

# A simple and efficient architecture for trainable activation functions

Andrea Apicella, Francesco Isgrò and Roberto Prevete

Dipartimento di Ingegneria Elettrica e delle Tecnologie dell'Informazione  
Università di Napoli Federico II

## Abstract

Learning automatically the best activation function for the task is an active topic in neural network research. At the moment, despite promising results, it is still difficult to determine a method for learning an activation function that is at the same time theoretically simple and easy to implement. Moreover, most of the methods proposed so far introduce new parameters or *adopt* different learning techniques. In this work we propose a simple method to obtain trained activation function which adds to the neural network local subnetworks with a small amount of neurons. Experiments show that this approach could lead to better result with respect to using a pre-defined activation function, without introducing a large amount of extra parameters that need to be learned.

**keywords:** neural networks, machine learning, activation functions, adaptive activation functions

## 1 Introduction

The success of deep learning approaches has led to an increase in interest in MultiLayer FeedForward (MLFF) neural networks. MLFF networks are composed of  $N$  elementary computing units (neurons), which are organized in  $L > 1$  layers. The first layer of a MLFF network is composed of  $d$  input variables. Each neuron  $i$  belonging to a layer  $l$ , with  $l > 1$ , **receives possibly connections** from all the neurons (or input variables in case of  $l = 2$ ) of the previous layer  $l - 1$ . Each connection is associated to a real value called *weight*. The flow of computation proceeds from the first layer to the last layer (*forward propagation*). The last neuron layer is called *output* layer, the remaining neuron layers are called *hidden* layers. The computation of a neuron  $i$  belonging to the layer  $l$  corresponds to a two-step process: first is computed the neuron activation  $a_i^l$  and then the neuron output  $z_i^l$ . The neuron activation  $a_i^l$  is usually constructed as a linear combination of the outputs of the previous layer:  $a_i^l = \sum_j w_{ij}^l z_j^{l-1} + b_i^l$  where  $w_{ij}^l$  is the weight of the connection going to the neuron  $j$  belonging to the layer  $l - 1$  to the neuron  $i$ ,  $b_i^l$  is a parameter said *bias*,  $l = 1, \dots, L$  and  $j$  runs on the indexes of the neurons of the layer  $l - 1$  which send connections to the neuron  $i$ . If  $l = 2$  the variables  $z^{l-1}$  correspond to the input variables. The neuron output  $z_i^l$  is usually computed by a differentiable, non linear activation function  $f(\cdot)$ :  $z_i^l = f(a_i^l)$ . The nonlinear functions  $f(\cdot)$  are generally chosen as simple *fixed* functions such as the *logistic sigmoid* or the *tanh* function.

Given a MLFF network with  $d$  input variables and  $c$  neurons in the output layer, it achieves a functional mapping from a  $d$ -dimensional space to a  $c$ -dimensional space. Thus, a MLFF network can be interpreted as a non-linear parametric function  $\mathbf{y} = \text{Net}(\mathbf{x}; \boldsymbol{\theta})$ , where the parameters  $\boldsymbol{\theta}$  are all the weights and biases of the network and  $\mathbf{y}$  is the response of the output layer. The approximation properties of MLFF networks have been widely studied [DeVore et al., 1996]. In a nutshell, a function approximation problem can be summarized as follows [Bishop, 2006, Ripley, 2007]: given an unknown function  $F : \mathbf{x} \in R^d \rightarrow \mathbf{y} = F(\mathbf{x}) \in R^c$  and a data set  $\{(\mathbf{x}^n, \mathbf{t}^n)\}_{n=1}^N$  representing a sampling of the unknown function, where  $\mathbf{t}^n = F(\mathbf{x}^n) + \epsilon$ , usually called *targets*, are the values assumed by  $F$  in  $\mathbf{x}^n$  with added an unknown noise  $\epsilon$ , the task is to find the appropriate values of the parameters of a parametric function  $M$  so as to get as close as possible to the unknown function  $F$ . In this context, there are two different problems, the first one concerns the expressive power of the parameterized function  $M$ , that is, if there are parameter values for which it is possible to approximate the unknown function  $F$ , the second one concerns the possibility of actually finding such parameter values. Interestingly, regarding the first problem, a MLFF network with a single hidden

layer, which is usually called *shallow* network, can approximate arbitrarily well any functional continuous mapping defined on a compact input domain, provided the number of hidden neurons is sufficiently large and the activation functions of the hidden neurons satisfy suitable properties, for example, to be sigmoidal or, more in general, not-polynomial functions [Sonoda and Murata, 2017, Bishop, 2006, Pinkus, 1999]. In other words, given a certain desired degree of approximation, it exists a set of parameters  $\bar{\theta}$  for which the neural network  $Net(\mathbf{x}; \bar{\theta})$  approaches the unknown function within this degree of approximation, supposed to have a sufficient number of hidden neurons and appropriate activation functions. In this sense, MLFF networks are said to be *universal approximators*.

However, the *key problem* is how to find these suitable network parameters, i.e., weights and biases. The process to determine the values of weights and biases on the basis of the data set is called *learning* or *training*. Importantly, although the choice of the non linear activation functions  $f(\cdot)$  does not affect the MLFF network’s universal approximator property, provided certain constraints are satisfied, this choice becomes a key aspect when network weights and biases are to be found during the training process. To clarify this aspect, let us briefly summarize what is a training process. The training process generally corresponds to the minimization of an *error function* with respect to the network parameters. The error function typically assumes the following form (although many other forms are possible [Bishop, 2006]):

$$E(\theta) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^c [y_k(\mathbf{x}^n; \theta) - t_k^n]^2$$

where  $y_k(\mathbf{x}^n; \theta)$  represents the output of the neuron  $k$  belonging to the output layer as a function of both the input  $\mathbf{x}^n$  and the network parameters  $\theta$ . The quantity  $t_k^n$  represents the target value for output neuron  $k$  when the input is  $\mathbf{x}^n$ . The solution for the network parameter values at the *global minimum* of the error function is usually found by iterating a gradient-based algorithm with the gradient computed through backpropagation [Bishop, 2006]. Since for MLFF networks the error function typically will be a highly non-linear function of the parameters (not-convex surface), there may exist many *local minima* or *saddlepoints*. Notice that parameter regions where the error function is very “flat” can mimick local minima insofar as the learning process is “trapped” for very long periods of time. In a learning process the main difficult is to avoid these stationary points or regions of the error function. Thus, the choice of the activation functions has a relevant impact on the shape of the error function and, consequently, on the performance of the learning process. Moreover, this choice can affect the number of hidden neurons and layers necessary to reach the desired degree of approximation [Guliyev and Ismailov, 2016, Eldan and Shamir, 2016].

For these reasons, recently, there is a very large literature proposing activation functions that differ from those standards such as sigmoid and tanh. In particular, the introduction of activation functions as ReLU [Nair and Hinton, 2010] and similar functions, such as Leaky ReLU and parametric ReLU described respectively in [Maas et al., 2013] and [He et al., 2015], has contributed to renew the interest of the scientific community for MLFF networks. The use of these new activation functions has been shown to improve significantly the networks in terms of performance and training speed, thanks to properties as no saturation, e.g. [Glorot and Bengio, 2010]. Another great improvement was given in [Clevert et al., 2015], where the learning is speeded up introducing the ELU activation function, and more recently in [Klambauer et al., 2017] with the introduction of SELU units.

Thus, finding alternative functions that can potentially improve further the results is still an *open field of research*. Consistent with this perspective, a number of recent papers compare neural architectures with different activation functions, as in [Pedomonti, 2018], or propose to search appropriate activation functions within a finite set of potentially interesting activation functions, as in [Ramachandran et al., 2018]. However, a very recent field of research focuses on the possibility to learn appropriate activation functions from data, in this way one has *adaptable* (or *trainable*) activation functions which are adjusted during the learning phase towards specific functions, allowing the network to exploit the data better (see, for example, [Qian et al., 2018]). Furthermore, any layer of the network could potentially have their own best activation function, increasing the number of choices to make at the design stage. On the other side, it is not guaranteed that fixing the same function for each layer is the best choice. Thus, a way to tackle the problem is to learn the activation functions from data, together with the other parameters of the network; the idea is to find the *good* activation functions that, together with the other network parameters, provides a *good* model for the data.

In this paper we introduce a new method for learning activation functions in the context of full-connected and convolutional MLFF neural networks. The impact of this method on the performance of the network

are experimentally assessed. The idea is built upon the possibility to obtain adaptable activation functions in terms of sub-networks with just one hidden layer. In a nutshell, each neuron with a non-linear activation function  $f$  can be *substituted* with a neuron with an *Identity* activation function which sends its output to a one-hidden layer sub-network with just one output neuron. This substitution enables us to obtain “any” activation function  $f$ , since an one-hidden layer neural network can approximate arbitrarily well any functional continuous mapping from one finite-dimensional space to another, provided the number of hidden neurons is sufficiently large and the activation functions of the hidden neurons satisfy suitable properties, for example, to be sigmoidal or, more in general, not-polynomial functions [Sonoda and Murata, 2017, Bishop, 2006, Pinkus, 1999]. Thus, our neural network architecture with variable activation functions is again a MLFF neural network. Importantly, this property means that any classical approach applicable to MLFF networks can also be directly applied to our architecture with trainable activation functions. Notably, as we will discuss in Section 2 and 3 our architecture represents a general framework in which several approaches recently proposed in literature can be included.

The paper is structured as follows. In the next section we critically discuss our approach with respect to the current literature. In Section 3 we introduce our architecture. Section 4 is dedicated to the experimental assessment. Finally, Section 5 is left to the conclusions.

## 2 Related work

Over the last years, ReLU has become the standard activation function for deep neural models, surpassing classic functions as sigmoid and tanh used in the past literature thanks to useful properties, such as the ability to avoid saturation issues [Nair and Hinton, 2010, Glorot et al., 2011]. Different variations of the ReLU have been proposed [Maas et al., 2013, Memisevic et al., 2014, Dugas et al., 2000]. All these functions are somehow different from ReLU, but once chosen they remain fixed, with the choice of which one to use taken during the design stage, typically in a heuristic way. A partial attempt to overcome this drawback moves in the direction of searching the best activation function from a predefined set [Liu and Yao, 1996, Yao, 1999, Ramachandran et al., 2018]. These techniques are limited by the fact that the functions are not learned, but just selected from a collection of standard functions. Thus, approaches by trainable activation functions propose more general frameworks. In this direction one can isolate three basic types of approaches: *parameterized standard activation function*, *linear combination of one-variable functions* and *ensemble of standard activation functions*. In Subsection 2.1, 2.2 and 2.3 we will discuss these three types of approaches. In Subsection 2.4 we will present other types of solutions. Our discussion will be mainly based on three dimensions: 1) how many new parameters are added to the network model, 2) the possibility or not to use standard techniques, within neural network context, for learning the new parameters, such as backpropagation for computing the error function gradient or sparse methods, 3) the expressive power of the trainable activation functions.

