APPLICATION OF ARTIFICIAL INTELLIGENCE TO WASTEWATER TREATMENT PLANT OPERATION

by

Praewa Wongburi

A dissertation submitted in partial fulfillment of the requirements for the degree of

Doctor of Philosophy
(Civil and Environmental Engineering)

2021

at the
UNIVERSITY OF WISCONSIN-MADISON

Date of final oral examination: 07/28/2021

The dissertation is approved by the following members of the Final Oral Committee:

Jae Park. Faculty, Professor, Civil and Environmental Engineering
Greg Harrington. Faculty, Professor, Civil and Environmental Engineering
Chin Wu. Faculty, Professor, Civil and Environmental Engineering
Andrea Hicks. Faculty, Assistant Professor, Civil and Environmental Engineering
Miaoyan Wang. Faculty, Assistant Professor, Statistics

ABSTRACT

In a wastewater treatment plant (WWTP), big data is collected from sensors installed in various unit processes, but limited data is used for operation and regulatory permit requirements. With the advancement in information technology, the data size in wastewater treatment systems has increased significantly. However, WWTPs have not used big data systematically to aid the operation and detect potential operational issues due to the lack of specialized analytical tools.

The objectives of the study were to: (1) develop analytics methods suitable for the management of big data generated in WWTPs, (2) interpret analytics results for extracting meaningful information, (3) implement a recurrent neural network (RNN) and Long Short-Term Memory (LSTM) to predict effluent water quality parameters and Sludge Volume Index (SVI), (4) apply an Explainable Artificial Intelligence (AI) algorithm to determine causes of predicted values, and (5) propose a real-time control using a predictive model to monitor and optimize the operation of WWTPs.

The predictive AI models in WWTPs were developed by applying big data analytics, statistical analysis, and RNN algorithms with an Explainable AI algorithm. The models successfully and accurately predicted the effluent water quality data and a key operational parameter, SVI. Furthermore, the Explainable AI algorithm provided insight into which influent parameters affected higher predicted effluent concentrations and SVI on a specific day, allowing operators to take corrective actions.

From a WWTP's operational data analysis, the RNN model successfully predicted the effluent concentrations of BOD₅, total nitrogen (TN) and total phosphorus (TP), and SVI. Furthermore, the

Explainable AI analysis found that higher influent NH₃N values lead to higher effluent BOD₅, and higher influent total suspended solids (TSS) and TP values resulted in lower effluent BOD₅, implying the importance of controlling dissolved oxygen (DO) in aeration basins. Since aeration is one of the major energy consumption sources in WWTPs, real-time prediction of the effluent water quality using the self-learning AI system developed in this study can be adopted to lower the energy cost significantly while improving effluent water quality. WWTPs must develop control methods based on the RNN prediction and Explainable AI analysis due to different operational conditions.

TABLE OF CONTENTS

1.	INTRODUCTION	1
	1.1 Background and Research Motivation	1
	1.2 Problem Statement	4
	1.2.1 Lack of Research on Big Data Management from WWTPs	4
	1.2.2 Lack of Comprehensive Statistical Analysis of Wastewater Treatment Data	5
	1.2.3 Lack of Predictive Models to Forecast Effluent Quality	5
	1.2.4 No Logistics for a Real-Time Model for WWTPs	6
	1.3 Research Objectives	6
	1.4 Research Scope and Methodology	6
	1.5 Organization of the Dissertation Proposal	8
2.	LITERATURE REVIEW	10
	2.1 Previous Studies of Big Data from WWTPs	10
	2.1.1 Big Data Basics	10
	2.1.2 Type of Big Data Analytics	11
	2.1.3 Data from WWTPs	13
	2.1.4 Big Data for Better WWTP Management.	15
	2.2 Previous Studies on Statistical Analysis in Wastewater Treatment Systems	16
	2.3 Previous Studies of the Development of Predictive Models in WWTPs	19

