# Continuation-Based Learning Algorithm for Discrete-Time Cellular Neural Networks

Holger Magnussen, Georgios Papoutsis and Josef A. Nossek

Institute for Network Theory and Circuit Design,
Technical University Munich, Munich, Germany
e-mail: holgi@nws.e-technik.tu-muenchen.de

**Abstract** – *The SGN-type nonlinearity of a standard Discrete-Time Cellular Neural Network (DTCNN) is replaced by a smooth, sigmoidal nonlinearity with variable gain. Therefore, the resulting dynamical system is fully differentiable. Bounds on gain of the sigmoidal function are given, so that the new, smooth system approximates the standard DTCNN within certain limits. A learning algorithm is proposed, which finds the template parameters for the standard DTCNN by gradually increasing the gain of the sigmoidal function.*

## 1 Introduction

The Discrete-Time Cellular Neural Network (DTCNN) was introduced in [1] as a discrete-time version of the Cellular Neural Network [2]. In a vector notation, the DTCNN equations can be written as follows

$$\begin{aligned} \mathbf{x}(t+1) &= \mathbf{A}\mathbf{y}(t) + \mathbf{B}\mathbf{u}(t) + \mathbf{i} \\ \mathbf{y}(t+1) &= \mathbf{s}(\mathbf{x}(t+1)), \quad \mathbf{y}(0) := \mathbf{y}_0 . \end{aligned} \tag{1}$$

Let $M_0$ be the number of (active and dummy) cells. $\mathbf{u}(t)$, $\mathbf{y}(t)$, and $\mathbf{x}(t)$ are the $M_0$-vectors of input signals, of cell output signals, and of cell states, respectively. Note that the dummy cells, which surround the array of active cells, are also included in the vectors. The $M_0 \times M_0$ matrices $\mathbf{A}$ and $\mathbf{B}$ contain the feedback and the feedforward template coefficients. Due to the local interconnection structure of the DTCNN and the translational invariance of the template coefficients, the two matrices are only sparse. $\mathbf{i}$ is the $M_0$-vector of (identical) cell bias values. The vector-valued function $\mathbf{s} : \mathbb{R}^{M_0} \to \{-1,1\}^{M_0}$ is equivalent to an element-wise application of the sign function $s : \mathbb{R} \to \{-1,1\}$ ($s(x) = +1$ for $x \geq 0$, and $-1$ otherwise) to each element of $\mathbf{x}(t)$ of the active cells and to a constant value for the dummy cells. Due to the SGN-type cell nonlinearity $s(x)$, the resulting dynamical system is not differentiable.

Since gradient information is not available, standard learning algorithms for differentiable dynamical systems [3], [4], [5] cannot be applied. Therefore, the SGN-type function $s(x)$ in Eq. 1 is replaced by a smooth, nonlinear, differentiable function $\hat{s}(\beta x)$ with

$\hat{s} : \mathbb{R} \rightarrow [-1, 1]$. $\hat{s}(\beta x)$ is a monotonous function with saturation regions, and it satisfies $\hat{s}(x) = -\hat{s}(-x)$, $\hat{s}(0) = 0$, and $\hat{s}(x \rightarrow \pm\infty) = \pm 1$. A typical choice is $\hat{s}(\beta x) := \tanh(\beta x)$. $\beta > 0$ is an adjustable gain factor. In the limit $\beta \rightarrow \infty$, the smooth function $\hat{s}(\beta x)$ is identical to SGN($x$) for all $x \neq 0$. For large, but finite $\beta$, the smooth function will approximate the SGN-function except for small values of $|x|$.

For the rest of this work, all cell states $x(t)$ and cell outputs $y(t)$ belonging to the system with the smooth nonlinearity will be distinguished from the standard DTCNN system by a hat "^". In analogy with Eq. 1, we define

$$
\begin{aligned}
\hat{x}(t+1) &= A\hat{y}(t) + Bu(t) + i \\
\hat{y}(t+1) &= \hat{s}(\beta\hat{x}(t+1)), \quad \hat{y}(0) := y(0) = y_0 .
\end{aligned}
\tag{2}
$$

This system is now fully differentiable with respect to the network parameters.

