

StressNAS: Affect State and Stress Detection Using Neural Architecture Search

Lam Huynh¹ Tri Nguyen² Thu Nguyen³ Susanna Pirttikangas² Pekka Siirtola⁴

¹Center for Machine Vision and Signal Analysis, University of Oulu ²Center for Ubiquitous Computing, University of Oulu

³Economics and Business Administration, University of Oulu ⁴Biomimetics and Intelligent Systems Group, University of Oulu

Abstract

Smartwatches have rapidly evolved towards capabilities to accurately capture physiological signals. As an appealing application, stress detection attracts many studies due to its potential benefits to human health. It is propitious to investigate the applicability of deep neural networks (DNN) to enhance human decision-making through physiological signals. However, manually engineering DNN proves a tedious task especially in stress detection due to the complex nature of this phenomenon. To this end, we propose an optimized deep neural network training scheme using neural architecture search merely using wrist-worn data from WESAD. Experiments show that our approach outperforms traditional ML methods by 8.22% and 6.02% in the three-state and two-state classifiers, respectively, using the combination of WESAD wrist signals. Moreover, the proposed method can minimize the need for human-design DNN while improving performance by 4.39% (three-state) and 8.99% (binary).

Keywords— Affect detection, Stress detection, Neural Architecture Search.

1. Introduction

Long-term stress can have negative effects on both mental and physical human's health. This can even lead to economic cost, such as absenteeism, diminished productivity at work, accidents, etc. [6]. Thus, detecting stress can greatly contribute to improving human's health and adding value to the economy. Therefore, detecting stress in a person has been widely discussed, and using physiological changes in the human body is one of the approaches in stress detection. However, the topic of detecting other affective states has not been seriously taken despite its contribution to human's emotion studies and commercial purposes.

Gjoreski et al. [5] is one of the first works that studied stress detection with a minimally intrusive approach. They used acceleration, blood volume pulse, heart rate data, galvanic skin response, and skin temperature recorded from a wrist-worn device to train a model for stress detection. Schmidt et al. [12] introduced a Multimodal Dataset for Wearable Stress and Affect Detection (WESAD). The data set contains physiological and motion data collected in a wrist-worn device and a chest-worn device.

They trained five machine learning models for detecting different affective states (baseline, stress and amusement): K-Nearest Neighbour (KNN), Linear Discriminant Analysis (LDA), Random Forest (RF), Decision Tree (DT), and AdaBoost Decision Tree (AB). The performance varies between the models and the types of classification problems. For the three-class classification problem (amusement vs. baseline vs. stress), the most well-performed approach when using wrist data is AB (75.21%). For binary classification problem (stress vs. non-stress), the most achievable accuracy is 87.12% obtained from RF using wrist data.

On the other hand, one might expect to improve affective detection performance by employing modern DNN architectures [13, 7, 9, 8, 14]. However, despite the efforts [4, 3], this remains a challenging task due to the complexity of human physiological signals. In this work, we introduced a novel framework, namely StressNAS, to optimize the deep neural network training using neural architecture search [1]. For comparison, we also implement other DNN architectures, including a multilayer perceptron, a fully convolutional network, and a residual-like DNN. Moreover, wrist-worn devices are less intrusive to the human body than chest-worn ones and are more common among users (e.g., smartwatches, wristbands). Potentially, affective state detection studies on data collected from wrist-worn devices can generate more user efficiency values. Therefore, we chose to evaluate our works in the three-state classifier (stress, baseline, and amusement) and binary classifier (stress and non-stress).

2. Data

The Wearable Stress and Affect Detection (WESAD) [12] is a high-quality multi-modal dataset aiming for human affective detection. Participants were asked to follow standard protocols to calibrate their states (neural, stress, amusement, meditation) before their physiology and motion signals were captured. The data were collected from 17 subjects; each took part in a 2-hour section. Unfortunately, due to device malfunction, data from subject ID #1 and #12 were discarded.

