

# Process State Estimation in a Wastewater Biological Treatment

IVÁN MACHÓN and HILARIO LÓPEZ and ANTONIO ROBLES

Departamento de Ingeniería Eléctrica, Electrónica de Computadores y Sistemas

Universidad de Oviedo

Edificio Departamental 2. Zona Oeste. Campus de Viesques s/n. 33204 Gijón (Asturias)

SPAIN

*Abstract:* - Using clustering techniques for data classification is very common. In this paper a Self-Organizing Map model is used to carry out an estimation of the process state in a wastewater biological treatment using clustering algorithms and validation indexes. The estimation is used to improve the efficiency of the treatment plant.

*Key-Words:* - Self-organizing mapping, clustering, validation, wastewater, biological treatment, chemical oxygen demand.

## 1 Introduction

This work is part of the KNOWATER II project “Implementation of a Knowledge Based System for Control of Steelworks Waste Water Treatment Plant”, which is sponsored by ECSC and their agreement number is 7210-PR-234. The contractors are Centro Sviluppo Materiali S.p.A., Corus RT&D, Betrieb Forschung Institut (BFI) and Universidad de Oviedo. The main objective of the KNOWATER II project was the development of plant supervision techniques for implementation in wastewater treatment plants. The present work was focused on the estimation of the process state in a wastewater biological treatment by means of the development of a neural network model that was obtained using Self-Organizing Map (SOM) and clustering algorithms. Once the model has been validated, a software tool was developed to supervise a biological wastewater treatment at a coke wastewater treatment plant (CWTP) of Arcelor in Avilés (Spain). The estimation of the current process state is calculated. Thus, important on-line knowledge is obtained.

## 2 Coke Wastewater Treatment Plant

The coke wastewater treatment plant (CWTP) consists of three zones: ammonia stripping towers, homogenization tank and biological reactor. The instrumentation of the CWTP was selected to control some key process variables (pHs and dissolved oxygen concentration) and to achieve the process monitoring by means of the designed software tool and its built-in AI technique.

Firstly, the influent stream is treated applying vapor in order to remove the ammonia in the stripping towers. A pHmeter is installed and a PID, which

controls a caustic soda dosing pump, controls the pH. The set-point is fixed at 12.

The second stage of the performance plant is the homogenization tank. At this point, a second pHmeter is established for the purpose of controlling pH by means of a PID controller with a set-point equal to 7. Sulphuric acid is added in order to neutralize the pH. Finally, the wastewater is treated biologically in a Sequencing Batch Reactor (SBR). The dissolved oxygen concentration is controlled by a third PID. Air is pumped into the reactor and a valve is regulated. The set-point is 3-5 mgO<sub>2</sub>/l.

An initial off-line study of the process was done and the PID controller output of the oxygen closed loop was connected and registered as one of the process variables to train the SOM network [1].

## 3 SOM

Self-Organizing Map (SOM) is a useful tool for process supervision and was used to construct a model that can be used as a pattern. The SOM [2] consists of a regular lattice typically defined in a two dimensional space composed of several neurons placed in the nodes of the lattice. SOM training implies assigning a set of coordinates in the input data space (prototype vector) to each neuron. Thus, each neuron is represented by a prototype vector and a correspondence is established between the coordinates of each neuron in the input space (data set) and their coordinates in the 2D-lattice or output space.

The present application was carried out compiling the SOM toolbox version 2.0 [3] developed at the HUT (Helsinki University of Technology). The steps taken to analyze the data are outlined in a previous work [4]. Firstly, the most significant process variables are

selected. These variables are described in table 1. Secondly, the data were normalized to a zero mean value and a unitary variance to make SOM treat them in the same way. After normalizing the SOM network was trained with these variables using batch training algorithm. Once the SOM has converged, it stores the most relevant information about the process in its prototype vectors. The visualization process allows all this information to be displayed in several ways: Interneuron distance matrix (Umatrix) that shows in gray or color levels the mean distance of each unit to its closest neighbors; the component planes that display the value of a given input variable throughout the whole data set using gray or color levels in the 2D lattice; the best clustering structure that allows the main process zones to be visualized [5].

