

On Robustness of the Explanatory Power of Machine Learning Models

Banamali Panigrahi¹, Saman Razavi¹, Lorne E Doig¹, Blanchard Cordell², Hoshin V. Gupta³, and Karsten Liber¹

¹University of Saskatchewan

²University of Saskatchewan

³The University of Arizona

March 15, 2024

Abstract

Machine learning (ML) is increasingly considered the solution to environmental problems where only limited or no physico-chemical process understanding is available. But when there is a need to provide support for high-stake decisions, where the ability to explain possible solutions is key to their acceptability and legitimacy, ML can come short. Here, we develop a method, rooted in formal sensitivity analysis (SA), that can detect the primary controls on the outputs of ML models. Unlike many common methods for explainable artificial intelligence (XAI), this method can account for complex multi-variate distributional properties of the input-output data, commonly observed with environmental systems. We apply this approach to a suite of ML models that are developed to predict various water quality variables in a pilot-scale experimental pit lake.

A critical finding is that subtle alterations in the design of an ML model (such as variations in random seed for initialization, functional class, hyperparameters, or data splitting) can lead to entirely different representational interpretations of the dependence of the outputs on explanatory inputs. Further, models based on different ML families (decision trees, connectionists, or kernels) seem to focus on different aspects of the information provided by data, although displaying similar levels of predictive power. Overall, this underscores the importance of employing ensembles of ML models when explanatory power is sought. Not doing so may compromise the ability of the analysis to deliver robust and reliable predictions, especially when generalizing to conditions beyond the training data.

Hosted file

Panigrahi_G-VARS_XAI paper_2024_03_09.docx available at <https://authorea.com/users/755307/articles/725108-on-robustness-of-the-explanatory-power-of-machine-learning-models>

Hosted file

Panigrahi_G-VARS_XAI paper_Supplement_2024_03_09.docx available at <https://authorea.com/users/755307/articles/725108-on-robustness-of-the-explanatory-power-of-machine-learning-models>

On Robustness of the Explanatory Power of Machine Learning Models

Banamali Panigrahi¹, Saman Razavi^{2,3,4,*}, Lorne E. Doig¹, Blanchard Cordell⁴, Hoshin V Gupta⁵, and Karsten Liber^{1,2}

¹Toxicology Centre, University of Saskatchewan, Canada.

²School of Environment and Sustainability, University of Saskatchewan, Saskatoon, Canada.

³Institute for Water Futures, Mathematical Sciences Institute, Australian National University, Canberra, Australia.

⁴Global Institute for Water Security, School of Environmental and Sustainability, University of Saskatchewan, Saskatoon, Canada.

⁵Department of Hydrology and Atmospheric Sciences, The University of Arizona, Tucson, Arizona, USA.

*Address correspondence to saman.razavi@usask.ca

Abstract:

Machine learning (ML) is increasingly considered the solution to environmental problems where only limited or no physico-chemical process understanding is available. But when there is a need to provide support for high-stake decisions, where the ability to *explain* possible solutions is key to their acceptability and legitimacy, ML can come short. Here, we develop a method, rooted in formal *sensitivity analysis* (SA), that can detect the primary controls on the outputs of ML models. Unlike many common methods for *explainable artificial intelligence* (XAI), this method can account for complex multi-variate distributional properties of the input-output data, commonly observed with environmental systems. We apply this approach to a suite of ML models that are developed to predict various water quality variables in a pilot-scale experimental pit lake.

A critical finding is that subtle alterations in the design of an ML model (such as variations in random seed for initialization, functional class, hyperparameters, or data splitting) can lead to entirely different representational interpretations of the dependence of the outputs on explanatory inputs. Further, models based on different ML families (decision trees, connectionists, or kernels) seem to focus on different aspects of the information provided by data, although displaying similar levels of predictive power. Overall, this underscores the importance of employing ensembles of ML models when explanatory power is sought. Not doing so may compromise the ability of the analysis to deliver robust and reliable predictions, especially when generalizing to conditions beyond the training data.

Key Points:

- We extend the sensitivity analysis (SA) paradigm to handle complex multivariate distributions encountered in machine learning (ML).
- We apply our new SA-based method to explain the controls of various ML models developed for water quality predictions.
- We show different how ML models may rely on entirely different predictors and data signals despite exhibiting comparable predictive power.

1. Introduction

Machine learning (ML) is increasingly used in various domains, and success stories abound for problems that do not carry significant risks of negative consequences. However, when erroneous predictions can have major adverse implications for societal and environmental well-being, ML often faces acceptability and legitimacy challenges (Lipton, 2017; Lakkaraju et al., 2019; Read et al., 2019; Rudin, 2019; Samek et

41 *al.*, 2019; Coyle and Weller, 2020; Lakkaraju et al., 2020; Roscher et al., 2020; Slack et al., 2021; Slack et
42 *al.*, 2023). In such cases, a prime concern is often the difficulty in explaining the reasoning behind an ML
43 model's predictions, which often leads decision makers to favor process-based models (PBMs) that rely
44 on representations that encode a physico-chemical understanding of the underlying system (Hipsey et al.,
45 2015; Fatichi et al., 2016; Read et al., 2019; Razavi, 2021; Razavi et al., 2022). However, for many
46 emerging, site-specific environmental problems, such PBMs are not yet readily available.

47 In recent years, there has been a growing emphasis on the need for *explainable artificial intelligence* (XAI)
48 methods. Various approaches to this have been developed, including *Partial Dependence Plots* (PDP;
49 Friedman, 1991), *Permutation Feature Importance* (PFI; Strobl et al., 2008), *Local Interpretable Model-*
50 *agnostic Explanations* (LIME; Ribeiro et al., 2016), and *SHapley Additive exPlanations* (SHAP; Lundberg and
51 Lee, 2017). These methods are designed to elucidate the specific contributions by which individual feature
52 instances lead to particular outputs (i.e., local interpretation), or to uncover how features collectively
53 influence model outputs across all instances (i.e., global interpretation).

54 While significant strides have recently been made in XAI, they continue to suffer from limitations. For
55 instance, the widely-used SHAP technique, which is based on game theory, is (1) limited in its scalability
56 to large, high-dimensional datasets due to computational constraints (Molnar, 2020; Molnar, 2022; Stein
57 et al., 2022), (2) unable to capture interactions (especially higher-order ones) between features, which
58 may limit its ability to provide comprehensive explanations (Kumar et al., 2020; Puy et al., 2022), and
59 further may (3) assign excessive importance to improbable instances, potentially leading to unreliable
60 outcomes (Molnar, 2022; Rudin et al., 2022). Various strategies to tackle these limitations have had
61 varying degrees of success (Owen, 2014; Strumbelj and Kononenko, 2014; Mase et al., 2019; Do and Razavi
62 et al., 2020; Frye et al., 2020; Janzing et al., 2020; Lundberg et al., 2020; Sheikholeslami et al., 2021; Aas
63 et al., 2021; Liu et al., 2024).

64 More recently, *sensitivity analysis* (SA) has emerged as an alternative approach to XAI (Razavi et al., 2021;
65 Scholbeck et al., 2023). SA is a relatively young discipline that aims to study how the outputs of a model
66 are related to, and are influenced by, its inputs and/or controlling factors (Saltelli et al., 2021). The
67 application of SA to ML models has gained traction in recent years, as evidenced by studies such as those
68 by Tunkiel et al. (2020), Paleari et al. (2021), Fel et al. (2021), Kuhnt and Kalka (2022), Ojha et al. (2022),
69 and Stein et al. (2022). Unlike conventional XAI methods like SHAP, which focus primarily on individual
70 data instances to evaluate feature importance, SA takes a broader approach by seeking to characterize
71 the entire '*response surface*' of a model – the hyperplane that maps the input variables onto the output
72 of interest. Consequently, the computational demand of SA can be independent of the dataset size used
73 for model training/calibration.

74 Various SA methods have been developed across different application disciplines. Broadly categorized in
75 Razavi et al. (2021), these methods fall under four main approaches: *derivative-based* (Morris, 1991;
76 Campolongo et al., 2007; Sobol' and Kucherenko, 2009; Lamboni et al., 2013; Rakovec et al., 2014;
77 Kucherenko and Iooss, 2016; Kucherenko and Song, 2016), *distribution-based* (Sobol', 1993; Owen, 1994;
78 Homma and Saltelli, 1996; Saltelli et al., 2008; Kucherenko and Song, 2016; Puy et al., 2021), *variogram-*
79 *based* (Razavi and Gupta, 2016b; Sheikholeslami and Razavi, 2020; Becker, 2020; Alipour et al., 2022), and
80 *regression-based* (Kleijnen, 1995; Kambhatla and Leen, 1997; Tonkin and Doherty, 2005; Shin et al., 2013).
81 These approaches offer different definitions of sensitivity, vary in computational demand, and exhibit
82 varying degrees of scalability to the input space dimension (Razavi and Gupta, 2015). While methods
83 under any of these approaches can, in theory, be applied to characterize the importance of inputs in an
84 ML model, a significant challenge arises when dealing with correlated inputs following complex multi-
85 variable distributions. This issue is prevalent across the majority, if not all, SA methods and poses a
86 considerable hurdle in using SA for explainable ML.

87 Here, we introduce an SA-based method specifically tailored for XAI. We accomplish this by extending the
 88 *Variogram Analysis of Response Surface* (VARS) framework (Razavi and Gupta, 2016b) to accommodate
 89 input-output datasets characterized by complex, multi-variable distributions commonly encountered in
 90 ML applications. VARS is a variogram-based method known for its high computational efficiency, even in
 91 high-dimensional problems (Becker, 2020). It stands out as the only method that considers crucial
 92 information regarding the structure of the response surface and perturbation scale (Haghnegahdar and
 93 Razavi, 2017). Consequently, our new SA method is well equipped to address the three challenges
 94 commonly encountered in XAI, as outlined above. It achieves this by being independent of the available
 95 data size, adept at handling correlated inputs with any complex marginal distributions, and capable of
 96 directly operating on the response surface, thereby ensuring robustness against improbable areas of input
 97 space.

98 We test this approach across a suite of ML models based on decision trees, connectionism, and kernels,
 99 as a possible solution to investigate the processes in a pilot pit lake in the Athabasca Oil Sands region of
 100 Western Canada. This pilot lake contains fluid tailings treated using the permanent aquatic storage
 101 structure process, capped with a blend of oil sands process-affected water and runoff water from the
 102 surrounding landscape. We show how the new SA approach illuminates the key controls of different ML
 103 models in this environmental system. We show, in particular, that while different ML models may
 104 demonstrate similar predictive power, they may do so based on fundamentally different signals and
 105 patterns extracted from the data. We also show that ML models based on the connectionism paradigm
 106 (including deep learning) may not necessarily be robust to randomness in their initialization, so that
 107 multiple replicates of the same model trained to the same data might utilize different underlying
 108 relationships to predict the output. We discuss how the understanding of hidden differences across
 109 different ML models is critical to enabling learning about physico-chemical processes in the systems under
 110 investigation, and to ensure that any decision made on this basis is supported by well-justified
 111 explanations.

112 2. The Sensitivity Analysis (SA) Method

113 We introduce a general approach, grounded in sensitivity analysis (SA), to illuminating the workings of
 114 any model, even a black box, by assessing the extent to which different inputs influence its outputs. This
 115 SA-based approach is especially applicable to the problem of “*explainability*” in artificial intelligence (AI),
 116 because it addresses three common challenges encountered in machine learning (ML) applications, as
 117 detailed in Section 1. In the rest of this section, we present an overview of the underlying VARS-based
 118 framework, highlighting its computational efficiency even for models with high-dimensional input spaces.
 119 Subsequently, we demonstrate how we extend this framework to handle models with correlated inputs
 120 characterized by complex marginal distributions.

121 2.1. Variogram Analysis of Response Surfaces (VARS)

122 The VARS framework, originally developed by Razavi and Gupta (2016a and b), aims to characterize the
 123 entire ‘*response surface*’ of a model by integrating the directional variograms of a model output over the
 124 entire input space and across the full range of ‘*perturbation scales*’, h , as follows:

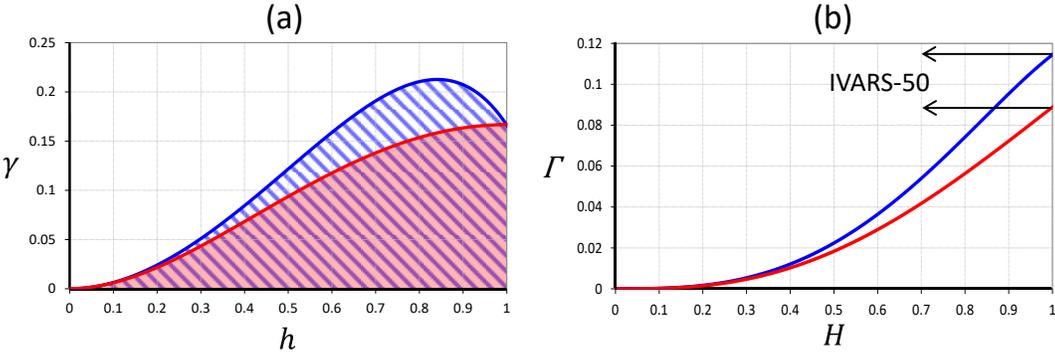
$$125 \quad \gamma(h) = \frac{1}{2}V[Z(\boldsymbol{\theta} + h) - Z(\boldsymbol{\theta})] \quad (1)$$

$$126 \quad \Gamma(H) = \frac{1}{2} \int_0^H V[Z(\boldsymbol{\theta} + h) - Z(\boldsymbol{\theta})]dh \quad (2)$$

127 where $\gamma(\cdot)$ and $\Gamma(\cdot)$ are directional variogram and integrated variogram functions, respectively, $\boldsymbol{\theta} =$
 128 $\{\theta_i \text{ for } i = 1, \dots, n\}$ and Z represent a point in the input space and its respective model output, n is the
 129 total number of inputs, $V[\]$ denotes the variance operator, and H is the range of perturbation scales of
 130 interest.

