Development of a predictive emissions model using a gradient boosting machine learning method

Minxing Si, Ke Du



 PII:
 S2352-1864(20)31328-6

 DOI:
 https://doi.org/10.1016/j.eti.2020.101028

 Reference:
 ETI 101028

To appear in: Environmental Technology & Innovation

Received date : 31 March 2020 Revised date : 30 June 2020 Accepted date : 2 July 2020

Please cite this article as: M. Si and K. Du, Development of a predictive emissions model using a gradient boosting machine learning method. *Environmental Technology & Innovation* (2020), doi: https://doi.org/10.1016/j.eti.2020.101028.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Elsevier B.V. All rights reserved.

**Journal Pre-proof** 

# **Development of a Predictive Emissions Model using a Gradient Boosting Machine Learning Method** Minxing Si<sup>a,b</sup> and Ke Du<sup>a\*</sup> <sup>a</sup> Department of Mechanical and Manufacturing Engineering, University of Calgary, 2500 University Drive NW, Calgary AB T2N 1N4 Canada <sup>b</sup> Tetra Tech Canada Inc., 140 Quarry Park Blvd Suite 110, Calgary, AB T2C 3G3 Canada \*Corresponding author: Ke Du (email: <u>kddu@ucalgary.ca</u>)

ABSTRACT: Predictive emissions monitoring systems (PEMSs) are alternatives to continuous emissions monitoring systems (CEMSs) for monitoring air pollutants, such as NO<sub>x</sub>. Existing PEMSs and related research have focused on applying artificial neural network (ANN) algorithms. However, ANN-based models are treated as "black boxes". Regulators and decision makers without a statistical background often have difficulty understanding these models, which poses a significant challenge for a broader application of PEMSs. In this study, we proposed a tree-based ensemble method with gradient boosting techniques for PEMS development. Compared to ANNs, tree-based methods are easier to understand and require less effort to preprocess data, fewer hyperparameters for model tuning, and less time for model training. We developed a predictive model using a gradient boosting machine learning library called XGBoost to monitor NOx emissions from a boiler located in Alberta, Canada. The model uses five process parameters as inputs and the predicted NO<sub>x</sub> emissions as output. We trained the model with 202,047 samples using random search methods to determine the best model and tested the model with 50,512 samples. We evaluated the test results against US EPA PEMS standards. The model passed all the statistical tests for precision outlined by US EPA Performance Specification 16. The Pearson correlation r value was 0.98 between the XGBoost-predicted NO<sub>x</sub> values and the CEMS-measured NO<sub>x</sub> values. The RMSE was 0.14, and the MAE was 0.09. We conclude that XGBoost is a good option for developing PEMSs. Facility operators can use the method provided in this study to develop PEMSs by themselves using the open-source library XGBoost at no cost. **KEYWORDS**: Predictive Emissions Monitoring System, PEMS, XGBoost, Gradient Boosting,

NO<sub>x</sub> Monitoring

Page 2 of 34

#### Introduction 1.0

The nitrogen oxides  $(NO_x)$  family comprises seven compounds, including nitric oxide (NO), nitrogen dioxide (NO<sub>2</sub>), and other oxides of nitrogen derivatives (US EPA, 1999). NO<sub>x</sub>, or oxides of nitrogen, commonly refers to NO and NO<sub>2</sub> together (Environment and Climate Change Canada, 2019; US EPA, 2019). NO<sub>x</sub> has significant adverse environmental and health effects, contributing to acid rain and smog and forming fine particles (PM) and ozone in the ambient air (tropospheric ozone or ground-level ozone) (Boningari and Smirniotis, 2016; Mauzerall et al., 2005). Industrial facilities with large stationary combustion sources are typically required to be equipped with one or more continuous emissions monitoring systems (CEMSs) to monitor NO<sub>x</sub> emissions for compliance with regulatory emission limits (Cozza and Faulkner, 1993; US EPA, 1994; White, 1993). However, CEMSs require high initial capital and operating costs, as well as frequent maintenance and operator training (Roth and Lawrence, 2010).

To address these issues with CEMSs and offer an alternative for monitoring NOx, a predictive emissions monitoring system (PEMS) was first developed in the 1970s (Hung, 1975). Since 1990, PEMS applications have been widely reported for a range of emission sources in various industries (Chien et al., 2010, 2005; Cooper and Andreasson, 1999; Hung and Langenbacher, 1995; Kamas and Keeler, 1995). The first generation of PEMSs was developed using first principles methods based on energy and mass balances (Chien et al., 2005; Faravelli et al., 2000; Harnevie et al., 1996). The second generation of PEMSs explored the statistical relationship between operating parameters and air emissions (Chien et al., 2010; Lee et al., 2005; Saiepour et al., 2006). When the operating parameters used to build statistical models were selected, thermodynamics was often considered. The second generation was referred to as statistical methods or statistical hybrid methods. The most recent PEMSs, or third generation, were developed using machine learning algorithms called artificial neural networks (ANNs) (Cuccu et al., 2017; Vanderhaegen et al., 2010). An ANN-

Page 3 of 34

based PEMS is more flexible than a PEMS of the first or second generations because it can use operating parameters beyond those related to NO<sub>x</sub> formation. Botros *et al.* (2011) developed ANN-based predictive models for NO<sub>x</sub> monitoring in gas turbines in natural gas compressor stations. The models used four process parameters and had one hidden layer with two units. The study found that the uncertainty of NO<sub>x</sub> prediction was  $\pm 2.5\%$  to  $\pm 6\%$ . Zain and Kien Kek Chua (2011) developed ANN-based models for CO, O<sub>2</sub>, NO<sub>x</sub>, and CO<sub>2</sub> monitoring in an incinerator. The ANN used three process parameters as the input layer and had two hidden layers. Each hidden layer had 5 units. The study concluded that the accuracy of the models was 98%. Tan et al. (2016) developed NO<sub>x</sub> predictive models using a single-hidden-layer feedforward ANN for process optimization and emission reductions in a 700 MW coal-fired boiler. The mean square error (MSE) was 62.1, and the correlation coefficient was 0.98. Si et al. (2019) developed ANN-based models for NOx monitoring in a cogeneration unit. The best model had a mean absolute error (MAE) of 0.60 and Pearson r value of 0.95.

The primary regulatory frameworks for using PEMSs as compliance monitoring and reporting tools were developed by the United States Environmental Protection Agency (US EPA) under Title 40 of the Code of Federal Regulations (CFR), specifically Performance Specification 16 (PS16) and Part 75 (US EPA, 2009). The European Committee for Standardization (CEN) published a technical specification (TS) for PEMS applicability, execution, and quality assurance (CEN/TS 17198:2018) in August 2018 (European Committee for Standardization, 2018). Many countries outside of the US and Europe, such as Gulf countries, use these guidelines for PEMS certification, especially the US EPA standards (ABB S.p.A., 2014). In the US, the first regulatory approval of a PEMS for emissions monitoring and reporting occurred in June 1993 (Kamas and Keeler, 1995). To <sub>56</sub> 101 date, a few hundred PEMSs have been installed worldwide and are used for compliance

Page 4 of 34

reporting, offline what-if analysis, and analyzer availability enhancement (Eisenmann et al.,2014; Swanson, 2018).

Although PEMSs were developed 50 years ago and have evolved from first principlesbased models to machine learning-based models to become more accurate, they have not been widely installed and used because facility operators did not realize significant cost savings when using commercial PEMSs, especially for the initial capital cost. For example, Arkansas Electric Cooperative Corporation installed seven PEMSs for emissions monitoring and regulatory reporting. They estimated that the initial capital costs of installing commercial PEMSs were similar to the costs of installing CEMSs (Bivens, 2019, 2017) because of the software licensing fee and consulting cost for model training and setup.

Si et al. were the first to introduce open-source machine learning libraries for PEMS 28 113 development. Their study built feedforward ANN-based models for NO<sub>x</sub> monitoring using Google's TensorFlow and Keras libraries and aimed to provide facility operators with 30 114 <sup>32</sup> 115 methods to develop predictive emissions models by themselves (Si et al., 2019). However, <sup>34</sup> 116 the models developed in the study inherited some common drawbacks of ANN algorithms. For example, ANN models are treated as "black boxes" because strong statistical knowledge is needed to understand ANN algorithms. Decision makers and regulators are often skeptical about the outputs produced from these black boxes due to the lack of statistical background (Guelman, 2012). 43 120

In this study, we developed a predictive model for NO<sub>x</sub> emission monitoring and
 regulatory reporting using a tree-based ensemble machine learning algorithm called
 XGBoost. XGBoost was implemented as an open-source machine learning library and
 developed based on gradient boosting techniques introduced by Friedman (2001). Gradient
 boosting has been successfully applied to tree models (Guelman, 2012; Krauss et al., 2017;
 Semanjski and Gautama, 2015; Zhang and Haghani, 2015). Gradient tree boosting adds
 multiple weak trees together and forms a strong learner to optimize predictive performance.