The data acquisition process utilized the chest-worn RespiBAN and wrist-worn Empatica E4 devices¹. The RespiBAN provided: respiration (RESP), electrocardiogram (ECG), electrodermal activity (EDA), electromyogram (EMG), skin temperature (TEMP), and three-axis acceleration (ACC) at 700 Hz. On the other hand, the Empatica E4 measured blood volume pulse (BVP, 64 Hz), electrodermal activity (EDA, 4 Hz), body temperature (TEMP, 4 Hz),

¹www.empatica.com/research/e4/

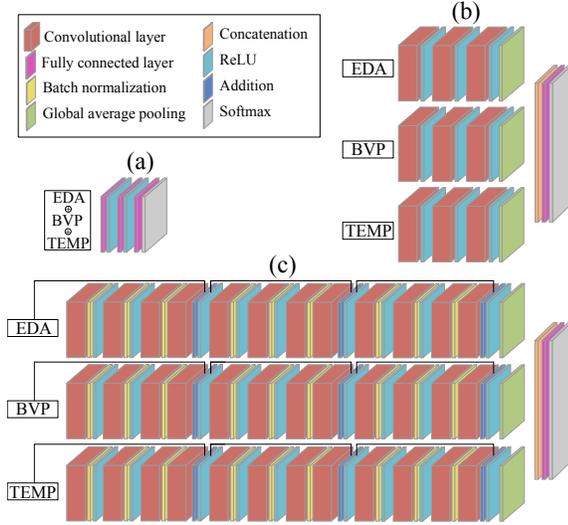


Figure 1. Manual design deep neural networks. (a) Multilayer Perceptron, (b) Fully Convolutional Network, and (c) Residual-like Deep Neural Network.

and three-axis acceleration (ACC, 32 Hz). The state conditions elicited from the protocol are referred to as the ground truth labels. As aforementioned, this work concentrates on the analysis of the WESAD wrist dataset.

3. Methods

This work uses a set of manual design DNNs with the WESAD dataset for preprocessing, training, and evaluating. Furthermore, we propose StressNAS, an automatic neural architecture search for affective state prediction.

3.1. Data processing

As mentioned earlier, we mainly utilized the data from the Empatica E4 for experimenting with the three-state (baseline, stress, and amusement) and the two-state (stress and non-stress) classification problems. We argue that training models with wrist data alone can be more challenging due to large variations in sampling rates of input modalities. Besides, we also transform the time-series data to filter banks for training DNNs.

The data were sampled using the sliding window technique. All experiments were conducted using a window length of 60 seconds with a 0.25-second shifting. As a result, approximately 132600 data samples were created in total.

A filter bank is a quadratic form of signal in the joint time-frequency domain that is a popular representation for training DNNs in speech processing. To obtain the filter banks, one can: 1) pass the signal through a pre-emphasis filter, 2) acquire overlapping frames from the filtered signal and then apply a windowing function (e.g. Hamming window), 3) take the Short-Time Fourier Transform to get the power spectrum, 4) apply triangular filters and mean normalization to calculate the filter banks.

3.2. Manual design deep neural networks

Figure 1 shows the architecture of a multilayer perceptron (MLP), a fully convolutional network (FCN) [10], and a residual-like DNN (ResNet) [7] used in our experiments. The MLP consists of three fully connected (FC) layers followed by nonlinear activation functions (rectified linear unit (ReLU) after layer 1-2 and Softmax for the last layer). The FCN contains several branches for separate modalities. Each branch has three convolutional layers (CONV), followed by ReLU with global average pooling (GAP) at the last layer. These features are then concatenated before feeding to an FC and Softmax for final prediction. The ResNet inherits a similar design principle of FCN while expanding the network using the residual blocks (Res-block). The Res-block contains four CONV-layers with batch normalization, ReLU, and residual connection. Each ResNet branch has three Res-block followed by a GAP before final concatenation and prediction.

3.3. StressNAS

The proposed network includes multiple DNNs that take in EDA, BVP, TEMP filter banks, and mixed features as inputs. Instead of manually constructing the DNN, we first randomly generated a set of network candidates (DNNs) for each input modality from the search space following [2] procedure. The core component of our candidates is the search cell. It is a directed acyclic graph that contains densely connected edges (feature transform operations, e.g., convolution, pooling, skip-connection) between its nodes (computed tensors).

