

Multi-Level Driver Workload Prediction Using Machine Learning and Off-The-Shelf Sensors

van Gent, Paul; Melman, T.; Farah, Haneen; van Nes, Nicole; van Arem, Bart

Publication date

2018

Document Version

Accepted author manuscript

Published in

Transportation Research Board Conference Proceedings 2018

Citation (APA)

van Gent, P., Melman, T., Farah, H., van Nes, N., & van Arem, B. (2018). Multi-Level Driver Workload Prediction Using Machine Learning and Off-The-Shelf Sensors. In *Transportation Research Board Conference Proceedings 2018* Transportation Research Board (TRB).

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Multi-Level Driver Workload Prediction Using Machine Learning and Off-The-Shelf Sensors

Paul van Gent

Delft University of Technology
Faculty of Civil Engineering and Geosciences
Stevinweg 1, 2628CN Delft, the Netherlands
Email: P.vanGent@tudelft.nl

Timo Melman

Delft University of Technology
Faculty of Mechanical, Maritime and Materials Engineering
Mekelweg 2, 2628CD Delft, the Netherlands
Email: T.Melman@tudelft.nl

Haneen Farah

Delft University of Technology
Faculty of Civil Engineering and Geosciences
Stevinweg 1, 2628CN Delft, the Netherlands
Email: H.Farah@tudelft.nl

Nicole van Nes

SWOV – Stichting Wetenschappelijk Onderzoek Verkeersveiligheid
Bezuidenhoutseweg 62, 2594AW The Hague, the Netherlands
Email: Nicole.van.Nes@swov.nl

Bart van Arem

Delft University of Technology
Faculty of Civil Engineering and Geosciences
Stevinweg 1, 2628CN Delft, the Netherlands
Email: B.vanArem@tudelft.nl

Word count: 6472 words text + 4 tables/figures x 250 words (each) = 7472 words

Revised Paper

Submission Date
31-07-2017

ABSTRACT

The present study aims to add to the literature on driver workload prediction using machine learning methods. The main aim is to develop workload prediction on a multi-level basis, rather than a binary high/low distinction as often found in literature. The presented approach relies on measures that can be obtained unobtrusively in the driving environment with off-the-shelf sensors, and on machine learning methods that can be implemented on low-power embedded systems.

Two simulator studies were performed, one inducing workload using realistic driving conditions, and one inducing workload with a relatively demanding lane-keeping task. Individual and group-based machine learning models were trained on both datasets and evaluated. For the group-based models the generalising capability, that is the performance when predicting data from previously unseen individuals, was also assessed.

Results show that multi-level workload prediction on the individual and group level works well, achieving high correct rates and accuracy scores. Generalising between individuals proved difficult using realistic driving conditions, but worked well in the high demanding lane-keeping task. Reasons for this discrepancy are discussed as well as future research directions.

Keywords: Driver workload, machine learning, workload prediction, random forest, support vector machine, embedded workload prediction

1 INTRODUCTION

2 Research into driver workload has been conducted for at least three decades (1, 2). Recently,
3 research efforts have shifted to using powerful Machine Learning (ML) methods, giving
4 promising results (5, 6). ML methods have been used for other driver-related classification
5 problems, such as driver distraction (7), driver interruptibility (8) or driver identification (9). The
6 present study aims to fill the gaps in the existing research on predicting driver workload using ML
7 methods in several ways, as will be explained in the next paragraphs.

8 First, ML studies into predicting driver workload often focus on a binary classification
9 problem (high workload vs. low workload). A more fine-grained prediction of workload may be
10 desirable to enable adaptive interfaces for in-vehicle advice systems (IVIS), systems that may
11 simplify their content (10), or driver assistance systems that may incrementally increase their level
12 of support based on the level of driver workload. The experiments described in this paper attempt
13 to predict workload on 7- and 10-point workload scales.

14 Second, studies to date often use intrusive sensors or measure variables (i.e.
15 electroencephalogram, EEG) that are not practical in the driving environment (see for example (5,
16 6)). Additionally, it is unknown how well results obtained by the high-grade intrusive sensors used
17 in experiments translate to low-cost sensors. This work uses low-cost sensors that can be
18 integrated into the real-world driving environment, and uses measures that can be obtained
19 non-intrusively. This is important, since especially low-cost sensors are likely to be integrated into
20 the driving environment in real-world applications.

21 Lastly, the models generated in most studies are not generally publicly available for use
22 by the research community. The models developed in this study will be made available for
23 scientific use after publication of results (<https://github.com/paulvangentcom>).
24

25 Research Objectives

26 The previous section outlined the main research gaps and ways to add to the present literature. This
27 led to the formulation of three criteria for predicting driver workload in the present work: The main
28 goal is to develop a workload algorithm that (A) has usable accuracy when predicting multiple
29 workload levels, while generalising among individuals, (B) uses data that can be measured with
30 available low-cost, sensors that can be integrated into the driving environment, and (C) is
31 implementable on embedded hardware (for example in a smart steering wheel).

32 The first criterion (A), predicting workload at a higher resolution than the binary low/high
33 found in previous literature while generalising among individuals, is addressed in the experimental
34 design and data analysis presented in subsequent sections.

35 The second criterion (B) entails using sensor inputs from readily available, low-cost
36 sensors that are easy to implement in the driving environment. By using low-cost sensors, which
37 are likely to present more noise in the signal compared to high-end sensors, results will give a
38 better reflection of real-world performance compared to studies using high-end sensors. Apart
39 from having been used successfully in other workload prediction studies, selected variables should
40 be measurable non-intrusively in the driving environment. This led to the selection of heart rate,
41 skin response, blink rate and several performance measures (for an overview of the selection
42 process, see (11)). This criterion ensures any results are directly applicable to in-car settings at a
43 low cost, and that results obtained are likely to translate well to real-world applications.

44 Criterion C, ensuring the model is implementable on an embedded system, means it must
45 be efficient both in memory use as well as computational requirements. Two machine learning
46 algorithms were selected that can satisfy this criterion: ‘Random Forest’ and ‘Support Vector
47 Machine’ algorithms. Random Forests (12) are computationally efficient (13) but can have a large

memory footprint. Solutions have been proposed that allow embedded implementations while maintaining performance (14), making it a suitable algorithm to use. Support Vector Machines (15) implementations can suffer from computational complexity, as well as high memory footprint for more complex models. Methods have been proposed, however, that achieve remarkable efficiency increases without sacrificing performance (16, 17), making SVM's also a suitable candidate algorithm.

