

FUZZY INFERENCE SYSTEM TO EVALUATE THE QUALITY OF GROUNDWATER IN MEXICAN WATER BODIES

LAURA I GARAY-JIMENEZ¹, ULISES MONTOYA CANALES², PILAR GOMEZ MIRANDA², ANA JUDITH MARMOLEJO RODRIGUEZ³, BLANCA TOVAR CORONA¹

¹INSTITUTO POLITÉCNICO NACIONAL, UNIDAD PROFESIONAL INTERDISCIPLINARIA EN INGENIERÍA Y TECNOLOGÍAS AVANZADAS, AV. INSTITUTO POLITÉCNICO NACIONAL 2580, LA LAGUNA TICOMÁN, GUSTAVO A. MADERO, 07340 CIUDAD DE MÉXICO, MEXICO.

²INSTITUTO POLITÉCNICO NACIONAL, UNIDAD PROFESIONAL INTERDISCIPLINARIA DE INGENIERÍA Y CIENCIAS SOCIALES Y ADMINISTRATIVAS, AV. TÉ 950, IZTACALCO, 08400, CIUDAD DE MÉXICO, MÉXICO. ³INSTITUTO POLITÉCNICO NACIONAL, CENTRO INTERDISCIPLINARIO DE CIENCIAS MARINAS, AVENIDA IPN, S/N COLONIA PLAYA PALO DE SANTA RITA, C.P. 23096, LA PAZ, BAJA CALIFORNIA SUR, MÉXICO.

EMAIL: ¹ lgaray@ipn.mx, ORCID ID: ¹ https://orcid.org/0000-0001-9478-4835 EMAIL: ²umontoya1800@egresado.ipn.mx, ORCID ID: ² https://orcid.org/0009-0003-7454-6356 EMAIL: ³ amarmole@ipn.mx, ORCID ID: ³ https://orcid.org/0000-0002-8913-2522

Corresponding Author*: Laura I Garay-Jiménez.

ABSTRACT:

Introduction: This paper presents the development of a fuzzy inference system to assess groundwater quality using data from 2012-2021 from the national water information system (SINA). The objective was to create a simplified semaphore based on fuzzy logic, classifying groundwater into CONAGUA's three traditional categories (green, yellow, red) while incorporating a degree of membership for each condition. Methodology: CONAGUA classifies water quality using 14 crisp variables, but we employed eight fuzzy variables as inputs to a Mamdani inference system. Results: Our fuzzy system achieved 84% similarity with CONAGUA's classification while providing an intraclass distribution for each semaphore color. A robustness evaluation using 2021 data showed comparable classification distribution (67% green, 62% yellow, and 49% red). The system effectively classifies gradual quality using key indicators: conductivity, hardness, total dissolved solids (TDS), and metal levels, aligning with CONAGUA's classical semaphore. Conclusion: Despite the existence of a superficial water semaphore, we propose using the groundwater semaphore instead. The superficial classification does not consider metals, yet preliminary multidisciplinary findings indicate metal presence in the Tampamachoco Lagoon. Therefore, the groundwater semaphore could be a suitable tool for assessing Tampamachoco Lagoon's water quality future studies.

KEYWORDS: Fuzzy logic, Fuzzy Inference System, groundwater, water quality, lagoon, open access database, water quality semaphore

INTRODUCTION

It was concluded by Arcega et al. and Ruiz et al. that human activities rather than environmental changes are the primary source of superficial and underground water pollution in lagoons in the Mexican Gulf. Industries are another primary source of pollution based on the levels of metals found in the lagoon's sediments higher than the recommended limits [1], [2],[3]. A standard metric should be used to compare the condition of different aquifers in Mexico and understand the relevance of these results. However, the water quality assessment is a complex nonlinear process because quality is a measurement that depends on many aspects, such as the environment of the aquifer and the chemical reaction between the lagoon water and its surrounding. It also depends on analyses of the qualitative and quantitative data reported by the researchers and the study areas. Even though, several proposed metrics and index methods are used for the water quality assessment [4].

In Mexico, Comisión Nacional del Agua (CONAGUA) and Secretaría de Medio Ambiente y Recursos Naturales (SEMARNAT) have established a national network to measure the quality of water (RENAMECA), which is responsible for measuring 5,000 sites across the country and started reporting in 2012. This program considers surface water, divided into lotic water bodies like streams and rivers and lentic water bodies like dams and



estuaries. It also studies groundwater bodies by measuring the quality of 2049 water samples taken directly from the water sources without processing and comparing the results with the drinking water references. A quality semaphore for groundwater was determined by the 14 parameters considered. In contrast, the semaphore developed for surface water considers 8 parameters for coastal zones and 12 for lotic and lentic ecosystems [5]. After ten years of the sensor net being in use, CONAGUA presents the statistics and geographic information maps of the quality of water and an evaluation of the sites based on three groups of the proposed semaphore, corresponding to the bodies that meet the established standards, bodies that partially meet the criteria and those that do not satisfy it [6][7][8]. The following criteria of the semaphore describe the water quality. The class is red if any of the eight parameters associated with metals, fecal coliforms, and fluorides fails to comply with the standards. Water bodies with alkalinity, hardness, and Total Dissolved Solids levels (TDS) associated with agricultural risks and salinization or with iron and manganese levels that exceed the standard thresholds are classified as yellow. Finally, if all 14 parameters are under the safe threshold, they are labeled as green. It has been observed that assessing water quality is a task that requires resources such as time, domain knowledge, and specialized tools to get samples of the study area. In some cases, it is hard to perform due to the location of water bodies because of lack of information and inaccuracy of measurements.

However, there is a need to assess water quality during research on natural environments to measure the anthropogenic impact on nature and explore the relationship between the components found in water bodies and the water quality. So, the data analysis of this information collected over ten years could be used as a reference to explore deeper and understand the ecological health of the lagoons [5].

