.06)	6 (0.72, 0.04)	2 (0.86, 0.09)
S.W.A.	2 (0.98, 0.05)	5 (0.99, 0.02)	4 (0.78, 0.08)	2 (0.99, 0.09)	4 (1.00, 0.00)	2 (0.92, 0.14)	3 (0.97, 0.05)
S.W.M.	3 (1.00, 0.00)	2 (0.96, 0.06)	4 (0.91, 0.05)	2 (1.00, 0.02)	2 (0.92, 0.09)	2 (0.96, 0.06)	2 (1.00, 0.00)
F.ACC.	4 (0.98, 0.05)	6 (0.93, 0.06)	2 (0.88, 0.09)	5 (0.87, 0.07)	2 (0.98, 0.05)	2 (0.82, 0.09)	2 (1.00, 0.00)
L.ACC.	3 (0.99, 0.04)	2 (0.83, 0.09)	2 (0.86, 0.10)	2 (0.92, 0.12)	2 (0.94, 0.08)	2 (0.80, 0.10)	2 (0.97, 0.08)

Table III: Optimal number of clusters for each combination of feature and signal as a result of the cross-validation process described in section III-D. In brackets, the value of mean and standard deviation referred to the optimal value as in Algorithm 1 .

strategy, which we call *contiguous subsampling*, a subsampling conditioned to spatial contiguity defined as follows. Given the vector $\mathbf{w}^{k,u}$ of dimension d , a random number $r \in \mathbb{N}$ is extracted uniformly in the interval $[1, d]$. Setting $l = \lfloor pd \rfloor$, where $p \in (0, 1)$ is the percentage of the elements to be subsampled, the vector $\mathbf{w}_S^{k,u}$ is constructed considering the elements of $\mathbf{w}^{k,u}$ with indexes from r to $(r + l) \bmod d$. In other words, the vector is subsampled taking, starting from a random element, its l consecutive elements, considering the vector with a circular structure.

For each of the two subsampling strategies defined, we propose an analysis that compares the clusterizations generated in two different ways: in the first, drivers are clustered upon all the data in the dataset, i.e. data coming from all the roads they have driven on; in the second, drivers are clustered upon only a portion of the data acquired. In this way, the first clustering can be considered somehow as a ground truth (being the result of all the data available to us), while the second is the result of a partial subsampling.

Figure 4C reports the results of the V-measure comparisons of the clusterings generated using all the data in the database with the clusterings generated by a subset of the data, for different sizes of subsets and for the two aforementioned subsampling methods. Every subsampling has been repeated 40 times with different random numbers and the K -means clusterings have been performed for each feature with the optimal value of K found earlier.

Results clearly show that the independent subsampling strategy performs better than the contiguous one, and for some features and signals it is possible to reduce the original dataset by a factor of 100 without impairing clustering performance. A comprehensive chart for all the combinations of signals and features can be found in Figure 6 in the Appendix.

V. CONCLUSIONS

In this paper, the problem of driving behavior analysis has been studied from a new point of view, that bridges the gap between driving behavior studies through uncontrolled experiments – leveraging only the GPS signal – and studies exploiting CAN bus data through very controlled experiments. This work proposes a methodology for delineating similarities among drivers using data collected in a completely uncontrolled experiment, through a clustering algorithm performed on seven different features of eight signals recorded by CAN

bus sensors, with a distributional approach. Moreover, it has been shown that, by properly choosing the subsampling strategy, it is possible to reduce the size of the dataset of as much as 99% without impairing clustering performance.

A. Discussion

Given the almost ontological question of what driver behavior is, this work attempts to define it through a data-driven approach. Without any external knowledge (ground truth), though, it is unclear how to define the boundary between the performance of the proposed method and the fuzziness and the unpredictability of human behavior. However, the promising results obtained in this study suggest that the present approach could be considered as a methodology for testing new signals, features and clustering methods which, coupled with additional field knowledge, may lead to pragmatic interpretations of the different clusters in terms of physical and behavioral characterization of driving styles.

It is important also to outline some limitations of this work: the number of users, 64 later reduced to 53 for data homogeneity reasons, likely does not offer a rich enough variety of driving behaviors to enable a comprehensive identification of common attitudes and outliers. Finally, an aspect that needs further investigation is the interaction of the different indicators and the signals directly in the clustering process.

