# Estimation of conductivity in hydraulic affluents through self-organizing maps (SOM)

*Predicción de la conductividad en afluentes hídricas mediante mapas auto-organizativos (SOM)*

Wilson Ricardo López Sánchez[1], Cesar Andrey Perdomo Charry[2],
Jorge Enrique Rodríguez Rodríguez[3]

INFORMACIÓN DEL ARTICULO

ABSTRACT

This paper shows the use of self-organizing maps applied to water tributaries prediction. Currently, the environment conservation and the efficient water use, are pretty relevant issues. However, the water quality is not easy to measure, due to the specialized equipment for determining parameters such us: total coliforms, PH and dissolved oxygen. To measure the water conductivity, it was used a database with parameters of water quality, through the correlation of many involved parameters in water quality.

RESUMEN:

En este artículo se presenta el uso de mapas auto-organizativos aplicados a la predicción de afluentes hídricas. Actualmente, uno de los principales problemas es la preservación del medio ambiente y el uso eficiente del recurso hídrico. Sin embargo, la medición de la calidad del agua no es una tarea fácil, debido a lo especializado de los equipos para la determinación de parámetros tales como: coliformes totales, PH y oxígeno disuelto. Se empleó una base de datos con parámetros de calidad del agua para calcular la conductividad de la misma, a través de la correlación de varios parámetros involucrados en la calidad del agua.

[1] BSc. In Electronic engineering, MSc. In Information and Communications Sciences, Universidad Distrital Francisco José del Caldas, Bogotá, Colombia. Current position: Universidad Distrital Francisco José del Caldas, Bogotá, Colombia. E-mail: wrlopezs@correo.udistrital.edu.co. ORCID: https://orcid.org/0000-0002-1377-0667

[2] BSc. In Electronic engineering, Universidad Surcolombiana, Colombia. MSc. In Information and Communications Sciences, Universidad Distrital Francisco José del Caldas, Bogotá, Colombia. Current position: Professor at Universidad Distrital Francisco José del Caldas, Bogotá, Colombia. E-mail: cperdomo@udistrital.edu.co. ORCID: https://orcid.org/0000-0001-7310-4618.

[3] BSc. In Systems engineering, Fundación Universitaria De Boyaca, Colombia. MSc. In Systems engineering, Universidad Nacional de Colombia, Colombia. Current position: Professor at Universidad Distrital Francisco José del Caldas, Bogotá, Colombia. E-mail: jerodriguezr@udistrital.edu.co. ORCID: https://orcid.org/0000-0002-88820948.

## 1.    Introduction

According to an ONU (United Nations) report posted on 2016, the water is an essential component of the national and local economy, and it is necessary to create and maintain the workplaces in all of the economic sectors. Half of the global manpower is utilized in eight sector that rely on water and natural resources: agriculture, forests, fishing, and energy, production with massive use of resources, recycling, building and delivery. Sustainable management of water, water's infrastructures and the access to a safe and affordable supply of water and suitable sanitation services improve living standards, spread local economies and promote the creation of respectable workplaces and also a greater social inclusion. Sustainable management of water is also an essential power for the green growth and sustainable development.

Therefore, the culture and weather diversity as well as the political and financial activities have led to the searching of new and better procedures to carry out the right integrated management of water resources. Therefore, the processes of hydrological globalization and regionalization enables to refine the studies details until reaching a suitable management for basins.

Otherwise, an artificial neural network is an information analysis system consisting of a great amount of processing elements, connected to each other through communication channels commonly unidirectional, operating over internal or external local information. This kind of networks has been used to recognize and classify patterns, complete a signal depart from partial values or reconstruct the right pattern from a distorted one, among other areas; this way, the networks are an important tool to carry out the data clustering, due to the capability to shape complex and multidimensional data. To predict the tributaries conductivity, the self-organizing maps (SOM), or the KOHONEN model were selected as topology of the artificial neural networks; this topology was presented on 1982 as a system with similar behavior of the human brain. It is a model of artificial neural network with the capability to build maps in a similar way as the human brain does; in the last one there are neurons which are organized in many zones, so that the feedback information, caught through the sensory organs, is internally represented as two-dimensional maps.