2.3.1 Autoregressive Integrated Moving Average (ARIMA)	19
2.3.2 Artificial Neural Network (ANN)	23
2.4 Conclusions and Recommendations	27
3. BIG DATA ANALYTICS FROM A WASTEWATER TREATMENT	28
3.1 Abstract	28
3.2 Introduction	29
3.2.1 Background	29
3.2.2 Previous Research	
3.2.3 Shortcoming of Previous Research	31
3.2.4 Study Objectives	32
3.3 Materials and Methods	32
3.3.1 Data Preprocessing	37
3.3.2 Statistical Analysis	39
3.4 Results and Discussions	41
3.4.1 Understanding of Data	41
3.4.2 Data Preparation.	41
3.4.3 Data Preprocessing	52
3.5 Conclusions and Recommendations	67
4. RECURRENT NEURAL NETWORKS (RNN) FOR MODEL PREDICTION.	69
4.1 Abstract	69
4.2 Introduction	70

	4.2.1 Background	70
	4.2.2 Ideal Predictive Model for a WWTP	71
	4.2.3 Shortcomings of Previous Predictive Models	75
	4.2.4 Study Objective	76
	4.3 Materials and Methods	77
	4.3.1 Data Preparation	77
	4.3.2 Development of Recurrent Neural Networks (RNNs) Models	81
	4.4 Results and Discussion	84
	4.5 Conclusions and Recommendations	116
	4.6 Future Research Plans	117
5.	LOGISTICS FOR A REAL-TIME PREDICTION MODEL	118
	5.1 Abstract	118
	5.2 Introduction	119
	5.2.1 Background	119
	5.2.2 SCADA Modernization with Python	120
	5.2.3 Explainable Artificial Intelligence	123
	5.2.4 Literature Review of Aeration Optimization in Wastewater Treatment Process	124
	5.2.5 Study Objectives	125
	5.3 Materials and Methods	125
	5.3.1 Logistics for a Real-Time Model in Wastewater Treatment Plant	125
	5.4 Results and Discussion	133

5.5 Conclusions and Recommendations	145
5.6 Future Research	146
6. PREDICTION OF SLUDGE VOLUME INDEX (SVI) IN A WASTE TREATMENT PLANT USING ARTIFICIAL INTELLIGENCE	
6.1 Abstract	148
6.2 Introduction	148
6.2.1 Background	
6.2.2 Sludge Volume Index	
6.2.3 Activated Sludge Process	151
6.2.4 Filamentous Bulking	151
6.3 Materials and Methods	152
6.4 Results and Discussion	160
6.5 Conclusions and Recommendations	173
7. REFERENCES	175
8. APPENDICES	185

LIST OF FIGURES

Figure 1.1 The growth of data from 2010 to 2020. (Source: Roser et al., 2015)
Figure 1.2 Number of deaths by risk factor. (Source: Ritchie & Roser, 2018)
Figure 1.3 Number of people with and without access to safe drinking water
Figure 1.4 Organization of the preliminary dissertation proposal9
Figure 2.1 The three V principles of big data. (Source: Su, 2018)
Figure 2.2 The iterative steps of the ARIMA approach. (Source: Box and Jenkins, 1970) 20
Figure 2.3 Frequency and trend of AI techniques applied to wastewater treatment during 1995-
2019. (Source: Zhao et al., 2020)
Figure 2.4 Classification tree of AI technology used in wastewater treatment
Figure 2.5 Basic structure of an Artificial Neural Network (ANN). (Source: Haider et al., 2019,
Figure 2.6 Basic structure of a Recurrent Neural Network (RNN). (Source: Haider et al., 2019,
Figure 2.7 Basic structure of a Long Short-Term Memory cell (LSTM).
Figure 3.1 The methodology of big data analytics in WWTPs
Figure 3.2 Nine Springs Wastewater Treatment Plant. (Source: McGowan & Wang, 2008) 34
Figure 3.3 The numbers of historical Nine Springs WWTP data from 1996-2019 34
Figure 3.4 Column name is 'MeasureCode', and 'LocationCode' contains various parameters, 35

