

# **WATER QUALITY ASSESSMENT WITH THE AID OF MACHINE LEARNING**

Author: Kalpana Baheti

## ***ABSTRACT***

The chemical compounds and their concentration in a water body change with the environment and with the addition of effluents. With several contaminants and nutrients in the matrix, it becomes very tedious to analyze which set of interdependent reactions will occur followed by the newly developed contaminants and the concentration of old contaminants in the water. Such information is necessary to a certain extent at the beginning of an industrial venture or a project plan. This foresight is hence obtained with the aid of a repository of snapshots of chemical tests performed on the water sample in question, prior to the start of the project. These snapshots are analyzed and the effect of each individual contaminant on certain chemicals known to be present in the water are extrapolated. These predicted concentrations of known chemicals then act as analytical indicators and return the entire chemical information of the water body in question at given time in the future. The set of contaminants present in the effluents can vary, hence pattern matching algorithms map the real-time sets to the snapshots in the database to minimize computation. In this manner, quantity of effluent release and schedule of release may be re-adjusted during project planning stage, and with enough certainty, toxicity limits will not be crossed.

## ***ANALYTICAL INDICATION***

Chemical indicators are substances that give a visible sign, usually a color change, of the presence or absence of a threshold concentration of a chemical species. Their limitations, however, are that most of them are one-to-one indicators and cannot confirm several presences at once, concentrations are not provided by many of them and most importantly, they can only be used in real-time data analysis, not predictive analysis.

Analytical indicators on the other hand, observes the relationship between every nutrient with the other, determines the rate of change of one nutrient with respect to another. Consider the complexity issues:

If there were  $N$  nutrients being analyzed, the total amount of comparisons to be made would be  ${}^NC_2$  which becomes considerably larger if  $N$  increases even a little. Hence, the search is filtered to analyzing sets of nutrients that have a correlation coefficient closer to 1 or -1 when compared to a single nutrient that has a strong dependence on time and physical environment and can be predicted with more certainty. To understand why we choose this arrangement, consider the diagrams below.

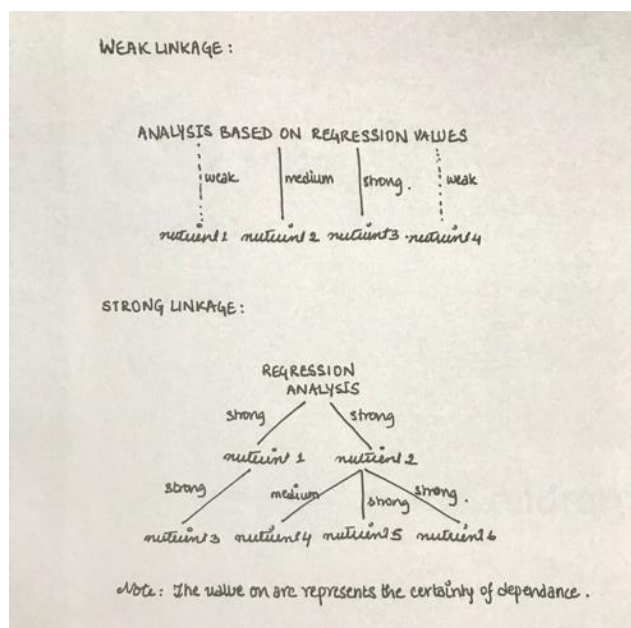


Diagram 2 has a much stronger logical linkage structure than the one in diagram 1. The correlation coefficients near 1 and -1 (selected within a range 1 to 0.75 and -0.75 to 1 for most cases) show that the dependency is linear in nature and is strong. The sign bits + and – refer to direct dependency and inverse dependency respectively. The closer they are to 1 or -1, the stronger is the dependency. Given below are marked nutrients that have more certainty of prediction based on regression analysis, and nutrients having high linear dependency:

