



HAL
open science

EMG-Based Automatic Gesture Recognition Using Robust Neural Networks (Extended version)

Ana-Antonia Neacșu, Jean-Christophe Pesquet, Corneliu Burileanu

► **To cite this version:**

Ana-Antonia Neacșu, Jean-Christophe Pesquet, Corneliu Burileanu. EMG-Based Automatic Gesture Recognition Using Robust Neural Networks (Extended version). 2022. hal-03751766v1

HAL Id: hal-03751766

<https://hal.science/hal-03751766v1>

Preprint submitted on 16 Aug 2022 (v1), last revised 17 Feb 2023 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

EMG-Based Automatic Gesture Recognition Using Robust Neural Networks (Extended version)

Ana Neacșu^{a,b}, Jean-Christophe Pesquet^a, Corneliu Burileanu^b

^a*Universite Paris-Saclay, CentraleSupélec, Inria, Centre de Vision Numérique, 3 Rue Joliot Curie, Gif-sur-Yvette, 91190, France*

^b*Speech and Dialogue Laboratory, University Politehnica of Bucharest, 313 Splaiul Independentei, Bucharest, 060042, Romania*

Abstract

This paper introduces a novel approach for building a robust Automatic Gesture Recognition system based on Surface Electromyographic (sEMG) signals, acquired at the forearm level. Our main contribution is to propose new constrained learning strategies that ensure robustness against adversarial perturbations by controlling the Lipschitz constant of the classifier. We focus on positive neural networks for which accurate Lipschitz bounds can be derived, and we propose different spectral norm constraints offering robustness guarantees from a theoretical viewpoint. Experimental results on two distinct datasets highlight that a good trade-off in terms of accuracy and performance is achieved. We then demonstrate the robustness of our models, compared to standard trained classifiers in three scenarios, considering both white-box and black-box attacks.

Keywords: Recognition, Machine Learning, Perturbations, Stability, Lipschitz regularity, Optimization, EMG

1. Introduction

In recent years, the concept of human-computer interaction (HCI) has been at the core of many scientific and sociological developments. Combined with the power of machine learning algorithms, it has led to some of the most outstanding achievements in nowadays technology which are used successfully in an ever-increasing number of areas impacting our lives, e.g. medicine [1],

autonomous driving [2], natural language processing [3], etc. Researchers all around the world focus on providing new intuitive and accurate ways of interacting with devices around, based on gesture, voice, or vision analysis [4]. Gestures constitute a universal and intuitive way of communication, with the potential of bringing the Internet of Things (IoT) experience to a different, more organic level [5]. Automatic gesture recognition (AGR) algorithms can be successfully used in various applications, from sign language recognition (SLR) [6] to VR games [7].

Various solutions for AGRs based on image or video stream analysis, leveraging on computer vision algorithms have been proposed; see for example [8, 9, 10]. Multi-modal approaches for gesture classification have been also studied [11]. Although a good performance is achieved on synthetic data, in real-life scenarios these systems may be sensitive to environmental conditions, e.g. light conditions, background, etc. Additionally, these systems are often computationally demanding and consequently not always suited for real-time applications. Accelerometers and electromyography (EMG) sensors provide an alternative low-cost technology for gesture sensing. *sEMG* stands for surface electromyography and represents the electrical manifestation of the neuromuscular activation related to the contraction of the muscles [12]. The applications of sEMG-based classification systems are focused on, but not limited to, assistive devices and rehabilitation or postural control therapy for physically impaired persons [13]. With the continuous development of more versatile signal processing techniques, the applications of EMG signal classification expanded to a wide range of domains including augmented reality, gaming industry, military applications, etc.[14, 15].

Two critical issues need to be addressed when developing AGR algorithms: fast enough inference to ensure real-time feeling for the end-user, and accurate and robust classification to guarantee that the gesture is correctly identified no matter the environmental conditions. Machine learning methods have become ubiquitous tools in a wide range of tasks including AGR, on account of their ability to solve a great variety of problems, from simple regressions to complex multi-modal classification.

However, deep neural networks, which are probably the most powerful methods, may appear as black boxes whose robustness is not always well-controlled. For real-life applications, it is mandatory to guarantee the reliability of such techniques. Nowadays, the main difficulty to overcome consists in developing high-performance systems that are also trustable and safe. An additional challenge is to avoid implementation heaviness during the learning phase.

In [16], the authors showed that slightly altering data inputs that were correctly classified by the network can lead to a wrong classification [17, 18, 19]. This finding was at the origin of the concept of adversarial inputs, which constitute malicious input data that can deceive machine learning models. For example, [18] shows how voice interfaces can be fooled by creating carefully crafted artificial audio inputs of unintelligible voice that are miss-classified as specific vocal commands by the system. Also, [20] introduces several methods for generating adversarial examples on ImageNet that are so close to the original data that differences are indistinguishable for the human eye.

It must be emphasized that adversarial inputs are not necessarily artifi-

cially created with the intention to sabotage the system. As other physiological signals, e.g. EEG or EKG, EMG signals have low frequency components (usually between 10 – 150Hz), and low amplitudes (≤ 10 mV Peak to Peak). This makes them very sensitive to noise and outside perturbations that can occur innately, under the form of noise stemming from acquisition devices, imperfect sensor contact, etc. Those can seriously flaw the performance of real-life applications based on pre-trained models [21].

As highlighted in [20], the Lipschitz behaviour of the network is tightly correlated with its robustness against adversarial attacks. The Lipschitz constant allows to upper bound the output perturbation knowing the magnitude of the input one, for a given metric [22]. Controlling this constant thus represents a feasible solution to limit the effect of adversarial attacks. Computing the exact Lipschitz constant of a neural network is however a very complex problem, so the main challenge is to find clever ways to approximate this constant effectively.

Recently, several techniques to ensure the Lipschitz stability of neural networks have been explored. For example, [19] proposes a novel weight spectral normalization technique applied to stabilize the training of the discriminator in Generative Adversarial Networks (GANs). The Lipschitz constant of the network is viewed as a hyper-parameter that can be tuned in the training process of the image generation task. Doing so leads to a model with improved generalization capabilities. In [23] norm-constraint GroupSort based architectures are proposed and it is shown that they can be used as universal Lipschitz function approximators. The authors apply gradient norm preservation to create Lipschitzian networks that offer adversarial robustness guar-

antees. In [24] the authors introduce Parseval networks, another approach for designing networks which are intrinsically robust to adversarial noise, by imposing the Lipschitz constant of each layer of the system to be less than 1. In [25] a convex optimization framework is introduced to compute tight upper bounds for the Lipschitz constant of Deep Neural Networks (DNNs). They make use of the observation that commonly used activation operators are gradients of convex functions. Semi-definite programming approaches to ensure robustness are also explored in [26].

The main contributions of this paper are:

- To propose a robust real-time Automatic Gesture Recognition system based on sEMG signals. The robustness is ensured by using a novel learning algorithm for training feedforward neural networks.
- To show that a good accuracy-robustness balance can be reached. To do so, we train the system under carefully crafted spectral norm constraints, allowing us to finely control its Lipschitz constant. A tight Lipschitz constant is efficiently estimated by focusing on neural networks with positive weights as in [27].
- To demonstrate the performance of the final architecture in real-life experiments where we show that the proposed robust model outperforms those trained conventionally.
- To analyze how our system behaves when the input is affected by different noise levels, simulating perturbations that may occur in real scenarios.
- To show the validity of our solution by experimenting on two distinct publicly available sEMG gestures datasets.

The rest of the paper is structured as follows. The theoretical background of our work is detailed in Section 2. In Section 3, we present the proposed optimization algorithm and we investigate the way of dealing with the constraints. The application and the results are discussed in Section 4, while Section 5 deals with how our model behaves when facing adversarial data. Finally, Section 6 contains some concluding remarks.

2. Robustness solutions in the context of nonnegative neural networks

2.1. Problem formulation

Any feedforward neural network is obtained by cascading m layers associated with operators $(T_i)_{1 \leq i \leq m}$. The neural network can thus be expressed as the following composition of operators:

$$T = T_m \circ \dots \circ T_1. \quad (1)$$

Each layer $i \in \{1, \dots, m\}$ has a real-valued vector input x_i of dimension N_{i-1} which is mapped to

$$T_i(x_i) = R_i(W_i x_i + b_i), \quad (2)$$

where $W_i \in \mathbb{R}^{N_i \times N_{i-1}}$, $b_i \in \mathbb{R}^{N_i}$ are the weight matrix and bias parameter respectively. $R_i: \mathbb{R}^{N_i} \rightarrow \mathbb{R}^{N_i}$ constitutes a non-linear activation operator which is applied component-wise (e.g., ReLU or Sigmoid) or globally (e.g., Softmax). Figure 1 shows a graphical representation of this concept.

Even though the choice of the activation R_i may differ depending on the task at hand, it has been shown in [28, 29] that most of them are actually α_i -averaged operators with $\alpha_i \in]0, 1]$. Recall that R_i is an α_i -averaged operator



Figure 1: Representation of a NN as a composition of operators

if, for every pair $(x_i, y_i) \in (\mathbb{R}^{N_i})^2$, the following inequality holds:

$$\|R_i(x_i) - R_i(y_i) - (1 - \alpha_i)(x_i - y_i)\| \leq \alpha_i \|x_i - y_i\|. \quad (3)$$

When $\alpha_i = 1/2$, R_i is said to be firmly nonexpansive. For standard choices of activation operators, R_i is firmly nonexpansive since it is the proximity operator of a proper, lower-semicontinuous function (see [29] for more details). Note that, in [25], it is assumed that R_i operates component-wise and is slope-bounded. The authors emphasize that the most common case corresponds to lower and upper slope values equal to 0 and 1, respectively. It follows from [30, Proposition 2.4] that a function satisfies this property if and only if it is the proximity operator of some proper lower-semicontinuous convex function, so that similar assumptions to those made in [29] are recovered.

