# Multiple Hidden Layered CEFYDRA: Cluster-first Explainable FuzzY-based Deep self-Reorganizing Algorithm

Javier Viaña[1[0000-0002-0563-784X]], Stephan Ralescu[2[0000-0002-3969-1342]],

Vladik Kreinovich[3[0000-0002-1244-1650]], Anca Ralescu[4[0000-0002-7564-3540]],

and Kelly Cohen[5[0000-0002-8655-1465]]

[1,2,3,4,5] University of Cincinnati, Cincinnati OH 45219, USA
[1] Massachusetts Institute of Technology, Cambridge MA 02139, USA
vianajr@mail.uc.edu
vianajr@mit.edu

**Abstract.** We propose a deep learning algorithm that breaks with the paradigm of weights and activation functions, CEFYDRA, a network of cluster-first fuzzy-based regression algorithms. In this paper, we cover the generalization of CEFYDRA to multilayered deep architectures. First, we provide the three laws that make this generalization possible: Displacement, Substitution and Composition. Then, we obtain the update formulas for the parameters of deep hidden layers. Finally, we show the pseudocode for prediction and update. We also briefly mention the reasons to believe that this algorithm, named CEFYDRA, is explainable and has the ability for plastic reorganization.

**Keywords:** Explainable AI, Neural Networks, Fuzzy Logic, Gradient Descent, Deep Learning.

## 1    Introduction

In the last year, multiple explainable neural network algorithms were created, e.g., [1-4]. Some of which have proved to be very successful in a variety of fields, such as the discovery of exoplanets [5] or industrial automation [6]. In fact, based on the number of publications and patents published over the last decade, it can be said that the neural networks have captured the attention of most of the research carried out in computer science. Despite the advances made, a new core methodology different from the classic mechanism of weights and activation functions is rarely proposed [7-8]. To a certain extent this implies a limitation in the type of problems that we can tackle.

The research on alternative and novel basic neural architectures should not decrease. Even so, the main techniques to improve the performance in a problem solved with deep learning continue to be the iterative (trial and error) modification of the number of layers, units, or the shape of the activation functions. Meanwhile, the avenue of varying the inner functioning of the neural unit itself remains partially unexploited.

## 2 On the Learning Formulas for the Hidden Layer

One of the goals of this research is to provide an alternative to the standard of weights, biases, and activation functions of a regular neural network. We proposed a model called CEFYDRA, an architecture whose units are based on the algorithms covered in [9-13]. At every unit of the network, the inputs are grouped in fuzzy clusters (whose degree of membership is defined through Cauchy membership functions), and mapped with multidimensional logistic functions. The output of a unit is obtained by merging the predictions of all the clusters using a Takagi-Sugeno-Kang approach. All the parameters are optimized following a gradient descent learning.

In this section we focus on how to update the parameters of deep hidden layers of a CEFYDRA.

In our formulation, the upper left index in parenthesis represents the layer reference (0 for output, 1 for the first hidden starting from the right), and the bottom index represents the unit within that layer. We use $\mathbf{x_i}, \mathbf{y_i}$, and $\hat{\mathbf{y}}_\mathbf{i}$ to denote the $i^{th}$ input, output, and predicted output, respectively. The symbols $t$ and $w$ represent the output and a generic parameter of a given unit. Finally, indexes $k$ and $j$ refer to the cluster and input dimension within the unit. The membership functions are obtained from Cauchy distributions,

$$
{}^{(\lambda)}_{h_\lambda}\mu_c(\mathbf{x_i}) = \frac{1}{1 + \left\| {}^{(\lambda)}_{h_\lambda}\mathbf{u_c} \cdot {}^{(\lambda+1)}\mathbf{t(x_i)} + {}^{(\lambda)}_{h_\lambda}\mathbf{v_c} \right\|^2} = \left[ 1 + \sum_{h=1}^{H} \left( {}^{(\lambda)}_{h_\lambda}u_{ch} \cdot {}^{(\lambda+1)}t_h(\mathbf{x_i}) + {}^{(\lambda)}_{h_\lambda}v_{ch} \right)^2 \right]^{-1}. \tag{1}
$$

In [7-8], the functions for the approximation of each cluster are linear, however, we want to constrain the output of each unit between 0 and 1. Thus, we use logistic functions instead, for a generic unit $h_\lambda$ within layer $\lambda$,

$$
{}^{(\lambda)}_{h_\lambda}r_c(\mathbf{x_i}) = \frac{1}{1 + \exp\left( -{}^{(\lambda)}_{h_\lambda}\mathbf{m_c} \cdot {}^{(\lambda+1)}\mathbf{t(x_i)} - {}^{(\lambda)}_{h_\lambda}n_c \right)}. \tag{2}
$$

Therefore, the parameters of each cluster are ${}^{(\lambda)}_{h_\lambda}\mathbf{u_c}, {}^{(\lambda)}_{h_\lambda}\mathbf{v_c}, {}^{(\lambda)}_{h_\lambda}\mathbf{m_c}$ and ${}^{(\lambda)}_{h_\lambda}n_c$.

With that set up, we start from the premise that between any output of a unit from the output layer and any other output of a unit from the first hidden layer exists a linear proportionality,

$$
\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial {}^{(1)}_{h_1}w_{kj}} = K \frac{\partial {}^{(1)}t_{h_1}(\mathbf{x_i})}{\partial {}^{(1)}_{h_1}w_{kj}}, \tag{3}
$$

where $K$ is a parameter that depends on the connection between these two units. It can be demonstrated that

$$
\begin{aligned}
K = \begin{bmatrix} {}^{(1)(0)}_{h_1\ \ l} \end{bmatrix} \phi(\mathbf{x_i}) \overset{\text{def}}{=} & \left[ \sum_{c=1}^{C} {}^{(0)}_l\mu_c(\mathbf{t_i}) \right]^{-1} \sum_{c=1}^{C} \left\{ {}^{(0)}_l\mu_c(\mathbf{t_i}) \left\{ -2 {}^{(0)}_l\mu_c(\mathbf{t_i}) \right. \right. \\
& \left. \cdot \left( {}^{(0)}_l u_{ch} \cdot {}^{(1)}t_{h_1}(\mathbf{x_i}) + {}^{(0)}_l v_{ch} \right) \cdot {}^{(0)}_l u_{ch} \cdot \left[ {}^{(0)}_l r_c(\mathbf{t_i}) - \hat{y}_l(\mathbf{x_i}) \right] \right. \\
& \left. \left. + {}^{(0)}_l m_{ch} \cdot {}^{(0)}_l \varpi_c(\mathbf{t_i}) \right\} \right\}.
\end{aligned} \tag{4}
$$

where $^{(0)}_l\varpi_c(\mathbf{t_i}) \stackrel{\text{def}}{=} \left[^{(0)}_l r_c(\mathbf{t_i})\right]^2 \cdot \exp\left(^{(0)}_l s_c(\mathbf{t_i})\right)$.

Due to this decoupling in the consecutive outputs, we can expect a generalized formula that relates any two non-consecutive outputs (from any two layers) of the type

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial^{(\lambda)}_{h_\lambda} w_{kj}} = \frac{\partial^{(\lambda)} t_{h_\lambda}(\mathbf{x_i})}{\partial^{(\lambda)}_{h_\lambda} w_{kj}} \prod \phi(\mathbf{x_i}). \tag{5}$$

We use the upper left index to distinguish between the hidden layers (Fig. 1). For a given layer $\lambda$, unit $h_\lambda$, the output variable is $^{(\lambda)} t_{h_\lambda}$.
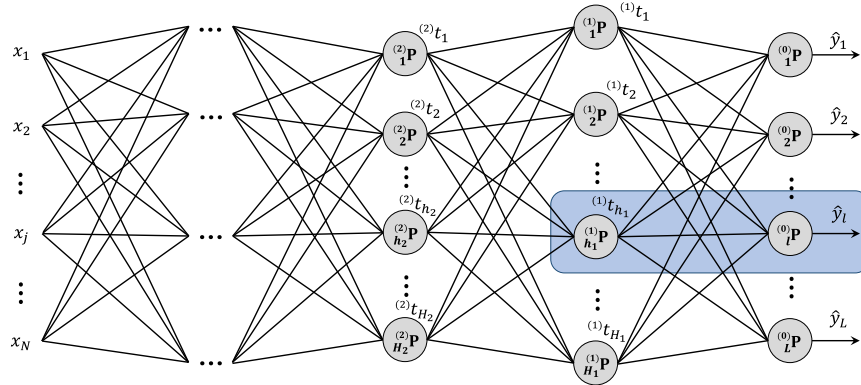


**Fig. 1.** Representation of a CEFYDRA with several hidden layers and an output layer of dimension $L$. The two units highlighted referred to equation (1).