131 Figure 1 provides an illustrative graphical representation of directional variograms and their integrated
 132 versions. Directional variograms characterize the variance of change in the response surface ($Z(\theta + h) -$
 133 $Z(\theta)$) as a function of perturbation scale (h). For small values of h , this variance of change resembles
 134 information akin to derivative-based SA for different inputs, while for larger values, variograms offer
 135 insights into the variance contribution of each input, akin to variance-based SA. Thus, VARS serves as a
 136 unifying theory bridging derivative-based and variance-based SA, while offering a spectrum of information
 137 on the response surface structure for all other values of h . When H is 50 percent of the input range, the
 138 respective integrated variogram (IVARS-50) is called the 'total-variogram effect'. This measure of input
 139 importance encapsulates sensitivity information across the full spectrum of perturbation scales; see
 140 [Razavi and Gupta \(2016a\)](#) and [Haghnegahdar and Razavi \(2017\)](#) for further details.

141



142

143

144

145

146

147

Figure 1. Illustrative example of (a) directional variograms and (b) integrated variograms for a hypothetical model with only two inputs. This example is adopted from Example 1a of [Razavi and Gupta \(2016a\)](#), where the range of inputs (θ) is two, resulting in a range of perturbation scales (h) of one, which is half of the input range.

148

149

150

151

152

153

154

Various studies have been shown VARS to be highly efficient and statistically robust (e.g., [Razavi and Gupta, 2016b](#); [Alipour et al., 2022](#); [Becker, 2020](#)). This efficiency is partly attributed to the estimation method used to calculate directional variograms, which relies on pairs of sample points rather than individual sample points ([Razavi and Gupta, 2016a & b](#)), thereby exploiting the fact that the number of pairs grows geometrically with an increase in the number of samples. As a result, VARS has been proven capable of effectively accommodating high-dimensional problems ([Sheikholeslami et al., 2019](#)). In the following sub-section, we expand upon this framework to adapt it for use in the context of XAI.

155 **2.2. VARS with complex input distributions**

156

157

158

159

160

161

162

A vast majority of SA methods and their applications in the literature operate under the assumption that model inputs are independent and uniformly distributed ([Do and Razavi, 2020](#)). This simplifying assumption is often made for the sake of computational convenience, as (1) accurately characterizing the multivariate distribution of inputs can be challenging or impractical in many cases, and (2) even if such distributions exist, incorporating them into the analysis may introduce computational complexity or feasibility issues. To address these limitations, [Do and Razavi \(2020\)](#) developed *Generalized VARS* (G-VARS), one of the first SA methods capable of accommodating correlated inputs.

163

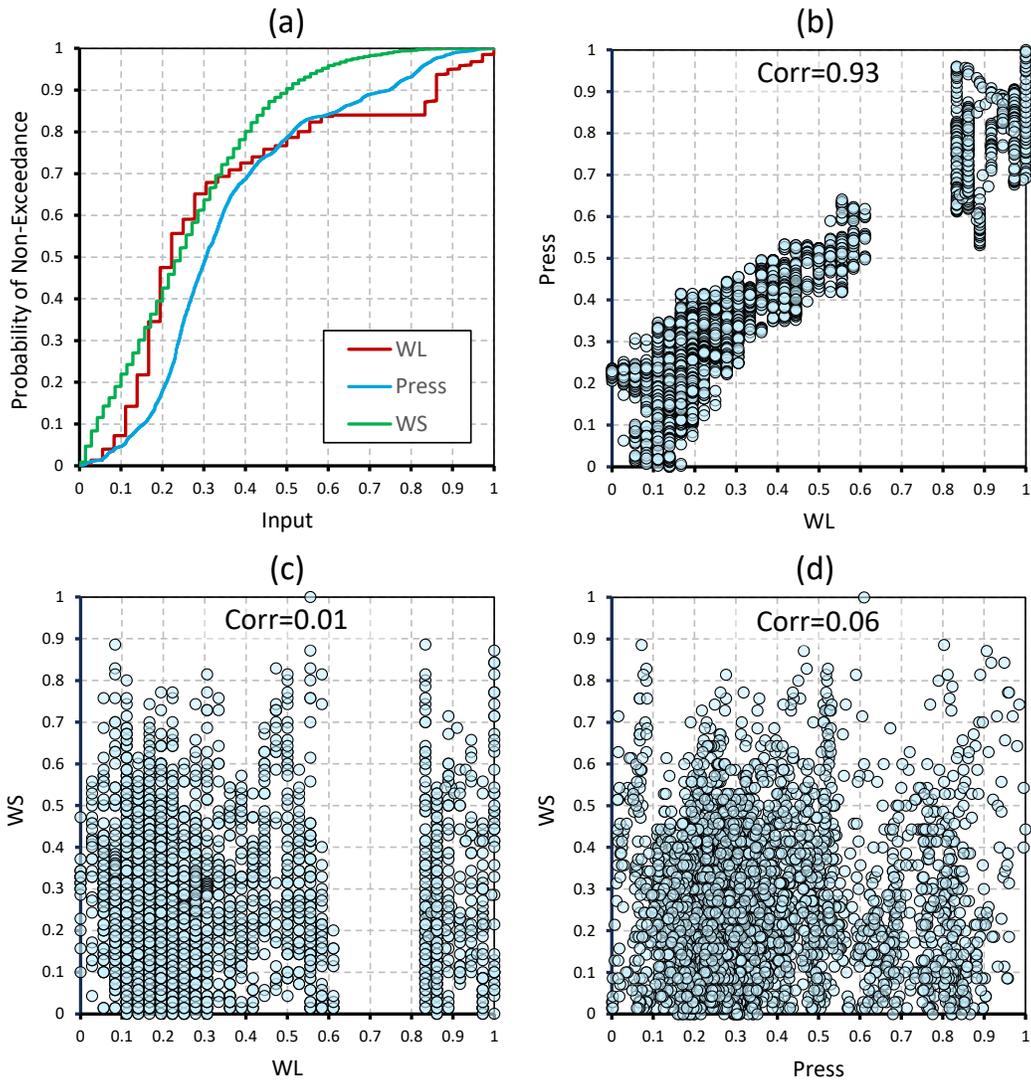
164

165

G-VARS involves novel sampling and estimation strategies that map the original input space onto a standard normal space, while accounting for pair-wise correlations between inputs using the *Nataf Isoprobabilistic Transformation*. However, G-VARS is constrained to handling simple, theoretical

166 multivariate distributions such as normal and triangular distributions. While these distributions can be
167 quite effective in characterizing various variables in modeling exercises, such as the parameters of a
168 hydrologic model as demonstrated by *Do and Razavi (2020)*, they may offer limited assistance in
169 characterizing highly complex datasets typical of ML inputs.

170 Figure 2 shows an illustrative three-variable example of datasets commonly encountered in practice. The
171 marginal distributions of real-world variables (e.g., Figure 2a) may be quite complex, exhibiting features
172 such as multiple modes, discontinuities, and long or heavy tails, while potentially being correlated with
173 one another (e.g., Figures 2b-c). Approximating such complex multivariate distributions with classic
174 theoretical distributions can often prove infeasible or impractical. Here, we developed a mathematical
175 construct, called G-VARS2, that adopts an *empirical* approach to represent the marginal distributional
176 properties of the data, as detailed below.



177
178 **Figure 2.** Example datasets used as inputs in machine learning. (a) Marginal cumulative
179 distribution functions (CDFs). (b-d) Scatter plots and Pearson correlation coefficients (Corr) of
180 pairs of inputs. WL, Press, and WS stand for water level, water pressure at the sediment-water
181 interface, and wind speed, respectively. The data shown are from the case study described in the
182 next section.

183

184 Suppose a model receives a set of n inputs, $\theta = \{\theta_i \text{ for } i = 1, \dots, n\}$, that follow a multivariate
185 distribution, $p(\theta)$, and produces an output, $Z(\theta)$. Let $\theta_{\tilde{i}}$ denote the set of all inputs except θ_i , and θ'_i
186 denote $\theta_i + h_i$, where h_i is a perturbation to θ_i . Now, Equation (1) can be rewritten as:

$$187 \quad \gamma(h) = \frac{1}{2}V[Z(\theta'_i | \theta_{\tilde{i}}, \theta_{\tilde{i}}) - Z(\theta_i | \theta_{\tilde{i}}, \theta_{\tilde{i}})] \quad (3)$$

188 where $\theta_{\tilde{i}}$ follows $p(\theta_{\tilde{i}})$, which is the marginal distribution derived from $p(\theta)$, and where $\theta'_i | \theta_{\tilde{i}}$ and
189 $\theta_i | \theta_{\tilde{i}}$ follow the conditional distributions when the inputs $\theta_{\tilde{i}}$ are set at specific values. Following the
190 proof in *Razavi and Gupta (2016a & b)*, the above equation can be numerically approximated by sampling
191 pairs of points from the model input-output space:

$$192 \quad \hat{\gamma}(h_i) = \frac{1}{2N_h} \sum_1^{N_h} [Z(\theta'_i | \theta_{\tilde{i}}, \theta_{\tilde{i}}) - Z(\theta_i | \theta_{\tilde{i}}, \theta_{\tilde{i}})]^2 \quad (4)$$

193 where N_h is the number of pairs of samples, spaced h_i apart, in the direction of θ_i .

194 Here, we adjust the sampling method of *Do and Razavi (2020)* to accommodate any ‘custom’ marginal
195 distribution of $p(\theta)$ that may exist in real-world data. The new method utilizes any available sample of
196 data for individual inputs to construct their empirical distributions, by calculating the frequencies of
197 different values or ranges of values from the data sample. Empirical distributions provide a data-driven
198 summary of the observed data’s distributional properties when the underlying theoretical distribution is
199 unknown, or is difficult to model accurately.

200 The new method processes a data sample for each input θ_i to derive its cumulative distribution function
201 (CDF) through the *Weibull* empirical approach, by sorting data entries in ascending order and assigning
202 each entry a probability of non-exceedance. Subsequently, the actual CDF of input θ_i , denoted as F_{θ_i} , is
203 estimated by linearly interpolating the points on the respective empirical CDF. The lower and upper
204 bounds of a custom-distributed input are assumed to be the minimum and maximum values of the
205 corresponding sample. Next, the inverse of F_{θ_i} for all inputs ($F_{\theta_i}^{-1}$) is incorporated into the G-VARS
206 framework through the following equation:

$$207 \quad \theta_i = F_{\theta_i}^{-1}[\phi(X_i)] \quad (5)$$

208 facilitating the transformation of samples between a standard normal space $X = \{X_i \text{ for } i = 1, \dots, n\}$ and
209 the original input space θ , where $\phi(\cdot)$ denotes the theoretical CDF of a standard normal distribution. The
210 software developed for G-VARS2 is accessible on GitHub at the following link: [https://github.com/vars-
211 tool/vars-tool](https://github.com/vars-tool/vars-tool).

212 3. Data used in ML: Pilot Scale Pit Lake

213 In the *Athabasca Oil Sands* (AOS) region of *Western Canada*, the accumulation of fluid tailings (FT) and Oil
214 Sands Process affected Water (OSPW) in tailing ponds has reached a concerning level that has attracted
215 global attention (*Gosselin et al., 2010; AEP, 2015; McNeill, 2017*). Despite decades of research on several
216 experimental reclamation techniques for fluid tailings management in the oil sand regions (*COSIA, 2012*),
217 there is still a need for an advanced pit lake technology under water capped fined deposit to handle the
218 substantial amount of FT (*Cossey et al., 2021*). To address this, one company in the AOS industries has
219 developed a pilot-scale experimental pit lake called *Lake Miwasin* as a prototypic precursor to large end
220 pit lakes.

221 Extensive measurements of water quantity and quality variables have been ongoing since the construction
222 of the lake to evaluate the performance of the system over time and its effectiveness in reclaiming
223 significant quantities of treated tailings materials stored onsite at the AOS. Our group has been using
224 wireless sensor technology to monitor water quality parameters at a high measurement frequency, to

225 gain a deeper understanding of the system's functioning, and to help guide further development and use
226 of pit lake strategies to mitigate the negative impacts of fluid waste on the environment.

227 Located on the east side of the *Athabasca River* (56° 53' 14" N and 111° 23' 7" W), Alberta, Canada, *Lake*
228 *Miwasin* is an engineered water body constructed using oil sand by-products as bottom substrate (TFT)
229 and overlying water (OSPW). The lake measures 70 m in width, 165 m in length and reaches a maximum
230 water depth of around 4.5 m. A littoral zone in the lake, comprising approximately 20% of the water
231 surface area, is present along the eastern periphery. This zone features a 0.2-5% slope that extends into
232 the surrounding upland, where limnetic zone has deep depth with substantial bottom substrate.

233 The wireless sensor network (WSN) was employed in both littoral and limnetic zones of the lake to monitor
234 (hourly) and relay information on key water quality parameters. The *Lake Miwasin* WSN utilized the
235 *Libelium™* smart water extreme and smart water ion models for monitoring. The sensor probes were
236 calibrated with standard solutions from *Libelium* (*Libelium, 2020*). Data was transmitted directly from the
237 sensor device to cloud storage through the *ThingSpeak™ Cloud* service and mobile devices using the local
238 4G network (via mobile SIM card). The sensor units, housed in custom acrylic boxes (30 × 25 × 25 cm) for
239 protection against field conditions, were affixed to high-density Styrofoam platforms (60 × 60 × 5 cm for
240 each sensor unit) using cinderblock anchors (Figure 3). Figure 3a illustrates one-time introduction of
241 TFT and OSPW stemming from the oil sands mining and extraction processes, Figure 3b shows
242 contemporary methods employed for treatment and reclamation of fluid tailings within AOS region,
243 and Figure 3c shows WSN technology was deployed in Lake Miwasin to monitor the water quality
244 conditions of the lake as the lake system ages.