Table 1. Training Variables

Name	Description
OXYGEN	Dissolved oxygen concentration (mgO <sub>2</sub> /litre)
CONTROLLER OUTPUT	Output of the PID controller of the oxygen closed loop (0-100)
TEMPERATURE_SBR	Temperature in the SBR (C)

#### 4 SOM validation

According to the properties of the SOM, the trained neural network must achieve the topology preservation of the data. Therefore the neighborhood on the output space and in the input space must be similar. If two prototype vectors close to each other in the input space are mapped wide apart on the grid, this is signaled by the situation where two closest best matching neurons of an input vector are not adjacent neurons. This kind of fold is considered as an indication of the topographic error in the mapping and does not verify the SOM property about training data topology preservation where neighbor neurons of the output space correspond to similar values of the process variables, i. e., regions of the output space represent working zones of the process

The topographic error [6] can be calculated by equation (1) as the proportion of sample vectors for which two best matching neurons are not adjacent.  $N$  is the number of samples,  $x_k$  is the  $k$ th sample of the data set and  $u(x_k)$  is equal to 1 if the first and second best matching neurons of  $x_k$  are not adjacent neurons, otherwise zero.

$$e_t = \frac{1}{N} \sum_{k=1}^N u(x_k) \quad (1)$$

The results of this error measurement are very easy to interpret and are also directly comparable between different models and even mapping of different data sets. Moreover, the prototype vectors approximate to the data set trying to substitute a data vector for a prototype vector of the SOM. A consequence of this approach is the quantization error. Equation (2) is usually used to calculate the average quantization error over the whole data set.  $N$  is the number of samples,  $x_i$  is the  $i$ th data sample and  $m_b$  is the prototype vector of the best matching neuron for  $x_i$ .

$$e_q = \frac{1}{N} \sum_{i=1}^N \|x_i - m_b\| \quad (2)$$

The SOM toolbox uses equations (3) and (4) to determine the output space size. The number of neurons of the output space is determined by equation (3).  $M$  is the number of neurons and  $N$  is the number of samples of the training data.

$$M = 5 \cdot \sqrt{N} \quad (3)$$

On the other hand, the criterion of the utilized toolbox to determine the ratio between the number of rows  $n_1$  and the number of columns  $n_2$  of the 2D grid or output space is calculated according to equation (4). The ratio between sidelengths of the map is the square root of the ratio between the two biggest eigenvalues of the training data. The highest eigenvalue is  $e_1$  and the second highest is  $e_2$ .

$$\frac{n_1}{n_2} = \sqrt{\frac{e_1}{e_2}} \quad (4)$$

Five data sets, which correspond to the aerobic phase of the SBR, are available to carry out the validation of the model. Pattern 1 before filtering is depicted in Fig. 1, whereas pattern 1 after filtering is showed in Fig. 2. Each sample is the mean value during 8 minutes and 20 seconds for each process variable. The training variables are showed in table 1. The objective is to find out the model that minimizes the quantization and topographic errors from several neural networks which have been trained using each of these available patterns and, at the same time, for different map sizes. Thus, a specific data set and an optimum map size must be selected. The validation method can be summarized in the following steps [7]:

1) A data set or pattern  $p_i$  is chosen to train the network. The data are normalized to a distribution with zero mean value and unitary variance.