### 2.1 Parameterized standard activation functions

With the expression parameterized standard activation functions we refer to all the functions with a shape that is very similar to a given standard activation function, but whose diversity from the latter comes from a set of trainable parameters. The addition of these parameters therefore requires changes, even minimal ones, in the learning algorithm, for example, in the case of using gradient-based methods with backpropagation, the partial derivatives of the error function respect to these new parameters are needed. A first attempt to have a parameterized activation function is given in [Hu, 1992] where the proposed activation function uses two trainable parameters  $\alpha, \beta$  to rule the function shape of a classic sigmoidal function. Similar works on sigmoidal and hyperbolic tangent functions are discussed in [Yamada and Yabuta, 1992a,b, Chen and Chang, 1996, Singh and Chandra, 2003, Chandra and Singh, 2004]. More recently, the authors in [He et al., 2015] introduce PReLU, a parametric version of ReLU, which modifies the function shape when the argument is negative. However, the resulting function remains basically a modified version of the ReLU function that can change its shape in a restricted domain. In [Clevert et al., 2015] ELU function is proposed, which outperforms the results obtained by ReLU on CIFAR100 dataset, becoming one of the best activation function currently known. However, it needs an external parameter  $\alpha$  to be set. In [Trottier et al., 2017] PELU unit is proposed, where the need to manually set the  $\alpha$  parameter is eliminated using two additional trainable parameters.

In all the approaches previously described, although the number of added parameters for each node is low, the expressive power of the trainable activation functions is limited.

## 2.2 Linear combination of one-variable functions

In this case, activation functions are modelled in terms of linear combinations of one-variable functions. These one-variable functions can in turn have additional parameters. For example, in [Agostinelli et al., 2014] each activation function is represented as a linear combination of  $S$  hinge-shaped functions. Each hinge-shaped function has just one parameter which regulates the location of the hinge. The number of additional parameters that must be learned when using this approach is  $2SM$ , where  $M$  is the total number of hidden units in the neural network. During the learning phase, the network can be trained using standard methods based on backpropagation. Any continuous piecewise-linear function can be approximated arbitrarily well provided the number  $S$  of hinge-shaped functions is sufficiently large.

A similar approach has been recently proposed by [Scardapane et al., 2018]. In this case, the activation function is modelled as a linear combination of  $S$  fixed functions, where the  $S$  fixed functions are defined in terms of parametric kernel functions. The parameters of the kernel functions are computed before the network learning phase by some heuristic procedure applied on the data set. During the network learning phase the number of additional parameters is  $SM$ , however for the kernel functions a number of  $KS$  parameters must be computed in a prior stage (where  $k$  is the number of parameters of the kernel functions). In case of a correct choice of the parameters of the kernel functions, any continuous one-to-one function defined on a compact set can be approximated arbitrarily well, provided the number of kernel functions is sufficiently large.

In [Ertuğrul, 2018], in the context of random weight artificial neural networks, a trainable activation function is proposed in terms of a polynomial function of degree  $p$ . The coefficients of the polynomial function are computed by linear regression. The number of added parameters corresponds to the number  $p + 1$  of coefficients for each neuron.

## 2.3 Ensemble of standard activation functions

In this type of approaches, activation functions are defined as an ensemble of a predetermined number of standard activation functions. For example, the authors of [Jin et al., 2016] designed an activation  $S$ -shape function composed by three linear functions taking inspiration by Webner-Fechner [Fechner, 1966] and Stevens law [Stevens, 1957], or in [Qian et al., 2018] a mixture of eLU and ReLU is presented. Interestingly, in [Sütfeld et al., 2018] the authors propose a trainable activation function in terms of a linear combination of  $n$  different, predefined and fixed functions such as hyperbolic tangent (tanh), ReLU and ELU. The added parameters are the  $n$  coefficients of the linear combination for each hidden neuron. A similar approach is proposed in [Harmon and Klabjan, 2017] where the authors model the trainable activation function as a linear combination of a predefined set of  $n$  normalized fixed activation functions. The added parameters are the coefficients of the linear combination and a set of offset parameters,  $\eta$  and  $\delta$ , which are used to dynamically offset the normalization range for each predefined function. Moreover, in order to force the network to choose amongst the predefined activation functions, during the learning process it is required that all the coefficients of the linear combination add to one. This then gives rise to another optimization process unrelated to the classic learning procedure for neural networks

## 2.4 Other approaches

Two interesting and successful approaches are Maxout[Goodfellow et al., 2013, Sun et al., 2018] and NIN[Lin et al., 2013]. However, despite the good performances, both approaches move away from the concept of trainable activation function as it has been previously discussed insofar as the adaptable function does not correspond to the neuron activation function by which the neuron output is computed on the basis of a scalar value (the neuron input) according to the standard two-stage process. In fact, in Maxout, instead of computing the input  $a_i$  of a neuron  $i$  and then assigning it as input to a trainable activation function,  $n$  input  $a_{ij}$  are computed, with  $j = 1, \dots, n$ , by  $n$  trainable linear functions, and then the maximum is taken over the output of these linear functions. NIN instead represents an approach used specifically in the case of convolutional

neural networks, wherein the nonlinear parts of the filters are replaced with a fully connected neural network acting on all channels simultaneously.

Another interesting way to tackle the problem is to use interpolating functions as in [Scardapane et al., 2017, Trentin, 2001]. For example, in [Scardapane et al., 2017] the authors propose an adaptable activation function by using a cubic spline interpolation, whose  $q$  control points for each neuron are adapted in the optimization phase. External methods to classic approaches in neural networks are needed to train the added parameters  $q * m$ , where  $m$  is the number of hidden neurons.

## 2.5 Summarizing

In all the known approaches, to the best of our knowledge, either the expressive power of the trainable activation functions is limited or they add new learning mechanisms, constraints and categories of parameters, by contrast in our approach we achieve a feed-forward neural network with *trainable* activation functions by a feed-forward neural network with *fixed* activation functions, thus leaving unaltered the classic learning mechanisms and categories of parameters. Thus, in our approach a number of attractive properties are simultaneously satisfied: *p1*) the trainable activation function can approximate arbitrarily well any continuous one-to-one mapping defined on a compact input domain, *p2*) any standard learning mechanism for neural network can be directly and easily applied, *p3*) no learning process in addition to those classically used for neural networks is added, *p4*) the added parameters are network weights or biases, therefore any classical regularization method can be used, including the possibility of imposing sparsity by using norms such as  $l_1$ .

None of the known approaches possess all these properties simultaneously. For example, property *p1* is non satisfied for all the approaches discussed in Section 2.1 and 2.3, the approaches discussed in Section 2.2 either do not satisfy property *p1* as in [Agostinelli et al., 2014] or property *p3* as in [Scardapane et al., 2018].

Interestingly, as we will discuss in Section 3 our architecture represents a general framework in which all the approaches described in Section 2.2 and some of the approaches in Section 2.3 can be included, insofar as any linear combination of  $m$  one-variable functions can be represented by a sub-network with  $m$  hidden neurons.

## 3 System architecture

### 3.1 Proposed model: Variable Activation Function Subnetwork

In general, as already introduced in Section 1, in a MLFF network the output of a neuron  $i$  belonging to the  $l$ -th layer is obtained by a two-step computation (see [Bishop, 2006], Chapter 4). The first step computes the input  $a_i^l = \sum_j w_{ij}^l z_j^{l-1} + b_i^l$ , where  $j$  runs over neuron's indexes (or network input values) of the previous

layer  $l - 1$ , which send connections to  $i$ ,  $z_j^{l-1}$  are the output of the neurons belonging to the previous layer (or network input values),  $w_{ij}^l$  are the connection weights going from the neurons  $j$  of the previous layer  $l - 1$  to the neuron  $i$ , and  $b_i^l$  is the bias associated to the neuron  $i$ . The output of the neuron  $i$  is then computed in a second step transforming the input  $a_i^l$  using a fixed activation function  $f$ , obtaining  $z_i^l = f(a_i^l)$ .

The key idea of our approach is to implement the second step of the computation by a "small" one-hidden layer sub-network, with  $k$  hidden neurons and just one input and one output neuron. Let us call it *Variable Activation Function* (henceforth, VAF) sub-network. So, a VAF for a neuron  $i$  can be described as a network composed by:

- an hidden layer, composed by  $k > 1$  neurons directly connected to the neuron  $i$  by a set of weights  $\alpha_h$ , with  $h = 1, 2, \dots, k$ ;
- a fixed activation function  $g$  for the  $k$  hidden neurons;
- an output layer composed by a single neuron connected to all the neurons of the hidden layer by a set of weights  $\beta_h$ , with  $h = 1, 2, \dots, k$ .

The computation of a VAF sub-network associated to a neuron  $i$  can be described as follows: VAF sub-network is fed with the input  $a_i$  of the neuron  $i$ , then the  $k$  neurons of the hidden layer compute  $k$  outputs as

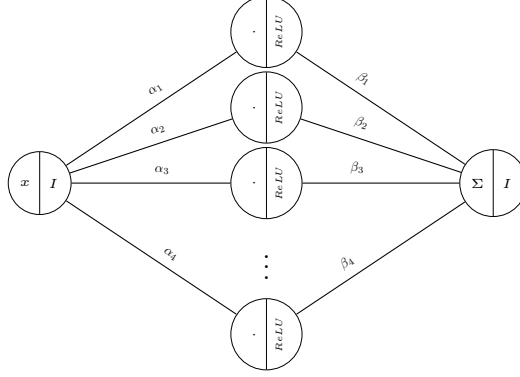


Figure 1: A general VAF scheme; with  $I$  we indicate the identity function

$y_h = g(\alpha_h a_i + \alpha_{0h})$  with  $h = 1, 2, \dots, k$ , while the output neuron computes  $z = \sum_h \beta_h y_h + \beta_0$ .  $\alpha_h$  and  $\alpha_{0h}$  are weights and biases of the hidden layer of the VAF sub-network, respectively, and  $\beta_h$  and  $\beta_0$  are weights and bias of the output layer of VAF sub-network, respectively. In this way the output  $z_i$  of the neuron  $i$  can be expressed as:

$$z_i = VAF(a_i) = \sum_{j=1}^k \beta_j g(\alpha_j a_i + \alpha_{0j}) + \beta_0 \quad (1)$$

$\alpha_j$ ,  $\alpha_{0j}$ ,  $\beta_j$  and  $\beta_0$  are the parameters to be learned from data during the training process.

