

## Modeling trihalomethanes concentrations in water treatment plants using machine learning techniques

Jongkwan Park<sup>a</sup>, Chan ho Lee<sup>a</sup>, Kyung Hwa Cho<sup>a</sup>, Seongho Hong<sup>b</sup>, Young Mo Kim<sup>c,\*</sup>, Yongeun Park<sup>d,\*</sup>

<sup>a</sup>School of Urban and Environmental Engineering, Ulsan National Institute of Science and Technology, Ulsan 44919, Korea

<sup>b</sup>Department of Chemical Engineering, Soongsil University, 369 Sangdo-Ro, Dongjak-Gu, Seoul 156-743, Korea

<sup>c</sup>School of Environmental Science and Engineering, Gwangju Institute of Science and Technology (GIST), 123 Cheomdan-gwagiro, Buk-gu, Gwangju 61005, Korea, email: youngmo@gist.ac.kr (Y.M. Kim)

<sup>d</sup>School of Civil and Environmental Engineering, Konkuk University, Seoul 05029, Korea, email: yepark@konkuk.ac.kr (Y. Park)

Received 25 November 2017; Accepted 17 January 2018

---

### ABSTRACT

Water disinfection process in a water treatment process results in the formation of disinfection by-products (DBPs), including total trihalomethanes (TTHMs). It takes a relatively long time to estimate TTHMs concentration level in the water treatment plants; thereby it is impossible to timely control operation parameters to reduce the TTHMs concentration. Here, we developed a predictive model to quantify TTHMs concentration using conventional water quality parameters from six water treatment plants in Han River. Before the developing the model, self-organizing map (SOM) and artificial neural network (ANN) restored missing values in input and output parameters. Then, an ANN model was trained to predict TTHMs by using relevant water quality parameters investigated from Pearson correlation. Pearson Correlation test selected six significant input parameters such as temperature, algae, pre-middle chlorine, post chlorine, total chlorine, and total organic carbon. Based on five-fold jackknife cross-validation, the ANN models built using different types of input data showed different performance in training (range of  $R^2$  from 0.62 to 0.92) and validation (range of  $R^2$  from 0.62 and 0.80) steps. This study can be a useful tool for predicting TTHMs concentrations using conventional water quality data in drinking water treatment plants. Machine learning models can be readily developed and utilized by managers working with drinking waters.

*Keywords:* Trihalomethanes (THMs); Drinking water treatment plant; Han River; Machine learning technique

---

\*Corresponding author.