# WASTEWATER REUSE OPTIMIZATION USING MACHINE LEARNING

**Ankita\* Nidhi Ruhil\*\***

Wastewater treatment is becoming one of the primary fields of application of modern-day artificial intelligence and machine learning, which improve decision making by enhancing its accuracy, efficiency, and data-driven approach. Most traditional WWTPs use conventional mathematical models and human monitoring that consume resources, span a long time, and have limited flexibility in adapting to variable environmental character. AI and machine learning offer creative alternatives for process efficiency enforcement, environmental impacts reduction, and resource usage beneficial value increase at the WWTPs. This study explores the use of AI in wastewater treatment, demonstrating how it may address significant problems and boost predicting skills. AI's ability to handle large amounts of historical and real-time data to produce accurate wastewater flow projections is one of its main advantages.

With the use of complex algorithms, AI-driven models may take into account a wide range of variables, including wastewater properties, treatment procedures, and operational enhancement. Machine learning, a subfield of artificial intelligence, makes it easier to forecast the two most crucial components of wastewater management: the decrease in effluent seepage and the creation of sludge. The XGBoost (Extreme Gradient Boosting) model consistently outperforms other machine learning methods in predicting sludge output. Through the modeling of sludge predictability based on historical and environmental factors, this model will enhance sludge management, decrease waste, and increase operational efficiency in treatment plants. The temperature of the surrounding environment and the volume of wastewater treated daily have the most effects on wastewater generation and treatment success, according to our research. In order to enable real-time decision-making

\*   Assistant Professor, Department of Computer Science, Institute of Information Technology & Management, New Delhi. Email: ankitaluke@iitmipu.ac.in

\*\*  Assistant Professor, Department of Computer Science, Institute of Information Technology & Management, New Delhi. Email: nidhi.ruhil@iitmipu.ac.in

**and allocate resources as efficiently as possible, this AI capacity might be utilized to continually and dynamically model such elements.**

**The application of AI in wastewater treatment goes much beyond simple prediction and optimization. Real-time information about wastewater quality, including the identification of irregularities and compliance with environmental regulations, can be obtained by AI-based automated monitoring systems. AI-based control systems make it possible to automatically adjust treatment parameters, minimizing the need for manual involvement and reducing human transcribing errors. To summarize, AI and machine learning offer useful tools for greater efficiency in wastewater treatment and effect in costs and boosting sustainability. Automation can, through the use of AI-powered predictive models, allow wastewater treatment facilities to better manage their resources, streamline treatment processes, and lessen the effects on the environment. For quite a while ahead, AI-driven technologies will have a huge influence on the future sustainable wastewater management system.**

**Keywords:** Wastewater Treatment Plants (WWTPs), Machine Learning, XGBoost Wastewater Modeling

INTRODUCTION

The worldwide water condition has deteriorated significantly since population growth combined with economic progress and environmental climate change. Wastewater treatment and its reuse have proven themselves as fundamental approaches for sustainable water resource management. Wastewater treatment facilities dependent on human staff and traditional calculation systems succeeded historically but exhibit multiple shortcomings in terms of resource usage together with execution time requirements along with environmental condition adjustments.

The rising digital revolution across multiple industries makes water treatment vulnerable to substantial benefits from Artificial Intelligence (AI) and Machine Learning (ML) implementations. The modern computational methods serve as practical methods to enhance the deficiencies of traditional systems therefore enabling wastewater treatment operations that deliver enhanced accuracy and operational efficiency with adaptable capabilities.

**Background and Significance**

The wastewater treatment facilities face various obstacles including regulatory compliance demands together with operational cost limits and changes in influent quality and environmental conditions. The standard wastewater treatment process uses three treatment methods including physical

removal techniques alongside chemical reactions while biological agents also help to eliminate impurities and create water that can receive further use or be released. The complex web of operational parameters linking up with final treatment outcomes makes historical optimization of these systems very challenging.