Recently, with the increased availability of open data sources, the rising environmental problems have been extensively studied by adopting fuzzy set theories. In general, transforming the partial crisp data to fuzzy information helps to develop a fair judgment by inferring based on previous knowledge and has been adopted for spatial analysis of water quality parameters [9], such as in Abidi et al. work[10] the scare samples set]were converted into fuzzy membership and produced dry and wet seasons maps.

Fuzzy Inference Systems (FIS) incorporate the knowledge of experts into systems, absorbing the complexity of ambiguity, and it deals with uncertainty to make decisions about a phenomenon. Another approach is a combination of neural networks with fuzzy logic named ANFIS system [11],[12], which could be trained with information of several data of the same place, but in this specific case where the site's measurements are integrated into a single mean without reported standard deviation is not enough information.

Fuzzy logic proves to be an excellent tool for generating approximations with certain levels of imprecision and helps to decrease ambiguity. The structure of a fuzzy system is simple to explain and represent. Therefore, we developed a system that would be easy to interpret and maintain, allowing researchers to evaluate water quality with tolerance or imprecision for promoting a deeper analysis of the semaphore results. Our fuzzy system was designed to evaluate groundwater quality and minimize the resources and time needed to perform this type of analysis. So, a system was implemented to approximate the semaphore data reported by CONAGUA, which classifies the water quality into three classes: green, yellow, and red, proposed by CONAGUA, but it includes a distribution inside the three main groups.

METHODOLOGY:

The process to develop the FIS started by gathering data related to Mexican water bodies, being our objective was to understand this type of natural environment.

We put our efforts into searching data through public databases and organizations that focus on studying water bodies. CONAGUA made available datasets containing the data collected as part of a study of groundwater bodies nationwide from 2012 to 2021.

We explored the data provided by CONAGUA to find patterns that would help us analyze groundwater bodies' quality. Then, it was required to clean and preprocess the datasets previously to perform a statistical analysis of the dataset. As part of this cleaning process, missing values were removed, and the data types of the parameters were changed to a suitable type for calculations.

Exploratory and statistical analysis was applied to the Groundwater quality 2012-2021 and 2021 datasets, and descriptive statistics of the numerical parameters were computed. The original analysis was focused on the parameters included in the datasets; however, it looked for the main contributors to the performance of the semaphore. After finishing the analysis, we could better understand the interactions between the 14 parameters and better select the FIS inputs.

Since 2012-2021 dataset was not available for each year but instead was averaged for the period, it was compared with the available last year's information. To accomplish this, we applied adversarial validation on the combined dataset. The process consists of the following steps: 1) Combine both datasets. 2)Add labels to identify the period they belong to, 0 for the 2012-2021 and 1 for the 2021 dataset. 3) Trained a simple classifier to infer the target class and 4) Infer the target class for the validation dataset and evaluate the model.

Then, the data was split into train, test, and validation datasets. Finally, a decision tree classifier was trained and evaluated. The classifier parameters were n estimators: 100, min sample split: 2, min sample leaf: 1, max features: sqrt, bootstrap: True and random state: 300. Because of the unbalanced classes associated to the semaphore and the size of the datasets, the adversarial validation process was combined with under-sampling for the large classes



and over-sampling for the smallest class as a strategy to balance the target class. Finally, cross-validation to test the model's performance was done.

The design phase of the fuzzy inference system (FIS) consisted of selecting the input parameters and defining the output. Then FIS type was selected, and the membership functions were proposed based on the data distributions and looking for consistency in the membership grades associated with each measured value of all possible values of the universe.

The dictionary of rules and fuzzy terms was built according to the knowledge of experts and the work reported by CONAGUA. Consequently, the groundwater quality fuzzy semaphore was implemented to inference the results. Finally, an evaluation of the inference results associated with the semaphore classes was performed and compared, based on classical metrics used in automatic classification.

(a) Data

The groundwater quality datasets were created by Comisión Nacional del Agua (CONAGUA) as part of a program that included 665 groundwater bodies and 14 parameters: Fluorides, Fecal Coliforms, Arsenic, Nitrate-Nitrogen, Cadmium, Chromium, Mercury, Lead, Alkalinity, Conductivity, Hardness, Total dissolved solids (TDS), Manganese and Iron. CONAGUA provides two available public datasets. The first one is the groundwater quality from 2012 to 2021 which contains the mean value of the measurements in the considered period and the semaphore class, including 2197 instances, and the second one is the Groundwater quality recorded in 2021, with 665 instances.

Both datasets include a class identification associated with the semaphore's color to classify a water body's quality. The quality of groundwater measured by CONAGUA follows national and international standards such as the National Waters Law and the Ecologic Equilibrium and Environment Protection Law [5]. A group of parameters determines the semaphore classes. If any of the parameters in the group exceeds the permitted levels, it affects the water body quality, and a label is assigned. The clustering is defined as follows: Green label: All 14 parameters lie within the permitted levels. Yellow label: either of these parameters surpasses the permitted levels, Alkalinity, Conductivity, Total dissolved solids (TDS), Manganese and Iron.

Red: either of these parameters surpasses the permitted levels, Fluorides, Fecal Coliforms, Nitrate-Nitrogen, Arsenic, Cadmium, Chromium, Mercury, and Lead.

Both datasets were converted to CSV format to process them using the Python library Pandas v1.5.3 for data manipulation the resulted data set sample is presented in Figure 1.

	Conductivity	Hardness	TDS	Semaphore
0	329.0	129.831	210.56	Green
1	615.0	221.7114	393.6	Green
2	636.0	221.7114	407.04	Green
3	379.0	141.8154	242.56	Green
4	354.0	139.818	226.56	Green

Figure 1. A sample of the transformed groundwater dataset 2012-2021

(b) Data Analysis

Before any statistical analysis of the datasets, it is necessary to clean them and prepare the format used to represent the data; In this case, some variables included nominal values containing less than and greater than symbols (i.e., 400>, 0.01<) instead of the float data expected. After casting the data type from string to float, rows with missing values were removed from the 2012-2021 and 2021 datasets having 129 and 361 rows removed, and 6% and 54% reduction, respectively. A statistical characterization was done to obtain each parameter's minimum, maximum, mean, standard deviation, and error from the 2012-2021 dataset. The analysis was repeated for the 2021 dataset and the distribution characteristics.