There are three laws that we can derive from (1):

**Displacement Law**. We can consider the equation (1) between every two consecutive layers. If we use the hidden layers 2 and 1, we get

$$\frac{\partial^{(1)} t_{h_1}(\mathbf{x_i})}{\partial^{(2)}_{h_2} w_{kj}} = \left[^{(2)(1)}_{h_2\ h_1}\right] \phi(\mathbf{x_i}) \frac{\partial^{(2)} t_{h_2}(\mathbf{x_i})}{\partial^{(2)}_{h_2} w_{kj}}, \tag{6}$$

where

$$\begin{aligned}
\left[^{(2)(1)}_{h_2\ h_1}\right]\phi(\mathbf{x_i}) \stackrel{\text{def}}{=} & \left[\sum_{c=1}^{C} {}^{(1)}_{h_1}\mu_c\left(^{(2)}\mathbf{t_i}\right)\right]^{-1} \sum_{c=1}^{C}\left\{^{(1)}_{h_1}\mu_c\left(^{(2)}\mathbf{t_i}\right)\right. \\
& \cdot\left\{-2\,^{(1)}_{h_1}\mu_c\left(^{(2)}\mathbf{t_i}\right)\cdot\left[^{(1)}_{h_1}u_{ch_2}\cdot{}^{(2)}t_{h_2}(\mathbf{x_i}) + {}^{(1)}_{h_1}v_{ch_2}\right]\cdot{}^{(1)}_{h_1}u_{ch_2}\right. \\
& \left.\left.\cdot\left[^{(1)}_{h_1}r_c\left(^{(2)}\mathbf{t_i}\right) - {}^{(1)}t_{h_1}(\mathbf{x_i})\right] + {}^{(1)}_{h_1}m_{ch_2}\cdot{}^{(1)}_{h_1}\varpi_c\left(^{(2)}\mathbf{t_i}\right)\right\}\right\}.
\end{aligned} \tag{7}$$

To obtain $\left[^{(2)(1)}_{h_2\ h_1}\right]\phi(\mathbf{x_i})$, we treat $^{(1)}t_{h_1}$ as the output $\hat{y}_l$ we used in (1) and $^{(2)}\mathbf{t_i}$ instead of the input $\mathbf{t_i}$. Similarly, we also adapt the parameters of $\left[^{(2)(1)}_{h_2\ h_1}\right]\phi(\mathbf{x_i})$ to those of the $h_1^{th}$ unit of the first hidden layer. Equation (6) is the particular form of the Displacement Law between the unit $h_1$ of layer 1 and the unit $h_2$ of layer 2 (Fig. 2).
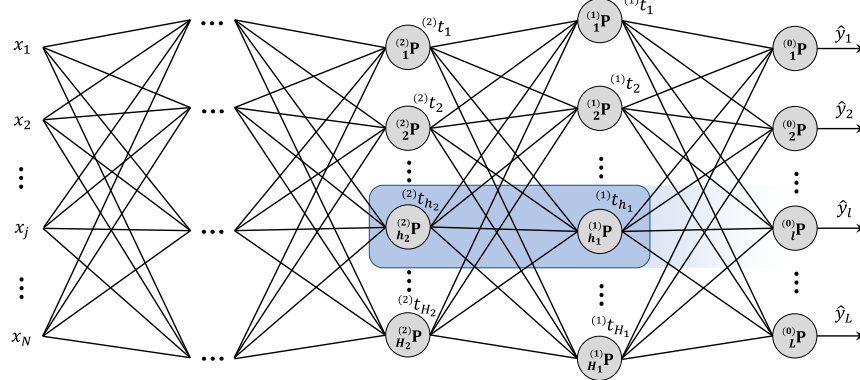
**Fig. 2.** Highlight of the two units within the CEFYDRA that refer to equation (6), i.e., the Displacement Law in its particular form between hidden layers 1 and 2.

We now translate the particular form of the Displacement Law (6) to its general form by considering the hidden layers $\lambda$ and, $\lambda - 1$. We focus on the $\alpha^{th}$ output of layer $\lambda$ ($^{(\lambda)}t_\alpha$) and the $\beta^{th}$ output of layer $\lambda - 1$ ($^{(\lambda-1)}t_\beta$) as show in Fig. 3.
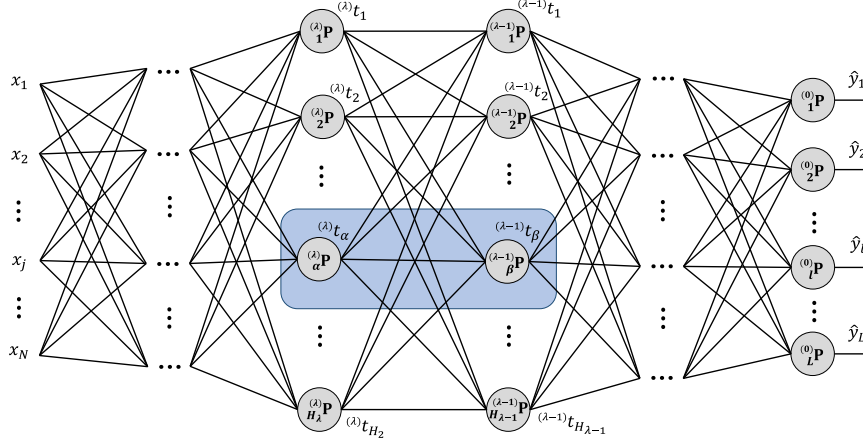


**Fig. 3.** Highlight of the two units within the CEFYDRA network that refer to equation (8), i.e., the Displacement Law in its general form between hidden layers $\lambda - 1$ and $\lambda$.

The resulting general form is

$$\frac{\partial^{(\lambda-1)}t_\beta(\mathbf{x_i})}{\partial_\alpha^{(\lambda)}w_{kj}} = \begin{bmatrix} (\lambda)(\lambda-1) \\ \alpha \quad \beta \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial_\alpha^{(\lambda)}w_{kj}}. \tag{8}$$

**Substitution Law.** This Law comes from the application of the chain rule in the derivatives of (1) to consider any other parameter of the system,

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial p} \frac{\partial p}{\partial_{h_1}^{(1)}w_{kj}} = \begin{bmatrix} (1)(0) \\ h_1 \quad l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(1)}t_{h_1}(\mathbf{x_i})}{\partial p} \frac{\partial p}{\partial_{h_1}^{(1)}w_{kj}}, \tag{9}$$

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial p} = \begin{bmatrix} (1)(0) \\ h_1 \quad l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(1)} t_{h_1}(\mathbf{x_i})}{\partial p}. \tag{10}$$

If we consider the parameter $p = \overset{(2)}{_{h_2}}w_{kj}$ for the layer 2 in formula (10), then

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}} = \begin{bmatrix} (1)(0) \\ h_1 \quad l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(1)} t_{h_1}(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}. \tag{11}$$

