



Autoassociative learning in relaxation labeling networks

Marcello Pelillo^{a,*}, Anna M. Fanelli^{b,1}

^a Dipartimento di Matematica Applicata e Informatica, Università "Ca' Foscari" di Venezia, Via Torino 155, 30173 Venezia Mestre, Italy

^b Dipartimento di Informatica, Università di Bari, Via Orabona 4, 70126 Bari, Italy

Received 20 June 1996

Abstract

We address the problem of training relaxation labeling processes, a popular class of parallel iterative procedures widely employed in pattern recognition and computer vision. The approach discussed here is entirely based on a theory of consistency developed by Hummel and Zucker, and contrasts with a recently introduced learning strategy which can be regarded as *heteroassociative*, i.e., what is actually learned is the association between patterns rather than the patterns themselves. The proposed learning model is instead *autoassociative* and involves making a set of training patterns *consistent*, in the sense rigorously defined by Hummel and Zucker; this implies that they become local attractors of the relaxation labeling dynamical system. The learning problem is formulated in terms of solving a system of linear inequalities, and a straightforward iterative algorithm is presented to accomplish this. The attractive feature of this algorithm is that it solves the system when it admits a solution, and automatically yields the best approximation solution when this is not the case. The learning model described here allows one to view the relaxation labeling process as a kind of asymmetric associative memory, the effectiveness of which is demonstrated experimentally. © 1997 Elsevier Science B.V.

Keywords: Learning; Relaxation labeling; Associative memories; Consistency; Neural networks;

1. Introduction

Relaxation labeling processes are a broad class of parallel cooperative procedures widely employed in computer vision and pattern recognition. They were pioneered by Rosenfeld, Hummel and Zucker (1976) in the mid-1970s and have since found applications in a variety of different problems (see (Kittler and Illingworth, 1985) for a review). The basic idea behind relaxation labeling is that complex, global computations can be done in a simple cooperative network of locally-interacting processing units, each representing

a specific hypothesis. This is commonly believed to be the style of computation actually implemented in the early stages of biological visual systems (e.g. (Ballard et al., 1983)). Due to the timing restrictions imposed by neuronal response, Marr was skeptical about this hypothesis. He raised the question that relaxation procedures "take too long and demand too much neural hardware to be implemented in any direct way" (Marr, 1982, p. 107). Recent experimental studies, however, contradict Marr's view, by showing that very few relaxation iterations are typically needed to arrive at a global "consistent" interpretation, provided that the initial processing stages are accurate enough (Zucker et al., 1988; Iverson, 1994).

One of the most important problems that arises in

* Corresponding author. E-mail: pelillo@dsi.unive.it.

¹ E-mail: fanelli@gauss.uniba.it.

applying a relaxation labeling algorithm to a practical task concerns finding a suitable set of *compatibility coefficients*. These are real-valued quantities that express the degree of agreement between different hypotheses, and effectively embody all the knowledge about the problem at hand. Traditional interpretations of compatibility coefficients have been in terms of statistical measures such as, for example, correlation (Rosenfeld et al., 1976) or mutual information (Peleg and Rosenfeld, 1978). The work by Peleg (1980) and Kittler and Hancock (1989) also supports this view. They developed new relaxation labeling algorithms, and provided probabilistic interpretations for the standard Rosenfeld et al.'s scheme, on the grounds of simple Bayesian arguments. The approach naturally leads to a statistical-based choice of the compatibility model. As pointed out by Hummel and Zucker (1983), even if Bayesian analysis provides much insight into the understanding of relaxation labeling algorithms, the approach is capable of accounting for at most one iteration of the process, and to understand its iterative behavior one has necessarily to resort to some approximation. Recently, however, the behavior of these relaxation labeling procedures has begun to be clarified, thanks to the work of Stoddart et al. (1995) who have uncovered certain interesting dynamical properties of the Kittler–Hancock probabilistic scheme.

In a recent work, we have tackled the problem from a radically different perspective (Pelillo and Refice, 1994). Instead of deriving the compatibilities of relaxation labeling by making recourse to probabilistic analysis, we try to *learn* them from a set of labeled data, in exactly the same way that neural networks do. Besides the obvious benefits of employing an automated procedure which is capable of determining the compatibility model in an “optimal” manner, the learning-based approach contributes to make the fields of relaxation labeling and neural networks closer. It also strengthens the biological plausibility of relaxation labeling processes as a mechanism of visual computation actually implemented in the brain, a hypothesis that has been forcefully put forward by Zucker and his collaborators (Zucker et al., 1989).

The approach described in (Pelillo and Refice, 1994) consists of deriving the compatibility strengths in such a way that the “distance” between what the relaxation process produces upon presentation of an input pattern and what is expected to produce

be as small as possible. Technically, this amounts to minimizing a certain cost function in the space of compatibility coefficients, a problem that can be solved by either classical gradient methods or more attractive global optimization procedures like genetic algorithms, as described in (Pelillo et al., 1995). The input to the learning algorithm in this case is a set of training exemplars in the form of ordered pairs of objects and labels. The objects are given as input to a local measurement process whose output will then initialize the relaxation process, while the labels represent the desired final responses. Borrowing the terminology from the neural network domain (Hertz et al., 1991), this can be regarded as a *heteroassociative* approach to learning, which means that what is actually learned (or stored) is the association between pairs of patterns, not the patterns themselves. A potential drawback of this kind of approach is that the derived compatibility coefficients *do* depend on the initial local measurement process. In some cases this may be advantageous, because in this way we effectively provide the learning algorithm with more problem-dependent information. However, in some circumstances it may result in poor generalization performance. This may be true when the initial local process is particularly noisy, thereby introducing noise in the learning process.