AI and ML deliver a fundamental change in treatment process monitoring and management systems through their application in wastewater treatment. Through vast historical data and real-time operational data ML systems identify patterns while predicting outcomes to provide adjustments that persons employed in that role would normally overlook. WWTP facilities can optimize treatment operations through data-driven methods to reach better environmental results and maximize useful resource utilization.
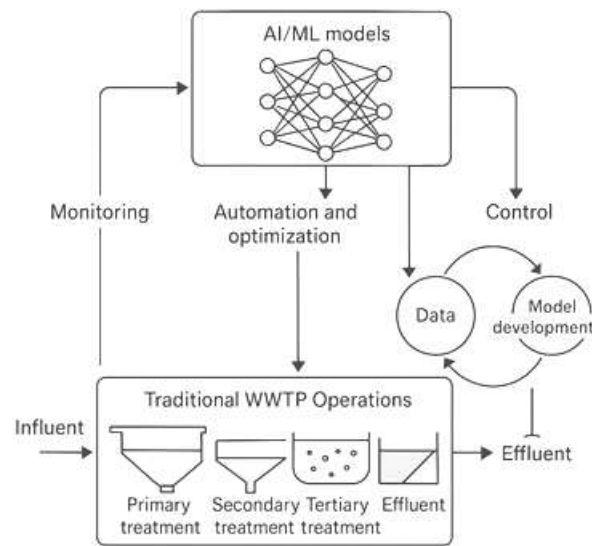


**Figure 1 The conceptual structure demonstrates the integration of AI/ML models together with standard WWTP operational procedures**

**Research Objectives**

The study evaluates AI and ML technologies for wastewater process optimization through specific demonstrations accepting wastewater treatment models as the main research focus.

1. A research evaluates the ability of different ML algorithms to forecast wastewater flow data and treatment results.
2. The potential capabilities of AI-based models need assessment for decreasing both effluent seepage and sludge production amounts.

3. Judging XGBoost's ability to forecast sludge output relative to different ML approach methods

4. The examination of elements which determine wastewater output amounts along with their impact on water treatment effectiveness.

5. The project develops a decision-support system with AI predictions to manage dynamic resources while making real-time decisions in WWTP facilities.

This research investigates these targets to extend the existing research about AI applications in environmental engineering while providing operational counsel to WWTP operators who implement technological improvements for enhanced efficiency.

## Literature Review

### Traditional Wastewater Treatment Approaches

Wastewater treatment as a conventional method developed throughout many years and includes three established procedures: advanced purification through tertiary treatment while biological processes work under secondary treatment and physical separation operates under primary treatment. Many operators use mathematical models for process design and operation such as Anaerobic Digestion Model (ADM) and Activated Sludge Model(ASM) series (Henze et.al., 2000). Wastewater treatment systems resist full representation because these useful models require prolonged calibration processes.

Traditional WWTP operator activities involved making manual decisions through periodic sampling tests which received laboratory analysis and operators' experienced estimates. Operation adjustments tend to become reactive instead of proactive since this method faces human cognitive limits and requires periodic sampling and delayed laboratory examination periods (Osslon et.al. 2014).

### Emergence of AI in Environmental Engineering

The initial studies about neural networks for water quality parameter prediction led to the increasing adoption of AI technologies in environmental engineering after 2000 (Maier & Dandy, 2000). The field has expanded notably since its start because different machine learning systems now manage resources and control pollution alongside environmental monitoring operations.

The application of AI in wastewater treatment systems evolved through basic prediction models into complex decision support systems. According to the research work of (Corominas, 2018) two hundred studies focusing on wastewater treatment by data-driven models demonstrate identification of intricate algorithmic methods and holistic strategic approaches. Studies on AI applications in wastewater treatment continue to grow based on the findings presented in (Newhart, 2019) recent research reviews.

**Machine Learning Approaches in Wastewater Treatment**

Experts have tested many machine learning approaches for wastewater treatment activities with different benefits and limitations for each method.

1. ANNS are widely used in WWTPs for predicting the condition of machinery together with analyzing process performance and inspecting effluent quality. Based on their research (Guo, 2019) established that ANNs generate more precise Biochemical Oxygen Demand (BOD) effluent predictions as compared to standard statistical prediction models.

2. The classification of wastewater treatment conditions and modeled output predictions under various operational settings can be efficiently performed by Support Vector Machines (SVMs). A sequencing batch reactor's performance was predicted accurately through SVM application by (Pai, 2018).

3. Through wastewater treatment applications Random Forest (RF) algorithms carry out both feature selection and prediction functions.