For each modality, we randomly search from 10000 architectures. We then rank these architectures based on their scores that is calculated as the covariance matrices of the gradient with respect to the input data at the initial time [11]. The best ten architectures with the highest-ranking are utilized for training. Output features of each modality are concatenated for final prediction. In total, the searching, ranking and training time of our proposed method on one Tesla-V100 is ~ 50 hours. Figure 2 shows the overview architecture of our proposed approach.

4. Experiments and Results

The experiments consist of the evaluation of manual design DNNs and StressNAS models in three-state and two-state classifiers. Besides, those DNNs' results are compared with traditional models by Schmidt et al. [12].

4.1. Experiment and Evaluation metric

The experiments are about the implementation of different DNN approaches. In particular, we implement and build MLP, FCN, ResNet, and especially our proposal StressNAS for a three-state classifier (stress, baseline, and amusement) and binary classifier (stress and non-stress). The binary classifier considers the stress records and non-stress records as the combination of baseline and amusement records. For evaluation, we separate our training and testing data using leave-one-out cross-validation on 15 subjects from the WESAD dataset to build user-invariant models. Particularly, each person, by turn, becomes the test set, and the other 14 samples' data are the training set. Hence, the balanced accuracy in the experiments is an average accuracy of 15 prediction results.

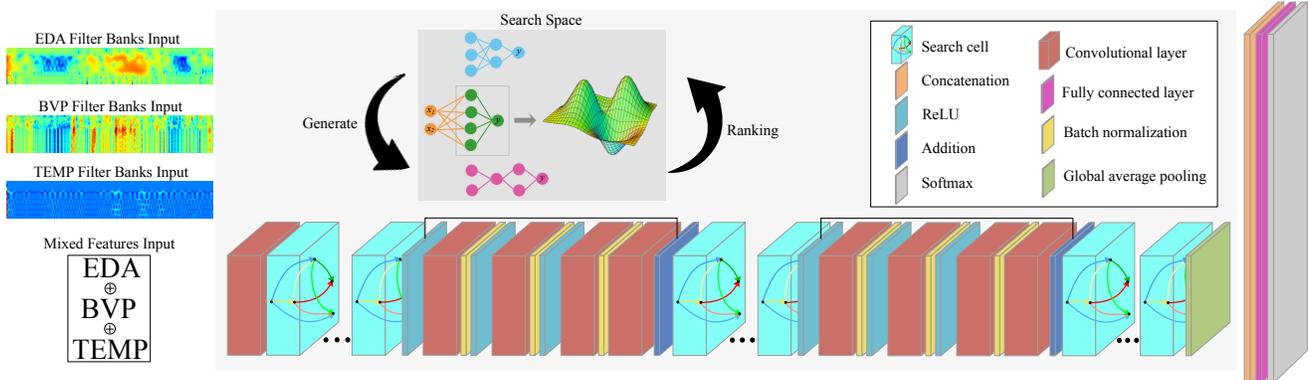


Figure 2. Overview architecture of StressNAS. The deep neural network (DNN) takes in EDA, BVP, TEMP filter banks, and mixed features. For each input, multiple DNNs are 1) randomly generated from the search space and 2) ranked based on their scores. Architectures with the highest-ranking are utilized for training. Output features of each modality are concatenated for final prediction.

Table 1. Results of different models and sensor combinations for classifying neural v. stress v. amusement states. Abbreviations: KNN = k-nearest neighbour, AB = AdaBoost DT, DT = Decision Tree, LDA = Linear discriminant analysis, RF = Random Forest, MLP = Multilayer Perceptron, FCN = Fully Convolution Network, ResNet = ResNet-like DNN, and StressNAS = our proposal.