As explained in [28], examples of activation operators R_i which are α_i -averaged with $\alpha_i > 1/2$ can be encountered. They basically correspond to over-relaxations of firmly nonexpansive operator. An example of such operators is the Swish activation function [31]. Another famous example is the group-sort operator:

$$\left(\forall x_i = \begin{bmatrix} x_{i,1} \\ \vdots \\ x_{i,M} \end{bmatrix} \in \mathbb{R}^{N_i} \right) \quad R_i(x_i) = \begin{bmatrix} x_{i,1}^\uparrow \\ \vdots \\ x_{i,M}^\uparrow \end{bmatrix}, \quad (4)$$

where the vector x_i has been decomposed in M subvectors $x_{i,j}$ with $j \in \{1, \dots, M\}$, of dimension B ($N_i = BM$) and $x_{i,j}^\uparrow$ designates the vector of components of $x_{i,j}$ sorted in ascending order. R_i is then purely nonexpansive, i.e. $\alpha_i = 1$. Note that max-pooling can be achieved by composing this group sort operation with a linear operator. Indeed, if $i < m$, $M = N_{i+1}$, and W_{i+1} is the matrix extracted from the $N_i \times N_i$ identity matrix Id_{N_i} by selecting the matrix rows with indices multiple of B , then $W_{i+1} \circ R_i$ corresponds to a max-pooling.

2.2. Lipschitz robustness certificate

Consider a neural network T as described in Fig. 1. let $x \in \mathbb{R}^{N_0}$ be the input of the network and let $T(x) \in \mathbb{R}^{N_m}$ be its associated output. By adding some small perturbation $z \in \mathbb{R}^0$ to the input, the perturbed input is $\tilde{x} = x + z$. The effect of the perturbation on the output of the system can be quantified by the following inequality:

$$\|T(\tilde{x}) - T(x)\| \leq \theta_m \|z\|, \quad (5)$$

where $\theta_m \geq 0$ denotes a Lipschitz constant of the network. θ_m represents thus an important parameter that allows us to assess and control the sensitivity of a neural network to various perturbations. It needs however to be accurately estimated to provide valuable information. A standard approximation to the Lipschitz constant [20] is given by

$$\theta_m = \prod_{i=1}^m \|W_i\|_S, \quad (6)$$

where $\|\cdot\|_S$ denotes the *spectral norm* of a matrix. Although simple to compute, this approximate bound is over-pessimistic. Different methods for

obtaining tighter estimates of the Lipschitz constant have been presented in the recent literature; see for example [22, 28, 25, 32, 33]. Local estimates of the Lipschitz constant can also be performed which may appear more relevant. But they are more complex to compute and, as we will see, controlling the global Lipschitz constant is usually sufficient to get a good performance. Estimating the global Lipschitz constant of the network is an NP (non-deterministic polynomial-time)-hard problem [22]. Although there exist efficient approaches to approximate an accurate bound [25, 32, 33], computing these estimates may be expensive for wide or deep networks. In addition, using these bounds within a training procedure is a difficult task [26]. In this work, we will make the following assumption.

Assumption 2.1 Let a neural network be given by (1) where the i -th layer with $i \in \{1, \dots, m\}$ is given by (2). We assume that

- (i) all the activation layers, except possibly the last one, consist of separable averaged operators, that is, for every $i \in \{1, \dots, m-1\}$, there exist averaged functions $(\rho_{i,k})_{1 \leq k \leq N_i}$ from \mathbb{R} to \mathbb{R} such that $R_i: (\xi_{i,k})_{1 \leq k \leq N_i} \mapsto (\rho_{i,k}(\xi_{i,k}))_{1 \leq i \leq k}$;
- (ii) at the last activation layer, R_m is an averaged operator.

Our approach will be grounded on the following result.

Proposition 2.2 [28] *Suppose that Assumption 2.1 holds. For every $i \in \{1, \dots, m\}$, let A_i be the matrix whose elements are the absolute values of those of W_i . Then,*

$$\vartheta_m = \|A_m \cdots A_1\|_s \tag{7}$$

is a Lipschitz constant of T . In addition

$$\|W_m \cdots W_1\| \leq \vartheta_m. \quad (8)$$

In particular if, for every $i \in \{1, \dots, m\}$, $W_i \in [0, +\infty[^{N_i \times N_{i-1}}$, ϑ_m is equal to the lower bound in (8).

Based on this proposition, the best estimate for the Lipschitz constant of a given feedforward neural network having nonnegative weights simplifies to the spectral norm of the product of all the weight matrices composing the network. More precisely, the obtained Lipschitz constant

$$\vartheta_m = \|W_m \cdots W_1\|_S$$

is the Lipschitz constant of a purely linear network, where all the non-linear activation operators have been replaced with the identity operator.

The above result is guaranteed to be valid only in the case when all the weights are nonnegative. In the general case of networks with weights having arbitrary signs, it can be proved that $\|W_m \cdots W_1\|_S$ represents only a lower bound of the Lipschitz constant established in [28]. It is also worth mentioning that the proposed results hold for any algebraic structure of the weight matrices $(W_i)_{1 \leq i \leq m}$.

3. Optimization methods for training robust feedforward networks

3.1. Stochastic gradient descent – projected variant

Standard training in neural networks consists in the minimization of a nonconvex cost function with respect to the model parameters by means of an iterative strategy. Let \mathcal{L} be the cost function defined as follows:

$$\mathcal{L}(\eta) = \sum_{k=1}^K \ell(z_k, \eta), \quad (9)$$

where $\eta = (\eta_i)_{1 \leq i \leq m}$ is a vector encompassing all the model parameters. For each layer $i \in \{1, \dots, m\}$, η_i denotes a vector of dimension $N_i(N_{i-1} + 1)$ that contains the scalar variables associated with the weight matrices W_i and the corresponding bias components b_i . The data information is represented by $(z_k)_{1 \leq k \leq K}$. For every $k \in \{1, \dots, K\}$, z_k is a pair consisting of an input of the system and the associated desired output (ground truth). Also, ℓ represents the loss function assumed to be differentiable (almost everywhere) with respect to η .

To ensure robustness, we shall impose spectral norm constraints on the weight matrices. In other words, the vector of parameters η is constrained to belong to a closed set \mathcal{S} that will be described in the next section. We propose to use an extension of a standard optimization techniques for training neural networks [34]. More specifically, we will implement a *projected stochastic gradient* algorithm. A momentum parameter is introduced in this algorithm to accelerate the convergence process.

Algorithm 1 describes the iterations performed at each epoch $n > 0$. We see that there are two nested loops: the outer loop operates on the batch index q and the second one on the layer index i . In this algorithm, $\gamma_n \in]0, +\infty[$ is the learning rate, while $\zeta_n \in [0, +\infty[$ denotes the inertia parameter for momentum. The algorithm is very similar to block-iterative techniques used in convex optimization [34]. The parameters of each layer are indeed updated successively by performing a gradient step on the data in the current mini-batch (which can be epoch-dependent). ∇_i represents the gradient, computed by standard backpropagation mechanism, with respect to η_i for each $i \in \{1, \dots, m\}$. This stochastic gradient step is followed by a

projection $\mathbb{P}_{\mathcal{S}_{i,n}}$ onto the constraint set $\mathcal{S}_{i,n}$. The definition of this set as well as the way of handling this projection are detailed in the following.

Algorithm 1: Projected SGD Algorithm

Partition $\{1, \dots, K\}$ into minibatches $(\mathbb{L}_{q,n})_{1 \leq q \leq Q}$ **foreach**

$q \in \{1, \dots, Q\}$ **do**

foreach $i \in \{1, \dots, m\}$ **do**

$$\left[\begin{array}{l} \Delta_{i,n} = (1 + \zeta_n)\eta_{i,n} - \zeta_n\eta_{i,n-1} \quad \tilde{\eta}_{i,n} = [(\eta_{j,n+1}^\top)_{j < i} \quad \Delta_{i,n}^\top \quad (\eta_{j,n}^\top)_{j > i}]^\top \\ \eta_{i,n+1} = \mathbb{P}_{\mathcal{S}_{i,n}} \left(\Delta_{i,n} - \gamma_n \sum_{k \in \mathbb{L}_{q,n}} \nabla_i \ell(z_k, \tilde{\eta}_{i,n}) \right) \end{array} \right.$$

where $\mathcal{S}_{i,n} = \{\eta_i \mid [(\eta_{j,n+1}^\top)_{j < i} \quad \eta_i^\top \quad (\eta_{j,n}^\top)_{j > i}]^\top \in \mathcal{S}\}$.