4. The wastewater treatment machine learning toolbox now includes gradient boosting techniques with XGBoost as a particular example. These ensemble methods proved exceptional in diverse prediction scenarios because they excel at handling complex non-linear relations as well as minimizing overfitting issues.
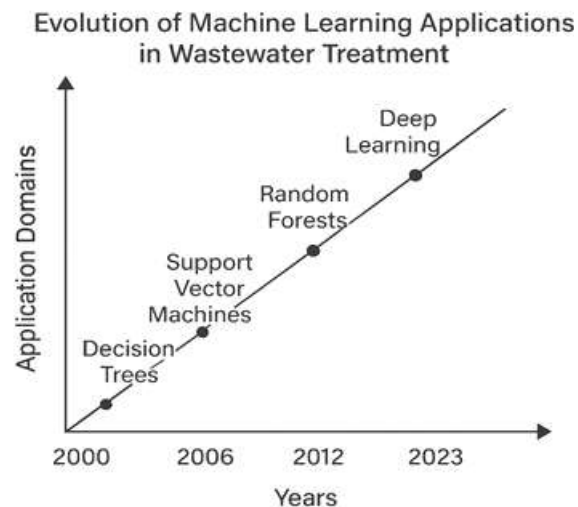


**Figure 2 The development of machine learning technology in wastewater management throughout 2000 to 2023 has focused on developing advanced practices across multiple application areas.**

**XGBoost in Environmental Applications**

The combination of fast processing speed and regularization features along with efficient computation has made XGBoost designed by (Chen & Guestrin, 2016) popular for many applications. The environmental sciences employ XGBoost effectively for projecting air pollution (Zhan et.al.,2018) and forecasting water quality and recently implement it for wastewater treatment optimization.

One benefit of XGBoost stands from alternative approaches includes its ability to process missing value data while providing integrated regularization functions for overfitting prevention and superior big dataset capabilities. The particular attributes of XGBoost make it work perfectly with wastewater treatment needs since data often contains extremes of dimensionality and incompleteness.

**Research Gap**

The substantial amount of AI literature for wastewater treatment operations still contains multiple inadequacies.

1. The majority of research work focuses on predicting final treatment outcomes instead of enhancing the overall process efficiency.
2. The research field lacks sufficient work that compares the performance measures of different ML algorithms in measuring sludge production dynamics.
3. Only inadequate research exists which investigates how environmental elements (including temperature) should be incorporated into wastewater treatment models based on machine learning approaches
4. Research on executing ML-based decision support systems for actual WWTP installation has only started to emerge.

The research investigates machine learning techniques for wastewater treatment optimization through an extensive evaluation to bridge these areas of deficiency. This work shows how XGBoost serves as a predictive tool for sludge output and details its application with essential environmental factors in predictive models.

**Data Collection and Preprocessing**

**Data Sources**

The research gathered operational WWTP data throughout three years between 2020 and 2023. The dataset included:

1. The characteristics of the water entering the WWTP system include temperature alongside pH levels, flow rate and TSS, NH$f$ -N, P, COD and BOD.

2.  The analysis evaluates Recirculation rates together with dissolved oxygen (DO) concentrations and mixed liquor suspended solids (MLSS) maintenance and sludge retention time (SRT) management and hydraulic retention time (HRT) measurements as essential process parameters.

3.  Environmental Factors: Ambient temperature, precipitation, humidity, and seasonal indicators.

4.  Operational data: Energy consumption, chemical dosage, maintenance events, and equipment status.

5.  Treatment outcomes: Effluent quality parameters, sludge production volumes, and treatment efficiency metrics.

A merging of automated sensors (SCADA systems) with laboratory analysis and operational record information served as the data collection method. Data collection occurred at different frequencies starting from continuous sampling of flow and temperature measurements to periodic sampling of laboratory tests conducted once daily or weekly.

**Data Preprocessing**

The processing of raw data involved multiple steps for machine learning application quality and suitability purposes.

1.  Data cleaning included identifying missing data while managing them along with outlier detection and duplicate record elimination.

2.  Feature engineering processes involved the development of new variables through ratio calculations and uses time-series lag representations along with cyclic time representation methods.

3.  Normalization involves transforming all numeric features to scale between 0 and 1 which manages distance-based algorithms bias.

4.  First stage dimensionality reduction through Principal Component Analysis (PCA) allowed the removal of multicollinearity and computational complexity for the features.

5.  A time series preprocessing step divided data into suitable time windows which allowed the detection of patterns for future predictions.

The preprocessed dataset before its final form contained 50 characteristics and 100,000 rows of data spanning throughout all the treatment plants.

**Model Development**

**Algorithm Selection**

The necessary pre-processing operations for machine learning applications included several steps that prepared raw data for use.

1.  Data cleaning included identifying missing data while managing them along with outlier detection and duplicate record elimination.