In this paper, we develop what may be called an *autoassociative* approach to the relaxation labeling learning task which, as a by-product, avoids this kind of problems. In this case, the learning problem consists of making a set of training patterns *consistent*, in the rigorous sense defined by Hummel and Zucker (1983). This is equivalent to saying that the learning data, viewed as points in state space, become attractive fixed points for the relaxation labeling dynamical system. We shall see that this amounts to solving a (sparse) system of linear inequalities, and shall describe an interesting iterative procedure to accomplish this, which has the distinguished feature of being able to always produce an “optimal” solution (in the sense of Chebyshev) when no exact solutions exist. It will also be seen how the proposed learning approach allows us to develop a powerful asymmetric associative memory model based on the dynamical properties of the relaxation labeling process. The validity of the model is confirmed experimentally.

We mention that essentially the same idea to de-

rive the compatibility coefficients of relaxation labeling was succinctly suggested by Hummel (1983). As a matter of fact, he did not suggest any practical way to accomplish this, and it seems that he did not recognize the potential applications of the model.

The outline of the paper is as follows. In Section 2, we introduce relaxation labeling algorithms and some of their fundamental properties. In Section 3, we formulate the learning problem as a system of linear inequalities, and in the subsequent section we describe an algorithm for solving it. Section 5 discusses the associative-memory application and presents some experimental results. In Section 6, we draw our conclusions.

2. Relaxation labeling processes and their properties

Relaxation labeling processes were developed to solve the so-called (continuous) labeling problem, where one has to assign labels to objects so as to satisfy a set of domain-specific constraints. Let $\mathbf{B} = \{b_1, \dots, b_n\}$ and $\mathbf{A} = \{\lambda_1, \dots, \lambda_m\}$ denote respectively the set of objects and the set of labels of the problem at hand. Moreover, let the constraints be quantitatively expressed in terms of a four-dimensional matrix of real-valued compatibility coefficients $R = \{r_{ij}(\lambda, \mu)\}$: the component $r_{ij}(\lambda, \mu)$ measures the strength of compatibility between the two hypotheses “ λ is on object b_i ” and “ μ is on object b_j ”. High values mean compatibility while low values mean incompatibility.

Let $p_i(\lambda)$ represent the degree of confidence of the hypothesis “label λ is on object b_i ”. It is assumed that $p_i(\lambda) \geq 0$ and $\sum_{\lambda} p_i(\lambda) = 1$, so that the m -dimensional vector $\bar{p}_i = (p_i(\lambda_1), \dots, p_i(\lambda_m))^T$ can be considered as the probability distribution of labels for object b_i . By putting together the \bar{p}_i 's we obtain a weighted labeling assignment for the objects of \mathbf{B} that will be denoted by \bar{p} , and will be conveniently considered as an $n \times m$ matrix. We will find it useful to introduce the space of weighted labeling assignments:

$$\mathcal{K} = \left\{ \bar{p} \in \mathbb{R}^{nm} : p_i(\lambda) \geq 0, i = 1, \dots, n, \lambda \in \mathbf{A} \right. \\ \left. \text{and } \sum_{\lambda} p_i(\lambda) = 1, i = 1, \dots, n \right\},$$

which is a linear convex set of \mathbb{R}^{nm} . Every vertex of \mathcal{K} represents an *unambiguous* labeling assignment which assigns exactly one label to each object. The set of these labelings will be denoted by \mathcal{K}^* :

$$\mathcal{K}^* = \{ \bar{p} \in \mathcal{K} : p_i(\lambda) = 0 \text{ or } 1, i = 1, \dots, n, \lambda \in \mathbf{A} \}.$$

Hummel and Zucker (1983) developed a general theory of consistency for the labeling problem which is the basis of the work reported here. The entire development of the theory is basically a generalization of the notion of consistency for unambiguous labelings which is more easily understood. Consider a labeling $\bar{p} \in \mathcal{K}$. The degree of agreement between the hypothesis that b_i is labeled with label λ and the context can be quantified by a linear *support* function

$$q_i(\lambda; \bar{p}) = \sum_j \sum_{\mu} r_{ij}(\lambda, \mu) p_j(\mu). \quad (1)$$

Now, let $\bar{p} \in \mathcal{K}^*$ be an unambiguous labeling assignment, and let $\lambda(i)$ denote the label assigned to b_i by \bar{p} (i.e., $p_i(\lambda(i)) = 1$). It seems reasonable to say that \bar{p} is consistent if and only if the assigned label of each object receives the greatest support at that object. This corresponds to having

$$q_i(\lambda; \bar{p}) \leq q_i(\lambda(i); \bar{p}), \quad i = 1, \dots, n, \lambda \in \mathbf{A}, \quad (2)$$

or, equivalently, $\sum_{\lambda} v_i(\lambda) q_i(\lambda; \bar{p}) \leq \sum_{\lambda} p_i(\lambda) q_i(\lambda; \bar{p})$, for all $\bar{v} \in \mathcal{K}^*$.