This method is also known as a gradient boosted regression tree (GBRT) and a gradient boosting machine (GBM). XGBoost was first released in 2014 (Labram, 2019) and has since become popular among machine learning practitioners. In 2015, 17 out of 29 winners in the Kaggle machine learning competitions used the XGBoost algorithm (Nielsen, 2016). XGBoost is easier to interpret than ANN-based algorithms, so decision makers and regulators can better understand how the predictive models work and how outcomes are produced. In addition, XGBoost requires less effort for data preprocessing and has fewer hyperparameters for tuning.

Existing PEMSs and related research have focused on the use of ANN algorithms. To the best of our knowledge, the gradient boosting approach has not been applied to predictive models for emissions monitoring and reporting, even though PEMSs have been **139** researched, developed, and installed for almost five decades. In addition, according to US EPA and EU PEMS standards, PEMSs need to be able to generate data for substitution when **140** <sup>32</sup> 141 one or more process parameters (sensors) fail, such as missing data. Previous studies <sup>34</sup> 142 removed missing data for model development. In this study, we test a method for data substitution to bridge this gap. Using the methods presented in this study, facility operators will have more complete knowledge to develop their emission models by themselves and obtain regulatory approval.

43 146 2.0

#### Material and Method

#### Tree-based Methods and Tree-based Ensemble Learning 2.1

Tree-based methods build predictive models by repeatedly splitting the predictor space into rectangles by some criteria and then fitting a model for each rectangle. Figure 1 shows that the space for two predictors, fuel gas flow and exhaust gas temperature, is divided into three rectangles, as summarized in Figure 2. R represents the divided areas in Figure 1 called the terminal regions or leaf nodes in the tree model shown in Figure 2. *j* represents

Page 6 of 34





A prediction method using tree-based ensemble learning is presented in Figure 3.



#### 169 Figure 3: Example of Tree-based Ensemble Learning

#### 0 2.2 Gradient Boosting and XGBoost

Boosting is an ensemble method that combines multiple weak learners to produce a power learner. A boosting method adds new learners  $f_x$  (tree or estimator) sequentially. At iteration k, the new estimator  $f_k(x_i)$  tries to correct the previous prediction  $\hat{y}_i^{(k-1)}$  or  $F^{(k-1)}(x)$ and generates a new prediction  $\hat{y}_i^{(k)}$  or  $F^{(k)}(x)$ , and this process can be written as Equation  $F^{(k-1)}(x)$  and  $F^{(k)}(x)$  represent the functions that produce the predicted values at iterations k - 1 and k, respectively. Iteration k also indicates that k trees are ensembled in the model.  $\hat{y}_i^{(k)} = \hat{y}_i^{(k-1)} + f_k(x_i)$ 

Algorithm 1: Gradient Tree Boosting AlgorithmNoteInput:Dataset 
$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$
 $x_i$  represents the values of the inputPage 8 of 34

(2)

$y_i$ represents the observed value for the i <sup>th</sup> sample. n represents the number of samples in the training dataset $L = \frac{(Observed - Predicted)^2}{2}$ F(x) is the function that produces the predicted values $\rho$ is an initial predicted value. $\rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. k is the k <sup>th</sup> tree
the i <sup>th</sup> sample. <i>n</i> represents the number of samples in the training dataset $L = \frac{(Observed - Predicted)^2}{2}$ <i>F</i> ( <i>x</i> ) is the function that produces the predicted values $\rho$ is an initial predicted value. $\rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. <i>k</i> is the k <sup>th</sup> tree
<i>n</i> represents the number of samples in the training dataset $L = \frac{(Observed - Predicted)^2}{2}$ $F(x)$ is the function that produces the predicted values $\rho$ is an initial predicted value. $\rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. <i>k</i> is the k <sup>th</sup> tree
in the training dataset $L = \frac{(Observed - Predicted)^2}{2}$ $F(x)$ is the function that produces the predicted values $\rho$ is an initial predicted value. $\rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. k is the k <sup>th</sup> tree
$L = \frac{(Observed - Predicted)^2}{2}$ $F(x) \text{ is the function that produces}$ the predicted values $\rho \text{ is an initial predicted value. } \rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. $k \text{ is the } k^{\text{th}} \text{ tree}$
F(x) is the function that produces the predicted values $\rho$ is an initial predicted value. $\rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. k is the k <sup>th</sup> tree
the predicted values $\rho$ is an initial predicted value. $\rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. k is the k <sup>th</sup> tree
$\rho$ is an initial predicted value. $\rho$ could be the mean of the observed values. In XGBoost, the default value is 0.5. k is the k <sup>th</sup> tree
could be the mean of the observed values. In XGBoost, the default value is 0.5. <i>k</i> is the k <sup>th</sup> tree
values. In XGBoost, the default value is 0.5. <i>k</i> is the k <sup>th</sup> tree
value is 0.5. $k$ is the k <sup>th</sup> tree
k is the k <sup>th</sup> tree
i is the i <sup>th</sup> sample in the training
data
j is the j <sup>th</sup> terminal region (leaf)
$\rho_k$ is the output of the $\mathbf{k}^{\mathrm{th}}$ estimator
(tree) $f_k(x)$

 $F^{(k-1)}(x)$  is the previous prediction

 $\sum_{j=1}^{J_k}$  summation is used when a

in Step c).

 $\rho_{jk}$  is the output from the tree made

single sample ends up in multiple leaf nodes 3 Output  $F^{(K)}(\mathbf{x})$  $\hat{y}_i = F^0(x) + \sum_{k=1}^{K} f_k(x_i)$ The learning steps by gradient boosting are illustrated in Figure 4.  $\hat{y}_{i}^{0} = F^{0}(x)$  $\hat{y}_i^1 = F^1(x) = F^0 + f_1(\checkmark)$  $\hat{y}_i^2 = F^2(x) = F^0 + f_1(-) + f_2(-)$  $\hat{y}_{i}^{k} = F^{k}(x) = F^{0} + f_{1}(\mathbf{x}) + f_{2}(\mathbf{x}) + \mathbf{o} + f_{k}(\mathbf{x})$ Figure 4: Illustration of Gradient Boosting Unlike traditional boosting, XGBoost adds a regularization function to prevent overfitting and optimizes the loss function by a Taylor expansion (Chen and Guestrin, 2016). The objective function J in XGBoost describes the model's performance and can be written as Equation 3:  $J = \sum_{i=0}^{n} L(y_i, \hat{y}_i) + \sum_{i=0}^{K} \Omega(f_k)$ (3) where n is the number of training samples and  $\Omega(f_k)$  is a regularization function.  $\Omega(f_k)$  is written as Equation 4:

Page 10 of 34

to determine the output values of each leaf

 $F^{(k)}(x) = F^{(k-1)}(x) + \eta \times \sum_{i=1}^{j_k} \rho_{jk} \mathbf{I}(x \in R_{jk})$ 

d) Update

$$\Omega(f_k) = \gamma \mathbf{T} + \frac{1}{2} \lambda \sum_{j=0}^{T} \mathbf{w}_j^2$$
(4)

where *T* is the number of leaf nodes and  $\gamma$  and  $\lambda$  are hyperparameters in XGBoost that users can define. *w* is the leaf weight (i.e., a predictive value in a terminal node). The objective function for the m<sup>th</sup> iteration can be written as Equation 5, where m represents elements of *K* iterations ( $m \in \{k_1, k_2, ..., K\}$ ).

$$J^{m} = \sum_{i=1}^{n} L(y_{i}, \hat{y}_{i}^{(m)}) + \sum_{k=1}^{m} \Omega(f_{m}(x_{i})) = \sum_{i=1}^{n} L(y_{i}, \hat{y}_{i}^{m-1} + f_{m}(x_{i})) + \sum_{k=1}^{m} \Omega(f_{m}(x_{i}))$$
(5)

To optimize the objective function, the XGBoost algorithm takes the following steps:

 Use a second-order Taylor expansion to approximate the derivatives, and write J<sup>m</sup> as Equation 6:

$$J^{m} \simeq \sum_{i=1}^{n} [L(y_{i}, \hat{y}_{i}^{(m-1)}) + g_{m}(x_{i})f_{m}(x_{i}) + \frac{1}{2}h_{m}(x_{i})f_{m}(x_{i})^{2}] + \sum_{k=1}^{m} \Omega(f_{m}(x_{i}))$$
(6)

where  $L(y_i, \hat{y}_i^{(m-1)})$  is the loss function for the previous prediction and  $g_m(x_i)$  represents the first derivative of the loss function  $L(y_i, \hat{y}_i^{(m-1)})$ , called the gradient.  $h_m(x_i)$  is the second derivative of the loss function  $L(y_i, \hat{y}_i^{(m-1)})$ , called the Hessian.  $g_m(x_i)$  and  $h_m(x_i)$ are written as follows:

$$g_m(x_i) = \frac{dL(y_i, \hat{y}_i^{(m-1)})}{d\hat{y}_i^{(m-1)}} \text{ and } h_m(x_i) = \frac{d^2L(y_i, \hat{y}_i^{(m-1)})}{d(\hat{y}_i^{(m-1)})^2}$$