Sensor Combinations	Schmidt et al. [12]					Manual Design DNNs			StressNAS
	AB	DT	RF	KNN	LDA	MLP	FCN	ResNet	
ACC+EDA+BVP+TEMP	75.21 ± 0.77	53.98 ± 1.79	74.85 ± 0.20	45.54	70.74	78.11	79.04	79.48	83.43
EDA+BVP+TEMP	73.62 ± 0.55	63.34 ± 1.00	76.17 ± 0.42	58.54	68.85	73.60	74.12	73.93	81.78
ACC	57.07 ± 0.57	53.71 ± 0.91	56.40 ± 0.16	45.54	47.73	52.16	46.53	45.25	55.81
EDA	59.42 ± 0.27	54.36 ± 0.27	56.57 ± 0.05	54.98	62.32	57.36	62.14	63.50	66.89
BVP	64.46 ± 0.21	57.57 ± 0.22	64.09 ± 0.12	59.44	70.17	62.43	65.42	68.11	71.24
TEMP	49.39 ± 0.23	47.42 ± 0.36	48.67 ± 0.21	44.32	58.96	55.14	56.32	61.35	62.15

Table 2. Results of different models and sensor combinations for classifying stress v. non-stress. Abbreviations: KNN = k-nearest neighbour, AB = AdaBoost DT, DT = Decision Tree, LDA = Linear discriminant analysis, RF = Random Forest, MLP = Multilayer Perceptron, FCN = Fully Convolution Network, ResNet = ResNet-like DNN, and StressNAS = our proposal.

Sensor Combinations	Schmidt et al. [12]					Manual Design DNNs			StressNAS
	AB	DT	RF	KNN	LDA	MLP	FCN	ResNet	
ACC+EDA+BVP+TEMP	83.98 ± 0.75	82.19 ± 0.44	87.12 ± 0.24	63.80	86.88	83.19	84.15	83.14	93.14
EDA+BVP+TEMP	88.05 ± 0.18	84.88 ± 0.11	88.33 ± 0.25	81.96	86.46	82.12	82.31	82.77	92.87
ACC	71.69 ± 0.45	64.08 ± 0.49	69.96 ± 0.42	63.80	60.02	65.15	67.88	66.85	72.15
EDA	79.71 ± 0.43	76.21 ± 0.27	76.29 ± 0.14	73.13	78.08	62.78	69.13	67.81	79.24
BVP	84.10 ± 0.13	81.39 ± 0.15	84.18 ± 0.11	82.06	85.83	69.97	72.15	69.15	81.16
TEMP	67.11 ± 0.34	68.22 ± 0.19	67.82 ± 0.11	64.46	69.24	55.17	68.12	62.54	71.46

4.2. Results

The results of the experiments are described in Table 1 and Table 2. In detail, the first column indicates the combination of sensors in the two first rows and the individual sensors in the remaining rows, while traditional approaches by Schmidt et al. [12] are illustrated in five first columns, and the other columns are DNNs' result. Remarkably, our proposal results with StressNAS are in the last column.

Table 1 presents the prediction results of four DNNs models in a three-state classifier with stress, baseline, and amusement. When comparing manual DNNs with traditional machine learning methods, ResNet performs better in three sensor combinations (ACC + EDA + BVP + TEMP, EDA, and TEMP), while MLP and FCN tend to have better performance with ACC + EDA + BVP + TEMP setting. Automatic DNN StressNAS provides bet-

ter prediction than traditional machine learning methods in most data combinations (except ACC data setting). Among different sensor combinations, both manual DNNs and StressNAS gain the highest accuracy with the combination of ACC, EDA, BVP, and TEMP. Specifically, all manual DNNs (MLP, FCN, and ResNet) and StressNAS obtain higher accuracy than traditional machine learning methods by 3-4% and 8%, respectively, in the combination of ACC, EDA, BVP, and TEMP. Moreover, Table 3 presents the best detection accuracies for each study subject for three states classification task. In general, the detection figures are evenly distributed among our subjects with a slight bias on subjects 3, 10, 13, and 16.

Table 2 shows the experimental results of the binary classifier (stress and non-stress detection). As is the case for the three-state classifier, both manual DNNs and StressNAS perform best with the sensor combination of ACC, EDA, BVP, and TEMP. How-

Table 3. Best detection accuracies for each study subject for three states classification task.

