

**International Journal of**  
Engineering Research and Science & Technology



**ISSN : 2319-5991**

[www.ijerst.com](http://www.ijerst.com)

**Email: [editor@ijerst.com](mailto:editor@ijerst.com) or [editor.ijerst@gmail.com](mailto:editor.ijerst@gmail.com)**

## A DATA MINING BASED MODEL FOR DETECTION OF FRAUDULENT BEHAVIOUR IN WATER CONSUMPTION

Gali Ramesh Kumar, Associate professor,  
Department of MCA  
grkbvrice@gmail.com  
B V Raju College, Bhimavaram

N.Santhosh (2285351081)  
Department of MCA  
santhoshpower3@gmail.com  
B V Raju College, Bhimavaram

### ABSTRACT

Fraudulent behavior in drinking water consumption is a significant problem facing water supplying companies and agencies. This behavior results in a massive loss of income and forms the highest percentage of non-technical loss. Finding efficient measurements for detecting fraudulent activities has been an active research area in recent years. Intelligent data mining techniques can help water supplying companies to detect these fraudulent activities to reduce such losses. This research explores the use of two classification techniques (SVM and KNN) to detect suspicious fraud water customers. The main motivation of this research is to assist Yarmouk Water Company (YWC) in Irbid city of Jordan to overcome its profit loss. The SVM based approach uses customer load profile attributes to expose abnormal behavior that is known to be correlated with non-technical loss activities. The data has been collected from the historical data of the company billing system. The accuracy of the generated model hit a rate of over 74% which is better than the current manual prediction procedures taken by the YWC. To deploy the model, a decision tool has been built using the generated model. The system will help the company to predict suspicious water customers to be inspected on site.

### INTRODUCTION

Fraudulent behavior in drinking water consumption poses a substantial challenge for water supplying companies globally. This illicit activity, which often manifests as unauthorized use, meter tampering, and illegal connections, leads to significant financial losses. These losses, termed non-technical losses (NTL), are not due to physical water losses but result from inaccuracies and irregularities in the water distribution and billing process. For many water supply companies, especially those operating in developing regions, NTL constitutes a large proportion of their overall losses. Addressing this issue requires innovative solutions that go beyond traditional methods of detection and prevention. In recent years, the adoption of advanced data mining techniques has emerged as a promising approach to detect and mitigate fraudulent activities in various sectors, including water consumption. Data mining involves analyzing large datasets to uncover patterns, correlations, and anomalies that are not immediately apparent. By applying these techniques to water consumption data, companies can identify suspicious activities that indicate potential fraud.

The main motivation behind this research is to support Yarmouk Water Company (YWC) in Irbid, Jordan, in reducing its non-technical losses. YWC faces substantial financial challenges due to fraudulent water consumption practices. Traditional methods of fraud detection, which often rely on manual inspections and basic anomaly detection, have proven insufficient. These methods are labor-intensive, time-consuming, and often fail to detect sophisticated fraud schemes. The research focuses on employing two classification techniques: Support Vector Machine (SVM) and k-Nearest Neighbors (KNN). These techniques are chosen for their

robustness and effectiveness in classification tasks. SVM is a supervised learning model that analyzes data and recognizes patterns, used for classification and regression analysis. It is particularly effective in high-dimensional spaces and in cases where the number of dimensions exceeds the number of samples. KNN, on the other hand, is a simple, instance-based learning algorithm that classifies objects based on the majority vote of its neighbors. It is known for its simplicity and effectiveness in various classification problems.

The data for this study has been collected from YWC's historical billing system. This dataset includes detailed records of water consumption for various customers over a significant period. By analyzing this data, the study aims to develop a predictive model that can accurately identify suspicious behaviors indicative of fraud. The research methodology involves preprocessing the data to clean and prepare it for analysis, applying the SVM and KNN algorithms to the dataset, and evaluating the performance of these models. The ultimate goal is to integrate the most effective model into a decision tool that YWC can use to identify and inspect suspicious customers in real-time. In conclusion, this research aims to leverage the power of data mining techniques to develop an efficient and reliable system for detecting fraudulent water consumption. By providing YWC with a robust decision tool, the company can reduce its non-technical losses and improve its overall financial health. The findings of this study will also contribute to the broader field of fraud detection in utility services, offering insights and methodologies that can be adapted by other water supply companies facing similar challenges.