## 2 Approximation Properties

If in the original DTCNN from Eq. 1 all cell state values at all time-steps $t$ are nonzero, then the behavior of the original DTCNN can be approximated with arbitrary precision by the smooth system from Eq. 2:

**Lemma 1 (Approximation properties of the smooth system):** *Let $x_\mu(t) \neq 0$ for all $t$ and all cells $\mu$. Let $0 < \delta < 1$. $a_\mu$ is the $\mu$-th row of the matrix A. Then*

$$
\left| x_\mu(t+1) - \frac{\delta}{2} a_\mu y(t) \right| \geq \frac{1}{\beta} \hat{s}^{-1}(1 - \delta) + \frac{\delta}{2} \cdot \|a_\mu\|_1 \qquad \forall\, t \geq 0, \forall\, \mu
\tag{3}
$$

*is a sufficient condition for*

$$
\|y(t) - \hat{y}(t)\|_\infty \leq \delta \qquad \forall\, t \geq 0 .
\tag{4}
$$

$\diamond$

*Proof:* by induction; omitted for lack of space $\qquad\qquad\square$

It is shown in [6] that the parameter space of the standard DTCNN system Eq. 1 is segmented into a very large number of convex cones, in which the behavior of the system is constant. The condition from Eq. 3 adds *transition regions* around the hyperplanes separating the convex cones, where the system Eq. 2 does not necessarily approximate system Eq. 1. This is illustrated for the case of a two-dimensional parameter space in Fig. 1, where the transition regions are shown in grey. The system behavior in the transition regions can be very complicated.

Lemma 1 implies that in the limit for very small $\delta$ and very large $\beta$, the smooth system Eq. 2 approximates the original DTCNN system Eq. 1 almost everywhere in the parameter space except for points where $|x_\mu(t)| \ll 1$ for some $t, \mu$, i.e. around the hyperplanes separating the convex cones in the parameter space.
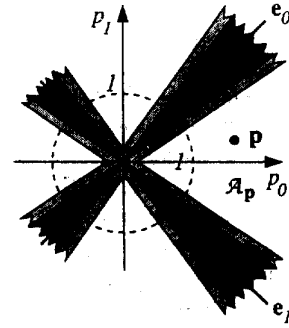


Fig. 1: Transition regions and approximation regions

172

Note that the SGN-type nonlinearity in a DTCNN hardware realisation is often re-
alised by a high gain amplifier [7]. Therefore, a hardware implementation is in fact a
smooth system of the form of Eq. 2. This implies that the template parameters of a
hardware implementation will have to obey the condition given in Lemma 1 in order to
work properly. This leads to even stricter accuracy conditions on the template parameters
than the ones derived in [6].

## 3  Learning by Continuation-Based Methods

Without loss of generality, we assume that the learning task consists of only one item. An
objective function is introduced, which measures the deviation of the actual behavior of
the network outputs $\hat{y}(t)$ from the desired behavior $\hat{d}(t) \in [-1,1]^{M_0}$ for $1 \leq t \leq T$. Let
p denote the *parameter vector* containing all independent template coefficients and the
bias. We then define the objective function $\hat{o}(p)$

$$\hat{o}(p) := \frac{1}{4M_0} \sum_{t=1}^{T} \alpha(t) \cdot \left\| \hat{y}(t) - \hat{d}(t) \right\|_2^2 \qquad \text{with } \alpha \in [0,1] , \sum_{t=1}^{T} \alpha(t) = 1 . \tag{5}$$

This objective function is differentiable with respect to the template parameters. Gradi-
ents can be computed using methods like the Williams/Zipser method [4] or the discrete-
time version of the method proposed by Pearlmutter [5]. Note that the network is only
operated for a fixed number of $T$ time-steps. The degree of freedom in systems Eq. 1 and
Eq. 2 ($\beta$ is counted as a system parameter as well) is eliminated by using a parameteri-
zation for p, which keeps the 2-norm of the parameter vector constant.