Two experiments were conducted to evaluate the feasibility of the previously defined criteria. First, a simulator experiment was performed, where workload was induced using realistic driving situations. Results of this experiment were explored further using a dataset obtained from another driving simulator experiment that induced workload with a demanding lane-keeping task. Finally, results of both experiments are discussed and future steps are outlined.

ESTIMATING WORKLOAD IN A REALISTIC DRIVING SCENARIO STUDY

To assess the feasibility of predicting driver workload in realistic driving settings, a simulator study was performed. The main goal was to evaluate the prediction of multi-level driver workload in realistic driving conditions.

Methods

Equipment

The study was performed in a fixed-base, medium-fidelity driving simulator. A dashboard mockup with three 4K-displays (resolution 4096*2160 px) provided roughly 180-degree vision. Actuators consisted of a Fanatec steering wheel and pedals, and a custom blinker control. The simulation ran in Unity3D. The simulated vehicle had an automatic gearbox and a top speed of 165 km/h.

FIGURE 1(A) illustrates the set-up.

Physiological data were recorded at 100Hz, using low-cost sensors powered by an Atmel ATmega328p embedded processor board. Heart rate was recorded using a photoplethysmographic (PPG) method (18) at the left index finger. Skin response was recorded at the middle and ring finger of the same hand (see FIGURE 1(B)). Additionally, blink data were recorded using a GoPro HERO+ camera on the dashboard, running at 1080p@30Hz. Simulator data were logged at 50Hz.

Simulator Scenarios

Two scenarios were created in Unity3D, one scenario with situations likely to induce high workload ('high' workload' scenario) and one with situations that are not likely to induce high workload ('low workload' scenario). Road geometry was based on a part of the Cooperative-ITS (C-ITS) corridor in the Netherlands: the A67, a two-lane highway between Eindhoven and Venlo with speed limit of 130km/h. Three weather conditions were designed for each scenario: clear weather, and two degrees of fog with visibility of approx. 150 meters ('light fog') and below 25 meters ('heavy fog'). This gave a total of six scenarios.

To accurately design the road geometry, CAD drawings of the road segments were secured from the open data program of the Dutch government (<https://data.overheid.nl>). Using Autodesk 3DS Max, the data in the CAD files were converted to 3D models and textured. The surrounding terrain was generated using height map data obtained from the Microsoft Bing Maps API (<https://www.bingmapsportal.com/>). Canals and wooded areas were extracted automatically from satellite imagery, and adjusted by hand where necessary. The location, shape, and content of traffic signs was inferred from Google Streetview, designed in 3DS Max and manually placed at the corresponding location in the scenario.

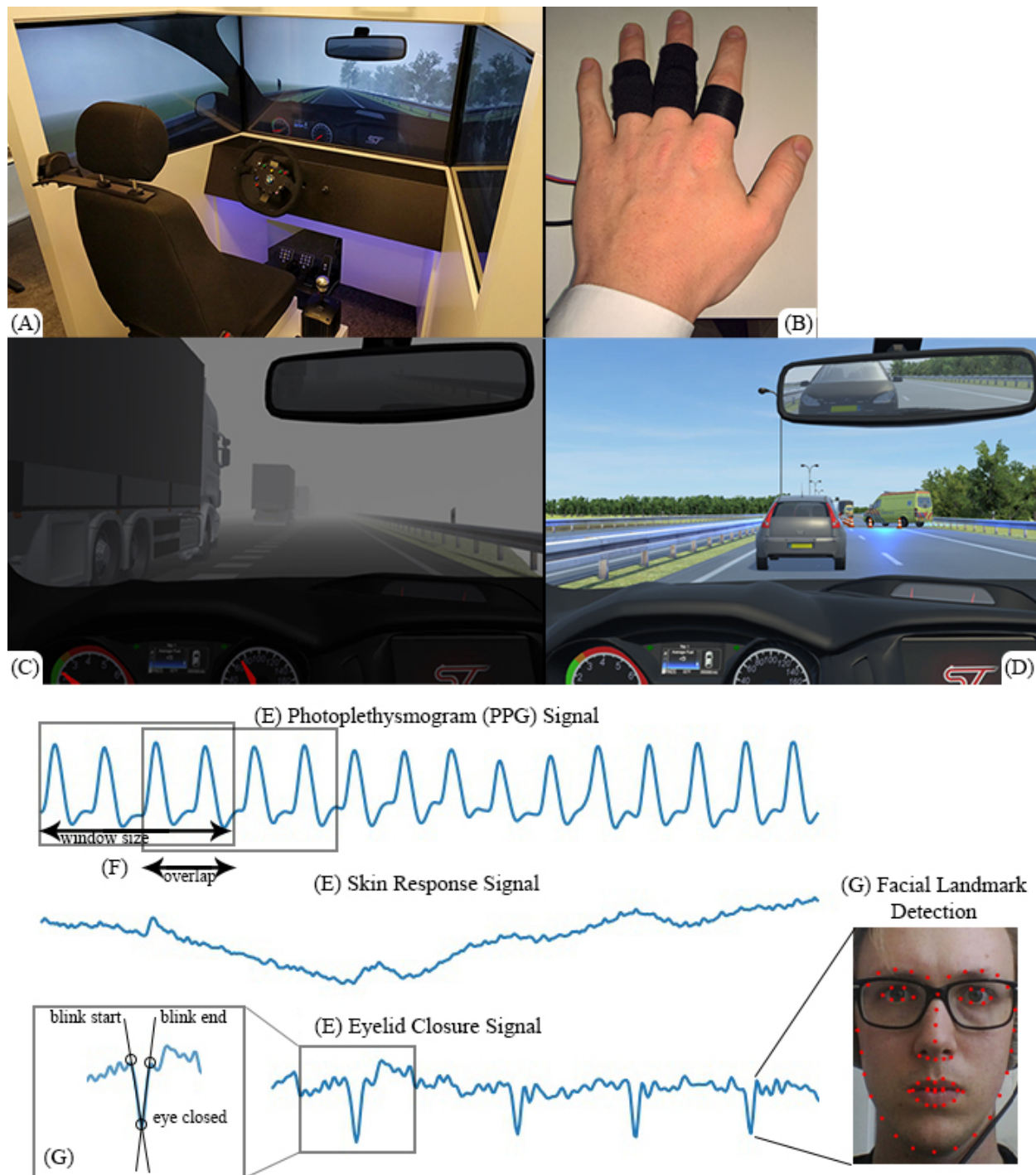


FIGURE 1 - Figure showing the simulator set-up (A), physiological sensors (B), the merging between a platoon of trucks in dense fog (C) and the accident site at the end of the 'high workload' scenario (D). Examples of the raw signal data are shown (E), the concepts of window size and overlap factor (F), an example of the facial landmark detection and the resulting process of analysing the blink rate signal (G).