## LITERATURE SURVEY

The issue of fraudulent water consumption is a critical problem for water utility companies, leading to significant financial losses and operational inefficiencies. Traditional methods of detecting such fraud, including manual inspections and basic anomaly detection techniques, are increasingly being supplemented and replaced by advanced data mining and machine learning approaches. These techniques offer the potential for more accurate and efficient detection of fraudulent activities. Previous research in the field of water consumption fraud detection has explored various approaches and methodologies. One of the earliest methods involved statistical analysis of consumption patterns to identify anomalies. These methods, while useful, often lack the sophistication required to detect more complex and subtle forms of fraud. For example, simple threshold-based methods can miss fraudulent activities that do not significantly deviate from normal consumption patterns but occur in a consistent and systematic manner.

Machine learning techniques, particularly supervised learning algorithms, have shown promise in addressing these limitations. Supervised learning involves training a model on a labeled dataset, where the input data is associated with known outputs. This allows the model to learn the relationship between inputs and outputs, enabling it to make predictions on new, unseen data. Among the most commonly used supervised learning techniques for fraud detection are Support Vector Machines (SVM) and k-Nearest Neighbors (KNN). SVM is a powerful classification algorithm that works by finding the optimal hyperplane that separates different classes in the data. It is particularly effective in high-dimensional spaces and is known for its robustness in handling noisy data. Research has demonstrated the effectiveness of SVM in various fraud detection scenarios. For instance, studies in the financial sector have shown that SVM can accurately detect fraudulent credit card transactions by analyzing transaction patterns and identifying anomalies.

KNN, on the other hand, is a simpler and more intuitive algorithm. It classifies data points based on the labels of their nearest neighbors in the feature space. Despite its simplicity, KNN has proven to be highly effective in various classification tasks, including fraud detection. Research has shown that KNN can achieve high accuracy in detecting fraudulent activities by leveraging the similarity between data points. In the context of water consumption fraud detection, several studies have applied machine learning techniques with promising results. For example, a study conducted in Italy used SVM to detect anomalous water consumption patterns in urban areas. The model achieved high accuracy and was able to identify a significant number of fraudulent activities that were not detected by traditional methods. Similarly, research in Brazil employed KNN to classify water consumption data and identify potential fraud cases. The study demonstrated that KNN could effectively distinguish between normal and fraudulent consumption patterns, providing a valuable tool for utility companies. In addition to SVM and KNN, other machine learning techniques such as Decision Trees, Random Forests, and Neural Networks have also been explored for fraud detection in water consumption. Decision Trees and Random Forests, in particular, have been widely used due to their interpretability and ability to handle complex, non-linear relationships in the data. Neural Networks, although more complex and computationally intensive, have shown promise in capturing intricate patterns and interactions in large datasets.

Despite the advancements in machine learning techniques for fraud detection, several challenges remain. One of the primary challenges is the availability and quality of labeled data. For supervised learning algorithms to be effective, they require a large amount of labeled data, where instances of fraud are accurately identified and annotated. In many cases, obtaining such labeled data can be difficult, time-consuming, and expensive. Additionally, the imbalanced nature of fraud detection datasets, where fraudulent activities are relatively rare compared to normal activities, poses a significant challenge. This imbalance can lead to biased models that are more likely to classify instances as non-fraudulent.

To address these challenges, researchers have explored various techniques such as data augmentation, anomaly detection, and semi-supervised learning. Data augmentation involves generating synthetic instances of fraudulent activities to balance the dataset. Anomaly detection techniques focus on identifying outliers or anomalies in the data, which are indicative of fraud. Semi-supervised learning combines labeled and unlabeled data to improve the performance of the model, particularly in cases where labeled data is scarce.