2) Remove the constant term  $\sum_{i=1}^{n} L(y_i, \hat{y}_i^{(m-1)})$  because it is not related to the output values  $f_m(x_i)$  and has no effect on optimizing the objective function. Equation 6 is then written as Equation 7:

$$J^{m} = \sum_{i=1}^{n} [g_{m}(x_{i})f_{m}(x_{i}) + \frac{1}{2}h_{m}(x_{i})f_{m}(x_{i})^{2}] + \sum_{k=1}^{m} \Omega(f_{m}(x_{i}))$$
(7)

3) Replace  $f_m(x_i)$  with the sum of the tree leaves written as Equation 8. The sum of the tree leaves is the sum of the outputs of the trees:

$$f_m(x_i) = \sum_{j=1}^T w_{jm} l(x \in R_{jm})$$
(8)

Replace Equation 7 with Equation 8, and 
$$J^m$$
 is written as Equation 9:

$$J^{m} = \sum_{i=1}^{n} [g_{m}(x_{i}) \sum_{j=1}^{T} w_{jm} + \frac{1}{2} h_{m}(x_{i}) \sum_{j=1}^{T} w_{jm}^{2}] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^{T} w_{jm}^{2}$$
(9)

The sums of  $g_m(x_i)$  and  $h_m(x_i)$  can be simplified as in Equation 10:

$$G_{jm} = \sum_{i \in I_{jm}} g_m(x_i), H_{jm} = \sum_{i \in I_{jm}} h_m(x_i)$$

$$(10)$$

where  $I_{jm}$  denotes the set of training samples (instances,  $x_i$ ) in the region (leaf)  $R_{jm}$ . 4) Replace  $g_m(x_i)$  and  $h_m(x_i)$  in Equation 9 with Equation 10. J<sup>m</sup> is written as Equation

20 215 

**227** 

$$J^{m} = \sum_{j=1}^{T} [G_{jm}w_{jm} + \frac{1}{2}(H_{jm} + \lambda)w_{jm}^{2}] + \gamma T$$
(11)

5) Minimize the objective  $J^m$  in Equation 11 by the weight  $w_{jm}$  for each leaf using the

derivative of  $J^m$  with respect to  $w_{jm}$ , as written in Equation 12.

$$\frac{\partial J^m}{\partial w_{jm}} = G_{jm} + (H_{jm} + \lambda)w_{jm} = 0$$
(12)

#### 6) Calculate the best weight $w_{jm}$ using Equation 13.

$$w_{jm} = -\frac{G_{jm}}{H_{jm} + \lambda} \tag{13}$$

For regression with MSE used for the loss function, the gradient  $G_{jm}$  is the sum of negative residuals, or  $-\sum_{i=1}^{n}(y_i - \hat{y}_i)$ , and the Hessian  $H_{jm}$  is the number of residuals

#### Therefore, Equation 13 is used to calculate the output value for a leaf.

7) Replace  $w_{jm}$  in Equation 11 with Equation 13; then, the objective  $J^m$  in Equation 11 **228** 

for

n.

#### the best tree structure can be written as Equation 14.

$$\mathbf{J}^m = -\frac{1}{2} \sum_{j=1}^{T} \left[ \frac{G_{jm}^2}{H_{jm} + \lambda} \right] + \gamma \mathbf{T}$$

(14)

#### 2.3 Data

б 

The predictive model was developed to monitor NO<sub>x</sub> emissions from a boiler. The NO<sub>x</sub> **234** emissions were continuously monitored by a CEMS unit. The CEMS unit was an in situ **235** monitoring system with path in situ analyzers manufactured by SICK. The NO<sub>x</sub> <sup>17</sup> 236 concentrations were measured by a differential optical absorption spectroscopy method. The NO<sub>x</sub> concentration range of the analyzer in the CEMS unit was 0-100 ppm by volume on a wet basis with 0-80 ppm for NO and 0-20 ppm for NO<sub>2</sub>. The CEMS unit measured NO<sub>x</sub> in ppm, exhaust temperature, and exhaust flow and then converted NO<sub>x</sub> ppm into kg/h. The NO<sub>x</sub> mass flow, exhaust temperature, and exhaust flow data were reported to Alberta's **241** provincial regulatory agency under the facility's operational permit. The facility CEMS report **242** included two parts: station status (SS), which records the boiler downtime, and record <sup>32</sup> 243 details (RD), which records the detailed NO<sub>x</sub> emissions data, including substitute data when the CEMS unit was offline or in maintenance.

We selected five process parameters as inputs for machine learning (Table 1). NOx in kg/h was the output of the model. The input parameters are also called PEMS sensors. We retrieved process data and NO<sub>x</sub> emissions data at one-minute intervals from 00:00 on January 1, 2019, to 23:59 on June 30, 2019, from a data storage system called the process **248 249** historian database (PHD). The data of the five process parameters and NO<sub>x</sub> emissions formed a structured dataset with 260,641 rows and 6 columns. Each row is called one sample or example. The five sensors were integrated into the facility's distributed control system and had routine preventative maintenance schedules for calibration and inspection to ensure that the instrument worked properly. The metered data obtained by the five **254** sensors were automatically transferred to the PHD.

- **255 Table 1: Model Inputs**

Sensor ID	Input (process) parameters	Unit
S0	fuel gas temperature	°C
S1	boiler feed water temperature	°C
S2	fuel gas flow	m³/h
S3	combustion air flow	10 <sup>3</sup> m <sup>3</sup> /day
S4	Exhaust gas temperature	°C
The following step • Step 1: Re	os were taken to process the raw data: move samples when the boiler was offline	e based on the SS records
because th	e	
facility doe	s not report $NO_x$ emissions in these period	ds (7,326 samples removed).
• Step 2: Re	move samples that have substitute $NO_x$ e	missions data (240 samples
removed).		
• Step 3: Re	move samples for which the fuel gas flow	was 0 (515 samples removed).
There was	a delay between the time that the fuel ga	is supply was stopped and the
time that t	he exhaust exited the stack. The CEMS re	corded small $NO_x$ readings for
approxima	tely 10 minutes after the fuel gas reading	s were 0. The $NO_x$ values ranged
from 0 to 5	5.3 kg/h and contributed to 0.007% of NC	D <sub>x</sub> emissions.
• Step 4: Sh	uffle the remaining dataset (252,559 sam	ples), and randomly split 80% o
the sample	es into the training dataset and 20% into t	the test dataset. The training
dataset ha	d 202,047 samples and was used for mod	lel building and training. The test
dataset ha	d 50,512 samples and was used to test th	ne model predictive power for an
unbiased e	valuation.	
2.4 Machi	ne Learning and Model Evaluation	
The paired proces	s data and CEMS-measured $NO_x$ emission	s data were used to train the
PEMS model using	the XGBoost library (version 0.90) in a F	Python (version 3.7.4)
environment. The	CEMS was the reference method in this s	tudy. Ten-fold cross-validation Page 14 of 3

2			
3 4 5	277	was us	sed to select the best model during training. The following hyperparameters were
6 7	278	tuned	by a random search technique:
8 9	279	•	Number of trees to fit (n_estimator)
10 11	280	•	Maximum depth of a tree (max_depth)
12 13	281	•	Step-size shrinkage used in update (learning_rate, $\eta$ )
14 15	282	•	Subsample ratio of columns when constructing each tree (colsample_bytree)
16 17	283	•	Minimum loss reduction required to make a further partition on a leaf node of the
19 20	284		tree (gamma, γ)
21 22	285	•	L2 regularization on weights (reg_lambda, $\lambda$ )
23 24	286	•	Minimum sum of instance weight needed in a child (min_child_weight)
25 26	287		A detailed explanation of each hyperparameter is provided in the XGBoost
27 28	288	docum	entation (XGBoost developers, 2019). We used a random search method for model
29 30	289	tuning	using the scikit-learn (0.21.3) library. "Random search" means that combinations of
31 32	290	the pa	rameters are randomly selected to find the best model structure.
34 35	291		We compared the $NO_x$ outputs from the model using the test dataset with the
36 37	292	corres	ponding CEMS-measured NO $_{\rm x}$ data. The results were evaluated against the criteria
38 39	293	outline	ed in the US EPA PS16 standard for PEMS precision. The criteria are presented in Table
40 41	294	2.	
42 43	295	Table	2: Precision Requirements for Predictive Models

Statistical test required by Criteria **PS16** Bias test  $d_{avg} \leq |cc|$ The mean difference  $d_{avg}$  between Reference Method (RM) values and predicted values is less than or equal to the absolute value of the confidence coefficient (CC) at a 97.5% one-sided confidence Page 15 of 34

	interval.
	Otherwise, a bias factor needs to be applied to the
	predicted values.
Pearson correlation coefficient	r ≥ 0.8
	The correlation between CEMS values and predicted
	values must be 0.8 or greater.
F-test	$F_{value} \leq F_{critical}$
	The variance ratio of predicted values and CEMS
	measured values must be less than or equal to the

The model was also assessed by root mean square error (RMSE), MAE and MSE using the **297** <sup>30</sup> 298 test dataset.