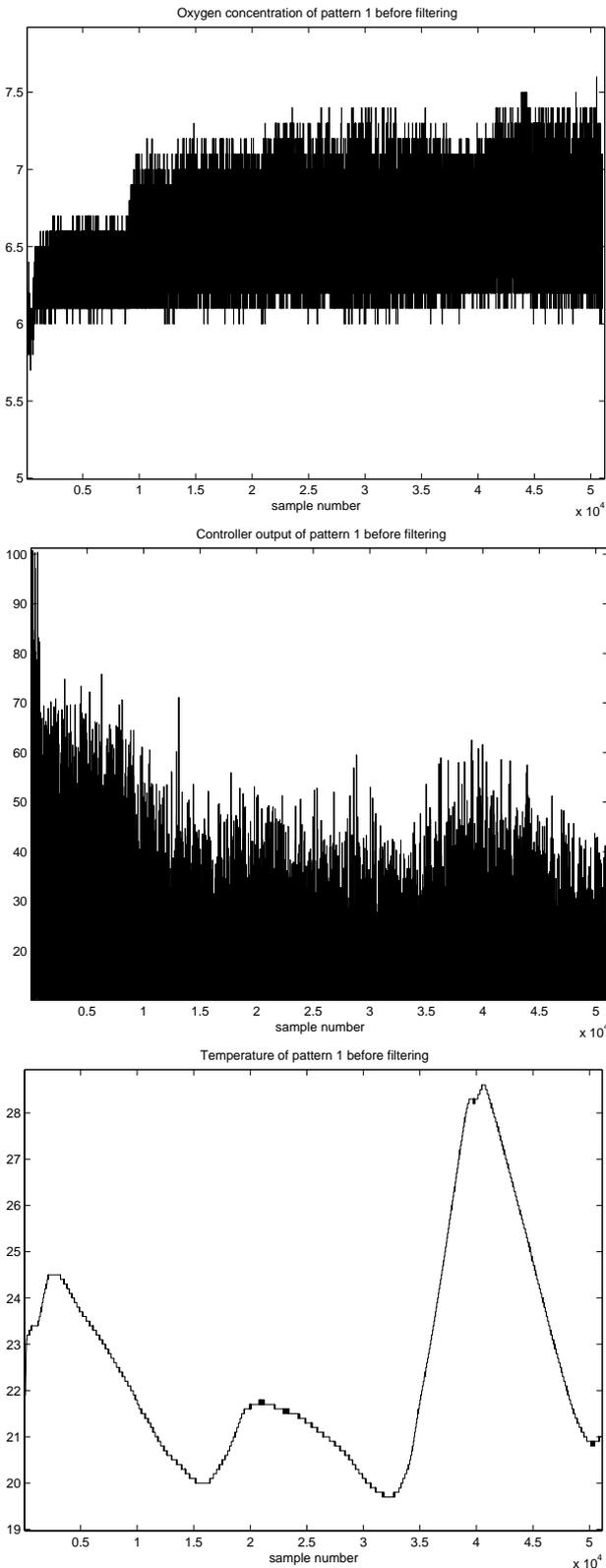


Fig. 1. Pattern 1 before filtering

- 2) Batch training is carried out on the SOM map whose sidelengths are calculated by means of equations (3) and (4) using pattern  $p_i$  as training data.
- 3) Once the trained model is obtained, the topographic and quantization errors are calculated for

the remaining patterns  $p_j$  which have not been used during the training. These patterns must also be previously normalized.

- 4) The size of this trained map is increased and reduced respecting the proportionality of its sidelengths (width and length). Once the size has been modified, the neural network is again trained using pattern  $p_i$ .

- 5) The third and fourth steps are repeated for different map sizes.

- 6) Steps 1 through 5 are repeated for the remaining patterns  $p_j$ , assuming each of these the role of pattern  $p_i$ .

Several map sizes have been trained using the five patterns. The mean values of the errors over the available patterns in function of the map number were calculated using each pattern as training data. The results for pattern 1 as training data are showed in Fig. 3. It can be seen that the larger the map size the lower the quantization error but the higher the topographic error. This is due to the neural network folds to reduce the quantization error. Moreover, the larger the map size the higher the computational cost. Therefore, there is compromise between the increase of the topographic error and the reduction of the quantization error. A curve, which represents the sum of both errors, has been added to the graphics. The model whose sidelengths have been calculated by means of equations (3) and (4) correspond to a horizontal axis value equal to 6 (map number equal to 6). The quantization error has been reduced and the topographic error has been incremented not very much. The final model to estimate the process state of the wastewater treatment was trained using pattern 1 because the values of the errors are the lowest for this pattern.