By analogy, a weighted labeling assignment $\bar{p} \in \mathcal{K}$ is said to be *consistent* provided that

$$\sum_{\lambda} v_i(\lambda) q_i(\lambda; \bar{p}) \leq \sum_{\lambda} p_i(\lambda) q_i(\lambda; \bar{p}), \\ i = 1, \dots, n, \quad (3)$$

for all $\bar{v} \in \mathcal{K}$. Furthermore, if strict inequalities hold in (3), for all $\bar{v} \neq \bar{p}$, then \bar{p} is said to be *strictly consistent*. It can be easily shown that for unambiguous labelings conditions (2) and (3) are equivalent (Hummel and Zucker, 1983). After defining the notion of consistency and proving some useful characterizations, Hummel and Zucker showed that when the compatibility matrix happens to be symmetric (i.e., $r_{ij}(\lambda, \mu) = r_{ji}(\mu, \lambda)$), then a sufficient condition for a labeling \bar{p} to be consistent is that it is a local minimum of the following “energy” function which is a measure of labeling’s (in)consistency

$$A(\bar{p}) = - \sum_{i,\lambda} \sum_{j,\mu} r_{ij}(\lambda, \mu) p_i(\lambda) p_j(\mu). \quad (4)$$

A relaxation labeling process takes as input an initial labeling assignment $\bar{p}^{(0)} \in \mathcal{K}$ and iteratively adjusts it taking into account the compatibility model, using an updating formula of the form

$$\bar{p}^{(t+1)} = f(\bar{p}^{(t)}, \bar{q}^{(t)}), \quad (5)$$

where t is the time index. The process evolves until (at least in theory) a fixed point is reached, which means that $\bar{p}^{(t+1)} = \bar{p}^{(t)}$. In practice, it is customary to stop the process when the distance between two successive labelings becomes negligible, or after a predetermined number of iterations.

The most popular form for the function f , which was also used in the experiments reported later in this work, is as follows:

$$p_i^{(t+1)}(\lambda) = \frac{p_i^{(t)}(\lambda) q_i^{(t)}(\lambda)}{\sum_{\mu} p_i^{(t)}(\mu) q_i^{(t)}(\mu)}, \quad (6)$$

provided that the compatibility coefficients are non-negative. This corresponds to the original nonlinear scheme developed heuristically by Rosenfeld et al. (1976). In the following, the relaxation algorithm will be best viewed as a continuous mapping \mathcal{T} of the assignment space \mathcal{K} onto itself. It starts out with $\bar{p}^{(0)}$ and iteratively produces a sequence of points $\bar{p}^{(0)}, \bar{p}^{(1)}, \bar{p}^{(2)}, \dots \in \mathcal{K}$, where $\bar{p}^{(t+1)} = \mathcal{T}(\bar{p}^{(t)})$, $t \geq 0$.

Recently it has been shown that despite its completely heuristic derivation, the original relaxation scheme (6) possesses a number of interesting properties (Pelillo, 1994). Firstly, when the compatibility matrix R is symmetric then A turns out to be a (strict) Liapunov function for the process, which means that it is monotonically decreasing along nonconstant trajectories. Secondly, and even more interestingly, it can be proven that strictly consistent labelings act as local attractors for the dynamical system defined by (6), whether or not the matrix R happens to be symmetric. This means that when started sufficiently close to a strictly consistent labeling \bar{e} , the relaxation process will tend to \bar{e} as time increases.² This property is formalized in the following theorem, which was

originally proven by Elfving and Eklundh (1982) in a slightly simplified form.

Theorem 1. *Let $\bar{e} \in \mathcal{K}^*$ be a strictly consistent labeling. Then \bar{e} is an asymptotically stable equilibrium point (and hence a local attractor) for the relaxation labeling scheme \mathcal{T} defined in formula (6).*

Proof. To prove the theorem we have to show that \bar{e} is a stable equilibrium point, and also a local attractor for \mathcal{T} . The latter condition was earlier proven in (Elfving and Eklundh, 1982, Theorem 10) by showing that the spectral radius of the Jacobian of \mathcal{T} evaluated at any strictly consistent labeling is less than 1. The fact that \bar{e} is a local attractor follows therefore immediately from a well-known result by Ostrowski (1966, Theorem 22.1) (we do not reproduce the proof here and refer to the original paper for technical details). It remains thus to see that \bar{e} is stable. Formally, this is expressed by the following condition: for any $\varepsilon > 0$ there exists a $\delta > 0$ such that $\|\bar{p} - \bar{e}\| < \delta$ implies $\|\mathcal{T}^t(\bar{p}) - \bar{e}\| < \varepsilon$ for all $t \geq 0$, where \mathcal{T}^t denotes the t th iterate of \mathcal{T} , i.e., $\mathcal{T}^0(\bar{p}) = \bar{p}$ and $\mathcal{T}^t(\bar{p}) = \mathcal{T}(\mathcal{T}^{t-1}(\bar{p}))$ for $t \geq 1$. First of all, since \bar{e} is a local attractor for \mathcal{T} , there must exist a $\delta' > 0$ such that $\lim_{t \rightarrow \infty} \mathcal{T}^t(\bar{p}) = \bar{e}$ whenever $\|\bar{p} - \bar{e}\| < \delta'$. Now, let $\varepsilon > 0$ be an arbitrary positive constant. There exists a nonnegative integer t_0 such that for all $t > t_0$ we have $\|\mathcal{T}^t(\bar{p}) - \bar{e}\| < \varepsilon$ whenever $\|\bar{p} - \bar{e}\| < \delta'$. Furthermore, since \mathcal{T} is continuous, so is \mathcal{T}^t for all $t \geq 0$. This means that, for any choice of t , there exists a $\delta_t > 0$ such that $\|\bar{p} - \bar{e}\| < \delta_t$ implies $\|\mathcal{T}^t(\bar{p}) - \mathcal{T}^t(\bar{e})\| = \|\mathcal{T}^t(\bar{p}) - \bar{e}\| < \varepsilon$ (recall that \bar{e} is a fixed point for \mathcal{T}). Therefore, by setting $\delta = \min\{\delta', \delta_1, \dots, \delta_{t_0}\}$ the condition for the stability of \bar{e} follows immediately, thereby proving the theorem. \square