#### 2.5 **Model Sensitivity Test**

<sub>35</sub> 300 A model sensitivity test was used to evaluate changes in model performance using substitute data when one or more sensors failed. As required by the EU PEMS TS and the US **301 302** EPA's PEMS standards, a PEMS needs to have a sensor validation system to identify sensor 41 303 failure hourly. The sensor validation system is responsible for generating substitute data, <sup>43</sup> 304 informing operators when sensors (process instruments, such as a fuel gas flow meter) need repair and indicating that the PEMS is out of control (US EPA, 2006). We followed the procedures for model sensitivity testing provided by the EU PEMS TS and US EPA CFR 40 Part 75 Subpart E. The procedures are summarized as follows: <sub>52</sub> 308 Select a set of reference sensor values and NO<sub>x</sub> emissions values. Ι. **309** II. Artificially fail one sensor, and then run the predictive model using substitute data. We assessed the effect on the model's accuracy by calculating the hourly percentage

difference between the reference NOx values in Step I and the predicted values using

Page 16 of 34

substitute data in Step II. Repeat this procedure for all the sensors used by the model individually.

314 III. Flag the outputs as invalid if the hourly percentage difference is greater than 10%.

IV. Perform a two-sensor failure test by repeating Step I, and then identify invalid data
by conducting Step III. We only performed the test for up to two sensor failures
because the US EPA requires a predictive model to have at least three input sensors,
and this study used five inputs for machine learning.

We used directly measured sensor data and NO<sub>x</sub> values from 00:00 on March 30,

320 2019, to 23:59 on April 5, 2019, for the set of reference sensor values and NO<sub>x</sub> emission

321 values as required in Step I. The boiler was in a normal operational condition in this period.

The data substitution method used for the sensitivity test is illustrated in Table 3. Data were

substituted by interpolating the time distance. We tested sensor failure for one hour (1H),

one day (1D), two days (2D), three days (3D), four days (4D), and five days (5D).

According to Alberta's CEMS code, a CEMS can only allow a sensor to be out of control or offline for up to five days.

327 Table 3: Example of Data Substitution by Interpolating the Time Distance

Timestamp	Fuel Gas Flow (10 <sup>3</sup> m <sup>3</sup> /h)	Substitute data
March 31, 2019 00:00	5	
March 31, 2019 00:01	Invalid data (sensor failure)	6
March 31, 2019 00:02	Invalid data (sensor failure)	7
March 31, 2019 00:03	Invalid data (sensor failure)	8
March 31, 2019 00:04	9	

## 29 3.0 Results and Discussion

## 30 3.1 Hyperparameter Tuning

The results of the random search for the hyperparameters are presented in Table 4. The gradient boosting model comprised 38 trees  $f_{K=38}(x)$ , and the learning rate  $\eta$  was 0.405. The

Page 17 of 34

333	two user-defined values $\gamma$ and $\lambda$ in the regularization function were 0.12 and 6	9.05,
334	respectively.	
335	Table 4: Random Search Results for Hyperparameters	
	XGBoost Hyperparameters	Values
	Number of trees to fit (n_estimator)	38
	Maximum depth of a tree (max_depth)	32
	Step-size shrinkage used in update (learning_rate, $\eta$ )	0.41
	Subsample ratio of columns when constructing each tree (colsample_bytree)	0.96
	Minimum loss reduction required to make a further partition on a leaf node	0.12
	of the tree (gamma, $\gamma$ )	0.12
	L2 regularization on weights (reg_lambda, $\lambda$ )	69.05
	Minimum sum of instance weight needed in a child (min_child_weight)	61.63

#### Evaluation 3.2

The results of the statistical tests are presented in Table 5. The model passed all three EPA requirements for precision. The RMSE was 0.14, and the MAE was 0.09.

#### **340**

#### **Table 5: Model Evaluation Results**

Statistical tests	Results
Bias test	$d_{avg} = -0.0011,  cc  = 0.0012$
	$d_{avg} <  cc $ ; the model passed the EPA requirement. Bias
	adjustment is not needed.
Pearson r value	r = 0.98
	r > 0.8; the model passed the EPA requirement.
F-test	$F_{value} = 0.96$ , $F_{critical} = 1.01$
	$F_{value} < F_{critical}$ ; the model passed the EPA requirement.
RMSE	0.14







#### **366** Figure 7: Importance of the Input Parameters

**365** 

#### **Model Sensitivity** 3.3

Figure 8 shows the absolute differences between the CEMS measured values and the **370** predicted values with substitute data when one sensor failed. The largest differences were found when the boiler feed water temperature sensor failed for five days (3.28%) and when **371** <sup>38</sup> 372 the exhaust gas temperature sensor failed for three days (3.01%). The differences were within the maximum allowable difference of 10%. The impact on the model performance of using substitute data for fuel gas temperature, fuel gas flow, and combustion air flow was less than 1% when one of the sensors failed for up to five days.

The absolute differences between the CEMS measured values and the predicted values with substitute data for two-sensor failure are presented in Figure 9. The largest difference was found when fuel gas temperature and boiler feed water temperature failed **378 379** together for five days (5.14%), but the difference was still within the 10% regulatory limit.



and sensor 1, boiler feed water temperature.

S04

S12

S13

IDs. For example, S01 represents the failures of sensor 0, fuel gas temperature,

Figure 9: Sensitivity Test for Two-Sensor Failure. The x-axis is a combination of sensor

S14

S23

S24

S34

Page 22 of 34

62 63 64

61

47 48 49

50 51

55 56 **386** 

57

58 **387** 59 60 0%

S01

S02

S03

In addition to the precision requirements that are outlined in Table 2, the US EPA (2018) **391** also has the following requirements to validate PEMSs for emissions monitoring (Table 6) in **392** CFR 40 PS16 and Part 75.

#### <sup>15</sup> 393 **Table 6: Other Regulatory Requirements for PEMSs**

Requirement	Criteria
Reliability	PEMS availability should be greater than or equal to 95%.
Quality Assurance	Input sensors must operate within the permitted ranges,
	such as model training ranges and manufacturing ranges.
	Daily check to ensure model is not modified
	Input sensors must be maintained in accordance with the
	manufacturer's recommendations.
	A PEMS should be equipped with an alarm system. The
	alarm system will inform facility operators when the PEMS is
	out of control, such as sensors out of permitted ranges and
	sensor failures.
	Routine relative accuracy test audits must be conducted to
	ensure a constant performance of a PEMS after initial
	certification.

#### 3.5 **Comparison with an Artificial Neural Network**

We also developed NO<sub>x</sub> predictive models using a feedforward neural network algorithm for comparison with the XGBoost model. The ANN model was trained using random search methods for the best model structure. The ANN model was trained and tested with the same 