From Lemma 1, it can be concluded that a minimization of the objective function
Eq. 5 does not make sense for large values of $\beta$, because in this case, Eq. 5 is a good
approximation of an objective function of system Eq. 1 defined analogously to Eq. 5. In
the case of large $\beta$, the objective function Eq. 5 has large areas with very low gradient
values, and small areas, most notably the ones excluded by Eq. 3, where the gradients
can become very large. The minimisation of such an objective function is very difficult
from a numerical point of view.

Therefore, the mathematical tool of *Contin-
uation* is applied [8]. Let $S^\beta$ denote a smooth
system of the form Eq. 2 with a gain parameter
$\beta$ and a parameter vector p satisfying $\|p\|_2 = 1$.
Let $\hat{o}^\beta(p)$ denote the corresponding objective
function according to Eq. 2. A PIDGIN AL-
GOL [9] description of the Continuation method
is given in Fig. 2.

The function $\texttt{optimize}(\hat{o}^{\beta_j}(p), p_{ini})$ mini-
mises the objective function $\hat{o}^{\beta_j}(p)$ constrained

```
begin Continuation method;
  randomize(p_ini);
  for j = 0 to j_max do
    P_opt = optimize(o^{βj}(p), p_ini);
    P_ini = new_init(P_opt);
  end
end Continuation method;
```

Fig. 2: Continuation-based method

to ($\|p\|_2 = 1$) using the *ParTan* algorithm, which is an optimisation algorithm especially
suited for functions with "narrow valleys" [10]. This minimisation starts from the initial
point $p_{ini}$. It returns the parameter vector $p_{opt}$ corresponding to a (local) optimum.
$\beta_0, \ldots, \beta_{j_{max}}$ is a monotonously increasing series of positive gain factors.

The algorithm works as follows: It starts with a small value $\beta_0$ and a random starting
point $p_{ini}$. Let $p_{opt}$ be the result of the optimisation of $\hat{o}^{\beta_j}(p)$. The new starting point
for the optimisation of $\hat{o}^{\beta_{j+1}}(p)$ is computed by the function $\texttt{new\_init}(p_{opt})$. In most

173

cases, $p_{ini} = p_{opt}$ is a reasonable strategy, but more elaborated extrapolation methods can be used as well. Then, $\beta$ is increased from $\beta_j$ to $\beta_{j+1}$. This is repeated, until the smooth system $S^{\beta_{j_{max}}}$ sufficiently approximates the standard DTCNN system for some large $\beta_{j_{max}}$. The idea behind the Continuation method is that by increasing $\beta$ from $\beta_j$ to $\beta_{j+1}$, where $\beta_{j+1} - \beta_j$ should be reasonably small, the location of minimum $p_{opt}$ will not move drastically, so that the new optimum is reached in relatively few steps.

The Continuation method can be visualized by a simple example. A one-dimensional DTCNN consisting of 8 cells with only an a-template is set up. This network is supposed to learn the 1D Connected Component Detector task. The objective function of the form (5) is summed over all possible $L = 256$ input patterns. Thus the network has 3 independent parameters. Due to the ($\|p\|_2 = 1$)-restriction, the parameter vector can be reduced to two dimensions by using a transformation of the form

$$p_1 = \sin v_1; \quad p_2 = \cos v_1 \sin v_2; \quad p_3 = \cos v_1 \cos v_2 \tag{6}$$

with $v_1 \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ and $v_2 \in [-\pi, \pi]$. Note that this transformation is cyclical in $v_2$. The objective function $\hat{o}^\beta(v_1, v_2)$ is shown in Fig. 3 for different values of $\beta$. The objective function value is coded as the greyscale value. Dark regions correspond to low objective function values. Note that the objective function values are scaled differently in the different plots. It can clearly be seen, how the objective function develops from a smooth surface for small $\beta$ to the plateau-type surface for large $\beta$. The location of the minima for the corresponding value of $\beta$ changes slowly with $\beta$. Due to the transformation in Eq. 6, the shape of the convex cones in the parameter space is distorted.