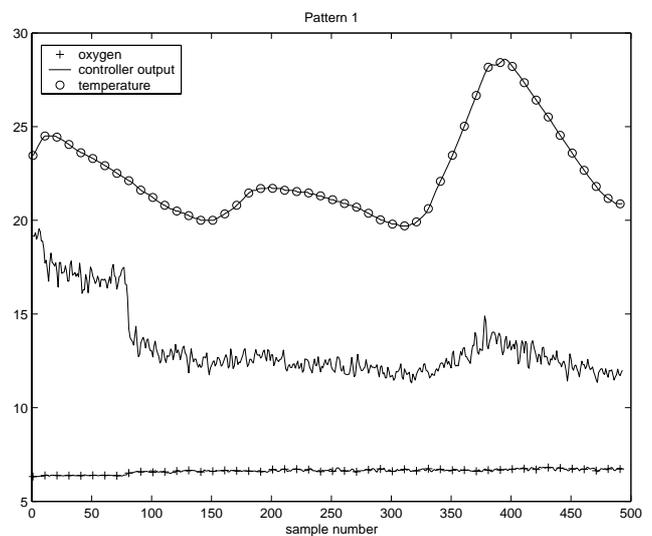


Fig. 2. Pattern 1 after filtering

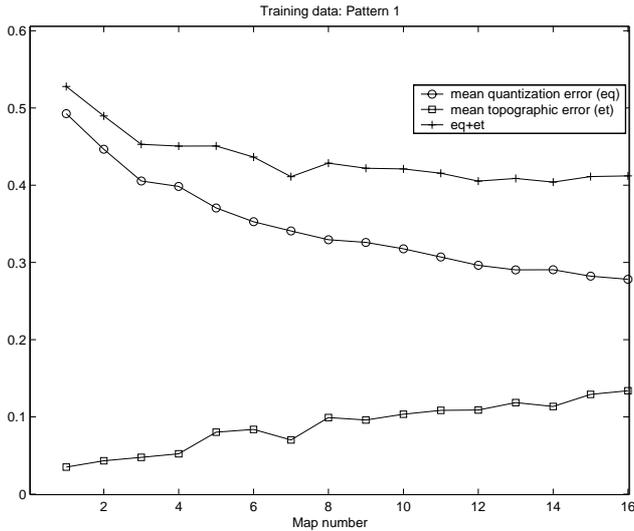


Fig. 3. Mean value of errors using pattern 1 as training data

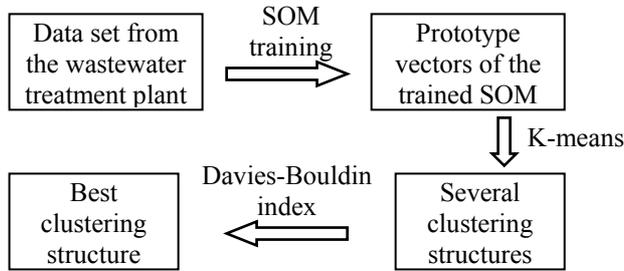


Fig. 4. Best Clustering Technique

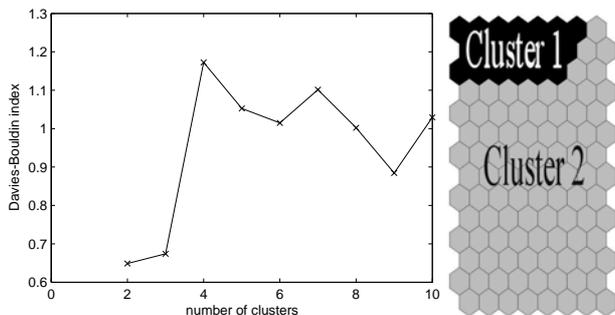


Fig. 5. Results of Davies-Bouldin Index and Best Clustering Structure

## 5 Clustering validation

The developed software tool carries out a clustering process which consists of a two-stage procedure [8] and is showed in Fig. 4. Firstly, the prototype vectors are obtained training the data of the aerobic phase using a SOM algorithm and then clustering them using a K-means algorithm, see [9]. Ten clustering structures were obtained varying the predefined number of clusters.