1 2						
3 4 5	399	training and test	datasets as the XGBoost	model. The best ANN model had the following		
6 7	400	hyperparameters	:			
8 9	401	Number of	hidden layers: 3			
10 11	402	Units in ea	ich hidden layer: 128			
12 13	403	• Learning r	ate: 0.0005			
14 15	404	L2 regular	ization: 0.01			
16 17	405	Optimizati	on: Nadam			
18 19	406	The ANN mod	el inputs and output (inp	ut and output layers) were the same as those in		
20	407	XGBoost.				
22 23 24	408	The Pearson r val	ue for the XGBoost mode	el was 0.98, compared to 0.92 for the ANN model.		
25 26	409	The RMSE for the	XGBoost model was 56%	6 less than that for the ANN model. The MAE for		
27 28	410	the XGBoost mod	el was 61% less than tha	at for the ANN model (Table 7). The higher r value	e	
29 30	411	and lower RMSE i	ndicated that the XGBoos	st model was a better model than the ANN model.		
31		Table 7: Comparison between VCPeers and the ANN weing the Test Peterst				
32	412	Table 7: Compa	rison between XGBoos	t and the ANN using the Test Dataset		
32 33 34	412	Table 7: Compa	rison between XGBoos	t and the ANN using the Test Dataset		
32 33 34 35	412	Table 7: Compa Statistical tests	rison between XGBoos XGBoost	t and the ANN using the Test Dataset		
32 33 34 35 36 37	412	Table 7: CompaStatistical testsPearson r	rison between XGBoos XGBoost 0.98	t and the ANN using the Test Dataset           ANN           0.92		
32 33 34 35 36 37 38 39 40	412	Table 7: CompaStatistical testsPearson rRMSE	rison between XGBoos XGBoost 0.98 0.14	t and the ANN using the Test Dataset          ANN         0.92         0.32		
32 33 34 35 36 37 38 39 40 41 42	412	Table 7: CompaStatistical testsPearson rRMSEMAE	rison between XGBoos XGBoost 0.98 0.14 0.09	t and the ANN using the Test Dataset          ANN         0.92         0.32         0.23		
32 33 34 35 36 37 38 39 40 41 42 43 44	412 413	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test data	rison between XGBoos XGBoost 0.98 0.14 0.09 taset contains 50,512 sa	ANN   0.92   0.32   0.23   mples at 1-minute intervals.		
32 33 34 35 36 37 38 39 40 41 42 43 44 45 46	412 413 414	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test da3.6Mode	rison between XGBoos XGBoost 0.98 0.14 0.09 taset contains 50,512 sa I Performance under	ANN   0.92   0.32   0.23   mples at 1-minute intervals. Non-normal Operating Conditions		
32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48	412 413 414 415	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test da3.6ModeUnder non-normal	rison between XGBoos XGBoost 0.98 0.14 0.09 taset contains 50,512 sat I Performance under	ANN         0.92         0.32         0.23         mples at 1-minute intervals.         Non-normal Operating Conditions         -NOCs), such as equipment startup and		
32 33 34 35 36 37 38 40 41 42 43 44 45 46 47 48 49 50	412 413 414 415 416	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test da3.6ModeUnder non-normalshutdown, the RM	XGBoost         0.98         0.14         0.09         taset contains 50,512 sat         I Performance under         I operating conditions (N         ISE increased from 0.13	ANN         0.92         0.32         0.23         mples at 1-minute intervals.         Non-normal Operating Conditions         -NOCs), such as equipment startup and         to 0.72 for XGBoost (Fig. 10) compared to that of	-	
32 33 34 35 37 38 40 41 423 44 45 46 47 48 49 512 52	412 413 414 415 416 417	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test da3.6ModeUnder non-normalshutdown, the RMthe normal operation	rison between XGBoos XGBoost 0.98 0.14 0.09 taset contains 50,512 sai I Performance under I operating conditions (N 1SE increased from 0.13 f	ANN         0.92         0.32         0.23         mples at 1-minute intervals.         Non-normal Operating Conditions         -NOCs), such as equipment startup and         to 0.72 for XGBoost (Fig. 10) compared to that of         e Pearson r value between the CEMS measured	=	
32 33 34 35 37 38 40 412 434 456 47 489 5512 553 552 553 552	412 413 414 415 416 417 418	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test data3.6ModeUnder non-normalshutdown, the RMthe normal operativalues and the m	XGBoost         0.98         0.14         0.09         taset contains 50,512 sat         Performance under         I operating conditions (N         1SE increased from 0.13 ft         ting condition (NOC). The         odel predicted values decord	ANN         0.92         0.32         0.23         mples at 1-minute intervals.         Non-normal Operating Conditions         -NOCs), such as equipment startup and         to 0.72 for XGBoost (Fig. 10) compared to that of         e Pearson r value between the CEMS measured         creased from 0.99 under the NOC to 0.91 under	Ŧ	
32 334 35 36 37 390 412 443 445 467 495 552 5555 555 555 5555 5555 5555 5555 5555 5555 5555	412 413 414 415 416 417 418 419	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test data3.6ModeUnder non-normalshutdown, the RMthe normal operativalues and the mN-NOCs for the X	XGBoost         0.98         0.14         0.09         taset contains 50,512 sat         Performance under         I operating conditions (N         ISE increased from 0.13 f         ting condition (NOC). The         odel predicted values dec         GBoost model. The increased	ANN         0.92         0.32         0.23         mples at 1-minute intervals.         Non-normal Operating Conditions         -NOCs), such as equipment startup and         to 0.72 for XGBoost (Fig. 10) compared to that of         e Pearson r value between the CEMS measured         creased from 0.99 under the NOC to 0.91 under         ase in RMSE and the decrease in the r value	F	
32 33 35 36 37 38 412 443 445 467 489 512 555	412 413 414 415 416 417 418 419 420	Table 7: CompaStatistical testsPearson rRMSEMAENote: The test data3.6ModeUnder non-normalshutdown, the RMthe normal operativalues and the mN-NOCs for the Xindicated that the	XGBoost         0.98         0.14         0.09         taset contains 50,512 sat         Performance under         I operating conditions (N         1SE increased from 0.13 ft         ting condition (NOC). The         odel predicted values dec         GBoost model. The increase         xGBoost model did not predicted did not predicted values did not predicted	ANN         0.92         0.32         0.23         mples at 1-minute intervals.         Non-normal Operating Conditions         -NOCs), such as equipment startup and         to 0.72 for XGBoost (Fig. 10) compared to that of         e Pearson r value between the CEMS measured         creased from 0.99 under the NOC to 0.91 under         ase in RMSE and the decrease in the r value         perform as well under N-NOCs as under the NOC.	F.	



# Figure 10: Non-Normal Operating Condition (N-NOC) vs. Normal Operating Condition (NOC). The NOC comprised 252,180 samples at 1-minute intervals. The N-NOC comprised 379 samples at 1-minute intervals.

Using 1-minute high-resolution data, the ANN model generated higher predicted values when the NO<sub>x</sub> mass flow rate was less than 2 kg/h and generated lower predicted values when the NO<sub>x</sub> mass flow rate was greater than 8 kg/h. In contrast, the XGBoost model performed better for flow rates less than 2 kg/h than the ANN model. However, the XGBoost model still produced values less than 12 kg/h when the CEMS values were greater than 12 kg/h (Fig. 11).



The facility reports its NOx emissions on an hourly basis, and the emission limit for the facility under its operation permit was set in kg/h. We averaged the 1-minute data to hourly data and compared the model performance. On an hourly basis, no emission rates were greater than 12 kg/h. Using the hourly data, the Pearson r was 0.93 for the ANN and 0.99 for XGBoost. All data provided by XGBoost fall close to the 1:1 ratio line (Fig. 12). The XGBoost model showed a higher accuracy than the ANN model when reporting emissions on an hourly basis.