(c) Fuzzy Inference System Design

A system with two Fuzzy inference subsystems was proposed. The primary FIS was designed to assess groundwater quality, looking to reduce the number of features. This FIS has four input parameters: Conductivity, Total Dissolved Solids (TDS), Hardness, and Metals level. The secondary FIS measures the levels of metals in groundwater; the result is an input to the primary subsystem. The proposed FIS for metals levels reduces the complexity and number of rules needed for the Groundwater Quality FIS.



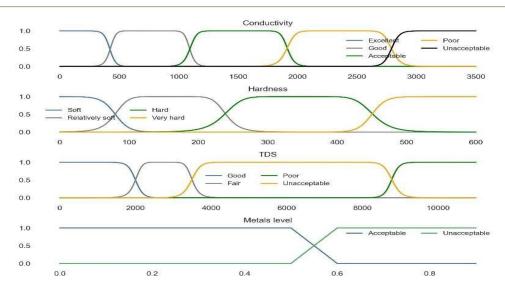


Figure 2. Membership functions (MFs) of the four input variables of the FIS for groundwater quality

Sigmoid and Gaussian membership functions (MF) were proposed for Conductivity, TDS, and Hardness. For the Metals levels, two trapezoidal functions were proposed. Figure 2 presents the four input variables of the Groundwater quality FIS.

The knowledge of expertise is summarized in the criterium of the thresholds for each variable and the primary interaction between them (Figure 3) So, the dictionary is provided to the Mamdani FIS, which contains 28 rules that define the system. Each rule evaluates the Antecedents (input variables) and finds an effect according to the rules generating a fuzzified response.

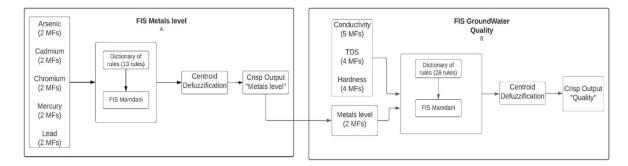


Figure 3. (a) FIS for metals levels. (b) FIS to assess Groundwater quality.

Trapezoidal functions were proposed for the Metals level FIS, and the membership functions are displayed in Figure 4. We considered Arsenic, Cadmium, Chromium, Mercury, and Lead as input in the system. Also, these five metals were considered by CONAGUA to impact the groundwater quality and modify the semaphore class to red if any of these parameters exceeded the acceptable level.

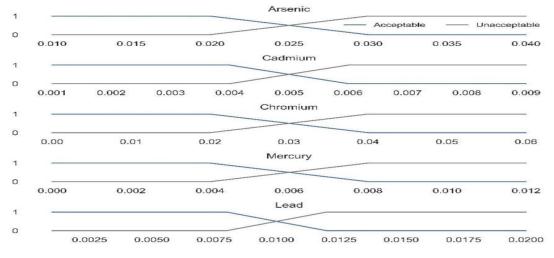


Figure 4. Memberships functions (MFs) of the five input variables of the FIS for metals concentration.



(d) FIS input selection

The Pearson correlation matrix was used to select the principal variables because this coefficient captured the linear correlation among pair of variables. It was observed that Conductivity, TDS, and Hardness have a high correlation compared to the relationship between the rest of the parameters used by CONAGUA to assess groundwater quality. Therefore, these indicators proved to have a positive relationship. The correlation coefficients matrix is shown in Figure 5.

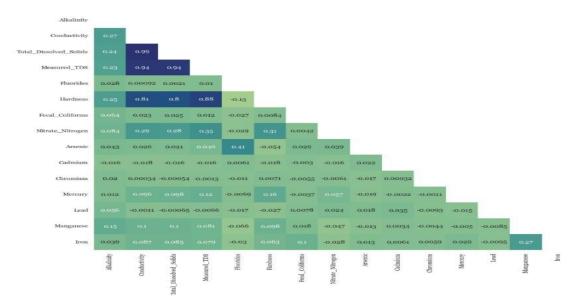


Figure 5. Correlation matrix. Groundwater 2012-2021 dataset

According to the Environmental Protection Agency (EPA), increasing amounts of TDS results in increased conductivity of water bodies. Figure 6 shows that the positive relation between TDS, Hardness, and Conductivity is visible, demonstrating the positive correlation of levels among the three parameters.

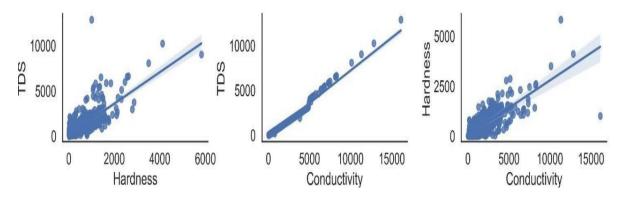


Figure 6. Relationship between Hardness, Conductivity and Total Dissolved Solids (TDS).

It is relevant to point out that in the study conducted by CONAGUA, if any of the following indicators: Alkalinity, Conductivity, TDS, Manganese, and Iron, exceed the acceptable levels, the water body would be directly classified as yellow[5].

(e) Membership Functions

The integrity of the membership function values associated with each variable along the range of the universe was tested and adjusted, so the sum of all membership degrees for each value of the four input values of the FIS, is always one as shown in (1) considering that each element only could belong from 0 to 100% to each function.

$$\mu_{i,j}(x_i) = \sum_{i=1}^{N} MF_i(x_i) = 1$$
 (1)

 $MF_{i}(x_{i})$ is the *i* membership function of N generated for each input fuzzy variable.

 $x_i \in X$ is the value of each input variable j in its universe.

j is each of the four input fuzzy variables.