Parameters	Free CO <sub>2</sub> r	DO r	BOD r	COD r	Sulphate r	Phosphate r	Nitrite r	Iron r	Sodium r
Free CO <sub>2</sub>	1	-	-	-	-	-	-	-	-
Dissolved Oxygen	0.349	1	-	-	-	-	-	-	-
BOD	-0.875 ✓	-0.836 ✓	1	-	-	-	-	-	-
COD	-0.651 ✓	-0.85 ✓	0.91 ✓	1	-	-	-	-	-
Sulphate	-0.369	0.165	0.529 ✓	0.478	1	-	-	-	-
Phosphate	-0.497	-0.497	0.572 ✓	0.405	0.94 ✓	1	-	-	-
Nitrite	0.197	-0.774 ✓	0.267	0.445	-0.121	-0.132	1	-	-
Iron	-0.395	-0.586 ✓	0.673 ✓	-0.883 ✓	0.282	0.221	0.7001 ✓	1	-
Sodium	-0.228	-0.678 ✓	0.622 ✓	0.669 ✓	0.3188	0.323	0.323	0.66 ✓	1
Calcium	0.897	-0.853 ✓	0.447	0.332	0.192	0.406	0.552 ✓	0.521 ✓	0.512 ✓
Magnesium	-0.25	-0.727 ✓	0.647 ✓	0.0689	0.435	0.452	0.882 ✓	0.799 ✓	0.923 ✓
Potassium	-0.213	-0.609 ✓	0.5479 ✓	0.509 ✓	0.234	0.272	0.234	0.799 ✓	0.933 ✓
Electrical conductivity	-0.257	-0.609	0.547	0.5006	0.234	0.276	-0.756	-0.714	-0.772
	0.228	0.356	0.245	0.23	0.125	0.218	0.132	0.125	0.256

	Temp	Transparency	pH	DO	BOD	Alkalinity	Chloride	TDS	TH	Nitrate	Phosphate
Temp.	-	-0.61* ✓	-0.45	-0.55 ✓	0.29	-0.11	0.56 ✓	0.01	0.76* ✓	0.09	0.31
Transparency		-	0.30	0.16	-0.35	0.18	-0.23	0.53 ✓	-0.55 ✓	-0.17	0.04
pH			-	0.14	-0.37	-0.07	0.09	-0.07	-0.08	-0.10	-0.22
DO				-	-0.36	0.08	-0.24	-0.25	-0.69* ✓	0.01	-0.05
BOD					-	-0.20	0.04	-0.08	0.43	0.44	0.25
Alkalinity						-	-0.05	0.39	-0.56 ✓	0.36	0.43
Chloride							-	0.13	0.48	0.52 ✓	0.37
TDS								-	-0.20	0.11	0.09
TH									-	-0.005	0.06
Nitrate										-	0.58* ✓
Phosphate											-

\*Significant at  $P < 0.05$  level

This relationship between nutrients is considerably more reliable than basing the relationship of nutrients on physical environment and time factor alone. Hence, the aim is to keep the set of nutrients predicted by regression analysis to the minimum. The effect of this fundamental set on secondary nutrients is carefully studied. Several instances of the concentrations at different times and circumstances are obtained, plotted and analyzed. These instances are called snapshots, they give information on quality of water in terms of nutrients. From these snapshots, multivariant statistical models are created and the algorithm acts upon these models to get correlation coefficients for pairs of nutrients. Then these nutrients are grouped using clustering.

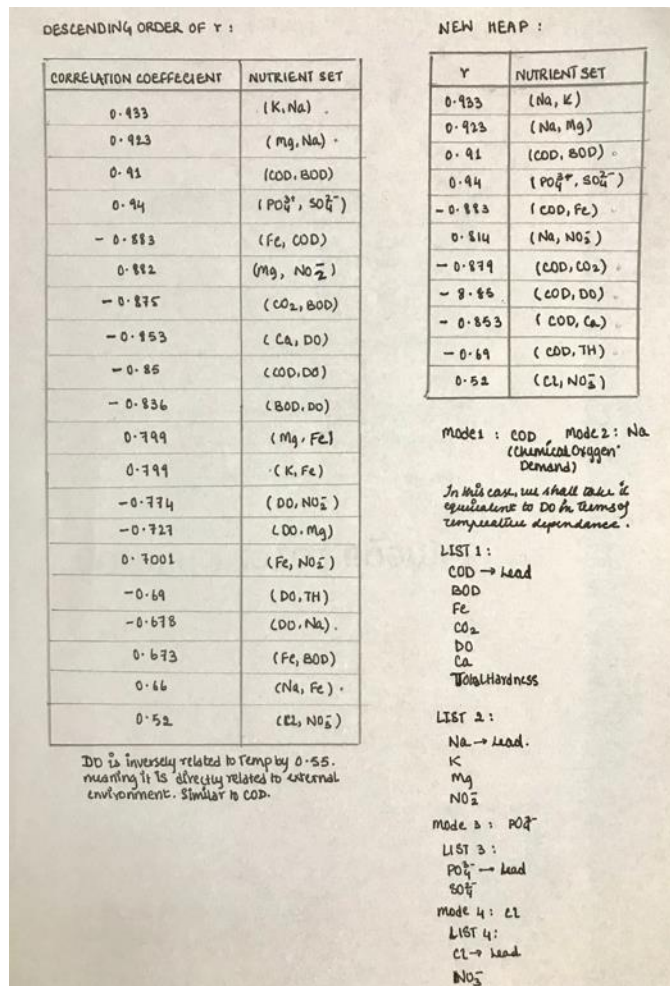
Examples of analytical indication would include the smooth inverse relation between concentration of dissolved oxygen and biological oxygen demand.