This is the analog to the fundamental local convergence result that Hummel and Zucker proved for a different relaxation algorithm (which makes use of a computationally expensive projection operator) that turns out to be an approximation of the standard Rosenfeld et al.'s scheme (Hummel and Zucker, 1983, Theorem 9.1). Its interestingness stems essentially from the fact that no restriction on the structure of the compatibility matrix is imposed. This result is the basis for the autoassociative learning approach described in this paper.

² See, e.g., (Luenberger, 1979) for an introduction to dynamical systems.

3. Problem formulation

The learning approach described in this paper is based on the assumption that we have access to a set of learning patterns

$$\Xi = \{\xi^{(1)}, \dots, \xi^{(P)}\}, \quad (7)$$

where each $\xi^{(\gamma)}$, $\gamma = 1, \dots, P$, is simply an ordered list of labels, i.e.,

$$\xi^{(\gamma)} = \xi_1^{(\gamma)} \xi_2^{(\gamma)} \dots \xi_n^{(\gamma)}, \quad (8)$$

with $\xi_i^{(\gamma)} \in \mathcal{A}$, for all $i = 1, \dots, n$. For each learning pattern $\xi^{(\gamma)}$, consider the corresponding unambiguous labeling assignment $\bar{p}(\xi^{(\gamma)}) \in \mathcal{K}^* \subset \mathbb{R}^{nm}$, which is just an alternative (but equivalent) representation of $\xi^{(\gamma)}$:

$$p_i(\lambda; \xi^{(\gamma)}) = \begin{cases} 1, & \text{if } \lambda = \xi_i^{(\gamma)}, \\ 0, & \text{otherwise,} \end{cases} \quad (9)$$

all $i = 1, \dots, n$, $\lambda \in \mathcal{A}$.

The key idea behind the proposed learning approach follows from the observation that since the learning patterns *are* instances of the problem we intend to solve, they must be "consistent" with the problem's constraints. This naturally suggests the following strategy: find a compatibility matrix R so that the learning patterns become consistent (in the following, we will find it useful to "linearize" the compatibility matrix and consider it as an $n^2 m^2$ -dimensional column vector \bar{r}). Indeed we want the learning patterns to be strictly consistent. In fact, from the preceding section (Theorem 1) we know that if started in the vicinity of one of them, the relaxation process will eventually converge toward it.

For the labeling $\bar{p}(\xi^{(\gamma)})$, $\gamma = 1, \dots, P$, to be strictly consistent the following relation must hold:

$$q_i(\lambda; \bar{p}(\xi^{(\gamma)})) < q_i(\xi_i^{(\gamma)}; \bar{p}(\xi^{(\gamma)})) \quad (10)$$

for all $\lambda \neq \xi_i^{(\gamma)}$ and $i = 1, \dots, n$. Simple algebraic manipulations yield

$$\sum_j r_{ij}(\lambda, \xi_j^{(\gamma)}) - \sum_j r_{ij}(\xi_i^{(\gamma)}, \xi_j^{(\gamma)}) < 0, \quad (11)$$

which is a system of $Pn(m-1)$ linear inequalities in the $n^2 m^2$ unknowns $\{r_{ij}(\lambda, \mu)\}$.

System (11) can be compactly represented as

$$C\bar{r} < \bar{0}, \quad (12)$$

where C is the $Pn(m-1) \times n^2 m^2$ matrix defined as

$$C(\gamma, i, \lambda; j, k, \mu, \eta) = \begin{cases} +1, & \text{if } j = i, \mu = \lambda, \eta = \xi_k^{(\gamma)}, \\ -1, & \text{if } j = i, \mu = \xi_j^{(\gamma)}, \eta = \xi_k^{(\gamma)}, \\ 0, & \text{otherwise} \end{cases}$$

(for notational convenience we use a three-component index for the rows and a four-component index for the columns), \bar{r} is the unknown compatibility vector, and $\bar{0}$ is the null vector.

In practical applications, it is customary when solving systems of linear inequalities to introduce a "margin" (Duda and Hart, 1973), which also ensures larger basins of attraction (Forrest, 1988). Accordingly, our system is rewritten as

$$\sum_j r_{ij}(\lambda, \xi_j^{(\gamma)}) - \sum_j r_{ij}(\xi_i^{(\gamma)}, \xi_j^{(\gamma)}) \leq \kappa \quad (13)$$

or, equivalently,

$$C\bar{r} \leq \kappa \bar{1}, \quad (14)$$

where κ is some predetermined negative constant (the margin), and $\bar{1}$ is the unity vector.

4. Solving the system: the Eremin algorithm

Once that the autoassociative learning problem has been cast in terms of solving a system of linear inequalities, the next step involves choosing a particular algorithm to solve it. Many algorithms have been developed for solving systems of linear inequalities. The most immediate approach (also suggested in (Hummel, 1983)) is to make use of the well-known simplex algorithm. Another popular method is the *relaxation* method developed by Agmon (1954) (not to be confused with our relaxation labeling process). A common characteristic of these methods, however, is that they require the system to be *compatible*, i.e., to admit at least a solution. Alternative procedures, like the one developed by Ho and Kashyap (1965) are able to automatically determine the solution of a system if

there is one, and just to report that no solutions exist when this is not the case.

