

aiNet: An Artificial Immune Network for Data Analysis

Leandro Nunes de Castro & Fernando José Von Zuben

{lnunes,vonzuben}@dca.fee.unicamp.br

<http://www.dca.fee.unicamp.br/~lnunes>

<ftp://ftp.dca.fee.unicamp.br/pub/docs/vonzuben/lnunes/DMHA.pdf>

DRAFT

In Data Mining: A Heuristic Approach

Hussein A. Abbass, Ruhul A. Sarker, and Charles S. Newton (Eds.)

Idea Group Publishing, USA

March, 2001

SUMMARY

INTRODUCTION.....	1
BASIC IDEAS AND RATIONALE	1
IMMUNE PRINCIPLES	3
IMMUNE NETWORK THEORY	4
CLONAL SELECTION AND AFFINITY MATURATION	6
AINET: AN ARTIFICIAL IMMUNE NETWORK MODEL FOR DATA ANALYSIS.....	8
AINET CHARACTERIZATION AND ANALYSIS.....	12
RELATED IMMUNE NETWORK MODELS.....	13
ANALYSIS OF THE ALGORITHM	15
KNOWLEDGE EXTRACTION AND STRUCTURE OF A TRAINED AINET	16
HIERARCHICAL CLUSTERING AND GRAPH THEORETICAL TECHNIQUES	17
AINET FUZZY CLUSTERING	22
EMPIRICAL RESULTS	23
TWO-SPIRALS PROBLEM: SPIR	25
THE CHAINLINK PROBLEM: CHAINLINK	27
FIVE NON-LINEARLY SEPARABLE CLASSES: 5-NLSC	29
SENSITIVITY ANALYSIS.....	31
ABOUT THE NUMBER OF CLUSTERS AND NETWORK PARAMETERS	34
CONCLUDING REMARKS.....	35
ACKNOWLEDGEMENTS.....	37
REFERENCES.....	37

aiNet: An Artificial Immune Network for Data Analysis

Leandro Nunes de Castro & Fernando José Von Zuben

{lnunes,vonzuben}@dca.fee.unicamp.br

State University of Campinas - UNICAMP

Department of Computer Engineering and Industrial Automation (DCA)

School of Electrical and Computer Engineering (FEEC)

Campinas - SP – Brazil

Introduction

This chapter explores basic aspects of the vertebrate immune system and proposes a novel artificial immune network model with the main goals of clustering and filtering crude data sets described by high-dimensional samples. We are not intent on reproducing with confidence any immune phenomenon, but to demonstrate that immune concepts can be used as inspiration to develop novel computational tools for data analysis. As important results of our model, the network evolved will be capable of reducing redundancy and describing data structure, including their spatial distribution and cluster inter-relations. Clustering is useful in several exploratory pattern analyses, grouping, decision-making and machine-learning tasks including data mining, knowledge discovery, document retrieval, image segmentation and automatic pattern classification. The data clustering approach was implemented in association with hierarchical clustering and graph theoretical techniques, and the network performance is illustrated using several benchmark problems. The computational complexity of the algorithm and a detailed sensitivity analysis of the user-defined parameters are presented. A trade-off among the proposed model for data analysis, connectionist models (artificial neural networks) and evolutionary algorithms is also discussed.

Basic Ideas and Rationale

The vertebrate immune system has several useful theories from the viewpoint of information processing. Among these, we can stress the immune network theory and the clonal selection and affinity maturation principles. The immune network theory hypothesizes the activities of the immune cells, the emergence of memory and the discrimination between our own cells (known as self) and external invaders (known as nonself). It also suggests that the immune system has an internal image of all existing pathogens (infectious nonself) to

which it might be exposed during its lifetime. On the other hand, the clonal selection principle proposes a description of the way the immune system copes with the pathogens to mount an adaptive immune response. The affinity maturation principle is used to explain how the immune system becomes increasingly better at its task of recognizing and eliminating these pathogens (antigenic substances). In this chapter, we will review these theories and show that many of their concepts and ideas can be used to develop an artificial immune network model, named aiNet, capable of solving pattern recognition tasks similarly to the vertebrate immune system.

The aiNet model will consist of a set of cells, named antibodies, interconnected by links with associated connection strengths. The aiNet antibodies are supposed to represent the network internal images of the pathogens (input patterns) contained in the environment to which it is exposed. The connections between the antibodies will determine their interrelations, providing a degree of similarity (in a given metric space) among them: the closer the antibodies, the more similar they are.

Based upon a set of unlabeled patterns $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$, where each pattern (object, or sample) \mathbf{x}_i , $i = 1, \dots, M$ is described by L variables (attributes or characteristics), a network will be constructed to answer questions like: (1) Is there a great amount of redundancy within the data set and, if there is, how can we reduce it? (2) Is there any group or subgroup intrinsic to the data? (3) How many groups are there within the data set? (4) What is the structure or spatial distribution of these data (groups)? (5) How can we generate decision rules to classify novel samples?

This chapter is organized as follows. In Section 2, the basic immunological principles to be employed are reviewed. The artificial immune network model, named aiNet, is described in Section 3, and analyzed in Section 4. Section 5 presents the hierarchical clustering and graph theoretical techniques used to define the network structure, and Section 6 presents the aiNet simulation results for several benchmark tasks, comparing with the Kohonen (1982) self-organizing map (SOM). Section 7 presents a sensitivity analysis of the proposed algorithm with relation to the most critical user-defined parameters. The chapter is concluded in Section 8 with a discussion of the network main characteristics, potential applications and future trends.

Immune Principles

As a first step, we shall sketch a few aspects of the human adaptive immune system. A number of concepts and technical terms will be introduced to make the reader familiar with the terminology. Master details about the immune network theory, clonal selection and affinity maturation principles will be given in dedicated sections. An interested reader shall refer to Janeway, Travers, Walport and Capra (2000) for a good introductory text in immunology and to de Castro and Von Zuben (1999a) for immunology under the computational intelligence perspective.

The immune system is a complex of cells, molecules and organs with the primary role of limiting damage to the host organism by pathogens, which elicit an immune response and thus are called antigens (Ag). One type of response is the secretion of antibody (Ab) molecules by B cells, or B lymphocytes. Antibodies are Y-shaped receptor molecules bound on the surface of a B cell with the primary role of recognizing and binding, through a complementary match, with an antigen. The antibody molecules recognize a portion of the antigen called epitope. Antibodies also present epitopes, which are named idiotopes. A set of idiotopes is called an idio type. While each B cell is known to have a single type of antibody, thus being called monospecific, antigens typically have several different types of epitopes, and can be recognized by several different antibodies. The antibody portion responsible for matching (recognizing) an antigen is called paratope, also known as V-region, for variable region. It is variable because it can alter its shape to achieve a better match (complementarity) with a given antigen. The strength and specificity of the Ag-Ab interaction is measured by the affinity (complementarity level) of their match. See Figure 1 for an illustration of an antigen with its many epitopes and an antibody with its paratope and idiotope.



Figure 1: B cell, antigen, antibody, epitopes, paratopes and idiotopes. (a) An antigen with its multiple epitopes recognized by different B cells. (b) Antibody combining site (V-region or paratope), and its idiotope.

When stimulated, the B cell proliferates and secretes its receptor molecules as free antibodies. Antibodies thus can either be free molecules or receptors attached to cells. Secretion requires that B cells become activated, undergo proliferation (cloning) and then finally differentiate into plasma and memory cells. A clone is a cell, or a set of cells, which is the progeny of the same cell. A plasma cell is the one capable of secreting antibody that presents high rates and a memory cell is the cell with high affinity with the antigen that will be rescued for a faster and stronger response to a previously seen (or related) antigen. Those cells that are valuable to the system, i.e. recognize antigens, grow in concentration and affinity (affinity maturation), while those that are not die out. This basic process of pattern recognition and selection is known as clonal selection (Burnet, 1978) and is similar to natural selection, except that it occurs on a rapid time scale on the order of days and weeks, within our bodies.

In order to be protective, the immune system must learn to distinguish between our own (*self*) cells and malefic external (*nonself*) invaders. This process is called self/nonself discrimination: those cells recognized as self do not promote an immune response, the system is said to be tolerant to them, while those that are not provoke a reaction resulting in their elimination.

Immune Network Theory

The immune network theory, as originally proposed by Jerne (1974a), hypothesized a novel viewpoint of lymphocyte activities, natural antibody production, pre-immune repertoire selection, tolerance and self/nonself discrimination, memory and the evolution of the immune system. It was suggested that the immune system is composed of a regulated network of cells and molecules that recognize one another even in the absence of antigens. The immune system was formally defined as an enormous and complex network of paratopes that recognize sets of idiotopes, and of idiotopes that are recognized by sets of paratopes, thus it could recognize as well as be recognized. The relevant events in the immune system are not only the molecules, but also their interactions. The immune cells can respond either positively or negatively to the recognition signal (antigen or other immune cell or molecule). A positive response would result into cell proliferation, cell activation and antibody secretion, while a negative response would lead to tolerance and suppression (see Figure 2(a)).

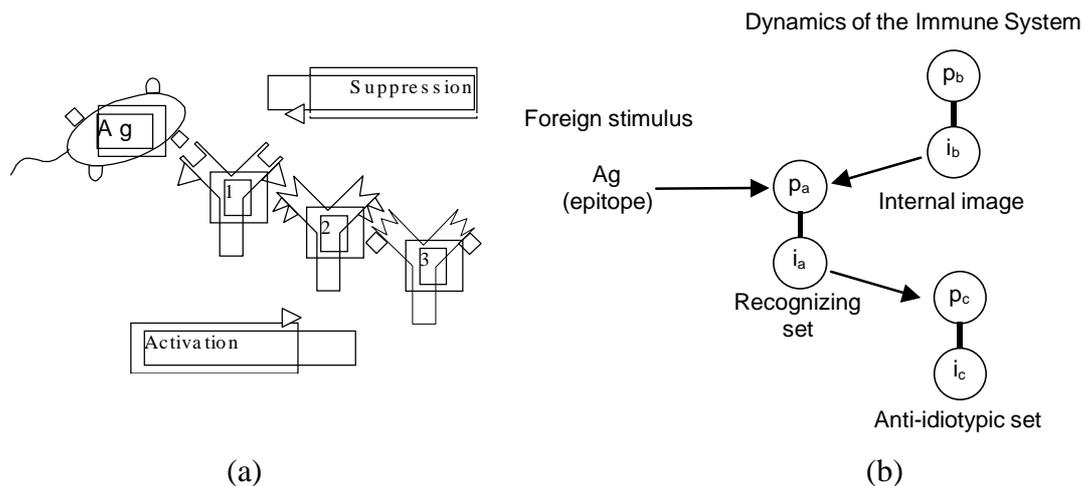


Figure 2: Idiotypic network representations. (a) An antigen stimulates the antibody production of class a, which stimulates class b, and so on. (b) More detailed view of idiotypic network (see text for details).

The network theory can be summarized as follows (see Figure 2(b)). When the immune system is primed with an antigen (Ag), its epitope is recognized (with various degrees of specificity) by a set of different paratopes, called p_a . The set i_b of idiotopes is called the internal image of the epitope (or antigen) because it is recognized by the same set p_a that recognized the antigen. The set i_b is associated with a set p_b of paratopes occurring on the molecules and cell receptors. Furthermore, each idiotope of the set i_a is recognized by a set of paratopes, so that the entire set i_a is recognized by an even larger set p_c of paratopes which occur together with a set i_c of idiotopes on antibodies and lymphocytes of the anti-idiotypic set. Following this scheme, we come to ever larger sets that recognize or are recognized by previously defined sets within the network. The arrows indicate a stimulatory effect when idiotopes are recognized by paratopes on cell receptors and a suppressive effect when paratopes recognize idiotopes on cell receptors.

There are several immune network models presented in the literature. Most of them are based upon a set of differential equations to describe the dynamics of the network cells and molecules (Jerne, 1974b; Bonna & Kohler, 1983; Farmer, Packard, & Perelson, 1986; Varela & Coutinho, 1991). The interactions, that could be either excitatory (network activation) or inhibitory (network suppression) between different types of elements would naturally lead to the network connectivity pattern and dynamics.

In the model proposed by Varela and Coutinho (1991), it is possible to stress three characteristics of the immune networks: 1) its structure, that describes the types of interaction among the network components, represented by matrices of connectivity; 2) its dynamics,

that accounts for the variation in time of the concentrations and affinities of its cells; and 3) its metadynamics, a property related to the continuous production of novel antibodies and death of non-stimulated or self-reactive cells. The central characteristic of the immune network theory is the definition of the individual's molecular identity (internal images), which emerges from a network organization followed by the learning of the molecular composition of the environment where the system develops.

The network approach is particularly interesting for the development of computational tools because it potentially provides a precise account of emergent properties such as learning and memory, self-tolerance, size control and diversity of cell populations. In general terms, the structure of most network models can be represented as (Perelson, 1989):

$$RPV = \begin{array}{c} \text{influx} \\ \text{of new cells} \end{array} - \begin{array}{c} \text{death of} \\ \text{unstimulated cells} \end{array} + \begin{array}{c} \text{reproduction of} \\ \text{stimulated cells} \end{array} \quad (1)$$

where RPV is the rate of population variation, and the last term includes Ab-Ab recognition and Ag-Ab stimulation.

Clonal Selection and Affinity Maturation

Learning in the immune system involves raising the population size and affinity of those lymphocytes that have proven themselves to be valuable by having recognized any antigen. During the lifetime of an individual, its immune system may be exposed to a given antigen repeatedly. The initial exposure to an antigen that stimulates an adaptive immune response is handled by a spectrum of small clones of B cells, each producing antibody with different affinity. The effectiveness of the immune response to secondary encounters is considerably enhanced by storing some high affinity antibody producing cells from the first infection (memory cells), so as to form a large initial high affinity clone for subsequent encounters (Ada & Nossal, 1987). This is an intrinsic scheme of a reinforcement learning strategy (Sutton & Barto, 1998), where the system is continuously improving its capability to perform its task (recognize antigens).