### *SNAPSHOTS AND CONCEPTS APPLIED TO THEM*

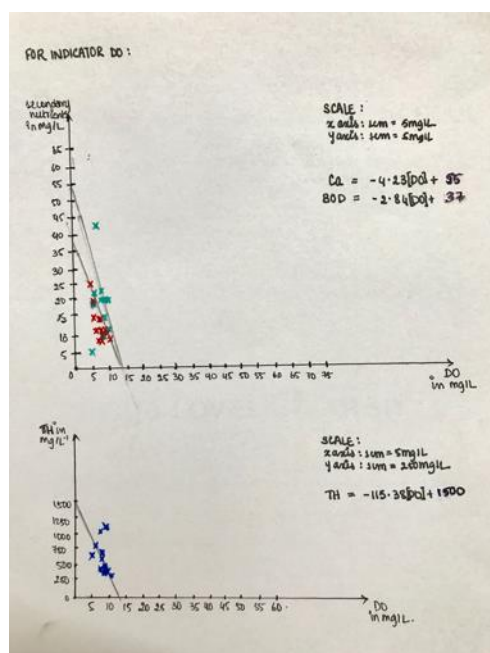
The first step as mentioned before is to categorize secondary nutrients based on their indicator nutrient. This will be done using the following partial derivative of Dijkstra's algorithm.

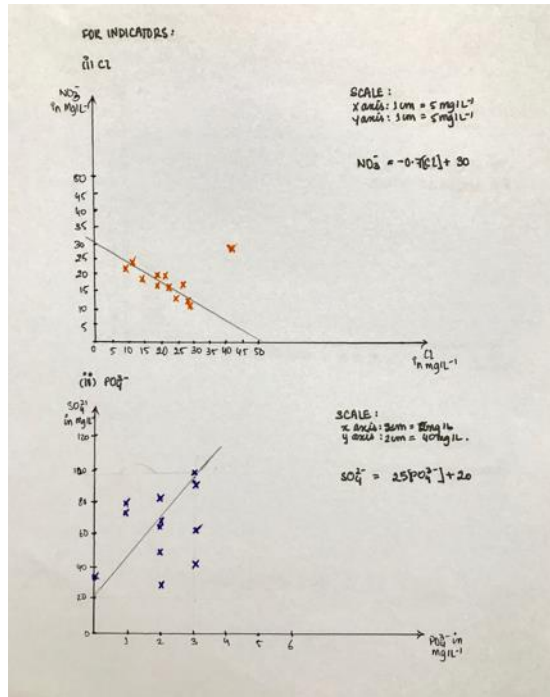
Algorithm:

1. For all positive selected correlation coefficients ( $r$ ), group them in descending order of  $|r|$  along with the concerned nutrients.
2. Initialize the count of each nutrient to 0 and flag to untraversed ( $u$ )  
 For all nutrient pairs, starting with highest coefficient;
  - If nutrient a. flag =  $u$  and nutrient b. flag =  $u$ 
    - Push onto heap
    - a. flag = b. flag = traversed ( $t$ )
    - a. count + = 1
    - b. count + = 1
  - Else if nutrient a. flag =  $u$  xor nutrient b. flag =  $u$ 
    - Consider if nutrient has already been traversed with nutrient c.
    - Create tuple ( $r_{c,b} = r_{b,a} * r_{a,c}$ , (nutrient c, nutrient a)) in favor of nutrient with Greater count.
    - Push tuple onto heap
    - b. flag =  $t$
    - a. count + = 1
    - b. count + = 1
  - Else
    - Pass
3. Assess the heap and find mode of the total nutrient set  $m$ .  
 This nutrient is an indicator nutrient. Push all tuples having  $m$  to another list and delete from the current heap.  
 Repeat step 3 on the heap till heap is empty.
4. For every list created mark the indicator.