The 'high workload' scenario was 15.9 km in length, and ran between Eindhoven and Someren. Participants would encounter several workload-inducing 'events' spread out across the scenario. After accelerating across an on-ramp, the first event was encountered: participants had to

merge into a dense platoon of trucks (4-5 meters headway, FIGURE 1(C)), a manoeuvre shown to increase workload on the driver (19). The second event was encountered two kilometres upstream, and consisted of a segment of slow moving traffic on the right lane, designed to nudge the participants to drive in the left lane. While passing the slow-moving traffic, an ambulance approached from behind exhibiting auditory and visual signals, travelling at the legally allowed max speed of 170km/h in the Netherlands (max. 40km/h difference with other traffic). This placed the participant in the demanding situation of quickly having to find a gap in the much slower moving lane to the right and perform a merging manoeuvre. The third event was a game of '20 questions' (20), intended to simulate an engaging (phone) conversation. By asking at most 20 polar (yes/no) questions, participants had to guess which animal, object or person the experimenter had in mind. The final event came near the end of the scenario. The right lane was closed off due to an accident, with slow moving (< 15 km/h) traffic on the left lane (FIGURE 1 (D)). The 20 questions game was played until the accident site was reached. If participants finished early, the game was restarted with a different subject. After this, participants took the next exit and stopped the car.

The 'low workload' scenario consisted of self-paced driving in light traffic for 20.5km. The simulated road was a replica of the A67 road between Someren and Venlo. There were no events. Participants drove until reaching a designated exit, where they stopped the car.

Experimental Procedure

Approval for the study was obtained from the ethics committee at Delft University of Technology. Participants drove the six scenarios spread out over three separate days, each day driving one randomly assigned 'high workload' and one 'low workload' scenario. This approach was taken because physiological measures can vary from day to day, as well as to avoid a fatigue effect from occurring when asking participants to drive six 10-15-minute scenarios consecutively.

In the 'high workload' scenario, participants were asked to rate their experienced mental effort and task difficulty on a 7-point scale after each event, leading to six workload data points per run. In the 'low workload' scenario, the questions were asked at fixed positions in the scenario, leading to four workload data points per run. The exact questions were '*How much mental effort did the driving task take in the last few moments, on a scale of 1-7?*' and '*How difficult was the driving task in the last few moments, on a scale of 1-7?*'. Scale labels ranged from very low/easy, to very high/difficult, and were explained to participants before the experiment started. Note that we did not use a standardised workload scale such as the NASA TLX or RSME, since we wanted to keep interaction time with and demands on the driver to a minimum.

Participants that registered for the experiment received a copy of the informed consent. It was signed and brought to the first session. After being seated in the simulator, a relaxation period of three minutes was given to the participants. This was to allow the physiological measures of each participant to return to its baseline. Sensors were attached, after which the signal quality was checked. A physiological baseline was recorded first. After the baseline, it was briefly explained to the participant that there would follow a drive on a segment of the A67 highway. Participants were instructed to drive at their own pace, but not exceed the speed limit as indicated on road-side signs. If a participant was unfamiliar with '20 questions', a test round was played to familiarise them with the game.

Data Analysis

Participants were asked to rate their mental effort and driving task difficulty on a 7-point scale. Since querying the driver might influence workload, the 'high workload' scenario was constructed in such a way that at least one minute of driving was between each two events, to allow signals to

return to baseline. The data recorded between two events were not used in the analysis. In the case of the 'low workload' scenario, one minute of data following each question were excluded from the analysis.

Preprocessing of Physiological Data

An algorithm was developed to extract the most commonly used features from the measured heart rate signal (21, 22), using a sliding window approach (see FIGURE 1F). The output measures are divided into time-domain (23) and frequency-domain measures (24). In the time-domain, the measures included are BPM (beats per minute), IBI (inter-beat interval), MAD (median absolute deviation of intervals between heart beats), SDNN (standard deviation of intervals between heart beats), RMSSD (root mean square of successive differences between neighbouring heart beat intervals), SDSD (standard deviation of successive differences between neighbouring heart beat intervals), and the pNN50 and pNN20 (proportion of differences between successive heart beats greater than 50ms and 20ms, resp.) In the frequency domain, included measures are LF (the low frequency band: 0,04-0,15Hz), which is related to short-term blood pressure variation, and HF (the high frequency band: 0,16-0,5Hz), which reflects breathing rate, and the LF/HF ratio, a measure of sympathetic-parasympathetic balance (24, 25).

Skin response consists of a tonic and phasic component (26). Tonic represents the long-term, slow variation in the signal, indicative of general psycho-physiological arousal (27). Phasic reflects relatively quick responses to discrete external stimuli, occurring generally between 1-3 seconds after stimulus onset (27). Power in the frequency spectrum of skin response between 0.03Hz-0.5Hz has been linked to short term workload changes (28). The mean, max-min difference, MAD (median absolute difference), and 0.03-0.5Hz frequency spectrum were extracted from the GSR signal, using the same window approach as for heart rate. Frequency spectra were extracted using a trapezoidal integration of the area under corresponding frequency bands in the power spectrum.

Blink data were detected offline from recorded video data. An algorithm was developed to extract blink number, blink duration and inter-blink-interval. It functioned by detecting 68 'facial landmarks' (29), then calculating eyelid distance for each frame. Blinks were detected in the resulting signal by finding large slopes, then finding the lowest point of reversal. The process is displayed visually in FIGURE 1 (G).

Driver Performance Data

Performance measures reflect how the control the driver exerts over the vehicle varies across conditions. We included steering wheel angle, steering wheel reversals, speed, variation in lateral and longitudinal position, and headway and time to collision when available (for more information, see (11)).

Generating Machine Learning Sets

Machine learning sets were generated from the raw data and labelled based on self-report data, by varying window size and overlap factor. Window size refers to how much data is used for the calculation of features, overlap factor refers to how much data any window W_i shares with the previous window W_{i-1} . Both concepts are visualised in FIGURE 1 (F). Window sizes of 5, 10 and 30 seconds, and overlap factors of 0% and 50% were used, leading to a total of 6 sets.