245 Before deployment at *Lake Miwasin*, a similar WSN methodology was tested in Canadian lakes to
246 delineate effluent exposure downstream of a Uranium Mill region (*Cupe-Flores et al., 2022*) and to
247 estimate selenium (Se) exposure using a site-specific threshold value (*Peixoto Mendes et al., 2023*). In
248 *Lake Miwasin*, deployment spanned from September 18th to October 10th in 2020; and from June 21st to
249 October 16th in 2021. Coinciding with visits to the lake for probe maintenance, we collected water samples
250 at each station, twice in 2020 and six times in 2021. For validation of sensor reading, samples were
251 collected manually in two replicates for each monitoring depth approximately 5 to 10 m apart from probes
252 using a Wildco® 2.2-L acrylic Van Dorn horizontal beta water sampler (Wildlife Supply, USA).

253 Prior to sampling, the sampler was thoroughly rinsed with lake surface water to avoid cross-contamination
254 between different station zones. The water sample was placed into pre-acid-washed ~250 mL and 30-mL
255 high-density polyethylene *Nalgene™* bottles (prewashed with 10% nitric acid and rinsed with distilled
256 water). Subsequently, the 30-mL bottle samples were filtered through a 0.45-µm polyether sulfone
257 membrane into two sets of 10-mL high-density polyethylene *Nalgene™* bottles using 5-mL syringes. Then
258 samples were refrigerated and transported in an ice-packed cooler to the *University of Saskatchewan*
259 *Toxicology Centre* (Saskatoon, SK, Canada) and kept at 4°C until laboratory analysis. A *Thermo Scientific™*
260 *Orion Star™ A329* portable multiparameter meter (Thermo Fisher Scientific, USA) was used to measure in
261 situ parameters (pH, EC, temperature, and DO) at a pre-defined monitoring depth. Similarly, the turbidity
262 was measured with a calibrated bench top turbidity meter (LaMotte®, 2020 meter). Field blanks
263 containing distilled water were included during sampling for quality control. Water quality parameters
264 measured by the sensor probes by our research team at University of Saskatchewan.

265 Additional climatological parameters included in this study were collected by Suncor Energy Inc. from the
266 meteorological station installed at the lake. In the context of studying this lake system, the focus was on
267 predicting key water quality parameters through the application of kernel-based, connectionist models,
268 and ensemble tree-based ML techniques.

269

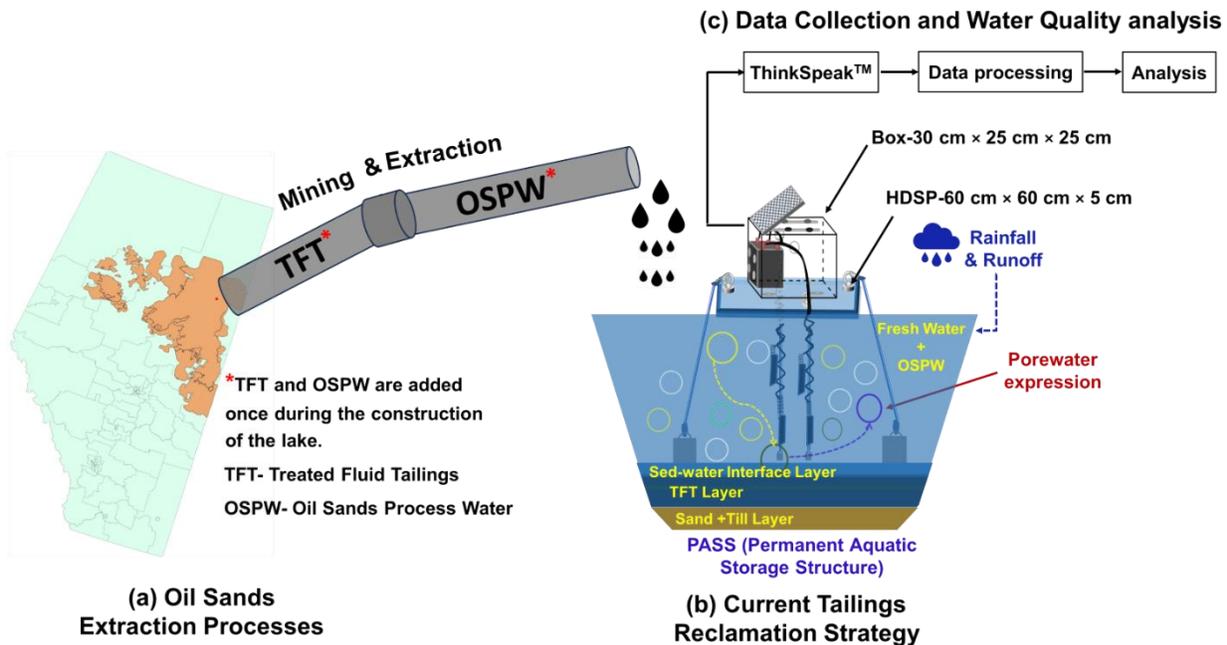


Figure 3. Graphical description of elements used for characterizing key water quality parameters in Lake Miwasin.

270
271
272
273

274 4. Design of Experiments

275 We evaluated the proposed SA method (G-VARS2) across various ML models, encompassing tree-based,
276 connectionist, and kernel-based models. Our objective was to gauge the robustness and consistency of
277 these models and to glean insights into controlling variables in the prediction process. This experiment
278 allowed us to not only develop accurate predictive ML models but also to provide transparent and
279 interpretable explanations regarding the inputs driving the predictions made by these models (Figure 4).
280 Similar to G-VARS, the user of G-VARS2 must set the following two algorithm parameters for sampling and
281 estimation: the number of star centers and the number of cross-sectional points. For this study, we
282 selected 100 and 10 for the former and latter, respectively. The selection of these numbers was informed
283 by our prior experience, and they have consistently shown robustness and stability in our initial
284 evaluations.

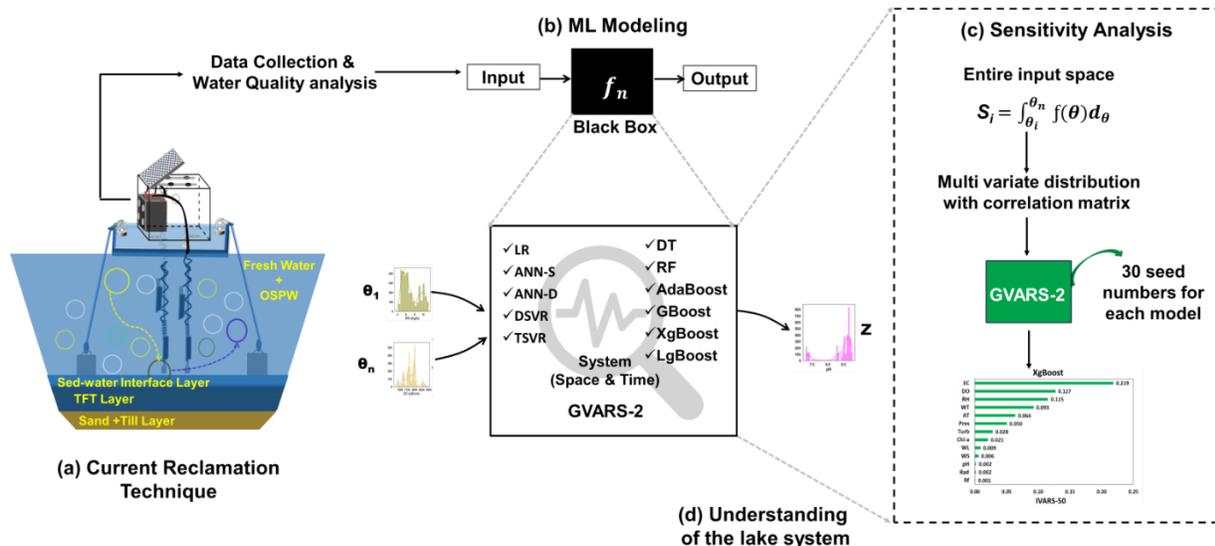
285 The tree-based ML models we used here include classic decision trees (DT), bagging-based random forest
286 (RF; [Breiman, 2001](#)), and boosting models including adaptive boosting (AdaBoost; [Freund and Schapire, 1997](#)),
287 gradient boosting (GBoost; [Friedman, 2001](#)), extreme gradient boosting (XgBoost; [Chen and Guestri, 2016](#)),
288 and light gradient boosting (LgBoost; [Ke et al., 2017](#)). Further, we used two artificial neural
289 networks (ANNs), one with a shallow, single hidden layer (ANN-S) and another with a deeper architecture
290 (ANN-D), and two support vector regression (SVR) models, one with default (DSVR) and another with
291 tuned parameterization (TSVR) (Table S1). We also used multiple linear regression (LR) to provide a
292 minimalist performance as a benchmark for model comparison.

293 We chose to assess the above ML models (Figure 4b) because they have been used widely in the realm of
294 surface water quality modeling ([Palani et al., 2008](#); [Ahmed, 2014](#); [Najah et al., 2014](#); [Ahmed et al., 2019a](#);
295 [Ahmed et al., 2019b](#); [Abobakr Yahya et al., 2019](#); [Banerjee et al., 2019](#); [Sinshaw et al. 2019](#); [Zou et al., 2019](#);
296 [Barzegar et al., 2020](#); [Yim et al., 2020](#); [Khullar and Singh, 2021](#); [Yamamoto et al., 2021](#); [Zhang et al., 2021](#);
297 [Zhou and Zhang, 2023](#)). Moreover, there is a growing literature on the application of advanced

298 bagging-boosting ensemble models based on DT in environmental sciences. These applications include
 299 estimating crop yields (*Liakos et al., 2018*), assessing energy performance (*Wang et al., 2018*), predicting
 300 water demand (*Wang et al., 2014*), estimating air particulate levels (*Brokampet et al., 2018*), quantifying
 301 climate and catchment control on hydrological drought (*Konapala and Mishra, 2019*), creating
 302 susceptibility maps for gully erosion (*Rahmati et al., 2017; Garosi et al., 2019*), mapping groundwater yield
 303 (*Sameen et al., 2019; Jelhouni et al., 2020; Mosavi et al., 2021*), predicting solar and wind energy (*Torres-*
 304 *Barrán et al. 2019*), forecasting water usage and rainfall (*Kim et al., 2020*). Although bagging-boosting
 305 models hold significant promise, their utilization in surface water quality research remains relatively
 306 limited (*Chen et al., 2020*). Examples include predicting biological oxygen demand, chemical oxygen
 307 demand (*Khullar and Singh, 2021*), turbidity (*Zhang et al., 2021*), Chlorophyll-a (*Savoy and Harvey, 2023*)
 308 and other water quality parameters such as permanganate index (COD_{Mn}), total phosphorus (TP), and total
 309 nitrogen (TN) (*Wang et al., 2021*), dissolved oxygen (DO), and ammonia (NH₃-N) (*Chen et al., 2020*).

310 In our analyses, we further accounted for the impact of data-splitting for training and testing of the
 311 different ML models as well as randomization in initializing the ML models. This is a very important, but
 312 often neglected step in model development, as described in *Maier et al. (2023)*. To do so, we developed
 313 30 replicates of every ML model (Figure 4c), each with a different random seed number for data-splitting
 314 and model initialization (Table S1). The detailed default and tuned hyperparameter values for all the ML
 315 models are presented in Table S2. We employed a standard approach to partition the datasets into two
 316 subsets for each ML model: 70% for the training dataset and 30% for the testing dataset. During the
 317 development phase, the training dataset was utilized as the foundation for constructing models, while the
 318 testing dataset served the critical role of evaluation by enabling performance comparisons among the
 319 developed models. The normalization and minimum-maximum scaling of all input and output variables
 320 were performed using the scikit-learn pre-processing library in Python (*Pedregosa et al. 2011*).

321



322

323 **Figure 4.** An illustration of the workflow, spanning from sampling to machine learning (ML)
 324 modeling, and subsequently to sensitivity analysis (SA). LR: Linear Regression; ANN-S: Simple
 325 Artificial Neural Network; ANN-D: Deep Artificial Neural Network; DSVR: Default Support Vector
 326 Regression; TSVR: Tuned Support Vector Regression; DT: Decision Tree; RF: Random Forest;
 327 AdaBoost: Adaptive Boosting; GBoost: Gradient Boosting; XgBoost: Extreme Gradient Boosting;
 328 LgBoost: Light Gradient Boosting.

329

330 To gauge the accuracy of our ML-based water quality prediction models, we utilized two performance
331 metrics: the coefficient of determination (R^2) and the root mean squared error (RMSE) as follows:

$$332 \quad R^2 = \left(\frac{\sum_{i=1}^n [(O_i - \bar{O})(P_i - \bar{P})]}{\sqrt{[(\sum_{i=1}^n (O_i - \bar{O})^2)(\sum_{i=1}^n (P_i - \bar{P})^2)]}} \right)^2 \quad (6)$$

$$333 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (O_i - P_i)^2} \quad (7)$$

334 where O_i and P_i are the observed and predicted values, respectively, \bar{O} and \bar{P} are the mean of the
335 observed and predicted values, respectively, and n is the number of data points.

336 Water quality monitoring in *Lake Miwasin* at the SWI involves collecting, measuring, and analysing water
337 samples to understand their chemical and biological attributes (Sinshaw et al. 2019). We used the
338 following sensor measured water quality variables: dissolved oxygen concentration (DO), pH, water
339 temperature (WT), conductivity (EC), chlorophyll-a (Chl-a), turbidity (Turb), ammonium (NH_4^+). We also
340 used meteorological variables, including water level (WL), wind speed (WS), solar radiation (Rad), water
341 pressure (Press), air temperature (AT), rainfall (Rf), relative humidity (RH). Using these variables, we
342 constructed ML-based models to predict four key water quality variables— NH_4^+ , Chl-a, DO, and pH—
343 utilizing all other variables as predictors (see summary statistics in Table S3). The selected target variables
344 are pivotal for lake monitoring, as they play a fundamental role in assessing the health of aquatic
345 ecosystems, influencing the growth and respiratory capabilities of aquatic life (Wetzel, 2001; Sánchez et
346 al., 2006; Pena et al., 2010; Risacher et al., 2018; Barzegar et al., 2020). Moreover, these variables provide
347 decision-makers with vital data to address environmental challenges in a sustainable manner (Wu and Liu
348 2012; Wu and Chen, 2013).