In the general form, we can take the general expression (8) obtained for the Displacement Law and consider any layer $\pi$ and any unit $\gamma$ within the layer $\pi$,

$$\frac{\partial^{(\lambda-1)} t_\beta(\mathbf{x_i})}{\partial \overset{(\pi)}{_\gamma}w_{kj}} = \begin{bmatrix} (\lambda)(\lambda-1) \\ \alpha \quad \beta \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(\lambda)} t_\alpha(\mathbf{x_i})}{\partial \overset{(\pi)}{_\gamma}w_{kj}}. \tag{12}$$

**Composition Law.** The Composition Law is achieved from the combination of both the Displacement and Substitution Laws. For the particular cases (6) and (10),

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}} = \begin{bmatrix} (2)(1) \\ h_2 \ h_1 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (1)(0) \\ h_1 \quad l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(2)} t_{h_2}(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}. \tag{13}$$

Equation (13) is certainly interesting, as it shows that $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}$ is not only dependent on $\begin{bmatrix} (1)(0) \\ h_1 \ l \end{bmatrix}\phi(\mathbf{x_i})$ and $\frac{\partial^{(2)} t_{h_2}(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}$ (which can be understood since in $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}$ we have both the information of the unit $h_2$ of the second hidden layer and the information of the unit $l$ of the output layer) but it is also dependent on $\begin{bmatrix} (2)(1) \\ h_2 \ h_1 \end{bmatrix}\phi(\mathbf{x_i})$, i.e., the information of the unit $h_1$ of the first hidden layer. This might seem confusing, since $h_1$ is not present in either $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}$ nor $\frac{\partial^{(2)} t_{h_2}(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}$. In fact, there is no specification for $h_1$, therefore $h_1$ could be any unit of the first hidden layer. This comes from the fact that the system is overdetermined due to the fully connected network. In other words, there are many paths that connect both $h_2$ and $l$, and they all provide a different result of $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}$. Thus, we suggest that $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}}$ should be calculated by averaging the value obtained from all the different paths (Fig. 4),

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}} = \frac{1}{H_1} \sum_{h_1=1}^{H_1} \left[ \begin{bmatrix} (2)(1) \\ h_2 \ h_1 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (1)(0) \\ h_1 \quad l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(2)} t_{h_2}(\mathbf{x_i})}{\partial \overset{(2)}{_{h_2}}w_{kj}} \right]. \tag{14}$$
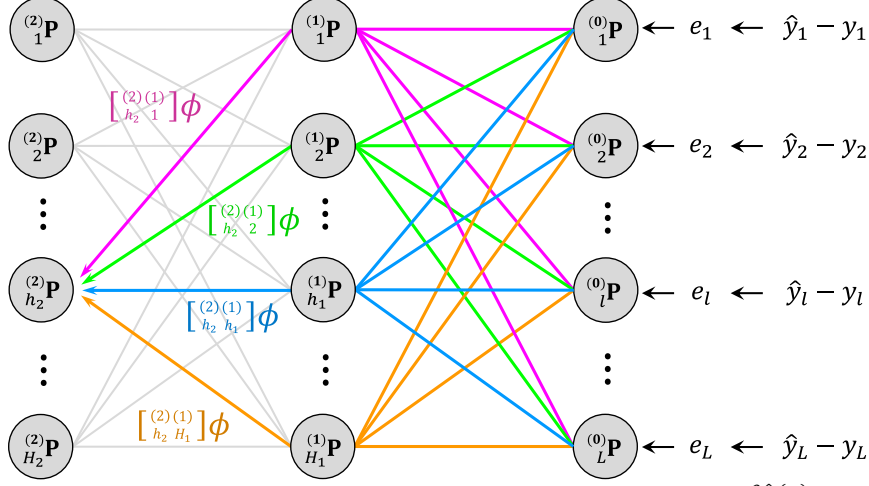
**Fig. 4.** Graphical representation of all the possible paths in the calculation of $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{h_2}^{(2)} w_{kj}}$. The term $e_{il}$ represents the error of each output unit $y_{il} - \hat{y}_l(\mathbf{x_i})$.

To calculate the general form of (13), we can now study the expression of $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}}$. We start from the expression of the particular Composition Law, and we apply the Substitution Law, i.e., instead of deriving with respect to $_{h_2}^{(2)} w_{kj}$, we derive with respect to $_{h_3}^{(3)} w_{kj}$, where 3 is the third hidden layer and $h_3$ is a generic hidden neuron, again, the factors $_l^{(0)}\phi(\mathbf{x_i})$ and $_{h_1}^{(1)}\phi(\mathbf{x_i})$ do not depend on $_{h_2}^{(2)} w_{kj}$ or $_{h_3}^{(3)} w_{kj}$, thus

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}} = \begin{bmatrix} (2)(1) \\ h_2 \ h_1 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (1)(0) \\ h_1 \ l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(2)} t_{h_2}(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}}, \tag{15}$$

now we consider the Displacement Law for layers 3 and 2, let us call $^{(3)}t_{h_3}$ the input before $^{(2)}t_{h_2}$,

$$\frac{\partial^{(2)} t_{h_2}(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}} = \begin{bmatrix} (3)(2) \\ h_3 \ h_2 \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(3)} t_{h_3}(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}}, \tag{16}$$

and we combine it with the result obtained in (15),

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}} = \begin{bmatrix} (3)(2) \\ h_3 \ h_2 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (2)(1) \\ h_2 \ h_1 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (1)(0) \\ h_1 \ l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(3)} t_{h_3}(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}}. \tag{17}$$

As it happened in (13), this particular case of the Composition Law (17), which defines the learning of the third layer's parameters, also depends on the path and the connections chosen in layers 1 and 2. Hence we use the average of all the possible paths as shown in Fig. 5,

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}} = \frac{1}{H_1 \cdot H_2} \sum_{h_2=1}^{H_2} \sum_{h_1=1}^{H_1} \left\{ \begin{bmatrix} (3)(2) \\ h_3 \ h_2 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (2)(1) \\ h_2 \ h_1 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (1)(0) \\ h_1 \ l \end{bmatrix} \phi(\mathbf{x_i}) \frac{\partial^{(3)} t_{h_3}(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}} \right\}, \tag{18}$$

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}} = \frac{\partial^{(3)} t_{h_3}(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}} \frac{1}{H_2} \sum_{h_2=1}^{H_2} \left\{ \begin{bmatrix} (3)(2) \\ h_3 \ h_2 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \frac{1}{H_1} \sum_{h_1=1}^{H_1} \left\{ \begin{bmatrix} (2)(1) \\ h_2 \ h_1 \end{bmatrix} \phi(\mathbf{x_i}) \cdot \begin{bmatrix} (1)(0) \\ h_1 \ l \end{bmatrix} \phi(\mathbf{x_i}) \right\} \right\}. \quad (19)$$
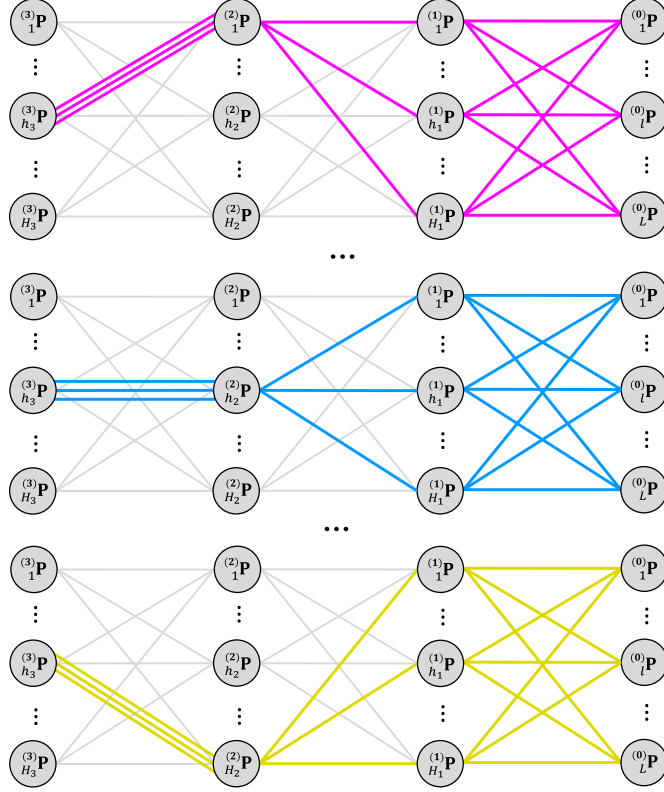


**Fig. 5.** Representation of all the paths to be considered for the calculation of $\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{h_3}^{(3)} w_{kj}}$.