The paper is organized in: problem, research methodology, approach to the basic theory of self-organizing maps, the implementation of the SOM model, results, discussion and conclusions.

## 2.    Problem

The computer systems' applications for the decision making and the natural events' modelling have significantly increased in the last few years, partly due to the progress in new models and the hardware capability to address complex problems [1]. Likewise, the obtained information from many natural events is analyzed with different computational techniques of machine learning and data mining [2]. The data extraction has been used in applications such as the searching of missing parameters and the parameters forecasting [3]. The data mining can be useful in the data acquisition or study the relation between parameters in order to obtain models and relations that cannot be obtained with conventional methods [4].

The mathematical models describing the behavior of physical phenomenon use a set of parameters and algorithms [5-6]. However, the most of the models must adapt to the parameters' alterations (climate changes). The models with the adaptive capacity to the parameters variation are so-called adaptive models [7-8]. The data extraction can determinate the correlation between features and estimate the lost values from a parameters set. [9] The data mining use in the natural events analysis plays an important role due to the results that have obtained with these methods [10], a sample of these studies is the fact of collecting models to research earthquakes [11], also the forecasting of water quality using satellite images [12- 14]. Many studies have been made for the parameters forecasting in the potable water sources, what has given as result the model generalization for different water sources. Some authors have created new adaptive models to predict natural events [15-16].

The two great limitations are the high cost associated with the data acquisition to measure the water quality and the missing of complete parameters through the time. The following paper shows a research with collected data in many rivers from Bogota city in different times of year, in order to assess the correlation between these data and the probability to make predictions about the water conductivity.

## 3.    Methodology

This section describes the type of scientific research used in the article along with the research method and the development methodology. The type of scientific research applied in this paper is descriptive-exploratory with an experimental approach. According to the formal research process, a hypothetical - deductive method was used in which a hypothesis was formulated, which through empirical validation was validated

through deductive reasoning. It was established, based on the experimentation, a mechanism for weighting the algorithm evaluation indicators in such a way that it was possible to evaluate said mechanism when changing the dataset.

The following tasks are defined to obtain the results, after applying the appropriate algorithms in feature selection:

1.      The basis for the development of this paper is based on the analysis and choice of algorithms for feature selection. In this stage, the case studies are of vital importance, the conclusions of this research and machine learning algorithms used in feature selection.

2.      Collection, integration, and data preprocessing. In this phase, the collection and integration of different datasets, data transformation; as the case may be, and cleaning in order to eliminate noise.

3.      Definition and application of tests of the algorithms used for feature selection. Based on the tests performed with the synthetic data.

4.      Evaluation of the imputation indicators were applied to apply the algorithms.

5.      Review of test results. The analysis of the data obtained in the allocation with the algorithms was performed. Likewise, the complexity of the algorithms was calculated, in order to determine their feasibility of implementation.

The instruments, on the basis of this research were developed, which are mainly cases of study and tests with synthetic data.

## 4.   Self-organizing maps – SOM

SOM was introduced by T. Kohonen on 1982 and it was a special kind of artificial neural networks of non-supervised learning which has been successfully implemented in data mining and in the Knowledge Discovery in Databases – KDD with a great variety of engineering
applications such as the pattern recognizing, image analysis, monitoring process and fails detection, and others [17].

This network or map finds out common features, regularities, correlations or categories in the input data, and integrates its intern structure of connections. The neurons must organize themselves according to inducements (data) from the outside [18-20].

These maps have the role of transforming patterns of arbitrary dimension as a response of the patterns of one or two fixing neurons dimensions which change by adapting themselves according to the characteristic input features. It can be claimed that the algorithms of selforganizing maps can be defined to display multidimensional data [20].