349 5. Results and Discussion

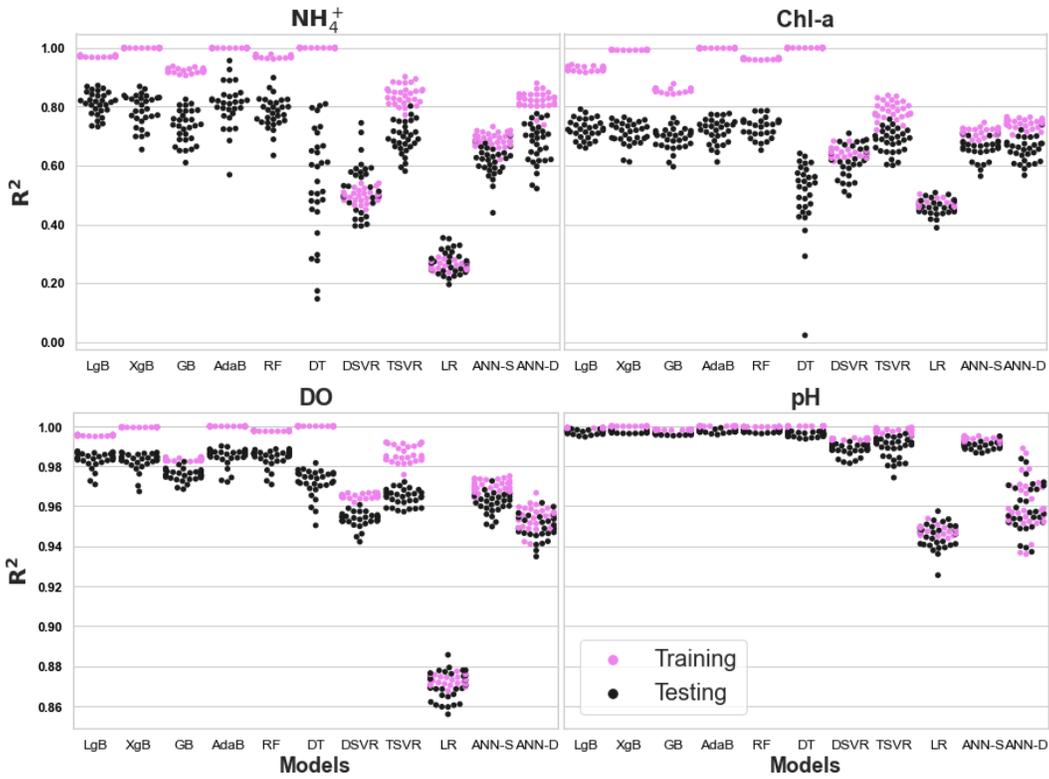
350 This section presents the outcomes of the designed experiments aimed at evaluating and explaining the
351 performance of various ML models using the developed SA method. We structure the results and
352 discussion around four key questions: (1) What is the predictive efficacy of different ML models? (2) Which
353 physico-chemical variables influence the predictions of ML models? (3) How robust and consistent is the
354 explanatory power of different ML models? (4) What novel insights do the ML models offer into the
355 underlying physico-chemical processes?

356 5.1. What is the predictive power of different ML models?

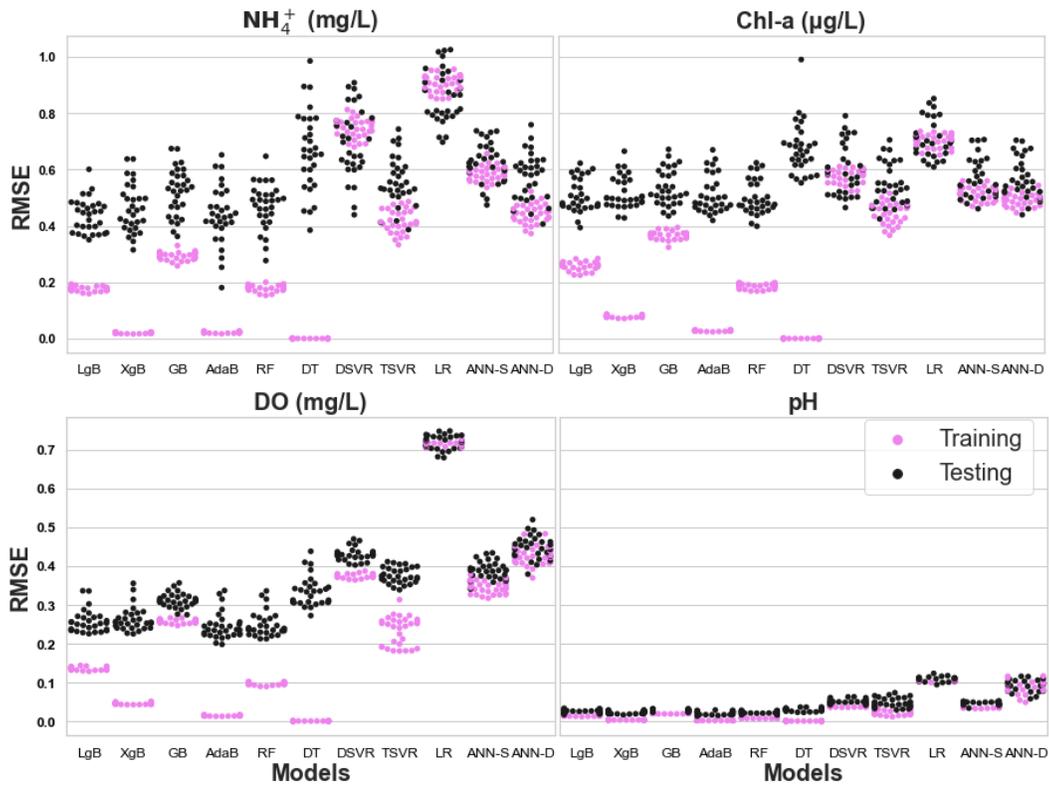
357 Figure 5 illustrates the predictive performance of various ML models across different target variables
358 during both the training and testing phases. Generally, these models demonstrated robust predictive
359 capabilities for pH, ranking second in performance for dissolved oxygen (DO), while exhibiting relatively
360 lower accuracy for predicting NH_4^+ and chlorophyll-a (Chl-a) levels. Notably, all ML models achieved
361 satisfactory performance for pH (with $R^2 > 0.93$ for linear regression (LR) and > 0.98 for other ML models)
362 and DO ($R^2 > 0.86$ for LR and > 0.9 for others). The LR model could only provide satisfactory performance
363 in the case of pH and DO prediction, suggesting that predicting NH_4^+ and Chl-a concentrations in the lake
364 rely predominantly on non-linear relationships, surpassing the capabilities of the LR model.

365

366



367



368

369

370

371

372

Figure 5. Performance of different ML models in training and testing across 30 replicates, each with a different random seed, according to the coefficient of determination (R^2) and root mean squared errors (RMSE). LgB: Light Gradient Boosting; XgB: Extreme Gradient Boosting; GB: Gradient Boosting;

373 AdaB: Adaptive Boosting; DT: Decision Tree; DSVR: Default Support Vector Regression; TSVR: Tuned
374 Support Vector Regression; L: Linear Regression; ANN-S: Simple Artificial Neural Network; ANN-D:
375 Deep Artificial Neural Network.

376
377 The results show that tree-based ML models outperformed the connectionist and kernel-based ML and
378 basic decision tree (DT) models. Throughout training, all tree-based models achieved R^2 scores remarkably
379 close to one. However, their performance declined notably on testing datasets (considerable reduction in
380 performance), particularly in NH_4^+ and Chl-a predictions. Importantly, all tree-based models with default
381 hyperparameter settings demonstrated acceptable performance, alleviating the need for extensive
382 hyperparameter tuning. This underscores the efficacy of these models in providing accurate predictions
383 for lake variables, albeit with some limitations in predicting certain parameters under testing conditions.

384 The variations observed in model performance depicted in Figure 5 across diverse seed numbers primarily
385 stem from the intrinsic randomness inherent in the model training processes, functional classes, and data
386 splitting. This susceptibility to initial conditions is a common trait among the ML models employed in this
387 study. For instance, the connectionist models such as ANN-S and ANN-D are influenced by variation in the
388 initial weights associated with different seeds, while variations in kernel-based SVR arise due to different
389 support vectors chosen in each model replicate. The basic DT model exhibits large variability in testing
390 due to the diverse tree structures resulting for different data splits. However, the ensemble tree-based
391 models such as bagging and boosting models were able to reduce those variations by building multiple
392 DTs and aggregating their predictions.

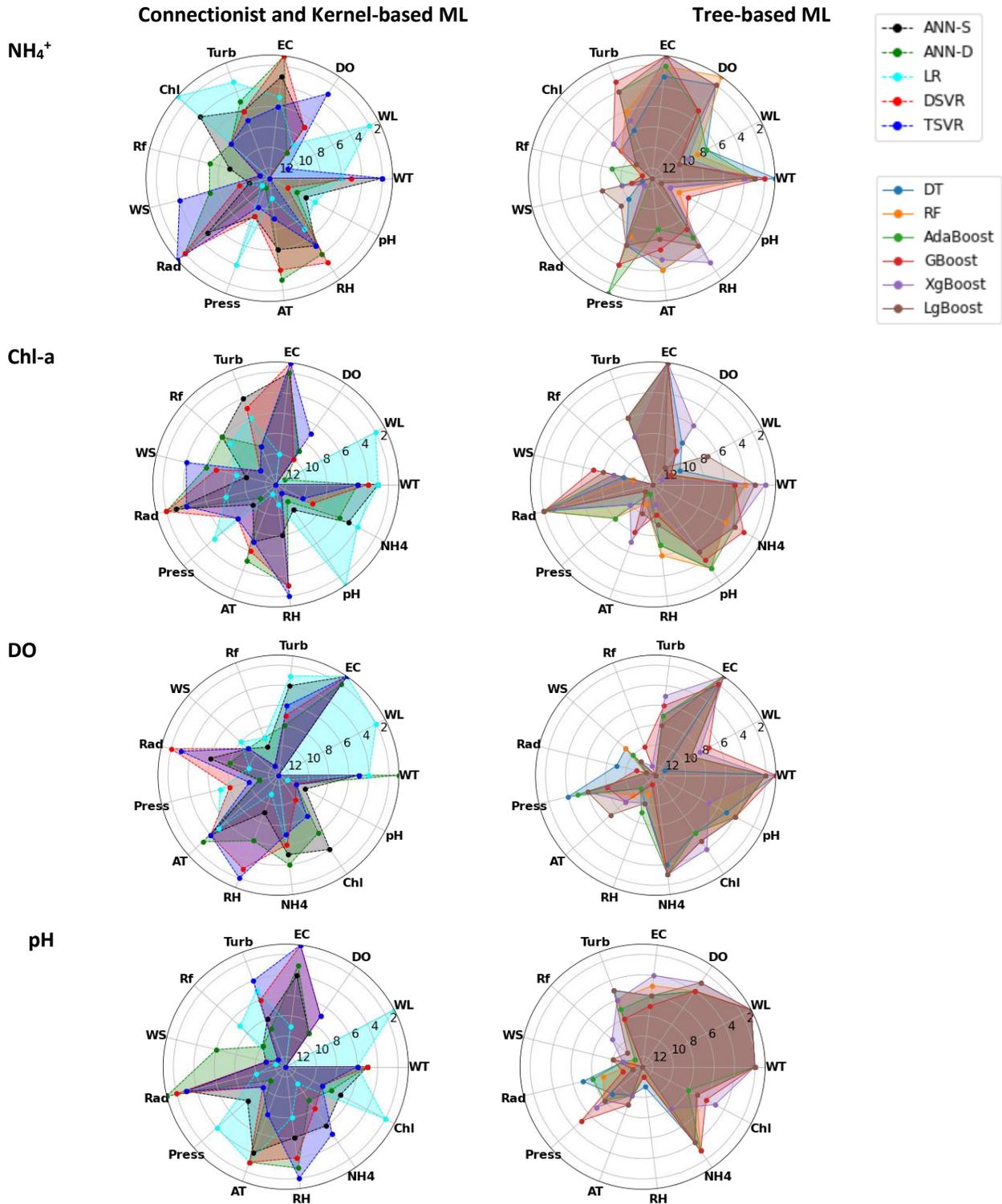
393 Another notable observation is the large reduction in model performance from training to testing for
394 certain models, particularly DT, AdaBoost, and XgBoost, underscoring the importance of thorough testing
395 and evaluation on datasets not utilized in training. More broadly, comprehending and addressing
396 potential variations in model performance resulting from various factors is crucial for recognizing the
397 uncertainty associated with the outcomes of ML models. The findings from the application of the
398 proposed SA-based method can provide further insights into this issue, as detailed in the subsequent
399 sections.

400 **5.2. What physico-chemical variables control the predictions of ML models?**

401 Most of the ML models showed success in mapping the inputs to the outputs in the pit lake, albeit to
402 varying degrees. In addition to providing predictive ability, these input-output mappings can potentially
403 offer a wealth of information on how the underlying physico-chemical processes work. However, such
404 mappings are typically comprised of very complex relationships that are hard to understand and explain.
405 Accordingly, use of SA to interrogate the ML models helps in characterizing the importance of the different
406 inputs on the functioning of the models to produce the output. Note that the SA method was run on each
407 of the 30 replicates of every ML model for each target output.

408 In Figure 6, star plots are to illustrate the overall importance (average over 30 replicates) of the inputs
409 into each of the ML models. The further a spoke extends outwards within the circle, the more influential
410 the respective input is in predicting the output. Accordingly, Rad and RH were identified as the most
411 influential inputs to the connectionist and kernel-based models for the predictions of NH_4^+ and Chl-a,
412 respectively. In contrast, for the tree-based models, DO and pH respectively turned out to be the most
413 influential inputs for prediction of NH_4^+ and Chl-a. For prediction of DO and pH, the connectionist and
414 kernel-based models were more sensitive to Rad, AT, and RH. The tree-based models, however, were
415 more sensitive to NH_4^+ , pH and Chl-a for the prediction of DO and to WL, WT and NH_4^+ for the prediction
416 of pH.