A general schema of a VAF unit is shown in figure 1. This schema enables one to approximate arbitrarily well any activation function provided that:

- the number  $k$  of hidden neurons in the VAF is sufficiently large;
- the activation function  $g$  of the hidden layer is a not-polynomial function.

As already discussed in Section 1, the first condition is given in [Hornik et al., 1989, Hornik, 1991], where it was shown that a shallow networks can approximate any continuous function provided that a sufficient number of hidden neurons are available and that the activation function is continuous, bounded and non-constant. This result was generalized in [Leshno et al., 1993], where it is proved that a shallow network can approximate any continuous function to any degree of accuracy if and only if the network's activation function is not polynomial. Therefore a VAF activation function can substitute any other network activation function without loss in generality, and having as overhead only an increase in the number of networks parameters, that is equal to  $N \cdot (3k + 1)$  with  $N$  total number of the hidden neurons of the network. Anyway, the number of required parameters can drop to  $L \cdot (3k + 1)$ , with  $L$  number of hidden layers, if we adopt the *shared weights principle*, so that the functions on the same layer share the same VAF weights. With this design choice, we reduce the number of parameters by making the reasonable assumption that if one function is good for a single neuron, then it should also be good for the other neurons of the same layer. This assumption can also be motivated, instead of under the profile of the sub-networks weights, in terms of activation function of a classic neural networks used in the neural network literature, where neurons on to same layer exhibit the same activation function. Summing up, under the shared weights principle for every network layer the only added hyper-parameters to set are:

- the number  $k$  of hidden neurons of the VAF subnetwork;
- the activation function  $g(\cdot)$  of the VAF hidden neurons.

It is worth to emphasize the fact that, in our approach, we have a neural network architecture which is still a MLFF network with fixed activation functions, without adding any external structure or parameters. Let us clarify this aspect (see also Figure 2). Given a neuron  $i$  belonging to  $l$ -th layer of a MLFF network  $Net$ , its output is computed as  $z_i^l = f(a_i^l)$ , in our approach we substitute the activation function  $f$  with the

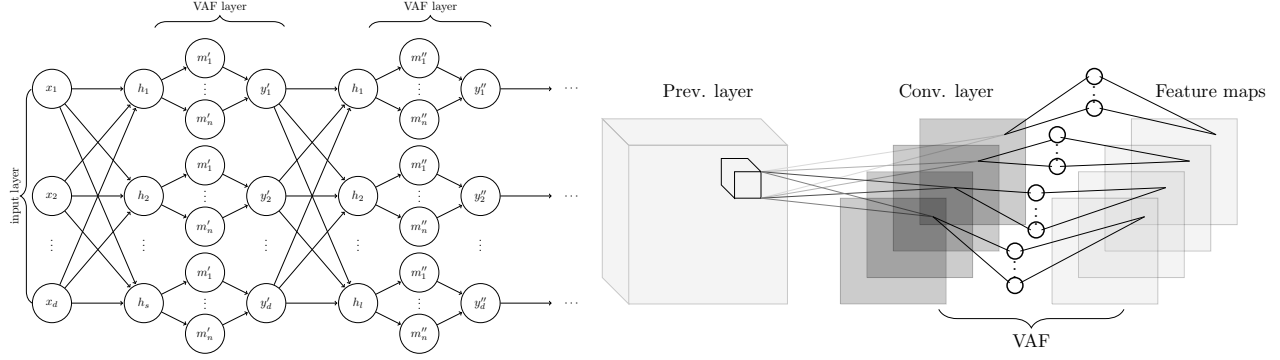


Figure 2: an example of VAF in a 2 full connected network (on the left) and in a convolutional layer (on the right)

Identity function, thus obtaining  $z_i^l = a_i^l$ . Then, we add a VAF sub-network which receives as input variable the output  $z_i^l$  of the neuron  $i$  and computes its output as defined in eq. 1. Finally, this output is sent as input of the next layer  $l + 1$  of  $Net$ . This procedure is uniformly performed for all the neurons of the MLFF network  $Net$ , but the output layer. Thus, one obtains a new neural network  $VafNet$  which is still a MLFF network with fixed activation functions, however it behaviours as  $Net$  equipped with trainable activation functions expressed in terms of eq. 1. Consequently, any standard training procedure can be left unaltered (e.g., Stochastic Gradient Descent).

Figure 2 shows how a VAF network can be integrated into a common multilayer full-connected neural network (on the left) and in a convolutional neural network (on the right).

Notably, given that a VAF subnetwork performs a linear combination of one-variable functions, any approach discussed in Section 2.2 can be included in this schema, provided to choose suitably the activation function  $g$  and the parameters  $\alpha$  and  $\beta$ .

### 3.2 VAF network learning

As discussed above, our neural architecture including VAF is a MLFF network, consequently it can be trained using any learning algorithm dedicated to MLFF network. However, in case of the same VAF acting uniformly for all neurons of a layer, then there is the constrain that the weights of VAF networks should be considered *shared weights*. From an implementation point of view this corresponds to consider a VAF network as a function convolving with the  $a_i^l$  values [Lin et al., 2013]. The weight values of the VAF, being few and connected to each unit, influence the behaviour of the entire network, therefore their behaviour must be taken into consideration during the training phase, and, in particular, the initial value of the VAF weights can be decisive. The training of a neural network usually starts initializing the weights and biases in a random way [Bishop, 2006], or using any initialization rule as for example [Glorot and Bengio, 2010]. Although these approaches can also be followed in our case, it is possible to choose different solutions for the VAF weights initialization. In particular, a possible alternative is to select the initial weights of the VAF so that at the beginning of the learning process the VAF networks approximate a fixed function. For example, we can select a classic activation function as ReLU or sigmoid, or the  $f$  activation function associated to the other hidden layers of the network. In this way hypothetically the function would start from a notoriously already valid form in which the training process should only modify it just enough to improve the performance of the network based on the training data. However, it should be kept in mind that this choice can affect negatively the solution generated by the learning process, given that the resulting VAF can be too similar

to the initial function.

---

**Algorithm 1:** Standard learning schema

---

**Input:**  $TS, VS, net, MaxEpochs$ :  $TS$  and  $VS$  are training and validation datasets, respectively;  $net$  is the network to be trained;  $MaxEpochs$  is the maximum number of epochs

**Output:**  $trainNet$ : Trained net

```

1  $net \leftarrow weightAndBiasInitialization(net)$  ;
2  $bestNet \leftarrow net$  ;
3  $n \leftarrow 0, minErr \leftarrow MAX$  ;
4 repeat
5    $n \leftarrow n + 1$  ;
6    $net \leftarrow learningAlgorithm(net, TS)$  ;
7    $errorT(n) \leftarrow Sim(net, TS)$ ;
8    $errorV(n) \leftarrow Sim(net, VS)$ ;
9   if  $errorV(n) < minErr$  then
10     $minErr \leftarrow errorV(n)$ ;
11     $bestNet \leftarrow net$ ;
12  end
13 until  $n > MaxEpochs$  OR  $earlyStoppingCriteria(errorT, errorV)$ ;
14  $trainNet \leftarrow bestNet$ ;

```

---

## 4 Experimental results

In this section we provide an experimental evaluation of the proposed trainable activation function architecture. In order to achieve a first clue on the validity of our approach, and some heuristic indications for the initialization strategies of VAF networks, in Section 4.1 we report some preliminary experiments on *Sensorless*, a relatively small classification dataset used as standard benchmark for supervised techniques.

On the basis of the results of these experiments, we performed two different series of experiments to test our approach on MLFF networks. In the former, we consider standard MLFF networks (Section 4.2), and in the latter convolutional MLFF networks (Section 4.3). In Section 4.2 we consider both classification and regression problems using 20 different datasets. In Section 4.3 we consider more large-scale dataset as MNIST, Fashion MNIST and CIFAR10.

### 4.1 VAF subnetworks: Activation functions, number of hidden neurons and weight initialization

For a preliminary analysis of the validity of our approach, and for defining some heuristic choices about VAF subnets such as the number and the activation functions of the hidden neurons, we perform a series of experiments on *Sensorless* dataset (see table 1 for details), partitioning it in a random sample of 60% for training, 20% for validation and another 20% for testing. According to what was also reported in [Scardapane et al., 2018], if one uses a standard shallow network, i.e., 1-hidden layer network, we found that  $\tanh$  is the best fixed activation function for this dataset. In particular, using a shallow network with 50 hidden neurons we obtained an accuracy on the test set very close to 100%. Thus, to better investigate the impact of our approach we chose a more “difficult situation” for a shallow network using network models with a small number of hidden neurons. More in detail, we selected three small shallow nets with 5, 10 and 20 hidden neurons.

For each model, We perform a set of experiments using different activation functions.

Firstly, we train these small networks using as fixed activation functions either  $\tanh$  or  $ReLU$ , then we repeat the same experiments substituting the fixed activation functions with VAF subnets as described in Section 3. We considered several scenarios: 1) different number  $k$  of VAF hidden neurons, in particular  $k \in \{3, 5, 7, 9, 11, 15\}$ ; 2)  $\tanh$  and  $ReLU$  as activation functions for VAF hidden neurons; 3) two different strategies for weight initialization of VAF subnets, both a classic random initialization and a weight initialization by which VAF subnets have a behaviour very similar to activation functions of the VAF hidden neurons, we will call the latter *specific weight initialization*; 4) as discussed in Section 3, we examine both

the case in which VAF subnets on the same layer share the weights (shared weights principle) and the case in which VAF subnets on the same layer can have different weights.

We trained all the networks using ADAM algorithm [Kingma and Ba, 2014] for 500 epochs. Furthermore, we repeat our experiments for 10 times.