B. Applications and future work

This paper projects the problem of driving behavior characterization using CAN bus technology from a research-oriented approach into an application-oriented technology that opens the way to wide scale and real-time implementations. In fact, as mentioned, the presence of the CAN bus data in almost every car could scale-up any possible application in a very broad and cost-effective way.

Car insurance companies, for example, are interested in assessing the risk of accidents for each user based on real data coming from their driving sessions. Users segmentation in fact, to the best of our knowledge, today is only performed – besides the accidents history – on general information like the geographical location, distance traveled, and velocity. More sophisticated concepts like “aggressiveness” or “nervousness” could be fully characterized. However, in order to do so, further studies have to be performed, comparing the insurance

companies drivers' profiles with the clustering obtained in this work, allowing their characterization based on a ground truth.

Another application is driver recognition, aiming to recognize a driver only upon the CAN bus data. This driver "fingerprint", already studied [28] but never tested in an uncontrolled experimental scenario, could let the car itself to identify the driver for security reasons or adapting settings for comfort or efficiency optimization.

Finally, integration of this modeling technique with physical detection technologies including sonar devices, stereo cameras, lasers and radar would allow to better understand and model driver behaviors, to improve the development of self driving cars and to have safer road networks.

Privacy disclaimer. The data reported herein was collected during experiments performed with drivers who were hired and were explicitly informed of the data collection process. In case the presented methodology should be used with consumer vehicles, it is fundamental to properly inform the customer about usage of data and the purpose of the collection. This needs to be done in order to comply with data privacy laws and regulations, but also to support customers' awareness and self-determination – especially in cases where the realization of an application requires providing personal data to third parties. It is the decision of the customer based on a declaration of consent, if personal data may be collected and for which purpose it may be used.