Model Development and Evaluation

Two different machine learning algorithms were used: A Random Forest Regressor (RFR), and a Support Vector Machines Regressor (SVR). The RFR creates an ensemble (forest) of regression trees in which each tree is trained on a random subset of the features. They have been used in for example (30). Support Vector Machines function by mapping the data to a higher dimensional space, and solving an optimization problem to identify a set of hyperplanes that separate the training data into classes. They have been used in for example (7, 9). With the SVR, the Polynomial kernel (SVR(poly)), and the Radial Basis Function kernel (SVR(rbf)) were evaluated. Algorithms that were used are taken from the SciKit-Learn repository (31).

The resulting models were evaluated using several metrics. Model error was evaluated using mean absolute error (AE_{μ}) and median absolute error ($AE_{\mu/2}$), both measures of the accuracy of the predictions. The coefficient of determination (R^2) was also computed as a goodness-of-fit measure. Performance for class-based predictions was also evaluated, expressed as correct rate.

Results

Participants

19 participants took part in the experiment. Data from one participant were excluded because of a failure to understand some tasks due to a language barrier. This left 18 participants, of which 12 were males and 6 were females. The average age was 34.56 years (SD 10.09). Of the 18 participants, 12 owned a car and reported using it three to four times a week on average, and travelling between 2500 and 15000km annually. All participants held a valid driver's license. No simulator sickness severe enough to terminate a driving session was reported. Reported mental effort and perceived difficulty correlated with weather conditions and with scenario type independently and in line with expectations, although no interaction effect was present (11).

Individual Models

The training and testing sets for the individual models were generated by dividing the dataset of each driver into training and testing sets with an 80%/20% split ratio, respectively. This split ratio was chosen to ensure sufficient training data, since individual datasets were relatively small.

The results indicated that the models functioned well, with the RFR outperforming the SVR. For all individual models with a window size of 5s and overlap of 0%, the AE_{μ} was 0.343, the $AE_{\mu/2}$ was 0.129, R^2 was 0.679, Correct Rate (CR) was 76.30% when predicting discrete classes, and 93.80% when miss-by-one errors were allowed (CR+/-1). This indicated that on average, predictions were off by 0.343, and that half the predictions had an error less than 0.129, from a total scale of 7 classes. See TABLE 1 for an overview of all results. Model performance increased with a larger overlap factor. This was expected, since a larger overlap creates a larger training set to fit the model to, and because a larger overlap factor indicates more shared variance between adjacent samples. Interestingly, an inverse relationship between window size and model performance was observed, contrary to what has been reported previously (5). Miss-by-one errors indicate predictions that are 'almost correct', and still contain enough information about the true workload states. For example, if workload is predicted as '6' while the true value is '7', the information in the prediction is still useful: in either case workload is on the high end.

Group Models

The second step was to estimate the model performance within the entire group. The dataset containing data from all drivers was split into training- and testing sets with a 60%/40% split ratio.

Since the size of the group dataset is much larger compared to individual dataset, a more stringent split ratio could be chosen while maintaining a sufficiently large training set.

Results indicated group models performed well. The AE_{μ} for the model with window size 5s and 0% overlap was 0.605, the $AE_{\mu1/2}$ 0.406, R^2 0.661, CR 57.40%, and CR \pm 1 90.60%.

TABLE 1 Performance Metrics RFR models

Window Size	5 sec		10 sec		30 sec	
Overlap Factor	0%	50%	0%	50%	0%	50%
Individual Model						
AE_{μ}	0.343	0.219	0.431	0.280	0.613	0.492
$AE_{\mu1/2}$	0.129	0.565	0.296	0.109	0.490	0.291
R^2	0.679	0.8716	0.590	0.794	0.071	0.306
CR	76.30%	85.21%	67.88%	80.77%	49.68%	60.82%
CR \pm 1	93.80%	97.61%	92.93%	96.13%	85.81%	89.55%
Group Model						
AE_{μ}	0.605	0.455	0.744	0.553	0.898	0.801
$AE_{\mu1/2}$	0.406	0.250	0.565	0.344	0.628	0.652
R^2	0.661	0.774	0.564	0.709	0.372	0.504
CR	57.40%	69.57%	46.12%	62.48%	40.47%	43.82%
CR \pm 1	90.60%	93.81%	87.02%	91.42%	80.60%	84.56%
Generalising Model						
AE_{μ}	1.522	1.536	1.457	1.519	1.375	1.424
$AE_{\mu1/2}$	1.163	1.201	1.199	1.253	1.174	1.230
R^2	-0.538	-0.623	-0.460	-0.602	-0.299	-0.396
CR	20.07%	20.05%	19.81%	20.21%	20.21%	20.47%
CR \pm 1	55.18%	55.19%	55.46%	54.94%	57.21%	55.89%

Performance metrics for the best performing (RFR) classifier. The table displays the mean (μ) and median ($\mu1/2$) absolute error metrics, the coefficient of determination (R^2), the correct rate (CR) and the miss-by-one correct rate (CR \pm 1).

Generalising Group Models

The last step was to assess how models would perform in a realistic setting, e.g. a setting where workload from an unknown driver is predicted based on data from a pool of other drivers. To achieve this, data were sampled using a k-fold approach, with $k = N_{\text{participants}}$. For every k_i , the training set consisted of all data except the held out participant k_i . Workload for participant k_i was then predicted and model performance evaluated. This method simulated how the trained models would perform when predicting data from previously unseen individuals. This obtained performance measure reflects real-world settings, where it is impractical for models to be trained on all possible drivers and generalising power is thus preferable.

Results showed that models did not perform well when generalizing to unknown drivers. The AE_{μ} for all individual models with window size 5s and 0% overlap was 1.522, $AE_{\mu1/2}$ was 1.163, R^2 was -0.538, CR 20.07%, and CR \pm 1 55.18%. The strongly negative coefficient of determination suggests unsatisfactory performance (the mean of the data is a better predictor than the trained model). The relatively low (though above chance level, not satisfactory) absolute error rates given R^2 are explained by a class imbalance in the dataset, where two classes (workload level 1 and 2) dominate. To assess whether this was a possible cause for the poor performance of the models, data were resampled using SMOTE (Synthetic Minority Over-Sampling Technique) (32).

This had little discernible effect on the model performance, and it was concluded that low performance was not due to the class imbalance in the dataset. It was also observed that R^2 increases slightly with increasing window size, in accordance with earlier studies (5) and contrary to the individual and group models in the present study.