Finally, the best clustering structure between the ten structures, which have been obtained from the K-means algorithm, is selected using the Davies-Bouldin index [10]. This index searches the model that minimizes the within-cluster distance and maximizes the between-clusters distance and is calculated according equation 5, where  $s_i$  and  $s_j$  are the within-cluster distances of clusters  $i$  and  $j$ ,  $d_{ij}$  is the distance between clusters  $i$  and  $j$  and  $n_c$  is the number of clusters.

$$DB_{nc} = \frac{1}{n_c} \sum_{i=1}^{n_c} \max_{i=1, \dots, n_c, i \neq j} \frac{(s_i + s_j)}{d_{ij}} \quad (5)$$

The Davies-Bouldin index is suitable for evaluation of K-means partitioning, because it gives low values indicating good clustering results for spherical clusters. Fig. 5 shows the Davies-Bouldin index after being applied to the data from the aerobic stage of the treatment. The best clustering corresponds to a number of two clusters and has been projected onto the SOM. It is displayed in Fig. 5. Cluster 1 corresponds to the first hours of the aerobic treatment where the values of the controller output are high due to the high chemical oxygen demand (COD). During this period the biological activity is high and the toxic substances are eliminated by means of the cellular metabolism, whereas cluster 2 represents the data collected after this high biological activity where the values of the controller output are lower because the COD has decreased. If this is the state of the treatment plant, the biological treatment of the aerobic stages can be finished improving the capacity of the plant.

## 6 Process state estimation

The developed software is running in a PC station that is connected to a data acquisition interface by means of Ethernet connection and TCP/IP protocol (MODBUS protocol in particular). The proposed AI techniques are integrated into this application to achieve the process monitoring and the process state estimation. Data from the aerobic stage is collected on-line automatically to train a SOM network. The plant operator can visualize the latest SOM network that corresponds to the latest aerobic treatment cycle of the plant, viewing the correlation between the process variables and the data classification can be obtained. The estimation of the current process state is also calculated by the validated pattern (pattern 1). The training data set must only contain the samples of the aerobic stage and is determined by the mean value of the controller output because this signal can be

considered as a key variable to estimate the states of the treatment, see [1] and [11].

The results of the latest SOM network can be visualized in Fig. 6 and they correspond to the latest aerobic stage of the wastewater treatment. They are the U-matrix, the component planes and the best clustering structure. Each component plane shows the value of each neuron to estimate the data variable of the input space. It is useful to determine the several zones where the variable value is high or low and to observe any correlation or relationship between process variables.

The correlations between variables can be observed, for example, between the controller output and the oxygen concentration. Also there is a correlation between the controller output and the temperature in the reactor due to the higher the temperature the lower the dissolved oxygen concentration and the controller must compensate this effect. As mentioned above, the best clustering structure is composed of 2 clusters and is calculated by means of the Davies-Bouldin index. A cluster corresponds to HIGH COD and the other is the LOW COD.

The cycles of the biological treatment at the sequencing batch reactor can be clearly observed in Fig. 8 and Fig. 9. The higher values correspond to the anoxic stage when the controller output is saturated and equal to 100%. The rest of the data corresponds to the aerobic stage (including sedimentation).

An important aspect appears: the end-point of the aerobic reaction. This end-point detection can be used to finalize the aerobic stage and in this way the duration of the cycle is shorter increasing the operating capacity of the plant. The estimation of the time of the main activity of the treatment (aerobic phase end-point) achieves operating cost savings and increases the plant performance; see [12], [13], [14] and [15]. The duration of the cycle was initially 48-72 hours one year ago, see Fig. 8, and it has been reduced to 24 hours as is showed in Fig. 9. In this way the operating capacity of the plant has been increased by reducing the retention time. In Fig. 7 the process state is estimated projecting the current values onto a SOM network by means of standing out the best matching neuron from the rest of the neurons. This SOM network is used as a pattern and is previously stored and validated using the validation method explained above (pattern 1). The projection is carried out onto the component planes and the best clustering structure. Thus, important on-line knowledge is obtained and the end of the main biological activity can be identified.