Even if in certain restricted circumstances our system can be proven to have a solution, in general this is certainly not the case. All these methods are therefore inappropriate in our case, for when the system is incompatible we would like at least to get a “good” approximation solution of it. Fortunately, this kind of algorithm is available. This was proposed in a rather obscure paper by the Russian mathematician Eremin in the early 1960s (Eremin, 1962). The attractive feature of Eremin’s method (which is a simple variant of Agmon’s relaxation algorithm) is that it *does* solve the system when it happens to be compatible, and yields the “best” approximation solution (in the sense of Chebyshev) when this is not the case. The major point is that it does not require any a priori knowledge about the system’s compatibility, being able to handle the two cases in a completely automatic manner.

Eremin’s method applies to systems of the general form

$$\sum_i a_{ji} x_i \leq b_j.$$

In our description, however, we will tailor it to our problem (13), or its equivalent formulation (14). Let $\bar{r} \in \mathbb{R}^{n^2 m^2}$ be an arbitrary compatibility vector, and put

$$\Delta(\bar{r}) = \max_{\gamma, i, \lambda} \left\{ \sum_j r_{ij}(\lambda, \xi_j^{(\gamma)}) - \sum_j r_{ij}(\xi_i^{(\gamma)}, \xi_j^{(\gamma)}) - \kappa \right\}. \quad (15)$$

Moreover, put

$$d(\bar{r}) = \begin{cases} \Delta(\bar{r}), & \text{if } \Delta(\bar{r}) \geq 0, \\ 0, & \text{otherwise,} \end{cases} \quad (16)$$

which is a continuous and convex function of \bar{r} . Eremin called the number $\varepsilon_0 = \min_{\bar{r}} d(\bar{r})$ the *defect* of system (13), which can be interpreted as the smallest ε for which the system

$$\sum_j r_{ij}(\lambda, \xi_j^{(\gamma)}) - \sum_j r_{ij}(\xi_i^{(\gamma)}, \xi_j^{(\gamma)}) - \kappa \leq \varepsilon \quad (17)$$

is compatible. Clearly if $\varepsilon_0 = 0$ then the original system (13) is compatible. When the system is incompat-

ible (i.e., $\varepsilon_0 > 0$) the number ε_0 is called the Chebyshev deviation of the system and any point \bar{r} satisfying the system

$$\sum_j r_{ij}(\lambda, \xi_j^{(\gamma)}) - \sum_j r_{ij}(\xi_i^{(\gamma)}, \xi_j^{(\gamma)}) - \kappa \leq \varepsilon_0 \quad (18)$$

will be called the point of its Chebyshev approximation. Furthermore, a sequence of points $\bar{r}_0, \bar{r}_1, \bar{r}_2, \dots \in \mathbb{R}^{n^2 m^2}$ will be called *solving* for the system (13) if it converges to some solution of the system (18).

Now, let $\{\lambda_k\}$ be a sequence of positive numbers such that $\lambda_k \rightarrow 0$ as $k \rightarrow \infty$, and $\sum_k \lambda_k = +\infty$ (e.g., $\lambda_k = 1/(k+1)$), and let \bar{r}_0 be an arbitrary initial point. Starting from \bar{r}_0 Eremin’s algorithm produces a sequence of points $\{\bar{r}_k\}$ according to the following scheme:

$$\bar{r}_{k+1} = \bar{r}_k - \lambda_{k+1} d(\bar{r}_k) C(\gamma_k, i_k, \lambda_k; \cdot)^T, \quad (19)$$

where $C(\gamma_k, i_k, \lambda_k; \cdot)$ is the row of the coefficient matrix C corresponding to the equation for which the value $d(\bar{r}_k)$ is attained (if more than one such equations exist, that with the least index is taken). Eremin (1962) proved the following result.

Theorem 2 (Eremin, 1962). *The sequence $\{\bar{r}_k\}$ defined in equation (19) is solving for the system (13).*

In other words, the preceding theorem states that if the system (13) happens to be compatible, then the sequence $\{\bar{r}_k\}$ will converge toward one of its solutions. Otherwise, it will converge toward a point of its Chebyshev approximation.

To conclude this section we need to clarify a final point. If we are to use the relaxation labeling algorithm defined in formula (6), then we need the compatibility coefficients to be nonnegative. Note, however, that there is no guarantee that the solution provided by the Eremin procedure will satisfy this constraint. A first possible approach to solve this problem is to add a set of “feasibility” constraints to our system (13). In our case, however, it is more convenient to initially derive arbitrary compatibility coefficients and then to scale them linearly so as to make them nonnegative (i.e., by simply adding the smallest negative coefficient). It can be readily seen that this has no effect on the structure of the space of consistent labelings (Pelillo, 1994).

5. An application: building an asymmetric associative memory

Following the seminal work of Hopfield (1982), there has been an increasing interest in the development of neural network models of associative memory (see e.g. (Hertz et al., 1991)). The basic idea behind this approach is that memory patterns can be stored as attractive fixed points of the system so that when started in their vicinity the memory will eventually “recall” the nearest one. Hopfield’s most valuable contribution was to show that certain highly-interconnected networks of neuron-like processing elements possess an energy function that drives their dynamical behavior toward low-energy states, provided that the connection weights between units are symmetric. Despite the manifested inspiration from neuroscience, however, the Hopfield model turns out to be unsatisfactory from a biological standpoint, because of the (essential) requirement that neurons be connected in a symmetrical fashion.