Antibodies present in a memory response have, on average, higher affinities than those of the early primary response. This phenomenon is referred to as the maturation of the immune response. This maturation requires that the antigen-binding sites of the antibody molecules in the matured response be structurally different from those present in the primary response. Random changes (mutations) are introduced into the variable region and occasionally one such change will lead to an increase in the affinity of the antibody. It is these high-affinity variants that are selected to enter the pool of memory cells. Those cells

carrying receptors with low antigenic affinity, or the self-reactive cells, must be efficiently eliminated (or become anergic). Instead of the expected clonal deletion of all self-reactive cells, B lymphocytes may undergo receptor editing: these B cells had deleted their self-reactive receptors and developed entirely new receptors. Some of the possible results of a meeting between a lymphocyte and an antigen are summarized in Figure 3.

George and Gray (1999) suggested that point mutations are good for exploring local regions, of an affinity landscape, while editing may rescue immune responses stuck on unsatisfactory local optima. A rapid accumulation of mutations (*hypermutation*) is necessary for a fast maturation of the immune response, but the majority of the changes will lead to poorer or non-functional antibodies. If a cell that has just picked up a useful mutation continues to be mutated at the same rate during the next immune responses, then the accumulation of deleterious changes may cause the loss of the advantageous mutation.

The selection mechanism may provide a means by which the regulation of the hypermutation process is made dependent on receptor affinity. Cells with low affinity receptors may be further mutated and eliminated if their affinity to the antigens remains small. However, in cells with high-affinity receptors, hypermutation may be gradually inactivated (Kepler & Perelson, 1993).

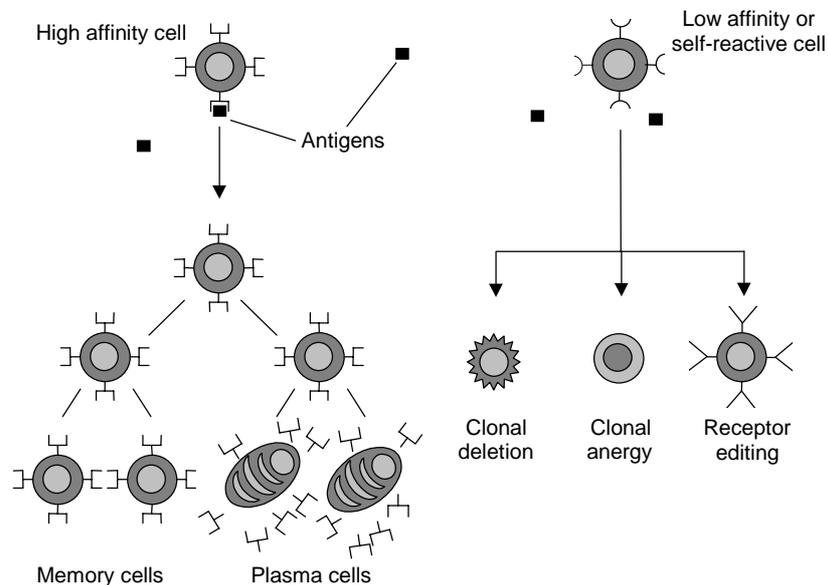


Figure 3: Antigenic interactions with lymphocytes. A minority of cells from the repertoire will recognize the antigen, and be activated by clonal selection. In cases the affinity between the cell receptor and the antigen is low or a self-reactive cell is detected, the lymphocyte in question might suffer apoptosis (programmed cell death), anergy (“freezing”), or a receptor editing process.

aiNet: An Artificial Immune Network Model for Data Analysis

In this section, we will present the aiNet learning algorithm focusing on its dynamics and metadynamics. A deeper analysis and several hierarchical clustering and graph theoretical techniques proposed to determine the final network architecture will be presented in the next section.

In order to quantify immune recognition, it is appropriate to consider all immune events as occurring in a shape-space S , which is a multi-dimensional metric space where each axis stands for a physico-chemical measure characterizing a molecular shape (Perelson & Oster, 1979). We will assume a problem dependent set of L measurements to characterize a molecular configuration as a point $s \in S$. Hence, a point in an L -dimensional space, called shape-space, specifies the set of features necessary to determine the antibody-antibody (Ab-Ab) and antigen-antibody (Ag-Ab) interactions. Mathematically, this shape (set of features that define either an antibody or an antigen) can be represented as an L -dimensional string, or vector. The possible interactions within the aiNet will be represented in the form of a connectivity graph. The proposed artificial immune network model can be formally defined:

Definition 1: The aiNet is an *edge-weighted graph*, not necessarily fully connected, composed of a set of nodes, called *antibodies*, and sets of node pairs called *edges* with an assigned number called *weight*, or *connection strength*, associated with each connected edge.

The aiNet clusters will serve as *internal images (mirrors)* responsible for mapping existing clusters in the data set into network clusters. As an illustration, suppose there is a data set composed of three regions with high density of data, according to Figure 4(a). A hypothetical network architecture generated by the learning algorithm to be presented is shown in Figure 4(b). The numbers within the cells indicate their labels (the total number is generally higher than the number of clusters and much smaller than the number of samples), the numbers next to the connections represent their strengths, and dashed lines suggest connections to be pruned, in order to detect clusters and define the final network structure. Notice the presence of three distinct clusters of antibodies, each of which with different number of antibodies, connections and strengths. These clusters map those of the original data set. Notice also that the number of antibodies in the network is much smaller than the number of data samples, characterizing an architecture suitable for data compression. Finally, the shape of the spatial distribution of antibodies follows the shape of the antigenic spatial distribution.

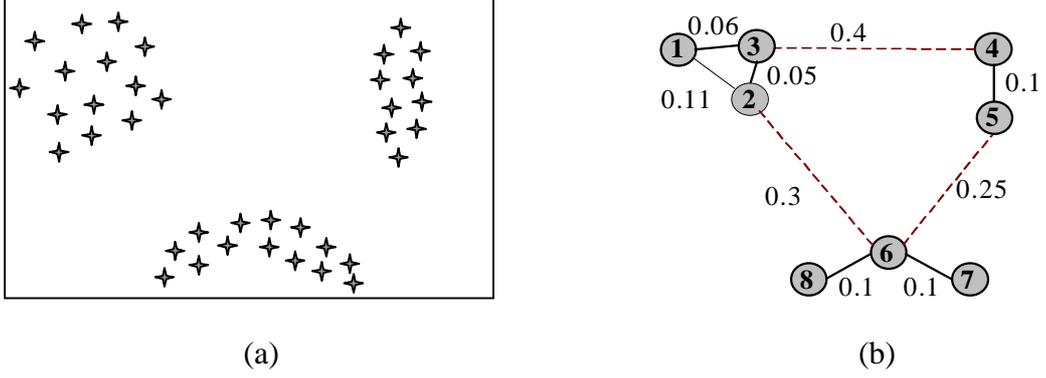


Figure 4: aiNet illustration. (a) Available data set with three clusters of high data density. (b) Network of labeled cells with their connection strengths assigned to the links. The dashed lines indicate connections to be pruned in order to generate disconnected sub-graphs, each characterizing a different cluster in the network.

Similarly to the models of Jerne (1974b) and Farmer et al. (1986), we make no distinction between the network cells and their surface molecules (antibodies). The Ag-Ab and Ab-Ab interactions are quantified through proximity (or similarity) measures. The goal is to use a distance metric to generate an antibody repertoire that constitutes the internal image of the antigens to be recognized, and evaluate the similarity degree among the aiNet antibodies, such that the cardinality of the repertoire can be controlled. Thus, the Ag-Ab affinity is inversely proportional to the distance between them: the smaller the distance, the higher the affinity, and vice-versa.

It is important to stress that, in the biological immune system, recognition occurs through a complementary match between a given antigen and the antibody. Nevertheless, in several artificial immune system applications (Hajela & Yoo, 1999; Hart & Ross, 1999; Oprea, 1999), and for the purposes of this model, the generation of an antibody repertoire with similar characteristics (instead of complementary) to the antigen set is a suitable choice.

As proposed in the original immune network theory, the existing cells will compete for antigenic recognition and those successful will lead to the network activation and cell proliferation (according to the clonal selection principle described in Section 3), while those who fail will be eliminated. In addition, Ab-Ab recognition will result in network suppression. In our model, suppression is performed by eliminating the self-recognizing antibodies, given a suppression threshold σ_s . Every pair Ag_j-Ab_i , $j = 1, \dots, M$, $i = 1, \dots, N$, will relate to each other within the shape-space S through the affinity $d_{i,j}$ of their interactions, which reflects the probability of starting a clonal response. Similarly, an affinity $s_{i,j}$ will be assigned to each pair Ab_j-Ab_i , $i, j = 1, \dots, N$, reflecting their interactions (similarity).

The following notation will be adopted:

- **Ab**: available antibody repertoire ($\mathbf{Ab} \in S^{N \times L}$, $\mathbf{Ab} = \mathbf{Ab}_{\{d\}} \cup \mathbf{Ab}_{\{m\}}$);
- $\mathbf{Ab}_{\{m\}}$: total memory antibody repertoire ($\mathbf{Ab}_{\{m\}} \in S^{m \times L}$, $m \leq N$);
- $\mathbf{Ab}_{\{d\}}$: d new antibodies to be inserted in **Ab** ($\mathbf{Ab}_{\{d\}} \in S^{d \times L}$);
- **Ag**: population of antigens ($\mathbf{Ag} \in S^{M \times L}$);
- f_j : vector containing the affinity of all the antibodies \mathbf{Ab}_i ($i = 1, \dots, N$) with relation to antigen \mathbf{Ag}_j . The affinity is inversely proportional to the Ag-Ab distance;
- **S**: similarity matrix between each pair $\mathbf{Ab}_i - \mathbf{Ab}_j$, with elements $s_{i,j}$ ($i, j = 1, \dots, N$);
- **C**: population of N_c clones generated from **Ab** ($\mathbf{C} \in S^{N_c \times L}$);
- **C***: population **C** after the affinity maturation process;
- d_j : vector containing the affinity between every element from the set **C*** with \mathbf{Ag}_j ;
- ζ : percentage of the mature antibodies to be selected;
- \mathbf{M}_j : memory clone for antigen \mathbf{Ag}_j (remaining from the process of clonal suppression);
- \mathbf{M}_j^* : resultant clonal memory for antigen \mathbf{Ag}_j ;
- σ_d : natural death threshold; and
- σ_s : suppression threshold.

The aiNet learning algorithm aims at building a memory set that recognizes and represents the data structural organization. The more specific the antibodies, the less parsimonious the network (low compression rate), whilst the more generalist the antibodies, the more parsimonious the network with relation to the number of antibodies (improved compression rate). The suppression threshold (σ_s) controls the specificity level of the antibodies, the clustering accuracy and network plasticity. In order to provide the user with important hints on how to set up the aiNet parameters, a sensitivity analysis of the algorithm with relation to the most critical user-defined parameters will be presented in Section 7.

The aiNet learning algorithm can be described as follows:

1. At each iteration, do:

1.1. For each antigenic pattern \mathbf{Ag}_j , $j = 1, \dots, M$, ($\mathbf{Ag}_j \in \mathbf{Ag}$), do:

1.1.1. Determine its affinity $f_{i,j}$, $i = 1, \dots, N$, to all \mathbf{Ab}_i . $f_{i,j} = 1/D_{i,j}$, $i = 1, \dots, N$:

$$D_{i,j} = \|\mathbf{Ab}_i - \mathbf{Ag}_j\|, \quad i = 1, \dots, N \quad (2)$$

1.1.2. A subset $\mathbf{Ab}_{\{n\}}$ composed of the n highest affinity antibodies is selected;

1.1.3. The n selected antibodies are going to proliferate (clone) proportionally to their antigenic affinity $f_{i,j}$, generating a set **C** of clones: the higher the affinity,

the larger the clone size for each of the n selected antibodies (see Equation (7));

- 1.1.4. The set \mathbf{C} is submitted to a directed affinity maturation process (guided mutation) generating a mutated set \mathbf{C}^* , where each antibody k from \mathbf{C}^* will suffer a mutation with a rate α_k inversely proportional to the antigenic affinity $f_{i,j}$ of its parent antibody: the higher the affinity, the smaller the mutation rate:

$$\mathbf{C}_k^* = \mathbf{C}_k + \alpha_k (\mathbf{A}\mathbf{g}_j - \mathbf{C}_k); \alpha_k \propto 1/f_{i,j}; k = 1, \dots, N_c; i = 1, \dots, N. \quad (3)$$

- 1.1.5. Determine the affinity $d_{k,j} = 1/D_{k,j}$ among $\mathbf{A}\mathbf{g}_j$ and all the elements of \mathbf{C}^* :

$$D_{k,j} = \|\mathbf{C}_k^* - \mathbf{A}\mathbf{g}_j\|, k = 1, \dots, N_c. \quad (4)$$

- 1.1.6. From \mathbf{C}^* , re-select $\zeta\%$ of the antibodies with highest $d_{k,j}$ and put them into a matrix \mathbf{M}_j of clonal memory;

- 1.1.7. *Apoptosis*: eliminate all the memory clones from \mathbf{M}_j whose affinity $D_{k,j} > \sigma_d$:

- 1.1.8. Determine the affinity $s_{i,k}$ among the memory clones:

$$s_{i,k} = \|\mathbf{M}_{j,i} - \mathbf{M}_{j,k}\|, \forall i, k. \quad (5)$$

- 1.1.9. *Clonal suppression*: eliminate those memory clones whose $s_{i,k} < \sigma_s$:

- 1.1.10. Concatenate the total antibody memory matrix with the resultant clonal memory \mathbf{M}_j^* for $\mathbf{A}\mathbf{g}_j$: $\mathbf{A}\mathbf{b}_{\{m\}} \leftarrow [\mathbf{A}\mathbf{b}_{\{m\}}; \mathbf{M}_j^*]$;

- 1.2. Determine the affinity among all the memory antibodies from $\mathbf{A}\mathbf{b}_{\{m\}}$:

$$s_{i,k} = \|\mathbf{A}\mathbf{b}_{\{m\}}^i - \mathbf{A}\mathbf{b}_{\{m\}}^k\|, \forall i, k. \quad (6)$$

- 1.3. *Network suppression*: eliminate all the antibodies such that $s_{i,k} < \sigma_s$:

- 1.4. Build the total antibody matrix $\mathbf{A}\mathbf{b} \leftarrow [\mathbf{A}\mathbf{b}_{\{m\}}; \mathbf{A}\mathbf{b}_{\{d\}}]$

2. Test the stopping criterion.

To determine the total clone size N_c generated for each of the M antigens, the following equation was employed:

$$N_c = \sum_{i=1}^n \text{round}(N - D_{i,j} \cdot N), \quad (7)$$

where N is the total amount of antibodies in $\mathbf{A}\mathbf{b}$, $\text{round}(\cdot)$ is the operator that rounds the value in parenthesis towards its closest integer and $D_{i,j}$ is the distance between the selected antibody i and the given antigen $\mathbf{A}\mathbf{g}_j$, given by Equation (2).