REFERENCES

- [1] "Driver Identification Using Automobile Sensor Data from a Single Turn," *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 953–958, 2016.
- [2] E. Massaro, C. Ahn, C. Ratti, P. Santi, R. Stahlmann, A. Lamprecht, M. Roehder, and M. Huber, "The Car as an Ambient Sensing Platform," *Proceedings of the IEEE*, vol. 105, no. 01, pp. 3–7, 2017.
- [3] U. Kiencke, S. Dais, and M. Litschel, "Automotive Serial Controller Area Network," *SAE Technical Paper*, 1986.
- [4] O. Moll, A. Zalewski, S. Pillai, S. Madden, M. Stonebraker, and V. Gadepally, "Exploring big volume sensor data with Vroom." [Online]. Available: <http://people.csail.mit.edu/spillai/vroom/vroom-proposal.pdf>
- [5] W. Wang, J. Xi, and H. Chen, "Modeling and recognizing driver behavior based on driving data: A survey," *Mathematical Problems in Engineering*, vol. 2014, 2014.
- [6] J. Grengs, X. Wang, and L. Kostyniuk, "Using GPS Data to Understand Driving Behavior," *Journal of Urban Technology*, vol. 15, no. 2, pp. 33–53, 2008.
- [7] J. Paefgen, F. Michahelles, and T. Staake, "GPS trajectory feature extraction for driver risk profiling," *Proceedings of the 2011 international workshop on Trajectory data mining and analysis*, pp. 53–56, 2011.
- [8] A. Pentland and A. Liu, "Modeling and Prediction of Human Behavior," *Neural Computation*, vol. 11, no. 1, pp. 229–242, 1999.
- [9] J. M. McNew, "Predicting cruising speed through data-driven driver modeling," *IEEE Conference on Intelligent Transportation Systems, Proceedings, ITSC*, pp. 1789–1796, 2012.
- [10] Y. Zhang, W. C. Lin, and Y. K. S. Chin, "A pattern-recognition approach for driving skill characterization," *IEEE Transactions on Intelligent Transportation Systems*, vol. 11, no. 4, pp. 905–916, 2010.
- [11] "Driver behavior recognition and prediction in a SmartCar," *Proceedings of SPIE*, vol. 4023, pp. 280–290, 2000.
- [12] J. Carmona, F. García, D. Martín, A. Escalera, and J. Armingol, "Data Fusion for Driver Behaviour Analysis," *Sensors*, vol. 15, no. 10, pp. 25 968–25 991, 2015.
- [13] "Analysis and classification of driver behavior using in-vehicle canbus information," *Biennial Workshop on DSP for In-Vehicle and Mobile Systems*, no. October 2015, pp. 17–19, 2007.
- [14] N. Li, J. J. Jain, and C. Busso, "Modeling of driver behavior in real world scenarios using multiple noninvasive sensors," *IEEE Transactions on Multimedia*, vol. 15, no. 5, pp. 1213–1225, 2013.
- [15] T. Wakita, K. Ozawa, C. Miyajima, and K. Takeda, "Parametric Versus Non-parametric Models of Driving Behavior Signals for Driver Identification," in *Audio- and Video-Based Biometric Person Authentication*, T. Kanade, A. Jain, and N. K. Ratha, Eds. Springer Berlin Heidelberg, 2005, pp. 739–747.
- [16] C. Miyajima, Y. Nishiwaki, K. Ozawa, T. Wakita, K. Itou, K. Takeda, and F. Itakura, "Driver modeling based on driving behavior and its evaluation in driver identification," *Proceedings of the IEEE*, vol. 95, no. 2, pp. 427–437, 2007.
- [17] C. Miyajima, Y. Nishiwaki, K. Ozawa, T. Wakita, K. Itou, and K. Takeda, "Cepstral Analysis of Driving Behavioral Signals for Driver Identification," *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*, vol. 5, no. October 2015, pp. 6–9, 2006.
- [18] N. Kawaguchi, K. Takeda, and F. Itakura, "Multimedia Corpus of In-Car Speech Communication," *The Journal of VLSI Signal Processing-Systems for Signal, Image, and Video Technology*, vol. 36, no. 2/3, pp. 153–159, 2004.
- [19] "Driver classification and driving style recognition using inertial sensors," *IEEE Intelligent Vehicles Symposium, Proceedings*, no. Iv, pp. 1040–1045, 2013.
- [20] D. A. Johnson and M. M. Trivedi, "Driving style recognition using a smartphone as a sensor platform," *IEEE Conference on Intelligent Transportation Systems, Proceedings, ITSC*, pp. 1609–1615, 2011.
- [21] J. Dai, J. Teng, X. Bai, Z. Shen, and D. Xuan, "Mobile phone based drunk driving detection," in *2010 4th International Conference on Pervasive Computing Technologies for Healthcare*. IEEE, 2010, pp. 1–8.
- [22] P. Angkitittrakul, J. H. Hansen, S. Choi, T. Creek, J. Hayes, J. Kim, D. Kwak, L. T. Noecker, and A. Phan, "UTDrive: The Smart Vehicle Project," in *In-Vehicle Corpus and Signal Processing for Driver Behavior*, H. Takeda, K. Erdogan, H., Hansen, J., Abut, Ed. Boston, MA: Springer US, 2009, pp. 55–67.
- [23] K. Takeda, S. Member, J. H. L. Hansen, P. Boyraz, C. Miyajima, H. Abut, and L. S. Member, "International Large-Scale Vehicle Corpora for Research on Driver Behavior on the Road," vol. 12, no. 4, pp. 1609–1623, 2011.
- [24] G. A. M. Meiring and H. C. Myburgh, "A review of intelligent driving style analysis systems and related artificial intelligence algorithms," *Sensors (Switzerland)*, vol. 15, no. 12, pp. 30 653–30 682, 2015.
- [25] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, 2nd ed. Springer, 2009.
- [26] P. J. Rousseeuw, "Silhouettes : a graphical aid to the interpretation and validation of cluster analysis," vol. 20, pp. 53–65, 1987.
- [27] A. Rosenberg and J. Hirschberg, "V-measure: A conditional entropy-based external cluster evaluation measure," *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pp. 410–420, 2007.
- [28] M. Enev, A. Takakuwa, K. Koscher, and T. Kohno, "Automobile Driver Fingerprinting," *Proceedings on Privacy Enhancing Technologies*, vol. 2016, no. 1, pp. 34–50, 2016.