3.2. Constraint sets

As mentioned before, this work revolves around feed-forward networks with positive weights. Thus, the first condition that we impose is nonnegativity for each layer $i \in \{1, \dots, m\}$, which is modeled by the constraint set

$$\mathcal{D}_i = \{W_i \in \mathbb{R}^{N_i \times N_{i-1}} \mid W_i \geq 0\} \quad (10)$$

Moreover, based on our standing assumptions and Proposition 2.2, we must impose a spectral norm constraint on the weight matrices to control the robustness of the system. This translates mathematically as the following upper bound constraint:

$$\|W_m \cdots W_1\|_S \leq \bar{\vartheta}, \quad (11)$$

where $\bar{\vartheta}$ represents the target maximum Lipschitz constant of the network. This bound constitutes a direct measure of the system level of robustness against adversarial inputs. We need to handle these two constraints simultaneously during the training process. Imposing nonnegativity is fairly easy

since (10) defines a simple convex constraint. By contrast, constraint (11) does not satisfy the convexity property. Since (11) corresponds to a closed set in the underlying space of weight matrices and this set has a nonempty intersection with $\mathcal{D}_1 \times \cdots \times \mathcal{D}_m$, the projection onto the intersection of the two sets can be defined but it is not guaranteed to be unique. To circumvent this difficulty, it can be noticed that (11) actually defines a multi-convex constraint in the sense that if, for every $i \in \{1, \dots, m\}$, $(W_j)_{1 \leq j \leq m, j \neq i}$ are given, then (11) imposes a convex constraint on W_i . This suggests to introduce the following closed and convex set:

$$\mathcal{C}_{i,n} = \{W_i \in \mathbb{R}^{N_i \times N_{i-1}} \mid \|A_{i,n}W_iB_{i,n}\|_S \leq \bar{\vartheta}\} \quad (12)$$

in order to control the Lipschitz constant. Hereabove, the matrices $A_{i,n}$ and $B_{i,n}$ represent the product of the weight matrices for the previous and the posterior layers, respectively. By adopting the convention that $A_{i,n} = \text{Id}$ if $i = m$ and $B_{i,n} = \text{Id}$ if $i = 1$, we define these matrix products as

$$A_{i,n} = W_{m,n} \cdots W_{i+1,n}, \quad B_{i,n} = W_{i-1,n+1} \cdots W_{1,n+1}, \quad (13)$$

where $(W_{j,n})_{1 \leq j \leq m}$ denote the estimates of the weight matrices at each iteration n , as it appears in Algorithm 1.

Thus, our objective will be to perform the projection onto the set $\mathcal{S}_{i,n} = \mathcal{D}_i \cap \mathcal{C}_{i,n}$, for each layer $i \in \{1, \dots, m\}$ and at each iteration n . Several algorithms can be envisaged to solve this convex optimization problem.

Before describing our proposed algorithmic solution, let us recall the expressions of the required elementary projections. For every $W \in \mathbb{R}^{S \times T}$, the projection of W onto $[0, +\infty]^{S \times T}$ is

$$\text{P}_{[0, +\infty]^{S \times T}}(W) = (\widetilde{W}_{s,t})_{1 \leq s \leq S, 1 \leq t \leq T}, \quad (14)$$

where, for every $s \in \{1, \dots, S\}$ and $t \in \{1, \dots, T\}$,

$$\widetilde{W}_{s,t} = \begin{cases} W_{s,t} & \text{if } W_{s,t} \geq 0 \\ 0 & \text{otherwise.} \end{cases} \quad (15)$$

Let $\mathcal{B}(0, \bar{\vartheta})$ be the closed spectral ball of center 0 and radius $\bar{\vartheta}$ defined as ¹

$$\mathcal{B}(0, \bar{\vartheta}) = \{W \in \mathbb{R}^{S \times T} \mid \|W\|_s \leq \bar{\vartheta}\}. \quad (16)$$

For every $W = (W_{s,t})_{1 \leq s \leq S, 1 \leq t \leq T} \in \mathbb{R}^{S \times T}$, let $U\Lambda V^\top$ be the singular value decomposition of W , where $U \in \mathbb{R}^{S \times R}$ and $V \in \mathbb{R}^{T \times R}$ are matrices such that $U^\top U = \text{Id}$ and $V^\top V = \text{Id}$, $R = \min\{S, T\}$, and $\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_R)$, $(\lambda_r)_{1 \leq r \leq R} \in [0, +\infty[^R$ being the singular values of W . Then the projection of W onto $\mathcal{B}(0, \bar{\vartheta})$ is expressed as

$$\text{P}_{\mathcal{B}(0, \bar{\vartheta})}(W) = U\widetilde{\Lambda}V^\top \quad (17)$$

where $\widetilde{\Lambda} = \text{Diag}(\widetilde{\lambda}_1, \dots, \widetilde{\lambda}_r)$ and

$$(\forall i \in \{1, \dots, r\}) \quad \widetilde{\lambda}_i = \begin{cases} \lambda_i & \text{if } \lambda_i \leq \bar{\vartheta} \\ \bar{\vartheta} & \text{otherwise.} \end{cases} \quad (18)$$

To compute the projection onto $\mathcal{S}_{i,n}$ of a matrix $\overline{W}_i \in \mathbb{R}^{N_i \times N_{i-1}}$, we propose to employ the FISTA (Fast Iterative Shrinkage-Thresholding Algorithm) version of a dual forward-backward method in Algorithm 2. This algorithm is based on a dual proximal approach [35] and constitutes an extension of the optimization method originally proposed in [36]. The rationale for this algorithm is given in the appendix.

¹To simplify our notation, $\mathcal{B}(0, \bar{\vartheta})$ will designate any spectral ball of this kind whatever the dimensions of the involved matrices.

Algorithm 2: FISTA-like accelerated version of DFB algorithm

Let $Y_0 \in \mathbb{R}^{N_m \times N_0}$

Set $\gamma = 1/(\|A_{i,n}\|_S \|B_{i,n}\|_S)^2$

Set $\alpha \in]2, +\infty[$

for $l = 0, 1, \dots$ **do**

$$\left\{ \begin{array}{l} \eta_l = \frac{l}{l+1+\alpha} \\ Z_l = Y_l + \eta_l(Y_l - Y_{l-1}) \\ V_l = \mathbf{P}_{D_i}(\bar{W}_i - A_{i,n}^\top Z_l B_{i,n}^\top) \\ \tilde{Y}_l = Z_l + \gamma A_{i,n} V_l B_{i,n} \\ Y_{l+1} = \tilde{Y}_l - \gamma \mathbf{P}_{\mathcal{B}(0, \bar{\vartheta})}(\gamma^{-1} \tilde{Y}_l) \end{array} \right.$$

return V_l

3.3. Handling looser constraints

The Lipschitz constant of the network can be controlled in multiple ways. Besides the solution formulated in Section 3.2, a more standard approach to control it [16] consists in imposing

$$\prod_{i=1}^m \|W_i\|_S \leq \bar{\vartheta}. \quad (19)$$

Two strategies have been implemented to enforce this constraint.

- (i) The first one consists in imposing a uniform bound on the spectral norm of each weight matrix $(W_i)_{1 \leq i \leq m}$, which leads to the following convex constraint sets:

$$(\forall i \in \{1, \dots, m\}) \quad \tilde{\mathcal{C}}_i = \{W_i \in \mathbb{R}^{N_i \times N_{i-1}} \mid \|W_i\|_S \leq \bar{\vartheta}^{1/m}\}. \quad (20)$$

- (ii) The second strategy aims at introducing more flexible bounds on the spectral norms of each layer. It is based on the following choice for the

individual convex constraint sets:

$$\begin{aligned}
 & (\forall n \in \mathbb{N} \setminus \{0\})(\forall i \in \{1, \dots, m\}) \\
 & \check{\mathcal{C}}_{i,n} = \left\{ W_i \in \mathbb{R}^{N_i \times N_{i-1}} \mid \|W_i\|_S \leq \|W_{i,n}\|_S \left(\frac{\vartheta}{\prod_{j=1}^m \|W_{j,n}\|_S} \right)^{1/m} \right\}.
 \end{aligned}$$

For every $i \in \{1, \dots, m\}$, projecting onto $\tilde{\mathcal{C}}_i$ or $\check{\mathcal{C}}_{i,n}$ is performed by truncating a singular value decomposition, similarly to the technique described at the end of Section 3.2. The projections onto $\tilde{\mathcal{C}}_i \cap \mathcal{D}_i$ and $\check{\mathcal{C}}_{i,n} \cap \mathcal{D}_i$ can then be computed by using the same iterative method as in Algorithm 2 with $A_{i,n} = B_{i,n} = \text{Id}$.

In all the proposed constrained optimization methods, the projection $\mathbf{P}_{\mathcal{B}(0, \tilde{\vartheta})}$ onto a spectral ball with radius $\tilde{\vartheta} > 0$ plays a prominent role. The ball radius depends on the handled constraint (11), (20), or (10). A complex operation such as a singular value decomposition may be very demanding in terms of computational resources when dealing with large size matrices. In that case, we propose to use an approximate projection [19] defined as

$$(\forall W \in \mathbb{R}^{S \times T}) \quad \mathbf{P}_{\mathcal{B}(0, \tilde{\vartheta})}(W) \simeq \begin{cases} W & \text{if } \|W\|_S \leq \tilde{\vartheta} \\ \frac{\tilde{\vartheta}}{\|W\|_S} W & \text{otherwise.} \end{cases} \quad (21)$$

Using this approximation in Algorithm 2 yields approximate projections $(\tilde{\mathbf{P}}_{\mathcal{C}_{i,n} \cap \mathcal{D}_i})_{1 \leq i \leq m, n > 0}$. Note however that we then lose the theoretical guarantees of convergence Algorithm 2, even if this issue was not observed in our implementation.