Conclusion

The results of this study showed that predicting self-reported workload in a simulated realistic environment was possible at the individual and group level, but proved difficult when generalising to unknown drivers. Several causes can be identified. The simulated scenarios might not have induced sufficient workload to be measurable with performance or physiological measures. Indeed, most participants indicated that driving in the simulator felt very different from actual driving, and was not that difficult at all. Since a self-report measure was used, which is a subjective measure, it is possible that different participants had biased response tendencies. Lastly, it might also be the case that different physiological response patterns to workload exist, in which case the sample size of 18 could have been too small to account for all occurring patterns.

This raises the question whether workload prediction is at all possible on non-binary scales, while generalising across drivers. To further explore this possibility, a dataset from a study with a lane-keeping task was obtained. This study and the results are discussed in the next section.

ESTIMATING WORKLOAD IN A FORCED-PACE SIMULATOR STUDY

A dataset was re-used from a previously executed study by Melman et al. (in press, (33)) to further assess multi-level workload prediction in drivers. The study featured a challenging lane-keeping task, which had the potential to induce higher workload than the previous study. The same physiological and performance measurements were used in as in the previously described simulator study.

Method

Equipment

The study was performed in a fixed-base driving simulator at the faculty of Aerospace Engineering, Delft University of Technology. The simulator consisted of a mockup dashboard with three LCD projectors (BenQ W1080ST 1080p) that provided roughly 180-degree vision. The simulated vehicle had an automatic gearbox and a top speed of 210 km/h.

Physiological data were logged using a biosignalsPlux wireless hub at 1000Hz. Heart rate was recorded using three pre-gelled Ag/AgCl electrodes at the heart's v3-node. Skin response was measured using the same pre-gelled electrodes, placed inside the palm and on the wrist of both hands. Simulator data were logged at 100Hz.

Scenarios

The scenarios used to induce workload in drivers each consisted of a 25km long, single-lane road. The road was divided into four 6km sections of different lane width (3.6m, 2.8m, 2.4m, 2.0m). Each section had seven curves, five with an inner radius of 750m and two with a 500m radius. Transitions between sections of different width always took place in a 750m radius curve, and were preceded by a road sign indicating a narrowing road. The four sections were identical, with the exception that the curves of segments 2 and 4 were mirrored with respect to section 1 and 3.

Cones were placed 8m apart on the road markings on both sides of the road. The main task was to hit as few cones as possible. A cone hit was indicated to the participant visually by a red dot on the side of the car where the cone was hit, and by a loud auditory beep. Extra difficulty in

lane-keeping was induced by a perturbation added to the vehicle's lateral motion. This perturbation was an unpredictable multi-sine signal with five frequencies between 0.067Hz and 0.25Hz, with a maximum summed amplitude of 1,000N. Without the perturbation, lane keeping (especially on straight segments) was not considered challenging enough. The width of the simulated vehicle was 1.8m.

Three runs were driven with the aim of inducing different levels of workload: a self-paced run and two forced-pace runs of 90km/h and 130km/h. In the self-paced run, participants had full longitudinal control over the car and could drive at their own pace. In the forced-pace conditions, however, the car's speed was automated and kept constant at 90km/h and 130km/h. This would push participants into curves at high speeds, with the goal of raising their workload significantly. The three runs were presented to the participants in randomised order.

Procedure

Participants read and signed an informed consent form, informing them of the purpose and procedure of the study. Participants were instructed that the main task was to minimise the total number of cone hits. Furthermore, participants were informed that during the experiment, a beep would sound every 20 seconds. At the sounding of this beep, participants were asked to verbally answer the question "From 0 to 10, how much effort does the current driving task take you?", with 0 being 'no effort', 5 being 'moderate effort' and 10 being 'a lot of effort'.

Before the experiment started, participants were familiarised with the simulator and the procedure by driving two 3.7km trial runs. The first trial run was self-paced, the second was forced-pace with speed at 110km/h. After the trial run, any question the participant had was answered. The electrodes were attached, and a one-minute baseline was recorded.

Analysis

Participants rated their mental effort on a scale of 0-10, every 20 seconds. This rating was annotated by the experimenter and added to the dataset. What data were logged, data preprocessing, ML set generation, model development and evaluation are identical to what has been described in the previous study.

Results

Participants

In total twenty-four participants took part in the experiment (17 male, 7 female). The average age was 24.6 years (SD 2.4). Participants reported driving multiple times a week (11 participants), at least once a month (7 participants) or less than one month (6 participants). All participants held a valid driving license. Reported mental effort was sensitive to the lane width variations, although regarding speed only to 130 km/h forced-pace condition (33).

Individual Models

As in the previous study, training and testing sets for the individual models were generated by dividing the dataset into two stratified sets. More data per participant were collected than in the previous experiment, so data were split with the more stringent 60%/40% split ratio.

Results were similar to the previous study, and indicated that the models performed well, with RFR outperforming SVR. An inverse relationship between model performance and overlap factor was observed, as well as increasing performance with increasing overlap factors, both as in the previous experiment. For all individual models with a window size of 5s and overlap of 0%, the AE_{μ} was 1.046, the $AE_{\mu/2}$ 0.662, R^2 0.635, CR 40.74%, and CR+/-1 77.31%. The relatively larger

absolute errors, compared to individual models in the previous study, might have resulted from the wider workload scale, the different nature of the driving task, or the more frequent reporting of mental workload. More information is displayed in table 2.

Group Models

To evaluate performance at the group level, data were split with a 60%/40% split ratio. Results indicated group models attained high performance. For the model with window size 5s and 0% overlap, the AE_{μ} was 0.904, the $AE_{\mu 1/2}$ 0.638, R^2 0.774, CR 41.61%, and CR+/-1 82.30%. TABLE 2 displays the full results. Performance increased with larger overlap factors, and again an (weak) inverse relationship between performance and window size was observed.

TABLE 2 Performance Metrics RFR Models

Window Size	5 sec		10 sec		30 sec	
Overlap Factor	0%	50%	0%	50%	0%	50%
Individual Model						
AE_{μ}	1.046	0.823	1.213	0.853	1.127	0.870
$AE_{\mu 1/2}$	0.662	0.511	0.833	0.518	0.959	0.694
R^2	0.635	0.763	0.600	0.675	0.561	0.735
CR	40.74%	50.31%	33.93%	45.83%	20.83%	40.28%
CR +/- 1	77.31%	84.34%	70.83%	81.94%	65.83%	81.48%
Group Model						
AE_{μ}	0.904	0.730	0.984	0.808	1.084	0.876
$AE_{\mu 1/2}$	0.638	0.482	0.722	0.546	0.792	0.663
R^2	0.774	0.830	0.740	0.802	0.718	0.811
CR	41.61%	51.30%	35.12%	46.44%	34.22%	37.87%
CR +/- 1	82.30%	88.18%	80.32%	85.88%	73.21%	82.41%
Generalising Model						
AE_{μ}	1.878	1.988	1.988	1.989	1.809	1.717
$AE_{\mu 1/2}$	1.831	1.844	1.718	1.741	1.680	1.568
R^2	0.118	0.079	0.196	0.177	0.411	0.433
CR	14.09%	13.45%	12.62%	13.44%	15.72%	15.21%
CR +/- 1	41.92%	40.70%	44.15%	42.29%	47.16%	46.32%

Performance metrics for the best performing (RFR) classifier. The table displays the mean (μ) and median ($\mu 1/2$) absolute error metrics, the coefficient of determination (R^2), the correct rate (CR) and the miss-by-one correct rate (CR +/- 1).