In the above algorithm, Steps 1.1.1 to 1.1.7 describe the clonal selection and affinity maturation processes as proposed by de Castro and Von Zuben (2000a) in their

computational implementation of the clonal selection principle. Steps 1.1.8 to 1.3 simulate the immune network activity.

As can be seen by the aiNet learning algorithm, a clonal immune response is elicited to each presented antigenic pattern. Notice also, the existence of two suppressive steps in this algorithm (1.1.9 and 1.3), that we call *clonal suppression* and *network suppression*, respectively. As far as a different clone is generated to each antigenic pattern presented, a clonal suppression is necessary to eliminate intra-clonal self-recognizing antibodies, while a network suppression is required to search for similarities between different sets of clones. After the learning phase, the network antibodies represent internal images of the antigens (or groups of antigens) presented to it.

The network outputs can be taken to be the matrix of memory antibodies' co-ordinates ($\mathbf{Ab}_{\{m\}}$) and their matrix of affinity (\mathbf{S}). While matrix $\mathbf{Ab}_{\{m\}}$ represents the network internal images of the antigens presented to the aiNet, matrix \mathbf{S} is responsible for determining which network antibodies are connected to each other, describing the general network structure.

To evaluate the aiNet convergence, several alternative criteria can be proposed:

1. Stop the iterative process after a pre-defined number of iteration steps;
2. Stop the iterative process when the network reaches a pre-defined number of antibodies;
3. Evaluate the average error between all the antigens and the network memory antibodies ($\mathbf{Ab}_{\{m\}}$) by calculating the distance from each network antibody to each antigen (this strategy will be useful for less parsimonious solutions), and stop the iterative process if this average error is larger than a pre-specified threshold; and
4. The network is supposed to have converged if its average error rises after k consecutive iterations.

aiNet Characterization and Analysis

The aiNet model can be classified as a *connectionist, competitive and constructive network*, where the antibodies correspond to the network nodes and the antibody concentration and affinity are their states. The learning mechanisms are responsible for the changes in antibody concentration and affinity. The connections among antibodies ($s_{i,k}$) corresponds to the physical mechanisms that measure their affinity, quantifying the immune network recognition. The aiNet graph representation describes its architecture, with the definition of the final number and spatial distribution of clusters. The dynamics governs the plasticity of the aiNet, while the metadynamics is responsible for a broader exploration of the

search-space and maintenance of diversity. The aiNet can also be classified as competitive, once its antibodies compete with each other for antigenic recognition and, consequently, survival. Antigenic competition is evident in Steps 1.1.2 and 1.1.6, while the competition for survival is performed in Step 1.1.7. Finally, the aiNet is plastic in nature, in the sense that its architecture, including number and shape of cells, is adaptable according to the problem.

The aiNet general structure is different from neural network models (Haykin, 1999) if one considers the function of the nodes and their connections. In the aiNet case, the nodes work as internal images of ensembles of patterns (thus representing the acquired knowledge), and the connection strengths describe the similarities among these ensembles. On the other hand, in the neural network case, the nodes are processing elements while the connection strengths may represent the knowledge.

As discussed by de Castro and Von Zuben (2000a), the immune clonal selection pattern of antigenic response can be seen as a microcosm of Darwinian evolution. The processes of simulated evolution (Holland, 1998) try to mimic some aspects of the original theory of evolution. Regarding the aiNet learning algorithm, it is possible to notice several features in common with simulated evolution (evolutionary algorithms). First, the aiNet is population based: an initial set of candidate solutions (antibodies) properly coded, is available at the beginning of the learning process. Second, a function to evaluate these candidate solutions has to be defined: an affinity measure as given by Equations (2) and (3). Third, the genetic encoding of the generated offspring (clones) is altered through a hypermutation mechanism. At last, several parameters have to be defined, like the number of highest affinity antibodies to be selected, and the natural death and suppression thresholds.

Related Immune Network Models

The proposed artificial immune network model also follows the general immune network structure presented in Equation (1), i.e. the rate of population variation is proportional to the sum of the network novel antibodies (Step 1.4), minus the death of unstimulated antibodies (Step 1.1.7), plus the reproduction of stimulated antibodies (Step 1.1.3). As a complement, we suppress self-recognizing antibodies (Steps 1.1.9 and 1.3). Nevertheless, the essence of the aiNet model is different from the existing ones in two respects. First, and most important, it is a discrete (iterative) instead of continuous model. Second, our network model may not be directly reproducing any biological immune phenomenon. The goal is to use the immune network paradigm, together with the clonal selection behavior of antigenic responses, as inspiration to develop an adaptive system

capable of solving complex information processing tasks, like data compression, pattern recognition, classification and clustering. Artificial immune networks (Dasgupta, 1999; de Castro & Von Zuben, 1999a) are problem dependent, in the sense that they are built according to the antigen set (problem).

Hunt and Cooke (1996) proposed an artificial immune system (AIS) model, based upon the immune network theory, to perform machine learning. The key features they tried to explore were a mechanism to construct the antibodies, a content addressable memory, the immune recognition (matching) mechanism and its self-organizing properties. Like in the aiNet case, they did not intend to provide a deep association between their proposed model and the vertebrate immune system. Instead, useful features were explored for the development of problem solving tools. Their model was based on the following elements with their respective roles and characteristics:

1. a bone marrow: generates antibodies, decides where in the network to insert the antigen, decides which B cell dies and triggers the addition of cells to the network;
2. B cells: carry the genetic information to build antibodies (and the antibodies themselves) along with their stimulation level;
3. antibodies: possess the paratope pattern;
4. antigen: possesses a single epitope; and
5. stimulation level: evaluates the strength of the Ag-Ab match and the affinity between different B cells.

Table 1 compares aiNet with the AIS model of Hunt and Cooke (1996).

Table 1: Trade-off between aiNet and the network model of Hunt and Cooke (1996).

	aiNet	Hunt & Cooke
Nodes	Antibodies	B cells
Coding	Real-valued vectors	Binary strings
Network initialization	Random with small influence in the final network (see Section 7)	Critical for the processing time
Antigenic presentation	To all the network	To a randomly chosen part of the network
Affinities	Euclidean distance	Proportional to the number of matching bits
Cell death	Suppressing antibodies with low antigenic and high antibody affinities	Suppressing B cell with low stimulation levels
Hipermutation	Controls learning	Promotes diversity

Timmis (2000) uncovered many problems of the model proposed by Hunt and Cooke (1996) and proposed a completely new domain independent algorithm with few parameters. His AIN (artificial immune network) is initialized as a network of B cell objects composed of a cross section of the data set to be learnt, while the remainder makes up the training data set (antigens). Each member of the antigen set is matched against each B cell in the AIN, with the similarity being determined by the Euclidean distance. B cells are stimulated by this matching and other connected B cells in the network. The stimulation level of a B cell determines its survival, as 5% of the weakest cells are removed at each iteration. In addition, the stimulation level indicates whether the B cell is going to be cloned. Table 2 presents the main differences between the aiNet and the AIN model of Timmis (2000).

Table 2: Main differences between aiNet and the AIN of Timmis (2000).

	aiNet	Timmis
Nodes	Antibodies	B cells
Network initialization	Random with small influence in the final network (see Section 7)	Cross section of the antigen set (training data set)
Cell death	Suppressing antibodies with low antigenic and high antibody affinities	Suppressing B cell with low stimulation levels
Hipermutation	Controls learning	Promotes diversity
Network connectivity	Defined by a hierarchical clustering or graph theoretical strategy	B cells with high affinities among each other (given a certain threshold) are linked to areas of the network with closest affinity

Analysis of the Algorithm

Analysis of an algorithm refers to the process of deriving estimates for the time and space needed during execution. Complexity of an algorithm refers to the amount of time and space required during execution (Johnsonbaugh, 1997). Determining the performance parameters of a computer program is a difficult task and depends on a number of factors such as the computer being used, the way the data are represented, how and with which programming language the code is implemented. The time needed to execute an algorithm is also a function of the input. Instead of dealing directly with the input, we use parameters that characterize its size, like the number (L) of variables of each input vector and the amount (M) of samples (patterns) available. In addition, the total number of network cells N , the number of cell clones N_c , and the network final number (m) of memory cells will be necessary to

evaluate its complexity. The most computational intensive step of the aiNet learning algorithm is the determination of the affinity between all the network antibodies (Step 1.2). The computation time required to compare all the elements of a matrix of size m is $O(m^2)$. Due to the asymptotic nature of the computational complexity, the total cost of the algorithm is taken to be $O(m^2)$. It is important to notice that m may vary along the learning iterations, such that at each generation the algorithm has a different computational cost.

Knowledge Extraction and Structure of a Trained aiNet

The aiNet memory antibodies $\mathbf{Ab}_{\{m\}}$ represent internal images of the antigens to which it is subject. This feature demands a representation in the same shape-space for the network of antibodies and for the antigens. Hence, visualizing the network for antigens (and antibodies) with $L > 3$ becomes a difficult task. In order to alleviate this difficulty, we suggest the use of several hierarchical clustering techniques to interpret the generated network. These techniques will help us to define the aiNet structure.

The aiNet structure could simply be determined by fully connecting all the network cells according to matrix \mathbf{S} , but it would make the network interpretation and knowledge extraction unfeasible processes. One way of simply reducing the complexity of a fully connected network of cells is to suppress all those connections whose strength extrapolates a pre-defined threshold. This idea though simple, does not account for any information within the network antibodies (indirectly in the data set) and might lead to erroneous interpretations of the resultant network. It is the main purpose here, to supply the user with formal and robust network interpretation strategies with high confidence levels. Explicitly speaking, the goals are to determine (1) the number of *clusters*, or classes (whenever a cluster corresponds to a class), (2) the spatial distribution of each cluster, and (3) the network antibodies belonging to each of the identified clusters. To do so, the network output is used, which is composed of the number m of memory antibodies, the matrix $\mathbf{Ab}_{\{m\}}$ of memory antibodies, and the upper triangular matrix \mathbf{S} of distances among these memory antibodies, along with some principles from cluster analysis. The problem is stated as follows.

Given a network with m memory antibodies (matrix $\mathbf{Ab}_{\{m\}}$), each of which being a vector of dimension L ($\mathbf{Ab}_{\{m\}} \in \mathfrak{R}^{m \times L}$) and their interconnections (matrix \mathbf{S}), devise a clustering scheme to detect inherent separations between subsets (clusters) of $\mathbf{Ab}_{\{m\}}$ in a metric space governed by a distance measure $d(x,y)$.

The algorithms to be presented here are well known from the statistical literature, but will be adapted and interpreted to the immune network paradigm. Under this perspective, the aiNet becomes responsible for extracting knowledge from the data set, while hierarchical cluster analysis techniques will be used to detect clusters in the resultant network, i.e. to interpret the aiNet. The network can be seen as a pre-processing for the cluster analysis technique, being a powerful tool to filter out redundant data from a given data set.

Hierarchical Clustering and Graph Theoretical Techniques

To illustrate the methods that will be used and the ones to be proposed, consider one of the simplest problems of data clustering presented in Figure 5(a). There are 50 samples subdivided into 5 clusters (non-overlapping classes) of 10 samples each. Figure 5(b) depicts the automatically generated network cells, considering the following aiNet training parameters: $n = 4$, $\zeta = 0.2$, $\sigma_d = 1.0$, $\sigma_s = 0.14$ and $d = 10$. The stopping criterion is a fixed number of generations: $N_{gen} = 10$. The resulting network contains only 10 cells, reducing the problem to 20% of its original complexity (size), corresponding to a compression rate $CR = 80\%$.

Hierarchical techniques may be subdivided into agglomerative methods, which proceed by a series of successive fusions of the m entities (antibodies) into groups, and divisive methods, which partition the set of m entities (antibodies) successively into finer partitions. The results of both agglomerative and divisive techniques may be represented in the form of a dendrogram, which is a two-dimensional diagram illustrating the fusions or partitions which have been made at each successive level (Everitt, 1993).