An additional advantage of Formula (21) is that it allows the nonnegativity of the elements of the input matrix to be kept. This allows us to derive cheap approximate versions of the projection onto $\tilde{\mathcal{C}}_i \cap \mathcal{D}_i$ with $i \in \{1, \dots, m\}$

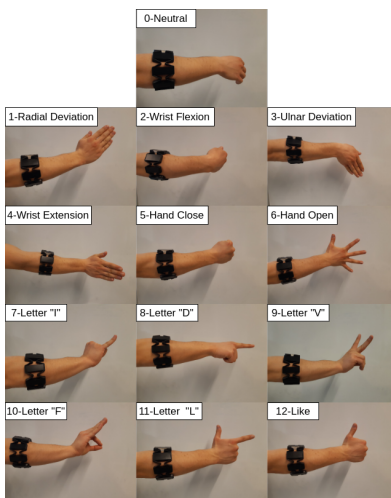


Figure 2: 13-gestures Dataset [37]

by first projecting onto \mathcal{D}_i and then applying the approximate projection onto $\tilde{\mathcal{C}}_i$. The resulting approximate projection is denoted by $(\tilde{\mathcal{P}}_{\tilde{\mathcal{C}}_i \cap \mathcal{D}_i})_{1 \leq i \leq m}$. A similar procedure can be followed to compute approximate projections $(\tilde{\mathcal{P}}_{\tilde{\mathcal{C}}_{i,n} \cap \mathcal{D}_i})_{1 \leq i \leq m, n > 0}$ onto $(\tilde{\mathcal{C}}_{i,n} \cap \mathcal{D}_i)_{1 \leq i \leq m, n > 0}$.

4. AGR Experimental Setup

4.1. sEMG datasets

We test our proposed training scheme on two online datasets containing EMG information of different hand gestures. Both were acquired using Myo armband, a device developed by Thalmic Labs, equipped with eight sEMG sensors displayed circularly. To the best of our knowledge these are the only gesture-based datasets currently acquired with this armband which are publicly available.

The first dataset, detailed in [38] contains EMG signals characterizing 7 hand gestures correlated to the primary movements of the hand. There are

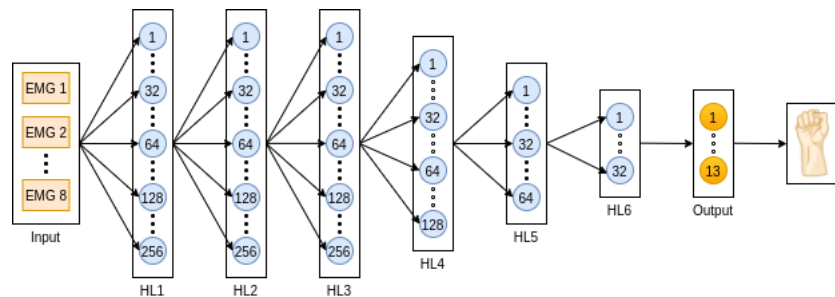


Figure 3: *Proposed neural network architecture for AGR. All the layers except the last one use ReLU activation functions; the last layer uses Softmax. The number of neurons considered for each layer is: 128, 128, 128, 64, 32, 16, in the case of 7-gestures dataset and 256, 256, 256, 128, 64, 32 in the case of 13-gesture dataset. The last layer has 7 or 13 neurons representing the gesture number being recognized. Each EMG box represents a column vector containing 8 time-descriptors.*

four mobility gestures (i.e., wrist flexion and extension, ulnar, and radial deviation) and two gestures used for grasping and releasing objects (i.e., spread fingers and close fist). The 7th gesture characterizes the neutral position, corresponding to the relaxation of the muscles.

The second dataset includes 13 gestures: the same 7 gestures described above, plus 6 additional classes. It contains gestures from 50 different subjects and two sets of trials per user. All 13 gestures are depicted in Figure 2. More details about the dataset can be found in [37].

We also validate our models in a real-context scenario. For the real-life predictions, we recorded the EMG activity associated with each gesture at forearm level using Myo armband. The information collected from each channel is transmitted to a computer via Bluetooth protocol where it is processed to extract relevant time domain features that will be used by the classifier to determine which gesture has been performed.

4.2. Proposed Architecture

The raw 8 channels EMG signal is split using a 250 ms sliding window, with 50% overlap. From each window of each channel a series of 8 time descriptors are extracted. The information from all the channels is then concatenated, forming a 64-dimensional vector. The 7-gestures dataset contains around 200k vector samples, while the 13-gestures dataset has around 59k vector samples. Those are split in training, validation, and test sets at user level according to the ratio: 70%, 20%, 10%. These vectors are fed to the network in mini-batches of size 2048. The considered architectures consists of a 6-hidden layer ($m = 6$) fully connected neural networks, with different parameters depending on the considered datasets, but the same core structure, as displayed in Figure 3. Let $x = (x_k)_{0 \leq k \leq K-1}$ be the vector of EMG samples acquired on a window from one channel. For this work we considered some of the most relevant features to describe sEMG data, as follows.

- (i) **Mean Absolute Value (MAV)** – the mean of the absolute values of the signal is given by

$$\text{MAV}(x) = \frac{1}{K} \sum_{k=0}^{K-1} |x_k|. \quad (22)$$

- (ii) **Zero Crossing Rate (ZCR)** – this feature counts the frequency at which the signal passes through zero. A threshold $\alpha \geq 0$ is used in order to lessen the noise effect. This feature can be computed in an incremental manner and it is defined as

$$\text{ZCR}(x) = \left| \left\{ k \in \{1, \dots, K-1\} \mid |x_k - x_{k-1}| \geq \alpha \text{ and } x_k x_{k-1} < 0 \right\} \right|. \quad (23)$$

- (iii) **Waveform Length (WL)** – this feature offers a simple characterization of the signal waveform. It corresponds to the following total variation seminorm:

$$\text{WL}(x) = \sum_{k=1}^{K-1} |x_k - x_{k-1}|. \quad (24)$$

- (iv) **Slope Sign Changes (SSC)** – measures the frequency at which the sign of the signal slope changes. It amounts in checking a condition on three consecutive samples x_k, x_{k-1}, x_{k+1} with $k \in \{2, \dots, K-2\}$:

$$\text{SSC}(x) = \left| \{k \in \{2, \dots, K-2\} \mid (x_k - x_{k-1})(x_k - x_{k+1}) \geq \alpha\} \right|, \quad (25)$$

where the threshold $\alpha > 0$ is employed to reduce the influence of the noise.

- (v) **Root Mean Square (RMS)** – this feature, also related to the quadratic mean or local energy of the signal is given by

$$\text{RMS}(x) = \sqrt{\frac{1}{K} \sum_{k=0}^{K-1} x_k^2}. \quad (26)$$

- (vi) **Hjorth parameters** – are a set of three features originally developed for characterizing electroencephalography signals and then successfully applied to sEMG signal recognition. The most relevant Hjorth activity parameter can be thought of as the integrated power spectrum and basically corresponds to the variance of the signal calculated as follows:

$$\sigma^2(x) = \frac{1}{K} \sum_{k=0}^{K-1} (x_k - \mu(x))^2, \quad (27)$$

where $\mu(x)$ represents the mean value of the signal. The standard deviation and $\text{RMS}(x)$ are equal when the mean of the signal is zero.

		Accuracy	75 %	80 %	85 %	90 %	95%
Lipschitz constant	$\tilde{\mathcal{C}}_i \cap \mathcal{D}_i$	$\tilde{\mathbb{P}}_{\tilde{\mathcal{C}}_i \cap \mathcal{D}_i}$	19.5	37.5	68.3	3.5×10^4	3.5×10^8
		$\mathbb{P}_{\tilde{\mathcal{C}}_i \cap \mathcal{D}_i}$	0.66	13.47	74.16	1.04×10^3	1.39×10^5
	$\check{\mathcal{C}}_{i,n} \cap \mathcal{D}_i$	$\tilde{\mathbb{P}}_{\check{\mathcal{C}}_{i,n} \cap \mathcal{D}_i}$	0.71	1.84	3.42	6.87	11.60
		$\mathbb{P}_{\check{\mathcal{C}}_{i,n} \cap \mathcal{D}_i}$	0.70	1.35	3.41	6.79	11.20
	$\mathcal{C}_{i,n} \cap \mathcal{D}_i$	$\tilde{\mathbb{P}}_{\mathcal{C}_{i,n} \cap \mathcal{D}_i}$	0.44	1.79	2.93	4.85	5.68
		$\mathbb{P}_{\mathcal{C}_{i,n} \cap \mathcal{D}_i}$	0.35	0.46	0.65	0.82	0.95

Table 1: Lipschitz constant obtained with various constrained optimization strategies for different accuracies – 7-gestures dataset.