Figure 3.5 Forms of data preprocessing. (Source: Han et al., 2012).	38
Figure 3.6 General layout of Nine Springs WWTP. (Source: McGowan & Wang, 2008)	42
Figure 3.7 The schematic of the liquid treatment process at the Nine Springs WWTP	43
Figure 3.8 The relationship between influent and effluent in TSS	50
Figure 3.9 The relationship between influent and effluent in TP.	51
Figure 3.10 The relationship between influent and effluent in TKN	51
Figure 3.11 The relationship between influent and effluent in NH ₃ N	51
Figure 3.12 The relationship between influent and effluent in BOD ₅	52
Figure 3.13 Correlation of effluent BOD5 to other parameters	55
Figure 3.14 Kurtosis and Skewness of normal distribution	58
Figure 3.15 Box plot of yearly and quarterly effluent BOD5	58
Figure 3.16 Normal probability distribution.	59
Figure 3.17 The means effluent BODs over a day, week, month, quarter, and year	60
Figure 3.18 The average means effluent BODs by year, quarter, month, and day	61
Figure 3.19 The patterns effluent BOD ₅ for each year	62
Figure 3.20 Time series plot of effluent BOD ₅ from 2015-2018	63
Figure 3.21 Box plot of effluent BODs by year and quarter from 2015-2018	63
Figure 3.22 The mean effluent BOD5 over a day, week, month, quarter, and year	64
Figure 3.23 The average means of effluent BOD5 over years, quarters, months, and days	65
Figure 3.24 The Dickey–Fuller test with hypothesis testing	60
Figure 4.1 A Recurrent Neural Networks. (Source: Olah, 2015)	72

Figure 4.2 Long Short-Term Memory networks (LSTM). (Source: Olah, 2015)	72
Figure 4.3 Forget gate architecture. (Source: Olah, 2015)	73
Figure 4.4 The cell state for the input gate. (Source: Olah, 2015)	74
Figure 4.5 Update each cell state. (Source: Olah, 2015)	74
Figure 4.6 Output gate architecture. (Source: Olah, 2015)	75
Figure 4.7 The overview of an RNN modeling process.	78
Figure 4.8 The architecture of the RNN model.	79
Figure 4.9 Five steps to develop a simple RNN model	82
Figure 4.10 Five steps to develop an LSTM model.	82
Figure 4.11 Steps to develop a time series model	84
Figure 4.12 Train and test loss over epoch for effluent BOD ₅ prediction models	87
Figure 4.13 The original data of effluent BOD5 from 2015 to 2018	91
Figure 4.14 The prediction of effluent BOD5 from 2015 to 2018 using the simple RNN mo	odel 92
Figure 4.15 The prediction of effluent BODs from 2015 to 2018 using the LSTM model	93
Figure 4.16 The prediction of effluent TP from 2015 to 2018 using the simple RNN model	·l 94
Figure 4.17 The prediction of effluent TP from 2015 to 2018 using the LSTM model	95
Figure 4.18 The prediction of effluent TKN from 2015 to 2018 using the simple RNN mod	del 96
Figure 4.19 The prediction of effluent TKN from 2015 to 2018 using the LSTM model	97
Figure 4.20 The prediction of effluent TSS from 2015 to 2018 using the simple RNN mode	'el 98
Figure 4.21 The prediction of effluent TSS from 2015 to 2018 using the LSTM model	99
Figure 4.22 The prediction of effluent NH ₃ N from 2015 to 2018 using the simple RNN mo	odel. 100

Figure 4.23 The prediction of effluent NH ₃ N from 2015 to 2018 using the LSTM model 101
Figure 4.24 The prediction of daily effluent BOD ₅ from 2015 to 2018 using the Simple RNN model.
Figure 4.25 The prediction of daily effluent BOD5 from 2015 to 2018 using the LSTM model. 103
Figure 4.26 The prediction of daily effluent TP from 2015 to 2018 using the Simple RNN model.
Figure 4.27 The prediction of daily effluent TP from 2015 to 2018 using the LSTM model 105
Figure 4.28 The prediction of daily effluent TKN from 2015 to 2018 using the Simple RNN model.
Figure 4.29 The prediction of daily effluent TKN from 2015 to 2018 using the LSTM model 107
Figure 4.30 The prediction of daily effluent TSS from 2015 to 2018 using the Simple RNN model.
Figure 4.31 The prediction of daily effluent TSS from 2015 to 2018 using the LSTM model 109
Figure 4.32 The prediction of daily effluent NH ₃ N from 2015 to 2018 using the Simple RNN model.
Figure 4.33 The prediction of daily effluent NH ₃ N from 2015 to 2018 using the LSTM model.111
Figure 4.34 The prediction of daily effluent SVI from 2015 to 2018 using the Simple RNN model.
Figure 4.35 The prediction of daily effluent SVI from 2015 to 2018 using the LSTM model 113
Figure 5.1 The diagram that shows a single ordering. (Source: Lundberg & Lee, 2017) 124
Figure 5.2 Three steps of an aeration optimization method
Figure 5.3 The general architecture of the SCADA system. (Source: Sosik, 2014)