The autoassociative learning model discussed in this paper naturally leads us to view the relaxation labeling process as a novel kind of multi-valued asymmetric associative memory. This consists of an $n \times m$ densely interconnected relaxation labeling network, where n is the word’s length and m is the number of possible values at each site. The connection strengths between units are determined by the compatibility matrix R , i.e., $r_{ij}(\lambda, \mu)$ is the weight on the connection between units (i, λ) and (j, μ) . The unit indexed (i, λ) updates its own state according to formula (6), and its activation value $p_i(\lambda)$ can therefore be thought of as the probability that λ be the correct value for word’s site i . From Theorem 1, we know that if started in the vicinity of a strictly consistent labeling, the network will eventually converge toward it, whether or not the connection strengths are symmetric. If the resulting labeling corresponds to a memory pattern, we say that the memory has “recalled” that pattern. The learning algorithm described in the previous section can therefore naturally be employed as a means to “store” a given set of patterns in the memory.

5.1. Experimental results

To assess the validity of this memory model some experiments were conducted aimed at testing its error-

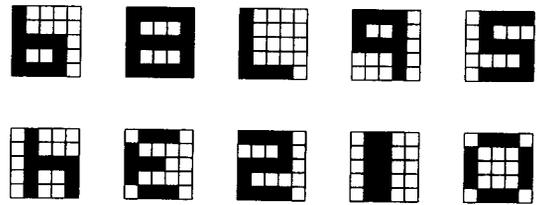


Fig. 1. Training set used in the experiments.

correction capabilities as well as evaluating its storage capacity. We concentrated on binary memories, i.e., $A = \{0, 1\}$. Note that, contrary to tradition, in our model each memory site has two computational units associated with it, namely one for the value “0” and the other for the “1”. It is interesting to observe that this kind of representation, which involves “opposing pairs”, is indeed the one preferred by many biological perceptual systems. Retinal ganglion cells with on- and off-center receptive fields are perhaps the most common example, but examples from other sensory systems abound (Hubel, 1988).

The learning set used in this study was taken from a recent paper by Zhuang et al. (1994) and consists of ten binary digits on a 5×5 matrix, as shown in Fig. 1. The training was carried out using the Eremin procedure described in the previous section, using different margin values $\kappa = -10, -20, -30, -50, -100$.

To experimentally test that the stored patterns had finite attraction basins, ten noisy versions of the corresponding unambiguous labelings were generated for each of them, using a continuous Gaussian noise with mean 0 and variance 0.1 (this was followed by a successive normalization step to ensure that the noisy labelings still belonged to \mathcal{K}). These 100 labelings were later given as input to the relaxation network which in the 100% of the cases quickly recalled the original ones.

Later, the error-correction capability of the model was tested. For each memory pattern, ten perturbed versions were generated by randomly flipping exactly d bits (for varying d), and the corresponding unambiguous labelings were obtained. In addition, since unambiguous labelings turn out to be fixed points for the relaxation scheme (6), a further Gaussian noise was inflicted as described before. The relaxation network was then allowed to run for (at most) 1,000 iterations, and the resulting (weighted) labelings were converted into binary patterns by a simple maxima selection criterion. The resulting binary patterns were

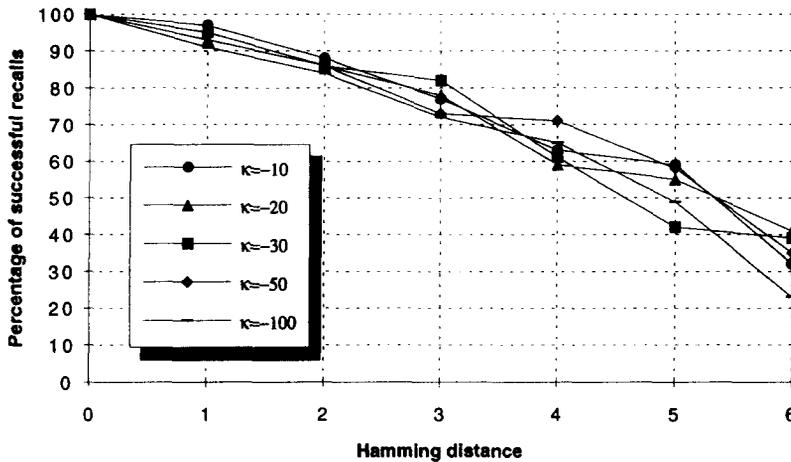


Fig. 2. Error-correction performance as a function of the Hamming distance, for various values of the margin κ .