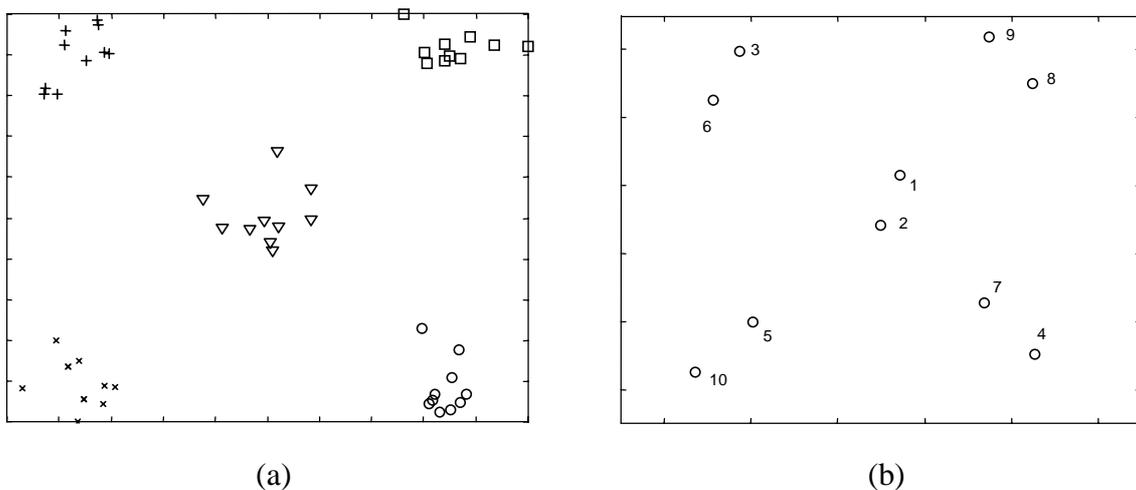


Figure 5: aiNet illustration. (a) Learning data. (b) Resulting network antibodies.

In this work, we will focus on the agglomerative methods, more specifically the nearest neighbor (or single link) method, and the centroid cluster analysis. As the aiNet may be seen as an interconnected graph of antibodies, it will be explored some graph-theoretical strategies for detecting and describing clusters, in particular the minimal spanning tree, MST.

The aiNet outputs are m , the matrix \mathbf{S} of dimension $S^{m(m-1)/2}$ and matrix $\mathbf{Ab}_{\{m\}}$ of dimension $S^{m \times L}$. Hence, the application of hierarchical methods for the construction of a dendrogram, like the nearest and furthest neighbor and centroid, is straightforward.

From \mathbf{S} and the methods listed above, we wish to construct a tree, or a nested set of clustering of the objects, in order to provide a striking visual display of similarity groupings of the network cells.

Definition 2: A *dendrogram* is defined as a rooted weighted tree where all terminal nodes are at the same distance (path length) from the root (Lapointe & Legendre, 1991).

We will not get into details on how to construct a dendrogram from a similarity matrix, the interested reader shall refer to Hartigan (1967) and Hubert, Arabie, and Meulman (1998). For the purposes of this chapter, three characteristics can adequately describe a dendrogram: its topology, labels and cluster heights (Lapointe & Legendre, 1995). Figure 6 illustrates the dendrogram representation for the centroid cluster strategy and the aiNet antibodies depicted in Figure 5(b). Notice the topology, cell labels, and cluster heights, representing the Ab-Ab affinities.

Virtually all clustering procedures provide little if any information as to the number of clusters present in data. Nonhierarchical procedures usually require the user to specify this parameter before any clustering is accomplished (the reason why we chose to use hierarchical methods instead of nonhierarchical ones) and hierarchical methods routinely produce a series of solutions ranging from m clusters to a solution with only one cluster present. As can be seen from Figure 6, the dendrogram can be broken at different levels to yield different clusterings of the network antibodies. In this case, the large variations in heights allow us to distinguish 5 clusters among the network antibodies, in accordance with the network depicted in Figure 5(b). This procedure is called *stepsizes* and involves examining the differences in *fusion values* between hierarchy levels. A broad review of several different methods for determining the number of clusters in a set of objects can be found in Milligan and Cooper (1985).

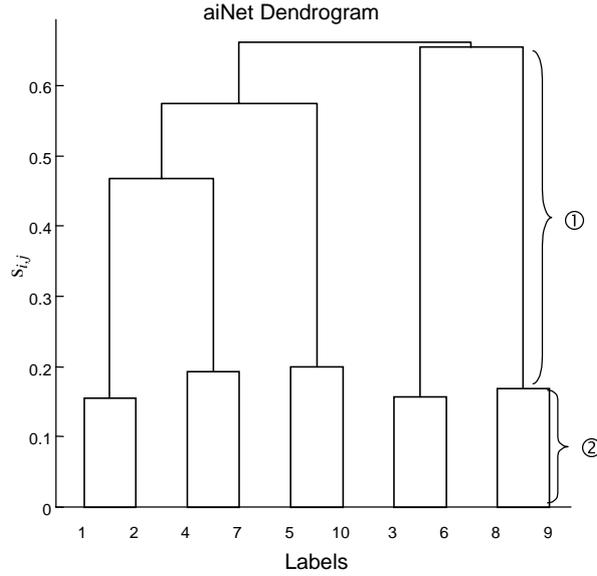


Figure 6: Dendrogram of the aiNet antibodies for the centroid method depicting large differences in the fusion values (① with relation to ②).

Keeping track of the nearest neighbor hierarchical clustering technique, we can find the *minimal spanning tree* (MST) of a graph to be a powerful mechanism to search for a locally adaptive interconnecting strategy for the network cells (Zahn, 1971). The MST will serve as another aid to detect and describe the structure of the aiNet clusters.

Definition 3: A tree is a *spanning tree* of a graph if it is a sub-graph containing all the vertices of the graph. A *minimal spanning tree* of a graph is a spanning tree with minimum weight. The weight of a tree is defined as the sum of the weights of its constituent edges (Leclerc, 1995).

Figure 7(a) depicts the minimal spanning tree (MST) for the constructed network. The visualization of this tree is only feasible for $L \leq 3$. By using the algorithm known as Prim's algorithm (Prim, 1957) to build the MST, we can draw a bar graph (see Figure 7(b)) representing the distances between neighboring cells.

Definition 4: A *minimax path* is the path between a pair of nodes that minimizes, over all paths, the cost, which is the maximum weight of the path (Carroll, 1995).

This definition is important in the aiNet context, once the preference for minimax paths in the MST forces it to connect two nodes i and j belonging to a tight cluster without straying outside the cluster. If the MST of a graph G is unique, then the set of minimax links of G defines this MST, else it defines the union of all MSTs of G .

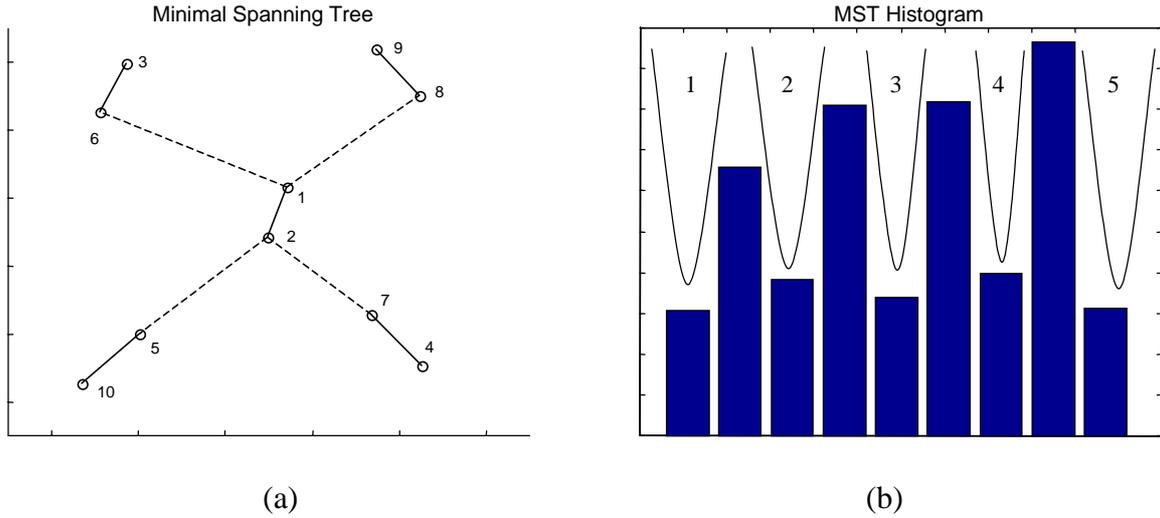


Figure 7: The minimal spanning tree and its histogram. (a) Edges to be removed (dashed lines) based upon the *factor* criteria, $r = 2$. (b) Number of clusters (Peaks + 1, or Valleys) for this MST.

Up to here, notice that the MST is used to define the number of network clusters, which will be equal to the number of higher peaks in the bar graph of Figure 7(b) plus one, indicating large variations in the minimax distances between cells. When the aiNet learning algorithm generates more than one antibody for each cluster (which is the case for this particular run), the number of clusters can also be measured as the number of valleys of the respective histogram. On the other hand, the dendrograms allow us not only to define the number of clusters but also to identify the elements (nodes or antibodies) belonging to each cluster. In order to automatically define the number and nodes composing each cluster of an MST, we can use some of the techniques proposed by Zahn (1971).

It is quite helpful that the MST does not break up the real clusters in $\mathbf{Ab}_{\{m\}}$, but at the same time neither does it force breaks where real gaps exist in the geometry of the network. A spanning tree is forced by its nature to span all the nodes in a network, but at least the MST jumps across the smaller gaps first.

There is the problem of deleting edges from an MST so that the resulting connected subtrees correspond to the observable clusters. In the example of Figure 7(a), we need an algorithm that can detect the appropriateness of deleting the edges 1-6, 1-8, 2-7 and 2-5. The following criterion is used.

An MST edge (i,k) whose weight $s_{i,k}$ is significantly larger than the average of nearby edge weights on both sides of the edge (i,k) should be deleted. This edge is called inconsistent.

There are two natural ways to measure the significance referred to. One is to see how many sample standard deviations separate $s_{i,k}$ from the average edge weights on each side. The other is to calculate the *factor* or *ratio* (r) between $s_{i,k}$ and the respective averages.

To illustrate this criterion, let us assume that all edges whose $s_{i,k}$ is greater than the average of nearby edges plus two standard deviations will be deleted, i.e., a factor $r = 2$ is chosen. Edges 1-6, 1-8, 2-7 and 2-5 will be selected for deletion (dashed lines in Figure 7(a)).

After determining the edges to be deleted, we can determine the number (p) of existing clusters ($c_i, i = 1, \dots, p$) in the aiNet and their respective components (antibodies). In this case, $c_1 = [6,3]$, $c_2 = [8,9]$, $c_3 = [1,2]$, $c_4 = [5,10]$ and $c_5 = [4,7]$.

The discussed criterion would fail to determine the correct number of network clusters in cases the network reaches its minimal size for the given data set and the clusters are approximately uniformly distributed over the search space. As an example, for the proposed problem (Figure 5(a)), suppose a minimal network with four cells was found. This would result in the detection of a single cluster by the MST factor criterion. On the other hand, all the remaining network antibodies could be seen as internal images of the data clusters, implying that its number is equal to the number of antibodies, each of which representing a single cluster.

As one last aspect of clustering to be discussed, consider the problem of cluster representation. Assume that each cluster can be uniquely represented by its *center of mass* ($v_k, k = 1, \dots, p$), or *centroid*, and the distance between clusters defined as the distance between the cluster centroids. Figure 8 depicts the resultant network antibodies defined by the aiNet learning algorithm and the network determined by the MST clustering algorithm described above for $r = 2$. The stars represent the centroids of each cluster. The use of the centroid to represent a cluster works well when the clusters are compact or isotropic. However, when the clusters are elongated or non-isotropic, this scheme fails to represent them properly, as will be discussed in the case of two examples presented in Section 6. Representing clusters by their centroids allow us to assign membership levels to each aiNet antibody with relation to the determined clusters, yielding a *fuzzy clustering* scheme.

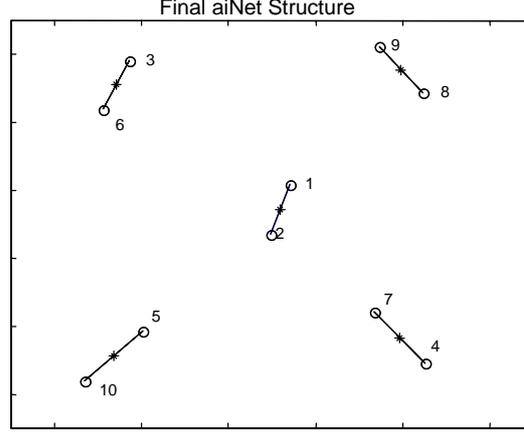


Figure 8: The resultant network is composed of five separate sub-graphs (sub-networks), each corresponding to a different cluster. The stars represent the centroids of each cluster.

aiNet Fuzzy Clustering

The presented clustering approaches generate partitions. In a partition, each cell belongs to one and only one cluster. Thus, the clusters in this hard clustering scheme are disjoint. Fuzzy clustering extends this notion to associate each cell (antibody) with every cluster using a membership function (Bezdek & Pal, 1992). The most well known fuzzy clustering techniques are the *fuzzy k-means* and the *fuzzy c-means* algorithms that iteratively update the cluster centers according to an actual proximity matrix (\mathbf{U}) until a small variation in \mathbf{U} is achieved. A brief exposition of fuzzy partition spaces is given by Bezdek and Pal (1992):

Let M be an integer, $1 < c < M$, and let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$ denote a set of M unlabeled feature vectors in \mathfrak{R}^L . Given \mathbf{X} , we say that p fuzzy subsets $\{u_i: \mathbf{X} \rightarrow [0,1]\}$ are a fuzzy p -partition of \mathbf{X} in case the (pM) values $\{u_{i,k} = u_i(\mathbf{x}_k), 1 \leq k \leq M, 1 \leq i \leq p\}$ satisfy three conditions:

$$0 \leq u_{i,k} \leq 1 \quad \text{for all } i, k; \quad (8a)$$

$$\sum u_{i,k} = 1 \quad \text{for all } k; \quad (8b)$$

$$0 < \sum u_{i,k} < 1 \quad \text{for all } i. \quad (8c)$$

Each set of (pM) values satisfying conditions (8a-c) can be arrayed as a $(p \times M)$ matrix $\mathbf{U} = [u_{i,k}]$. The set of all such matrices is the nondegenerate fuzzy c -partitions of \mathbf{X} .