Figure 5.4 The main components of the SCADA system. (Modified from Manda et al., 2018)	. 128
Figure 5.5 General WWTP layout. (Source: Richards, 2020)	. 129
Figure 5.6 The flow chart of the real-time logistics system.	. 130
Figure 5.7 The data collection diagram in a wastewater treatment facility.	. 131
Figure 5.8 The automation control system of a WWTP. (Source: Du et al., 2019)	. 132
Figure 5.9 The structure of the modern SCADA. (Modified from Manda et al., 2018)	. 132
Figure 5.10 The modern SCADA with Python.	. 133
Figure 5.11 The SHapley summary plot of the model that includes flow rate	. 135
Figure 5.12 The SHapley summary plot of the model that includes organic loading	. 135
Figure 5.13 The SHapley summary plot of the model that includes flow rate and organic loa	ding.
	. 136
Figure 5.14 The force plot of the first observation including flow rate	. 136
Figure 5.15 The force plot of the first observation including organic loading	. 136
Figure 5.16 The force plot of the first observation including flow rate and organic loading	. 137
Figure 5.17 The collective force plot of all input variables.	. 137
Figure 5.18 The inputs values in the collective force plot with different inputs	. 139
Figure 5.19 The mean of inputs values in the collective force plot by selecting one input v	alue.
	. 140
Figure 5.20 SHAP dependence plot of influent NH ₃ N to influent BOD ₅ .	. 141
Figure 5.21 SHAP dependence plot of influent TKN to influent BOD ₅	. 142
Figure 5.22 The diagram of the real-time logistics for wastewater treatment operation	. 143

Figure 5.23 The real-time logistics for the prediction models in WWTPs	144
Figure 6.1 Sludge Volume Index data from 1996 to 2020	154
Figure 6.2 Box plot of yearly SVI from 1996 to 2020.	154
Figure 6.3 Sludge Volume Index data from 2001 to 2020.	155
Figure 6.4 Box plot of yearly SVI from 2001 to 2020.	155
Figure 6.5 Sludge Volume Index data from 2010 to 2020.	156
Figure 6.6 Box plot of yearly SVI from 2010 to 2020	157
Figure 6.7 The Recurrent Neural Networks model prediction in Python	159
Figure 6.8 Correlation coefficient between each parameter	160
Figure 6.9 Normal distribution of the first dataset from 1996 to 2020	161
Figure 6.10 Normal distribution of the second dataset from 2001 to 2020	162
Figure 6.11 Normal distribution of the third dataset from 2010 to 2020	162
Figure 6.12 Normal probability plot of the first dataset.	163
Figure 6.13 Normal probability plot of the second dataset	164
Figure 6.14 Normal probability plot of the third dataset	164
Figure 6.15 Mean Absolute Error between each epoch of the model	165
Figure 6.16 The Sludge Volume Index prediction model of the first set of data (1996	to 2020).167
Figure 6.17 The Sludge Volume Index prediction model of the second set of data (20	001 to 2020).
	168
Figure 6.18 The Sludge Volume Index prediction model of the third set of data (20)10 to 2020).
	169

Figure 6.19 SHapley interpretation plots for the first model	170
Figure 6.20 SHapley interpretation plots for the second model	171
Figure 6.21 SHapley interpretation plots for the third model	172
Figure 8.1 A simple RNN model in Python.	188
Figure 8.2 A simple RNN model in Python (Continued).	189
Figure 8.3 A simple RNN model in Python (Continued).	190
Figure 8.4 An LSTM model in Python.	191
Figure 8.5 An LSTM model in Python (Continued).	
Figure 8.6 An LSTM model in Python (Continued)	193