A SOM model is consisting of the following items [21]:

-       Neurons matrix: input layer (formed by $N$ neurons, one by each input variable) it receives and transmits the information from the outside to the output layer. The hidden layers (formed by M neurons) has the role to process the information and form the features map. Usually, they organize as a two-dimensional map which can be rectangular or hexagonal.

-       Relationship between neurons: there is proximity relationship between neurons which is the key to shape the map during the training stage.

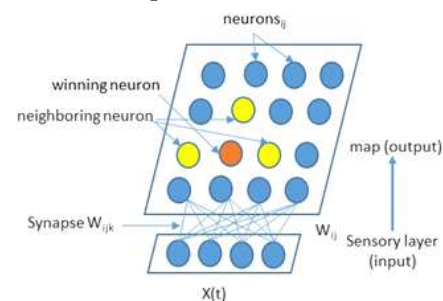In Figure 1, are shown the items composing an organizational map.



**Figure 1.** Application of self-organizing maps to the data display [21]

The SOM learning process is:

**Step 1.** A $x$ vector is randomly selected from the data set and the distance to the codebook's vectors is calculated, using, for instance, the Euclidean distance:

$$\|x - m_m\| = min\{\|x - m_j\|\}$$ (1)

**Step 2.** Once the closest vector is found, the other codebook's vectors is updated. The closest vectors and their neighborhoods (in a topological way) move close of $x$ vector in the data space.

The magnitude of such attraction is ruled by the learning rate. While the updating process is producing and new vectors are assigned to the map, the learning rate gradually decrease to zero. Along with this one, the proximity ratio also decreases. The updating rule for the reference vector *i* is the next one:

$$m_j(t+1) = \begin{cases} m_j(t+1) + \propto (t)\big(x(t) - m_j * t\big) \\ m_j(t) \end{cases} \quad (2)$$

Where $J \in N_c(t)$.

Steps 1 and step 2 repeat until the training is over. The number of training steps must be stablished a priori, to calculate the convergence rate in terms of proximity and the learning rate. Once the training is over, the map must be ordered in a topological way: **n** vectors topologically close are applied in *n* adjacent neurons or even the same neuron [18].

The network working is relatively simple when the input information $E_K = (e_1^k, \dots, e_n^k)$ is done, each one of the *M* neurons of the output layer is received through the feedback, which are the weight connections $W_{IJ}$. In the same way, the neurons receive the lateral connections inputs or outputs which influence will rely on the distance they are at.

In this way, the generated output by a output layer *j* in light of an input vector $E_K$ will be:

$$S_1(t+1) = f\left(\sum_{i=1}^{N} W_{ij} e_i^{(k)} + \sum_{p=1}^{M} Int_{pj} S_p^{(k)}\right) \quad (3)$$

Where $Int_{pj}$ is a function that represents the lateral influence of the *p* neuron over the *j* neuron. It is clear that it is a competitive network, due to, by submitting an input *Ek* the network evolve until a stable situation where an output layer turns on, the winner. So, the mathematical equation can get simplified to the following term:

$$S_1 = MIN\|E_k - W_j\| = \sqrt{\sum_{i=1}^{n}\left(e_i^{(k)} - W_{ij}\right)^2} \quad (4)$$

Where $\|E_k - W_j\|$ is a measure (for instance Euclidean distance) of the difference between the input $E_k = \left(e_1^{(k)}, \dots, e_n^{(k)}\right)$ and the weight vector $E_k = \left(e_1^{(k)}, \dots, e_n^{(k)}\right)$ of the conections between each one of the input layers and output layers *j*. In in the network during the learning process. It is aim to find the most similar data to the input data to enquire which neuron will get activated and where in the output two-dimensional space it will be [20].

The SOM advantages consist in they are able to preserve the topology of data space, they project highly dimensional data to a low dimensional scheme of representation and if they have the ability to find similarities in the data [17].