After the number and members of each (hard) cluster are defined, and the network clusters are represented by their centers of mass, it is possible to apply a fuzzy clustering concept to the aiNet, where each antibody will have a measurable membership value to each of the determined clusters (centroids). In the aiNet context, the fuzzy clustering relaxes the membership of the network antibodies to the cluster centers, $\mathbf{U} = [u_{i,k}]$, which in this case can assume any value over the $[0,1]$ interval. Conditions (8b) and (8c) are also relaxed, so that the sum of memberships is not required to be one.

Matrix \mathbf{U} for the aiNet can be determined by calculating the distance between all the network memory antibodies $\mathbf{Ab}_{\{m\}}^i$, $i = 1, \dots, m$, and the centroids of the clusters v_k , $k = 1, \dots, p$, \mathbf{U}^* , normalizing its rows over the $[0,1]$ interval and then passing it through a squashing function, such that the smaller the distance between the aiNet antibodies and their respective centroids, the closer its membership value to unity. This can be achieved by applying a sigmoidal function to $\mathbf{U} = 1./\mathbf{U}^*$, where the $./$ operator means that each value of the \mathbf{U} matrix will be determined by dividing one by the respective element of \mathbf{U}^* .

The proximity matrix \mathbf{U} assigning the membership of each network antibody to the determined centroids is presented in Table 2. From this table it is possible to see that cells c_1 and c_2 belong to cluster v_1 with membership 1.0, to cluster v_2 with membership 0.58 and 0.63 ($u_{2,1}$ and $u_{2,2}$), respectively, and so on.

Table 3: Membership values for each cell c_i , $i = 1, \dots, 10$, with relation to each cluster centroid v_j , $j = 1, \dots, 5$.

	c_1	c_2	c_3	c_4	c_5	c_6	c_7	c_8	c_9	c_{10}
v_1	1.00	1.00	0.67	0.71	0.76	0.69	0.84	0.75	0.71	0.66
v_2	0.58	0.63	0.50	1.00	0.68	0.50	1.00	0.60	0.56	0.64
v_3	0.67	0.50	0.63	0.56	0.50	0.58	0.57	1.00	1.00	0.50
v_4	0.50	0.55	0.54	0.62	1.00	0.57	0.64	0.50	0.50	1.00
v_5	0.60	0.50	1.00	0.50	0.59	1.00	0.50	0.60	0.56	0.64

Empirical Results

In order to evaluate the performance of the aiNet, three benchmark problems were considered: SPIR, CHAINLINK, 5-NLSC, according to Figure 9. Note that, though the samples are labeled in the picture they are unlabeled for the aiNet. Each task has its own particularity and will serve to evaluate several network features, among which we can stress

cluster partition and representation, and its potential to reduce data redundancy. Table 1 presents the aiNet training parameters for all problems. The stopping criterion was a maximum number of generations N_{gen} .

In the three cases, the aiNet performance was compared to that of the Kohonen (1982) self-organizing map (SOM), which is also an unsupervised technique broadly used in clustering tasks. The SOM was implemented with a 0.9 geometrical decreasing learning rate (at each five iterations) with an initial value $\alpha_0 = 0.9$ and final value $\alpha_f = 10^{-3}$ as the stopping criterion. The weights were initialized using a uniform distribution over the interval $[-0.1, 0.1]$. The output grid is uni-dimensional with a variable output number of neurons according to the problem under evaluation. At the end of the SOM learning phase, all those output neurons that do not classify any input datum will be pruned from the network.

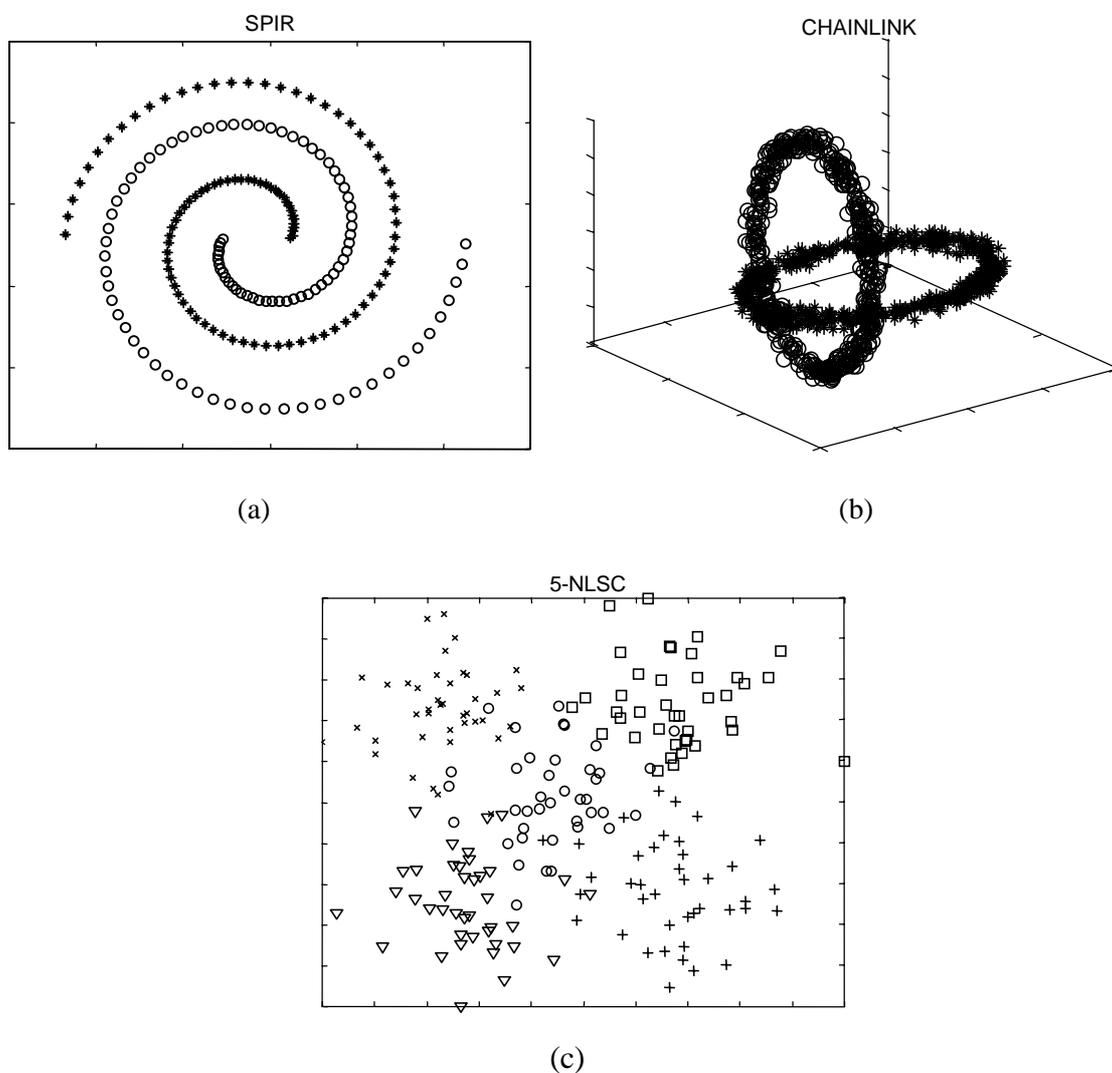


Figure 9: Test problems for the aiNet, where M is the number of samples. (a) SPIR, $M = 190$. (b) CHAINLINK, $M = 1000$. (c) 5-NLSC, $M = 200$.

Table 4: Training parameters for the aiNet learning algorithm.

	σ_s	σ_d	n	d	$\zeta(\%)$	N_{gen}
SPIR	0.07	1.0	4	10	10	40
CHAINLINK	0.15				10	40
5-NLSC	0.20				20	10

Two-Spirals Problem: SPIR

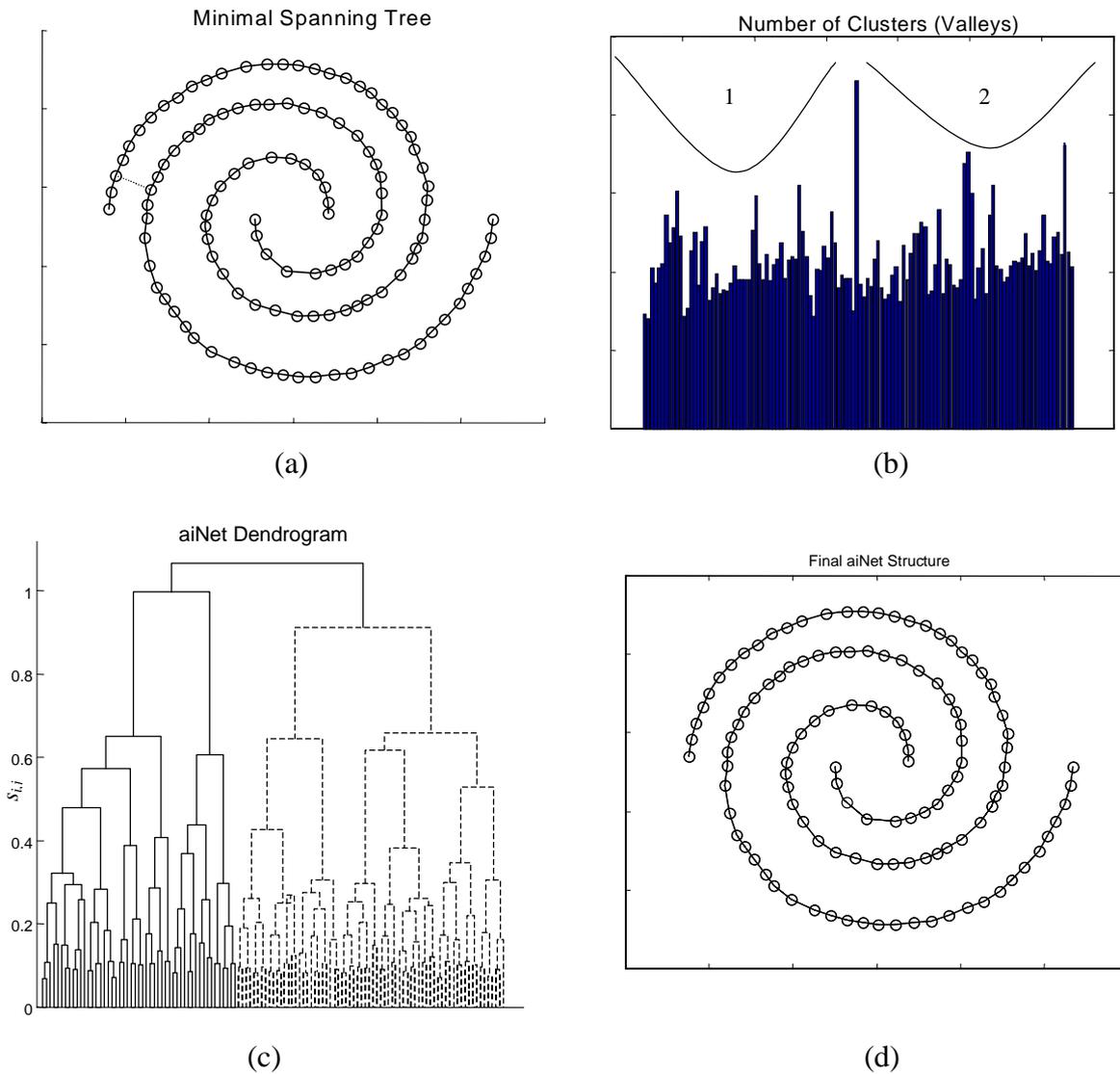


Figure 10: aiNet applied to the 2-spirals problem. (a) Minimal spanning tree, in which the dashed line represents the connection to be pruned. (b) MST histogram indicating two clusters. (c) aiNet dendrogram. (d) Final network structure, determining the spatial distribution of the 2 clusters.

The first problem was the so-called 2-spirals problem, illustrated in Figure 9(a). This training set is composed of 190 samples in \mathfrak{R}^2 . This task aims at testing the aiNet capability to detect non-linearly separable clusters.

Figures 10(a) and (b) depict the MST and its corresponding histogram. From the histogram we can detect the existence of two different clusters in the network, which are automatically obtained using a factor $r = 2$. In this case, the resultant memory matrix was composed of $m = 121$ antibodies, corresponding to a $CR = 36.32\%$ reduction in the sample set size. Note that the compression was superior in regions where the amount of redundancy is larger, i.e., the centers of the spirals (see Figure 9(a)). The network dendrogram also allows us to detect two large clusters of data, as differentiated by the solid and dashed parts of Figure 10(c).

Figures 11(a) and (b) present the results for the SOM with an initial number of 50 output neurons; a single output neuron was pruned after learning ($m = 49$). The network final configuration and the resultant U-matrix (Ultsch, 1995) indicates the way the neurons were distributed, and according to their neighborhood the SOM would not be able to appropriately solve this problem, since nothing can be inferred from the U-matrix plot.