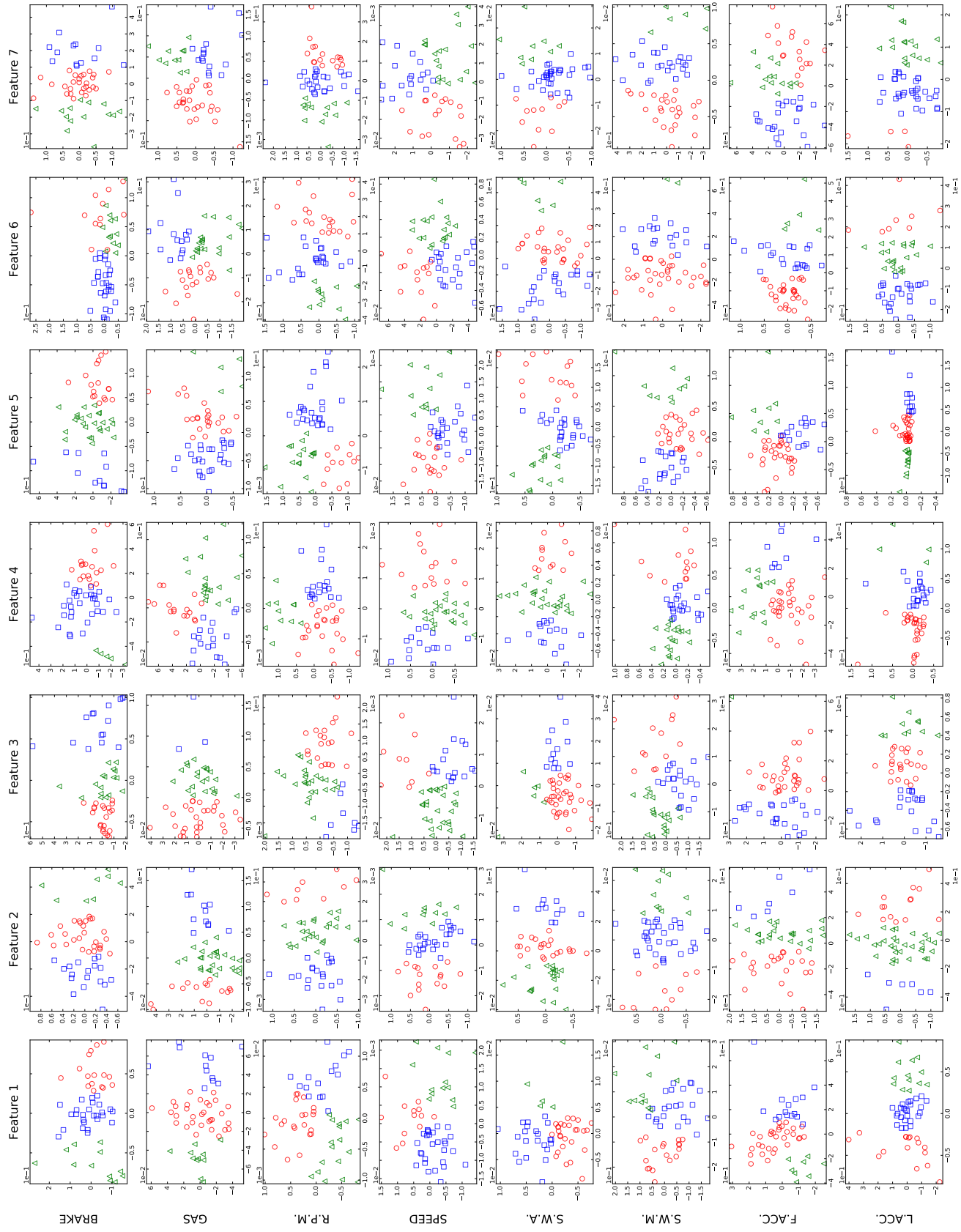


Figure 5: Drivers clusterings for different signals and features. The K -means algorithm has been run on all data in the database and for the optimal values of K as in Table III