Generalising Group Models

Model performance when generalising to unknown individuals was then assessed, which did not perform well in the first simulator experiment. Data sampling methods were identical to the previous study.

Results indicated models performed moderately well. For the best performing model with window size 30s and 50% overlap, the AE_{μ} was 1.717, the $AE_{\mu 1/2}$ 1.568, R^2 0.433, CR 15.21%, CR+/-1 46.32%. Although model absolute error is relatively large, the coefficient of determination indicated a moderate relationship between model and data. FIGURE 2 below displays the predicted and true values for the first four participants. Individual model performance varied, with workload being predicted well for some participants, while for others showed a correct trend but with a constant offset error. These offset errors inflated the absolute error rates and deflated the

predictive accuracy despite good model performance. Generally, a decreased performance with increased overlap factor was observed (except for the largest window size of 30s), as well as increased performance with increased window size. The effect is similar to results for the model generalisation step in the previous study, but more pronounced. The effect also corresponds with what has been reported before (5).

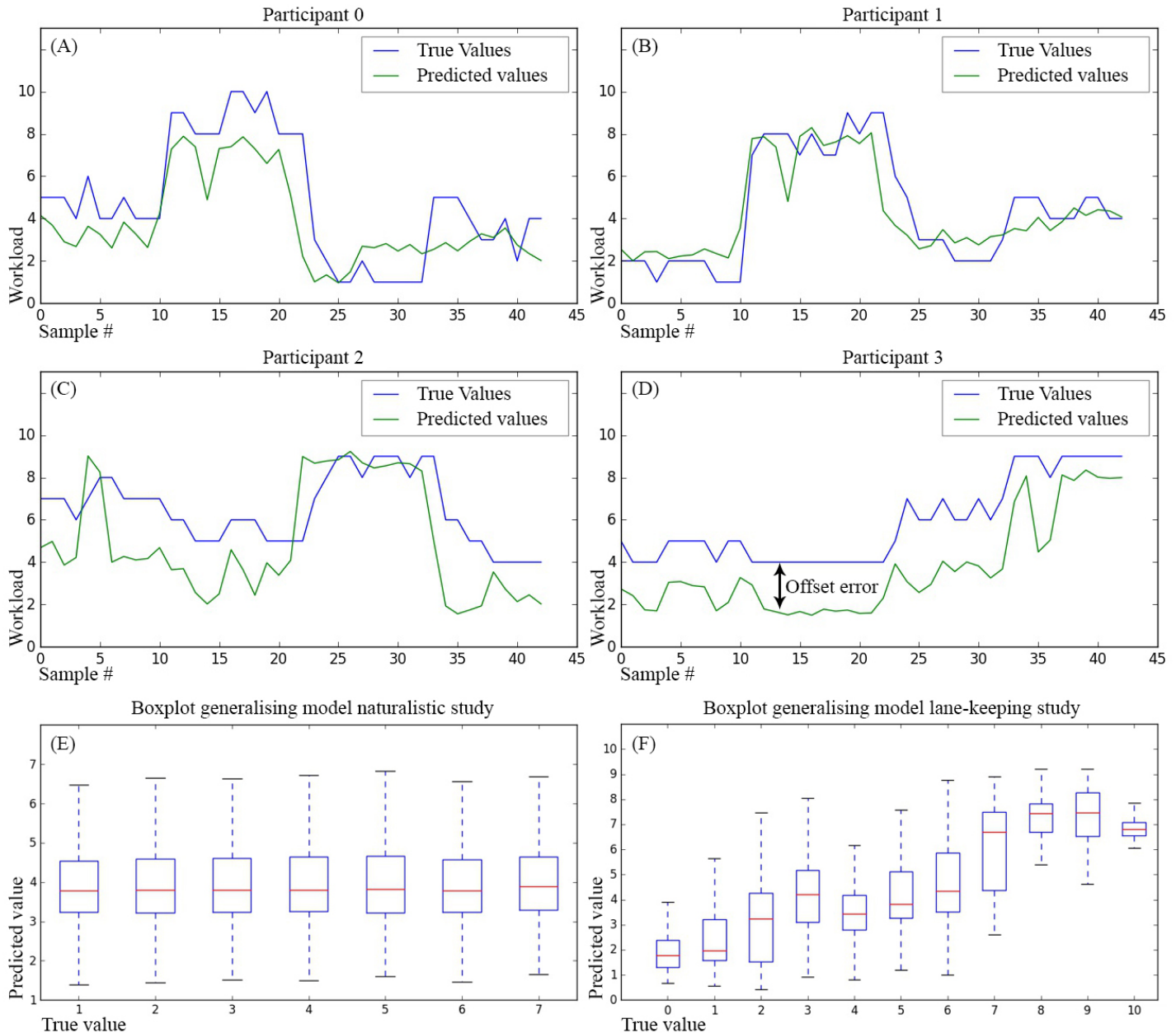


FIGURE 2 – The top four windows (A-D) show plots from the first four participants indicate that the models performed well, with the relatively large absolute errors likely resulting from individual scaling problems in the predictions. These offset errors are indicated in (D): the general trend is predicted well but there is a constant offset error. The last two windows (E-F) show box plots, further exploring the generalising models from both studies.

Conclusion

The results of this study show similarities with the previous study for individual and group-based models. Additionally, this second experiment shows that, when predicting multi-level workload (11 classes), generalising performance was satisfactory, although still with room for improvement.

This study seems to indicate that indeed non-binary workload prediction that generalises to unknown individuals is possible using ML methods. Although models generalising between individuals showed variations in performance based on which individual's workload was being predicted, including constant offset errors in several participants, overall performance was promising.