417



418 **Figure 6.** Input importance of different ML models (connectionist and tree-based) characterized
 419 through the proposed sensitivity analysis that is based on Integrated Variogram Across a Range of
 420 Scales (IVARS-50). Shown for each ML model is the median of input rankings across the 30 replicates.
 421 Rank 1 indicates the most influential input, rank 2 the second most influential input and so on.
 422

423 Furthermore, the importance of some inputs turned out to be quite different in the case of different ML
424 models. For example, EC was identified to be the most important input by DSVR to predict NH_4^+ , 3rd most
425 important by ANN-S, and the 7th most important by TSVR. Moreover, the rankings based on the linear
426 regression model were largely inconsistent with those based on other models, suggesting the existence
427 of strong non-linearity in the problems at hand. Overall, such considerable differences in input variable
428 importance indicate that, while different ML models may show comparable performance in terms of
429 predictive power, they may do so by relying on entirely different signals embedded in the training data
430 (Figures S1-S4).

431 We should note that prior to doing any modeling we observed from histogram plots and correlations that
432 water quality parameters such as conductivity (EC), Total Dissolved Solids (TDS), and Salinity have very
433 similar data distributions and are highly correlated (~ 0.99). A similar pattern was observed for DO and
434 Saturated oxygen concentration. If not addressed, this high level of correlation can lead to what is known
435 as "*substitution effects*," where different variables can essentially serve as substitutes for each other in
436 explaining the outcome. To tackle the problem of collinearity in the input dataset, we selected only EC
437 and DO to use as inputs and discarded the parameters that were highly correlated to them, finally
438 retaining only a few of the parameters (pH, Press and WL) that had a moderately high correlation of ~ 0.90 .
439 This choice to retain EC and DO is based on existing theory in the AOS region. It is worth noting that the
440 discarded parameters were computed empirically by the WSN system.

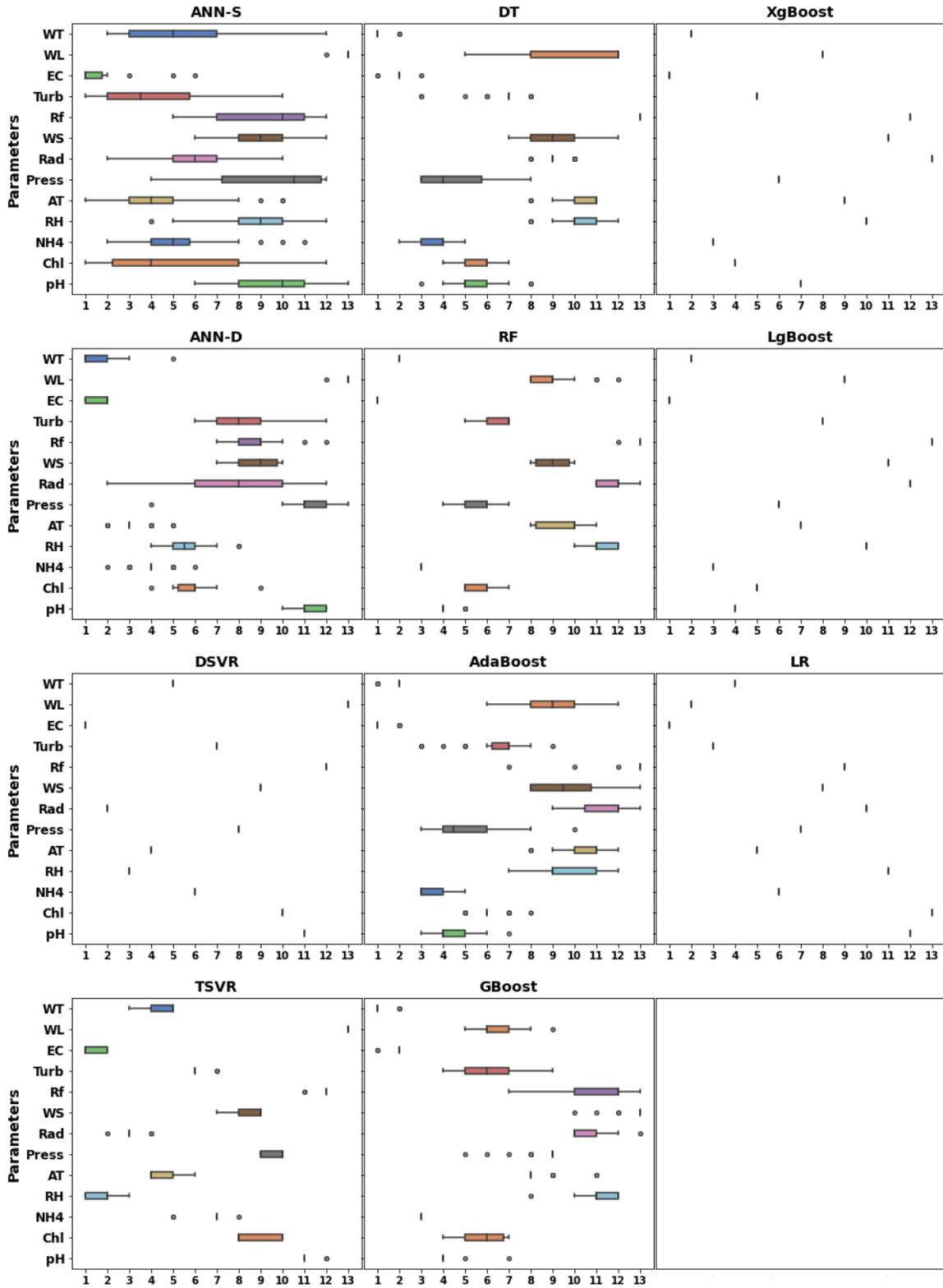
441 Note, however, that the novel reclamation technique being used in Lake Miwasin ages over time, causing
442 the lake to differ from its natural state in various ways, so we cannot solely rely on existing scientific
443 evidence. Moreover, our ML models, particularly RF and all Boosting models inherently address
444 substitution effects through regularization and ensemble techniques. Further, the various ML models
445 differ in their choice of explanatory variables, which is relevant to understanding these substitution
446 effects. Certain ML strategies are more susceptible to this issue while others are not when selecting
447 explanatory variables.

448 **5.3. How robust is the explanatory power of different ML models?**

449 By robustness, we refer here to low sensitivity of the ML results to the random seed chosen for the
450 randomization of their initial weights and for the splitting of available data to training and testing subsets.
451 Overall, the results of a more robust model are expected to be more reliable, particularly in the cases
452 where available data are limited. Here, we assessed model robustness by checking the dispersion of SA
453 results across 30 replicates. The more dispersed the input importance is across replicates, the less robust
454 is the model. Figure 7 shows the distribution of input ranks as perceived by different ML models across
455 the 30 replicates in predicting DO. Overall, the ANN-based models demonstrated the lowest robustness,
456 frequently relying upon different inputs as their primary controls for driving an output.

457 For example, ANNs showed a wide range of sensitivity to WT, Turb, and Chl when predicting DO, almost
458 as if the model may randomly pick up different predictive signals in the data each time it is set up. The
459 level of robustness demonstrated by ANN-D may be deemed comparable to that of ANN-S, although they
460 showed different predictive power in the previous sub-section. On the other hand, some ML models
461 showed high robustness, consistently identifying the same/similar inputs as their primary drivers for
462 prediction. For example, the variants of gradient boosting, such as GBoost, XgBoost and LgBoost models,
463 show comparatively less dispersion in behaviour over the 30 replicates, indicating more consistency in
464 their results.

465



466

467

468

469

Figure 7. Box plots showing the dispersity (range and distribution) of input importance for DO across the 30 replicates; Rank 1 indicates the most influential the respective input.

470 In general, the boosting models yielded very similar sets of “most influential” inputs but with slightly
471 different orders in their ranking. AdaBoost was an exception, showing large variability in rankings, due to
472 its random weight initialization process for model training. An interesting observation was the behaviour
473 of support vector regression with default hyperparameters (DSVR) versus that with hyper-parameter
474 tuning (TSVR). The former showed no variability in input rankings, while the latter showed some level of
475 variability due to hyperparameter tuning separately for each of the 30 replicates. This indicates that the
476 robustness of an ML model may be partly a manifestation of its inherent mechanisms in fitting data, and
477 partly of the way the user sets them up. Figure S3 in supplementary materials shows similar results for
478 ML models predicting the other target variables.

479 **5.4. What new insights do the ML models provide into the underlying physico-chemical processes?**

480 The SA-based explanation approach assessed the extent to which one variable can be impacted by other
481 variables, with the strength of these influences reflecting the underlying physical processes in *Lake*
482 *Miwasin*. Here, we compare the ML results for predicting NH_4^+ , Chl-a, DO and pH with established scientific
483 evidence from the oil sands region, primarily drawn from literature reviews and our research background.
484 By answering this research question, our objective was to highlight the explanations provided by ML
485 models that align with conventional wisdom and those that do not. Of course, it is always possible that
486 some new explanations may be discovered that identify relationships we are either unaware of or that
487 are not recognized by existing theories. This exploration may yield two possible scenarios, one where the
488 ML model may be providing the right answer for wrong reason and another where it challenges the
489 validity of existing theory.

490 From the SA results, we see that NH_4^+ concentrations at the sediment water interface in the lake are
491 influenced by parameters like EC, WT, AT, DO, turbidity, and Press (as seen in Figure 5). The water content
492 present in parent untreated fluid tailings contains high concentrations of dissolved constituents, including
493 Na, Cl, organics, and NH_3 (Dompierre et al., 2016). These dissolved constituents are released from tailings
494 due to the upward movement of water associated with tailings densification (Dompierre and Barbour
495 2016), providing a mechanistic explanation for the associations between EC and NH_3 (detectable NH_4^+).
496 AT affects WT, which in turn influences the rate of microbial respiration, with elevated WT promoting
497 biological oxygen demand and production of NH_3 (Stumm and Morgan, 2012). Warmer temperatures and
498 declining DO can also increase sediment bioturbation rate by chironomid (invertebrate) larvae (Roskosch
499 et al., 2012), further promoting bio-irrigation-mediated benthic fluxes of salt, NH_3 , and other dissolved
500 constituents.

501 Further, there is an interplay between NH_4^+ and DO, as NH_4^+ could be an oxygen consuming constituent in
502 oil sands end pit lakes (Risacher et al. 2018). The DO parameter was well captured by the tree-based
503 model, while the connectionist ML models were unable to identify DO as an important parameter for
504 NH_4^+ . Moreover, the interplay between overland water flow and bioturbation enhances metal flux from
505 low permeability sediment beds (Amato et al., 2016; Xie et al., 2018). In particular, it is possible that this
506 bioturbation process can partly destabilise the sediment bed in *Lake Miwasin* and cause a temporary
507 remobilisation of suspended particles and particulate organic matter that can yield to overall fluctuation
508 in turbidity levels. Fluid tailings is a source of particulate organic matter (Sasar et al., 2022) and may
509 increase turbidity values during periods of water column stratification. Overall, the influential input
510 parameters (EC, WT, AT, DO, Turb) are very well captured by tree-based models in our SA for prediction
511 of NH_4^+ .

512 The SA indicates that Chl-a prediction using ML models are mostly influenced by WT, EC, Rad (light), pH,
513 NH_4^+ . Evidence showed that Chl-a could be an indicator of (a) photosynthesis (affected by DO, solar
514 radiation, and temperature; Shammas et al., 2009; Wallace et al., 2016); (b) nutrient status (affected by
515 pH since algae grow better at higher pH values by taking up more nutrients and CO_2 under alkaline

516 conditions; *Veeresh et al., 2010; Wallace et al., 2016*), and (c) the growth and distribution of
517 phytoplankton species composition (affected by solar light, DO and WT; *Harrison et al., 2018; Bouffard et*
518 *al., 2018, Liu and Georgakakos, 2021*) in the lake. Without any ambiguity, our SA method captured most
519 of the important sensitive parameters for Chl-a prediction using the tree-based models, but it failed to
520 identify DO as an important input.

521 For prediction of DO, the suggested variables EC, WT, NH_4^+ , and pH are frequently signaled by different
522 models in the SA. Similarly, the prediction of pH can be influenced by WL, WT, NH_4^+ , and DO (Figure 6).
523 Most of these parameters are indirectly or directly involved in different processes occurring within the
524 lake, such as release of pore water from TFT and fluctuation of DO due to the release of oxygen consuming
525 constituents (e.g., NH_4^+) (*Dompierre and Barbour, 2016; Risacher et al. 2018*). Based on our SA, the tree-
526 based models identified most of these as sensitive parameters for prediction of DO and pH whereas the
527 connectionist ML models did not (Figure 6).

528 6. Conclusions

529 This study developed a new approach to XAI through the lens of SA. This approach has the conceptual
530 strength that it characterizes the entire response surface of an ML model – whereas other methods
531 typically look only at the model response in the region of the available input-output data points. This
532 approach was used to investigate the primary controls on the physico-chemical processes of a major
533 environmental problem, as determined by a suite of connectionist, kernel-based, and tree-based ML
534 models. The analyses enabled efficient detection of important explanatory variables, thereby guiding
535 long-term monitoring programs with reduced data collection cost.

536 Notably, although most of the ML models showed similar levels of predictive power, they tended to base
537 their predictions on different explanatory variables (inputs). In particular, the connectionist ML models
538 such as neural networks showed a large degree of variability in how their outputs depended on the various
539 inputs. Different replicates of the same connectionist model were often primarily driven by different
540 inputs, suggesting that the model may pick up different signals in the data to provide similar levels of
541 predictability. Interestingly, the important inputs of the tree-based ML models were more consistent with
542 each other, while tending to be somewhat different from those of the connectionist and kernel-based
543 models.