 Page 26 of 34





465       Figure 14: Model evaluation of a process upset event using data at 1-hou         466       intervals. The data period is from 00:00 on May 1, 2019, to 23:         467       1, 2019.         468       3.7         469       3.7         470       Although a predictive model is the core component of a PEMS, facility operators of         471       consider the following when building a PEMS on their own:         11       1)       Deploying predictive models to existing process control networks. Process         472       1)       Deploying predictive models to existing process control networks. Process         473       devices, such as distributed control systems, on process networks may no         474       sufficient computational power to run complicated predictive models.         475       2)       Integrating predictive models into existing data acquisition systems and rus         476       systems.       3)         477       3)       Multiple models may be needed to replace one CEMS unit. For example, fa         478       Alberta, Canada, are required to report exhaust flow and temperature, in         479       mass emission rates.         480       Predictive modeling, including ANN and XGBoost algorithms, has shown great periodictive modeling.         481       for emission monitoring. Computer-based emission monitoring methods can	r DO on May weed to control
6intervals. The data period is from 00:00 on May 1, 2019, to 23:84671, 2019.114683.7Limitations and Future Study124693.7Limitations and Future Study14470Although a predictive model is the core component of a PEMS, facility operators of16471consider the following when building a PEMS on their own:184721) Deploying predictive models to existing process control networks. Process19473devices, such as distributed control systems, on process networks may not21474sufficient computational power to run complicated predictive models.24474sufficient computational power to run complicated predictive models.254752) Integrating predictive models into existing data acquisition systems and re264773) Multiple models may be needed to replace one CEMS unit. For example, fa27478Alberta, Canada, are required to report exhaust flow and temperature, in28479mass emission rates.29478small stationary combustion sources that do not have regulatory requirements for28481for emission monitoring. A PEMS provides a more accurate emissions reporting methods29483sources than engineering estimation methods with generic emission factor39484small sources than engineering estimation methods with generic emission factor484484small sources than engineering estimation methods with generic emission factor484484small sources than engineerin	<b>DO on May</b> weed to control
8       467       1, 2019.         468       3.7       Limitations and Future Study         470       Although a predictive model is the core component of a PEMS, facility operators of consider the following when building a PEMS on their own:         471       consider the following when building a PEMS on their own:         18       472       1) Deploying predictive models to existing process control networks. Process devices, such as distributed control systems, on process networks may not sufficient computational power to run complicated predictive models.         24       2474       sufficient computational power to run complicated predictive models.         25       475       2) Integrating predictive models into existing data acquisition systems and response.         26       477       3) Multiple models may be needed to replace one CEMS unit. For example, for systems.         26       478       Alberta, Canada, are required to report exhaust flow and temperature, in mass emission rates.         27       478       Alberta, Computer-based emission monitoring methods can be in small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting methods with generic emission factor normally used for regulatory reporting. Predictive modeling also offers the potent of normally used for regulatory reporting. Predictive modeling also offers the potent of mortally used for regulatory reporting. Predictive modeling also offers the potent of the following winhout high capital cost. However, further research is	eed to
468 <b>3.7</b> Limitations and Future Study         470       Although a predictive model is the core component of a PEMS, facility operators of consider the following when building a PEMS on their own:         471       consider the following when building a PEMS on their own:         11       Deploying predictive models to existing process control networks. Process devices, such as distributed control systems, on process networks may not sufficient computational power to run complicated predictive models.         473       2) Integrating predictive models into existing data acquisition systems and response.         474       3) Multiple models may be needed to replace one CEMS unit. For example, for Alberta, Canada, are required to report exhaust flow and temperature, in mass emission rates.         479       mass emission rates.         481       for emission monitoring. Computer-based emission monitoring methods can be in small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting methods with generic emission factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emisting without high capital cost. However, further research is needed because of the follower for the follow	eed to control
4470Although a predictive model is the core component of a PEMS, facility operators of6471consider the following when building a PEMS on their own:94721) Deploying predictive models to existing process control networks. Process11) Deploying predictive models to existing process control networks. Process11) Deploying predictive models to existing process control networks. Process12) Integrating predictive models into existing data acquisition systems and resisting predictive models into existing data acquisition systems and resistence13) Multiple models may be needed to replace one CEMS unit. For example, fa1478Alberta, Canada, are required to report exhaust flow and temperature, in1mass emission rates.6480Predictive modeling, including ANN and XGBoost algorithms, has shown great per1481for emission monitoring. Computer-based emission monitoring methods can be in2483small stationary combustion sources that do not have regulatory requirements for2484small sources than engineering estimation methods with generic emissions factor2486normally used for regulatory reporting. Predictive modeling also offers the potent487without high capital cost. However, further research is needed because of the followed for continuous for the followed for the followed for the followed for the followed for continuous for continuous for continuous for continuous for the followed for the followe	eed to control
<ul> <li>471 consider the following when building a PEMS on their own:</li> <li>472 1) Deploying predictive models to existing process control networks. Process 473 devices, such as distributed control systems, on process networks may not 474 sufficient computational power to run complicated predictive models.</li> <li>475 2) Integrating predictive models into existing data acquisition systems and run 476 systems.</li> <li>477 3) Multiple models may be needed to replace one CEMS unit. For example, for 478 Alberta, Canada, are required to report exhaust flow and temperature, in 479 mass emission rates.</li> <li>480 Predictive modeling, including ANN and XGBoost algorithms, has shown great per 481 for emission monitoring. Computer-based emission monitoring methods can be in 482 small stationary combustion sources that do not have regulatory requirements for 483 continuous monitoring. A PEMS provides a more accurate emissions reporting met 484 small sources than engineering estimation methods with generic emissions factor 485 normally used for regulatory reporting. Predictive modeling also offers the potent 486 construct a high-resolution network for continuously monitoring point source emi 487 without high capital cost. However, further research is needed because of the foll</li> </ul>	control
<ol> <li>Deploying predictive models to existing process control networks. Process devices, such as distributed control systems, on process networks may not sufficient computational power to run complicated predictive models.</li> <li>Integrating predictive models into existing data acquisition systems and run systems.</li> <li>Integrating predictive models into existing data acquisition systems and run systems.</li> <li>Multiple models may be needed to replace one CEMS unit. For example, for Alberta, Canada, are required to report exhaust flow and temperature, in mass emission rates.</li> <li>Predictive modeling, including ANN and XGBoost algorithms, has shown great per for emission monitoring. Computer-based emission monitoring methods can be in small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting met small sources than engineering estimation methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emit without high capital cost. However, further research is needed because of the foll</li> </ol>	control
<ul> <li>devices, such as distributed control systems, on process networks may not sufficient computational power to run complicated predictive models.</li> <li>2) Integrating predictive models into existing data acquisition systems and run systems.</li> <li>3) Multiple models may be needed to replace one CEMS unit. For example, for Alberta, Canada, are required to report exhaust flow and temperature, in mass emission rates.</li> <li>Predictive modeling, including ANN and XGBoost algorithms, has shown great per for emission monitoring. Computer-based emission monitoring methods can be in small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting methods small sources than engineering estimation methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent without high capital cost. However, further research is needed because of the follow</li> </ul>	
<ul> <li>474 sufficient computational power to run complicated predictive models.</li> <li>475 2) Integrating predictive models into existing data acquisition systems and response</li> <li>476 systems.</li> <li>477 3) Multiple models may be needed to replace one CEMS unit. For example, for</li> <li>478 Alberta, Canada, are required to report exhaust flow and temperature, in</li> <li>479 mass emission rates.</li> <li>480 Predictive modeling, including ANN and XGBoost algorithms, has shown great per</li> <li>481 for emission monitoring. Computer-based emission monitoring methods can be in</li> <li>482 small stationary combustion sources that do not have regulatory requirements for</li> <li>483 continuous monitoring. A PEMS provides a more accurate emissions reporting methods</li> <li>484 small sources than engineering estimation methods with generic emissions factor</li> <li>485 normally used for regulatory reporting. Predictive modeling also offers the potent</li> <li>486 construct a high-resolution network for continuously monitoring point source emis</li> <li>487 without high capital cost. However, further research is needed because of the follow</li> </ul>	t have
<ul> <li>2) Integrating predictive models into existing data acquisition systems and response systems.</li> <li>3) Multiple models may be needed to replace one CEMS unit. For example, for Alberta, Canada, are required to report exhaust flow and temperature, in mass emission rates.</li> <li>Predictive modeling, including ANN and XGBoost algorithms, has shown great performed for emission monitoring. Computer-based emission monitoring methods can be in small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emission without high capital cost. However, further research is needed because of the follower and the presence in the follower and the presence in the follower and the presence is needed because of the follower and the presence is neaded because of the follower and t</li></ul>	
<ul> <li>476 systems.</li> <li>477 3) Multiple models may be needed to replace one CEMS unit. For example, fa</li> <li>478 Alberta, Canada, are required to report exhaust flow and temperature, in</li> <li>479 mass emission rates.</li> <li>480 Predictive modeling, including ANN and XGBoost algorithms, has shown great per</li> <li>481 for emission monitoring. Computer-based emission monitoring methods can be in</li> <li>482 small stationary combustion sources that do not have regulatory requirements for</li> <li>483 continuous monitoring. A PEMS provides a more accurate emissions reporting met</li> <li>484 small sources than engineering estimation methods with generic emissions factor</li> <li>485 normally used for regulatory reporting. Predictive modeling also offers the potent</li> <li>486 construct a high-resolution network for continuously monitoring point source emi</li> <li>487 without high capital cost. However, further research is needed because of the follow</li> </ul>	porting
<ul> <li>3) Multiple models may be needed to replace one CEMS unit. For example, fa</li> <li>Alberta, Canada, are required to report exhaust flow and temperature, in</li> <li>mass emission rates.</li> <li>Predictive modeling, including ANN and XGBoost algorithms, has shown great per</li> <li>for emission monitoring. Computer-based emission monitoring methods can be in</li> <li>small stationary combustion sources that do not have regulatory requirements for</li> <li>continuous monitoring. A PEMS provides a more accurate emissions reporting met</li> <li>small sources than engineering estimation methods with generic emissions factor</li> <li>normally used for regulatory reporting. Predictive modeling also offers the potent</li> <li>construct a high-resolution network for continuously monitoring point source emi</li> <li>without high capital cost. However, further research is needed because of the following is predictive in the potent is source in the potent is source in the potent is needed because of the following is predictive.</li> </ul>	
<ul> <li>Alberta, Canada, are required to report exhaust flow and temperature, in</li> <li>mass emission rates.</li> <li>Predictive modeling, including ANN and XGBoost algorithms, has shown great per</li> <li>for emission monitoring. Computer-based emission monitoring methods can be in</li> <li>small stationary combustion sources that do not have regulatory requirements for</li> <li>continuous monitoring. A PEMS provides a more accurate emissions reporting met</li> <li>small sources than engineering estimation methods with generic emissions factor</li> <li>normally used for regulatory reporting. Predictive modeling also offers the potent</li> <li>construct a high-resolution network for continuously monitoring point source emi</li> <li>without high capital cost. However, further research is needed because of the following and the potent of the potent potent of the potent potent of the potent of the potent of the potent potent of the potent potent of the potent p</li></ul>	cilities in
479 mass emission rates. 480 Predictive modeling, including ANN and XGBoost algorithms, has shown great per 481 for emission monitoring. Computer-based emission monitoring methods can be in 482 small stationary combustion sources that do not have regulatory requirements for 483 continuous monitoring. A PEMS provides a more accurate emissions reporting met 484 small sources than engineering estimation methods with generic emissions factor 485 normally used for regulatory reporting. Predictive modeling also offers the potent 486 construct a high-resolution network for continuously monitoring point source emi 487 without high capital cost. However, further research is needed because of the follow	addition to
Predictive modeling, including ANN and XGBoost algorithms, has shown great per for emission monitoring. Computer-based emission monitoring methods can be in small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting met small sources than engineering estimation methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emi without high capital cost. However, further research is needed because of the foll	
for emission monitoring. Computer-based emission monitoring methods can be in small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting methods small sources than engineering estimation methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emi without high capital cost. However, further research is needed because of the foll	formance
small stationary combustion sources that do not have regulatory requirements for continuous monitoring. A PEMS provides a more accurate emissions reporting me small sources than engineering estimation methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emi without high capital cost. However, further research is needed because of the foll	stalled on
continuous monitoring. A PEMS provides a more accurate emissions reporting me small sources than engineering estimation methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emi without high capital cost. However, further research is needed because of the foll	r
small sources than engineering estimation methods with generic emissions factor normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emi without high capital cost. However, further research is needed because of the foll	thod for
normally used for regulatory reporting. Predictive modeling also offers the potent construct a high-resolution network for continuously monitoring point source emi without high capital cost. However, further research is needed because of the following th	s that are
<ul><li>486 construct a high-resolution network for continuously monitoring point source emi</li><li>487 without high capital cost. However, further research is needed because of the following the following point source emi</li></ul>	ial to
487 without high capital cost. However, further research is needed because of the following the following the following the second seco	ssions
	owing
88 limitations:	
• Predictive models tend to be equipment-specific. A model trained for one	piece of
equipment may not generate the same performance when applied to anot	her piece
	Page 29 of 34