OVERALL CONCLUSIONS AND DISCUSSION

The present study tried to model driver workload using machine learning techniques that can run on embedded systems, with data collected from low-cost-sensors. Results have shown that individual models and within-group models functioned well in both a realistic driving setting as well as an artificial lane-keeping task setting. When generalising to unknown drivers, only the lane-keeping study produced usable results. As displayed in FIGURE 2 (E-F), in the first study the generalised model learns to predict values around the mean to optimize accuracy, in the second study the model learns to predict based on the reported workload.

Since the data we gathered in the study are time-series human physiological and performance data, it likely exhibits strong autocorrelation from one sample to the next. This might be a potential explanator for the higher performance in the individual and group models in both studies. Since with random sampling, shared variance between samples from the training set and the prediction set might bias the classifier towards a higher accuracy. To better assess performance, training cases were included where the models had to generalize to unknown individuals. These give a more accurate indication of performance, since with this approach there is no shared variance between training set (all participants minus participant k) and the testing set (participant k). As such, only the generalizing training case offers a reliable index of performance. This is an important distinction, since it shows that although using machine learning to predict driver workload can lead to promising results, care must be taken when interpreting the results. Without care in selecting the sampling techniques used, model performance might be inflated.

Possible reasons for the discrepancy in generalizing performance between both studies could include that the workload induced in the realistic settings was too low to be reflected in the physiological or performance signals, that workload induced by artificial tasks is more easily measurable than that induced by more realistic tasks, or that different physiological response patterns to workload might exist and that the sample in the first study was either too small or contained too much individual variation.

Possible limitations of the present study are that we employed a self-report measure as ground truth of the experienced mental workload of the drivers. We did not employ standardised workload scales such as NASA TLX, to keep interaction time and demand with the driver to a minimum. However, this may have contributed to lower model performance through participant response tendencies, and leaves some doubt as to what degree the data captures workload. In addition to this, we did not look at compensatory behaviour drivers might employ to manage their workload, such as reducing speed in complex or demanding situations.

Future directions are planned. These include feature space normalisation of the dataset to attempt to reduce the offset errors observed in some individuals, as well as exploring additional feature extraction methods. After this, on-road testing is planned to explore model performance in real-world driving settings. Lastly, development of an embedded variant of the model is planned.

REFERENCES

1. de Waard, D. *The Measurement of Drivers' Mental Workload*. Drukkerij Haasbeek, Alphen aan den Rijn, 1996.
2. Aasman, J., G. Mulder, and L. J. M. Mulder. Operator Effort and the Measurement of Heart-Rate Variability. *Human Factors*, Vol. 29, No. 2, 1987, pp. 161–170.
3. Matthews, G., L. E. Reinerman-Jones, D. J. Barber, and J. Abich. The Psychometrics of Mental Workload: Multiple Measures Are Sensitive but Divergent. *Human Factors: The Journal of the Human Factors and Ergonomics Society*, Vol. 57, No. 1, 2015, pp. 125–143. <https://doi.org/10.1177/0018720814539505>.
4. Mehler, B., B. Reimer, and J. F. Coughlin. Sensitivity of Physiological Measures for Detecting Systematic Variations in Cognitive Demand From a Working Memory Task: An On-Road Study Across Three Age Groups. *Human Factors: The Journal of the Human Factors and Ergonomics Society*, Vol. 54, No. 3, 2012, pp. 396–412. <https://doi.org/10.1177/0018720812442086>.
5. Solovey, E. T., M. Zec, E. A. Garcia Perez, B. Reimer, and B. Mehler. Classifying Driver Workload Using Physiological and Driving Performance Data. *Proceedings of the 32nd annual ACM conference on Human factors in computing systems - CHI '14*, 2014, pp. 4057–4066. <https://doi.org/10.1145/2556288.2557068>.
6. Jarvis, J., F. Putze, D. Heger, and T. Schultz. Multimodal Person Independent Recognition of Workload Related Biosignal Patterns. *Proceedings of the 13th international conference on multimodal interfaces - ICMI '11*, 2011, p. 205. <https://doi.org/10.1145/2070481.2070516>.
7. Liang, Y., M. L. Reyes, and J. D. Lee. Real-Time Detection of Driver Cognitive Distraction Using Support Vector Machines. *IEEE Transactions on Intelligent Transportation Systems*, Vol. 8, No. 2, 2007, pp. 340–350. <https://doi.org/10.1109/TITS.2007.895298>.
8. Kim, S., J. Chun, and A. K. Dey. Sensors Know When to Interrupt You in the Car. *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems - CHI '15*, 2015, pp. 487–496. <https://doi.org/10.1145/2702123.2702409>.
9. Moreira-matias, L., and H. Farah. On Developing a Driver Identification Methodology Using In-Vehicle Data Recorders. *IEEE Transactions on Intelligent Transportation Systems*, Vol. submitted, 2017. <https://doi.org/10.1109/TITS.2016.2639361>.
10. Birrel, S., M. Young, N. Staton, and P. Jennings. Using Adaptive Interfaces to Encourage Smart Driving and Their Effect on Driver Workload. Vol. 484, 2017, p. 764. <https://doi.org/10.1007/978-3-319-41682-3>.
11. Gent, P. Van, H. Farah, N. Van Nes, and B. Van Arem. Towards Real-Time, Nonintrusive Estimation of Driver Workload: A Simulator Study. *Proceedings of the Road Safety and Simulation Conference 2017*, 2017.
12. Breiman, L. Random Forests. *Machine learning*, Vol. 45, No. 1, 2001, pp. 5–32. <https://doi.org/10.1023/A:1010933404324>.
13. Sventnik, V., A. Liaw, C. Tong, J. C. Culberson, R. P. Sheridan, and B. P. Feuston. Random Forest: A Tool for Classification and Regression in Compound Classification and QSAR Modeling. *Journal of Chemical Information and Computer Science*, Vol. 43, No. 6, 2003, pp. 1947–1958. <https://doi.org/10.1016/j.rse.2008.02.011>.
14. Mishina, Y., R. Murata, Y. Yamauchi, T. Yamashita, and H. Fujiyoshi. Boosted Random Forest. *IEICE Transactions on Information and Systems*, Vol. E98, No. D, 2015, pp. 1630–1636.
15. Cortes, C., and V. Vapnik. Support-Vector Networks. *Machine Learning*, Vol. 20, No. 3, 1995, pp. 273–297. <https://doi.org/10.1023/A:1022627411411>.
16. Theodorides, T., and S. Member. Embedded Hardware-Efficient Real-Time Vector Machines. *IEEE transactions on neural networks and learning systems*, Vol. 27, No. 1, 2016, pp. 99–112. <https://doi.org/10.1109/TNNLS.2015.2428738>.
17. Bajaj, N., G. T. C. Chiu, and J. P. Allebach. Reduction of Memory Footprint and Computation Time for Embedded Support Vector Machine (SVM) by Kernel Expansion and Consolidation. *IEEE International Workshop on Machine Learning for Signal Processing, MLSP*, 2014.