Study subject	Accuracy \uparrow
2	81.07
3	72.55
4	82.82
5	79.84
6	89.50
7	81.52
8	81.84
9	88.43
10	74.55
11	85.61
13	75.55
14	87.11
15	91.40
16	95.55
17	84.20

ever, interestingly, different from the three-class classifier, manual DNNs’ results are similar or lower than traditional models’ results, especially LDA and AB’s results. In contrast, the proposed method StressNAS provides a substantial improvement for almost all settings (except EDA and BVP signals). Notably, with the data combination of ACC, EDA, BVP, and TEMP, StressNAS achieves better prediction than traditional machine learning methods by 6%.

5. Conclusion

With the swift development of technologies, sensor devices have become more popular and wide uses in health care. Stress detection is a crucial topic supporting and improving human health by monitoring and analyzing body signals’ changes. Due to the interest in stress detection through physiological signals, many studies have concentrated on different activities, including constructing datasets, applying signal processing, and utilizing machine learning models. Despite existing studies related to stress detection, applying deep learning approaches still lacks consideration from the community. Therefore, the paper details using deep learning approaches to detect stress state through physiological signals. In detail, the paper utilizes a set of deep neural approaches (MLP, FCN, and ResNet) and proposes StressNAS, an automatically detected neural network architecture, to analyze affective states from WESAD’s wrist dataset. The evaluation is based on leave-one-out cross-validation to gain the result for comparison among different approaches in a three-state classifier (stress, baseline, and amusement) and binary classifier (stress and non-stress). Manual design DNNs are somewhat borderline with classical machine learning methods (CLASS), which suggests that carefully designing network architecture is a tedious task. On the other hand, DNN generated by neural architecture search (StressNAS) yields enhanced performance across all sensor combinations. Also, StressNAS tends to achieve better accuracy than CLASS in both the classifiers.

Acknowledgement

This work is supported by the Academy of Finland 6Genesis Flagship (grant 318927), the Vision-based 3D perception for mixed reality applications project and the TrustedMaaS project by the Infotech institute of the University of Oulu.

References

- [1] Zoph Barret and V Le Quoc. Neural architecture search with reinforcement learning. In *International conference on learning representations*, 2017. 1
- [2] Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. In *International Conference on Learning Representations*, 2019. 2
- [3] Maciej Dzieżyc, Martin Gjoreski, Przemysław Kazienko, Stanisław Saganowski, and Matjaž Gams. Can we ditch feature engineering? end-to-end deep learning for affect recognition from physiological sensor data. *Sensors*, 20(22):6535, 2020. 1
- [4] Martin Gjoreski, Matjaž Gams, Mitja Luštrek, Pelin Genc, Jens-U Garbas, and Teena Hassan. Machine learning and end-to-end deep learning for monitoring driver distractions from physiological and visual signals. *IEEE Access*, 8:70590–70603, 2020. 1
- [5] Martin Gjoreski, Hristijan Gjoreski, Mitja Luštrek, and Matjaž Gams. Continuous stress detection using a wrist device: in laboratory and real life. In *proceedings of the 2016 ACM international joint conference on pervasive and ubiquitous computing: Adjunct*, pages 1185–1193, 2016. 1
- [6] George Halkos and Dimitrios Bousinakis. The effect of stress and satisfaction on productivity. *International Journal of Productivity and Performance Management*, 2010. 1
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 2
- [8] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 1
- [9] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016. 1
- [10] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 2
- [11] Joseph Mellor, Jack Turner, Amos Storkey, and Elliot J. Crowley. Neural architecture search without training, 2020. 2
- [12] Philip Schmidt, Attila Reiss, Robert Duerichen, Claus Marberger, and Kristof Van Laerhoven. Introducing wesad, a multimodal dataset for wearable stress and affect detection. In *Proceedings of the 20th ACM International Conference on Multimodal Interaction*, pages 400–408, 2018. 1, 2, 3

- [13] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [1](#)
- [14] A Vaswani, N Shazeer, N Parmar, J Uszkoreit, L Jones, and AN Gomez. & polosukhin, i.(2017). attention is all you need. *Advances in neural information processing systems*, pages 5998–6008, 2017. [1](#)