LIST OF TABLES

Table 2.1 The definitions of big data related to four themes. (Source: Riahi & Riahi, 2018)	12
Table 2.2 Three stages for gaining confidence in sensors and analyzers. (Source: Shaw, 2017)	7)14
Table 2.3 The classification of concentration in domestic wastewater	16
Table 3.1 The selection of parameters from 'SiteCode' and 'LocationCode'	47
Table 3.2 The influent table from 'MTR VLT' location	48
Table 3.3 The effluent table from 'EFF BLDG' location	49
Table 3.4 Number of parameters in 'MeasureCode'	49
Table 3.5 The clean dataset from the Nine Springs WWTP big data	50
Table 3.6 The table of descriptive statistics	53
Table 3.7 Correlation coefficient between parameters	54
Table 3.8 Heatmap of the correlation coefficient	55
Table 3.9 Program for the result after removing missing values	56
Table 3.10 Program for processing normal distribution test.	57
Table 4.1 Applications of AI models for operation management in WWTPs	77
Table 4.2 The result from data scaling.	80
Table 4.3 The values of water quality parameters in wastewater treatment from 2015 to 2018.	. 85
Table 4.4 Performances of the effluent BOD5 prediction models	86
Table 4.5 The summary of the accuracy of effluent parameters models	89
Table 4.6 General descriptive statistics of the daily data from 2015 to 2018	90

Table 4.7 Comparison of the model accuracy between the discrete big data and the daily dataset.	
	114
Table 4.8 The average R2 score between the discrete big data and the daily dataset	115
Table 5.1 The dataset of the model including output and input parameters.	134
Table 6.1 Nine Springs Wastewater Treatment Dataset	152
Table 6.2 Dataset from 1996 to 2020	153
Table 6.3 Normalization inputs and output values of the models	158
Table 6.4 General statistics for the dataset from 2010 to 2020	160
Table 8.1 Application of AI technologies to pollutant removal in WWTPs	185
Table 8.2 Application of AI technologies to pollutant removal in WWTPs (Continued)	186
Table 8.3 Application of AI models for operation management during wastewater treatme	ent 187

1. INTRODUCTION

1.1 Background and Research Motivation

Due to the demands for lower operation and maintenance cost, energy cost, stringent compliance requirements of water quality parameters, and lower greenhouse gas emission, WWTPs (WWTPs) must be in a smart management mode. Due to the extensive use of water quality sensors, big data is generated every day or even every second. In the wastewater treatment sector and many other fields such as finance, marketing, stocks, health care, and so on, big data plays an essential role. Data generation has been tremendously increasing since 2010, and 90% of the world's data has been created in the past two years (Figure 1.1).

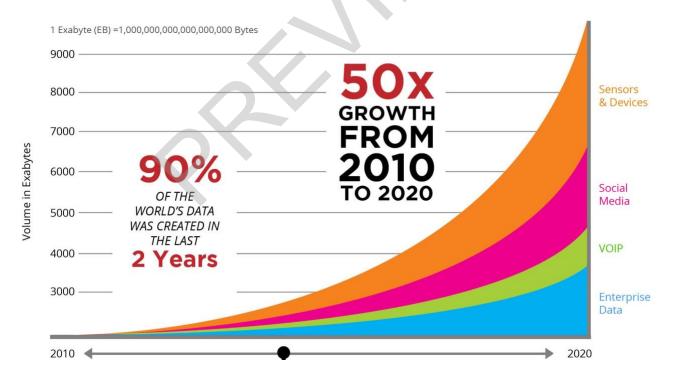


Figure 1.1 The growth of data from 2010 to 2020. (Source: Roser et al., 2015)

Big data is currently most preferably exploited in organizations, companies, and businesses. However, a massive amount of data results in the difficulties of storing, monitoring, analyzing, and visualizing data for further processing. Therefore, data is often stored and then underutilized. A good understanding of the data dynamics in WWTPs is vital for reliable monitoring and control activities. However, the dynamical behavior of the data is usually complicated and uncertain due to nonlinearity, variations from the environmental conditions, strong interactions between the process variables involved, and changes in the flow rate and concentration of the composition of the influent (Harrou et al., 2018). Finding insight from historical and real-time data can directly improve traditional operational systems in a better direction.

A series of wastewater treatment processes remove pollutants from wastewater to be safely reused or discharged into natural water resources. Treated water can be recycled and redistributed for agricultural, industrial, and other purposes or safely released back into the natural resources without causing any adverse effects (Grant et al., 2012). The effluent from a WWTP must meet the National Pollutant Discharge Elimination System (NPEDES) permit to protect the environment and public health (Siegrist, 2017). Lack of access to safe water leads to a risk factor for infectious diseases such as cholera, diarrhea, and dysentery. According to the Global Burden of Disease study, 1.2 people died prematurely in 2017 due to unsafe water (Figure 1.2). This number was three times the number of homicides in 2017 and approximately equal to the amount that died in road accidents globally. Besides, only 71% of the world population has access to safe drinking water, which means that 29% of the world population does not have access to safe water. It equates to 2.1 billion people globally (Figure 1.3).