Generalizing the solution (17), we can then conclude that

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{\alpha}^{(\lambda)} w_{kj}} = \frac{\partial^{(\lambda)} t_{\alpha}(\mathbf{x_i})}{\partial_{\alpha}^{(\lambda)} w_{kj}} \prod_{\pi=0}^{\lambda-1} \begin{bmatrix} (\pi+1)(\pi) \\ \delta \quad \gamma \end{bmatrix} \phi(\mathbf{x_i}), \quad (20)$$

$$\frac{\partial \hat{y}_l(\mathbf{x_i})}{\partial_{\alpha}^{(\lambda)} w_{kj}} = \frac{\partial^{(\lambda)} t_{\alpha}(\mathbf{x_i})}{\partial_{\alpha}^{(\lambda)} w_{kj}} \prod_{\pi=1}^{\lambda-1} \left\{ \begin{bmatrix} (\pi+1)(\pi) \\ \delta \quad \gamma \end{bmatrix} \phi(\mathbf{x_i}) \right\} \cdot \begin{bmatrix} (1)(0) \\ h_1 \ l \end{bmatrix} \phi(\mathbf{x_i}). \quad (21)$$

where $\gamma$ and $\delta$ represent the units within the layers $\pi$ and $\pi + 1$ respectively, and

$$\left[\begin{smallmatrix}(\pi+1)(\pi)\\\delta\quad\gamma\end{smallmatrix}\right]\phi(\mathbf{x_i}) = \left[\sum_{c=1}^{C} {}_{\gamma}^{(\pi)}\mu_c\big({}^{(\pi+1)}\mathbf{t_i}\big)\right]^{-1} \sum_{c=1}^{C}\Big\{{}_{\gamma}^{(\pi)}\mu_c\big({}^{(\pi+1)}\mathbf{t_i}\big)$$
$$\cdot\Big\{-2{}_{\gamma}^{(\pi)}\mu_c\big({}^{(\pi+1)}\mathbf{t_i}\big)\cdot\big[{}_{\gamma}^{(\pi)}u_{c\delta}\cdot{}^{(\pi+1)}t_{\delta}(\mathbf{x_i}) + {}_{\gamma}^{(\pi)}v_{c\delta}\big]\cdot{}_{\gamma}^{(\pi)}u_{c\delta}$$
$$\cdot\big[{}_{\gamma}^{(\pi)}r_c\big({}^{(\pi+1)}\mathbf{t_i}\big) - {}^{(\pi)}t_{\gamma}(\mathbf{x_i})\big] + {}_{\gamma}^{(\pi)}m_{c\delta}\cdot{}_{\gamma}^{(\pi)}\varpi_c\big({}^{(\pi+1)}\mathbf{t_i}\big)\Big\}\Big\}. \tag{22}$$

For the case where $\pi$ is identifying the last hidden layer (the first layer starting from the left), the inputs that receive the units are indeed the inputs of the system, so ${}^{(\pi+1)}\mathbf{t_i}\big|_{\pi=\Pi} = \mathbf{x_i}$, where $\Pi$ represents the last hidden layer. The formula (20), aligns with our preliminary reasoning and the expected form (5).

Again we consider the averaging of all the possible paths to correct (20),

$$\frac{\partial\hat{y}_l(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}} = \frac{1}{\prod_{\pi=1}^{\lambda-1}H_{\pi}}\sum_{h_1,\dots,h_{\lambda-1}=1}^{H_1,\dots,H_{\lambda-1}}\left\{\frac{\partial{}^{(\lambda)}t_{\alpha}(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}}\prod_{\pi=1}^{\lambda-1}\Big\{\big[\begin{smallmatrix}(\pi+1)(\pi)\\\delta\quad\gamma\end{smallmatrix}\big]\phi(\mathbf{x_i})\Big\}\cdot\big[\begin{smallmatrix}(1)(0)\\h_1\quad l\end{smallmatrix}\big]\phi(\mathbf{x_i})\right\}, \tag{23}$$

$$\frac{\partial\hat{y}_l(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}} = \frac{\partial{}^{(\lambda)}t_{\alpha}(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}}\frac{1}{\prod_{\pi=1}^{\lambda-1}H_{\pi}}\sum_{h_1,\dots,h_{\lambda-1}=1}^{H_1,\dots,H_{\lambda-1}}\left\{\prod_{\pi=1}^{\lambda-1}\Big\{\big[\begin{smallmatrix}(\pi+1)(\pi)\\\delta\quad\gamma\end{smallmatrix}\big]\phi(\mathbf{x_i})\Big\}\cdot\big[\begin{smallmatrix}(1)(0)\\h_1\quad l\end{smallmatrix}\big]\phi(\mathbf{x_i})\right\}. \tag{24}$$

Note that we cannot replace the sum of the product, $\sum_{h_1,\dots,h_{\lambda-1}=1}^{H_1,\dots,H_{\lambda-1}}\big\{\prod_{\pi=1}^{\lambda-1}\big[\begin{smallmatrix}(\pi+1)(\pi)\\\delta\quad\gamma\end{smallmatrix}\big]\phi(\mathbf{x_i})\big\}$, for the product of the sum, $\prod_{\pi=1}^{\lambda-1}\big\{\sum_{h_{\pi}=1}^{H_{\pi}}\big[\begin{smallmatrix}(\pi+1)(\pi)\\\delta\quad\gamma\end{smallmatrix}\big]\phi(\mathbf{x_i})\big\}$, since the terms $\big[\begin{smallmatrix}(\pi+1)(\pi)\\\delta\quad\gamma\end{smallmatrix}\big]\phi$ are not independent, i.e., all the intermediate units considered must create a connected path. Instead, we can group the series in the following way

$$\frac{1}{\prod_{\pi=1}^{\lambda-1}H_{\pi}}\sum_{h_1,\dots,h_{\lambda-1}=1}^{H_1,\dots,H_{\lambda-1}}\left\{\prod_{\pi=1}^{\lambda-1}\Big\{\big[\begin{smallmatrix}(\pi+1)(\pi)\\\delta\quad\gamma\end{smallmatrix}\big]\phi(\mathbf{x_i})\Big\}\cdot\big[\begin{smallmatrix}(1)(0)\\h_1\quad l\end{smallmatrix}\big]\phi(\mathbf{x_i})\right\} = \frac{1}{H_{\lambda-1}}\sum_{h_{\lambda-1}=1}^{H_{\lambda-1}}\left\{\big[\begin{smallmatrix}(\lambda)(\lambda-1)\\\alpha\quad\beta\end{smallmatrix}\big]\phi(\mathbf{x_i})\right.$$
$$\left.\cdot\left\{\dots\frac{1}{H_2}\left\{\sum_{h_2=1}^{H_2}\big[\begin{smallmatrix}(3)(2)\\h_3\quad h_2\end{smallmatrix}\big]\phi(\mathbf{x_i})\frac{1}{H_1}\left\{\sum_{h_1=1}^{H_1}\Big\{\big[\begin{smallmatrix}(2)(1)\\h_2\quad h_1\end{smallmatrix}\big]\phi(\mathbf{x_i})\cdot\big[\begin{smallmatrix}(1)(0)\\h_1\quad l\end{smallmatrix}\big]\phi(\mathbf{x_i})\Big\}\right\}\right\}\right\}\right\}. \tag{25}$$