 $\mu_{i,j}(x_j)$ is the sum of the grades of membership associated to each the value x



(f) Inference Rules

The inferences rules for the FIS were defined by the observations obtained as part of the data analysis, domain knowledge of the water quality indices, and the work developed by CONAGUA. Twenty-eight inference rules were defined for the FIS rules dictionary. Usually, the number of rules in a FIS is expressed by an exponential formula (2). This formula is as follows:

 $\mathbf{r} = \mathbf{N}^{\mathbf{i}} \tag{2}$

Where N is the number of fuzzy variables considered in the input.

r is the number of possible rules.

i is the number of membership functions are associated to the linguistic terms.

If the five metal parameters were used directly in the main FIS, 32,768 rules would have been needed. So, in this proposal, the Groundwater quality FIS required a maximum of 1024 rules, and the second FIS required 25 rules. After analyzing the results, the number of rules was reduced and finally was set to 28. Some rules examples are presented in Table 1.

Table 1 Example of the Inference rules defined in the FIS.

If Conductivity is Excellent ar	d TDS is Good and Hardness	s is Soft and Metals Concentration is
acceptable, then Quality is Gree	n	

If Conductivity is Good and TDS is Good and Hardness is Soft and Metals Concentration is Acceptable, then Quality is Green

If Conductivity is Acceptable and TDS is Poor and Hardness is Hard and Metals Concentration is acceptable, then Quality is Yellow

If Conductivity is Acceptable and TDS is Unacceptable and Hardness is Hard and Metals Concentration is Unacceptable, then Quality is Red.

If Conductivity is Excellent and TDS is Fair and Hardness is Very Hard and Metals Concentration is Unacceptable, then Quality is Red.

(g) FIS output

The output of the FIS is a value between 0 and 1 representing the quality of water. The system uses the centroid method as defuzzification. Each value between 0 and 1 has a membership value for each of the membership functions of the output variable as shown in Figure 7. We associated the output crisp value with one of the three semaphore classes, using the maximum membership value on the three possible fuzzy classes as the criterium.

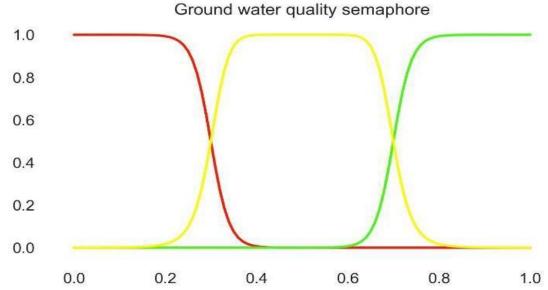


Figure 7. Membership functions of the FIS output variable. The output is a crisp value in the range of 0.0-1.0 that represents the quality of groundwater. The output variable emulates the water quality semaphore implemented by CONAGUA, with 3 classes: green, yellow and red.

CONAGUA provides the standards to assess the quality of groundwater bodies in Mexico. In Table 2, we observe the 14 parameters considered, the standards thresholds, and the semaphore class associated in case the parameters exceed the acceptable levels. So, the dictionary of rules, the membership functions for each input parameter as well as the output of the FIS design were guided by the CONAGUA standards.



Table. 2 Groundwater quality standards proposed by CONAGUA as part of their study to assess groundwater quality, including the three semaphore classes [5]

Parameter	Ranges						Semaphore	
	Accep			able Unacc			able	
Conductivity	C<=250	250	0 <c<=7:< td=""><td colspan="2">750 750<c<=2000< td=""><td>2000<c<=3000< td=""><td>Cond>3000</td><td>Yellow</td></c<=3000<></td></c<=2000<></td></c<=7:<>	750 750 <c<=2000< td=""><td>2000<c<=3000< td=""><td>Cond>3000</td><td>Yellow</td></c<=3000<></td></c<=2000<>		2000 <c<=3000< td=""><td>Cond>3000</td><td>Yellow</td></c<=3000<>	Cond>3000	Yellow
	Excellent		Good		Acceptable	Poor	Unacceptable	
Hardness	H<=60)	60 <h<< td=""><td colspan="2"><=120 120<h<=500< td=""><td colspan="2">H>500</td><td>Yellow</td></h<=500<></td></h<<>	<=120 120 <h<=500< td=""><td colspan="2">H>500</td><td>Yellow</td></h<=500<>		H>500		Yellow
	Soft		Relati	Relatively Hard		Very hard		
	so		ft					
TDS	TDS<	<=100	0	100	0 <tds<=2000< td=""><td>2000<tds<=10000< td=""><td>TDS>10000</td><td>Yellow</td></tds<=10000<></td></tds<=2000<>	2000 <tds<=10000< td=""><td>TDS>10000</td><td>Yellow</td></tds<=10000<>	TDS>10000	Yellow
	Good					Poor	Unacceptable	
Arsenic	As<	As<=0.01		0.01 <as<=0.025< td=""><td colspan="2">As>0.025</td><td>Red</td></as<=0.025<>		As>0.025		Red
	Exc	Excellent		Good		Unacceptable		
Cadmium	$Cd \le 0.003$		0.003 <cd<=0.005< td=""><td colspan="2">Cd>0.005</td><td>Red</td></cd<=0.005<>		Cd>0.005		Red	
	Excellent		Good		Unacceptable			
Chromium	Cr<=0		0.05	05 Cr>0.05		5	Red	
	Excellent					Unacceptable		
Mercury	Mercury Hg<=0		.006		Hg>0.006		Red	
	Excell		lent		Unacceptable			
Lead	Pb<=().01		Pb>0.01		Red	
	Excelle			lent		Unaccept	able	

(h) Defuzzification

The FIS uses the Centroid as a defuzzification method, so the fuzzy variables turn in crisp values. The Centroid method calculates the center of the area under the curve of the output membership function (MF) obtained with the Mamdani inference method.