## 4  Experimental Results

The Continuation-based algorithm described in Section 3 was first applied to easy learning problems with low-dimensional parameter vectors. In most of these cases the algorithm could easily find the expected global minimum of the described objective function $\hat{o}(p)$ without being very sensitive to the specific strategy of increasing the gain parameters $\beta_j$.

However, the performance of the method deteriorates for more complicated problems with parameter vectors of a higher dimensionality. This happens, because the optimization algorithm can get stuck in undesirable local minima of the objective function Eq. 5. These local minima lie in the transition regions excluded by Eq. 3, where at least one cell output $\hat{y}_\mu(t)$ can have very small values, so that the signal levels are no longer "quasi-binary". The local minima can be very narrow and very deep, especially for large values of $\beta$. They correspond to solutions that cannot be performed by the original DTCNN system from Eq. 1 due to its binary-valued cell outputs. Actually, it can be shown that it is possible to chose a binary-valued desired output $\hat{d}(t)$, so that the objective function Eq. 5 has a local minimum in a transition region [11].

As an example, Fig. 4 shows $\hat{o}(p)$ (solid line), the objective function of the smooth system according to Eq. 5, and a corresponding objective function $o(p)$ of the standard DTCNN system from Eq. 1 (dotted line). Both functions are shown for parameter vector values taken from a short line segment through a transition region in the parameter space, symbolically denoted by $p$. Here, $\beta$ is very large, and the local minimum, which only exists for the smooth system, is clearly visible.

In the case of more complicated problems, the choice of the strategy for increasing $\beta$ becomes very important. By changing $\beta$, the shape of the objective function and the
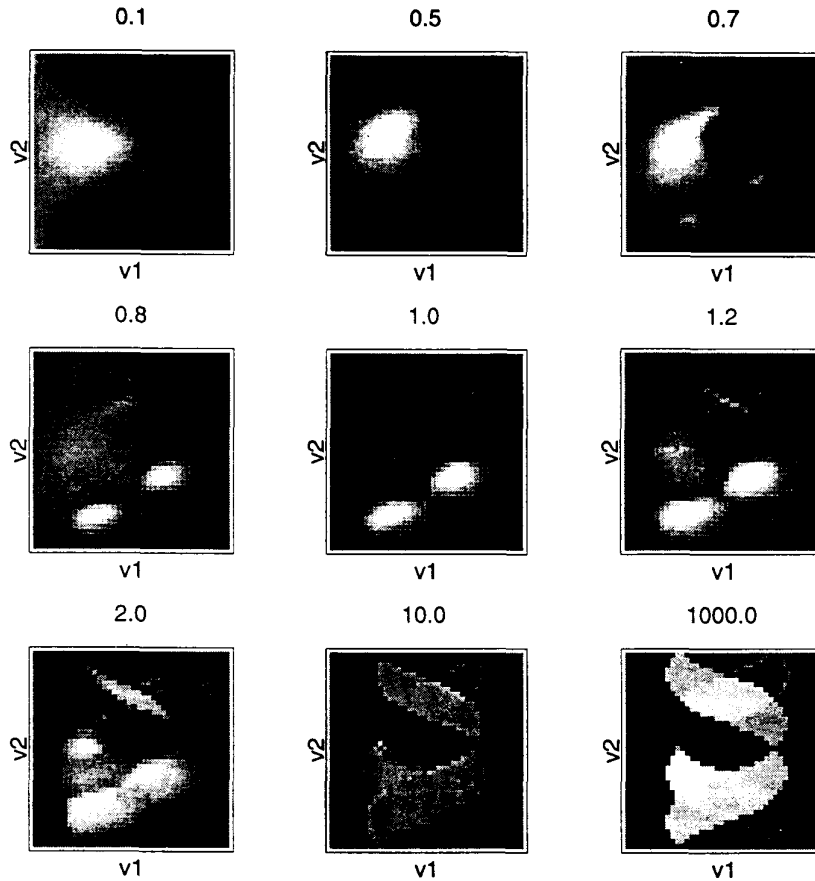
174