Note that in (25), the term $\big[\begin{smallmatrix}(1)(0)\\h_1\quad l\end{smallmatrix}\big]\phi(\mathbf{x_i})$ is not the argument of a hypothetical final series because all the paths have to end in the output unit $l$. However, when we study the partial derivative of the loss $J$ with respect to the generic parameter ${}_{\alpha}^{(\lambda)}w_{kj}$, and we integrate the definition of $\frac{\partial\hat{y}_l(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}}$ given by (24), we can see how this new formula has one more series whose argument is indeed $\big[\begin{smallmatrix}(1)(0)\\h_1\quad l\end{smallmatrix}\big]\phi(\mathbf{x_i})$ (and the error term $e_{il} = y_{il} - \hat{y}_l(\mathbf{x_i})$),

$$\frac{\partial J(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}} = -\frac{1}{L}\sum_{l=1}^{L}\left\{[y_{il} - \hat{y}_l(\mathbf{x_i})]\frac{\partial\hat{y}_l(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}}\right\} = -\frac{1}{L}\sum_{l=1}^{L}\left\{e_{il}\frac{\partial\hat{y}_l(\mathbf{x_i})}{\partial{}_{\alpha}^{(\lambda)}w_{kj}}\right\} = \tag{26}$$

$$= -\frac{\partial^{(\lambda)} t_\alpha(\mathbf{x_i})}{\partial^{(\lambda)}_\alpha w_{kj}} \frac{1}{\prod_{\pi=1}^{\lambda-1} H_\pi} \sum_{h_1,\dots,h_{\lambda-1}=1}^{H_1,\dots,H_{\lambda-1}} \left\{ \prod_{\pi=1}^{\lambda-1} \left\{ \left[ \begin{smallmatrix} (\pi+1)(\pi) \\ \delta \quad\; \gamma \end{smallmatrix} \right] \phi(\mathbf{x_i}) \right\} \right.$$

$$\left. \cdot \frac{1}{L} \sum_{l=1}^{L} \left\{ \left[ \begin{smallmatrix} (1)(0) \\ h_1 \;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot e_{il} \right\} \right\}.$$

So, following with the development of (25), we can replace the term $\left[ \begin{smallmatrix} (1)(0) \\ h_1 \;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i})$ for $\sum_{l=1}^{L} \left\{ \left[ \begin{smallmatrix} (1)(0) \\ h_1 \;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot e_{il} \right\}$ in order to obtain the definition of $\frac{\partial J(\mathbf{x_i})}{\partial^{(\lambda)}_\alpha w_{kj}}$,

$$\frac{1}{\prod_{\pi=1}^{\lambda-1} H_\pi} \sum_{h_1,\dots,h_{\lambda-1}=1}^{H_1,\dots,H_{\lambda-1}} \left\{ \prod_{\pi=1}^{\lambda-1} \left\{ \left[ \begin{smallmatrix} (\pi+1)(\pi) \\ \delta \quad\; \gamma \end{smallmatrix} \right] \phi(\mathbf{x_i}) \right\} \cdot \sum_{l=1}^{L} \left\{ \left[ \begin{smallmatrix} (1)(0) \\ h_1 \;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot e_{il} \right\} \right\} =$$

$$\frac{1}{H_{\lambda-1}} \sum_{h_{\lambda-1}=1}^{H_{\lambda-1}} \left\{ \left[ \begin{smallmatrix} (\lambda)(\lambda-1) \\ \alpha \quad\;\; \beta \end{smallmatrix} \right] \phi(\mathbf{x_i}) \right. \tag{27}$$

$$\left. \cdot \left\{ \dots \cdot \frac{1}{H_2} \left\{ \sum_{h_2=1}^{H_2} \left[ \begin{smallmatrix} (3)(2) \\ h_3\; h_2 \end{smallmatrix} \right] \phi(\mathbf{x_i}) \frac{1}{H_1} \left\{ \sum_{h_1=1}^{H_1} \left\{ \left[ \begin{smallmatrix} (2)(1) \\ h_2\; h_1 \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot \frac{1}{L} \sum_{l=1}^{L} \left\{ \left[ \begin{smallmatrix} (1)(0) \\ h_1\;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot e_{il} \right\} \right\} \right\} \right\} \right\}.$$

At this point, we can define a recursive function as follows,

$$\begin{cases} \overset{(\pi+1)}{\underset{\delta}{\Gamma}}(\mathbf{x_i}) = \dfrac{1}{H_\pi} \sum_{\gamma=1}^{H_\pi} \left\{ \left[ \begin{smallmatrix} (\pi+1)(\pi) \\ \delta \quad\; \gamma \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot \overset{(\pi)}{\underset{\gamma}{\Gamma}}(\mathbf{x_i}) \right\} & \text{if } \pi > -1 \\[3mm] \overset{(0)}{\underset{l}{\Gamma}}(\mathbf{x_i}) = e_{il} = y_{il} - \hat{y}_l(\mathbf{x_i}) \end{cases} \tag{28}$$

so that

$$\overset{(1)}{\underset{h_1}{\Gamma}}(\mathbf{x_i}) = \frac{1}{L} \sum_{l=1}^{L} \left\{ \left[ \begin{smallmatrix} (1)(0) \\ h_1\;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot \overset{(0)}{\underset{l}{\Gamma}}(\mathbf{x_i}) \right\} = \frac{1}{L} \sum_{l=1}^{L} \left\{ \left[ \begin{smallmatrix} (1)(0) \\ h_1\;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot e_{il} \right\}. \tag{29}$$

Thus, using the recursive function defined at (28), we can redefine equation (27) as

$$\frac{1}{\prod_{\pi=1}^{\lambda-1} H_\pi} \sum_{h_1,\dots,h_{\lambda-1}=1}^{H_1,\dots,H_{\lambda-1}} \left\{ \prod_{\pi=1}^{\lambda-1} \left\{ \left[ \begin{smallmatrix} (\pi+1)(\pi) \\ \delta \quad\; \gamma \end{smallmatrix} \right] \phi(\mathbf{x_i}) \right\} \cdot \sum_{l=1}^{L} \left\{ \left[ \begin{smallmatrix} (1)(0) \\ h_1\;\; l \end{smallmatrix} \right] \phi(\mathbf{x_i}) \cdot e_{il} \right\} \right\} = \overset{(\lambda)}{\underset{\alpha}{\Gamma}}(\mathbf{x_i}). \tag{30}$$

Finally, the formula of $\frac{\partial J(\mathbf{x_i})}{\partial^{(\lambda)}_\alpha w_{kj}}$ becomes

$$\frac{\partial J(\mathbf{x_i})}{\partial^{(\lambda)}_\alpha w_{kj}} = -\overset{(\lambda)}{\underset{\alpha}{\Gamma}}(\mathbf{x_i}) \frac{\partial^{(\lambda)} t_\alpha(\mathbf{x_i})}{\partial^{(\lambda)}_\alpha w_{kj}}. \tag{31}$$

This is coherent with what we have for the case of a single hidden layer, i.e., equation (1). Note that $\overset{(\lambda)}{\underset{\alpha}{\Gamma}}(\mathbf{x_i})$ does not depend on the cluster $k$ nor on the dimension $j$ despite the fact that the parameter $\overset{(\lambda)}{\underset{\alpha}{w}}_{kj}$ of equation (31) does.