(i) FIS implementation

We used MATLAB® version R2022b and the Fuzzy Logic Toolbox to design the system. Then it was migrated to a PythonTM version to be deployed as a Web service or a portable online system. In this case, the FIS was built using the library ScikitFuzzy v0.4.2s, Python 3.10.9, and NumPy 1.22.4. In Figure 8 is presented an inference result with the FIS using a chosen data point from the 2012-2021 dataset. The input values used for this simulation were: 641 (Conductivity), 410.24 (TDS), 219.714 (Hardness), and 0.77 (Metals level). For this selected point, the quality was 0.1522, so it is mapped to the "Red" class of the semaphore. The vertical dashed line in the figure marks the intersection with the membership associated with the output variable, so the crisp output value obtained is the maximum of the membership grades at that point: 0.1522.

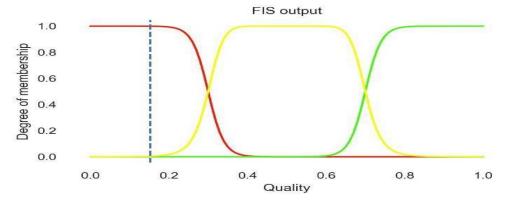
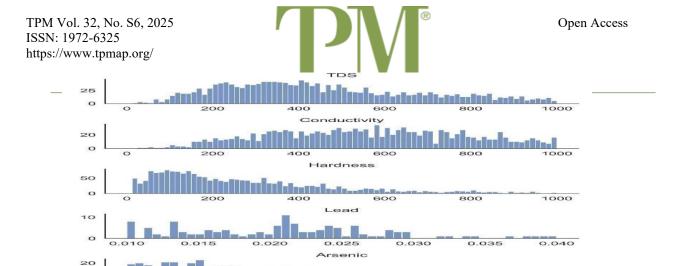


Figure 8. FIS inference example using the data point (641,410.24,219.71,0.77) from the 2012-2021 dataset.

(j) FIS testing

The 2012-2021 and 2021 datasets were used to test the FIS. The first part of the testing was to fetch the data to the Metals FIS to obtain the output "Metals level", one of the inputs to the Groundwater quality FIS. Once the Metals level parameter was obtained, we fetched the four inputs to the FIS (conductivity, TDS, hardness, and Metals level). The inference and defuzzification of the Groundwater quality FIS were evaluated with the dataset. Finally, outputs were mapped to one of the three semaphore classes based on each output's membership value. So, fuzzy semaphore and a crisp semaphore were compared.



0.025

0.016

0.030

0.017

0.0004908

0.0004906

0.040

0.0004910

100

0.015

0.0004904

(k) FIS Evaluation

Then a confusion matrix was used as a metric to evaluate the similarities of the proposed system with the crisp CONAGUA semaphore as a reference. In the confusion matrix, the actual classes reported by CONAGUA and the obtained classes by the FIS are compared with each other to analyze the performance of the FIS using both datasets.

Because the 2012-2021 dataset was used to construct the FIS, then the 2021 dataset was used to evaluate it. Metrics computed to evaluate the FIS were accuracy, precision, recall, and F1 score. To calculate each score, we separately considered the semaphore labels obtained for the 2012-2021 and 2021 datasets.

RESULTS:

To analyze the data distribution visually, we plotted the histograms for each parameter for both datasets, Groundwater quality 2012-2021 and 2021 (Figure 9 and Figure 10). The histograms show a positive-skewed distribution for the TDS, conductivity, and Hardness; therefore, an exponential distribution would be better to approximate these parameters.

On the other hand, metals such as Cadmium, Lead, Arsenic, Mercury, and Chromium do not have significant variation. Descriptive statistics are presented in Table 3.

If just the 2021 dataset is considered, the distribution characteristics are summarized in Table 4, and the histograms shown in Figure 10. The 2012-2021 dataset results were summarized in a confusion matrix presenting the classes assigned by the proposed FIS (Figure 11A) compared with the classes included in the dataset defined by CONAGUA.

Table 3 Descriptive statistics of the features considered in Groundwater dataset for 2012-2021.

Parameter	Min	Max	Mean	Standard	Standard	Kurtosis	Skewness
				deviation	error		
TDS	7.68e+01	3148.80	6.13e+02	4.92e+02	2.82e+01	4.97	2.02
Conductivity	1.20e+02	4920.00	9.58e+02	7.69e+02	4.41e+01	4.97	2.02
Hardness	1.99e+01	1637.86	2.81e+02	2.28e+02	1.30e+01	10.08	2.52
Cadmium	2.90e-03	0.00	2.90e-03	4.01e-05	2.30e-06	304.00	17.43
Lead	4.90e-03	0.02	5.09e-03	1.89e-03	1.08e-04	126.71	11.03
Arsenic	9.00e-03	0.31	1.92e-02	2.92e-02	1.67e-03	42.09	5.60
Mercury	4.90e-04	0.001	5.07e-04	9.25e-05	5.31e-06	49.81	6.63
Chromium	4.90e-03	4.25	2.42e-02	2.44e-01	1.39e-02	300.50	17.28

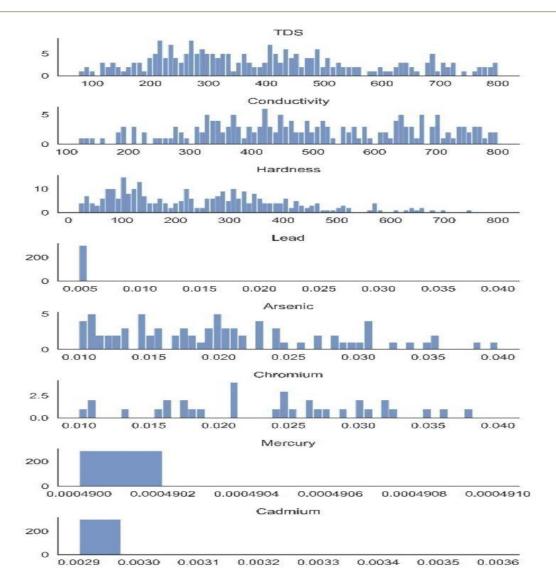


Figure 9. Parameters distribution. 2012-2021 Groundwater Quality dataset

The 2012-2021 dataset results were summarized in a confusion matrix presenting the classes assigned by the proposed FIS (Figure 11A) compared with the classes included in the dataset defined by CONAGUA. CONAGUA classified 2197 groundwater bodies all over Mexico. After the data cleaning process, 6% of rows were removed from the 2012-2021 dataset leaving information about 2068 water bodies left. Out of the 2068 water bodies registered in the 2012-2021 dataset, 888 were classified as Green, 802 as Red, and 378 as Yellow by CONAGUA. On the other hand, for the 2021 dataset, the FIS classified 918 water bodies as Green, 598 as Red, and 552 as Yellow.