Number of deaths by risk factor, World, 2017

Total annual number of deaths by risk factor, measured across all age groups and both sexes.



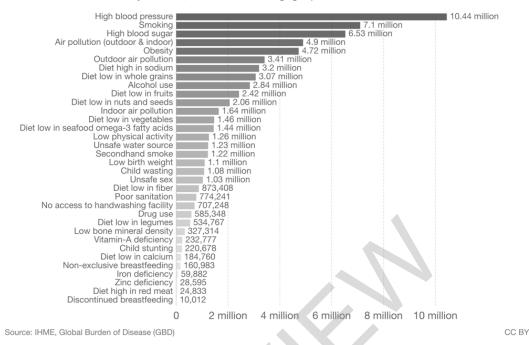


Figure 1.2 Number of deaths by risk factor. (Source: Ritchie & Roser, 2018)

Number of people with and without access to safe drinking water, World ■ With access to safe 7 billion drinking water Without access to safe drinking water 6 billion 5 billion 4 billion 3 billion 2 billion 1 billion 2008 2000 2002 2006 Source: Our World in Data based on WHO, WASH OurWorldInData.org/water-access • CC BY

Figure 1.3 Number of people with and without access to safe drinking water.

(Source: Ritchie & Roser, 2018)

Even though there is an increasing amount of historical data in WWTPs, most information in the data will remain unexploited. According to the operator's point of view, the main reason for this is the high dimensionality of the data, where traditional analysis tools cannot be used (Durrenmatt, 2011). As a result, a large amount of data is lost. Various methods and tools such as big data analytics, statistical analysis, deep learning, and artificial intelligence (AI) applications extract information hidden in the data. Furthermore, this would assist the operator in further optimization of the WWTP, improve the effluent quality, reduce the human errors in operating processes, foster the operator's knowledge of the plant processes, and provide other supporting information. Finally, the main goal of a WWTP is to ensure the efficiency of wastewater treatment operation, which is tremendously essential for community health and the environment. Advanced tools and techniques to improve wastewater treatment operations were introduced in this study for analyzing, modeling, optimizing, and forecasting wastewater treatment quality.

1.2 Problem Statement

A literature review on big data management in wastewater treatment operation indicates four principle wastewater management problems that need to be solved in this study: (1) lack of research on big data management from WWTP, (2) lack of comprehensive statistical analysis of the data, and (3) lack of real-time predictive models to forecast effluent quality.

1.2.1 Lack of Research on Big Data Management from WWTPs

After researching publications in big data management in wastewater treatment facilities from various resources such as Google Scholar®, Scopus®, and Web of Science®, there have been limited studies (Ghernaout et al., 2018; Romero et al., 2017). A large amount of data is generated from various wastewater treatment operations every day. Big data should be exploited to enhance the operating systems. Unfortunately, due to the lack of specialized tools, operators and engineers

cannot extract meaningful and valuable information from the massive amount of highdimensional data.

1.2.2 Lack of Comprehensive Statistical Analysis of Wastewater Treatment Data

Several studies have analyzed energy flow and influent in wastewater treatment facilities to monitor, assess, and model the WWTPs (Martin & Vanrolleghem, 2014). In addition, many publications have illustrated the usefulness of statistical analysis models for WWTP optimization (Cheng et al., 2019; Newhart et al., 2019); operation (Garbowski et al., 2018); analysis (Pantsar-Kallio et al., 1999; Taheriyoun & Moradinejad, 2015; Zhang et al., 2019) and control (Harrou et al., 2018; Maiza et al., 2013). However, few studies have been performed on finding patterns, determining the relationship of each parameter, and selecting meaningful information from big data in wastewater treatment operations with the use of advanced analytics. As a result, big data is underutilized.