(vii) **Skewness** – measures the overall asymmetry of probability distribution of the data:

$$\text{Skew}(x) = \frac{1}{K} \sum_{k=0}^{K-1} \left(\frac{x_k - \mu(x)}{\sigma(x)} \right)^3. \quad (28)$$

(viii) **Integrated Square-root EMG (ISEMG)** – is a feature returning the sum of the fully-rectified signal:

$$\text{ISEMG}(x) = \sum_{k=0}^{K-1} \sqrt{|x_k|}. \quad (29)$$

4.3. Performance analysis in terms of accuracy and robustness

The performance of our AGR system trained conventionally achieves state-of-art performance [43, 37], of over 99% accuracy for both datasets. Since in this case the weights are not guaranteed to be positive, the lower bound introduced in Proposition 2.2 does not constitute a valid Lipschitz constant. Computing the exact Lipschitz constant θ_m of the system is a very difficult task [28], but we can easily bound θ_m between the estimate given by (6) and the spectral norm of the product of all the weight ma-

		Accuracy	75 %	80 %	85 %	90 %	95%
Lipschitz constant	$\tilde{\mathcal{C}}_i \cap \mathcal{D}_i$	$\tilde{\mathbb{P}}_{\tilde{\mathcal{C}}_i \cap \mathcal{D}_i}$	20.2	41.8	145.2	2.2×10^5	1.21×10^{11}
		$\mathbb{P}_{\tilde{\mathcal{C}}_i \cap \mathcal{D}_i}$	0.85	20.47	112.3	1.62×10^4	2.31×10^8
	$\check{\mathcal{C}}_{i,n} \cap \mathcal{D}_i$	$\tilde{\mathbb{P}}_{\check{\mathcal{C}}_{i,n} \cap \mathcal{D}_i}$	0.84	2.08	4.23	7.54	12.02
		$\mathbb{P}_{\check{\mathcal{C}}_{i,n} \cap \mathcal{D}_i}$	0.81	2.01	4.12	7.50	11.92
	$\mathcal{C}_{i,n} \cap \mathcal{D}_i$	$\tilde{\mathbb{P}}_{\mathcal{C}_{i,n} \cap \mathcal{D}_i}$	0.54	1.87	3.38	4.20	5.78
		$\mathbb{P}_{\mathcal{C}_{i,n} \cap \mathcal{D}_i}$	0.49	0.53	0.75	0.92	1.25

Table 2: Lipschitz constant obtained with various constrained optimization strategies for different accuracies – 13-gestures dataset.

		Accuracy	75%	80%	85%	90%	95%
Model	empirical std		0.65	1.22	0.56	1.35	1.10
Variation	robust std		1.02	0.94	0.53	0.87	1.07

Table 3: Standard deviation of accuracy computed on 15 epochs, after convergence, on the test set for constrained models – 7-gestures dataset

trices from the network. We found that the Lipschitz constant estimate $\theta_m \in [1.56 \times 10^{12}, 1.59 \times 10^{14}]$ for both datasets. This suggests that despite the high performance of the classifiers, their robustness is poorly controlled, leaving the systems vulnerable to adversarial perturbations.

A first step towards controlling the Lipschitz constant of the classifica-

		Accuracy	75%	80%	85%	90%	95%
Model	empirical std		0.65	1.05	0.75	0.75	0.72
Variation	robust std		0.77	0.81	0.72	0.97	0.59

Table 4: Standard deviation of accuracy computed on 15 epochs, after convergence, on the test set for constrained models – 13-gestures dataset

		Accuracy	75%	80%	85%	90%	95%
Lipschitz constant	\mathcal{C}_i	$\tilde{\mathcal{P}}_{\tilde{\mathcal{C}}_i}$	72.03	127.5	1296	8.75×10^4	5.43×10^9
		$\mathcal{P}_{\tilde{\mathcal{C}}_i}$	52.06	102.49	905.45	7.23×10^4	8.14×10^8
		<i>Deel-lip</i> [39]	75.81	126.9	1283.6	8.70×10^4	5.43×10^9

Table 5: Lipschitz constant for networks trained with arbitrary signs – 7-gestures dataset.

		Accuracy	75%	80%	85%	90%	95%
Lipschitz constant	\mathcal{C}_i	$\tilde{\mathcal{P}}_{\tilde{\mathcal{C}}_i}$	76.59	125.20	1016	2.03×10^4	4.3×10^8
		$\mathcal{P}_{\tilde{\mathcal{C}}_i}$	61.22	99.74	740	1.26×10^4	6.7×10^7
		<i>Deel-lip</i> [39]	77.21	125.63	1120	2.04×10^4	4.5×10^8

Table 6: Lipschitz constant for networks trained with arbitrary signs – 13-gestures dataset.

tion algorithm and implicitly its robustness is to impose the nonnegativity condition associated with constraint \mathcal{D} . Training under such a nonnegativity constraint is shown to improve the network operation interpretability [27] and acts as a regularization, reducing overfitting. On the other hand, it can affect its approximation capability and potentially lead to a performance decay. Training the proposed system subject to constraint \mathcal{D} results in an overall accuracy of 96.92 % and 95.87% for the case of 7 and 13 classes, respectively. The performance decay was balanced by an increase in the robustness, since the Lipschitz constant, computed as indicated in Proposition 2.2, equals $\theta_m = 9.69 \times 10^{10}$ for 7 classes and $\theta_m = 9.73 \times 10^{10}$ for 13 classes. We observed that the accuracy reduction can be overcome by adding additional layers to the architecture. Indeed, we were able to obtain a similar accuracy to the baseline by adding an extra layer to the existing architecture and retraining both systems subject to \mathcal{D} . Furthermore, compared to the unconstrained models, we managed to maintain a high performance

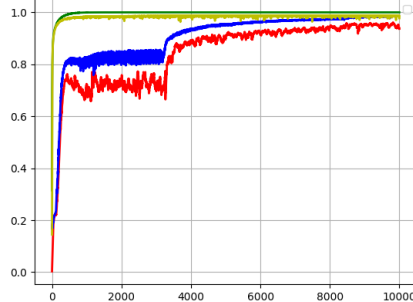


Figure 4: Accuracy vs. Iterations – constrained and unconstrained models in the context of 7-gesture dataset. The training and validation curves are displayed in green and yellow, respectively, for the unconstrained model. The training and validation curves are displayed in blue and red, respectively, in the case of constrained training, with the bound $\bar{\vartheta} = 0.95$.

while improving the robustness with respect to unconstrained training, i.e. $\theta_m = 1.02 \times 10^{11}$ for the 7-classes dataset and $\theta_m = 9.96 \times 10^{10}$ for the 13-classes dataset. We can however conclude from these tests that imposing the nonnegativity of the weight coefficients is not sufficient to reach satisfactory robustness.

To further control the robustness of the systems, we have to manage the Lipschitz constant of the networks by training them under an additional spectral norm constraints, as described by (11). Searching for the optimal accuracy robustness trade-off, we trained several models considering each of the three aforementioned constraints, namely $(\mathcal{C}_{i,n})_{1 \leq i \leq m, n \in \mathbb{N}}$ in (12), $(\tilde{\mathcal{C}}_i)_{1 \leq i \leq m}$ in (20), and $(\check{\mathcal{C}}_{i,n})_{1 \leq i \leq m, n \in \mathbb{N}}$ in (10). By adjusting the upper bound $\bar{\vartheta}$, we were able to assess the effect of a robustness constraint on the overall performance of the neural network-based classifiers, and finally to achieve the optimal trade-off. All our models were trained using Algorithm 1 as the optimizer.

The obtained results are summarized in Tables 1 and 2 for the two

Accuracy [%]						
Attack	robust model			baseline model		
	adversarial	non-adversarial		adversarial	non-adversarial	
FGSM [20]	91.75			76.48	71.21	
C&W ℓ_2 [40]	90.09	92.95		48.03	99.78	45.85
PGD [41]	91.92			59.36	56.38	
JSMA [42]	91.10			89.37	81.27	

Table 7: White-box attack results. We consider our best constrained model, having a Lipschitz constant $\theta = 0.97$ and two models trained conventionally: the best baseline [43] and another one having similar performance as the constrained one.

Movement	User #1		User #2		User #3		User #4		User #5		User #6		User #7		User #8		User #9		User #10	
	C	U	C	U	C	U	C	U	C	U	C	U	C	U	C	U	C	U	C	U
up	2	2	1	3	0	0	0	1	0	0	0	2	0	2	1	2	0	2	1	3
down	1	1	0	2	2	3	0	0	2	4	1	0	2	3	1	1	0	1	0	1
right	0	4	0	0	0	1	0	1	1	1	0	2	0	0	0	0	0	1	1	2
left	3	5	1	4	0	1	0	1	2	5	0	0	0	1	2	3	1	2	0	1
fist	0	2	2	4	0	0	1	0	0	3	0	1	1	1	0	2	1	1	1	3
spread	0	3	2	5	3	4	2	4	1	0	0	0	1	2	1	0	0	1	0	3
Sum	6	17	6	18	5	9	3	7	6	13	1	5	4	9	6	7	2	8	3	3
Error rate (%)	5	14	5	15	4.1	7.5	2.5	5.7	5	10.7	0.7	4.1	3.3	7.5	5	5.7	1.6	6.6	2.5	10.7

Table 8: Experiment results

datasets, respectively. As expected, obtaining a good robustness-accuracy trade-off requires paying attention to the way we design our constrained networks. We show that using tight constraints to approximate the Lipschitz constant accurately has a positive impact on the overall performance of both networks. We also evaluated the use of an inexact projection, designated by \tilde{P} (see Section 3.3). It can be observed that using an exact projection yields significantly better results. By combining tight constraints and exact projection techniques, we observe that the robustness of the network can be properly ensured while keeping a good accuracy in both cases. Indeed, we