## Results

In Figure 3a, 3b and 3c are reported the results with respect to the shallow networks with 5, 10 and 20 hidden neurons, respectively, in the case in which VAF subnets on the same layer do not share the weights. In Figure 4a, 4b and 4c are reported the results in the case in which VAF subnets on the same layer share the weights.

Notably, one can observe that all the models equipped with VAF subnets outperform the corresponding shallow networks. Interestingly, these results support the possibility of using a shared VAF approach with a fairly low number of VAF hidden neurons, thus having a lower number of parameters to be learned. In fact, the two approaches, non-shared (Figure 3) and shared (Figure 4) VAF subnets, exhibit a very similar behaviour, and although in all cases accuracy tends to increase as the number of neurons of the VAF subnets increases, this increase is not always very relevant. The two types of VAF weight initialization seem to give similar results, with slightly better performances for random initialization. The use of tanh or ReLU as activation function of VAF hidden neurons, on the other hand, seems to significantly change the network performance. In fact, the accuracy obtained by networks with ReLU activation function for the VAF hidden neurons is uniformly lower than those obtained with the tanh activation function. We suppose that this result is due to the fact that we are always using shallow nets.

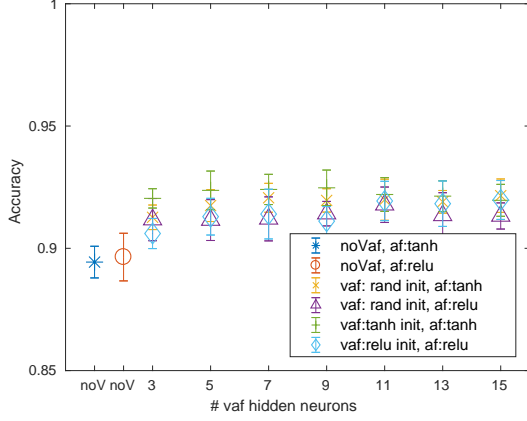
In Figure 5a and 5b are reported the output values of trained VAF subnetworks when a random or a specific weight initialization is chosen, respectively. One can note that the resulting activation functions are often strongly different from the classic tanh and ReLU, and that they exhibit similarly a high degree of non-linearity.

## 4.2 Full-connected MLFF networks: classification and regression

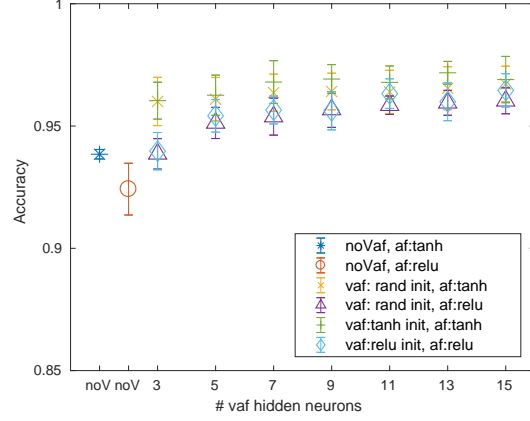
In this experimental scenario we focus on evaluating the impact of both VAF subnetworks and VAF weight initialization using full-connected MLFF networks with 1 or 2 hidden layers trained on 20 public datasets (see Table 1). 10 of these datasets are suitable for classification problems, and 10 for regression problems. The number of hidden neurons varies in the set  $\{10, 25, 50, 100\}$ , but for neural networks with 2 hidden layers we only selected neural networks with a number of hidden neurons belonging to the first layer larger than the number of hidden neurons of the second layer. ReLU was selected as activation function  $g$  of the hidden neurons of VAF sub-networks. Thus, for each dataset we obtained 4 network models with 1-hidden layer, and 6 with 2-hidden layers. Let us call  $net_{m_1}$  and  $net_{m_1, m_2}$  the 1-hidden and 2-hidden layer networks, respectively, with  $m_1, m_2 \in \{10, 25, 50, 100\}$ . On the basis of what was discussed in Section 3, to each network  $net_{m_1}$  ( $net_{m_1, m_2}$ ) it is possible to associate a neural network  $vnet_{m_1}^k$  ( $vnet_{m_1, m_2}^k$ ) equipped with VAF subnetworks, where  $k$  is the number of hidden neurons of VAF subnetworks.

On the basis of the results discussed in Section 4.1, we considered VAF subnets shared on each layer, and  $k = 3$ . In Table 2 we report the neural network architectures used in this series of experiments. Neural network architectures were sorted in ascending order according to their complexity. Networks were trained according to an usual learning approach, described in Algorithm 1. In particular, we used a batch approach, RProp [Riedmiller and Braun, 1992], with “small” datasets, i.e., when the number of examples was less than  $5 \cdot 10^3$ , otherwise we used a mini-batch approach, RMSProp [Tieleman and Hinton, 2012]. Moreover, networks with VAF subnetworks were trained using both a random initialization and an specific weight initialization such that they approximate a ReLU function. All the network models, i.e.,  $net_{m_1}, net_{m_1, m_2}, vnet_{m_1}^k$  and  $vnet_{m_1, m_2}^k$  were compared in a  $K$ -fold cross validation schema (see Algorithm 2), with  $K = 10$ .

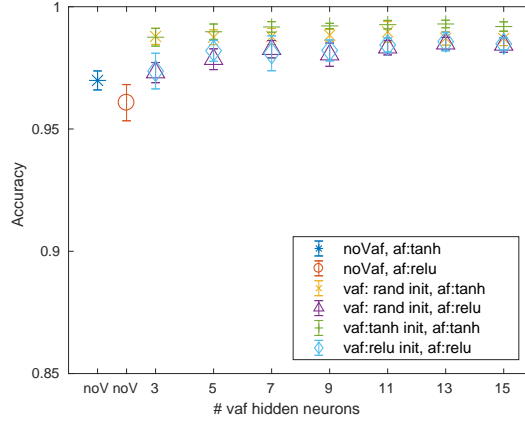
Note that Learning Rate (LR) in RMSProp spans in the range  $[0.0001, 0.1]$ , considering 10 equispaced values, while in RProp  $\eta^+$  was selected equal to 1.01, and  $\eta^-$  equal to 0.5. In Table 3 are summarized the



(a) Shallow Network with 5 hidden neurons

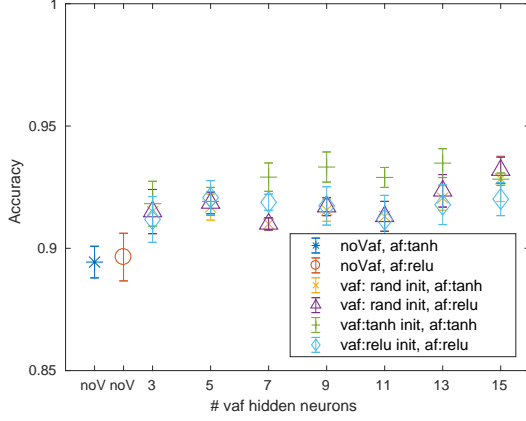


(b) Shallow Network with 10 hidden neurons

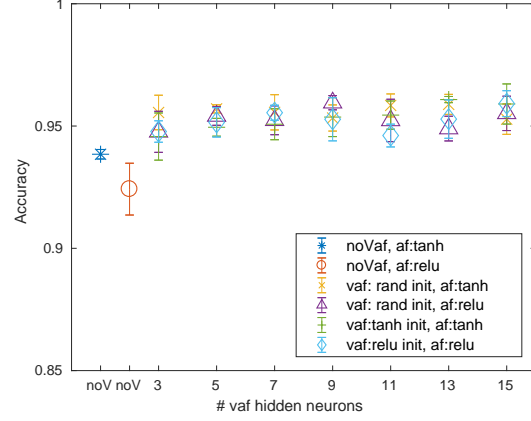


(c) Shallow Network with 20 hidden neurons

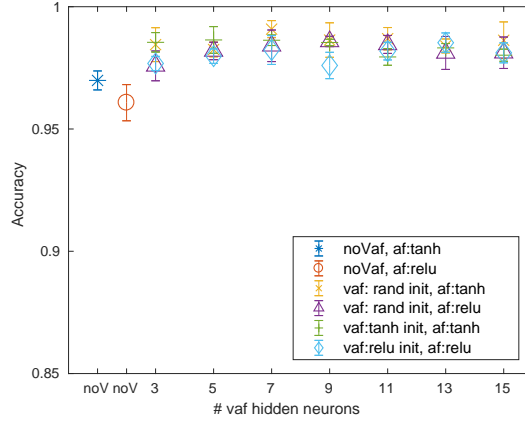
Figure 3: Test accuracy of networks with different VAF subnets on each layer. Using *Sensorless* dataset, we trained three small shallow networks composed of 5, 10 and 20 hidden neurons with fixed activation functions corresponding to either *tanh* or *ReLU*. In figure such networks are referred as *noVaf*. Then we repeated the same experiments substituting the fixed activation functions with VAF subnets. The number VAF hidden neurons ranges in  $k \in 3, 5, 7, 9, 11, 15$ , the possible activation functions for VAF hidden neurons are *tanh* and *ReLU*. Weight initialization of VAF subnets is either a classic random initialization or a weight initialization by which VAF subnets have a behaviour very similar to activation functions of the VAF hidden neurons. VAF subnets on the same layer can have different weights.



(a) Shallow Network with 5 hidden neurons



(b) Shallow Network with 10 hidden neurons



(c) Shallow Network with 20 hidden neurons

Figure 4: Test accuracy of networks with shared VAF subnets on each layer. Using *Sensorless* dataset, we trained three small shallow networks composed of 5, 10 and 20 hidden neurons with fixed activation functions corresponding to either *tanh* or *ReLU*. In figure such networks are referred as *noVaf*. Then we repeated the same experiments substituting the fixed activation functions with VAF subnets. The number VAF hidden neurons ranges in  $k \in \{3, 5, 7, 9, 11, 15\}$ , the possible activation functions for VAF hidden neurons are *tanh* and *ReLU*. Weight initialization of VAF subnets is either a classic random initialization or a weight initialization by which VAF subnets have a behaviour very similar to activation functions of the VAF hidden neurons. VAF subnets on the same layer share the weights.