1.2.3 Lack of Predictive Models to Forecast Effluent Quality

The water quality predictions in WWTPs were attempted using advanced machine learning tools and techniques. Nonetheless, without the pretreatment of big data, the forecast may not be accurate. Also, previous studies focused on traditional deterministic modeling methodology (Boyd et al., 2019; Huang et al., 2016; Khademikia et al., 2016; Pisa et al., 2019). The novel deep learning method, an RNN, is rarely applied. The conventional approaches might be accurate in predicting if the system is fixed and all the parameters are determined. Furthermore, it is time-consuming, intrusive, and limited by small data analysis.

1.2.4 No Logistics for a Real-Time Model for WWTPs

Real-time data reflects the status of an operational system. The common characteristic of the data is the strict time constraint (Wu et al., 2006). Real-time information is associated with a timestamp and life cycle, and they are only valid for the responding sampling time. Deep learning algorithms with real-time modeling would be a reliable tool for early warning of potential operational upsets and subsequent effluent quality permit violations. There is no study on the development of logistics for real-time modeling using RNNs to help operators for monitoring, detecting fault operation systems in WWTPs.

1.3 Research Objectives

The main objective of this paper is to enhance the operation and performance of WWTPs through big data management and model prediction with AI applications. The study has the following four specific objectives:

- (1) To develop analytics for big data from a WWTP;
- (2) To interpret analytics results for extracting meaningful information;
- (3) To develop deep learning models for forecasting effluent quality from historical wastewater treatment data, which include train and evaluate the models to acquire the most effective algorithms; and
- (4) To establish logistics for a real-time self-learning AI system for monitoring and detecting problems during WWTP operation.

1.4 Research Scope and Methodology

Four principal tasks in this study are discussed below.

Task 1: Big data analytics frameworks in wastewater treatment operation.

This task involves collecting data, understanding the processes in wastewater treatment, visualizing data, selecting meaningful information, and developing big data analytics procedures.

Task 2: Statistical analysis techniques to obtain a pattern and meaningful information.

This task applies various statistical methods such as descriptive statistical analysis, correlation coefficient, box plot, normal distribution, and hypothesis testing to find a relationship between parameters, a pattern of data, and insight of information.

Task 3: DNNs model development for prediction.

This task is to develop a predictive model by implementing advanced modeling techniques, Recurrent Neural Networks (RNNs). The methods include preparing the data, selecting the train and test dataset, build a simple RNNs model and the Long Short-Term Memory (LSTM) model with a different number of hidden layers, train the models to predict an output result, and compare and evaluate the models. The parameters to be predicted will be Total Phosphorus (TP), Total Suspended Solids (TSS), and NH₃/Total Nitrogen (TN) in addition to BOD₅ (Biochemical Oxygen Demand) and lastly, Sludge Volume Index (SVI). Control of dissolved oxygen (DO), sludge and energy management logistics may be attempted from the big data.

<u>Task 4: Propose logistics for a real-time prediction model to detect fault errors in wastewater treatment operations.</u>

Unless a WWTP allows access to real-time big data and incorporates the AI model into their Supervisory Control and Data Acquisition (SCADA) system, it is impossible to implement a real-time prediction model to the existing system. Therefore, this task proposes logistics for

real-time model development in WWTPs to help a WWTP's operator for failure detections beforehand.

1.5 Organization of the Dissertation Proposal

The thesis proposal is organized as follows:

Chapter 2 presents a comprehensive review of previous studies on big data in the wastewater treatment sector, statistical analysis in wastewater treatment systems, and the development of prediction models in WWTPs.

Chapter 3 presents big data management with statistical analysis. Big data analytics processes include data collection, data understanding, data preparation, data mining, evaluation, and deployment. Data was collected from the Nine Springs WWTP in Madison Metropolitan Sewerage District (MMSD), Wisconsin. This data is comprehensively studied through data preprocessing techniques contain data cleaning, data integration, data transformation, and data reduction. Finally, statistical analysis techniques extracted meaningful patterns and information to obtain the appropriate dataset.

Chapter 4 presents the development of model predictions using RNNs. After the big data is extracted to a manageable size and the statistical preprocessing technique is implemented to select meaningful information, the prediction efficiency of wastewater effluent quality will improve significantly. The results of traditional RNN models and RNN-LSTM models were compared and evaluated to choose an optimization algorithm for implementation.

Chapter 5 proposes the subsequent work plan to develop an RNN model using different parameters, develop logistics for monitoring, detecting, and proactive maintenance, assist better decision making, and optimize wastewater treatment facilities.