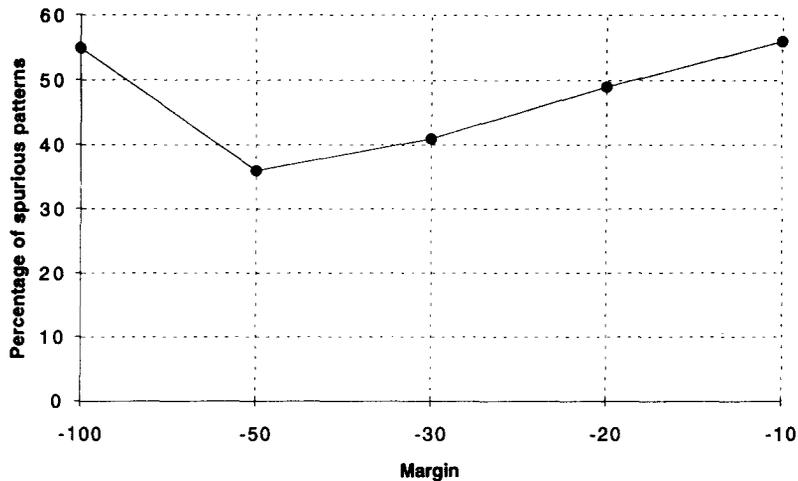


Fig. 3. Percentage of spurious patterns as a function of the margin κ .

then compared with the original ones, and a *success* was recorded when a perfect match was achieved. Fig. 2 shows the percentage of successful recalls as a function of d , for various values of the margin κ . As can be seen, there is no significant difference in performance, but it seems that there is no advantage in going beyond a certain limit which, in our case, may be taken as $\kappa = -50$. Our results compare favorably with those obtained on the same training set by Zhuang et al. (1994), who developed a sophisticated learning algorithm for Hopfield memories (cf. their Fig. 4(a)). Their results, in turn, were by far superior to those obtained using the standard Hebb rule originally introduced by Hopfield (1982).

We also experimentally estimated the number of spurious patterns in the proposed memory. We generated a hundred 25-bit random vectors and gave each of them as input to the relaxation labeling network. The converged patterns were then compared with the memory vectors shown in Fig. 1 and those not contained in the training set were considered as spurious. Fig. 3 shows the percentage of such patterns for the various margin values employed. Again, we note how decreasing the margin κ below -50 results in a deterioration of results. Unfortunately, we cannot compare in a fair manner our results with those presented in (Zhuang et al., 1994) since they considered a converged pattern as a spurious one if it differed

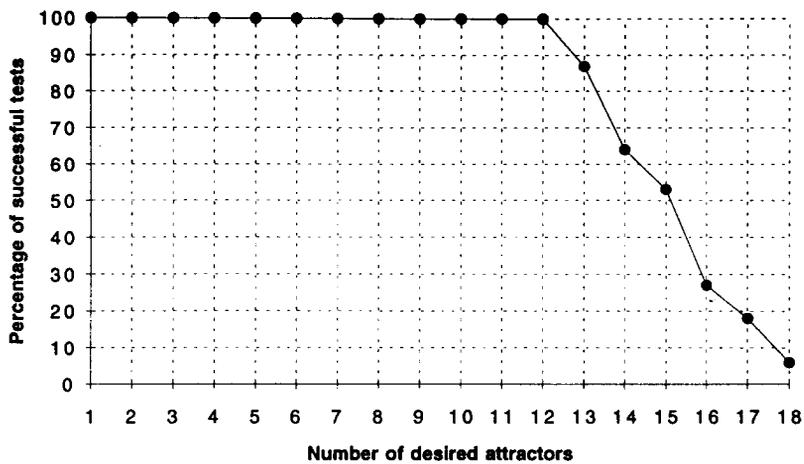


Fig. 4. Storage capacity of the relaxation labeling associative memory.

not only from the memory vectors $\xi^{(\gamma)}$, but also from $-\xi^{(\gamma)}$, for $\gamma = 1, \dots, P$ (they considered networks with $-1/+1$ states); storing a pattern ξ in a Hopfield network, in fact, implies storing also $-\xi$. They found about 20% spurious memories, a figure that is comparable to our best one (corresponding to $\kappa = -50$) if we take into account the criterion they employed to obtain it. Hebb-trained Hopfield networks produced instead a percentage of spurious memories very close to 100% (Zhuang et al., 1994).

Our next goal was to estimate the storage capacity of the proposed memory. To do so, the following experiment was conducted. We fixed the word's length at $n = 10$ and tried to store s desired memories, for increasing values of s . In this case, only the margin value $\kappa = -50$ was used. For each value of s , 100 such tests were performed each involving exactly s randomly generated patterns. After the learning phase, we tested if the desired memories were made strictly consistent by the learning algorithm. As seen, this is equivalent to saying that the patterns were stored into the memory. A test was considered successful if the percentage of stored memories was greater than or equal to 90%. We then calculated the percentage of successful tests for each value of s , and Fig. 4 shows the results obtained.

From this study we can say that the capacity of our model is approximately $1.2n$. This figure is by far superior to the capacity of Hebb-trained associative memories, which is $0.15n$ (Hopfield, 1982), and is also substantially larger than the estimated capacity

of the powerful model developed by Zhuang et al. (1994), which was $0.8n$ for $n = 10$.³ These results seem also to be superior to those obtainable with the asymmetric memory model described in (Michel and Farrell, 1990), which is capable of effectively storing a number of (linearly independent) patterns of the order of $0.5n$.

6. Conclusions

A crucial problem in applying a relaxation labeling algorithm to a practical task is to derive the so-called compatibility coefficients, which embody all the knowledge about the problem being solved. In analogy with neural network models, the idea of learning has recently been introduced into the relaxation labeling domain. This involves deriving the compatibility coefficients by means of an incremental algorithm, in such way that the performance of the relaxation process gradually improves over time. Borrowing the terminology from the neural network field, we have distinguished between heteroassociative and autoassociative learning procedures. In the first case, the task is to learn an association between two different patterns of stimuli. Our previous learning approach (Pelillo and Refice, 1994; Pelillo et al., 1995) falls into this class. One potential drawback of this approach, is that it is dependent on the initial local measurements, and this

³ In their experiments, however, they performed 1,000 tests for each value of s .

may potentially result in poor generalization performance.