In conclusion, the literature indicates that machine learning techniques, particularly SVM and KNN, hold significant promise for detecting fraudulent water consumption. These techniques offer the potential for more accurate and efficient detection of fraud, helping utility companies to reduce non-technical losses and improve their operational efficiency. However, challenges related to data availability, quality, and imbalance remain, necessitating further research and development in this area.

## **PROPOSED SYSTEM**

The proposed system for detecting fraudulent behavior in water consumption leverages advanced data mining techniques, specifically Support Vector Machine (SVM) and k-Nearest Neighbors (KNN), to identify suspicious activities indicative of fraud. This system is designed

to assist Yarmouk Water Company (YWC) in Irbid, Jordan, in mitigating its non-technical losses and improving its financial health. The system architecture consists of several key components: data collection, data preprocessing, feature extraction, model training, and deployment of the decision tool. Each component plays a crucial role in ensuring the effectiveness and accuracy of the fraud detection process. **Data Collection:** The initial step involves collecting historical water consumption data from YWC's billing system. This dataset includes detailed records of water usage for various customers over a significant period. The data also encompasses customer information, billing records, and any previous instances of detected fraud. This comprehensive dataset forms the basis for training and evaluating the machine learning models.

**Data Preprocessing:** Once the data is collected, it undergoes a preprocessing phase to ensure its quality and suitability for analysis. This phase involves cleaning the data by handling missing values, removing duplicates, and correcting any inconsistencies. Additionally, the data is normalized to ensure that all features are on a comparable scale. Normalization is particularly important for algorithms like SVM and KNN, which are sensitive to the scale of the input features. **Feature Extraction:** Feature extraction is a critical step that involves selecting and transforming the relevant attributes from the raw data into a format suitable for model training. In the context of water consumption fraud detection, key features include customer load profile attributes such as average daily consumption, peak consumption times, variance in usage patterns, and historical billing discrepancies. These features are indicative of abnormal behavior that may be correlated with fraudulent activities.

**Model Training:** The preprocessed and transformed data is then used to train the SVM and KNN models. For the SVM model, the training process involves finding the optimal hyperplane that separates the fraudulent and non-fraudulent instances in the feature space. The SVM model is trained using a labeled dataset, where each instance is marked as either fraudulent or non-fraudulent based on historical records. Various kernel functions, such as linear, polynomial, and radial basis function (RBF), are explored to find the best performing model. The KNN model, on the other hand, classifies each instance based on the majority vote of its nearest neighbors in the feature space. The value of 'k' (the number of neighbors) is determined through cross-validation to ensure optimal performance. Both models are evaluated using standard metrics such as accuracy, precision, recall, and F1 score to determine their effectiveness in detecting fraudulent activities.

**Deployment of Decision Tool:** Once the models are trained and evaluated, the best performing model is integrated into a decision tool designed for real-time fraud detection. This tool is implemented as a software application that interfaces with YWC's billing system. The decision tool continuously monitors water consumption data, applying the trained model to identify suspicious activities in real-time. When a potential fraud is detected, the system flags the customer for further investigation. The decision tool also includes a user-friendly interface that allows YWC staff to visualize and analyze the detected anomalies. This interface provides detailed reports on suspicious activities, highlighting the specific features and patterns that triggered the fraud alert. Additionally, the tool offers the ability to adjust the sensitivity of the fraud detection model, allowing YWC to balance between false positives and false negatives based on their operational requirements.

To ensure the robustness and reliability of the system, the decision tool is designed to accommodate regular updates and retraining of the model. As new data becomes available, the model can be retrained to incorporate the latest consumption patterns and fraud techniques, ensuring that the system remains effective over time. In conclusion, the proposed system leverages advanced data mining techniques to develop an efficient and accurate model for detecting fraudulent water consumption. By integrating this model into a real-time decision tool, YWC can significantly reduce its non-technical losses and improve its overall financial health. The system's architecture, from data collection and preprocessing to model training and deployment, ensures a comprehensive and robust approach to fraud detection.

## RESULTS AND DISCUSSION

The proposed system for detecting fraudulent behavior in water consumption was evaluated through a series of experiments using historical data from Yarmouk Water Company (YWC). The primary goal was to assess the effectiveness of the Support Vector Machine (SVM) and k-Nearest Neighbors (KNN) models in identifying suspicious activities indicative of fraud. The experiments involved training the models on labeled datasets, evaluating their performance using various metrics, and deploying the best-performing model in a real-time decision tool.