Figure 10. Parameters distribution of the 2021 Groundwater dataset **Table 4** Descriptive statistics of the parameters of Groundwater quality 2021 dataset.

Parameter	Min	Max	Mean	Standard	Standard	Kurtosis	Skewness
				deviation	error		
TDS	2.49e+01	12880	7.01e+02	8.09e+02	1.77e+01	53.20	5.72
Conductivity	2.77e+01	16100	1.07e+03	1.10e+03	2.41e+01	35.49	4.50
Hardness	1.99e+01	5828.68	3.38e+02	3.71e+02	8.16e+00	36.55	4.27
Cadmium	2.90e-03	0.15	3.02e-03	3.57e-03	7.85e-05	1576.02	38.42
Lead	4.90e-03	0.08	6.65e-03	7.10e-03	1.56e-04	36.53	5.48
Arsenic	9.00e-03	0.41	2.12e-02	3.42e-02	7.52e-04	42.07	5.62
Mercury	4.90e-04	0.02	5.18e-04	4.46e-04	9.82e-06	1627.45	38.46
Chromium	4.90e-03	2.14	9.44e-03	7.60e-02	1.67e-03	704.93	26.34



The total number of groundwater bodies labeled by the FIS is 2068, accounting for 94% of the 2012-2021 dataset. And the Groundwater quality FIS inferred the classes presented in Figure 12, reporting a percentage of the same labeled groundwater bodies: 75% Green, 69% Yellow, and 67% Red. The results obtained by the FIS using the Groundwater Quality 2021 dataset are shown in the confusion matrix in Figure 11B, presenting the CONAGUA semaphore and the FIS classes.

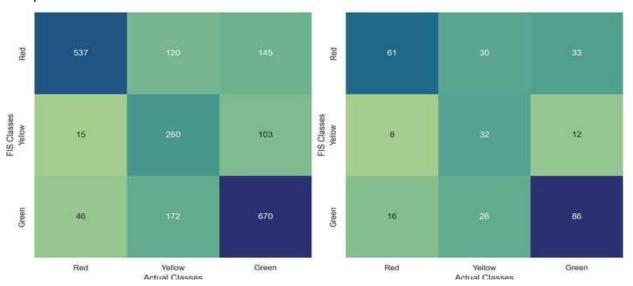


Figure 11. Confusion Matrix of the FIS and CONAGUA classes of Groundwater quality. A) 2012-2021, B) 2021 dataset.

As for the 2021 dataset, CONAGUA classified 665 groundwater bodies. After removing the rows with missing values, we ended up with 304 water bodies representing 46% of the 665 water bodies; 128 were classified as Green, 52 as Yellow, and 124 as Red by CONAGUA.

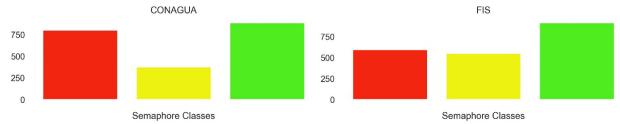


Figure 12. Semaphore classes. Groundwater Quality 2012-2021 dataset.

The FIS inferred the classes shown in Fig. 13, reporting 131 as Green, 88 as Yellow, and 85 water bodies as Red. The percentage of the same labeled water bodies per class is 67% Green, 62% Yellow, and 49% Red.

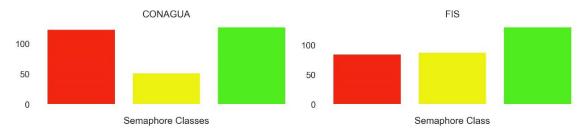


Figure 13. Semaphore classes. Groundwater Quality 2021 dataset

This result shows differences mainly in the red and yellow classes, so the fuzzy values were analyzed. The distribution of the FIS output for the 2012-2021 and 2021 datasets is shown in Figure 14. The FIS output named "Quality" is in the range of 0 to 1. The semaphore class colors the membership grade values of the output variable. One first observation is that the 2012-2021 data represents the average values of 9 years. Nevertheless, the distribution of the semaphore classes is like the 2021 dataset.

Therefore, we can tell that the behavior of groundwater bodies in 2021 is not far from the general behavior as is reported in the adversarial evaluation.



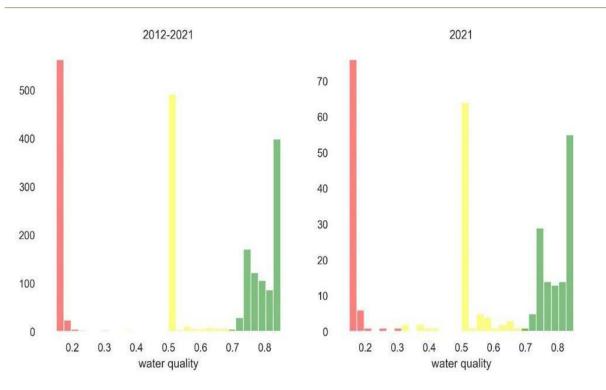


Figure 14. FIS output values colored by their corresponding semaphore class

Finally, the comparison of the dataset in the fuzzy semaphore shows a slightly difference in the internal distribution into the semaphore that is not possible to observe in the crisp version proposed by CONAGUA.

For our study, the weighted F1 better approximates how well the FIS performed, given that this metric considers both the precision, the recall, and the imbalance of the semaphore classes. In Table 5, the results are presented for both datasets classification.

Table 5 Scores of the metrics used to evaluate the FIS semaphore with controversial evaluation.