## 5. Application of the SOM neural network to estimation of conductivity

The database analysis was made according to KDD [16], this database was measured with 2177 samples of 131 parameters with the following layout: 122 station features along Bogota city. Taking into account Tunjuelo, Fucha, Salitre and Torca rivers, as well as another parameters such as: the sample number, the station ID, the river identification, data, type of simple outline and an outline about the weather with the starting hour and the sample of ending hour, it was taken during eight months. On the other hand, from the database, some registers of parameters regarding to the quality of water where they were registered and some other were deleted and how they used with the following parameters, which are shown in Table 1 as follows:

| Features | Units |
|---|---|
| pH | -- |
| Temperature | [°C] |
| Conductivity | [μS.cm-1] |
| Dissolved oxygen [OD Winkler] | [mg O2.L-1] |
| Fats and oils | [mg.L-1] |
| Total Coliforms | [NMP.100 mL-1] |
| Faecal Coliforms | [NMP.100 mL-1] |

**Table 1.** Parameters related to water quality. Source: own.

A quantification process for the feature analysis was necessary. A numerical value or ID was assigned to each station where the samples have remarks about the weather, the data and the moment they were taken.

The same process was made for the remarks of weather, for this date, few values between 1 and 8 were taken, in this way, 1 fits to the month of November 2009 and July 10, 2010 for the case of time of day value of 0 hours less than 12 noon and 1 for others. The patterns with missing values were deleted. After the purification the base has 825 data to do a statistical analysis: estimating the mean, the variance and the standard deviation. With the results from the statistical analysis, the outliers data are deleted, and with a box–plot graph (graph 1), the data with a high dispersion before the purification of the database are shown, after this, goodness of fit test were performed, it is plotted a histogram by features to notice if the data have any distribution function established. Finally, can be noticed the correlation between the features through the artificial neural networks.

Turning on different neurons in the SOM neural network determines if there is any relationship between the features and the importance of these ones, so, this

way, the most suitable features for the training of a neuronal neural network are selected, as of an algorithm LevenbergMarquardt which does the prediction of the water conductivity. Finally, considering the obtained results, many neural networks are trained by changing the number of neurons in the hidden layer and by comparing the estimating error; in this way, it is selected the most suitable network to determine the water conductivity in water tributaries, fore the present research instance. For the neural networks training, the data set was divided in three sub-sets: 70% for the network training, 15% for the validation test and 15% for the tests applied to the neural networks.

## 6. Results

Here below, in Table 2, the mean, the variance and the standard deviation of data related to water quality, are shown.

| Features | Mean | Standard deviation | Variance |
|---|---|---|---|
| pH | 7,33 | 7,1 | 0,31311753 |
| Temperature | 16,5 | 13,5 | 8,31266956 |
| Conductivity | 334 | 308 | 122719,063 |
| Dissolved oxygen [OD Winkler] | 0,8 | 0,1 | 5,65387026 |
| Fats and oils | 9 | 3,6 | 7980,47612 |
| Total Coliforms | 2400000 | 11000000 | 1,63E+19 |
| Faecal Coliforms | 900000 | 4300000 | 4,98E+18 |

**Table 2.** Mean, standard deviation and variance of data taken from the stations in water tributaries. Source: own.

Data with the highest dispersion in fats and oils, total coliforms and faecal coliforms features are deleted, in Figure 2. In which is submitted the box-plot charter for the three features.
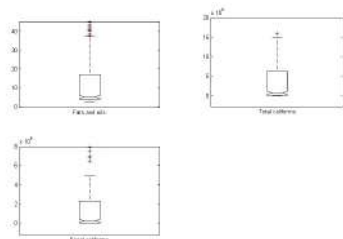


**Figure 2.** Box plot charter of fats and oils, total coliforms and faecal coliforms debugged features. Source: own.

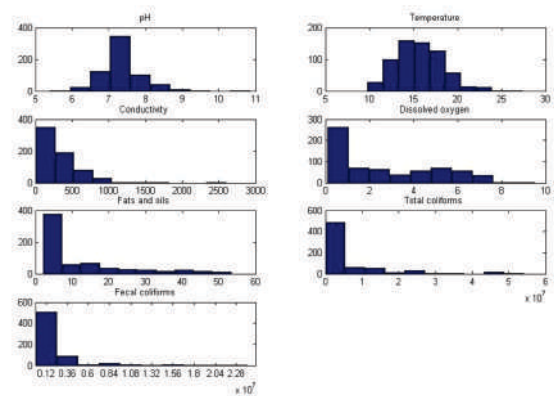In Figure 3 the histogram of each one of the features is submitted.