In this paper, we have proposed a new, autoassociative algorithm to train relaxation labeling processes. The approach is based on a formal theory of consistency developed in a landmark paper by Hummel and Zucker (1983). After introducing the basic relaxation labeling formulas and discussing their fundamental dynamical properties, we have formulated the autoassociative learning problem as one of solving a system of linear inequalities. Among the many existing algorithms for solving such systems we have chosen a relaxation-style procedure developed by the Russian mathematician Eremin. This has the attractive feature of being able to automatically yield the best approximation solution when no exact solution exists. The proposed learning algorithm has been tested over an associative memory application, where memory patterns are stored as attractive fixed points of the relaxation labeling dynamical system. The results obtained indicate that our model exhibits error-correction performance that are competitive with sophisticated Hopfield-style memories, and turns out to have a substantially higher storage capacity.

References

- Agmon, S. (1954). The relaxation method for linear inequalities. *Canad. J. Math* 6, 382–392.
- Ballard, D.H., G.E. Hinton and T.J. Sejnowski (1983). Parallel visual computation. *Nature* 306, 21–26.
- Duda, R.O. and P.E. Hart (1973). *Pattern Classification and Scene Analysis*. Wiley, New York.
- Elfving, T. and J.-O. Eklundh (1982). Some properties of stochastic labeling procedures. *Comput. Graph. Image Process.* 20, 158–170.
- Eremin, I.I. (1962). Iteration method for Chebyshev approximations for sets of incompatible linear inequalities. *Soviet Math. Doklady* 3, 570–572.
- Forrest, B.M. (1988). Content addressability and learning in neural networks. *J. Phys. A: Math. Gen.* 21, 245–255.
- Hertz, J., A. Krogh and R.G. Palmer (1991). *Introduction to the Theory of Neural Computation*. Addison-Wesley, Redwood City, CA.
- Ho, Y.-C. and R.L. Kashyap (1965). An algorithm for linear inequalities and its applications. *IEEE Trans. Elect. Comp.* 14 (5), 683–688.
- Hopfield, J.J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci. USA* 79, 2554–2558.
- Hubel, D.H. (1988). *Eye, Brain, and Vision*. Scientific American Library, New York.
- Hummel, R.A. (1983). A design method for relaxation labeling applications. In: *Proc. Natl. Conf. Artificial Intell. (AAAI-83)*, Washington, DC, 1983, 168–171.
- Hummel, R.A. and S.W. Zucker (1983). On the foundations of relaxation labeling processes. *IEEE Trans. Pattern Anal. Machine Intell.* 5 (3), 267–287.
- Iverson, L.A. (1995). Toward discrete geometric models for early vision. PhD Thesis, McGill University, Montréal, Canada.
- Kittler, J. and E.R. Hancock (1989). Combining evidence in probabilistic relaxation. *Internat. J. Pattern Recognition Artificial Intell.* 3 (1), 29–51.
- Kittler, J. and J. Illingworth (1985). Relaxation labeling algorithms – A review. *Image Vision Comput.* 3, 206–216.
- Luenberger, D.G. (1979). *Introduction to Dynamic Systems*. Wiley, New York.
- Marr, D. (1982). *Vision*. Freeman, New York, 1982.
- Michel, A.M. and J.A. Farrell (1990). Associative memories via artificial neural networks. *IEEE Control Syst. Mag.*, April 1990, 6–17.
- Ostrowski, A.M. (1966). *Solution of Equations and Systems of Equations*. Academic Press, New York.
- Peleg, S. (1980). A new probabilistic relaxation scheme. *IEEE Trans. Pattern Anal. Machine Intell.* 2 (4), 362–369.
- Peleg, S. and A. Rosenfeld (1978). Determining compatibility coefficients for curve enhancement relaxation processes. *IEEE Trans. Syst. Man Cybernet.* 8 (7), 548–555.
- Pelillo, M. (1994). Nonlinear relaxation labeling as growth transformation. In: *Proc. 12th Internat. Conf. Pattern Recognition*, Jerusalem, Israel, 1994, 201–206.
- Pelillo, M. and M. Refice (1994). Learning compatibility coefficients for relaxation labeling processes. *IEEE Trans. Pattern Anal. Machine Intell.* 16 (9), 933–945.
- Pelillo, M., F. Abbattista and A. Maffione (1995). An evolutionary approach to training relaxation labeling processes. *Pattern Recognition Letters* 16 (10), 1069–1078.
- Rosenfeld, A., R.A. Hummel and S.W. Zucker (1976). Scene labeling by relaxation operations. *IEEE Trans. Syst. Man Cybernet.* 6 (6), 420–433.
- Stoddart, A.J., M. Petrou and J. Kittler (1995). Probabilistic relaxation as an optimizer. In: *Proc. British Machine Vision Conf.*, Birmingham, UK, 1995, 613–622.
- Zhuang, X., Y. Huang and F.A. Yu (1994). Design of Hopfield content-addressable memories. *IEEE Trans. Signal Process.* 42, 492–495.
- Zucker, S.W., C. David, A. Dobbins and L. Iverson (1988). The organization of curve detection: Coarse tangent fields and fine spline coverings. In: *Proc. 2nd Internat. Conf. Computer Vision*, Tampa, FL, 1988, 568–577.
- Zucker, S.W., A. Dobbins and L. Iverson (1989). Two stages of curve detection suggest two styles of visual computation. *Neural Computation* 1, 68–81.