The calculation of $\frac{\partial \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})}{\partial \overset{(\lambda)}{\underset{\alpha}{w}}_{kj}}$ is similar to that studied in [7-8] but instead of having the real input $\mathbf{x_i}$ as input of the layer, we have $\overset{(\lambda+1)}{\mathbf{t}}(\mathbf{x_i})$, i.e., we substitute $x_{ij}$ for $\overset{(\lambda+1)}{t_j}(\mathbf{x_i})$. We also denote $\overset{(1)}{\underset{h}{\xi}}_k(\mathbf{x_i}) \overset{\text{def}}{=} \overset{(1)}{\underset{h}{\mu}}_k(\mathbf{x_i}) \cdot \left[\sum_{c=1}^{C} \overset{(1)}{\underset{h}{\mu}}_c(\mathbf{x_i})\right]^{-1}$. After a short straight-forward development, we obtain the partial derivatives of $\overset{(\lambda)}{t_\alpha}(\mathbf{x_i})$,

$$\frac{\partial \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})}{\partial \overset{(\lambda)}{\underset{\alpha}{m}}_{kj}} = \overset{(\lambda)}{\underset{\alpha}{\xi}}_k(\mathbf{x_i}) \cdot \overset{(\lambda)}{\underset{\alpha}{\varpi}}_k(\mathbf{x_i}) \cdot \overset{(\lambda+1)}{t_j}(\mathbf{x_i}), \tag{32}$$

$$\frac{\partial \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})}{\partial \overset{(\lambda)}{\underset{\alpha}{n}}_{k}} = \overset{(\lambda)}{\underset{\alpha}{\xi}}_k(\mathbf{x_i}) \cdot \overset{(\lambda)}{\underset{\alpha}{\varpi}}_k(\mathbf{x_i}), \tag{33}$$

$$\frac{\partial \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})}{\partial \overset{(\lambda)}{\underset{\alpha}{u}}_{kj}} = -2\overset{(\lambda)}{\underset{\alpha}{\xi}}_k(\mathbf{x_i}) \cdot \overset{(\lambda)}{\underset{\alpha}{\mu}}_k(\mathbf{x_i}) \cdot \left[\overset{(\lambda)}{\underset{\alpha}{r}}_k(\mathbf{x_i}) - \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})\right] \cdot \left(\overset{(\lambda)}{\underset{\alpha}{u}}_{kj} \cdot \overset{(\lambda+1)}{t_j}(\mathbf{x_i}) + \overset{(\lambda)}{\underset{\alpha}{v}}_{kj}\right)$$
$$\cdot \overset{(\lambda+1)}{t_j}(\mathbf{x_i}), \tag{34}$$

$$\frac{\partial \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})}{\partial \overset{(\lambda)}{\underset{\alpha}{v}}_{kj}} = -2\overset{(\lambda)}{\underset{\alpha}{\xi}}_k(\mathbf{x_i}) \cdot \overset{(\lambda)}{\underset{\alpha}{\mu}}_k(\mathbf{x_i}) \cdot \left[\overset{(\lambda)}{\underset{\alpha}{r}}_k(\mathbf{x_i}) - \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})\right] \cdot \left(\overset{(\lambda)}{\underset{\alpha}{u}}_{kj} \cdot \overset{(\lambda+1)}{t_j}(\mathbf{x_i}) + \overset{(\lambda)}{\underset{\alpha}{v}}_{kj}\right). \tag{35}$$

Using the learning rate and (31) we can then get

$$\Delta \overset{(\lambda)}{\underset{\alpha}{w}}_{kj} = -\eta \frac{\partial J(\mathbf{x_i})}{\partial \overset{(\lambda)}{\underset{\alpha}{w}}_{kj}} = \eta \cdot \overset{(\lambda)}{\underset{\alpha}{\Gamma}}(\mathbf{x_i}) \frac{\partial \overset{(\lambda)}{t_\alpha}(\mathbf{x_i})}{\partial \overset{(\lambda)}{\underset{\alpha}{w}}_{kj}}, \tag{36}$$

where the update rule of $\overset{(\lambda)}{\underset{\alpha}{w}}_{kj}$ is as follows,

$$\overset{(\lambda)}{\underset{\alpha}{w}}_{kj}^{new} = \Delta \overset{(\lambda)}{\underset{\alpha}{w}}_{kj} + \overset{(\lambda)}{\underset{\alpha}{w}}_{kj}^{old}. \tag{37}$$

## 3 On How to Perform the Update While Minimizing the Calculations

The update of the parameters requires the calculation of all the combinations of the paths that link the unit of the parameter and the output units. This would be a heavy calculation if it had to be recalculated at each layer. However, if we are performing the update from right to left (starting with the output layer), then we can recycle the calculations as it happens in the backpropagation method of neural networks.

In this section we explain how this recycling works and which variables should be stored at each layer in order to be retrieved for the updates of the consecutive layers.

First, we need to input an observation $(\mathbf{x_i}, \mathbf{y_i})$ and calculate the prediction with the current parameters of the CEFYDRA, $\hat{\mathbf{y}}_i = f(\mathbf{x_i})$. Then, we start with the backpropagation phase. As mentioned, we move from the right to left.

We start with the calculation of all the errors of the output units $\{e_{i1}, e_{i2}, \dots, e_{iL}\}$ and we store these values.

We move on to the first hidden layer, and we go through each of the units within the layer. At each unit, we obtain all the $\phi(\mathbf{x_i})$ variables that link the current unit of the hidden layer with all the units of the output layer, calculating the value of $\overset{(1)}{\underset{h_1}{\Gamma}}(\mathbf{x_i}) = \frac{1}{L}\sum_{l=1}^{L}\left\{e_{il} \cdot \left[\overset{(1)(0)}{\underset{h_1}{}}{}_l\right]\phi(\mathbf{x_i})\right\}$. Note that we do not need to recalculate the errors at each $\overset{(1)}{\underset{h_1}{\Gamma}}(\mathbf{x_i})$; instead we retrieve the values from the memory. Then, we store the values of $\left\{\overset{(1)}{\underset{1}{\Gamma}}(\mathbf{x_i}), \overset{(1)}{\underset{2}{\Gamma}}(\mathbf{x_i}), \dots, \overset{(1)}{\underset{H_1}{\Gamma}}(\mathbf{x_i})\right\}$ so that we can retrieve them in the next layer. The time complexity of this step is $O(H_1 \times L)$, and the space complexity $O(H_1)$.

Before moving on to the next layer, we calculate the partial derivatives $\frac{\partial \overset{(1)}{}t_{h_1}(\mathbf{x_i})}{\partial \overset{(1)}{h_1}w_{kj}}$ for all the hidden units of the first hidden layer, for all the parameters in each unit. In a fully connected network, at each unit of the first hidden layer we have a total of $C \times [(H_2 \times 3) + 1]$ parameters, where $C$ is the number of clusters, $H_2$ the number of input dimensions in the first hidden layer, 3 represents the variables $\overset{(1)}{h_1}u_{kj}$, $\overset{(1)}{h_1}v_{kj}$, and $\overset{(1)}{h_1}m_{kj}$, and 1 represents the independent variable $\overset{(1)}{h_1}n_k$. Finally, we get the values of $\frac{\partial J(\mathbf{x_i})}{\partial \overset{(1)}{h_1}w_{kj}}$ and then using the learning rates we calculate the deltas $\Delta \overset{(1)}{h_1}w_{kj}$, which need to be stored for the update phase. The time complexity of this step is $O(H_1 \times C \times [(H_2 \times 3) + 1])$ (if all the units have the same number of clusters), and the space complexity is $O(H_1 \times C \times [(H_2 \times 3) + 1])$ which comes from the storage of the deltas.

Next, we move to the second hidden layer. At this point, the process is the same with the exception that for the calculation of $\left\{\overset{(2)}{\underset{1}{\Gamma}}(\mathbf{x_i}), \overset{(2)}{\underset{2}{\Gamma}}(\mathbf{x_i}), \dots, \overset{(2)}{\underset{H_2}{\Gamma}}(\mathbf{x_i})\right\}$, instead of using the error terms as we did for the first hidden layer, we use the values of $\left\{\overset{(1)}{\underset{1}{\Gamma}}(\mathbf{x_i}), \overset{(1)}{\underset{2}{\Gamma}}(\mathbf{x_i}), \dots, \overset{(1)}{\underset{H_1}{\Gamma}}(\mathbf{x_i})\right\}$.