Figure 3: Objective function $\partial^\beta(v_1, v_2)$ for the smooth system and the CCD task

location of the minima changes. If $\beta$ is increased too fast, then the starting point for the optimization of system $\mathcal{S}^{\beta_j}$, i.e. the result of the optimization of system $\mathcal{S}^{\beta_{j-1}}$, is too far away from the actual location of the minimum of system $\mathcal{S}^{\beta_j}$. On the other hand, if $\beta$ is increased too slowly, then computation time is wasted, because the function optimize() has to be called too many times. In addition, it becomes more likely that the algorithm tracks one of the narrow, undesired minima in the transition regions until $\beta_{max}$ is reached. This solution is worthless, since it does not reflect the behavior of the standard DTCNN system. For moderate increases of $\beta$, the algorithm can "jump out" of the undesired local minima while at the same time it can track the much broader minima corresponding to a mapping performable by the standard DTCNN system.
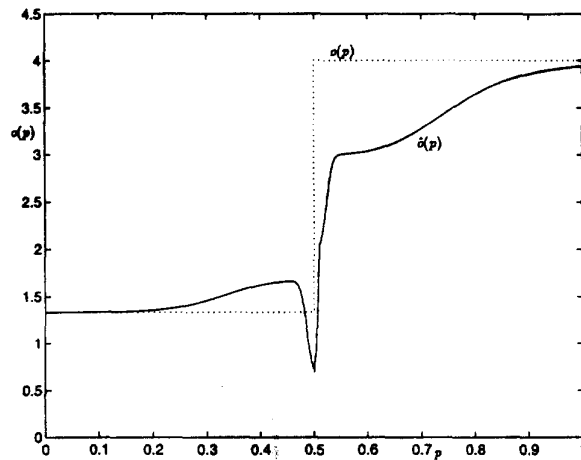
175

Figure 4: Objective function $\hat{o}(\mathbf{p})$ and $o(\mathbf{p})$ in a transition region

## References

[1] H. Harrer and J. A. Nossek, "Discrete-Time Cellular Neural Networks," *International Journal of Circuit Theory and Applications*, vol. 20, pp. 453–467, Sept. 1992.

[2] L. O. Chua and L. Yang, "Cellular neural networks: Theory," *IEEE Trans. on Circuits and Systems*, vol. 35, pp. 1257–1272, Oct. 1988.

[3] F. J. Pineda, "Dynamics and architecture for neural computation," in *Artificial Neural Networks: Paradigms, Applications and Hardware Implementation* (E. Sánchez-Sinencio and C. Lau, eds.), pp. 58–81, IEEE press, 1992.

[4] R. J. Williams and D. Zipser, "A learning algorithm for continually running fully recurrent neural networks," *Neural Computation*, vol. 1, pp. 270–280, 1989.

[5] B. A. Pearlmutter, "Learning state space trajectories in recurrent neural networks," *Neural Computation*, vol. 1, pp. 263–269, 1989.

[6] H. Magnussen and J. A. Nossek, "A geometric approach to properties of the discrete-time cellular neural network *(to be published)*," *IEEE Transactions on Circuits and Systems — I: Fundamental Theory and Applications*, 1994.

[7] H. Harrer, J.A. Nossek and R. Stelzl "An analog implementation of Discrete-Time Cellular Neural Networks," *IEEE Transactions on Neural Networks*, vol. 3, pp. 466–476, May 1992.

[8] R. Seydel, "Tutorial on continuation," *Int. J. of Bifurcation and Chaos*, vol. 1, no. 1, pp. 3–11, 1991.

[9] C. H. Papadimitriou and K. Steiglitz, *Combinatorial Optimization: Algorithms and Complexity*. Prentice-Hall, 1. ed., 1982.

[10] D. G. Luenberger, *Introduction to Linear and Nonlinear Programming*. Addison-Wesley, 1. ed., 1973.

[11] G. Papoutsis, "Implementierung eines Optimierungsalgorithmus als Lernalgorithmus für Zeitdiskrete Zellulare Neuronale Netze *(in german)*." Diplomarbeit, Dec. 1993.

176