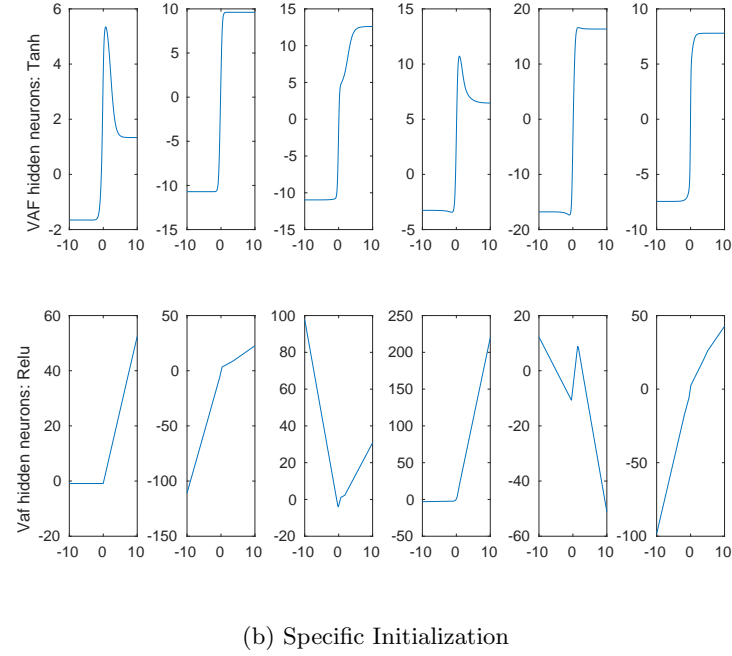
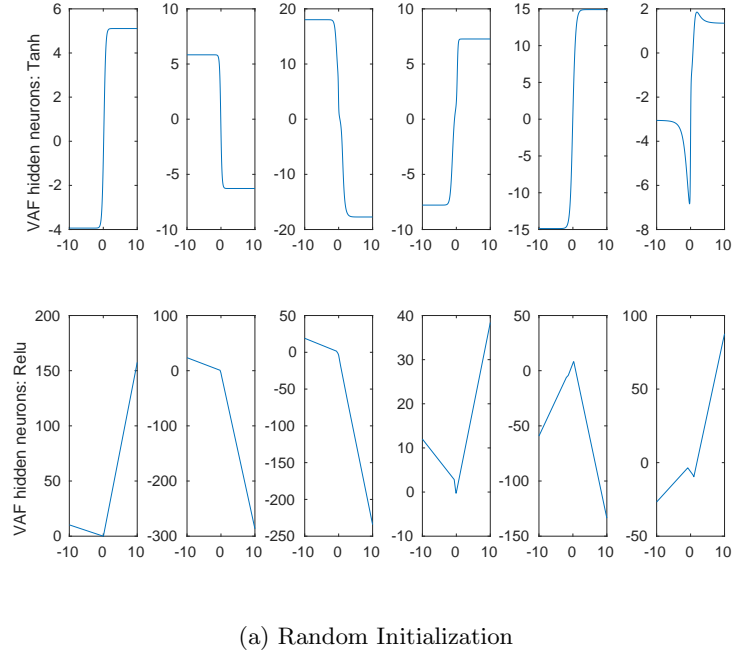


Figure 5: Examples of trained VAF subnetworks. On the y-axis we plot the output value of the VAF. In 5a are plotted trained VAF subnetworks when a random weight initialization is chosen. In 5b when a specific weight initialization is chosen

Name	Istances	Input Dim.	N. classes	Task	Neural Network Arch.	Ref.
CPU-Small	8192	12	-	Regress.	MLFF	Dheeru and Karra Taniskidou [2017]
DeltaElevator	9517	6	-	Regress.	MLFF	<a href="https://www.dcc.fc.up.pt">https://www.dcc.fc.up.pt</a> [2009]
Elevators	16599	18	-	Regress.	MLFF	<a href="https://www.dcc.fc.up.pt">https://www.dcc.fc.up.pt</a> [2009]
Kinematics	8192	8	-	Regress.	MLFF	<a href="https://www.dcc.fc.up.pt">https://www.dcc.fc.up.pt</a> [2009]
Puma-8NH	8192	8	-	Regress.	MLFF	<a href="https://www.dcc.fc.up.pt">https://www.dcc.fc.up.pt</a> [2009]
Puma-32NH	8192	32	-	Regress.	MLFF	<a href="https://www.dcc.fc.up.pt">https://www.dcc.fc.up.pt</a> [2009]
Servo	197	4	-	Regress.	MLFF	Dheeru and Karra Taniskidou [2017]
Energy Cooling	768	8	-	Regress.	MLFF	<a href="https://www.dcc.fc.up.pt">https://www.dcc.fc.up.pt</a> [2009]
Energy Heating	768	8	-	Regress.	MLFF	<a href="https://www.dcc.fc.up.pt">https://www.dcc.fc.up.pt</a> [2009]
Yatch	308	7	-	Regress.	MLFF	Dheeru and Karra Taniskidou [2017]
Sensorless	58509	49	11	Classif.	MLFF	Dheeru and Karra Taniskidou [2017]
Liver	345	7	2	Classif.	MLFF	Dheeru and Karra Taniskidou [2017]
Wine	178	13	3	Classif.	MLFF	Dheeru and Karra Taniskidou [2017]
Statlog Image Segmentation	2310	19	7	Classif.	MLFF	Dheeru and Karra Taniskidou [2017]
Statlog Landsat Satellite	6435	36	7	Classif.	MLFF	Dheeru and Karra Taniskidou [2017]
Cardiotocography	2126	22	3	Classif.	MLFF	Dheeru and Karra Taniskidou [2017]
Seismic bumps	2584	18	2	Classif.	MLFF	Sikora and Wróbel [2010]
Dermatology	336	35	3	Classif.	MLFF	Dheeru and Karra Taniskidou [2017]
Diabetic retinopathy debrecen	1151	19	2	Classif.	MLFF	Antal and Hajdu [2014]
QSAR biodegradation	1055	41	2	Classif.	MLFF	Mansouri et al. [2013]
Climate model simulation	540	18	2	Classif.	MLFF	Lucas et al. [2013]
MNIST	70000	$28 \times 28$	10	Classif.	CNN	LeCun and Cortes [2010]
Fashion MNIST	70000	$28 \times 28$	10	Classif.	CNN	Xiao et al. [2017]
Cifar10	60000	$32 \times 32 \times 3$	10	Classif.	CNN	Krizhevsky and Hinton [2009]

Table 1: Properties of the datasets used for the experiments, and architectures of the neural network applied to the data.

parameters of this series of empirical evaluations.

---

**Algorithm 2:**  $K$ -fold cross validation procedure

---

**Input:** Dataset  $D$ , network model  $mnet$ , number of folds  $k$ , hyper-parameters values  $\{p_1, p_2, \dots, p_n\}$  with  $p_i = \{ \text{possible values for } i\text{-th parameter} \}$  with  $1 \leq i \leq n$

- 1  $FoldResults = []$ ;
- 2 split  $D$  in a  $k$ -partition  $P^k(D)$  ;
- 3 **forall**  $1 \leq i \leq k$  **do**
- 4      $TestSet \leftarrow P_i^k(D)$ ;
- 5      $R \leftarrow P^k(D) \setminus \{TestSet\}$ ;
- 6     split  $R$  in a 2-partition  $P^2(R)$  ;
- 7      $TrainSet \leftarrow P_1^2(R)$ ;
- 8      $ValSet \leftarrow P_2^2(R)$ ;
- 9      $bestParams \leftarrow \emptyset$ ;
- 10     $bestResults \leftarrow \emptyset$ ;
- 11    **forall**  $h \in p_1 \times p_2 \times \dots \times p_n$  **do**
- 12        $model \leftarrow Train(mnet, TrainSet, ValSet, h)$ ;
- 13        $results \leftarrow Sim(model, TestSet)$ ;
- 14       **if**  $results$  **better than**  $bestResults$  **then**
- 15            $bestResults \leftarrow results$ ;
- 16            $bestParams \leftarrow h$ ;
- 17            $bestModel \leftarrow model$ ;
- 18       **end**
- 19    **end**
- 20     $FoldResults[i] \leftarrow bestResults$ ;
- 21 **end**
- 22 **return**  $Average(FoldResults)$

---

## Results

In Table 4 and 5 are showed mean and standard deviations of RMSE and accuracy for regression and classification datasets, respectively, by using a K-fold cross-validation approach. The best results are displayed

Table 2: Neural network architectures used in the first experimental scenario. See text for further details.

	#1	#2	#3	#4	#5	#6	#7	#8	#9	#10
Stand	$net_{10}$	$net_{25}$	$net_{50}$	$net_{100}$	$net_{25,10}$	$net_{50,10}$	$net_{100,10}$	$net_{50,25}$	$net_{100,25}$	$net_{100,50}$
VAF	$vnet_{10}^3$	$vnet_{25}^3$	$vnet_{50}^3$	$vnet_{100}^3$	$vnet_{25,10}^3$	$vnet_{50,10}^3$	$vnet_{100,10}^3$	$vnet_{50,25}^3$	$vnet_{100,25}^3$	$vnet_{100,50}^3$

Table 3: Parameters of the first experimental scenario. See text for further details.

$m_1, m_2$	k	VAF initialization	Learning approaches	# maximum epochs	K
{10, 25, 50, 100}	{3}	{Random, ReLU}	{RMSPProp, RProp}	300	10

in bold.

In case of the regression datasets, VAF approach uniformly overcomes standard approach. Only in one case we obtain the best result with a standard approach. For four datasets (DeltaElevator, Elevators, Puma-32H and Yatch) RMSE 's mean obtained by VAF networks results much smaller than RMSE obtained by neural network without VAF subnetwork. For example in DeltaElevator dataset RMSE's mean was reduced by two (VAF init random) and one (VAF Init ReLU) order of magnitude. Moreover standard deviations remain comparable or lower than those without VAF subnetworks. This suggests that the training process of network with VAF subnetworks is sufficiently stable.

Similar results were obtained with classification datasets (see Table 5). Neural networks with VAF outperforms neural networks without VAF. Only in two datasets (20%) neural networks without VAF outperforms neural networks with VAF. Also in this case, standard deviations remain comparable or lower than those without VAF.