1

of equipment for air emissions monitoring. Model generalization is needed for wider installation and application.

The performance of predictive models for low and high emission rates (or low and high emission concentrations) is not as good as that for normal emission rates. **495** Accurately predicting high emission rates is critical because facility operators must demonstrate through monitoring methods that their operations meet the maximum **496** <sup>17</sup> **497** emission limits set by the regulator. In this study, the XGBoost model generated <sup>19</sup> 498 lower predictive values than the CEMS measured values for high emission rates.

#### 3.8 Conclusions

<sub>24</sub> 500 Open-source libraries can play a critical role in the wider installation of PEMSs because **501** commercial PEMSs reportedly have initial capital costs similar to CEMSs. Compared with **502** ANN-based predictive models, tree-based machine learning methods, such as XGBoost, are <sup>30</sup> 503 easier to understand, require less effort for data preprocessing, have fewer <sup>32</sup> 504 hyperparameters for model tuning, and need less time for model training.

In this study, we demonstrated that the XGBoost machine learning algorithm can be used to build predictive models for NO<sub>x</sub> emissions monitoring. The NO<sub>x</sub> emission values **507** predicted by the best model structure using random search techniques meet all the 41 508 statistical requirements outlined by US EPA PS16. The Pearson correlation r value between **509** the XGBoost-predicted NO<sub>x</sub> values and the CEMS-measured NO<sub>x</sub> values was 0.98 at 1-<sup>45</sup> 510 minute intervals. The RMSE and MAE using 1-minute interval data in the test dataset were 0.14 and 0.09, respectively. 

In addition, we proposed a data imputation method for sensor failure and tested model sensitivity using this method. Although there are many more complex data **514** imputation methods, such as K-nearest neighbors, K-means, and multiple imputations by **515** chained equations (Schmitt et al., 2015), we demonstrated a simple imputation by <sup>58</sup> 516 interpolating that the time distance is enough to meet regulatory requirements for

Page 30 of 34

predictive emissions modeling. The model inputs are univariate time series data, and linear interpolation is effective for missing data imputation because the precision of the XGBoost model used in this study changed only up to 5.14% using substitute data when two sensors

Further research is still needed to improve model precision for non-normal operating conditions, especially when emission rates (concentrations) are high, because the precision determines if facilities meet the emission limits required by regulations.

The authors thank the Environment and Water group and Jessica Coles from Tetra Tech

Canada Inc. for editing the manuscript.

This study was supported by the Natural Sciences and Engineering Research Council

(NSERC) of Canada (fund number CRDPJ535813-18) and by Mitacs through the Mitacs

Accelerate program (fund number IT18400).

#### **Author Contributions**

MS: Formal analysis; Methodology; Writing - original draft. DK: Project administration;

Funding acquisition; Writing - original draft.

#### **Declaration of Competing Interest**

The authors declare they have no actual or potential competing interests.

ABB S.p.A., 2014. Predictive Emission Monitoring Systems The new approach for monitoring <sup>54</sup> 540 emissions from industry [WWW Document]. URL <sup>55</sup> 541 https://library.e.abb.com/public/4128e88396a14d83b8116a2a71b5b88f/PB\_PEMS-<sup>56</sup> 542 EN.pdf (accessed 8.22.18).

Bivens, T., 2019. Arkansas Electric Cooperative Corporation. Little Rock, Arkansas. Personal communication

Page 31 of 34

1	
2	
3	
<sup>4</sup> 545	Bivens, T., 2017. CEMS Alternatives – A Different Way.
<sup>5</sup> 546	https://www.jasupra.com/legainews/cems-alternatives-a-different-way-tim-
- 547 7 E49	42097/(dccessed 8.23.18). Regingeri T. Smirniotic R.C. 2016 Impact of nitrogen evides on the environment and
8 540	build
<sup>9</sup> 550	13 133–141 https://doi.org/10.1016/j.coche 2016.09.004
$^{10}_{11}$ 551	Botros, K.K., Williams-Gossen, C., Makwana, S., Siarkowski, L., 2011, Predictive Emission
12 <b>552</b>	Monitoring (PEM) Systems Development and Implementation.
13 553	Chen, T., Guestrin, C., 2016. XGBoost: A Scalable Tree Boosting System, in: Proceedings of
14 554	the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data
15 <b>555</b>	Mining - KDD '16. Presented at the the 22nd ACM SIGKDD International Conference,
16 <b>556</b>	ACM Press, San Francisco, California, USA, pp. 785–794.
17 557	https://doi.org/10.1145/2939672.2939785
18 558	Chien, I.W., Chu, H., Hsu, W.C., Tu, Y.Y., Isai, H.S., Chen, K.Y., 2005. A performance
19 <b>559</b>	study of PEMS applied to the Hsinta power station of Talpower. Atmos. Environ. 39,
20 560	223-230.  Inttps://doi.org/10.1010/j.dtmosenv.2004.09.002
<sup>22</sup> 562	Feasibility Study of a Predictive Emissions Monitoring System Applied to Tainower's
<sup>23</sup> 563	Nannu and Hsinta Power Plants, 1. Air Waste Manag, Assoc. 60, 907–913.
<sup>24</sup> 564	https://doi.org/10.3155/1047-3289.60.8.907
<sup>25</sup> 565	Cooper, D.A., Andreasson, K., 1999. Predictive NOx emission monitoring on board a
<sup>26</sup> <sub>27</sub> 566	passenger ferry. Atmos. Environ. 33, 4637–4650. https://doi.org/10.1016/S1352-
<sup>2</sup> / <sub>28</sub> 567	2310(99)00239-3
29 <b>568</b>	Cozza, A., Faulkner, K.F., 1993. Acid rain program offers free-market incentives, portends
<sub>30</sub> 569	future regulation. Hazmat World U. S. 6:5.
31 570	Cuccu, G., Danafar, S., Cudre-Mauroux, P., Gassner, M., Bernero, S., Kryszczuk, K., 2017.
32 571	A data-univen approach to predict NOX-emissions of gas turbines, in: 2017 IEEE International Conference on Big Data (Big Data), Presented at the 2017 IEEE
33 J72 24 573	International Conference on Big Data (Big Data), IFEE Boston, MA, np. 1283–1288
35 574	https://doi.org/10.1109/BigData.2017.8258056
36 575	Eisenmann, T., Bianchin, D.R., Triebel, D., 2014. Predictive Emission Monitoring (PEM):
37 576	Suitability and Application in View of U.S. EPA and European Regulatory Frameworks.
<sup>38</sup> 577	Presented at the The 11th International Conference & Exhibition on Emissions
<sup>39</sup> 578	Monitoring (CEM) 2014, Istanbul, Turkey, p. 15.
<sup>40</sup> 579	Environment and Climate Change Canada, 2019. Canada's Air Pollutant Emissions Inventory
<sup>42</sup> 580	Report: annex 1 [WWW Document]. aem. URL
<sup>43</sup> 581	https://www.canada.ca/en/environment-climate-change/services/air-
44 583	1 4 20)
<sup>45</sup> 584	European Committe for Standardization 2018 Stationary source emissions - Predictive
$\frac{46}{47}$ 585	Emission Monitoring Systems (PEMS) - Applicability, execution and quality
<sup>47</sup> 586	assurance.
49 <b>587</b>	Faravelli, T., Bua, L., Frassoldati, A., Antifora, A., Tognotti, L., Ranzi, E., 2000. A new
<sub>50</sub> 588	procedure for predicting NOx emissions from furnaces. Comput. Aided Chem. Eng.,
51 <b>589</b>	European Symposium on Computer Aided Process Engineering-10 8, 859–864.
52 <b>590</b>	https://doi.org/10.1016/S1570-7946(00)80145-5
53 <b>591</b>	Friedman, J.H., 2001. Greedy Function Approximation: A Gradient Boosting Machine. Ann.
55 502	Stat. 29, 1189-1232. Gualman J. 2012. Cradient boosting trees for puts incurance lass cost modeling and
56 504	nrediction Expert Syst Appl 39 3659-3667
57 595	https://doi.org/10.1016/i.eswa.2011.09.058
58	
59	
6U 61	Page 32 of 34
62	
63	
64	