**Figure 3.** Features histogram. Source: own.

It is estimated the mean, standard deviation and variance for each one of the features (Table 3), after debugging the database.

| Features | Mean | Standard deviation | Variance |
|---|---|---|---|
| pH | 7,29 | 0,5397471 | 0,29132693 |
| Temperature | 15,6 | 2,71309025 | 7,36085868 |
| Conductivity | 228,5 | 268,071 | 71862,0611 |
| Dissolved oxygen [OD Winkler] | 1,8 | 2,40174399 | 5,76837417 |
| Fats and oils | 4,95 | 12,4506106 | 155,017706 |
| Total Coliforms | 770000 | 8937255,88 | 7,9875E+13 |
| Faecal Coliforms | 230000 | 3816388,03 | 1,4565E+13 |

**Table 3.** Mean, standard deviation and variance of debugged data. Source: own.

It is trained a SOM network of $10x10$ neurons having as an input each one of the features from the database to observe the activating level of each neuron and to look for the correlation between the features (Figure 4).

The neural networks are trained as of the algorithm Levenberg-Marquardt by varying the number of neurons in the hidden layer from 3 until 30 neurons, getting as a result that the network with lowest MSE (Mean squared error) is the pertinent to 10 neurons in the hidden layer, this experiment was made with 102 times.

In Figure 5, it is evident that the lowest MSE reached by the validation is between 5 and 10 times. Also, that as the number of times increases, so the test error and the validation increases, what means that the network is over-trained, therefore the network is trained again for a maximum of 15 times of training; the MSE of this network is shown in charter 5.
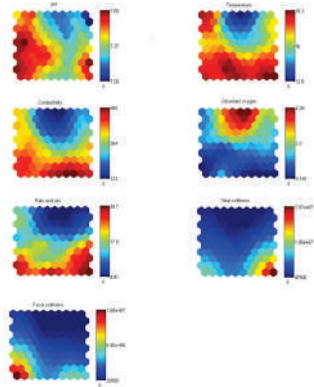
**Figure 4**. Neural activation in SOM network. Activation of neurons in SOM network of the seven features related to water quality. Source: own.
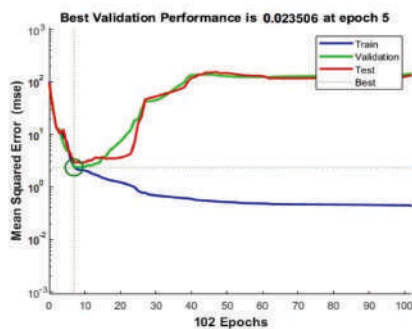


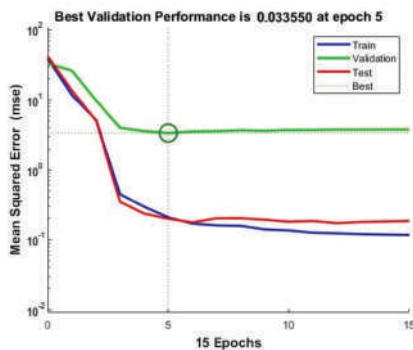**Figure 5.** Network efficiency for 102 times. Source: own.