The Algorithms 1, 2, and 3 provide the pseudocode for all the three phases, prediction, backpropagation, and update, respectively.

---

**Algorithm 1:** Prediction of $\hat{\mathbf{y}}_\mathbf{i} = f(\mathbf{x}_\mathbf{i})$ where $f(\cdot)$ is the entire CEYDRA. Time complexity $O(\Pi \times H_\lambda)$ where $H_\lambda = \max_{\pi \in \{0,1,\dots,\Pi\}} H_\pi$. Space complexity $O(\Pi \times H_\lambda)$.

---

**Input:** The observation vector $\mathbf{x}_\mathbf{i}$ of dimension $N$.
**Output:** The prediction vector $\hat{\mathbf{y}}_\mathbf{i}$ of dimension $L$.

1   $^{(\Pi+1)}\mathbf{t_i} \leftarrow \mathbf{x_i}$
2   Descending order of layers: $\mathbf{\Lambda} \leftarrow \{\Pi, \Pi - 1, \dots, 1, 0\}$
3   **for** each **layer** $\lambda$ **in** $\mathbf{\Lambda}$ **do**
4     **for** each unit $\alpha$ **in** layer $\lambda$ **do**

5       $^{(\lambda)}_\alpha\mathbf{P}(\mathbf{t_i}) \leftarrow \left[\sum_{c=1}^{C} {}^{(\lambda)}_\alpha\mu_c\big(^{(\lambda+1)}\mathbf{t_i}\big) \cdot {}^{(\lambda)}_\alpha r_c\big(^{(\lambda+1)}\mathbf{t_i}\big)\right] \cdot \left[\sum_{c=1}^{C} {}^{(\lambda)}_\alpha\mu_c\big(^{(\lambda+1)}\mathbf{t_i}\big)\right]^{-1}$

6       $^{(\lambda)}_\alpha t_i \leftarrow {}^{(\lambda)}_\alpha\mathbf{P}(\mathbf{t_i})$

7       $^{(\lambda)}\mathbf{t_i}[\alpha] \leftarrow {}^{(\lambda)}_\alpha t_i$

8     **end**
9   **end**
10   $\hat{\mathbf{y}}_\mathbf{i} \leftarrow {}^{(0)}\mathbf{t_i}$
11   **return** $\hat{\mathbf{y}}_\mathbf{i}$

---

---

**Algorithm 2:** Backpropagation process to calculate the values of $\Delta^{(\lambda)}_\alpha w_{kj}$ of the CEFYDRA. The time complexity is $O\big(\Pi \times H_\lambda \times (H_{\lambda-1} + C \times H_{\lambda+1})\big)$ where $H_\lambda \times H_{\lambda-1} = \max_{\pi \in \{1,\dots,\Pi\}} H_\pi \times H_{\pi-1}$ or $H_{\lambda+1} \times H_\lambda = \max_{\pi \in \{0,1,\dots,\Pi\}} H_{\pi+1} \times H_\pi$, thus we can simplify the previous expressions to $O(\Pi \times C \times H_{\lambda+1} \times H_\lambda)$ with $H_{\lambda+1} \times H_\lambda = \max_{\pi \in \{0,1,\dots,\Pi\}} H_{\pi+1} \times H_\pi$. The space complexity is also $O(\Pi \times C \times H_{\lambda+1} \times H_\lambda)$, which comes from the space needed by $\Delta\mathbf{W}$.

---

**Inputs:** $\mathbf{x_i}, \mathbf{y_i}, \hat{\mathbf{y}}_\mathbf{i}$
**Outputs:** The update structure $\Delta\mathbf{W}$, organized in a [Layer, Unit] fashion, which has as many entries as the number of units in the system, and each of those entries $\Delta\mathbf{W}[\lambda, \alpha]$ is a three-dimensional matrix that stores the values of $\Delta^{(\lambda)}_\alpha w_{kj}$.

1   Initialize with null values: $\Delta\mathbf{W} \leftarrow \mathbf{0}$
2   **for** each unit $l$ **in layer** 0 **do**
3     $\overset{(0)}{\underset{l}{\prod}}(\mathbf{x_i}) \leftarrow y_{il} - \hat{y}_l(\mathbf{x_i})$
4   **end**
5   Ascending order of layers: $\mathbf{\Lambda} \leftarrow \{1, \dots, \Pi - 1, \Pi\}$;
6   **for** each **layer** $\lambda$ **in** $\mathbf{\Lambda}$ **do**
7     **for** each unit $\alpha$ **in** layer $\lambda$ **do**

---

8       **for** each unit $\beta$ **in layer** $\lambda - 1$ **do**

9
$$\left[ {}^{(\lambda)(\lambda-1)}_{\ \alpha\quad\ \beta} \right] \phi(\mathbf{x_i}) \leftarrow \left[ \sum_{c=1}^{C} {}^{(\lambda-1)}_{\ \ \ \ \beta}\mu_c\left({}^{(\lambda)}\mathbf{t_i}\right) \right]^{-1} \sum_{c=1}^{C} \left\{ {}^{(\lambda-1)}_{\ \ \ \ \beta}\mu_c\left({}^{(\lambda)}\mathbf{t_i}\right) \right.$$
$$\cdot \left\{ -2\,{}^{(\lambda-1)}_{\ \ \ \ \beta}\mu_c\left({}^{(\lambda)}\mathbf{t_i}\right) \cdot \left[ {}^{(\lambda-1)}_{\ \ \ \ \beta}u_{c\alpha} \cdot {}^{(\lambda)}t_\alpha(\mathbf{x_i}) + {}^{(\lambda-1)}_{\ \ \ \ \beta}v_{c\alpha} \right] \cdot {}^{(\lambda-1)}_{\ \ \ \ \beta}u_{c\alpha} \right.$$
$$\left. \left. \cdot \left[ {}^{(\lambda-1)}_{\ \ \ \ \beta}r_c\left({}^{(\lambda)}\mathbf{t_i}\right) - {}^{(\lambda-1)}t_\beta(\mathbf{x_i}) \right] + {}^{(\lambda-1)}_{\ \ \ \ \beta}m_{c\alpha} \cdot {}^{(\lambda-1)}_{\ \ \ \ \beta}\varpi_k\left({}^{(\lambda)}\mathbf{t_i}\right) \right\} \right\}$$

10       **end**

11       $H_\lambda \leftarrow \text{len}(\textbf{layer } \lambda)$

12
$$\overset{(\lambda)}{\underset{\alpha}{\Gamma}}(\mathbf{x_i}) \leftarrow \frac{1}{H_{\lambda-1}} \sum_{\beta=1}^{H_{\lambda-1}} \left\{ \left[ {}^{(\lambda)(\lambda-1)}_{\ \alpha\quad\ \beta} \right] \phi(\mathbf{x_i}) \cdot \overset{(\lambda-1)}{\underset{\beta}{\Gamma}}(\mathbf{x_i}) \right\}$$

13       **for** each **cluster** $k$ **in** unit $\alpha$ **do**

14
$$ {}^{(\lambda)}_{\alpha}\xi_k(\mathbf{x_i}) \leftarrow {}^{(\lambda)}_{\alpha}\mu_k(\mathbf{x_i}) \cdot \left[ \sum_{c=1}^{C} {}^{(\lambda)}_{\alpha}\mu_c(\mathbf{x_i}) \right]^{-1} $$

15
$$\frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}n_k} \leftarrow {}^{(\lambda)}_{\alpha}\xi_k(\mathbf{x_i}) \cdot {}^{(\lambda)}_{\alpha}\varpi_k(\mathbf{x_i})$$

16       Update the entry of the parameter $n$ in $\mathbf{W}$:

17       $\mathbf{W}[\lambda, \alpha]_k\{n\} \leftarrow \dfrac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}n_k}$