1	
2	
4 596 5 597	Harnevie, H., Sarkoezi, L., Trenkle, S., 1996. Predictive emission monitoring system (PEMS) for emission control in biomass fired plants (No. SVF575). Stiftelsen foer
<sup>7</sup> 598	Hung, W.S.Y., 1975. An Experimentally Verified NOx Emission Model for Gas Turbine
<sup>°</sup> <sub>9</sub> 600	Combustors V01BT02A009. https://doi.org/10.1115/75-GT-71
$10^{10}$ 601	Hung, W.S.Y., Langenbacher, F., 1995. PEMS: Monitoring NOX Emissions From Gas
<sup>11</sup> 602	Structures and Dynamics: Controls, Diagnostics and Instrumentation: Education:
12 603	IGTI Scholar Award Presented at the ASME 1995 International Gas Turbine and
14 605	Aeroengine Congress and Exposition, ASME, Houston, Texas, USA, p. V005T15A016.
15 606	https://doi.org/10.1115/95-GT-415
16 <b>607</b>	Kamas, J., Keeler, J., 1995. Predictive emissions monitoring systems: a low-cost alternative
17 <b>608</b>	for emissions monitoring [in cement industry], in: 1995 IEEE Cement Industry
18 <b>609</b>	Technical Conference. 37th Conference Record. Presented at the 1995 IEEE Cement
19 610	Industry Technical Conference. 37th Conference Record, IEEE, San Juan, Puerto
20 611 21 612	Rico, pp. 497–509. https://doi.org/10.1109/CITCON.1995.514350
<sup>22</sup> 613	random forests: Statistical arbitrage on the S&P 500 Fur 1 Oper Res 259 689-
<sup>23</sup> 614	702. https://doi.org/10.1016/i.ejor.2016.10.031
<sup>24</sup> 615	Labram, A., 2019. Article: Fitting data with XGBoost   Institute and Faculty of Actuaries
<sup>25</sup> 616	[WWW Document]. URL https://www.actuaries.org.uk/news-and-
$\frac{20}{27}$ 617	insights/news/article-fitting-data-xgboost (accessed 1.23.20).
28 618	Lee, YH., Kim, M., Han, C., 2005. Application of Multivariate Statistical Models to
29 619	Prediction of NUX Emissions from Complex Industrial Heater Systems. J. Environ.
30 020	Mandel 1 S.P. 2015 A Comparison of Six Methods for Missing Data Imputation, 1 Biom
32 622	Biostat. 06. https://doi.org/10.4172/2155-6180.1000224
33 623	Mauzerall, D., Sultan, B., Kim, N., Bradford, D., 2005. NO emissions from large point
34 624	sources: variability in ozone production, resulting health damages and economic
35 <b>625</b>	costs. Atmos. Environ. 39, 2851–2866.
36 626	https://doi.org/10.1016/j.atmosenv.2004.12.041
38 628	Technology Trondheim Norway
<sup>39</sup> 629	Roth, M., Lawrence, P., 2010. A cost-effective alternative to continuous emission monitoring
<sup>40</sup> 630	systems. Environ. Sci. Eng. Mag. 56–57.
<sup>41</sup> 631	Saiepour, M., Schofield, N., Leden, B., Niska, J., Link, N., Unamuno, I., Gomes, J., 2006.
<sup>42</sup> <sub>43</sub> 632	Development and Assessment of Predictive Emission Monitoring Systems (PEMS)
<sup>13</sup> 633	Models in the Steel Industry, in: Iron and Steel Technology Conference Proceedings.
45 634	Presented at the AISTech 2006, Association for Iron & Steel Technology, Cleveland,
$\frac{46}{100}$ 636	Semaniski I Gautama S 2015 Smart City Mobility Application—Gradient Boosting Trees
4 <sup>1</sup> / 637	for Mobility Prediction and Analysis Based on Crowdsourced Data. Sensors 15,
49 638	15974–15987. https://doi.org/10.3390/s150715974
<sub>50</sub> 639	Si, M., Tarnoczi, T.J., Wiens, B.M., Du, K., 2019. Development of Predictive Emissions
51 640	Monitoring System Using Open Source Machine Learning Library – Keras: A Case
52 641	Study on a Cogeneration Unit. IEEE Access 7, 113463-113475.
53 042 54 643	Swanson B 2018 CMC Solutions Pleasant Hill IA Personal communication Nov 2018
55 644	Tan, P., Xia, J., Zhang, C., Fang, O., Chen, G., 2016. Modeling and reduction of NOX
<sup>56</sup> 645	emissions for a 700 MW coal-fired boiler with the advanced machine learning
<sup>57</sup> 646	method. Energy 94, 672–679. https://doi.org/10.1016/j.energy.2015.11.020
<sup>58</sup> 647	US EPA, 2019. What Pollutants are Included in "oxides of nitrogen" in MOVES?   MOVES and
59 <b>648</b> 60	Other Mobile Source Emissions Models   US EPA [WWW Document]. URL
61	Page 33 of 34
62	
63	

1	
2	
3	
<sup>4</sup> 649	https://www.epa.gov/moves/what-pollutants-are-included-oxides-nitrogen-moves
<sup>5</sup> 650	(accessed 12.31.19).
<sup>6</sup> 651	US EPA, 2018. Performance specification 16—specifications and test procedures for
έg652	predictive emission monitoring systems in stationary sources.
<sub>ິ</sub> 653	US EPA, 2009. Performance Specification 16 for Predictive Emissions Monitoring Systems
10 <b>654</b>	and Amendments to Testing and Monitoring Provisions. Fed. Regist., Rules and
11 655	Regulations 74, 12575–12591.
<sub>12</sub> 656	US EPA, 2006. Approval of the Predictive Emission Monitoring System Installed on Unit
13 657	BL2100 at Dearborn Industrial Generation (Facility ID (ORISPL) 55088).
14 658	US EPA, 1999. Nitrogen Oxides (NOX), Why and How They are Controlled.
15 659	US EPA, 1994. An Operator's Guide To Eliminating Bias in CEM Systems (No. EPA/430/R-94-
16 660	U10). Vanderbeggen E. Deneve M. Laget II. Faniel N. Mertene J. 2010. Predictive Emissione.
10 662	Vandernaegen, E., Deneve, M., Lagel, H., Faniel, N., Mertens, J., 2010. Predictive Emissions
19 662	https://doi.org/10.1115/GT2010-22800
20 664	White 1 R 1993 CEMs turn monitoring giant Pollut Eng. U.S. 25:13
21 665	XGBoost developers 2019 XGBoost Parameters — xgboost 1.0.0-SNAPSHOT
<sup>22</sup> 666	documentation [WWW Document] URI
<sup>23</sup> 667	https://xgboost.readthedocs.io/en/latest/parameter.html (accessed 1.24.20).
<sup>24</sup> 668	Zain, S.M., Kien Kek Chua, 2011. Development of a neural network Predictive Emission
<sup>25</sup> 669	Monitoring System for flue gas measurement, in: 2011 IEEE 7th International
<sup>26</sup> 670	Colloquium on Signal Processing and Its Applications. Presented at the 2011 IEEE 7th
28 <sup>/</sup> 671	International Colloquium on Signal Processing and its Applications, pp. 314–317.
<sub>29</sub> 672	https://doi.org/10.1109/CSPA.2011.5759894
<sub>30</sub> 673	Zhang, Y., Haghani, A., 2015. A gradient boosting method to improve travel time prediction.
31 6/4	Iransp. Res. Part C Emerg. Technol. 58, 308–324.
32 0/5	https://doi.org/10.1016/j.trc.2015.02.019
33 0/0	
34 25	
36	
37	
38	
39	
40 41	
±⊥ 42	
43	
44	
45	
46	
47 49	
49	
50	
51	
52	
53	
54 FF	
55 56	
57	
58	
59	
60 C1	Pane 34 of 34
o⊥ 62	
63	
64	
65	

- Gradient Boosting
- Machine learning on environmental monitoring and modeling
- Predictive emissions monitoring system
- PEMS
- XGBoost

**Minxing Si**: Formal analysis; Methodology; Role/Writing - original draft.

**Ke Du**: Project administration; Funding acquisition; Role/Writing - review & editing.

#### **Declaration of interests**

 $\boxtimes$  The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

□ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