**Figure 6.** Artificial neural network efficiency for 10 times. Source: own.

In Figure 7 is shown a comparison between the real conductivity data and the obtained data by the neural network, in this case, there is a sample of 15 data.
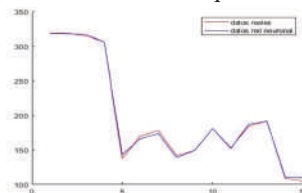


**Figure 7.** Comparison between the real conductivity data and the obtained data by the neural network. Source: own.

## 7. Discussion

The mean, the variance and the standard deviation of the data were calculated. From the features in table 2 can be seen that the conductivity, fats and oils, total coliforms and faecal coliforms have a high variance, which implicates a data dispersion. By comparing the mean and the standard deviation in table 2 it is observed that: the features of fats and oils, total coliforms and faecal coliforms is around $1 \times 10^{10}$, which states a huge value spectrum or a high dispersion.

The dispersion of the features fats and oils, total coliforms and faecal coliforms decreases by deleting the most disperse ones regarding each feature. In charter 1 is submitted a box-plot charter about the three features where can be seen the dispersion decrease for the three features. It is submitted a histogram of the features (chapter 2) except for the PH and temperature, it is evident a data concentration for values close to cero.

The obtained database after the cleaning and debugging has 12 features and 654 data. To determine the correlation of the seven related features to the water quality, it was applied the SOM network of 10x10 neurons, in charter 3 is shown that the activation of the first 5 features' neurons has a similar behavior. The temperature and conductivity show a similar behavior, by comparing the activation of the first rows neurons the feature PH has similarities to the last two features. The temperature and conductivity behavior regarding to the dissolved oxygen in water is inverse. By applying a mirror effect to any feature related to the Coliforms, their behavior is similar to the other one.

Finally, some artificial neural networks are trained, getting a better result with the model using 10 neurons in the hidden layer for the prediction of water conductivity considering some features as network input. This model is trained with 100 iterations, reaching stability before the tenth iteration, it is done a new network training, stopping the training when it gets to 10 iterations or to 95 validations; in charter 5 it can be seen the efficiency the network got with a optimal point of work for 5 iterations.

## 8. Conclusions

Data preprocessing enabled a better efficiency of SOM network to find out different correlations between the database features, without outliers which could generate wrong answers or noise in the neural network and finally, to decrease the problem magnitude.

The neural network was trained for 50 iterations avoiding an over training of the network that restricted the prediction by introducing new data.

For future studies it is desirable to do a data analysis by discriminating the different stations where samples were taken; to do a study about the data standardization so the networks training can give better results.

## References

[1] J. Refonaa, M. Lakshmi and V. Vivek, "Analysis and prediction of natural disaster using spatial data mining technique", International Conference on Circuit, Power and Computing Technologies, 2015. https://doi.org/10.1109/ICCPCT.2015.7159 379.

[2] C. Doña, N. Chang, V. Caselles, J. M. Sánchez, A. Camacho, J. Delegido and B. W. Vannah, "Integrated satellite data fusion and mining for monitoring lake water quality status of the Albufera de Valencia in Spain", *Journal of Environmental Management*, vol. 151, 2015, pp. 416-426. https://doi.org/10.1016/j.jenvman.2014.12.0 03.

[3] G. Ssali and T. Marwala, "Computational intelligence and decision trees for missing data estimation", IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence), 2008. https://doi.org/10.1109/IJCNN.2008.463379 0.

[4] M. K. Sohrabi and S. Akbari, "A comprehensive study on the effects of using data mining techniques to predict tie strength", *Computers in Human Behavior*, vol. 60, 2016, pp. 534-541. https://doi.org/10.1016/j.chb.2016.02.092.

[5] Z. Yu, F. Haghighat and B. Fung, "Advances and challenges in building engineering and data mining applications for energy-efficient communities", *Sustainable Cities and Society*, vol. 25, 2016, pp. 33-38. http://dx.doi.org/10.1016/j.scs.2015.12.001

[6] S. Chapra, "Surface Water Quality Modelling", Mc Graw Hill. Brown, L.C., and Barnwell, pp. 189, 1987.

[7] S. Chapra, G. Pelletier and H. Tao, "QUAL 2K: A Modeling Framework for Simulating River and Stream Water Quality", 2005. [Online]. Available at: https://www.hydroreform.org/node/3420.