18       Auxiliar variable 1:

19       $ {}^{(\lambda)}_{\alpha}\vartheta_k(\mathbf{x_i}) \leftarrow -2\,{}^{(\lambda)}_{\alpha}\xi_k(\mathbf{x_i}) \cdot {}^{(\lambda)}_{\alpha}\mu_k(\mathbf{x_i}) \cdot \left[ {}^{(\lambda)}_{\alpha}r_k(\mathbf{x_i}) - {}^{(\lambda)}t_\alpha(\mathbf{x_i}) \right]$

20       **for** each **dimension** $j$ **in** ${}^{(\lambda+1)}\mathbf{t}(\mathbf{x_i})$ **do**

21       Auxiliar variable 2: $ {}^{(\lambda+1)}_{\alpha}\varrho_k(\mathbf{x_i}) \leftarrow {}^{(\lambda)}_{\alpha}u_{kj} \cdot {}^{(\lambda+1)}t_j(\mathbf{x_i}) + {}^{(\lambda)}_{\alpha}v_{kj}$

22
$$\frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}m_{kj}} \leftarrow {}^{(\lambda+1)}t_j(\mathbf{x_i}) \frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}n_k}$$

23
$$\frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}v_{kj}} \leftarrow {}^{(\lambda+1)}_{\alpha}\varrho_k(\mathbf{x_i}) \cdot {}^{(\lambda)}_{\alpha}\vartheta_k(\mathbf{x_i})$$

24
$$\frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}u_{kj}} \leftarrow {}^{(\lambda+1)}t_j(\mathbf{x_i}) \frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}v_{kj}}$$

25       Update the entries in $\Delta\mathbf{W}$:

26
$$\Delta\mathbf{W}[\lambda, \alpha]_{kj}\{m, a, b\} \leftarrow \left\{ \frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}m_{kj}}, \frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}a_{kj}}, \frac{\partial\, {}^{(\lambda)}t_\alpha(\mathbf{x_i})}{\partial\, {}^{(\lambda)}_{\alpha}b_{kj}} \right\}$$

27       **end**

28       **end**

29       $\Delta\mathbf{W}[\lambda, \alpha] \leftarrow \overset{(\lambda)}{\underset{\alpha}{\Gamma}}(\mathbf{x_i}) \cdot \Delta\mathbf{W}[\lambda, \alpha]$

30     **end**

31    **end**

32    $\Delta\mathbf{W} \leftarrow \eta \cdot \Delta\mathbf{W}$

33    **return** $\Delta\mathbf{W}$

---

**Algorithm 3:** Update process to calculate the parameters of the next epoch. Time complexity $O(\Pi \times H_\lambda)$ where $H_\lambda = \max_{\pi \in \{0,1,\dots,\Pi\}} H_\pi$. Space complexity $O(1)$ if we modify the pre-existing structure $\mathbf{W}$.

---

**Inputs:** The update structure $\Delta\mathbf{W}$, the initial structure of parameters $\mathbf{W}^0$.
**Outputs:** The updated parameter final structure $\mathbf{W}^f$.

1  Initialize with the initial structure: $\mathbf{W}^f \leftarrow \mathbf{W}^0$
2  Order of the layers: $\Lambda \leftarrow \{\Pi, \Pi - 1, \dots, 1, 0\}$
3  **for** each **layer** $\lambda$ **in** $\Lambda$ **do**
4    **for** each **unit** $\alpha$ **in layer** $\lambda$ **do**
5      $\mathbf{W}^{new}[\lambda, \alpha]_{kj} \leftarrow \Delta\mathbf{W}[\lambda, \alpha]_{kj} + \mathbf{W}^{old}[\lambda, \alpha]_{kj}$
6    **end**
7  **end**
8  **return** $\mathbf{W}^{new}$

---

Both the explainability and the self-reorganization concepts of CEFYDRA will be covered in detail in a different paper. However, we make a note on these two so the reader can understand further the implications of this technique.

The ability to explain the predictions generated (i.e., explainability) is a direct consequence of using membership functions to cluster the information and later approximate with different expressions. At each unit of the CEFYDRA, the dominant individual functions with highest membership values can be collected and aggregated to generate a final linear approximation of the inputs and outputs, which can then be displayed to the end user as an explanation of the prediction. Note that this explanations refer to a specific instance.

The ability for self-reorganization of CEFYDRA comes from the fact that during the learning, the membership functions of some clusters might get reduced to the point where they become negligible. At this point, the network can be modified by suppressing such cluster and reassigning it in a different unit.

## 4    Conclusions

We proposed an alternative core mathematical mechanism for the units of a neural network. The CEFYDRA is an algorithm that leverages a fuzzy Takagi-Sugeno-Kang inference to map the inputs of each unit via Cauchy membership functions and multi-dimensional logistic functions. In this paper we focused our efforts on demonstrating the learning formulation for the deep hidden layers of the CEFYDRA. We covered the three laws that make this generalization possible, i.e., Displacement, Substitution and Composition. Finally, we also propose a method to minimize the calculations (induced from the fully connected system) and reduce the complexity in the training.

## Acknowledgements

## References

1. Yang, Z., Zhang, A., Sudjianto, A.: GAMI-Net: An explainable neural network based on generalized additive models with structured interactions. Pattern Recognition, 120 (2021).
2. Lee, Y. Extraction of Competitive Factors in a Competitor Analysis Using an Explainable Neural Network. Neural Processing Letters 53, 1979–1994 (2021).
3. Fauvel, K., Lin, T., Masson, V., Fromont, É., Termier, A.: XCM: An Explainable Convolutional Neural Network for Multivariate Time Series Classification. Mathematics 9 (23), (2021).
4. Sasaki, H., Hidaka, Y., Igarashi, H.: Explainable Deep Neural Network for Design of Electric Motors. IEEE Transactions on Magnetics 57(6), 1-4 (2021).
5. Valizadegan, H. et al.: ExoMiner: A Highly Accurate and Explainable Deep Learning Classifier that Validates 301 New Exoplanets. arXiv:2111.10009 (2021).
6. Kim, M. S., Yun, J. P., Park, P.: An Explainable Convolutional Neural Network for Fault Diagnosis in Linear Motion Guide. In IEEE Transactions on Industrial Informatics, 17 (6), 4036-4045 (2021).
7. Gallego, V., Ríos Insua, D.: Current Advances in Neural Networks. Annual Review of Statistics and Its Application 9(1), (2022).
8. Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., Müller, K. -R.: Explaining Deep Neural Networks and Beyond: A Review of Methods and Applications. Proceedings of the IEEE 109(3), 247-278 (2021).
9. Viaña, J., Ralescu, S., Cohen, K., Ralescu, A., Kreinovich, V: Localized Learning A Possible Alternative to Current Deep Learning Techniques. In: Melin, P., Castillo, O. (eds.) Studies in Computational Intelligence (2021).
10. Viaña, J., Cohen, K.: Fuzzy-Based, Noise-Resilient, Explainable Algorithm for Regression. In: Explainable AI and Other Applications of Fuzzy Techniques. 1st edn. Springer, Cham (2022).
11. Viaña, J., Ralescu, S., Cohen, K., Ralescu, A., Kreinovich, V.: Extension to multi-dimensional problems of a fuzzy-based explainable and noise-resilient algorithm". In: Proceedings of the 14th International Workshop on Constraint Programming and Decision Making CoProd'2021, Szeged, Hungary (2021).
12. Viaña, J., Ralescu, S., Cohen, K., Ralescu, A., Kreinovich, V.: Why Cauchy Membership Functions: Reliability. Advances in Artificial Intelligence and Machine Learning, (To Appear).
13. Viaña, J., Ralescu, S., Cohen, K., Ralescu, A., Kreinovich, V.: Why Cauchy Membership Functions: Efficiency. Advances in Artificial Intelligence and Machine Learning, 1(1), 81-88 (2021).