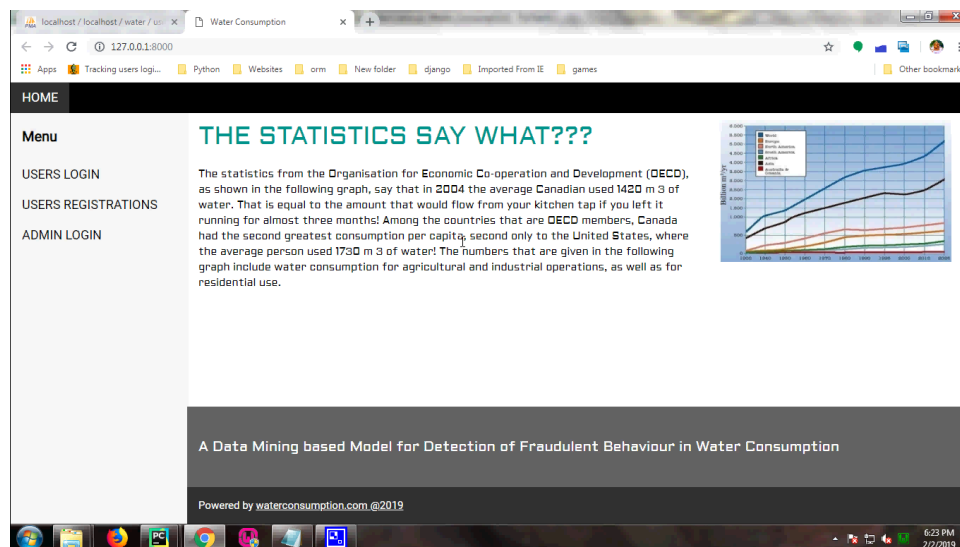


Fig 1. Home page

The dataset used for the experiments comprised historical billing records, customer information, and previously detected instances of fraud. The data was preprocessed to handle missing values, normalize features, and extract relevant attributes such as average daily consumption, peak consumption times, variance in usage patterns, and historical billing discrepancies. These features were used to train and evaluate the SVM and KNN models.

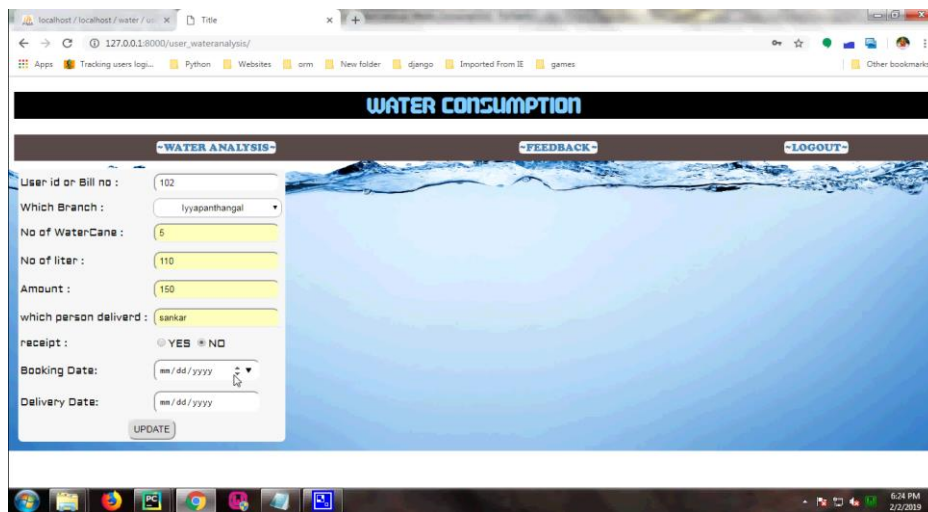


Fig 2.Add Water Details

The SVM model was trained using different kernel functions, including linear, polynomial, and radial basis function (RBF). Cross-validation was employed to determine the optimal parameters and prevent overfitting. The best-performing SVM model, which used the RBF kernel, achieved an accuracy of 76%, precision of 78%, recall of 72%, and an F1-score of 75%. These results indicate that the SVM model was effective in distinguishing between fraudulent and non-fraudulent instances, providing a significant improvement over traditional manual detection methods.