Metric	2012-2021 Scores	2021 scores
Accuracy	0.71	0.59
Precision	0.75	0.63
Recall	0.71	0.59
F1	0.72	0.60

DISCUSSION

The Mexican states with better groundwater quality evaluation are associated with more green assignments to their water bodies classified by CONAGUA. They are Sonora (136), Guanajuato (132), and Durango (109); the FIS semaphore obtained the same states as follows Guanajuato (157), Sonora (131), Durango (114). On the other hand, states with more red water bodies identified by CONAGUA are Durango (241), Guanajuato (109), and Coahuila (77). Similarly, the results obtained by the FIS semaphore are Durango (195), Guanajuato (79), and Coahuila (50).

Lastly, the top states with yellow water bodies classified by CONAGUA are Sonora (44), Yucatan (42), and Coahuila (31). In contrast, the FIS semaphore obtained the following: Yucatan (82), Sonora (68), and Tamaulipas (56). The FIS obtained similar results in all three classes, accounting for the top 3 states per class.

Over the years, several indices have been developed to assess surface and groundwater quality. These works define the critical parameters concerning water quality and provide a partial solution in water assessments [13]. Some of the worldwide representatives are the National Sanitation Foundation Water Quality Index (NSFWQI), the Canadian Council of Ministers of the environment water quality index (CCMEWQI), the Oregon Water Quality Index (OWQI), and the Weight Arithmetic Water Quality Index (WAWQI). All these indexes include TDS, Hardness, Alkalinity, Biochemical Oxygen Demand (BOD), and Dissolved Oxygen (DO) as the usual parameters to calculate the water quality. The parameters' selection differences are related to the water's final use. CONAGUA uses worldwide information such as guidelines for drinking water quality[14], agriculture water



quality, and other countries regulations as Bolivia, Chile, Spain, New Zealand, Canada, Malaysia[15]. Nevertheless, the final Mexican thresholds are defined by Mexican regulations, considering information from the Federal Laws in the section: National Applied Disposition in the matter of water [6], Health Secretary [7] and the Urban Development and Ecology Secretary [8]. Hence, it is the reference of this study. An updated norm NOM-001-SEMARNAT-2021 has been identified, but it was not considered in this model proposal because data were classified with the previous version. However, if required, the system could be upgraded with new values in a future version.

This work pretends to generate a methodology for implementing the fuzzy inference system based on the available national information. One of the first steps is to understand the data provided by the datasets and gather all the information associated with the problem. Then preprocessing the information and finally designing and implementing the system. The results supported by the computed scores show a good performance of the Proposed FIS for 2012-2021. However, the performance was drastically reduced, evaluating only the 2021 dataset.

The reduction of the measurement available, the presence of missing values, and changes in the distribution of the data measured in 2021 could explain the poor performance. Nonetheless, the distribution of the semaphore classes is similar, and the extensive adversarial test proves no difference. However, the overall performance of the classifier used for adversarial validation is 0.6132, which means that the distribution of both datasets is somehow different but not significant enough for the classifier to distinguish them furthermore.

Similar to Sajib et al. [16], the relevance of heavy metal measurements is confirmed in the search for understanding the relationship between Groundwater and its context. The semaphore is observed to be highly affected by existing geology, quality of recharge, degree of chemical weathering, level of Groundwater, and some surface elements. The result of the complex interaction between these processes is reflected in the metrics considered [10], [17]. As others author reported, it is observed that the quality of Groundwater is subjected to the interaction between geological and hydrological processes.

CONCLUSION:

A FIS to assess the Groundwater quality was proposed using a reduced number of parameters from the CONAGUA semaphore. A similar semaphore was implemented using the FIS results to classify groundwater bodies. However, the FIS uses 8 of the 14 original variables, and five of those eight variables were grouped in a new indicator named Metals level, which monitors the non-allowed metal levels. The advantage of using two FIS in parallel is that the number of combinations in the final evaluation is considerably reduced. If all variables were used, there would be 1024 rules according to the exponential formula used to calculate the total number of rules. The number of rules used is 28, which represents only 3% of the total number of rules. The parameters with the highest correlation are Conductivity, TDS, and Hardness. These parameters have a positive correlation higher than 0.80 and were selected as inputs in the principal subsystem. On the other hand, Cadmium, Chromium, and Mercury present the highest skewness and kurtosis; therefore, we can tell that the dataset contains several outliers representing extreme values for these parameters. Even ANFIS could be included in each site to define the water quality in an intelligent sensor net. In this case, FIS was selected because of the database information structure. The classification was not the primary goal of this paper; even though the FIS managed to correctly classify above 50% of data points for each semaphore class, the FIS can improve its performance by fine-tuning the chosen inputs and increasing the dictionary of rules. The main goals were to follow the standard thresholds provided by CONAGUA and give a degree of membership in the classes to analyze the distribution inside each one and find a methodology to prepare and define the rules.

This work is doing into a multidisciplinary project about multidimensional health of Tampamachoco Lagoon. Although there is a superficial water semaphore, we consider using the groundwater semaphore instead because the superficial does not consider metals as an indicator of water quality due to the preliminary results of the multidisciplinary team; Because Metals are presented in the Lagoon then the groundwater semaphore could be suitable for evaluating the Tampamachoco Lagoon water quality using the proposed model in future work.

Acknowledgement:

Special thanks to the FRESA group, whose work is focused on the analysis of the risk factors involved in the environmental health assessment of coastal bodies with different degrees of contamination, supported by the grant number 20243987. Their insights and recommendations have been invaluable in shaping our proposal and fostering our interest in the field of multivariable automatic analysis in this topic.

Author contributions:

All authors contributed to the study conception and design. Inference System implementation was done by Ulises Montoya Canales. Material preparation, data collection and analysis were performed by Ulises Montoya Canales, Pilar Gomez Miranda and Laura Ivoone Garay Jiménez. The first draft of the manuscript was written by Ulises Montoya Canales, Laura Ivoone Garay Jiménez and Blanca Tovar Corona, and it was revised by Ana Judith Marmolejo Rodriguez and Pilar Gomez Miranda. All authors commented on previous versions of the manuscript. All authors read and approved the final manuscript.