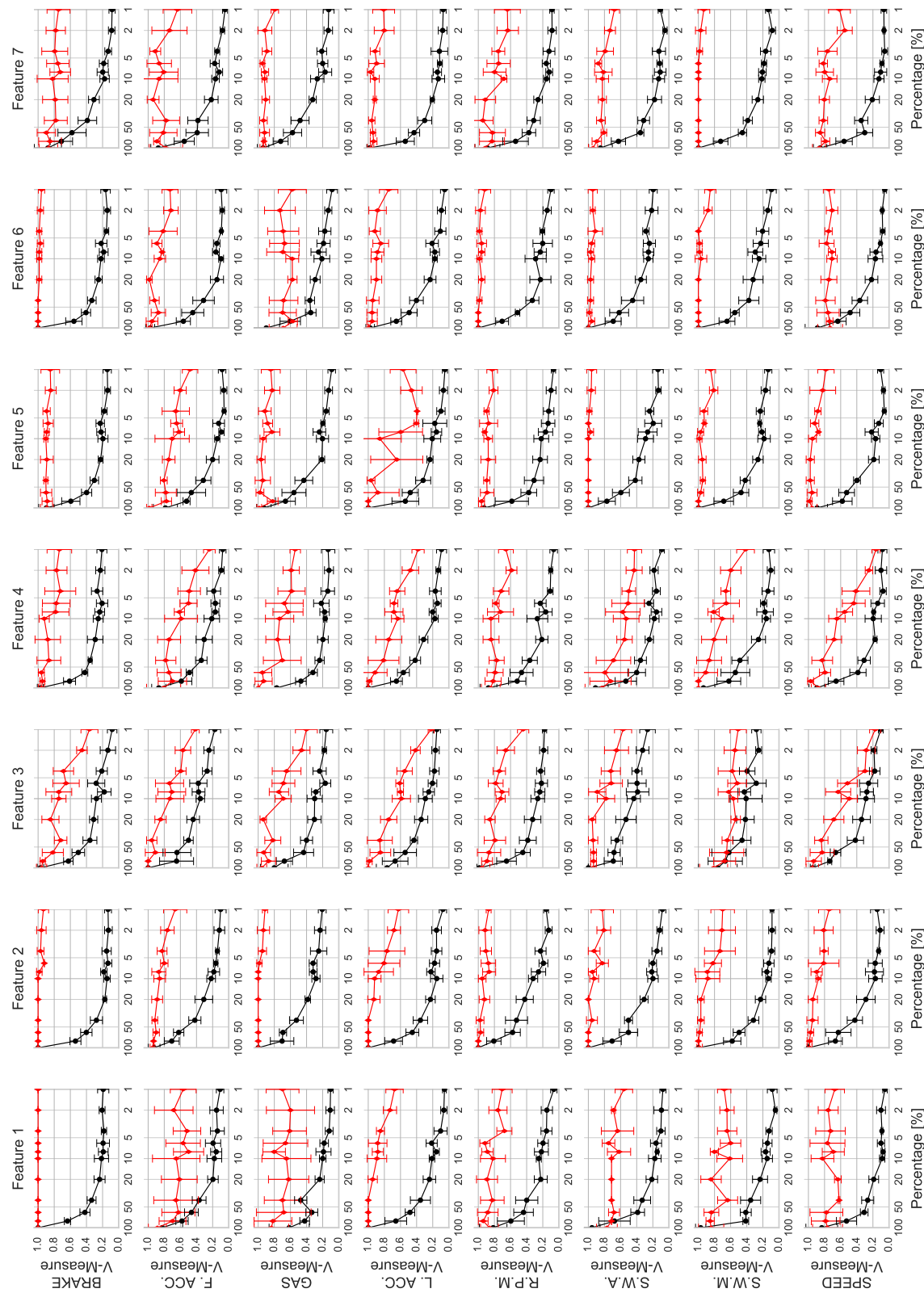


Figure 6: Comparison of different subsampling methods: *independent subsampling* (red line, diamonds) and *contiguous subsampling* (black line, circles). V-measures of the comparisons of the K -means clusters generated using all the data in the database, with the clusters generated by a subset of the data (validation set), for different sizes of the validation set (100%, 50%, 20% 10%, 5%, 2%, 1% of the original data). The clusterings use the optimal values of K as in Table III.

Umberto Fugigliando is a Research Fellow at MIT Senseable City Lab. He received his Bachelor degree (2013) and Master degree (2016) in Applied Mathematics from Politecnico di Torino (Italy), with a thesis on driving behavior. He is also a ASP Alta Scuola Politecnica fellow and he has spent a semester at KTH Royal Institute of Technology in Stockholm (Sweden). His research interests are in the area of digital technology and data science with applications to mobility, acoustics and human behavior characterization.

Paolo Santi is Research Scientist at MIT Senseable City Lab where he leads the MIT/Fraunhofer Ambient Mobility initiative, and a Senior Research at the Istituto di Informatica e Telematica, CNR, Pisa. Dr. Santi holds a "Laurea" degree and PhD in computer science from the University of Pisa, Italy. Dr. Santi is a member of the IEEE Computer Society and has recently been recognized as Distinguished Scientist by the Association for Computing Machinery. His research interest is in the modeling and analysis of complex systems ranging from wireless multi hop networks to sensor and vehicular networks and, more recently, smart mobility and intelligent transportation systems. In these fields, he has contributed more than 120 scientific papers and two books. Dr. Santi has been involved in the technical and organizing committee of several conferences in the field, and he is/has been an Associate Editor of the IEEE Transactions on Mobile Computing, the IEEE Transactions on Parallel and Distributed Systems, and Computer Networks. Dr. Santi was Guest Editor of the Proceeding of the IEEE special issue on Vehicular Communications: Ubiquitous Networks for Sustainable Mobility in 2011, to which he also contributed a paper.