544 Overall, our analysis reveals an important issue that is arguably a critical takeaway message of this paper.
545 Subtle alterations in the design of ML models (such as variations in the random seed used for initialization,
546 functional classes, hyperparameters, or data splitting) can lead to entirely different representational
547 interpretations of the dependence of the outputs on explanatory variables (inputs). This strongly
548 reinforces the importance of utilizing ensembles of ML models when explanatory power is a desirable
549 outcome (see also *De La Fuente et al., 2023*). Such ensembles could be generated via multiple replicates
550 of the same model, or by employing diverse types of ML models, or some combination of both. Failure to
551 do so could mean that the analysis cannot be relied upon to guarantee the delivery of robust and reliable
552 predictions, particularly when using the developed models to generalize to conditions beyond the
553 region(s) of the data used for model training.

554 **Acknowledgement:** We thank Suncor for their invaluable on-site assistance at Lake Miwasin and providing
555 access to their database platform.

556 **Funding:** We acknowledge support by the Suncor Energy Inc.

557 **Competing interests:** The authors declare no conflict of interest.

558 **Data and software availability:** The software developed here for G-VARS2 has been made available on
559 GitHub at <https://github.com/vars-tool/vars-tool>. The data used will be made available on a public
560 repository.

561 **REFERENCES:**

- 562 Aas, K., Jullum, M., Løland, A. (2021). Explaining individual predictions when features are dependent: More
563 accurate approximations to Shapley values. *Artificial Intelligence*, 298, 103502.
564 <https://doi.org/10.1016/j.artint.2021.103502>.
- 565 AEP (2015). Total Area of the Oil Sands Tailings Ponds over Time, Oil Sands Information Portal, Alberta
566 Environment and Parks, Edmonton, AB. Dataset no. 542. 4 March 2015. Available:
567 <http://osip.alberta.ca/library/Dataset/Details/542>.
- 568 Abobakr Yahya, A. S., Najah, A., Othman, F. B., Ibrahim, R. K., Afan, H. A., El-Shafie, A., Fai, C. M., Hossain, M. S.,
569 Ehteram, M., & Elshafie, A. (2019). Water quality prediction model-based support vector machine model
570 for ungauged river catchment under dual scenarios. *Water*, 11(6), 1231.
571 <https://doi.org/10.3390/w11061231>.
- 572 Ahmed, A. A. M. (2014). Prediction of dissolved oxygen in Surma River by biochemical oxygen demand and
573 chemical oxygen demand using the artificial neural networks (ANNs). *Journal of King Saud University-
574 Engineering Sciences*, 29(2), 151-158. <https://doi.org/10.1016/j.jksues.2014.05.001>.
- 575 Ahmed, A.N., Othman, F. B., Afan, H. A., Ibrahim, R. K., Fai, C. M., Hossain, M. S., Hossain, M. S., Ehteram, M., &
576 Elshafie, A. (2019a). Machine learning methods for better water quality prediction. *Journal of Hydrology*,
577 578, 124084. <https://doi.org/10.1016/j.jhydrol.2019.124084>.
- 578 Ahmed, U., Mumtaz, R., Anwar, H., Shah, A. A., Irfan, R., & García-Nieto, J. (2019b). Efficient water quality
579 prediction using supervised machine learning. *Water*, 11(11), 2210. <https://doi.org/10.3390/w11112210>.
- 580 Alipour, A., Jafarzadegan, K., & Moradkhani, H. (2022). Global sensitivity analysis in hydrodynamic modeling and
581 flood inundation mapping. *Environmental Modelling & Software*, 152, 105398.
582 <https://doi.org/10.1016/j.envsoft.2022.105398>.
- 583 Amato, E. D., Simpson, S. L., Remaili, T. M., Spadaro, D. A., Jarolimek, C. V., & Jolley, D. F. (2016). Assessing the
584 effects of bioturbation on metal bioavailability in contaminated sediments by diffusive gradients in thin
585 films (DGT), *Environ. Sci. Technol.* 50(2016) 3055–3064. <https://doi.org/10.1021/acs.est.5b04995>.
- 586 Banerjee, A., Chakrabarty, M., Rakshit, N., Bhowmick, A.R., & Ray, S. (2019). Environmental factors as indicators
587 of dissolved oxygen concentration and zooplankton abundance: Deep learning versus traditional
588 regression approach. *Ecological indicators*, 100, 99-117. <https://doi.org/10.1016/j.ecolind.2018.09.051>.
- 589 Barzegar, R., Aalami, M. T., & Adamowski, J. (2020). Short-term water quality variable prediction using a hybrid
590 CNN–LSTM deep learning model. *Stochastic Environmental Research and Risk Assessment*, 34(2), 415-433.
591 <https://doi.org/10.1007/s00477-020-01776-2>.
- 592 Becker, W. (2020). Meta functions for benchmarking in sensitivity analysis. *Reliability Engineering & System
593 Safety*, 204, 107189. <https://doi.org/10.1016/j.res.2020.107189>.
- 594 Bouffard, D., Kiefer, I., Wüest, A., Wunderle, S., & Odermatt, D., (2018). Are surface temperature and
595 chlorophyll in a large deep lake related? An analysis based on satellite observations in synergy with
596 hydrodynamic modelling and in-situ data. *Remote Sens. Environ.* 209, 510–523.
597 <https://doi.org/10.1016/j.rse.2018.02.056> <https://doi.org/>.
- 598 Breiman, L. (2001). Random forests. *Machine learning*, 45, 5-32.
- 599 Brokamp, C., Jandarov, R., Hossain, M., & Ryan, P. (2018). Predicting daily urban fine particulate matter
600 concentrations using a Random Forest model. *Environ. Sci. Technol.* 52, 4173–4179.
601 <https://doi.org/10.1021/acs.est.7b05381>.
- 602 Campolongo, F., Cariboni, J., & Saltelli, A. (2007). An effective screening design for sensitivity analysis of large
603 models. *Environ. Modelling & Software, Modelling, Comput. Assisted Simul. Map. Dangerous
604 Phenomena for Hazard Assess.* 22, 1509–1518. <https://doi.org/10.1016/j.envsoft.2006.10.004>.

605 Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system, in: Proceedings of the 22nd Acm
606 Sigkdd International Conference on Knowledge Discovery and Data Mining, ACM, 2016, pp. 785–794.
607 <https://doi.org/10.1145/2939672.2939785>.

608 Chen, K., Chen, H., Zhou, C., Huang, Y., Qi, X., Shen, R., & Ren, H. (2020). Comparative analysis of surface water
609 quality prediction performance and identification of key water parameters using different machine
610 learning models based on big data. *Water Res.* 171, 115454.
611 <https://doi.org/10.1016/j.watres.2019.115454>.

612 COSIA (2012). Technical Guide for Fluid Fine Tailings Management. Canadian Oil Sands Innovation Alliance,
613 Edmonton, AB. August 2012. Report, 1-131. [Available at [http://www.cosia.ca /initiatives /project-](http://www.cosia.ca/initiatives/project-research)
614 [research](http://www.cosia.ca/initiatives/project-research), last accessed 27 September 2023.]

615 Cossey, H. L., Batycky, A. E., Kaminsky, H., & Ulrich, A. C. (2021). Geochemical stability of oil sands tailings in
616 mine closure landforms. *Minerals*, 11(8), 830. <https://doi.org/10.3390/min11080830>.

617 Coyle, D., & Weller, A. (2020). “Explaining” machine learning reveals policy challenges. *science*, 368(6498),
618 1433-1434. <https://doi.org/10.1126/science.aba9647>.

619 Cupe-Flores, B., Mendes, M., Panigrahi, B., & Liber, K. (2022). Delineating effluent exposure and cumulative
620 ecotoxicological risk of metals downstream of a Saskatchewan uranium mill using autonomous sensors.
621 *Environmental Toxicology and Chemistry* 41, 1765-1777. <http://dx.doi.org/10.1002/etc.5341>.

622 De la Fuente, L. A., Gupta, H. V., & Condon, L. E. (2023). Toward a Multi-Representational Approach to Prediction
623 and Understanding, in Support of Discovery in Hydrology. *Water Resources Research*, 59(1),
624 e2021WR031548.

625 Do, N. C., & Razavi, S., (2020). Correlation effects? A major but often neglected component in sensitivity and
626 uncertainty analysis. *Water Resources Research*, 56, e2019WR025436.
627 <https://doi.org/10.1029/2019WR025436>.

628 Dompierre, K. A., Lindsay, M. B. J., Cruz-Hernández, P., & Halferdahl, G. M. (2016). Initial geochemical
629 characteristics of fluid fine tailings in an oil sands end pit lake, *Sci. Total Environ.*, *Science of the Total*
630 *Environment* 556, 196-206. [https://doi.org/10.1016/ j.scitotenv.2016.03.002](https://doi.org/10.1016/j.scitotenv.2016.03.002).

631 Dompierre, K.A., & Barbour, S.L. (2016). Characterization of physical mass transport through oil sands fluid fine
632 tailings in an end pit lake: A multi-tracer study, *J. Contam. Hydrol.*, 189(1), 12–26.
633 <https://doi.org/10.1016/j.jconhyd.2016.03.006>.

634 Fatichi, S., Vivoni, E. R., Ogden, F. L., Ivanov, V. Y., Mirus, B., Gochis, D., Downer, C. W., Camporese, M., Davison,
635 J. H., Ebel, B., Jones, N., Kim, J., Mascaro, G., Niswonger, R., Restrepo, P., Rigon, R., Shen, C., Sulis, M., &
636 Tarboton, D. (2016). An overview of current applications, challenges, and future trends in distributed
637 process-based models in hydrology. *Journal of Hydrology*, 537, 45–60. [https://doi.org/10.1016/J.](https://doi.org/10.1016/J.JHYDROL.2016.03.026)
638 [JHYDROL.2016.03.026](https://doi.org/10.1016/J.JHYDROL.2016.03.026).

639 Fel, T., Cadene, R., Chalvidal, M., Cord, M., Vigouroux, D., & Serre, T. (2021) Look at the variance! Efficient
640 black-box explanations with Sobol-based sensitivity analysis. In Beygelzimer, A., Dauphin, Y., Liang, P.,
641 and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, 2021. URL
642 <https://openreview.net/forum?id=hA-PHQGOjqQ>.

643 Freund, Y., & Schapire, R. E. (1997). A desicion-theoretic generalization of on-line learning and an application
644 to boosting *J. Comput. Syst. Sci.*, 55 (1997), pp. 23-37.

645 Friedman, J. H. (1991). Multivariate adaptive regression splines. *The annals of statistics*, 19(1), 1-67.
646 <https://doi.org/10.1214/aos/1176347963>.

647 Friedman J. H. (2001). Greedy function approximation: a gradient boosting machine, *Ann. Statist.* (2001) 1189–
648 1232. <https://www.jstor.org/stable/2699986>.

649 Frye, C., Rowat, C., & Feige, I. (2020). Asymmetric shapley values: incorporating causal knowledge into model-
650 agnostic explainability. *Advances in Neural Information Processing Systems*, 33, 1229-1239.

651 Garosi, Y., Sheklabadi, M., Conoscenti, C., Pourghasemi, H. R., & Van Oost, K. (2019). Assessing the performance
652 of GIS-based machine learning models with different accuracy measures for determining susceptibility
653 to gully erosion. *Science of the Total Environment*, 664, 1117-1132.
654 <https://doi.org/10.1016/j.scitotenv.2019.02.093>.

655 Gosselin, P., Hrudehy, S. E., Naeth, M. A., Plourde, A., Therrien, R., Kraak, G. V. D., & Xu, Z. (2010). Environmental
656 and Health Impacts of Canada's Oils Sands Industry, R. Soc. of Canada, Ottawa, Canada.

657 Haghnegahdar, A., & Razavi, S. (2017). Insights into sensitivity analysis of Earth and environmental systems
658 models: On the impact of parameter perturbation scale. *Environmental Modelling & Software*, 95, 115-
659 131. <http://dx.doi.org/10.1016/j.envsoft.2017.03.031>.

660 Harrison, J. W., Beecraft, L., & Smith, R. E. H. (2018). Implications of irradiance exposure and non-
661 photochemical quenching for multi-wavelength (bbe FluoroProbe) fluorometry. *J. Photochem.*
662 *Photobiol. B Biol.* 189, 36–48. <https://doi.org/10.1016/j.jphotobiol.2018.09.013>.

663 Hipsey, M. R., Hamilton, D. P., Hanson, P. C., Carey, C. C., Coletti, J. Z., Read, J. S., Ibelings, B. W., Valesini, F. J.,
664 & Brookes, J. D. (2015). Predicting the resilience and recovery of aquatic systems: A framework for model
665 evolution within environmental observatories. *Water Resources Research*, 51, 7023–7043.
666 <https://doi.org/10.1002/2015WR017175>.

667 Homma, T., & Saltelli, A. (1996). Importance measures in global sensitivity analysis of nonlinear models.
668 *Reliability Engineering & System Safety*, 52(1):1–17. [https://doi.org/10.1016/0951-8320\(96\)00002-6](https://doi.org/10.1016/0951-8320(96)00002-6).

669 Janzing, D., Minorics, L., & Blöbaum, P. (2020). Feature relevance quantification in explainable AI: A causal
670 problem. In *International Conference on artificial intelligence and statistics*, PMLR, 2020, pp. 2907-2916.

671 Jeihouni, M., Toomanian, A., & Mansourian, A. (2020). Decision tree-based data mining and rule induction for
672 identifying high quality groundwater zones to water supply management: a novel hybrid use of data
673 mining and GIS. *Water Resour. Manag.* 34, 139–154. <https://doi.org/10.1007/s11269-019-02447-w>.

674 Kambhatla, N., & Leen, T. K. (1997). Dimension reduction by local principal component analysis. *Neural*
675 *Comput.* 9, 1493–1516. <https://doi.org/10.1162/neco.1997.9.7.1493>.

676 Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., & Liu, T. Y. (2017). LightGBM: A highly efficient
677 gradient boosting decision tree, in: *Advances in Neural Information Processing Systems*, pp. 3146–3154.