In Figure 5 and 6 are reported some VAF subnetwork behaviours at the end of the learning process.

	standard Relu RMSE: mean + St.Dev	VAF Init random RMSE: mean + St.Dev	VAF Init ReLU RMSE: mean + St.Dev
CPUsmall	$0.0616 \pm 0.0016$ ( $net_{50,10}$ )	<b><math>0.0593 \pm 0.0017</math></b> ( $vnet_{100}^3$ )	$0.0606 \pm 0.0032$ ( $vnet_{100,10}^3$ )
DeltaElevator	$0.1355 \pm 0.0055$ ( $net_{25}$ )	<b><math>0.0030 \pm 0.0006</math></b> ( $vnet_{100,50}^3$ )	$0.0414 \pm 0.0393$ ( $vnet_{10}^3$ )
Elevators	$1.2746 \pm 0.3741$ ( $net_{50,10}$ )	<b><math>0.0068 \pm 0.0003</math></b> ( $vnet_{50,25}^3$ )	$0.1915 \pm 0.1658$ ( $vnet_{10}^3$ )
Kinematics	$0.1090 \pm 0.0040$ ( $net_{100}$ )	$0.1315 \pm 0.0185$ ( $vnet_{25}^3$ )	<b><math>0.0935 \pm 0.0048</math></b> ( $vnet_{50,25}^3$ )
Puma-8NH	$0.1336 \pm 0.0023$ ( $net_{25,10}$ )	<b><math>0.1316 \pm 0.0018</math></b> ( $vnet_{100,10}^3$ )	$0.1331 \pm 0.0025$ ( $vnet_{25,10}^3$ )
Puma-32H	$0.2372 \pm 0.1473$ ( $net_{25,10}$ )	<b><math>0.0273 \pm 0.0005</math></b> ( $vnet_{100,10}^3$ )	$0.0317 \pm 0.0039$ ( $vnet_{25,10}^3$ )
Servo	$0.0946 \pm 0.0158$ ( $net_{100}$ )	<b><math>0.0896 \pm 0.0276</math></b> ( $vnet_{10}^3$ )	$0.0961 \pm 0.0271$ ( $vnet_{25}^3$ )
Energy Cooling	$0.0417 \pm 0.0014$ ( $net_{100,10}$ )	$0.0461 \pm 0.0024$ ( $vnet_{100,25}^3$ )	<b><math>0.0400 \pm 0.0026</math></b> ( $vnet_{50,25}^3$ )
Energy Heating	<b><math>0.0206 \pm 0.0026</math></b> ( $net_{100,10}$ )	$0.0304 \pm 0.0237$ ( $vnet_{25,10}^3$ )	$0.0213 \pm 0.0027$ ( $vnet_{100,10}^3$ )
Yatch	$0.2442 \pm 0.1146$ ( $net_{25}$ )	$0.3435 \pm 0.2186$ ( $vnet_{100}^3$ )	<b><math>0.1481 \pm 0.0553</math></b> ( $vnet_{25}^3$ )

Table 4: RMSE for the experiments on the regression datasets. We used a K-Fold Cross-validation evaluation. In bold the best results. The best neural architecture for each case is between parentheses.

### 4.3 Convolutional MLFF networks

In order to evaluate experimentally the impact of VAF on Convolutional Neural Networks (CNN), we consider standard CNN networks with 2 and 3 convolutional layers, and run experiments on three different dataset: MNIST, Fashion MNIST and CIFAR10 (see Table 1 for further details). As discussed in Section 3.2 and Section 4.2, a key aspect is the initialization of the VAF networks. Thus, also in this case, we chose to initialize the weights of the VAF subnetworks either randomly or to approximate a ReLU function. To this aim, we build two CNN architectures (similar to the basic network used in [Lin et al., 2013]), the first one composed of 2-layer CNN networks used for MNIST and Fashion-MNIST and the second one composed of 3-layers trained and tested with the more complex CIFAR10 dataset. Let us call  $cnet_{A_2}$  and  $cnet_{A_3}$  respectively the 2-layer CNN and the 3-layer CNN; as stated in Section 3, it is possible to associate to each

	standard Relu	VAF Init random	VAF Init ReLU
	Accuracy: mean + St.Dev	Accuracy: mean + St.Dev	Accuracy: mean + St.Dev
Liver	$0.6203 \pm 0.0474$ ( $net_{25,10}$ )	$0.6290 \pm 0.0378$ ( $vnet_{100,10}^3$ )	<b><math>0.6348 \pm 0.0375</math></b> ( $vnet_{25}^3$ )
Wine	$0.8879 \pm 0.0516$ ( $net_{10}$ )	<b><math>0.9552 \pm 0.0371</math></b> ( $vnet_{50}^3$ )	$0.9162 \pm 0.0434$ ( $vnet_{10}^3$ )
Image segmentation	<b><math>0.9463 \pm 0.0128</math></b> ( $net_{25}$ )	$0.9351 \pm 0.0179$ ( $vnet_{50}^3$ )	$0.9381 \pm 0.0079$ ( $vnet_{50}^3$ )
Satellite image	$0.8821 \pm 0.0101$ ( $net_{100,50}$ )	$0.8856 \pm 0.0028$ ( $vnet_{100}^3$ )	<b><math>0.8875 \pm 0.0080</math></b> ( $vnet_{100,25}^3$ )
CTG	$0.8979 \pm 0.0263$ ( $net_{100}$ )	<b><math>0.9040 \pm 0.0073</math></b> ( $vnet_{100}^3$ )	$0.8984 \pm 0.0261$ ( $vnet_{50,25}^3$ )
Seismic bumps	<b><math>0.9346 \pm 0.0009</math></b> ( $net_{10}$ )	$0.9234 \pm 0.0074$ ( $vnet_{10}^3$ )	$0.9342 \pm 0.0001$ ( $vnet_{10}^3$ )
Dermatology	$0.9749 \pm 0.0116$ ( $net_{50,25}$ )	$0.9692 \pm 0.0182$ ( $vnet_{10}^3$ )	<b><math>0.9750 \pm 0.0248</math></b> ( $vnet_{100}^3$ )
Diabetic	$0.7254 \pm 0.0290$ ( $net_{100}$ )	$0.7315 \pm 0.0238$ ( $vnet_{10}^3$ )	<b><math>0.7333 \pm 0.0231</math></b> ( $vnet_{50}^3$ )
Biodegradation	$0.8635 \pm 0.0336$ ( $net_{10}$ )	<b><math>0.8673 \pm 0.0225</math></b> ( $vnet_{100,10}^3$ )	$0.8569 \pm 0.0108$ ( $vnet_{50}^3$ )
Climate simulation	$0.9500 \pm 0.0140$ ( $net_{50,25}$ )	$0.9519 \pm 0.0211$ ( $vnet_{10}^3$ )	<b><math>0.9556 \pm 0.0240</math></b> ( $vnet_{100}^3$ )

Table 5: Accuracies for the experiments on the classification datasets. We used a K-Fold Cross-validation evaluation. In bold the best results. The best neural architecture for each case is between parentheses.

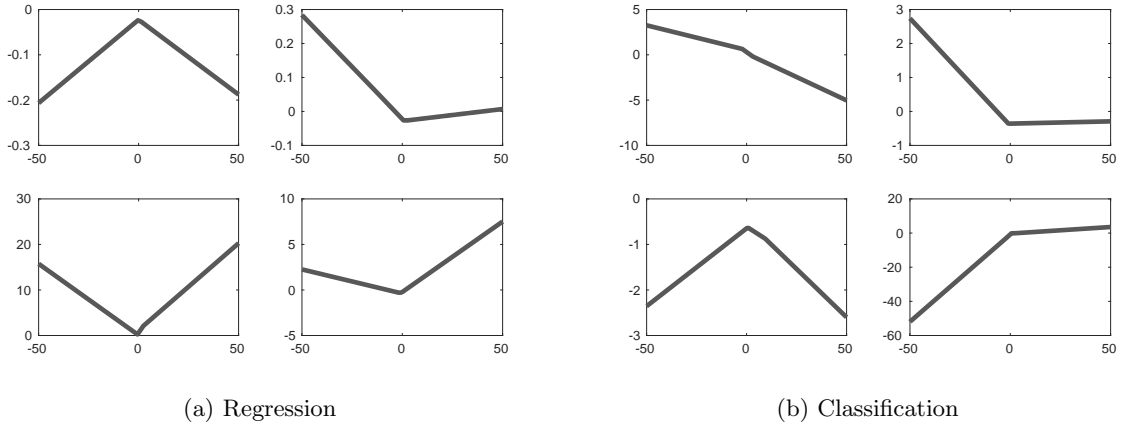


Figure 6: Plots of some VAF behaviours at the end of the learning process. In 6a for regression datasets, in 6b for classification datasets.

$cnet$  a neural network  $vnet^k$  equipped with VAF sub-networks having  $k$  hidden units. The experiments were preformed using a 10-fold cross validation schema as described in 2. Networks were trained using Stochastic Gradient Descent (SGD) method with mini-batching.

Furthermore, we compare our architecture with two other neural architectures also equipped with trainable activation functions. The first one is KAFnet, a very recent and promising approach proposed in [Scardapane et al., 2018] and already discussed in Section 2.2. The second one is Network in Network (NIN), a successful approach proposed in [Lin et al., 2013] and already discussed in Section 2.4. To this aim, we used the same experimental settings described in [Scardapane et al., 2018], i.e., a convolutional MLFF network composed by two convolutional layers, each of these followed by a  $3 \times 3$  maxpooling layer and a dropout layer of 0.25 (see Table 6). To distinguish it from the others models, we will call this network  $cnet_B$ . Starting from  $cnet_B$ , we obtained three different types of neural networks with trainable activation functions according to three different procedures proposed for KAFnet, NIN and VAF. We use the classic CIFAR10 data configuration (50000 training samples + 10000 test samples) to train the three types of obtained networks. The network  $cnet_B$  with fixed activation function corresponding to ReLU is also considered as baseline. Finally, we repeat the same setup using MNIST and Fashion-MNIST dataset.

Properties of the used CNN architectures and learning process are summarised in table 6.

## Results

In Table 7 are shown mean and standard deviations of accuracies for the three datasets Cifar10, MNIST and Fashion MNIST, using a 10-fold cross-validation approach for the neural architecture summarized in the first

Table 6: Parameters of the second experimental scenario. See text for further details.