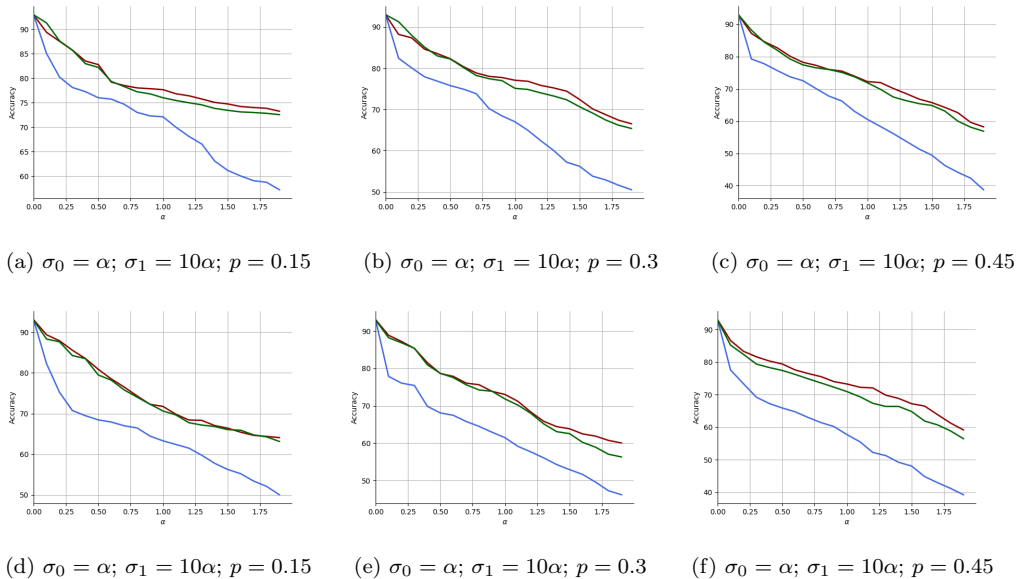


Figure 5: Accuracy vs. α in the context Noisy Inputs training. First row: 7-gesture dataset; Second row: 13-gestures dataset. Red line: robust model; Blue line: baseline model; Green line: adversarial trained model

succeeded in ensuring a Lipschitz constant around 1 for a 95% accuracy. The observed loss in accuracy with respect to a standard training is consistent with the “no free lunch theorem” [44].

Training neural networks subject to tight spectral norm constraints can be challenging,² and the cost of obtaining a good performance is the training time. Figure 4 shows the training curves for both validation and training sets in the context of the unconstrained baseline model (yellow and green lines), and in the case of training a constrained version (red and blue lines) using the optimal projection $\mathbb{P}_{\mathcal{C}_{i,n} \cap \mathcal{D}_i}$, with $\bar{\vartheta}_m = 0.95$. Even though it requires more

²A code in TensorFlow will be made available upon the acceptance of the paper.

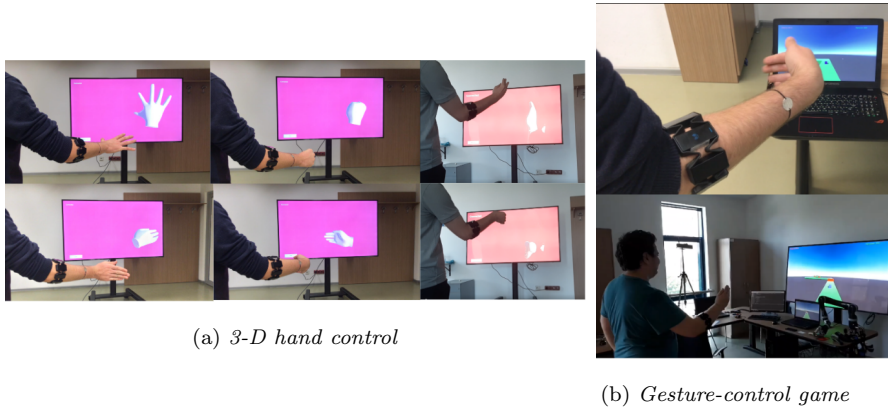


Figure 6: *Real-life experimental setup*

iterations, the constrained model is capable of reaching an accuracy comparable with the baseline, while providing a robustness certificate. Since the training curves may show some slight variations, we measured the accuracy variations in two ways: by computing the classical *standard deviation* (*std*), and by employing *median absolute deviation* (*mad*). For a vector $(x_i)_{1 \leq i \leq I}$, it is expressed as $\text{MAD} = \text{median}(|x_i - \zeta(x)|_{1 \leq i \leq I})$, where $\zeta(x)$ represents the median of the vector components. From this quantity, we can derive an empirical estimate of the standard deviation by multiplying MAD with a factor equal to 1.4826. The latter estimate is known to be more robust to outliers for Gaussian distributed data, especially in the case of small populations. The results are summarized in Tables 3 and 4. It can be observed that the empirical standard deviation is below 1.6% and the robust estimate of it is below 1.1% for both datasets.

We have also evaluated how the positivity constraint impacts the overall accuracy of our system. We trained a robust network by allowing the weights to have arbitrary signs. For this purpose, we control individually the Lips-

chitz constant of each layer $i \in \{1, \dots, m\}$ to be less than a given value $\bar{\vartheta}^{1/m}$. The exact projection onto $\tilde{\mathcal{C}}_i$, $\mathbb{P}_{\tilde{\mathcal{C}}_i}$, as well as the approximate one $\tilde{\mathbb{P}}_{\tilde{\mathcal{C}}_i}$ were computed as described previously. In this case $\bar{\vartheta}$ represents an upper bound on the Lipschitz constant of the system. Tables 5 and 6 summarize the results for different values of $\bar{\vartheta}$, for the two respective datasets. We compare our method for dealing with Lipschitz constraints with the approach proposed in [39]. This approach which is implemented in the *deel-lip* library allows the user to train robust networks in a convenient manner, offering a robustness certificate by performing a spectral normalization for each layer. It can be observed on these datasets that our method yields similar results when using the approximate projection, but better ones when using the exact projection.

5. Robustness validation

In this section we investigate to what extent the theoretical concepts described in the previous sections help in improving the robustness of the classifier in different settings. To this goal, we consider the following three scenarios. In the first one, we examine the impact of adversarial attacks on the performance of the classifier. The second scenario takes into account the effect of noise in the acquisition process. In the case of sEMG signals, this noise may come from imperfect skin-sensor contact caused by hairs or drops of sweat. In the last scenario we perform a real-life experiment using 10 able-bodied volunteers.

5.1. Sensitivity to adversarial attacks

We evaluate our robust model on purposely designed perturbations, by studying their influence on the overall performance of the system. We lead at-

tacks on our best robust model in terms of accuracy and robustness achieving 92.95% accuracy and a Lipschitz constant $\bar{\vartheta} = 0.87$ for the 7-gesture dataset. We compare the results with two conventionally trained models: the best one in terms of performance, which achieves 99.78% prediction accuracy on non-adversarial data, and another one trained to have similar performance as our robust model reaching 92.99% accuracy on the original test set.

To create the adversarial samples we used some of the most popular white-box attackers, namely: *Fast gradient sign method (FGSM)* [20] – generates adversarial data based on the gradient of the cost function with respect to the input data; *Jacobian Saliency Map Attacker (JSMA)* [42] – computes a perturbation based on ℓ_0 distance metric by iteratively selecting the input sample that will increase the chance of miss-classification; *Projected gradient descent (PGD)*[41] – uses local first order information about the network to create adversarial examples; *Carlini and Wagner (C&W)* [40] – utilizes ℓ_2 distance to compute the optimal adversarial perturbation.

The results summarised in Table 7 show the performance obtained for the 7-gesture test set. Note that the robust model performance is barely affected by the adversarial perturbations, whereas the baseline models show a huge drop in accuracy. It can be concluded that our goal of building a robust AGR system is achieved. It also proves that precisely controlling the Lipschitz constant of the system contributes greatly in improving its robustness against adversarial samples.

5.2. Noisy input behaviour

To simulate the effect of underlying noise generated during the acquisition process, we added synthetic noise directly to the raw sEMG data, prior to

the feature extraction step. The noise is chosen independent and identically distributed according to a Gaussian mixture law $(1-p)\mathcal{N}(0, \sigma_0^2) + p\mathcal{N}(0, \sigma_1^2)$. The mixture comprises a background component, corresponding to the intrinsic electronic noise in the armband, such as thermal or quantization noise, and an impulsive component accounting for outliers. Those may be related to imperfect wiring that can generate impulse-like artifacts. In our experiments, we consider background and impulse noises with standard deviations $\sigma_0 = \alpha$ and $\sigma_1 = 10\alpha$ with $\alpha \in [0, +\infty[$. We generate different levels of noise, by varying the parameter α . The probability of peaks $p \in [0, 1]$ is also adjusted to simulate more or less severe scenarios in terms of outliers.

From the resulting noisy signals, we extract the features described in Section 4 and pass them to the classifier, using our robust models reaching an accuracy of 92.95% ($\bar{\vartheta} = 0.87$) for the 7-gestures dataset, and 93.05% ($\bar{\vartheta} = 0.98$) in the case of the 13-gestures dataset, trained with non-altered data. We compared the results achieved with our robust training with those obtained with i) classical training and ii) adversarial training. The adversarial training was performed by generating an extended dataset, containing the original data and corrupted versions of them by additive noise following the Gaussian mixture law described above, where the parameters p and α were drawn randomly in a uniform manner on $[0.15, 0.45]$ and $[0, 2]$, respectively. In the absence of noise, a similar performance in terms of accuracy was obtained: 7-gestures dataset – 92.99%, and 92.97%, 13-gestures dataset – 93.03% and 92.98% for baseline and the adversarial training, respectively.

The experimental results obtained on both datasets are depicted in Figure 5. The red, blue, and green lines correspond to the unconstrained, con-

strained, and adversarial models, respectively. We observe that the constrained model is significantly less affected by the presence of noise in the inputs than the one trained without robustness guarantees. It is also worth noting that training with adversarial inputs also leads to satisfactory results, although usually slightly less accurate. The Lipschitz lower and upper bounds computed for the networks trained in an adversarial manner are indeed much lower than those with standard training, but they remain quite large ((1845.23, 79534.2) for 7-gestures dataset and (1754.74, 64595.8) for 13-gestures dataset).