678 Khullar, S., & Singh, N. (2021). Water quality assessment of a river using deep learning Bi-LSTM methodology:
679 forecasting and validation. *Environmental Science and Pollution Research*, 29(9), 12875-12889.
680 <https://doi.org/10.1007/s11356-021-13875-w>.

681 Kim, H., & Kim, B. (2020). Analysis of major rainfall factors affecting inundation based on observed rainfall and
682 Random Forest. *Korean Soc. Hazard Mitig.* 20, 301–310.
683 <https://doi.org/10.9798/KOSHAM.2020.20.6.301>.

684 Kleijnen, J. P. C. (1995). Sensitivity analysis and optimization of system dynamics models: regression analysis
685 and statistical design of experiments. *Syst. Dynam. Rev.* 11, 275–288.
686 <https://doi.org/10.1002/sdr.4260110403>.

687 Konapala, G., & Mishra, A. (2020). Quantifying climate and catchment control on hydrological drought in the
688 continental United States. *Water Resources Research*, 56 (1). <https://doi.org/10.1029/2018WR024620>.

689 Kucherenko, S., & Iooss, B. (2016). Derivative-based global sensitivity measures. In *Handbook of Uncertainty*
690 *Quantification*, pp. 1–24. Springer International Publishing, Cham. https://doi.org/10.1007/978-3-319-11259-6_36-1.

691

692 Kucherenko, S., & Song, S. (2016). Derivative-based global sensitivity measures and their link with Sobol'
693 sensitivity indices. *Monte Carlo and Quasi-Monte Carlo Methods*, pp. 455–469.
694 https://doi.org/10.1007/978-3-319-33507-0_23.

695 Kuhnt, S. & Kalka, A. (2022). Global sensitivity analysis for the interpretation of machine learning algorithms,
696 pp. 155–169. Springer International Publishing, Cham, 2022. https://doi.org/10.1007/978-3-031-07155-3_6.

698 Kumar, I. E., Venkatasubramanian, S., Scheidegger, C., & Friedler, S. (2020). Problems with Shapley-value-based
699 explanations as feature importance measures. In *International Conference on Machine Learning* (pp.
700 5491-5500). *Proceedings of Machine Learning Research* (PMLR).
701 <https://proceedings.mlr.press/v119/kumar20e/kumar20e.pdf>.

702 Lakkaraju, H., Kamar, E. Caruana, R., & Leskovec, J. (2019). Faithful and customizable explanations of black box
703 models. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 131–138.

704 Lakkaraju, H. & Bastani, O. (2020). How do I fool you? Manipulating user trust via misleading black box
705 explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society* (pp. 79-85).

706 Lamboni, M., Iooss, B., Popelin, A.L., & Gamboa, F. (2013). Derivative-based global sensitivity measures: general
707 links with Sobol' indices and numerical tests. *Math. Comput. Simulat.* 87, 45–54.
708 <https://doi.org/10.1016/j.matcom.2013.02.002>.

709 Liakos, K. G., Busato, P., Moshou, D., Pearson, S., & Bochtis, D. (2018). Machine learning in agriculture: A review.
710 *Sensors*, 18(8), 2674. <https://doi.org/10.3390/s18082674>.

711 Libelium (2020). *Technical Guide-Smart Water Extreme*, Document version: v7.3- 02/2020.
712 <https://www.scribd.com/document/375986458/Waspmote-Technical-Guide>.

713 Lipton, Z. C. (2017). The Mythos of Model Interpretability: In Machine Learning, the Concept of Interpretability
714 Is Both Important and Slippery. *Queue*, 16, 31–57. <https://arxiv.org/pdf/1606.03490.pdf>.

715 Liu, H., Clark, M. P., Gharari, S., Sheikholeslami, R., Freer, J., Knoben, W. J. M., Marsh, C. B., & Papalexioiu, S. M.
716 (2024). An improved copula-based framework for efficient global sensitivity analysis. *Water Resources*
717 *Research*, 60, e2022WR033808. <https://doi.org/10.1029/2022WR033808>.

718 Liu, X., & Georgakakos, A. P. (2021). Chlorophyll a estimation in lakes using multi-parameter sonde data, *Water*
719 *Research*, Volume 205, 117661, ISSN 0043-1354, <https://doi.org/10.1016/j.watres.2021.117661>.

720 Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in neural*
721 *information processing systems*, 30. <https://doi.org/10.48550/arXiv.1705.07874>.

722 Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., &
723 Lee, S. I. (2020). From local explanations to global understanding with explainable AI for trees, *Nature*
724 *Machine Intelligence*, Vol 56-67. <https://doi.org/10.1038/s42256-019-0138-9>.

725 Maier, H. R., Galelli, S., Razavi, S., Castelletti, A., Rizzoli, A., Athanasiadis, I. N., Sanchez-Marre, M., Acutis, M., Wu,
726 W., & Humphrey, G. B. (2023). Exploding the myths: An introduction to artificial neural networks for
727 prediction and forecasting. *Environmental modelling & software*, 105776.
728 <https://doi.org/10.1016/j.envsoft.2023.105776>.

729 Mase, M., Owen, A. B., & Seiler, B. (2019). Explaining black box decisions by shapley cohort refinement. *arXiv*
730 preprint arXiv:1911.00467.

731 McNeill, J. (2017). Alberta Government continues its mismanagement of fluid tailings with approval of CNRL's
732 tailings plan: Pembina Institute reacts to Alberta Energy Regulator's decision.
733 <https://www.pembina.org/media-release>.

734 Molnar, C., König, G., Herbringer, J., Freiesleben, T., Dandl, S., Scholbeck, C. A., Casalicchio, G., Grosse-Wentrup,
735 M. & Bischl, B. (2020). Pitfalls to avoid when interpreting machine learning models.

736 Molnar, C. (2022) Interpretable Machine Learning. 2nd edition, 2022. URL [https://christophm.github.io/](https://christophm.github.io/interpretable-ml-book)
737 interpretable-ml-book.

738 Morris, M. D. (1991). Factorial sampling plans for preliminary computational experiments. *Technometrics* 33,
739 161–174. <https://doi.org/10.1080/00401706.1991.10484804>.

740 Mosavi, A., Hosseini, S. F., Choubin, B., Goodarzi, M., Dineva, A. A., & Sardooi, R. E. (2021). Ensemble boosting
741 and bagging based machine learning models for groundwater potential prediction. *Water Resources*
742 *Management*, 35, 23-37. <https://doi.org/10.1007/s11269-020-02704-3>.

743 Najah, A., El-Shafie, A., Karim, O.A., & El-Shafie, A.H. (2014). Performance of ANFIS versus MLP-NN dissolved
744 oxygen prediction models in water quality monitoring. *Environ. Sci. Pollut. Res.* 2014, 21, 1658–1670.
745 <https://doi.org/10.1007/s11356-013-2048-4>.

746 Nossent, J., Elsen, P., & Sobol', B.W. (2011). Sensitivity analysis of a complex environmental model.
747 *Environmental Modelling & Software*, 26(12):1515–1525, 2011.
748 <https://doi.org/10.1016/j.envsoft.2011.08.010>.

749 Ojha, V., Timmis, J., & Nicosia, G. (2022). Assessing ranking and effectiveness of evolutionary algorithm
750 hyperparameters using global sensitivity analysis methodologies. *Swarm and Evolutionary Computation*,
751 74,101130. <https://doi.org/10.1016/j.swevo.2022.101130>.

752 Owen, A. (1994). Lattice Sampling Revisited: Monte Carlo Variance of Means Over Randomized Orthogonal
753 Arrays. *Ann. Statist.* 22, 930–945. <https://doi.org/10.1214/aos/1176325504>.

754 Owen, A. (2014). Sobol' indices and Shapley value. *SIAM/ASA Journal on Uncertainty Quantification*, 2:245–
755 251, 01 2014. <https://doi.org/10.1137/130936233>.

756 Paleari, L., Movedi, E., Zoli, M., Burato, A., Cecconi, I., Errahouly, J., Pecollo, E., Sorvillo, C., & Confalonieri, R.
757 (2021). Sensitivity analysis using Morris: Just screening or an effective ranking method? *Ecological*
758 *Modelling*, 455: 109648, 2021. <https://doi.org/10.1016/j.ecolmodel.2021.109648>.

759 Palani, S., Liong, S.Y., & Tkalich, P. (2008). An ANN application for water quality forecasting. *Mar. Pollut. Bull.*
760 56, 1586–1597. <https://doi.org/10.1016/j.marpolbul.2008.05.021>.

761 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss,
762 R., Dubourg, V., Vanderplas, J., Passos, A., & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python.
763 *the Journal of machine Learning research*, 12, 2825-2830.

764 Peixoto Mendes, M., Cupe-Flores, B., Panigrahi, B., & Liber, K. (2023). Application of autonomous sensor
765 technology to estimate selenium exposure and a site-specific selenium threshold in a Canadian boreal
766 lake. *Integrated Environmental Assessment and Management* 19, 395-411.
767 <http://dx.doi.org/10.1002/ieam.4644>.

768 Pena, M., Katsev, S., Oguz, T., & Gilbert, D. (2010). Modeling dissolved oxygen dynamics and hypoxia.
769 *Biogeosciences*, 7 (3), 933–957. <https://doi.org/10.5194/bg-7-933-2010>.

770 Puy, A., Becker, W., Lo Piano, S., & Saltelli, A. A. (2021). comprehensive comparison of total-order estimators
771 for global sensitivity analysis. *International Journal for Uncertainty Quantification*.
772 <https://doi.org/10.1615/int.j.uncertaintyquantification.2021038133>.

773 Puy, A., Beneventano, P., Levin, S. A., Lo Piano, S., Portaluri, T., & Saltelli, A. (2022). Models with higher effective
774 dimensions tend to produce more uncertain estimates. *Science Advances*, 8(42),
775 <https://doi.org/10.1126/sciadv.abn945>.

776 Rahmati, O., Tahmasebipour, N., Haghizadeh, A., Pourghasemi, H. R., & Feizizadeh, B. (2017). Evaluation of
777 different machine learning models for predicting and mapping the susceptibility of gully erosion.
778 *Geomorphology*, 298, 118-137. <https://doi.org/10.1016/j.geomorph.2017.09.006>.

779 Rakovec, O., Hill, M. C., Clark, M. P., Weerts, A. H., Teuling, A. J., & Uijlenhoet, R. (2014). Distributed Evaluation
780 of Local Sensitivity Analysis (DELSA), with application to hydrologic models: distributed evaluation of
781 local sensitivity analysis. *Water Resour. Res.* 50, 409–426. <https://doi.org/10.1002/2013WR014063>.

782 Razavi, S., & Gupta, H.V. (2015). What do we mean by sensitivity analysis? The need for comprehensive
783 characterization of “global” sensitivity in Earth and Environmental systems models: a Critical Look at
784 Sensitivity Analysis. *Water Resour. Res.* 51, 3070–3092. <https://doi.org/10.1002/2014WR016527>.

785 Razavi, S., & Gupta, H. V. (2016a). A new framework for comprehensive, robust, and efficient global sensitivity
786 analysis: 1. Theory. *Water Resour. Res.* 52, 423–439. <https://doi.org/10.1002/2015WR017558>.

787 Razavi, S., & Gupta, H. V. (2016b). A new framework for comprehensive, robust, and efficient global sensitivity
788 analysis: 2. Application. *Water Resour. Res.* 52, 440–455. <https://doi.org/10.1002/2015WR017559>.

789 Razavi, S. (2021). Deep learning, explained: Fundamentals, explainability, and bridgeability to process-based
790 modelling. *Environmental Modelling & Software*, 144:105159, 2021. doi:
791 <https://doi.org/10.1016/j.envsoft.2021.105159>.

792 Razavi, S., Jakeman, A., Saltelli, A., Prieur, C., Iooss, B., Borgonovo, E., Plischke, E., Piano, L. S., Iwanaga, T., Becker,
793 W., Tarantola, S., Guillaume, J. H., Jakeman, J., Gupta, H., Melillo, N., Rabitti, G., Chabridon, V., Duan, Q.,
794 Sun, X., Smith, S., Sheikholeslami, R., Hosseini, N., Asadzadeh, M., Puy, A., Kucherenko, S., & Maier, H. R.
795 (2021). The Future of Sensitivity Analysis: An essential discipline for systems modeling and policy support,
796 *Environmental Modelling and Software* 137 (2021) 104954,
797 <https://doi.org/10.1016/j.envsoft.2020.104954>.

798 Razavi, S., Hannah, D. M., Elshorbagy, A., Kumar, S., Marshall, L., Solomatine, D. P., Dezfouli, A., Sadegh, M., &
799 Famiglietti, J. (2022). Coevolution of machine learning and process-based modelling to revolutionize
800 Earth and environmental sciences: A perspective. *Hydrological Processes*, 36(6).
801 <https://doi.org/10.1002/hyp.14596>.

802 Read, J. S., Jia, X., Willard, J., Appling, A. P., Zwart, J. A., Oliver, S. K., Karpatne, A., Hansen, G. J. A., Hanson, P.
803 C., Watkins, W., Steinbach, M., & Kumar, V. (2019). Process-guided deep learning predictions of lake
804 water temperature. *Water Resources Research*, 55, <https://doi.org/10.1029/2019WR024922>.

805 Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). Why should i trust you? Explaining the predictions of any
806 classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and*
807 *data mining* (pp. 1135-1144).

808 Risacher, F. F., Morris, P. K., Arriaga, D., Goad, C., Nelson, T. C., Slater, G. F., & Warren, L. A. (2018). The interplay
809 of methane and ammonia as key oxygen consuming constituents in early stage development of Base
810 Mine Lake, the first demonstration oil sands pit lake, *Applied Geochemistry*, Volume 93, 2018, Pages 49-
811 59, ISSN 0883-2927, <https://doi.org/10.1016/j.apgeochem.2018.03.013>.

812 Roscher, R., Bohn, B., Duarte, A., & Garcke, J. (2020). Explainable Machine Learning for Scientific Insights and
813 Discoveries 10.1109/ACCESS.2020.2976199.