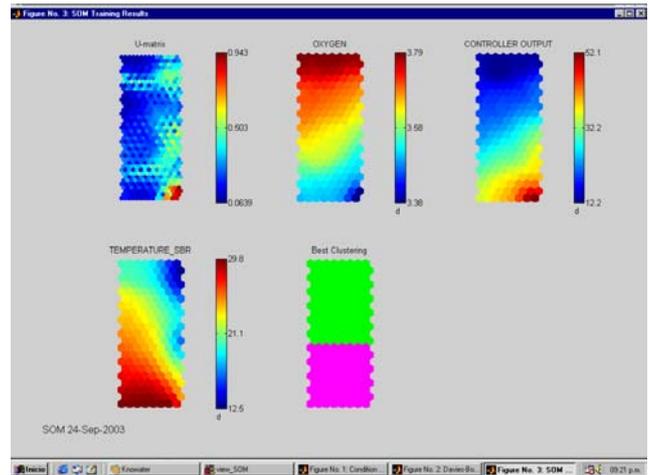


Fig. 6. SOM Training Results

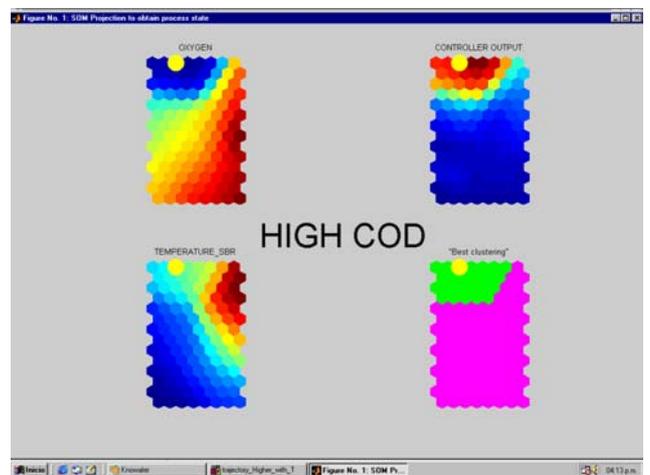


Fig. 7. Process State Estimation

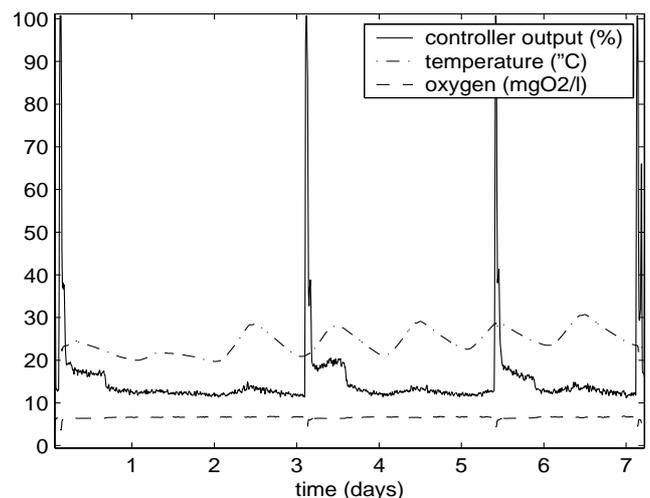


Fig. 8. Process Values of one year ago

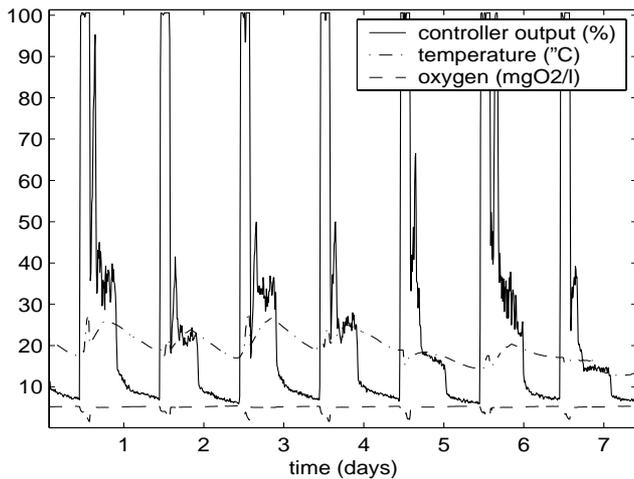


Fig. 9. Current Process Values

## 7 Conclusions

The results verify that the developed neural network model can achieve the estimation of the process state of the biological treatment. After model validation a software tool was developed to supervise a coke wastewater treatment plant (CWTP) and is a stand-alone application which is composed of the data acquisition system from the CWTP and the proposed AI technique. The data set of the aerobic stage is collected to train automatically a SOM network. The data classification is obtained using K-means algorithm as partitive clustering algorithm and Davies-Bouldin index for clustering validation. The estimation of the current process state can be assigned calculating the best matching neuron that corresponds to the current process values. The endpoint of the aerobic reaction can be detected by this AI technique. So, operating cost savings are achieved and the plant performance is increased. In this way, total retention time was reduced from 48-72 hours to 24 hours.

### References:

- [1] López H. and I. Machón. 2004a. "Biological wastewater treatment analysis using som and clustering algorithms," in *Proc. 12th Mediterranean Conference on Control and Automation*, Kusadasi.
- [2] Kohonen T. 2001. *Self-Organizing Maps*. New York: Springer-Verlag.
- [3] Vesanto J.; E. Alhoniemi; J. Himberg; K. Kiviluoto and J. Parviainen. 1999. "Self-organizing map for data mining in matlab: the som toolbox," *Simulation News Europe*, pp. 25–54.
- [4] López H.; I. Machón and S. Roces. 2003. "Waste treatment monitoring using self-organizing map and condition achievement maps," in *Proc. IFAC 5th Symposium on Intelligent Components and Instruments for Control Applications*, Aveiro.