Name	Layers	VAF initialization	Learning approaches	# maximum epochs
$cnet_{A_2}$	2× (Conv. 192 + Maxout + Dropout)	{Random, ReLU}	SGD	300
$cnet_{A_3}$	3× (Conv. 192 + Maxout + Dropout)	{Random, ReLU}	SGD	300
$cnet_B$	2× (Conv. 150 + Maxout + Dropout)	Random	Adam	300

two rows of Table 6. The best results are reported in bold style. One can note that VAF approach uniformly outperforms the standard approach, especially when using a random initialization scheme. Also in this experimental scenario the standard deviations obtained by networks with VAF remain comparable or lower than those without VAF subnetworks. Especially for the CIFAR10 dataset, we obtain a considerable improvement.

In Figures 7 and 8 are shown some examples of trained activation functions respectively in  $vcnn_{A_2}$  and  $vcnn_{A_3}$ ; it should be noted the influence of VAF initialization on the trained activation function: it seems that, in case of initialization as ReLU, the initial shape remains mostly unchanged, giving a resulting function that looks like a PReLU/Leaky ReLU. A more interesting behaviour is given by random initialization, where every VAF unit seems to exhibit greater changes respect to the initial function. This greater variability given by random initialization respect to ReLU initialization seems to give an improvement in accuracy results as shown in Table 7.

In Table 8 we show the performances of KAFnet, NIN and VAF network on the two datasets Fashion-MNIST and CIFAR10 in terms of accuracy. VAF network outperforms KAF and NIN on both the dataset CIFAR10 and Fashion MNIST. We do not report the MNIST results because are all very similar between them (over the 99% of accuracy). Notably, therefore, also with respect to two other two approaches with trainable activation functions known in literature, our approach results in better performance.

	standard ReLU Acc. + St.Dev	VAF Init random Acc. + St.Dev	VAF Init ReLU Acc. + St.Dev
Cifar10	$0.857 \pm 0.002$ ( $cnet_{A_3}^5$ )	<b><math>0.875 \pm 0.003</math></b> ( $vcnet_{A_3}^5$ )	$0.860 \pm 0.002$ ( $vcnet_{A_3}^5$ )
MNIST	$0.991 \pm 0.001$ ( $cnet_{A_2}^5$ )	<b><math>0.994 \pm 0.001</math></b> ( $vcnet_{A_2}^5$ )	$0.993 \pm 0.002$ ( $vcnet_{A_2}^5$ )
Fashion MNIST	$0.923 \pm 0.001$ ( $cnet_{A_2}^5$ )	<b><math>0.935 \pm 0.002</math></b> ( $vcnet_{A_2}^5$ )	$0.934 \pm 0.001$ ( $vcnet_{A_2}^5$ )

Table 7: Results of the convolutional networks with a 10-fold cross Validation with  $cnet_A$ .

	standard ReLU Accuracy	VAF(M=5) Accuracy	KAF(D=20) Accuracy	NIN Accuracy
Cifar10	0.707	<b>0.812</b>	0.802	0.763
MNIST	0.995	0.995	0.995	<b>0.996</b>
Fashion MNIST	0.920	<b>0.935</b>	0.929	0.925

Table 8: Comparison between different activation functions on different dataset using the standard division on  $cnet_B$ .

## 5 Conclusion

In this work, we proposed a simple and direct way to obtain adaptable activation functions in feed-forward neural networks. In particular, we proposed to modify a feed-forward neural network by adding Variable Activation Functions (VAF) in terms of one-hidden layer subnetworks (see Section 3). The resulting network is still a feed-forward neural network. The proposed architecture doesn't need many more parameters than networks using not adaptable activation functions as ReLU, and the learning process follows standard approaches (see Section 3.2). Importantly, VAF subnetworks can approximate arbitrarily well any activation functions, provided that the number of hidden neurons is sufficiently large (see Section 3).

It is worth to remark that our approach distinguishes from other approaches proposed in literature insofar as it satisfies simultaneously the properties  $p1 - p4$  as described in Section 2.5. These properties include a high expressive power of the trainable activation functions, no external parameter or learning process in addition to the classical ones for neural networks, and the possibility to use classical regularization methods.

Interestingly, as we discussed in Section 3 our architecture represents a general framework in which all the approaches described in Section 2.2 and some of the approaches in Section 2.3 can be included.

We experimentally evaluated our architecture on three different sets of experiments. In the former (see Section 4.1, we tested our approach using small shallow networks for defining some heuristic choices about VAF subnets. Notably, all the models equipped with VAF subnets outperform the corresponding shallow networks, and the results support the possibility of using a shared VAF approach with a fairly low number of VAF hidden neurons. In the second series of experiments (see Section 4.2), we considered full-connected Multi-Layered Neural Network (MLFF) networks. More specifically, we selected 10 networks with 1 or 2 hidden layers. A correspondent network with VAF subnetworks was built for each of these 10 networks (see Section 3 and 4.2). We obtained a total of 20 different neural network architectures. These neural architectures were evaluated and compared using a  $K$ -Fold Cross-Validation procedure (see Algorithm 2) on 20 different datasets (see Table 1), either for classification tasks or regression tasks. The results show that the networks with VAF subnetworks are uniformly more performing than the ones without VAF networks. In particular, our approach outperforms that without VAF networks on the 85% of the datasets. Only on three datasets our approach had worse results.

In the last set of experiments, we considered Convolutional Neural Networks with 2 and 3 layers and correspondent networks with VAF units and we evaluate them using 3 image datasets (MNIST, Fashion MNIST and CIFAR10) for classification. Also in this case the VAF subnetworks outperform networks with static units and selected state-of-the-art neural architectures (KAFNet and NIN) equipped with trainable activation functions.

In conclusion, VAF units have been tested using traditional MLNN networks and CNN networks with various datasets and give better results compared with networks with similar design both with traditional ReLU functions and trainable activation functions. We showed that is possible to obtain encouraging results without the need to use complex designs, particular initialization schemes or learning process in addition to those classically used for neural networks.

## Acknowledgements

The work has been partially supported by the national project Perception, Performativity and Cognitive Sciences - PRIN2015 Cod. 2015TM24JS\_009.

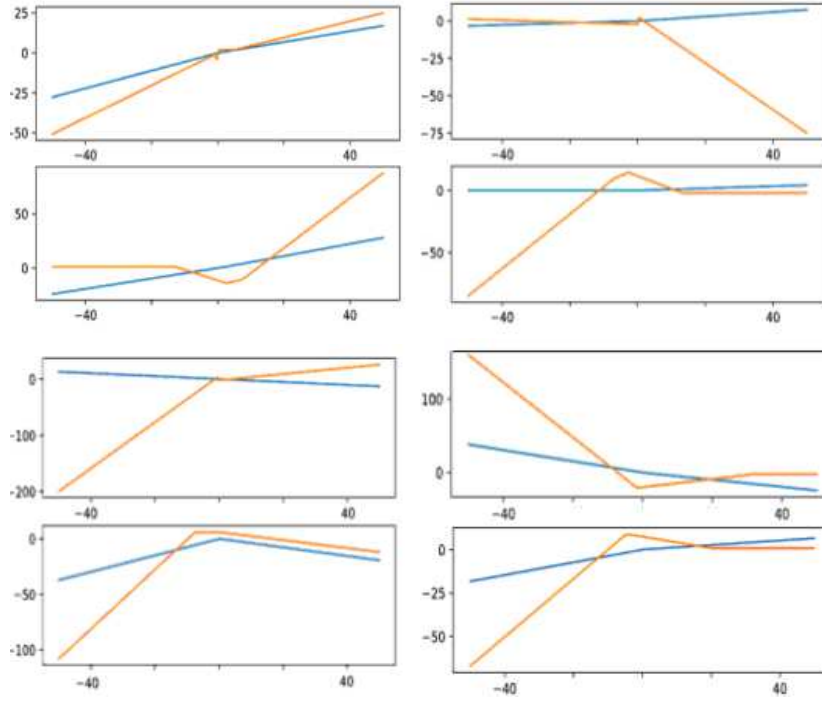
## References

- Ronald A DeVore, Konstantin I Oskolkov, and Pencho P Petrushev. Approximation by feed-forward neural networks. *Annals of Numerical Mathematics*, 4:261–288, 1996.
- C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- Brian D Ripley. *Pattern recognition and neural networks*. Cambridge university press, 2007.
- Sho Sonoda and Noboru Murata. Neural network with unbounded activation functions is universal approximator. *Applied and Computational Harmonic Analysis*, 43(2):233–268, 2017.
- Allan Pinkus. Approximation theory of the mlp model in neural networks. *Acta Numerica*, 8(January): 143–195, 1999.
- Namig J Guliyev and Vugar E Ismailov. A single hidden layer feedforward network with only one neuron in the hidden layer can approximate any univariate function. *Neural computation*, 28(7):1289–1304, 2016.
- Ronen Eldan and Ohad Shamir. The power of depth for feedforward neural networks. In *Conference on Learning Theory*, pages 907–940, 2016.