814 Roskosch, A., Hette, N., Hupfer, M., & Lewandowski, J. (2012). Alteration of Chironomus plumosus ventilation
815 activity and bioirrigation-mediated benthic fluxes by changes in temper ature, oxygen concentration,
816 and seasonal variations. *Fresh water Science* 31:269–281.

817 Rudin, S. (2019). Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use
818 Interpretable Models Instead. *Nat. Mach. Intell*, 1, 206–215.

819 Rudin, C., Chen, C., Chen, Z., Huang, H., Semenova, L., & Zhong, C. (2022). Interpretable machine learning:
820 Fundamental principles and 10 grand challenges. *Statistics Surveys*, 16:1 – 85, 2022.
821 <https://doi.org/10.1214/21-SS133>.

822 Saltelli, A., Jakeman, A., Razavi, S., & Wu, Q. (2021). Sensitivity analysis: A discipline coming of age.
823 *Environmental Modelling & Software*, 146, 105226. <https://doi.org/10.1016/j.envsoft.2021.105226>.

824 Sameen, M. I., Pradhan, B., & Lee, S. (2019). Self-learning random forests model for mapping groundwater yield
825 in data-scarce areas. *Natural Resources Research*, 28, 757-775. [https://doi.org/10.1007/s11053-018-](https://doi.org/10.1007/s11053-018-9416-1)
826 [9416-1](https://doi.org/10.1007/s11053-018-9416-1).

827 Samek, W., Montavon, G., Vedaldi, A., Hansen, L. K. & Müller, K. R. (2019). *Explainable AI: Interpreting,*
828 *Explaining and Visualizing Deep Learning; Lecture Notes in Artificial Intelligence, Lect.Notes*
829 *ComputerState-of-the-Art Surveys; Springer: Berlin/Heidelberg, Germany, ISBN 978-3-030-28953-9.*

830 Sánchez, E., Colmenarejo, M. F., Vicente, J., Rubio, A., García, M. G., Travieso, L., & Borja, R. (2006). Use of the
831 water quality index and dissolved oxygen deficit as simple indicators of watersheds pollution. *Ecol. Indic.*
832 2007, 7 (2), 315–328. <https://doi.org/10.1016/j.ecolind.2006.02.005>.

833 Sasar M., Johnston C. T., & Santagata M. (2022). Characterization and Dynamics of Residual Organics in Oil
834 Sands Fluid Fine Tailings. *Energy Fuels* 2022, 36, 6881–689.
835 <https://doi.org/10.1021/acs.energyfuels.2c01315>.

836 Savoy, P., & Harvey, J. W. (2023). Predicting daily river chlorophyll concentrations at a continental scale. *Water*
837 *Resources Research*, 59, e2022WR034215. <https://doi.org/10.1029/2022WR034215>.

838 Scholbeck, C. A., Moosbauer, J., Casalicchio, G., Gupta, H., Bischl, B., & Heumann, C. (2023). Position Paper:
839 Bridging the Gap Between Machine Learning and Sensitivity Analysis. arXiv preprint arXiv:2312.13234.

840 Shammas, N., Wang, L., & Wu, Z. (2009). Waste Stabilization Ponds and WSPs, in Volume 8: Biological
841 Treatment Processes, *Handbook of Environmental Engineering*, Humana Press, Totowa, New Jersey, pp.
842 315–370.

843 Sheikholeslami, R., Gharari, S., Papalexioiu, S. M., & Clark, M. P. (2021). VISCOUS: A variance-based sensitivity
844 analysis using copulas for efficient identification of dominant hydrological processes. *Water Resources*
845 *Research*, 57, e2020WR028435. <https://doi.org/10.1029/2020WR028435>.

846 Sheikholeslami, R., & Razavi, S. (2020). A fresh look at variography: measuring dependence and possible
847 sensitivities across geophysical systems from any given data. *Geophysical Research Letters*, 47(20).
848 <https://doi.org/10.1029/2020GL089829>.

849 Sheikholeslami, R., Razavi, S., Gupta, H.V., Becker, W., & Haghnegahdar, A. (2019). Global sensitivity analysis
850 for high-dimensional problems: how to objectively group factors and measure robustness and
851 convergence while reducing computational cost. *Environ. Model. Software* 111, 282–299.
852 <https://doi.org/10.1016/j.envsoft.2018.09.002>.

853 Sheikholeslami, R., Yassin, F., Lindenschmidt, K. E., & Razavi, S. (2017). Improved understanding of river ice
854 processes using global sensitivity analysis approaches. *Journal of Hydrologic Engineering*, 22(11).
855 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0001574](https://doi.org/10.1061/(ASCE)HE.1943-5584.0001574).

856 Shin, M.-J., Guillaume, J. H., Croke, B. F., & Jakeman, A. J. (2013) Addressing ten questions about conceptual
857 rainfall– runoff models with global sensitivity analyses in R. *Journal of Hydrology*, 503:135–152, 2013.
858 <https://doi.org/10.1016/j.jhydrol.2013.08.047>.

859 Sinshaw, T.A., Surbeck, C.Q., Yasarer, H., & Najjar Y. (2019). Artificialneural network for prediction of total
860 nitrogen and phosphorus in US lakes. *J Environ Eng* 145(6):04019032.
861 [https://doi.org/10.1061/\(ASCE\)EE.1943-7870.0001528](https://doi.org/10.1061/(ASCE)EE.1943-7870.0001528).

862 Slack, D., Hilgard, A., Singh, S., & Lakkaraju, H. (2021). Reliable post hoc explanations: Modeling uncertainty in
863 explainability. *Advances in neural information processing systems*, 34, 9391-9404.

864 Slack, D., Krishna, S., Lakkaraju, H., & Singh, S. (2023). Explaining machine learning models with interactive
865 natural language conversations using TalkToModel. *Nature Machine Intelligence*, 5(8), 873-883.
866 <https://doi.org/10.1038/s42256-023-00692-8>.

867 Sobol', I. M. (1993). Sensitivity analysis for non-linear mathematical models. *Mathematical Modelling and*
868 *Computational Experiment, Translated from Russian: I.*

869 M. Sobol', Sensitivity estimates for nonlinear mathematical models. *Matematicheskoe Modelirovanie* 2,
870 407–414 (1990) 112–118 1.

871 Sobol, I. M. (2001). Global sensitivity indices for nonlinear mathematical models and their Monte Carlo
872 estimates. *MATH COMPUT SIMULAT*,55(1–3),271-280, [https://doi.org/10.1016/S0378-4754\(00\)00270-](https://doi.org/10.1016/S0378-4754(00)00270-6)
873 [6](https://doi.org/10.1016/S0378-4754(00)00270-6).

874 Sobol', I. M., & Kucherenko, S. (2009). Derivative based global sensitivity measures and their link with global
875 sensitivity indices. *Math. Comput. Simulat.* 79, 3009–3017.
876 <https://doi.org/10.1016/j.matcom.2009.01.023>.

877 Song, X., Zhang, J., Zhan, C., Xuan, Y., Ye, M., & Xu, C. (2015). Global sensitivity analysis in hydrological modeling:
878 Review of concepts, methods, theoretical framework, and applications. *Journal of Hydrology*, 523:739–
879 757, 2015. <https://doi.org/10.1016/j.jhydrol.2015.02.013>.

880 Song, E., Nelson, B. L., & Staum, J. (2016). Shapley effects for global sensitivity analysis: Theory and
881 computation. *SIAM/ASA Journal on Uncertainty Quantification*, 4(1), 1060-1083.
882 <https://doi.org/10.1137/15M1048070>.

883 Stein, B. V., Raponi, E., Sadeghi, Z., Bouman, N., Van Ham, R. C. H. J., & Back, T. (2022). A comparison of global
884 sensitivity analysis methods for explainable AI with an application in genomic prediction. *IEEE Access*,
885 10:103364–103381.

886 Strobl, C., Boulesteix, A. L., Kneib, T., Augustin, T., & Zeileis, A. (2008). Conditional variable importance for
887 random forests. *BMC bioinformatics*, 9, 1-11. <https://doi:10.1186/1471-2105-9-307>.

888 Strumbelj, E., & Kononenko, I. (2014). Explaining prediction models and individual predictions with feature
889 contributions. *Knowledge and Information Systems*, 41(3): 647–665, Dec 2014.

890 Stumm, W., & Morgan, J. J. (2012). *Aquatic chemistry: chemical equilibria and rates in natural waters*. John
891 Wiley & Sons.

892 Tonkin, M. J., & Doherty, J. (2005). A hybrid regularized inversion methodology for highly parameterized
893 environmental models: hybrid regularization methodology. *Water Resour. Res.* 41
894 <https://doi.org/10.1029/2005WR003995>.

895 Torres-Barrán, A., Alonso, Á., & Dorronsoro, J. R. (2019). Regression tree ensembles for wind energy and solar
896 radiation prediction. *Neurocomputing*, 326, 151-160. [http://dx.doi.](http://dx.doi.org/10.1016/j.neucom.2017.05.104)
897 [org/10.1016/j.neucom.2017.05.104](http://dx.doi.org/10.1016/j.neucom.2017.05.104).

898 Tunkiel, A. T., Sui, D., & Wiktorski, T. (2020). Data-driven sensitivity analysis of complex machine learning
899 models: A case study of directional drilling. *Journal of Petroleum Science and Engineering*, 195:107630,
900 2020. <https://doi.org/10.1016/j.petrol.2020.107630>.

901 Veeresh, M., Veeresh, A. V, Huddar B. D., & Hosetti, B. B. (2010). Dynamics of industrial waste stabilization
902 pond treatment processes, *Environ. Monit. Assess.*, 169, 55–65. [https://doi.org/10.1007/s10661-009-](https://doi.org/10.1007/s10661-009-1150-z)
903 [1150-z](https://doi.org/10.1007/s10661-009-1150-z).

904 Wallace, J., Champagne, P. & Hall, G. (2016). Time series relationships between chlorophyll-a, dissolved oxygen,
905 and pH in three facultative wastewater stabilization ponds, *Environ. Sci.: Water Res. Technol.*, 2016, 2,
906 1032. <https://doi.org/10.1039/C6EW00202A>.

907 Wang, F., Wang, Y., Zhang, K., Hu, M., Weng, Q., & Zhang, H. (2021). Spatial heterogeneity modeling of water
908 quality based on random forest regression and model interpretation. *Environmental Research*, 202,
909 111660. <https://doi.org/10.1016/j.envres.2021.111660>.

910 Wang, H., Jasim, A., & Chen, X. (2018). Energy harvesting technologies in roadway and bridge for different
911 applications—A comprehensive review. *Appl. Energy* 212, 1083–1094.
912 <https://doi.org/10.1016/j.apenergy.2017.12.125>.

913 Wang, P., Lu, B., Zhang, H., Zhang, W., Su, Y., & Ji, Y. (2014). Water demand prediction model based on random
914 forests model and its application. *Water Resour. Prot.* 30, 34–37.

915 Wetzel, R. G. (2001). *Limnology: Lake and river ecosystems*. 3rd ed., Academic press, San Diego. 1006p.

916 Wu, Y., & Chen, J. (2013). Investigating the effects of point source and nonpoint source pollution on the water
917 quality of the East River (Dongjiang) in South China. *Ecol Indic* 32:294–304.
918 <https://doi.org/10.1016/j.ecolind.2013.04.002>.

919 Wu, Y., & Liu, S. (2012). Modeling of land use and reservoir effects on nonpoint source pollution in a highly
920 agricultural basin. *J Environ Monit* 14(9):2350–2361. <https://doi.org/10.1039/C2EM30278K>.

921 Xie, M., Wang, N., Gaillard, J. F., & Packman, A. I. (2018). Interplay between flow and bioturbation enhances
922 metal efflux from low-permeability sediments. *J. Hazard Mater.* 341, 304–312.
923 <https://doi.org/10.1016/j.jhazmat.2017.08.002>.

924 Yim, I., Shin, J., Lee, H., Park, S., Nam, G., Kang, T., Cho, K.H., & Cha, Y. (2020). Deep learning-based retrieval of
925 cyanobacteria pigment in inland water for in-situ and airborne hyperspectral data. *Ecological Indicators*,
926 110, 105879. <https://doi.org/10.1016/j.ecolind.2019.105879>.

927 Yamamoto, R., Harada, M., Hiramatsu, K., & Tabata, T. (2021). Three-layered Feedforward artificial neural
928 network with dropout for short-term prediction of class-differentiated Chl-a based on weekly water-
929 quality observations in a eutrophic agricultural reservoir. *Paddy and Water Environment*, 1-18.
930 <https://doi.org/10.1007/s10333-021-00874-3>.

931 Zhang, W., Han, S., Zhang, D., Jin, X., Cao, E., Shan, B., & Wei, D. (2022). Preliminary Study on the Dissolved
932 Oxygen Recovery Process in Freshwater Ecosystems under the Coupling Effect of Oxygen-Consuming
933 Pollutants and Temperature. *ACS ES&T Water* 2, 1639-1646.
934 <http://dx.doi.org/10.1021/acsestwater.2c00150>

935 Zhang, Y., Yao, X., Wu, Q., Huang, Y., Zhou, Z., Yang, J., & Liu, X. (2021). Turbidity prediction of lake-type raw
936 water using random forest model based on meteorological data: A case study of Tai lake, China. *Journal*
937 *of Environmental Management*, 290, 112657. <https://doi.org/10.1016/j.jenvman.2021.112657>.

938 Zhou, C., & Zhang, J. (2023). Simultaneous measurement of chemical oxygen demand and turbidity in water
939 based on broad optical spectra using backpropagation neural network. *Chemometrics and Intelligent*
940 *Laboratory Systems*, 237, 104830. <https://doi.org/10.1016/j.chemolab.2023.104830>.

941 Zou, X. Y., Lin, Y. L., Xu, B., Guo, Z. B., Xia, S. J., Zhang, T. Y., & Gao, N. Y. (2019). A novel event detection model
942 for water distribution systems based on data-driven estimation and support vector machine
943 classification. *Water Resour. Manag.* 33 (13), 4569–4581. [https://doi.org/10.1007/s11269-019-02317-](https://doi.org/10.1007/s11269-019-02317-5)
944 5.