This experiment emphasizes that controlling the Lipschitz constant of a network improves its robustness not only against targeted adversarial attacks, as shown previously, but also in the case of black-box attacks, where no prior information about the model is used.

5.3. *Real-life scenario validation*

To illustrate the practical applicability of our findings, we proceed to validate our model in a real-life context. For this purpose, we designed an experiment to compare a conventionally trained model with the constrained one. We integrated both models in a real-time application that controls a 3D hand on a screen, as well as a game that can be controlled by gestures, to give the user a tangible feedback. We asked 10 volunteers (males and females) to test both models by performing each gesture 20 times. We emphasize that the user had no prior knowledge about what model was implemented, since it was randomly selected at the beginning of each new trial. Pictures of the experimental setup are provided in Figure 6. Table 8 details on a user level, how many (out of the 20) trials were erroneously classified. U and C

denote the Unconstrained and the Constrained models, respectively. Note that, despite obtaining very good results on the test set, the unconstrained model loses a lot in terms of performance (up to 15%) when facing real-life data. We can observe that training a positive neural networks subject to Lipschitz constraints improves the overall robustness of the classifier against adversarial perturbations, not only from a theoretical viewpoint, but also practically by leading to more reliable systems with greater generalization power.

6. Conclusion

This work has shown the usefulness of designing robust feed-forward neural networks for automatic gesture recognition based on sEMG physiological signals. More precisely, we proposed to finely control the Lipschitz constant of these nonlinear systems by considering positively weighted neural architectures. To offer robustness certificates, we also developed new optimization techniques for training classifiers subject to spectral norm constraints on the weights. We studied various constrained formulations and showed that robustness can be secured without sacrificing accuracy when using a combination of tight constraints and exact projections. Experiments on two distinct datasets illustrated the good performance of our approach. We further demonstrated the effectiveness of our robust classifier, compared to classically trained ones, when facing white-box and black-box attacks, as well as in real-life usage.

In future works, it would be interesting to apply such a robust training procedure to other applications in pattern recognition involving data acquired in real-time.

Appendix – Accelerated DFB algorithm

Let $n \in \mathbb{N} \setminus \{0\}$ and $i \in \{1, \dots, m\}$. Computing the projection of a matrix $\overline{W}_i \in \mathbb{R}^{N_i \times N_{i-1}}$ onto $\mathcal{D}_i \cap \mathcal{C}_{i,n}$ is equivalent to solve the following matrix optimization problem:

$$\underset{W_i \in \mathbb{R}^{N_i \times N_{i-1}}}{\text{minimize}} \quad \iota_{\mathcal{D}_i}(W_i) + \iota_{\mathcal{B}(0, \overline{\vartheta})}(A_{i,n}W_iB_{i,n}) + \frac{1}{2}\|W_i - \overline{W}_i\|_{\mathbb{F}}^2 \quad (30)$$

where $\|\cdot\|_{\mathbb{F}}$ is the Frobenius norm and $\iota_{\mathcal{S}}$ denotes the indicator of a set \mathcal{S} (this function is equal to 0 on this set and $+\infty$ otherwise.) The dual optimization problem associated to this strongly convex minimization problem reads

$$\underset{Y \in \mathbb{R}^{N_m \times N_0}}{\text{minimize}} \quad f^*(-A_{i,n}^\top Y B_{i,n}^\top) + \iota_{\mathcal{B}(0, \overline{\vartheta})}^*(Y), \quad (31)$$

where for a given function g , g^* denotes its Fenchel-Legendre conjugate. In our case $f = \iota_{\mathcal{D}_i} + \frac{1}{2}\|\cdot - \overline{W}_i\|_{\mathbb{F}}^2$. From standard conjugation rules [35], f^* is equal to

$$(\forall W_i \in \mathbb{R}^{N_i \times N_{i-1}}) \quad f^*(W_i) = \widetilde{\iota}_{\mathcal{D}_i}(W_i + \overline{W}_i), \quad (32)$$

where $\widetilde{\iota}_{\mathcal{D}_i}$ is the Moreau envelope of $\iota_{\mathcal{D}_i}^*$ given by

$$\widetilde{\iota}_{\mathcal{D}_i}(W_i) = \inf_{W'_i \in \mathbb{R}^{N_i \times N_{i-1}}} \iota_{\mathcal{D}_i}^*(W'_i) + \frac{1}{2}\|W'_i - W_i\|_{\mathbb{F}}^2. \quad (33)$$

The Moreau envelope of a proper lower-semicontinuous convex function is differentiable. Thus f^* is differentiable and its gradient is [45, Example 17.33]

$$\nabla f^*(W_i) = \text{P}_{\mathcal{D}_i}(W_i + \overline{W}_i). \quad (34)$$

We deduce that the gradient of $Y \mapsto f^*(-A_{i,n}^\top Y B_{i,n}^\top)$ is

$$-A_{i,n} \text{P}_{\mathcal{D}_i}(\overline{W}_i - A_{i,n}^\top Y B_{i,n}^\top) B_{i,n}.$$

Since $\mathbf{P}_{\mathcal{D}_i}$ is a nonexpansive operator, the latter function has a Lipschitz gradient with constant $\beta = \|A_{i,n}\|_{\mathbb{S}}^2 \|B_{i,n}\|_{\mathbb{S}}^2$. The dual problem (31) thus corresponds to the minimization of the sum of a smooth convex function and a proper lower-semicontinuous function. Consequently, it can be minimized by a proximal algorithm. Such a strategy will require to calculate the proximity operator of $\gamma \mathcal{L}_{\mathcal{B}(0,\bar{\vartheta})}^*$ for some scaling parameter $\gamma \in]0, +\infty[$. By using Moreau's formula [45], this proximity operator is expressed as

$$(\forall Y \in \mathbb{R}^{N_m \times N_0}) \quad \text{prox}_{\gamma \mathcal{L}_{\mathcal{B}(0,\bar{\vartheta})}^*}(Y) = Y - \gamma \mathbf{P}_{\mathcal{B}(0,\bar{\vartheta})}(\gamma^{-1}Y). \quad (35)$$

A classical solution for solving the dual problem consists in using the standard forward-backward algorithm [46, 30]. This leads to Algorithm 3 [36]. Another solution consists in using the FISTA-like algorithm in [47], which leads to the accelerated version in Algorithm 2. The sequences $(Y_\ell)_{\ell \in \mathbb{N}}$ generated by these two algorithms is guaranteed to converge to a solution \widehat{Y} to the dual problem. In addition, from Kuhn-Tucker conditions, the solution to the primal problem $\widehat{W}_i = \mathbf{P}_{\mathcal{S}_{i,n}}(\overline{W}_i)$ is equal to $\nabla f^*(-A_{i,n}^\top \widehat{Y} B_{i,n}^\top)$. It follows from (34) and the continuity of $\mathbf{P}_{\mathcal{D}_i}$ that the sequence $(V_\ell)_{\ell \in \mathbb{N}}$ converges to \widehat{W}_i .

Algorithm 3: Dual Forward-backward algorithm

Let $Y_0 \in \mathbb{R}^{N_m \times N_0}$

Set $\epsilon \in]0, 1/(\|A_{i,n}\|_{\mathbb{S}}\|B_{i,n}\|_{\mathbb{S}})^2[$

for $l = 0, 1, \dots$ **do**

Set $\gamma_\ell \in [\epsilon, 2/(\|A_{i,n}\|_{\mathbb{S}}\|B_{i,n}\|_{\mathbb{S}})^2 - \epsilon]$ $V_\ell = \mathbf{P}_{\mathcal{D}_i}(\overline{W}_i - A_{i,n}^\top Y_\ell B_{i,n}^\top)$
 $\widetilde{Y}_\ell = Y_\ell + \gamma_\ell A_{i,n} V_\ell B_{i,n}$ $Y_{\ell+1} = \widetilde{Y}_\ell - \gamma_\ell \mathbf{P}_{\mathcal{B}(0,\bar{\vartheta})}(\gamma_\ell^{-1} \widetilde{Y}_\ell)$

Acknowledgements

This project was supported by the BRIDGEABLE ANR Research and Teaching Chair in AI and by the Romanian National Authority for Scientific Research and Innovation, UEFISCDI, project EGMP-AI, agreement 425PED/2020, grant PN-III-P2-2.1-PED-2019-2392.