- 1 <https://doi.org/10.1109/MLSP.2014.6958875>.
- 2 18. Jae Baek, H., H. Bit Lee, J. Soo Kim, J. Min Choi, K. Keun Kim, and K. Suk Park. Nonintrusive
- 3 Biological Signal Monitoring in a Car to Evaluate a Driver's Stress and Health State. *Telemedicine*
- 4 *and e-HEALTH*, Vol. 15, No. 2, 2009, pp. 182–189.
- 5 19. de Waard, D., A. Kruizinga, and K. A. Brookhuis. The Consequences of an Increase in Heavy Goods
- 6 Vehicles for Passenger Car Drivers' Mental Workload and Behaviour: A Simulator Study. *Accident*
- 7 *Analysis and Prevention*, Vol. 40, No. 2, 2008, pp. 818–828.
- 8 <https://doi.org/10.1016/j.aap.2007.09.029>.
- 9 20. Kun, A. L., A. Shyrovkov, and P. a. Heeman. Interactions between Human-Human Multi-Threaded
- 10 Dialogues and Driving. *Personal and Ubiquitous Computing*, Vol. 17, No. 5, 2013, pp. 825–834.
- 11 <https://doi.org/10.1007/s00779-012-0518-1>.
- 12 21. van Gent, P. Analyzing a Discrete Heart Rate Signal Using Python.
- 13 [http://www.paulvangent.com/2016/03/15/analyzing-a-discrete-heart-rate-signal-using-python-part-](http://www.paulvangent.com/2016/03/15/analyzing-a-discrete-heart-rate-signal-using-python-part-1/)
- 14 [1/](http://www.paulvangent.com/2016/03/15/analyzing-a-discrete-heart-rate-signal-using-python-part-1/).
- 15 22. van Gent, P. Python Heart Rate Analysis Toolkit. *GitHub Repository*.
- 16 https://github.com/paulvangentcom/heart_rate_analysis_python.
- 17 23. Reimer, B., B. Donmez, M. Lavallière, B. Mehler, J. F. Coughlin, and N. Teasdale. Impact of Age
- 18 and Cognitive Demand on Lane Choice and Changing under Actual Highway Conditions. *Accident*
- 19 *Analysis and Prevention*, Vol. 52, 2013, pp. 125–132. <https://doi.org/10.1016/j.aap.2012.12.008>.
- 20 24. Montano, N., A. Porta, C. Cogliati, G. Costantino, E. Tobaldini, K. R. Casali, and F. Iellamo. Heart
- 21 Rate Variability Explored in the Frequency Domain: A Tool to Investigate the Link between Heart
- 22 and Behavior. *Neuroscience and Biobehavioral Reviews*, Vol. 33, No. 2, 2009, pp. 71–80.
- 23 <https://doi.org/10.1016/j.neubiorev.2008.07.006>.
- 24 25. Billman, G. E. Heart Rate Variability - A Historical Perspective. *Frontiers in Physiology*, Vol. 2
- 25 NOV, No. November, 2011, pp. 1–13. <https://doi.org/10.3389/fphys.2011.00086>.
- 26 26. Lim, C. L., C. Rennie, R. J. Barry, H. Bahramali, I. Lazzaro, B. Manor, and E. Gordon.
- 27 Decomposing Skin Conductance into Tonic and Phasic Components. *International Journal of*
- 28 *Psychophysiology*, Vol. 25, No. 2, 1997, pp. 97–109.
- 29 [https://doi.org/10.1016/S0167-8760\(96\)00713-1](https://doi.org/10.1016/S0167-8760(96)00713-1).
- 30 27. Seitz, M., T. J. Daun, A. Zimmermann, and M. Lienkamp. Measurement of Electrodermal Activity
- 31 to Evaluate the Impact of Environmental Complexity on Driver Workload. *Proceedings of the*
- 32 *FISITA 2012 World Automotive Congress*, 2012, pp. 245–256.
- 33 <https://doi.org/10.1007/978-3-642-33741-3>.
- 34 28. Shimomura, Y., T. Yoda, K. Sugiura, A. Horiguchi, K. Iwanaga, and T. Katsuura. Use of Frequency
- 35 Domain Analysis of Skin Conductance for Evaluation of Mental Workload. *Journal of physiological*
- 36 *anthropology*, Vol. 27, No. 4, 2008, pp. 173–177. <https://doi.org/10.2114/jpa2.27.173>.
- 37 29. Köstinger, M., P. Wohlhart, P. M. Roth, and H. Bischof. Annotated Facial Landmarks in the Wild: A
- 38 Large-Scale, Real-World Database for Facial Landmark Localization. *Proceedings of the IEEE*
- 39 *International Conference on Computer Vision*, 2011, pp. 2144–2151.
- 40 <https://doi.org/10.1109/ICCVW.2011.6130513>.
- 41 30. Miyaji, M., M. Danno, H. Kawanaka, and K. Oguri. Driver's Cognitive Distraction Detection Using
- 42 Adaboost on Pattern Recognition Basis. *Proceedings of the 2008 IEEE International Conference on*
- 43 *Vehicular Electronics and Safety, ICVES 2008*, 2008, pp. 51–56.
- 44 <https://doi.org/10.1109/ICVES.2008.4640853>.
- 45 31. Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.
- 46 Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M.
- 47 Perrot, and É. Duchesnay. Scikit-Learn: Machine Learning in Python. *Journal of Machine Learning*
- 48 *Research*, Vol. 12, 2012, pp. 2825–2830. <https://doi.org/10.1007/s13398-014-0173-7.2>.
- 49 32. Chawla, N. V., K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. SMOTE: Synthetic Minority
- 50 over-Sampling Technique. *Journal of Artificial Intelligence Research*, Vol. 16, 2002, pp. 321–357.
- 51 <https://doi.org/10.1613/jair.953>.
- 52 33. Melman, T., D. A. Abbink, M. M. van Paassen, E. R. de Boer, and J. C. F. Winter. What Determines

1 Drivers' Speed? An Empirical Investigation of Three Behavioural Adaptation Models. *Manuscript*
2 *submitted for publication.*
3
4
5