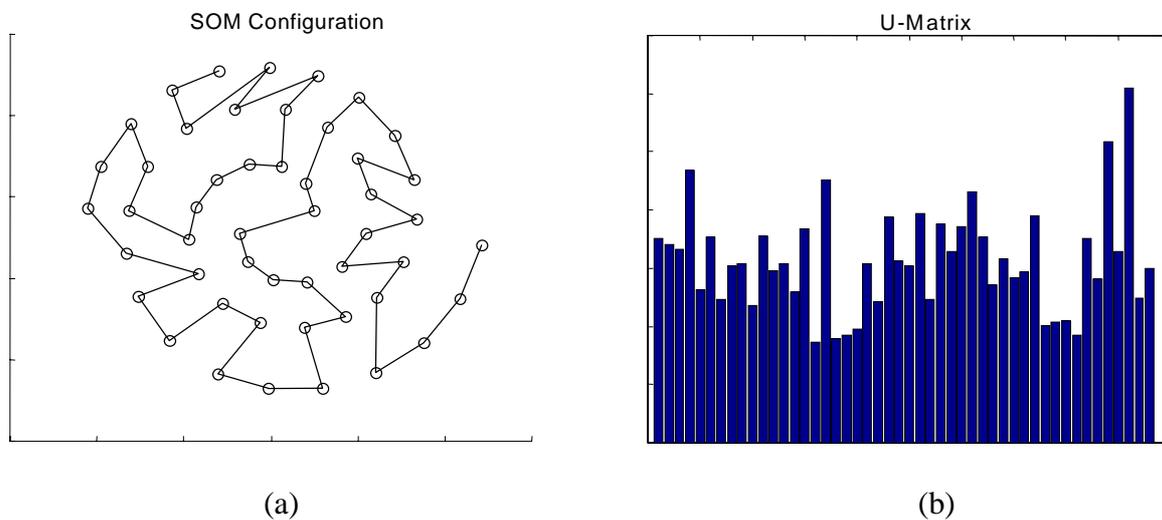


Figure 11: Results obtained by the application of the SOM to the SPIR problem. (a) Final network configuration and weight neighborhood, $m = 49$. (b) U-matrix.

The ChainLink Problem: CHAINLINK

1000 data points in the \mathfrak{R}^3 -space were arranged such that they form the shape of two intertwined 3-D rings, of whom one is extended along the x - y direction and the other one along the x - z direction. The two rings can be thought of as two links of a chain with each one consisting of 500 data points. The data is provided by a random number generator whose values are inside two toroids with radius $R = 1.0$ and $r = 0.1$ (see Figure 9(b)).

Statistical analysis of the data set shows that it has the following properties: the different components (x,y,z) of the data are uncorrelated (using Pearson's correlation coefficient) and the distribution of each component may be regarded as normal. Principal Component Analysis shows that there is no lower inherent dimensionality in the data. There exists however no (hyper-)plane that can separate the data set correctly into the two subsets, i.e., the data set is non-linearly separable. From the way the data set is created one can clearly state, however, that it consists of two distinguished subsets. The innersubset distance of a data point is by an order of magnitude smaller than the intersubset distance (Ultsch, 1995).

Figures 12(a) and (b) depict the MST and its corresponding histogram when the aiNet is applied to the CHAINLINK problem. From the histogram we can detect the existence of two different clusters in the network, which are automatically obtained using a factor $r = 2$.

Note that, in this case, the evaluation of the fusion values (*stepsize*) of the network dendrogram (Figure 12(c)) represents a difficult task, and may lead to an incorrect clustering. The compression rate of this problem was at the order of $CR = 94.5\%$ ($m = 55$).

Figures 13(a) and (b) depict the final SOM configuration taking into account the neurons' neighborhood and the U-matrix, respectively, for $m = 46$ output neurons (4 output neurons were pruned after learning, since they represent no input datum). In this case, note that the U-matrix is composed of 5 valleys, indicating 5 different clusters, what is not in accordance with the correct number of clusters. Notice either, from Figure 13(a), that the 5 clusters can be obtained by drawing 5 hyperplanes cutting each of the rings in its respective parts.

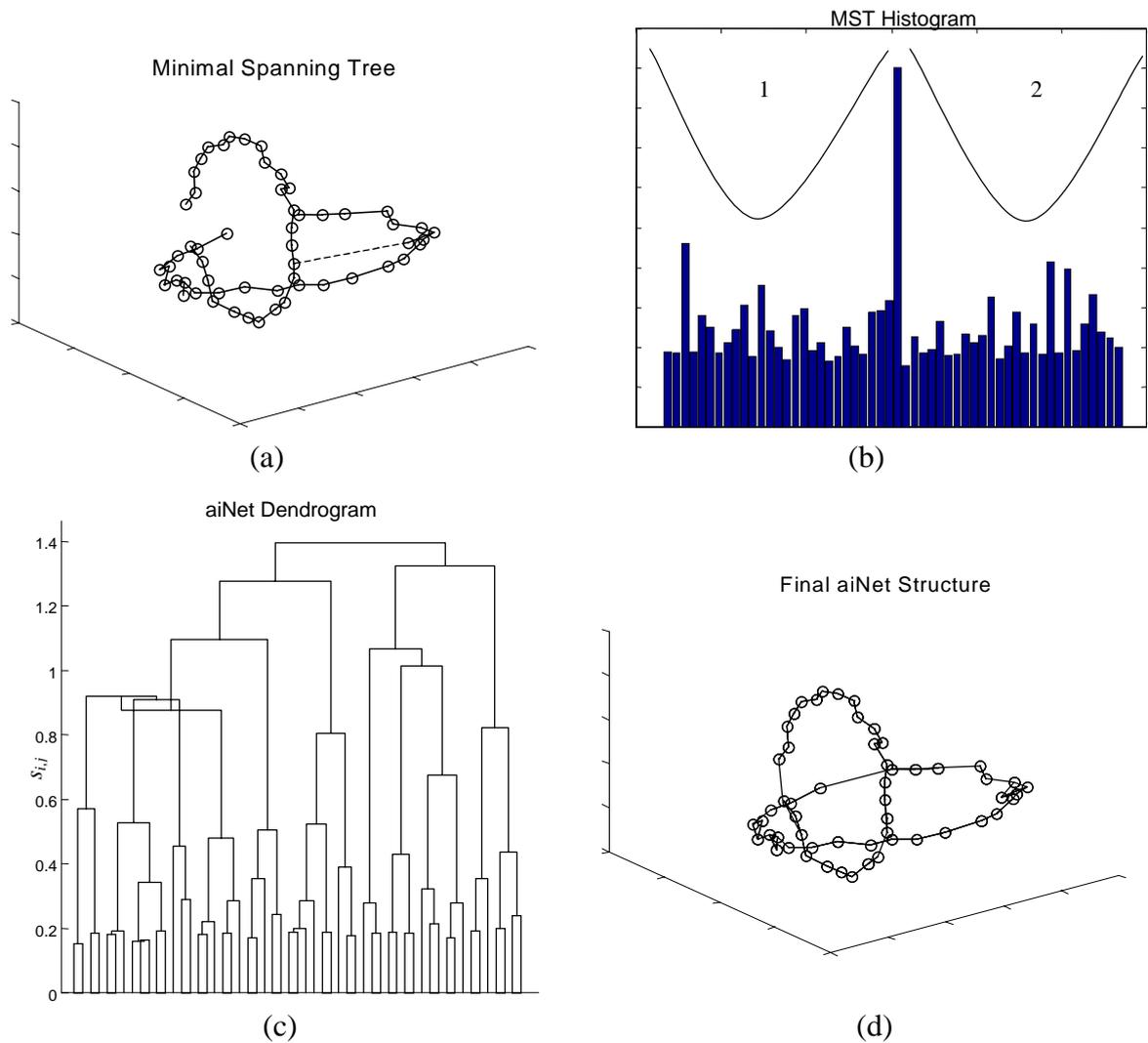


Figure 12: aiNet application to the CHAINLINK problem. (a) Minimal spanning tree with the dashed connection to be pruned. (b) MST histogram indicating the presence of two well separated clusters. (c) aiNet dendrogram. (d) Final network architecture.

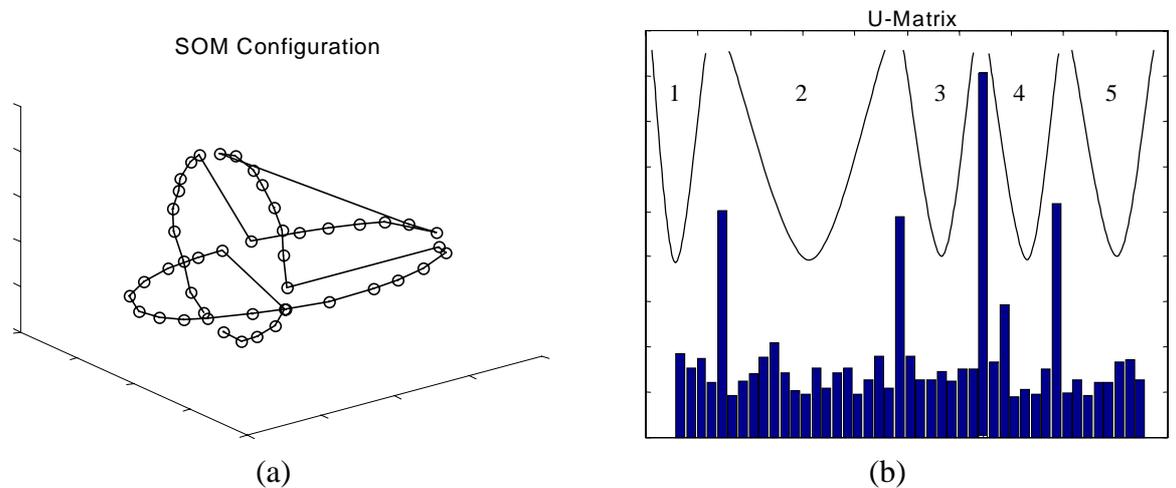


Figure 13: Uni-dimensional SOM applied to the CHAINLINK problem. (a) Final network configuration and weight neighborhood. (b) U-matrix.

Five Non-Linearly Separable Classes: 5-NLSC

As a last example, consider the problem illustrated in Figure 9(c). This example is particularly interesting, because the distinction among all the classes is not clear, even for a human observer. Note that, as in the previous examples, though the samples are labeled in the picture, they are unlabeled for the aiNet. This example has already been used by de Castro & Von Zuben (1999b) to evaluate the performance of a pruning method for the Kohonen SOM, named PSOM.

Figures 14(a) and (b) depict the MST and its corresponding histogram, respectively. The aiNet presented a compression rate $CR = 96\%$, with a final memory size $m = 8$.

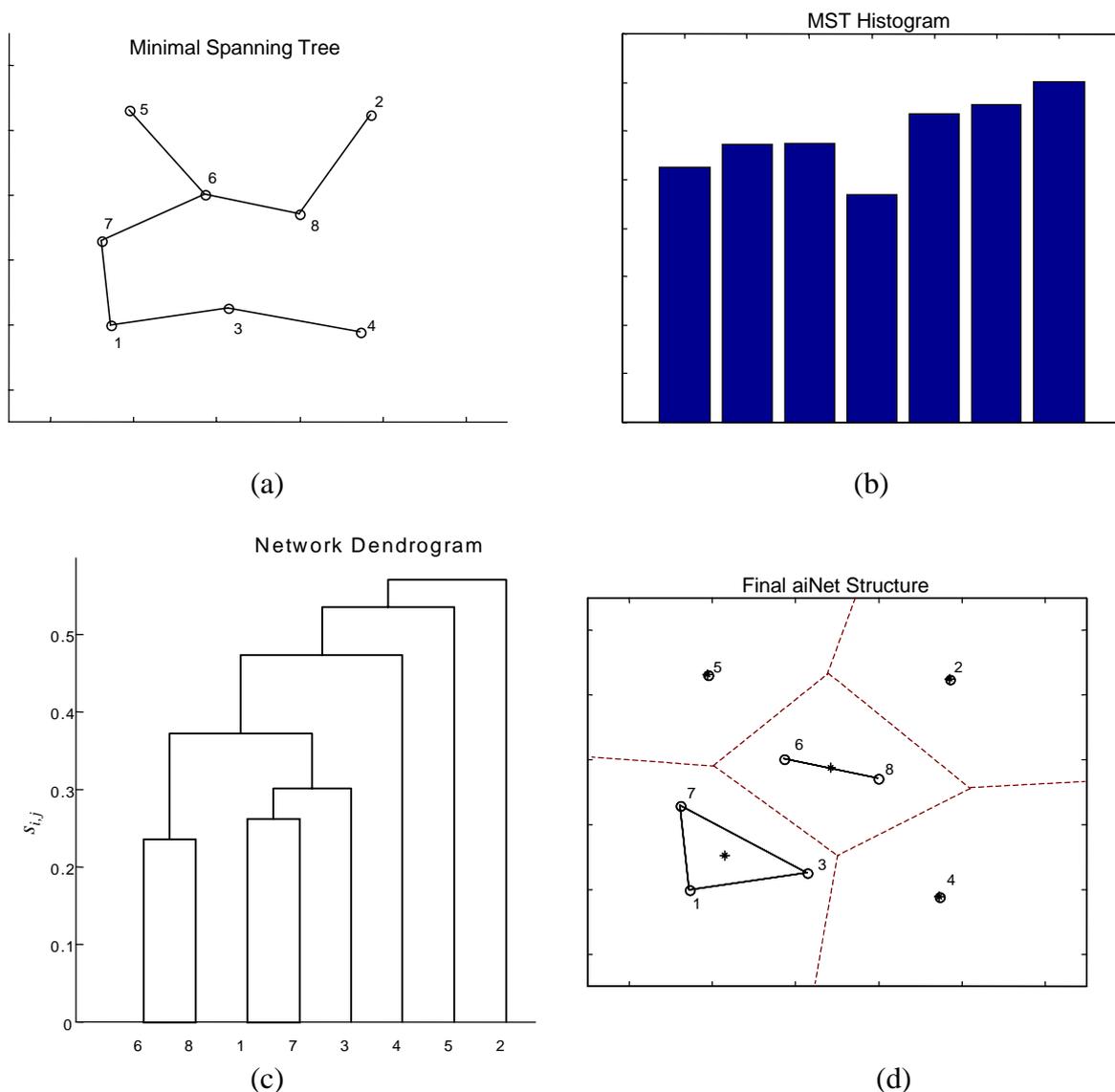


Figure 14: aiNet applied to the 5-NLSC problem. (a) MST and its corresponding histogram (b). (c) Network dendrogram. (d) Final network structure with the centroids depicted, and the Voronoi diagram with relation to the centers of the clusters.