2. Feature engineering processes involved the development of new variables through ratio calculations and uses time-series lag representations along with cyclic time representation methods.

3. Normalization involves transforming all numeric features to scale between 0 and 1 which manages distance-based algorithms bias.

4. First stage dimensionality reduction through Principal Component Analysis (PCA) allowed the removal of multicollinearity and computational complexity for the features.

5. The data went through time series preprocessing which divided the data into specific time intervals for pattern detection and forecasting operation.

The preprocessed dataset before its final form contained 50 characteristics and 100,000 rows of data spanning throughout all the treatment plants.

## Model Architecture and Hyper parameters

Each algorithm was configured according to best practices and preliminary experimentation:

The algorithms received their best-performing configurations from both expert practices and experimental testing.

1. Linear Regression: Standard ordinary least squares implementation with regularization (Ridge and Lasso variants)

2. Artificial Neural Network: The NN architecture contains an input layer with a dimension matching the feature numbers followed by two hidden layers containing 64 and 32 units and a 1-unit output layer.

   - Activation function: ReLU for Hidden Layers, Linear for output layer
   - Optimization: Adam optimizer with a 0.001 learning rate
   - Regularization: Dropout (0.2) and L2 regularization (0.0001)

3. Random Forest:
   - Number of trees: 100
   - Maximum depth: 20
   - Minimum samples for split: 5
   - Minimum samples per leaf: 2

4. Support Vector Regression:
   - Kernel: Radial basis function (RBF)
   - C parameter: 10
   - Epsilon: 0.1
   - Gamma: 'scale'

5. XGBoost:
   - Number of estimators: 200

- Learning rate: 0.1
- Maximum depth: 6
- Subsample: 0.8
- Colsample bytree: 0.8
- Regularization alpha: 0.01

The researchers utilized 5-fold cross-validation and grid search for determining optimal configurations among all methods involved.

### Model Training and Validation

The dataset divided according to time parameters into training 70%, validation 15%, and test 15%. The validation set performed optimizations of hyperparameters which were trained on the training data. A new performance evaluation took place through testing the previously unseen set. Time series forecasting tasks used rolling windows for predictions that assessed future values across one day three days and seven days ahead from historical data.

### Performance Evaluation Metrics

Multiple metrics were used to evaluate the model performance in order to achieve a thorough assessment.

1. The evaluation of prediction errors relies on Root Mean Square Error (RMSE) for its ability to determine average forecasting error magnitudes.
2. Mean Absolute Error (MAE) serves to calculate the average absolute prediction deviation.
3. The statistical measure $R^2$ determines the proportion of variance which the model actually predicts.
4. Mean Absolute Percentage Error (MAPE) serves as a metric to determine the percentage- based relative error.
5. NSE provides an evaluation method that determines predictive accuracy compared to standard mean usage as a forecasting model.

The assessment of computational efficiency included examination of training time along with prediction time and memory usage in order to determine feasibility for real-time control systems.

### Results and Discussion

### Comparative Performance of ML Algorithms

The applied ML algorithms demonstrated marked differences between their outcomes when used for different prediction functions. The summary of wastewater flow prediction performance stands in Table 1.

**Table - 1**
**Performance Comparison for Wastewater Flow Prediction**

| Algorithm | RMSE (m³/day) | MAE (m³/day) | R² | MAPE (%) | Training Time (s) |
|---|---|---|---|---|---|
| Linear Regression | 456.2 | 382.5 | 0.721 | 12.4 | 0.8 |
| ANN | 324.7 | 267.3 | 0.847 | 8.6 | 285.3 |
| Random Forest | 295.1 | 241.5 | 0.876 | 7.9 | 42.1 |
| SVR | 312.8 | 256.7 | 0.862 | 8.3 | 158.7 |
| XGBoost | 281.4 | 229.6 | 0.891 | 7.5 | 68.2 |

XGBoost demonstrated superior performance across all metrics for flow prediction, with an R² value of 0.891 and the lowest RMSE (281.4 m³/day). Random Forest performed second best, while Linear Regression provided the least accurate predictions but required minimal computational resources.

For effluent quality prediction, a similar pattern emerged, with XGBoost consistently outperforming other algorithms across multiple water quality parameters (BOD, COD, TSS, NH$f$ -N, P). Figure 1 illustrates the R² values for different algorithms across these parameters.