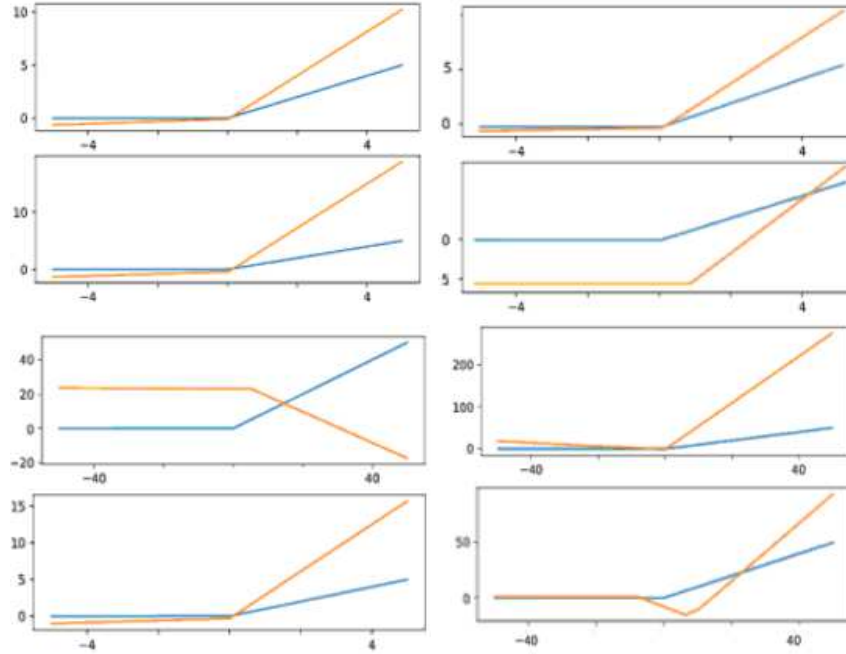
- Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, pages 807–814, 2010.
- Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, 2013.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, volume 9, pages 249–256, 2010.
- Djork-Arn Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *CoRR*, abs/1511.07289, 2015.
- Gnter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. *CoRR*, abs/1706.02515, 2017.
- Dabal Pedamonti. Comparison of non-linear activation functions for deep neural networks on mnist classification task. *CoRR*, abs/1804.02763, 2018.
- Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for activation functions. In *6th International Conference on Learning Representations (ICLR), 2018*, 2018.
- Sheng Qian, Hua Liu, Cheng Liu, Si Wu, and Hau-San Wong. Adaptive activation functions in convolutional neural networks. *Neurocomputing*, 272:204–212, 2018.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323, 2011.
- Roland Memisevic, Kishore Reddy Konda, and David Krueger. Zero-bias autoencoders and the benefits of co-adapting features. *CoRR*, abs/1402.3337, 2014.
- Charles Dugas, Yoshua Bengio, Franois Blisle, Claude Nadeau, and Ren Garcia. Incorporating second-order functional knowledge for better option pricing. In *NIPS*, pages 472–478, 2000.
- Yong Liu and Xin Yao. Evolutionary design of artificial neural networks with different nodes. In *International Conference on Evolutionary Computation*, pages 670–675, 1996.
- X. Yao. Evolving artificial neural networks. *Proceedings of the IEEE*, 87:1423–1447, 1999.
- Shao H. Hu, Z. The study of neural network adaptive control systems. *Control and Decision*, 7(2):361–366, 1992.
- T. Yamada and T. Yabuta. Neural network controller using autotuning method for nonlinear functions. *IEEE Transactions on Neural Networks*, 3(4):595–601, 1992a.
- T. Yamada and T. Yabuta. Remarks on a neural network controller which uses an auto-tuning method for nonlinear functions. In *[Proceedings 1992] IJCNN International Joint Conference on Neural Networks*, volume 2, pages 775–780 vol.2, 1992b.
- Chyi-Tsong Chen and Wei-Der Chang. A feedforward neural network with function shape autotuning. *Neural networks*, 9(4):627–641, 1996.
- Yogesh Singh and Pravin Chandra. A class+ 1 sigmoidal activation functions for ffnns. *Journal of Economic Dynamics and Control*, 28(1):183–187, 2003.
- Pravin Chandra and Yogesh Singh. An activation function adapting training algorithm for sigmoidal feed-forward networks. *Neurocomputing*, 61:429–437, 2004.

- Ludovic Trottier, Philippe Gigu, Brahim Chaib-draa, et al. Parametric exponential linear unit for deep convolutional neural networks. In *Machine Learning and Applications (ICMLA), 2017 16th IEEE International Conference on*, pages 207–214. IEEE, 2017.
- Forest Agostinelli, Matthew Hoffman, Peter J. Sadowski, and Pierre Baldi. Learning activation functions to improve deep neural networks. *CoRR*, abs/1412.6830, 2014.
- Simone Scardapane, Steven Van Vaerenbergh, Simone Totaro, and Aurelio Uncini. Kafnets: Kernel-based non-parametric activation functions for neural networks. *Neural Networks*, 2018.
- Ömer Faruk Ertuğrul. A novel type of activation function in artificial neural networks: Trained activation function. *Neural Networks*, 99:148–157, 2018.
- Xiaojie Jin, Chunyan Xu, Jiashi Feng, Yunchao Wei, Junjun Xiong, and Shuicheng Yan. Deep learning with s-shaped rectified linear activation units. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 1737–1743. AAAI Press, 2016.
- Gustav Fechner. *Elements of psychophysics*. New York, Holt, Rinehart and Winston, 1966.
- Stanley S Stevens. On the psychophysical law. *Psychological review*, 64(3):153, 1957.
- Leon René Sütthof, Flemming Brieger, Holger Finger, Sonja Füllhase, and Gordon Pipa. Adaptive blending units: Trainable activation functions for deep neural networks. *arXiv preprint arXiv:1806.10064*, 2018.
- Mark Harmon and Diego Klabjan. Activation ensembles for deep neural networks. *CoRR*, abs/1702.07790, 2017.
- Ian J. Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout networks. In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, ICML’13*, pages III–1319–III–1327, 2013.
- Weichen Sun, Fei Su, and Leiquan Wang. Improving deep neural networks with multi-layer maxout networks and a novel initialization method. *Neurocomputing*, 278:34 – 40, 2018.
- Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *CoRR*, abs/1312.4400, 2013.
- Simone Scardapane, Michele Scarpiniti, Danilo Comminiello, and Aurelio Uncini. Learning activation functions from data using cubic spline interpolation. In *Italian Workshop on Neural Nets*, pages 73–83. Springer, 2017.
- Edmondo Trentin. Networks with trainable amplitude of activation functions. *Neural Networks*, 14(4-5): 471–493, 2001.
- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359 – 366, 1989. ISSN 0893-6080. doi: [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8). URL <http://www.sciencedirect.com/science/article/pii/0893608089900208>.
- Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- Moshe Leshno, Vladimir Ya Lin, Allan Pinkus, and Shimon Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural networks*, 6(6):861–867, 1993.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 12 2014.
- Martin Riedmiller and Heinrich Braun. Rprop - a fast adaptive learning algorithm. Technical report, Proc. of ISCIS VII), Universitat, 1992.

- Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- Dua Dheeru and Efi Karra Taniskidou. UCI machine learning repository, 2017.
- <https://www.dcc.fc.up.pt>. <https://www.dcc.fc.up.pt>, June 2009.
- Marek Sikora and Lukasz Wróbel. Application of rule induction algorithms for analysis of data collected by seismic hazard monitoring systems in coal mines. *Archives of Mining Sciences*, 55(1):91–114, 2010.
- Bálint Antal and András Hajdu. An ensemble-based system for automatic screening of diabetic retinopathy. *Knowledge-based systems*, 60:20–27, 2014.
- Kamel Mansouri, Tine Ringsted, Davide Ballabio, Roberto Todeschini, and Viviana Consonni. Quantitative structure-activity relationship models for ready biodegradability of chemicals. *Journal of Chemical Information and Modeling*, 53(4):867–878, 2013.
- D. D. Lucas, R. Klein, J. Tannahill, D. Ivanova, S. Brandon, D. Domyancic, and Y. Zhang. Failure analysis of parameter-induced simulation crashes in climate models. *Geoscientific Model Development*, 6(4):1157–1171, 2013.
- Yann LeCun and Corinna Cortes. MNIST handwritten digit database. <http://yann.lecun.com/exdb/mnist/>, 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- A Krizhevsky and G Hinton. Learning multiple layers of features from tiny images. In *Computer Science Department, University of Toronto, Tech. Rep*, volume 1, 01 2009.



Random init.



ReLU init.

Figure 7: Examples of changes in a VAF in a 2 layer conv. network using random (top) and ReLU initialization (bottom). The blue line is the start function, the orange line is the learned function.

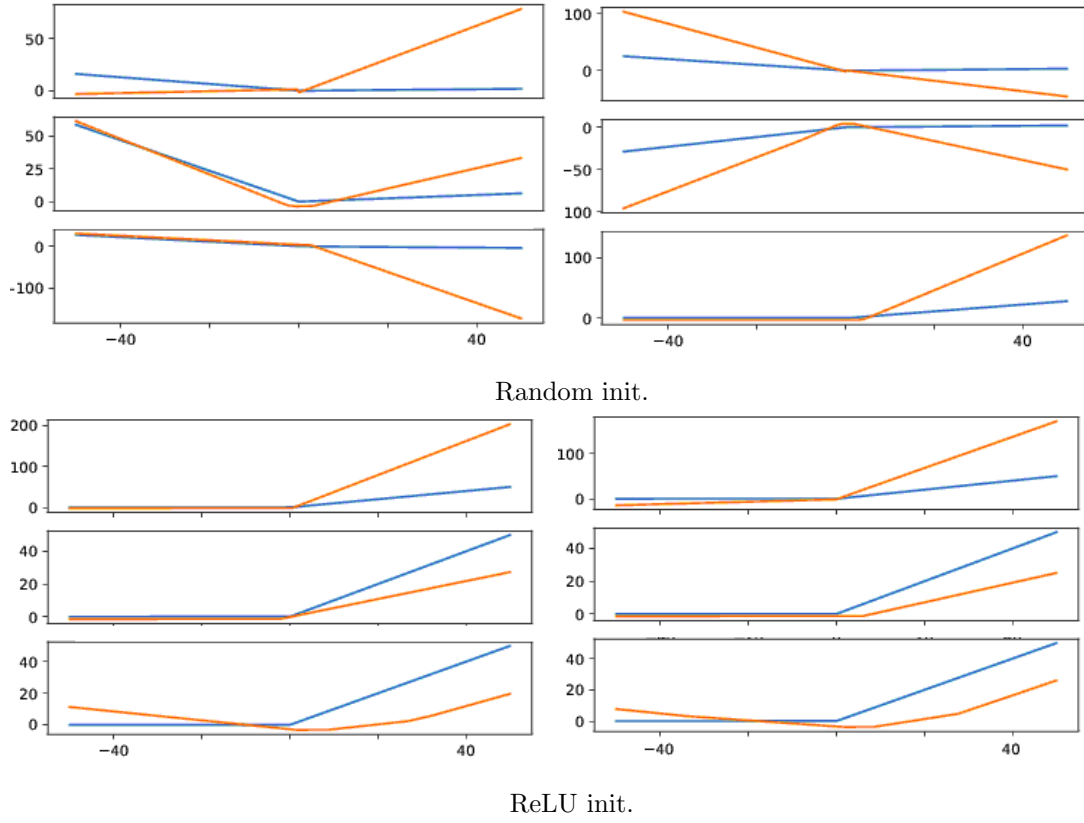


Figure 8: Examples of a resulting VAF in a 3 layer conv. using random (top) and ReLU initialization (bottom). The blue line is the start function, the orange line is the learned function.