References

- [1] J. He, S. L. Baxter, J. Xu, J. Xu, X. Zhou, K. Zhang, The practical implementation of artificial intelligence technologies in medicine, *Nat. med.* 25 (1) (2019) 30–36.
- [2] P. Li, X. Chen, S. Shen, Stereo R-CNN based 3D object detection for autonomous driving, in: *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, Long Beach, CA, USA, 2019, pp. 7644–7652.
- [3] N. Widiastuti, Convolution neural network for text mining and natural language processing, in: *IOP Conf.: Mater. Sci. Eng.*, Vol. 662, Kazimierz Dolny, Poland, 2019.
- [4] A. Kurakin, Z. Zhang, Z. Liu, A real time system for dynamic hand gesture recognition with a depth sensor, in: *Proc. IEEE European Signal Processing Conf.*, Bucharest, Romania, 2012, pp. 1975–1979.
- [5] J. Qi, G. Jiang, G. Li, Y. Sun, B. Tao, Intelligent human-computer interaction based on surface EMG gesture recognition, *IEEE Access* 7 (2019) 61378–61387.
- [6] M. J. Cheok, Z. Omar, M. H. Jaward, A review of hand gesture and sign language recognition techniques, *Int. J. Mach. Learn. Cyber.* 10 (1) (2019) 131–153.
- [7] F. Wen, Z. Sun, T. He, Q. Shi, M. Zhu, Z. Zhang, L. Li, T. Zhang, C. Lee, Machine learning glove using self-powered conductive superhydrophobic triboelectric textile for gesture recognition in VR/AR applications, *Adv. Sci.* 7 (14) (2020) 2000261.
- [8] X. Gao, Y. Jin, Q. Dou, P.-A. Heng, Automatic gesture recognition in robot-assisted surgery with reinforcement learning and tree search, in: *Proc. IEEE Int. Conf. Robot. Autom.*, Paris, France, 2020, pp. 8440–8446.
- [9] O. Kopuklu, N. Kose, G. Rigoll, Motion fused frames: Data level fusion strategy for hand gesture recognition, in: *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, Utah, USA, 2018, pp. 2103–2111.
- [10] T. J. Darrell, I. A. Essa, A. P. Pentland, Task-specific gesture analysis in real-time using interpolated views, *IEEE Trans. Pattern Anal. Mach. Intell.* 18 (12) (1996) 1236–1242.
- [11] N. Neverova, C. Wolf, G. Taylor, F. Nebout, Moddrop: adaptive multi-modal gesture recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 38 (8) (2015) 1692–1706.
- [12] M. Atzori, A. Gijsberts, C. Castellini, B. Caputo, A.-G. M. Hager, S. Elsig, G. Giatsidis, F. Bassetto, H. Müller, Electromyography data for non-invasive naturally-controlled robotic hand prostheses, *Sci. data* (140053) (23 Dec 2014).

- [13] R. N. Khushaba, S. Kodagoda, Electromyogram (EMG) feature reduction using mutual components analysis for multifunction prosthetic fingers control, in: *Int. Conf. Control Autom. Robotics & Vision*, Guangzhou, China, 2012, pp. 1534–1539.
- [14] J. Kim, S. Mastnik, E. André, EMG-based hand gesture recognition for realtime biosignal interfacing, in: *Proc. Int. Conf. Intell. User Interfac.*, Canaria, Spain, 2008, pp. 30–39.
- [15] A. D. Orjuela-Cañón, A. F. Ruíz-Olaya, L. Forero, Deep neural network for EMG signal classification of wrist position: Preliminary results, in: *IEEE Latin American Conf. Comput. Intell.*, Arequipa, Peru, 2017, pp. 1–5.
- [16] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, R. Fergus, Intriguing properties of neural networks, in: *Proc. Int. Conf. Learn. Represent.*, Banff, Canada, 2014.
- [17] A. Kurakin, I. Goodfellow, S. Bengio, Adversarial machine learning at scale, in: *Int. Conf. Learn. Represent.*, Toulon, France, 2017.
- [18] N. Carlini, P. Mishra, T. Vaidya, Y. Zhang, M. Sherr, C. Shields, D. Wagner, W. Zhou, Hidden voice commands, in: *USENIX Security Symp.*, Austin, TX, USA, 2016, pp. 513–530.
- [19] M. Takeru, K. Toshiki, K. Masanori, Y. Yuichi, Spectral normalization for generative adversarial networks, in: *Int. Conf. Learn. Represent.*, Vancouver, Canada, 2018.
- [20] I. J. Goodfellow, J. Shlens, C. Szegedy, Explaining and harnessing adversarial examples, in: *Int. Conf. Learn. Represent.*, San Diego, CA, USA, 2015.
- [21] A. Neacșu, J.-C. Pesquet, C. Burileanu, Accuracy-robustness trade-off for positively weighted neural networks, in: *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, Barcelona, Spain, 2020, pp. 8389–8393.
- [22] K. Scaman, A. Virmaux, Lipschitz regularity of deep neural networks: analysis and efficient estimation, in: *Proc. Ann. Conf. Neur. Inform. Proc. Syst.*, Montreal, Canada, 2018, pp. 3839–3848.
- [23] C. Anil, J. Lucas, R. Grosse, Sorting out Lipschitz function approximation, in: *Proc. Int. Conf. Mach. Learn.*, Long Beach, California, USA, 2019, pp. 291–301.
- [24] M. Cisse, P. Bojanowski, E. Grave, Y. Dauphin, N. Usunier, Parseval networks: Improving robustness to adversarial examples, in: *Proc. Int. Conf. Mach. Learn.*, Sydney, NSW, Australia, 2017, pp. 854–863.
- [25] M. Fazlyab, A. Robey, H. Hassani, M. Morari, G. Pappas, Efficient and accurate estimation of lipschitz constants for deep neural networks, in: *Adv. Neural Info. Process. Syst.*, Vancouver, Canada, 2019, pp. 11423–11434.
- [26] P. Pauli, A. Koch, J. Berberich, P. Kohler, F. Allgower, Training robust neural networks using Lipschitz bounds, *IEEE Control Syst. Lett.* 6 (2021) 121–126.
- [27] J. Chorowski, J. M. Zurada, Learning understandable neural networks with nonnegative weight constraints, *IEEE Trans. Neural Netw. Learn. Syst.* 26 (1) (2015) 62–69.
- [28] P. L. Combettes, J.-C. Pesquet, Lipschitz certificates for layered network structures driven by averaged activation operators, *SIAM J. on Math. Data Sci.* 2 (4) (2020) 529–557.

- [29] P. Combettes, J.-C. Pesquet, Deep neural network structures solving variational inequalities, *Set-Valued Var. Anal.* (2020) 1–28.
- [30] P. Combettes, J.-C. Pesquet, Proximal thresholding algorithm for minimization over orthonormal bases, *SIAM J. on Optim.* 18 (4) (2008) 1351–1376.
- [31] P. Ramachandran, B. Zoph, Q. V. Le, Searching for activation functions, *Proc. Int. Conf. Learn. Represent.* (30 Apr.–03 May 2018).
URL <https://arxiv.org/pdf/1710.05941.pdf>
- [32] F. Latorre, P. Rolland, V. Cevher, Lipschitz constant estimation of neural networks via sparse polynomial optimization, in: *Int. Conf. on Learning Representations*, 2020.
URL https://openreview.net/forum?id=rJe4_xSFDB
- [33] T. Chen, J.-B. Lasserre, V. Magron, E. Pauwels, Semialgebraic optimization for Lipschitz constants of ReLU networks (2020) 19189–19200.
- [34] P. L. Combettes, J.-C. Pesquet, Fixed point strategies in data science, *IEEE Trans. Sig. Proc.* 69 (2021) 3878–3905.
- [35] N. Komodakis, J.-C. Pesquet, Playing with duality: An overview of recent primal-dual approaches for solving large-scale optimization problems, *IEEE Signal Process. Mag.* 32 (6) (2015) 31–54.
- [36] P. L. Combettes, D. Düng, B. C. Vũ, Dualization of signal recovery problems, *Set-Valued Var. Anal.* 18 (3–4) (2010) 373–404.
- [37] C. Andronache, M. Negru, A. Neacșu, G. Cioroiu, A. Rădoi, C. Burileanu, Towards extending real-time EMG-based gesture recognition system, in: *Proc. IEEE Int. Conf. Telecomm. Signal Process.*, Milan, 2020, pp. 301–304.
- [38] U. Côté-Allard, C. L. Fall, A. Drouin, A. Campeau-Lecours, C. Gosselin, K. Glette, F. Laviolette, B. Gosselin, Deep learning for electromyographic hand gesture signal classification using transfer learning, *IEEE Trans. Neural Syst. Rehabilitation Eng.* 27 (4) (2019) 760–771.
- [39] M. Serrurier, F. Mamalet, A. González-Sanz, T. Boissin, J.-M. Loubes, E. del Barrio, Achieving robustness in classification using optimal transport with hinge regularization, in: *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, Nashville, USA, 2021, pp. 505–514.
- [40] N. Carlini, D. Wagner, Towards evaluating the robustness of neural networks, in: *IEEE Symp. Security Privacy*, San Jose, CA, USA, 2017, pp. 39–57.
- [41] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu, Towards deep learning models resistant to adversarial attacks, in: *Proc. Int. Conf. Learning Representations*, Vancouver, BC, Canada, 2018.
- [42] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, A. Swami, The limitations of deep learning in adversarial settings, in: *IEEE Symp. Security Privacy*, Saarbrücken, Germany, 2016.
- [43] A. A. Neacsu, G. Cioroiu, A. Radoi, C. Burileanu, Automatic EMG-based hand gesture recognition system using time-domain descriptors and fully-connected neural networks, in: *Int. Conf. Telecommunications Signal Process.*, Budapest, Hungary, 2019, pp. 232–235.
- [44] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, A. Madry, There is no free lunch in adversarial robustness (but there are unexpected benefits), *CoRR*, abs/1805.12152, 2018 (2018).

- [45] H. Bauschke, P. Combettes, *Convex analysis and monotone operator theory in Hilbert spaces*, 2019.
- [46] P. L. Combettes, V. R. Wajs, Signal recovery by proximal forward-backward splitting, *Multiscale Model. Simul.* 4 (4) (2005) 1168–1200.
- [47] A. Chambolle, C. Dossal, On the convergence of the iterates of the fast iterative shrinkage/thresholding algorithm, *J. Optim. Theory Appl.* 166 (3) (2015) 968–982.