- [5] López H. and I. Machón. 2004b. "Self-organizing map and clustering for wastewater treatment monitoring," *Engineering Applications of Artificial Intelligence*, vol. 17, no. 3, pp. 215–225.
- [6] Kiviluoto K. 1996. "Topology preservation in self-organizing maps," in *IEEE International Conference on Neural Networks*, vol. 1, pp. 294–299.
- [7] Machón I. and H. López. 2004. "An application for on-line control of a sequencing batch reactor," in *Proc. IFAC Workshop on Modelling and Control for Participatory Planning and Managing Water Systems*, Venice.
- [8] Vesanto J. and E. Alhoniemi. 2000 "Clustering of the self-organizing map," *IEEE Trans. Neural Networks*, vol. 11, no. 3, pp. 586–600.
- [9] McQueen J. 1967. "Some methods for classification and analysis of multivariate observations," in *5-th Berkeley Symposium on mathematics, Statistics and Probability*, no. 1, pp. 281–298.
- [10] Davies D. and D. Bouldin. 1979. "A cluster separation measure," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 1, no. 2, pp. 224–227.
- [11] López H. and I. Machón. 2004c. "An introduction to biological wastewater treatment explained by som and clustering algorithms," in *Proc. IEEE International Symposium on Industrial Electronics*, Ajaccio.
- [12] Andreottola G.; P. Foladori and M. Ragazzi. 2001. "On-line control of a sbr system for nitrogen removal from industrial wastewater," *Water Science and Technology*, vol. 43, no. 3, pp. 93–100.
- [13] Cho B.; S. Liaw; C. Chang; R. Yu; S. Yang and B. Chiou. 2001. "Development of a real-time control strategy with artificial neural network for automatic control of a continuous-flow sequencing batch reactor," *Water Science and Technology*, vol. 44, no. 1, pp. 95–104.
- [14] Paul E.; S. Plisson-Saune; M. Mauret and J. Cantet. 1998. "Process state evaluation of alternating oxic-anoxic activated sludge using orp, ph and do," *Water Science and Technology*, vol. 38, no. 3, pp. 299–306.
- [15] Yu R.; S. Liaw; B. Cho and S. Yang. 2001. "Dynamic control of a continuous-inflow sbr with time-varying influent loading," *Water Science and Technology*, vol. 43, no. 3, pp. 107–114.