As discussed previously, when the final number of network antibodies is close to its minimal size ($m = 5$ in this case), the MST method might not be able to produce an accurate cluster separation and the network dendrogram might serve as an alternative. This is clear in this example, where by looking at Figures 14(a) and (b) we can not conclude anything about the final number of network clusters. In this case, the aiNet dendrogram (Figure 14(c)) served the purpose of correctly determining the number and members of each cluster. Figure 14(d) presents the final network configuration, the centroids of the clusters and the Voronoi diagram plotted with relation to the centroids of the clusters.

Table 5 shows the membership values for each network cell ($c_i, i = 1, \dots, 8$) with relation to the five clusters ($v_j, j = 1, \dots, 5$), and Figure 15 depicts the Voronoi diagram of Figure 14(d) together with the data set of Figure 9(c).

Table 5: Membership values for each cell $c_i, i = 1, \dots, 8$, with relation to each cluster centroid $v_j, j = 1, \dots, 5$.

	c_1	c_2	c_3	c_4	c_5	c_6	c_7	c_8
v_1	1.00	0.50	1.00	0.50	0.50	0.86	1.00	0.63
v_2	0.80	0.50	1.00	0.50	0.50	1.00	0.97	1.00
v_3	0.65	0.50	0.55	0.50	1.00	0.89	0.94	0.50
v_4	0.71	0.50	1.00	1.00	0.50	0.50	0.59	0.72
v_5	0.50	1.00	0.50	0.50	0.50	0.63	0.50	0.80

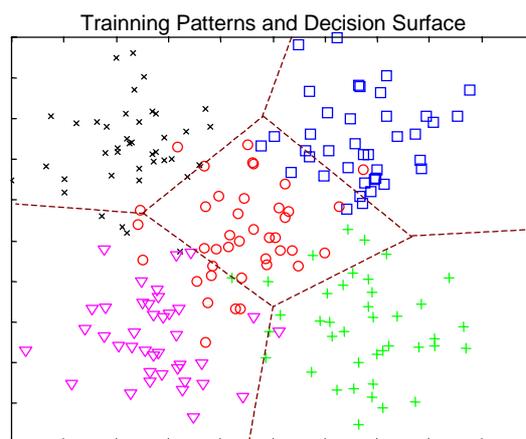


Figure 15: Decision surface (Voronoi diagram), taken from Figure 14(d), for the data set of Figure 9(c).

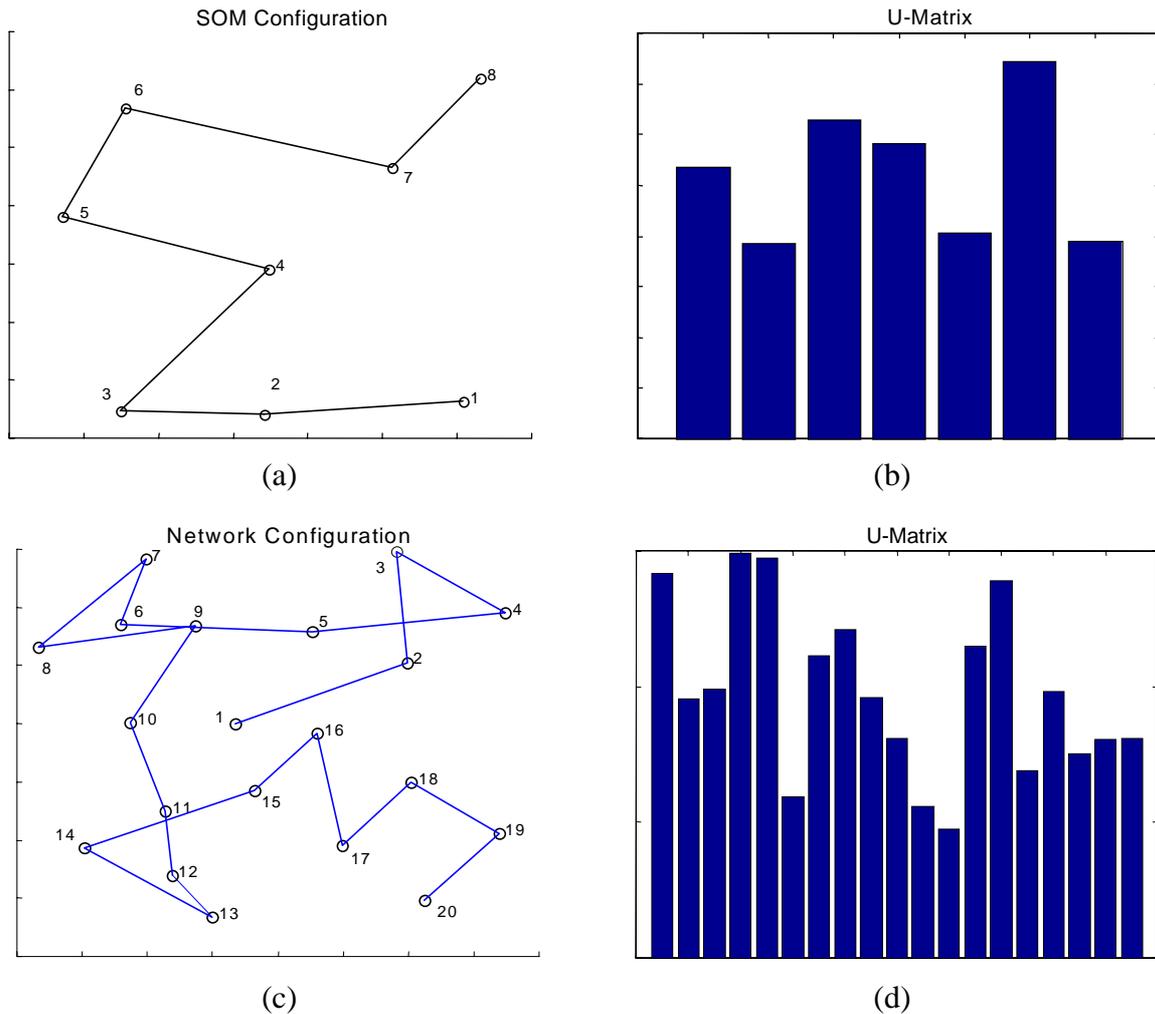


Figure 16: Results of the SOM to the 5-NLSC problem. (a) Network configuration and neighborhood, $m = 8$. (b) U-matrix, $m = 8$. (c) Network configuration and neighborhood, $m = 20$. (d) U-matrix, $m = 20$.

We used a SOM to solve the 5-NLSC problem with the same parameters as used in all the other examples. The number of output units was chosen to be $m = 8$ and $m = 20$, and the results are presented in Figure 16. Note that, from the U-matrix, we can not infer anything about the number of clusters in the resultant SOM.

Sensitivity Analysis

To apply the aiNet to any problem, a number of parameters have to be defined by the user, as can be seen from Table 3. In this section, we intend to discuss and analyze how sensitive the aiNet is to some of these user-defined parameters. In particular, it will be studied the influence of the parameters σ_s , σ_d , n and ζ in the convergence speed, final network size and recognition accuracy.

Figure 17 shows the trade-off between the suppression threshold σ_s and the final number N of output cells in the aiNet for the SPIR and CHAINLINK problems. As discussed previously, σ_s controls the final network size and is responsible for the network plasticity. Larger values for σ_s indicate more generalist antibodies, while smaller values result in highly specific antibodies. This parameter is critical, because the choice of a high value for σ_s might yield a misleading clustering. For the problems tested, the limiting values for σ_s that lead to correct results are $\sigma_s = 0.08$, $\sigma_s = 0.2$ and $\sigma_s = 0.2$, respectively. Higher values resulted in wrong clustering for some trials.

The pruning threshold (σ_d) is responsible for eliminating antibodies with low antigenic affinity. Without loss of generality, if we consider the illustrative problem presented in Section 5, we can evaluate the relevance of this parameter for the aiNet learning. Table 5 shows the amount of antibodies pruned from the network at the first generation. The results presented were taken from 10 different runs. In all runs, this parameter pruned network antibodies only at the first generation, and in the following generations no antibody was pruned by σ_d . This can be explained by the fact that the initial population of antibodies is randomly generated, but after the first generation, some of these antibodies were already selected, reproduced and matured to recognize the antigens (input patterns). Hence, we can conclude that the selection pressure and learning imposed by the algorithm are strong enough to properly guide the initial network towards a reasonable representation of the antigens in a single generation.

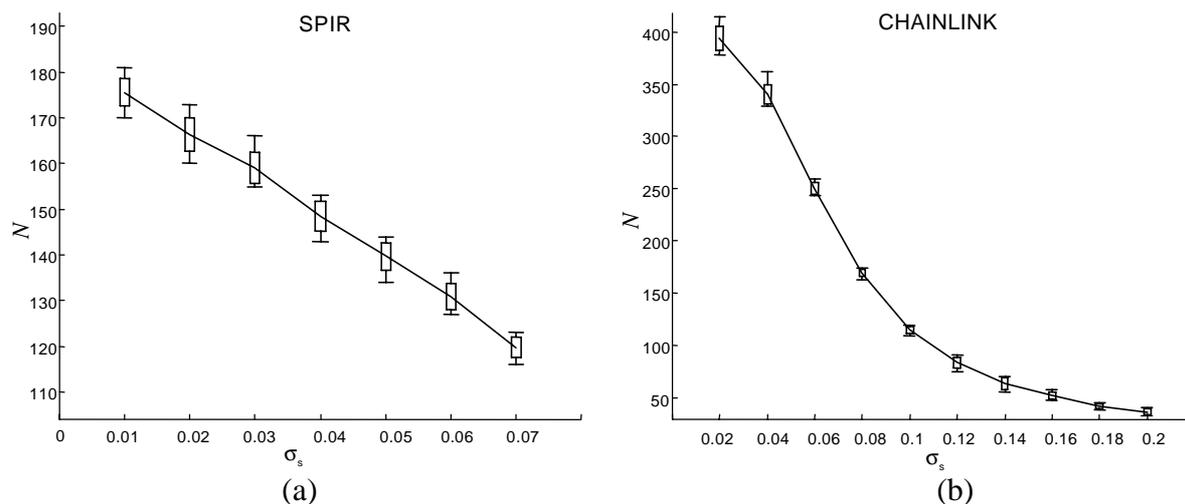


Figure 17: Trade-off between the final number of output units (N) and the suppression threshold (σ_s). The results are the maximum, minimum, mean and standard deviation taken over 10 runs.

Table 6: Number of antibodies pruned (Np) from the aiNet, at the first generation, for problem 5-LSC along ten runs.

<i>Run</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>Average</i>
Np	74	87	73	60	81	54	55	58	62	77	68.1±11.7

It is known that immune recognition is performed by a complementary Ag-Ab match. On the other hand, if we suppose that the aiNet main goal is to reproduce (build internal images of) the antigens to be recognized, instead of creating complementary antibodies, it is possible to define as the stopping criterion an average distance between the aiNet antibodies and the antigens and try to minimize this distance. In Steps 1.1.2 and 1.1.6 we minimize the Ag-Ab distance in order to maximize their affinity.

To properly study the aiNet sensitivity with relation to parameters n and ζ , it was chosen a value for the suppression threshold that would lead to a final network with approximately 50 antibodies. Based on this idea, one can test the aiNet potential to appropriately learn the antigens by simply defining as the stopping criterion (SC) a small value for the Ag-Ab average distance (10^{-2} , for example).

While evaluating the aiNet sensitivity to n , the following parameters were chosen: $\sigma_s = 0.01$, $\sigma_d = 1.0$, $n = 1..10$, $\zeta = 10\%$ and $SC = 10^{-2}$. Figure 18(a) depicts the trade-off between n and N (final network size), and Figure 18(b) illustrates the trade-off between n and the final number of generations for convergence. Note that the larger n , the larger the network size N , indicating that n has a direct influence on the network plasticity (see Figure 18(a)). On the other hand, from Figure 18(b) we can conclude that the larger n , the smaller the number of generations required for convergence (learning). The results presented are the maximum, minimum and mean taken over 10 runs.

Finally, to study the aiNet sensitivity to ζ , we chose the parameters $\sigma_s = 0.01$, $\sigma_d = 1.0$, $n = 4$, and $SC = 10^{-2}$; ζ was varied from 2% to 24% with steps of size 2%. Figure 19 shows the trade-off between ζ , N and the final number of generations for convergence (mean value taken over 10 runs). From this picture we can notice that ζ does not have a great influence on the final network size, but larger values of ζ imply in slower convergence.