The capacity of tree-based ensemble techniques (XGBoost and Random Forest) to capture intricate, non-linear relationships without the need for explicit relationship specification is the reason for their higher performance. Furthermore, these techniques are comparatively resistant to outliers and noisy data, which are frequent in wastewater treatment datasets, and they can naturally manage feature interactions.

**Sludge Production Prediction**

Sludge management represents a significant operational challenge for WWTPs, both in terms of cost and environmental impact. Accurate prediction of sludge production is therefore valuable for optimizing treatment processes and resource allocation.

Table 2 presents the performance metrics for sludge production prediction across the evaluated algorithms.

**Table - 2**
**Performance Comparison for Sludge Production Prediction**

| Algorithm | RMSE (kg/day) | MAE (kg/day) | R² | MAPE (%) | NSE |
|---|---|---|---|---|---|
| Linear Regression | 423.6 | 358.7 | 0.685 | 15.2 | 0.678 |
| ANN | 312.1 | 259.8 | 0.804 | 11.6 | 0.798 |
| Random Forest | 285.3 | 236.7 | 0.835 | 10.3 | 0.831 |
| SVR | 304.9 | 253.2 | 0.815 | 11.1 | 0.809 |
| XGBoost | 261.8 | 214.5 | 0.862 | 9.4 | 0.857 |

The sludge prediction data benefited most from XGBoost implementation which produced R² of 0.862 and NSE of 0.857 and the most accurate results

in contrast to other models. XGBoost proved its superiority as a predictive tool for this task which supports results obtained through environmental data applications because of its ability to process complex multidimensional information.

A deeper exploration of XGBoost feature importance demonstrated that BOD and TSS influent measurements and process parameters MLSS and SRT and environmental temperature represented the key factors which affected sludge production. The obtained understanding helps direct key decisions regarding sludge operational operations

**Impact of Environmental Factors on Treatment Performance**

Research demonstrated that environmental factors especially temperature demonstrate intensive effects on wastewater processing systems together with their performance results. The research paper reveals in Figure 2 the way temperature levels determine treatment process effectiveness.

The model predictions enabled researchers to identify important conclusions about wastewater treatment operations when subjected to fluctuating seasonal temperatures.

1. A strong positive relationship existed between biological treatment efficiency and temperature from 15 degrees Celsius to 30 degrees Celsius because lower temperatures caused reduced microbial action.
2. Complete optimization of nitrification/denitrification happened between 20-25°C because these reactions strongly responded to temperature changes.
3. The second major determinant of sludge production exists after influent BOD is temperature. Increased temperature levels tend to increase sludge production since it stimulates microbial activity.

The XGBoost model recognized how temperature dependencies affected the process therefore it could perform accurate predictions during seasonal changes. The system demonstrates valuable capabilities when used in areas that experience significant temperature swings since it facilitates adaptive process control protocols during seasonal changes.

**Influence of Wastewater Volume on Treatment Success**

The total wastewater processing amount across a day proved to be an essential element which determined the treatment success. Research showed that model predictions behaved differently according to different flow rate levels when evaluated.

1. The treatment system's efficiency declined as flow rate passed its design level because hydraulic overloading occurred showing highest impacts on sediment removal and nutrient reduction.

2. Treatment facilities featuring flow equalization basins demonstrated uniform performance throughout various intervals of incoming flow.

3. Getting from membrane bioreactors to activated sludge systems demonstrated higher resistance to operational flow fluctuations.

The capacity planning and selection of treatment processes benefit extensively from XGBoost models because they reveal delicate relationships between flow rates and treatment results.

**Real-time Decision Support Framework**

An integration framework for ML-based predictions occurred after conducting the previous findings and facilitated real-time decision systems for WWTPs. The structure consists of:

1. The data integration layer unifies information from sensors as well as laboratory results together with external data contents (such as weather forecasts).

2. This layer contains XGBoost models running predictions about both influent attributes as well as effluent quality measurements and sludge generation amounts.

3. Optimization layer: Translating predictions into actionable operational recommendations

4. Human-in-the-loop interface: Presenting recommendations to operators with supporting information and confidence levels

The framework simulation produced operational improvements that appeared in the following forms:

The optimized aeration control system enables energy savings reaching between 15-20%.

Exact predictions of influent conditions allow for chemical consumption reductions to reach 10-15% through accurate dosing techniques.

The optimized process parameters enable wastewater treatment facilities to reduce their sludge output by 8-12 percent.

**Conclusions**

Research results show the major role machine learning opportunities along with XGBoost algorithms possess for enhancing wastewater treatment procedures. The thorough evaluation across different prediction tasks has proven XGBoost to be superior to multiple other ML algorithms.