Data Availability statement:

The data that support the findings of this study are available from Sistema Nacional de Información del Agua (SINA) acceded on https://www.gob.mx/conagua/es/articulos/indicadores-de-calidad-delagua?idiom=es#:~:text=Los%20Indicadores%20subterr%C3%A1neos%20son%2014,Cond_elec)%2C%20Dure za%20Total%20(Dur_Tot

Conflict of interest: "The authors declare that there is no conflict of interest".

Funding Statement:

This work has been supported by the project "Multidimensional models of temporal series associated to the anthropic contamination in marine organisms consumed by humans and its effect on their overall health", SIP 20211164, 20220701, 20230872 funded by Instituto Politécnico Nacional of Mexico within 2021–2023.

REFERENCES:

- [1] Arcega-Cabrera F, Garza-Perez R, Noreña-Barroso E, Oceguera-Vargas I, Impacts of Geochemical and Environmental Factors on Seasonal Variation of Heavy Metals in a Coastal Lagoon Yucatan, Mexico. Bull Environ Contam Toxicol, 2015, 58–65. https://doi.org/10.1007/s00128-014-1416-1
- [2] Ruiz-Fernandez AC, Rangel-Garcia M, Perez-Bernal LH, Lopez-Mendoza PG, Gracia A,Schwing P, Hollander D, Paez-Osuna F, Cadorso-Mohedano JG, Cuellar-Martinez T, Sanchez-Cabeza JA, Mercury in sediment cores from the southern Gulf of Mexico: Preindustrial levels and temporal enrichment trends. Marine Pollution Bulletin, 2019, 149. doi:https://doi.org/10.1016/j.marpolbul.2019.110498
- [3] Martínez ML, Silva R, Lithgow D, Mendoza E, Flores P, Martínez R, Cruz C, Human impact on coastal resilience along the coast of Veracruz, Mexico. In: Martinez, M.L.; Taramelli, A., and Silva, R. (eds.), Coastal Resilience: Exploring the Many Challenges from Different Viewpoints. Journal of Coastal Research, 2017, Special Issue 77: 143–153. https://doi.org/10.2112/SI77-015.1
- [4] Chidiac SE, A comprehensive review of water quality indices (WQIs): history, models, attempts and perspectives. Reviews in Environmental Science and Bio/Technology, 2023, 22:349-395. 634 https://doi.org/10.1007/s11157-023-09650-7
- [5] CONAGUA, Water Quality (in spanish), 2022. Consulted on https://www.gob.mx/cms/uploads/attachment/file/925192/Generalidades_Indicadores_de_calidad_del_agua.pdf. Accessed march 2023. Updated version 2024.
- [6] CONAGUA, Federal Rights Law. Provisions applicable to national waters, 2015, https://www.gob.mx/conagua/acciones-y-programas/situacion-de-los-recursos-hidricos Accessed 20 February 2023 (in spanish)
- [7] Secretaría de Salud, Monitoring of primary contact water in seawater from beaches and freshwater bodies (in Spanish), 2015, https://www.gob.mx/cofepris/documentos/manual-operativo-monitoreo-de-agua-de-contacto-primario-en-el-agua-de-mar-de-playas-y-cuerpos-de-agua-dulce Accessed 17 February 2024
- [8] Secretaría de Desarrollo Urbano y Ecología (SEDUE), AGREEMENT establishing the Mexican Ecological Water Quality Criteria CE-CCA-001/89. Diario Oficial de la Federación,1989.
- https://www.dof.gob.mx/nota_detalle.php?codigo=4837548&fecha=13/12/1989#gsc.tab=0 Accessed 2 march 2024
- [9] Chidambaram S, Prasanna MV, Ventramanan S, Nepolian M, Pradeep K, Banajarani P, Thivya C, Thilagavathi R, Groundwater quality assessment for irrigation by adopting new suitability plot and spatial analysis based on fuzzy logic technique. Environmental Research, 2022, 204. https://doi.org/10.1016/j.envres.2021.111729
- [10]ABIDI, Jamila Hammami, et al. Evaluation of groundwater quality indices using multi-criteria decision-making techniques and a fuzzy logic model in an irrigated area. Groundwater for Sustainable Development, 2024, vol. 25, p. 101122. https://doi.org/10.1016/j.gsd.2024.101122Get rights and content
- [11] Shwetank, SJ, Suhas, Jitendra KC, Hybridization of ANFIS and fuzzy logic for groundwater quality assessment. Groundwater or Sustainable Development, 2022, 18. doi:https://doi.org/10.1016/j.gsd.2022.100777
- [12] Jha MK, Assessing groundwater quality for drinking water supply using hybrid fuzzy-GIS-based water quality index. Water Research, 2020, 179. doi:https://doi.org/10.1016/j.watres.2020.115867
- [13] Vigueras-Velázquez, M. E., Carbajal-Hernández, J. J., Sánchez-Fernández, L. P., Vázquez-Burgos, J. L., & Tello-Ballinas, J. A. Weighted fuzzy inference system for water quality management of Chirostoma estor estor culture, Aquaculture Reports, 2020, 18, 100487. https://doi.org/10.1016/j.aqrep.2020.100487
- [14] OMS, Guidelines for Drinking-water Quality., 2008, 1. https://apps.who.int/iris/handle/10665/42852 Accesed 20 march 2024.
- [15] FAO, Water quality for agriculture. FAO Irrigation and Drainage, 1994. https://www.fao.org/3/t0234e/T0234E00.htm#TOC, Accessed 23 september 2024
- [16] SAJIB, Abdul Majed, et al. Developing a novel tool for assessing the groundwater incorporating water quality index and machine learning approach. Groundwater for Sustainable Development, 2023, vol. 23, p. 101049.
- [17] PATEL, Neha; BHATT, Darshana. Insights of ground water quality assessment methods—A review. Materials Today: Proceedings, 2024. https://doi.org/10.1016/j.matpr.2024.04.045