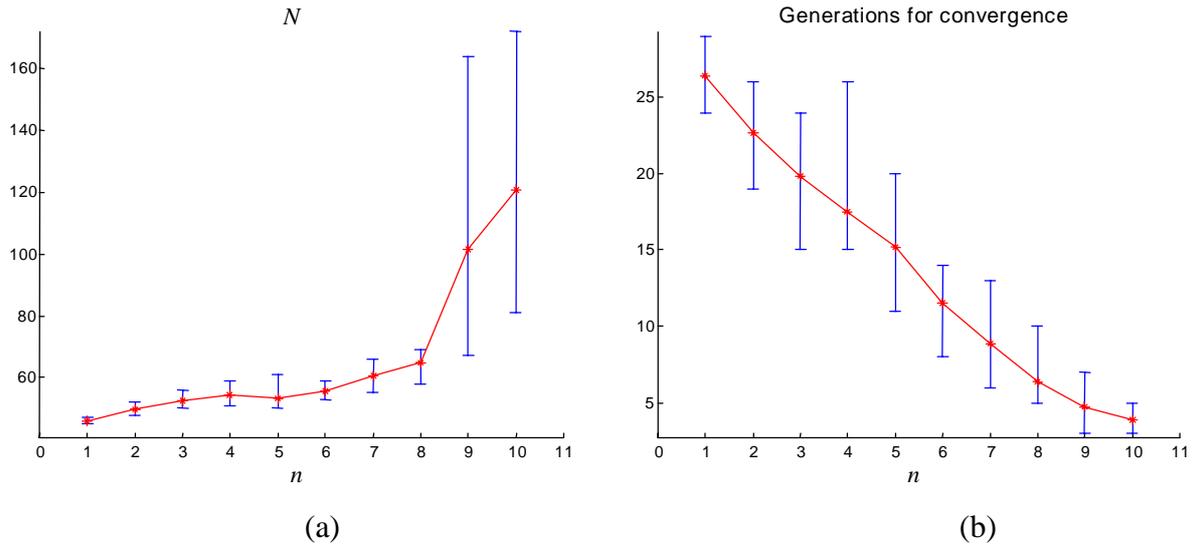


Figure 18: aiNet sensitivity to the number n of highest affinity cells to be selected for the next generation, Maximum, minimum and mean taken over ten runs. (a) Trade-off $n \times N$. (b) Trade-off $n \times gen$.

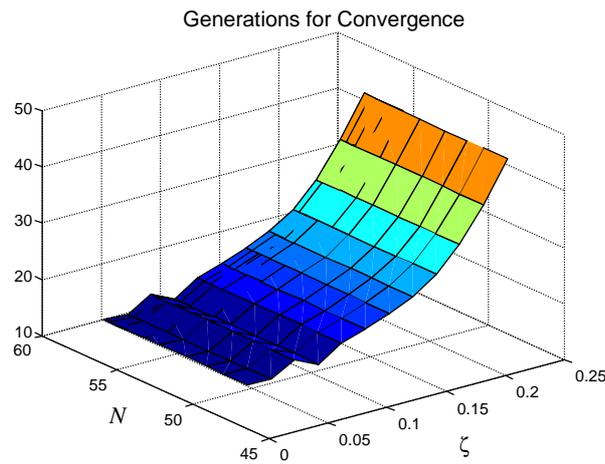


Figure 19: Trade-off between ζ , N and the number of generations for convergence (average over ten runs).

About the Number of Clusters and Network Parameters

The definition of a suppression threshold (σ_s) parameter is crucial in the determination of the final network size and consequently the number and shapes of the final cluster generated by the minimal spanning tree. This parameter has been determined in an ad hoc fashion, and as a further extension of this model we suggest the co-evolution of σ_s together with the network antibodies.

The amount of highest affinity antibodies to be selected for reproduction (n) also demonstrated to be decisive for the final network size. Nevertheless, the authors kept most of the parameters fixed for all problems, as can be seen from Table 1. In the most computationally intensive problems (SPIR and CHAINLINK), ζ was set to 10% in order to increase the learning speed. The only parameter that seemed to be really critical for the network clustering was the suppression threshold. It is also important to mention that in all the problems tested, the network demonstrated to be insensitive to the initial antibody repertoire, i.e., initial conditions.

In the SOM network case, the a priori definition of the network size may impose a network architecture not capable of correctly mapping the input data into the output nodes. Several models have been proposed to overcome this drawback (Fritzke, 1994; Cho, 1997, de Castro & Von Zuben, 1999b). If the clusters are non-linearly separable, but there is still a small distance between them (e.g., problems SPIR and CHAINLINK), the aiNet demonstrated to be capable of solving these problems, while the SOM network failed.

Concluding Remarks

In this chapter, an artificial immune network model, named aiNet, was proposed to solve data clustering problems. The resulting learning algorithm was formally described, and related to other connectionist and evolutionary models. In addition, the aiNet was applied to several benchmark problems and the obtained results compared to those of the Kohonen self-organizing neural network. As there is a great amount of user-defined parameters associated with the aiNet training, a sensitivity analysis was also performed.

The general purposes of the aiNet are: the automation of knowledge discovery, the mining of redundant data and the automatic clustering partition. This way, we can make use of antibodies and input patterns to be recognized (antigens) of the same dimension. One of the main reasons to take this decision, is that the aiNet can maintain the original topology of the classes, which is usually lost when it is promoted a dimensionality reduction.

On the one hand the two main goals of the SOM are to reduce data dimensionality and to preserve the metric and topological relationships of the input patterns (Kohonen, 1982). On the other hand, the aiNet reduces data redundancy, not dimensionality, and allows the reconstruction of the metric and topological relationships. Due to the possibility of reproducing the topological relationships, similar information (based upon a distance metric)

are mapped into closer antibodies, eventually the same one, characterizing the quantization and clustering of the input space.

By the time the immune network theory was proposed, the selective view of immune recognition (clonal selection principle) was well established and accepted. This new paradigm was in conflict with the selective theory, and network models did not take into account a clonal selection pattern of antigenic response. The network model presented in this chapter is different from the existing ones in the sense that it is discrete, instead of continuous, and it brings together the two originally conflicting theories: clonal selection and immune network. Nevertheless, the aiNet model takes into account the same processes covered by most of the dynamical models found in the literature (Equation (1)), but does not aim at directly mimicking any immune phenomenon.

In the aiNet model, clonal selection controls the amount and shapes of the network antibodies (its dynamics and metadynamics), while hierarchical and graph-theoretical clustering techniques are used to define the final network structure. The learning algorithm is generic, but the resultant networks are problem dependent, i.e., the set of patterns (antigens) to be recognized will guide the search for the network structure and shape of clusters. As its main drawbacks, we can mention its high number of user-defined parameters and its high computational cost per iteration, $O(m^2)$, with relation to the number, m , of memory antibodies.

As can be seen from Figures 12(d) and 14(d), the resulting aiNet clusters present a very peculiar spatial distribution. If we tried to represent these clusters by their respective centers of mass to perform the aiNet fuzzy clustering, the membership value of most of the antibodies would be incorrect, once the centroids are not representative of the real distribution of the classes.

As possible extensions and future trends we can stress the application of the aiNet to real-world benchmark problems of dimension $L > 3$, its application to the n-TSP (n-route travelling salesman) problem, the treatment of feasibility in the shape-space and its possible hybridization with local search techniques. In addition, the aiNet can be augmented to take into account adaptive parameters, aiming at reducing the amount of user defined parameters.

Acknowledgements

Leandro Nunes de Castro would like to thank FAPESP (Proc. n. 98/11333-9) for the financial support. Fernando Von Zuben would like to thank FAPESP (Proc. n. 98/09939-6) and CNPq (Proc. n. 300910/96-7) for their financial support.

References

- Ada, G. L., & Nossal, G. (1987). The Clonal Selection Theory, Scientific American, 257(2), 50-57.
- Bezdek, J. C. & Pal, S. K. (1992). Fuzzy Models for Pattern Recognition: Methods that Search for Structures in Data, New York, IEEE.
- Bonna, C. A., & Kohler, H. (Eds.). (1983). Immune Networks, Annals of the New York Academy of Sciences, 418.
- Burnet, F. M. (1978). Clonal Selection and After, In G. I. Bell, A. S. Perelson, & G. H. Pimbley Jr. (Eds.), Theoretical Immunology (pp. 63-85). Marcel Dekker Inc.
- Carrol, J. D. (1995). 'Minimax Length Links' of a Dissimilarity Matrix and Minimum Spanning Trees, Psychometrika, 60(3), 371-374.
- Cho S.-B. (1997). Self-Organizing Map with Dynamical Node Splitting: Application to Handwritten Digit Recognition, Neural Computation, 9, 1345-1355.
- Dasgupta, D. (Ed.). (1999c). Artificial Immune Systems and Their Applications, Springer-Verlag.
- De Castro, L. N., & Von Zuben, F. J. (2000a). The Clonal Selection Algorithm with Engineering Applications, In Workshop Proceedings of the GECCO 2000, 36-37. [On-Line]. Available: <http://www.dca.fee.unicamp.br/~lnunes/immune.html>
- De Castro, L. N., & Von Zuben, F. J. (1999a). Artificial Immune Systems: Part I – Basic Theory and Applications, (Tech. Rep. – RT DCA 01/99). Campinas, SP: State University of Campinas, Brasil. [On-Line]. Available: ftp://ftp.dca.fee.unicamp.br/pub/docs/vonzuben/tr_dca/trdca0199.pdf
- De Castro, L. N., & Von Zuben, F. J. (1999b). An Improving Pruning Technique with Restart for the Kohonen Self-Organizing Feature Map, In Proceedings of International Joint Conference on Neural Networks, 3 (pp. 1916-1919). Washington D.C., USA.
- Everitt, B. (1993). Cluster Analysis, Heinemann Educational Books.
- Farmer, J. D., Packard, N. H., & Perelson, A. S. (1986). The Immune System, Adaptation, and Machine Learning, Physica 22D, 187-204.
- Fritzke B. (1994). Growing Cell Structures – A Self-Organizing Network for Unsupervised and Supervised Learning, Neural Networks, 7(9), 1441-1460.

- George, A. J. T., & Gray, D. (1999). Receptor Editing during Affinity Maturation, Immunology Today, 20(4), p. 196.
- Hajela, P., & Yoo, J. S. (1999). Immune Network Modelling in Design Optimization, In D. Corne, M. Dorigo, & F. Glover (Eds.). New Ideas in Optimization (pp. 203-215). McGraw Hill, London.
- Hart, E., & Ross, P. (1999). The Evolution and Analysis of a Potential Antibody Library for Use in Job-Shop Scheduling, In D. Corne, M. Dorigo, & F. Glover (Eds.). New Ideas in Optimization (pp. 185-202). McGraw Hill, London.
- Hartigan, J. A. (1967). Representations of Similarity Matrices by Trees, Journal of the American Statistical Association, 62, 1440-1158.
- Haykin S. (1999). Neural Networks – A Comprehensive Foundation (2nd ed.). Prentice Hall.
- Holland, J. H. (1998). Adaptation in Natural and Artificial Systems (5th ed.). MIT Press.
- Hubert, L., Arabie, P., & Meulman, J. (1998). Graph-Theoretic representations for Proximity Matrices Through Strongly-Anti-Robinson or Circular Strongly-Anti-Robinson Matrices, Psychometrika, 63(4), 341-358.
- Hunt, J. E., & Cooke, D. E. (1996). Learning Using an Artificial Immune System, Journal of Network and Computer Applications, 19, 189-212.
- Janeway, C. A., P. Travers, Walport, M., & Capra, J. D. (2000). Immunobiology The Immune System in Health and Disease (4th ed.). Artes Médicas (in Portuguese).
- Jerne, N. K. (1974a). Towards a Network Theory of the Immune System, Ann. Immunol. (Inst. Pasteur) 125C, 373-389.
- Jerne, N. K. (1974b). Clonal Selection in a Lymphocyte Network. In G. M. Edelman (Ed.). Cellular Selection and Regulation in the Immune Response (p. 39). Raven Press, New York.
- Johnsonbaugh, R. (1997). Discrete Mathematics (4th ed.). Prentice Hall.
- Kepler, T. B., & Perelson, A. S. (1993). Somatic Hypermutation in B Cells: An Optimal Control Treatment, Journal of Theoretical Biology, 164, 37-64.
- Kohonen T. (1982). Self-Organized Formation of Topologically Correct Feature Maps, Biological Cybernetics, 43, 59-69.
- Lapointe, F-J., & Legendre, P. (1995). Comparison Tests for Dendrograms: A Comparative Evaluation, Journal of Classification, 12, 265-282.
- Lapointe, F-J., & Legendre, P. (1991). The Generation of Random Ultrametric Matrices Representing Dendrograms, Journal of Classification, 8, 177-200.
- Leclerc, B. (1995). Minimum Spanning Trees for Tree Metrics: Abridgements and Adjustments, Journal of Classification, 12, 207-241.
- Milligan, G. W., & Cooper, M. C. (1985). An Examination of Procedures for Determining the Number of Clusters in a Data Set, Psychometrika, 50(2). 159-179.

Oprea, M. (1999). Antibody Repertoires and Pathogen Recognition: The Role of Germline Diversity and Somatic Hypermutation (Ph.D. Dissertation, University of New Mexico, Albuquerque, New Mexico, USA).

Perelson, A. S. (1989). Immune Network Theory, Immunological Review, 110, 5-36.

Perelson, A. S., & Oster, G. F. (1979). Theoretical Studies of Clonal Selection: Minimal Antibody Repertoire Size and Reliability of Self-Nonself Discrimination, Journal of Theoretical Biology, 81, 645-670.

Prim, R. C. (1957). Shortest Connection Networks and Some Generalizations, Bell System Technology Journal, 1389-1401.

Sutton, R. S., & Barto, A. G. (1998). Reinforcement Learning an Introduction, A Bradford Book.

Timmis, J. I. (2000). Artificial Immune Systems: A Novel Data Analysis Technique Inspired by the Immune Network Theory (Ph.D. Dissertation, University of Wales, Aberystwyth, UK)

Ultsch, A. (1995). Self-Organizing Neural Networks Perform Different from Statistical k-means, Gesellschaft für Klassifikation.

Varela, F. J., & Coutinho, A. (1991). Second Generation Immune Networks, Immunology Today, 12(5), 159-166.

Zahn, C. T. (1971). Graph-Theoretical Methods for Detecting and Describing Gestalt Clusters, IEEE Transactions on Computers, C-20(1), 68-86.