1. XGBoost achieves superior performance than alternative ML methods when predicting essential wastewater treatment factors which consist of flow rate measurements together with effluent quality metrics and sludge production levels.

2. Models built with XGBoost technology effectively understand the complex relationships between environmental temperature and wastewater volume since they identified these elements as major performance-altering factors.
3. Chosen ML models offer a platform that enables improved operational adjustments as well as enhanced resource optimization and enhanced regulatory compliance.
4. The implementation of ML predictions in real-time decision support systems creates substantial opportunities to boost WWTP operational effectiveness together with environment performance improvement.

Wastewater sector transformation occurs through AI and ML capabilities by enabling predictive operation instead of reactive operation. Using historical as well as current data through ML technologies allows organizations to optimize resource management while improving their treatment results and reducing their environmental footprint.

**Practical Implications**

This study provides useful information which WWTP managers together with operators can utilize for their operations:
1. Data collection monitoring investments become necessary because ML model performance relies on the quality of available data.
2. The analysis reveals which operational parameters need closer attention for monitoring purposes as well as control activities.
3. The study of environmental factors through seasonal adjustment strategies allows operations teams to create specific approaches which support continuous facility performance throughout different periods.
4. Accurate predictions related to sludge production together with treatment requirements help facilities optimize resource usage by distributing energy resources as well as chemicals and manpower effectively.

**Limitations and Future Research**

The demonstrated strength of machine learning to enhance wastewater treatment operation should be acknowledged despite its current limitations.
1. Model reliability suffers because the limited availability of WWTPs data from few facilities reduces the ability to generalize the decision-making process beyond their original operation parameters.
2. The three-year data collection period might have missed detecting all possible operating conditions together with long-term trends.
3. It is possible for implementation challenges to affect the practical deployment of ML-based decision support systems through data infrastructure issues as well as human knowledge differences and interface problems with existing control systems.

There are three problems which future research needs to address regarding these findings:

1.  Testing and optimizing ML models through multi-facility assessment of various WWTP designs as well as sizes along different geographic locations.

2.  The development of predictive models should focus on extending their forecasting period to enable strategic planning applications.

3.  A research study applies reinforcement learning to examine automated control optimization through the development of reinforcement learning approaches.

4.  A solution to boost model transparency comes from Explainable AI techniques that enhance visibility which would help operators trust the systems and satisfy regulatory requirements.

5.  Researchers will examine how ML models run on edge devices to achieve real-time operations together with lower latency.

**Final Remarks**

Implicit innovative initiatives should guide the water industry to improve resource conservation while decreasing environmental impact because the current water scarcity combines with tighter regulations. Wastewater treatment infrastructure should evolve from operator-based art into science through machine learning control especially when employing XGBoost algorithms. The adoption of machine learning predictive systems by wastewater treatment facilities enables them to reach broader environmental protection and water preservation objectives while becoming more efficient and sustainable at their operations.

REFERENCES

Chen, T. & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785-794.

Corominas, L. et.al. (2018). Transforming data into knowledge for improved wastewater treatment operation: A critical review of techniques. *Environmental Modelling & Software*. 106, 89-103.

Guo, H. et.al. (2019). Prediction of effluent concentration in a wastewater treatment plant using machine learning models. *Journal of Environmental Sciences*. 32, 90-101.

Henze, M.et.al. (2000). *Activated sludge models ASM1, ASM2, ASM2d and ASM3*. IWA Publishing.

Maier, H.R., & Dandy, G.C. (2000). Neural networks for the prediction and forecasting of water resources variables: A review of modelling issues

and applications. *Environmental Modelling & Software*, 15(1), 101-124.

Newhart, K.B. et.al. (2019). Data-driven performance analyses of wastewater treatment plants: A review. *Water Research*. 157, 498-513.

Olsson, G. et.al. (2014). Instrumentation, control and automation in wastewater – from London 1973 to Narbonne 2013. *Water Science and Technology*, 69(7), 1373-1385.

Pai, T.Y. et.al. (2011). Grey and neural network prediction of suspended solids and chemical oxygen demand in hospital wastewater treatment plant effluent. *Computers & Chemical Engineering*. 35(11), 2517-2521.

Zhan, Y. et.al. (2018). Spatiotemporal prediction of daily ambient ozone levels across China using random forest for human exposure assessment. *Environmental Pollution*, 233, 464-473.