Emanuele Massaro, PhD is a Postdoctoral Research Fellow at the MIT Senseable City Lab. He received both his Bachelor (2006) and his Master (2009) in Environmental Engineering from the University of Florence (Italy). He then received his PhD in Complex Systems and Nonlinear Dynamics in 2014 from the Department of Information Engineering and Department of Physics and Astronomy at the University of Florence. He came to the United States in March 2014 to conduct his postdoctoral research where he worked for one year as Postdoctoral Associate at the Department of Civil and Environmental Engineering Carnegie Mellon University and also as a contractor for the Risk and Decision Science Team of US Army Corps of Engineer. He joined the Massachusetts Institute of Technology in March 2015. His broad research interests are in the areas of socio-technical systems and computational social science: he aims to understand the theory of, and quantify the interplay among physical infrastructures, information, and human (societal) activities.

Sebastiano Milardo received his Bachelor degree in 2011 and his Master degree in 2013, both in Computer Engineering from the University of Catania. From January 2014 to April 2015 he worked in the Italian National Consortium of Telecommunications (CNIT), as Researcher within the NEWCOM# and SIGMA Projects. Since 2015 he is currently a Ph.D. student in Information and Communication Technologies at the University of Palermo. His research interests include Software Defined Networking, Sensor Networks, network protocols for the Internet of Things and Big Data analysis.

Kacem Abida, PhD, is a senior engineer at the Volkswagen Group of America Electronics Research Lab (ERL). Dr. Abida is currently leading the big data projects at ERL. He holds a PhD degree in Electrical and Computer Engineering from the University of Waterloo, Canada. His areas of interest include speech and natural language technologies, as well as machine learning based big data analytics.

Rainer Stahlmann received his diploma in electrical engineering and computer science from University of Applied Sciences Ingolstadt, Germany, in 2009. Since then he has been working for AUDI AG in Ingolstadt, Germany, where he is currently in the Department of Data Strategy and Analytic Services. In cooperation with the Chair for Computer Networks and Communication Systems at University of Erlangen, Germany, he is working toward his Ph.D. degree. His research is focused on vehicular data processing and analytics as well as on technical evaluation of V2X communication systems.

Dr. Ing. **Florian Netter** received his diploma in mechanical engineering from the technical university of Munich, Germany, in 2010. Since then he has been working for AUDI AG in Ingolstadt, Germany, where he received his Ph.D. degree in cooperation with the Karlsruhe Institute of Technology,

Germany, in 2015. During his research he focused on complexity adaptation of simulation models in entire system simulations to identify quantification attributes for a high goodness of fit and in spite of increasing computational power still maintaining a short simulation period. Currently he is working at AUDI AG in the Department for Platform Development and Data Analytics taking care of vehicular data stream processing in cloud computing environments.

Carlo Ratti is the founder and Director of the MIT Senseable City Lab. An architect and engineer by training, Carlo Ratti practices in Italy and teaches at the Massachusetts Institute of Technology. He graduated from the Politecnico di Torino and the cole Nationale des Ponts et Chaussées in Paris, and later earned his MPhil and PhD at the University of Cambridge, UK. Ratti has co-authored over 200 publications and holds several patents. His work has been exhibited worldwide at venues such as the Venice Biennale, the Design Museum Barcelona, the Science Museum in London, GAFTA in San Francisco and The Museum of Modern Art in New York. His Digital Water Pavilion at the 2008 World Expo was hailed by Time Magazine as one of the Best Inventions of the Year. He has been included in Esquire Magazine Best and Brightest list, in Blueprint Magazine 25 People Who Will Change the World of Design, and in Forbes Magazine Names You Need To Know in 2011. Ratti was a presenter at TED 2011 and is serving as a member of the World Economic Forum Global Agenda Council for Urban Management. He is a regular contributor to the architecture magazine Domus and the Italian newspaper Il Sole 24 Ore. He has also written as an op-ed contributor for BBC, La Stampa, Scientific American and The New York Times.