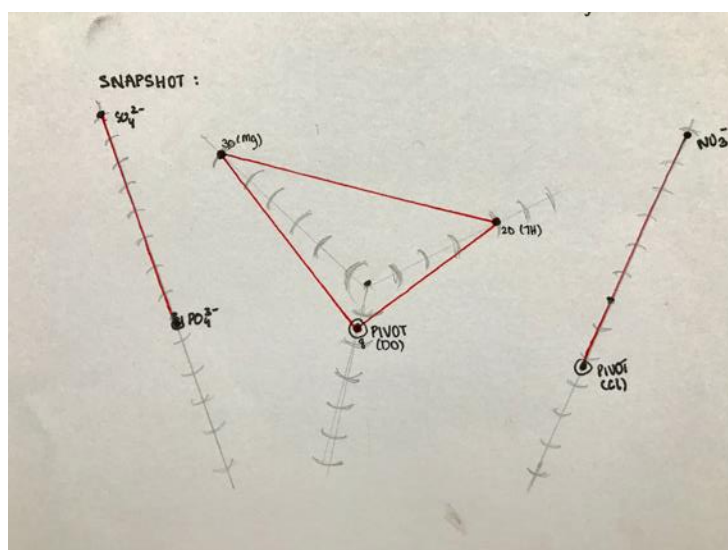
The next step involves sample data analysis with the use of snapshot algorithms based on variation of indicators which has now reduced our computation. Arrange the indicator samples in ascending order of concentration. Since the rest of chemicals have strong linear dependency, the graphs plotted, with the help of curve fitting, will give gradients of the dependency, that is, relative rate of change.





### The Snapshot Algorithm:

1. Obtain future values of indicator nutrients with the aid of regression analysis and plot the points on the set of concentric circles representing concentrations. Then join all points to form a polygon.
2. Name the indicator's vertex as pivot. Attach the rate relative gradient (b) of all the other vertices as an attribute. Perform this for all indicator graphs.
3. Now study the effect of a set of contaminants from an effluent on the indicator set of nutrients, since these are of smaller size, they may be handled reaction by reaction.
4. Deduce new values of indicator concentrations.
5. Push the pivots of each graph to new concentrations.
  - a. Nutrient concentration = (pivot concentration)\*b + C
  - b. Move the vertices to the new locations on graph.



## RESULTS AND PROOF OF CONCEPT

After sampling these results against a set of contaminants (Heptachlor, Dieldrin, Aldrin, Diazinon), that are periodically drained in the lake, the propagated values of the nutrients are:

1. Chloride = 120 ppm
2. Dissolved Oxygen = 8 ppm
3. Phosphate = 0.5 ppm
4. Sulphate = 32.5 ppm
5. Nitrate = 0.5 ppm
6. Calcium = 21.16 ppm

Comparing it with the table given below, it is proven that the method works with required accuracy.

## STRENGTHS AND DRAWBACKS

1. This method of computation does not require in depth knowledge of chemical redox reactions that occur in water, nor does it need a reaction-by-reaction computation caliber. It is based on observations and data analysis, and the understanding of how to maximize use of strong dependences, eliminate weak dependences and return more accurate results.
2. The drawback is that several times perfect analytical indicators are not found. Also, there may always be a possibility of error as all predictions are subjected to, due to inclusion of outliers while calculating gradients, incorrect correlation coefficients and such misfortunes.

## SCOPE IN ENVIRONMENTAL CHANGE

A slightly more accurate method can result in better planning of projects involving effluent release, while simultaneously ensuring their obligation to prescribed environmental water standards. This may be generated as a template and handle interrelated resources' chemical quality in a reinforced manner. The method may be applied to set of contaminants too, in such a manner that when multiple effluent loads are released, even though each one individually might tip the environmental standards, the net effect of the aggregated release is considerably harmless to the environment. This will also help project planners decide when to release which effluents and how to schedule their production year.

## BIBLIOGRAPHY

1. D. Kamal, A.N. Khan, M.A. Rahman and F. Ahamed, 2007. Study on the Physicochemical Properties of Water of Mouri River, Khulna, Bangladesh, *Pakistan Journal of Biological Sciences*, 10: 710-717.  
DOI: [10.3923/pjbs.2007.710.717](https://doi.org/10.3923/pjbs.2007.710.717)  
URL: <https://scialert.net/abstract/?doi=pjbs.2007.710.717> for correlation coefficient data Udaisagar Lake, Udaipur.
2. <http://www.moef.nic.in/sites/default/files/nlcp/P%20-%20World%20Case%20Studies/P-30.pdf> for comparison with actual data.
3. <http://www.environmentaljournal.org/1-3/ujert-1-3-8.pdf> for correlation coefficient data and training data sets of Udaisagar Lake, Udaipur.
4. [https://www.researchgate.net/figure/Water-chemistry-of-Lake-Udaisagar-values-are-mean-n-24-AE1-SE\\_fig3\\_321575458](https://www.researchgate.net/figure/Water-chemistry-of-Lake-Udaisagar-values-are-mean-n-24-AE1-SE_fig3_321575458) for training data sets of Udaisagar Lake, Udaipur.
5. Oxidation of Diazinon by Aqueous Chlorine: Kinetics, Mechanisms, and Product Studies  
Qi Zhang and and Simo O. Pehkonen\*

*Journal of Agricultural and Food Chemistry* **1999** 47 (4), 1760-1766  
DOI: 10.1021/jf981004e to study effect of contaminants on nutrients.