[8] T. Garrido Reyes and J. Mendoza Crisosto, "Characterization of Dissolved Organic Matter in River Water by Conventional Methods and Direct Sample Analysis-Time of Flight-Mass Spectrometry", *Journal of Chemistry*, vol. 2016, 2016. http://dx.doi.org/10.1155/2016/1537370.

[9] F. B. Wagner, P. Borch Nielsen, R. Boe-Hansen and H. Albrechtsen, "Copper deficiency can limit nitrification in biological rapid sand filters for drinking water production", *Water Research*, vol. 95, 2016, pp. 280-288. http://dx.doi.org/10.1016/j.watres.2016.03. 025.

[10] T. Umair Ali, Z. Ahmed and D. Kim, "Estimation of N2O emission during wastewater nitrification with activated sludge: Effect of ammonium and nitrite concentration by regression analysis", *Journal of Industrial and Engineering Chemistry*, vol. 20, no. 4, 2014, pp. 2574-2579. http://dx.doi.org/10.1016/j.jiec.2013.10.04 2.

[11] A. Erturk, M. Gurel, A. Ekdal, C. Tavsan, A. Ugurluoglu, D. Zafer Seker, A. Tanik and I. Ozturk, "Water quality assessment and meta model development in Melen watershed – Turkey", *Journal of Environmental Management*, vol. 91, no. 7, 2010, pp. 1526-1545. https://doi.org/10.1016/j.jenvman.2010.02. 021.

[12] S. L. Nimmagadda and H. Dreher, "Ontology based data warehouse modeling and mining of earthquake data: prediction analysis along Eurasian-Australian continental plates," 5th IEEE International Conference on Industrial Informatics, Vienna, 2007, pp. 597-602. https://doi.org/10.1109/INDIN.2007.43848 25.

[13]   M. Bonansea, M. C. Rodriguez, L. Pinotti and S. Ferrero, "Using multi-temporal Landsat imagery and linear mixed models for assessing water quality parameters in Río Tercero reservoir (Argentina)", *Remote Sensing of Environment*, vol. 158, 2015, pp. 28-41. https://doi.org/10.1016/j.rse.2014.10.032.

[14]   E. Therese Harvey, S. Kratzer and P. Philipson, "Satellite-based water quality monitoring for improved spatial and temporal retrieval of chlorophyll-a in coastal waters", *Remote Sensing of Environment*, vol. 158, 2015, pp. 417-430. https://doi.org/10.1016/j.rse.2014.11.017.

[15]   X. Wang, J. Zhang and V. Babovic, "Improving real-time forecasting of water quality indicators with combination of process-based models and data assimilation technique", *Ecological Indicators*, vol. 66, 2016, pp. 428-439. http://dx.doi.org/10.1016/j.ecolind.2016.02.016.

[16]   G. Piatetsky-Shapiro, "Advances in knowledge discovery and data mining", USA: American Association for Artificial Intelligence Menlo Park, 1996.

[17]   R. Salas, "Mapas Autoorganizativas de Kohonen (SOM)", 2004. [Online]. Available at: http://www.inf.utfsm.cl/~rsalas/Pagina_Investigacion/docs/Apuntes/Redes%20SOM.pdf.

[18]   J. Marín, "Mapas Auto-organizativos de Kohonen SOM", 2008. [Online]. Available at: http://halweb.uc3m.es/esp/Personal/personas/jmmarin/esp/DM/tema5dm.pdf.

[19]   A. Díaz, "Redes neuronales no supervisadas con topología dinámica de la segmentación de imágenes a color", 2010. [Online]. Available at. http://agt.cie.uma.es/~adiaz/Publications/PFCAntonioDiazRamos.pdf.

[20]   J. Rodríguez, "Fundamentos de la minería de datos", Bogotá D.C, Colombia: Universidad Distrital Francisco José de Caldas, 2010.

[21]   F. Serrano, "Aplicación de mapas autoorganizados (SOM) a la visualización de datos", 2009. [Online]. Available at: https://docplayer.es/9172902-Aplicacion-de-mapas-