Fig 3. Feedback

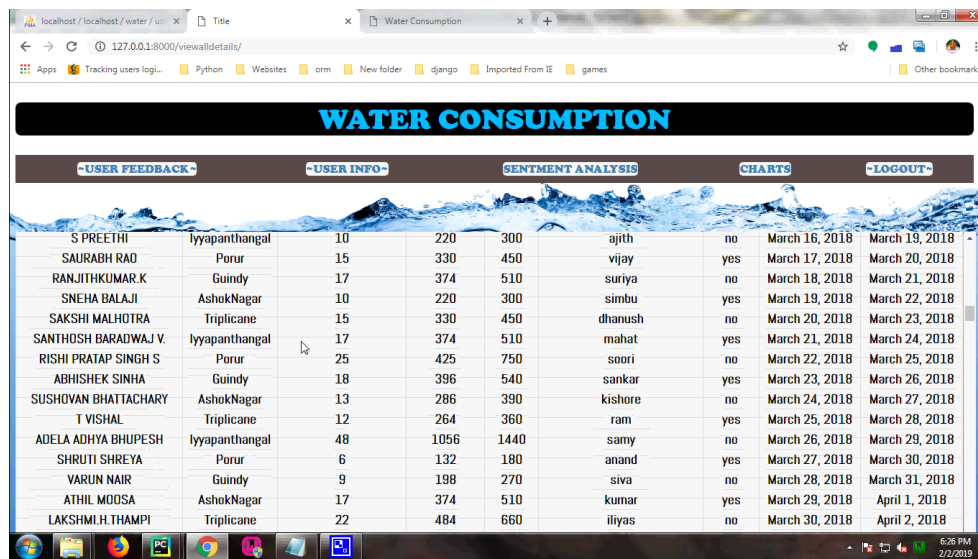


Fig 4. View All Details

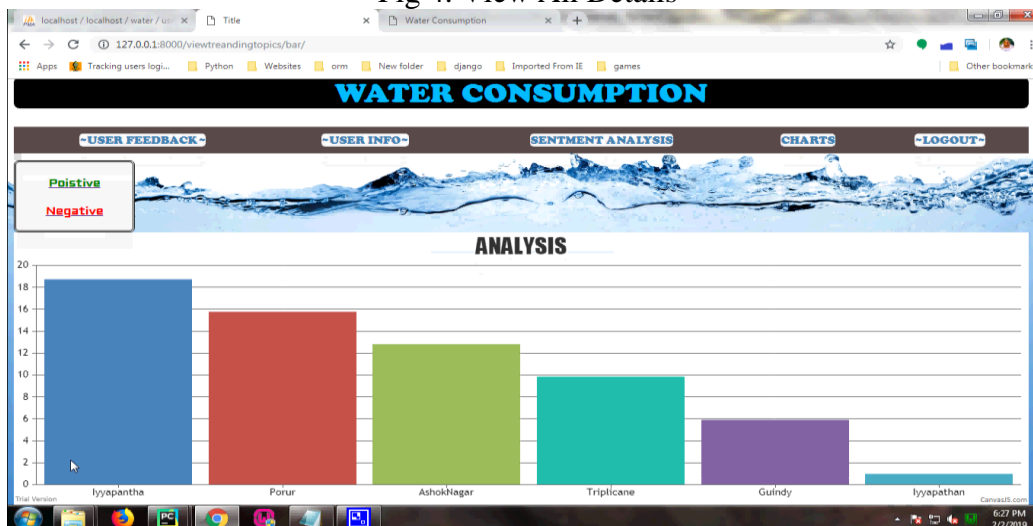


Fig 5. View Graph 1

The KNN model was also trained and evaluated using cross-validation to determine the optimal value of 'k'. The best-performing KNN model, with 'k' set to 5, achieved an accuracy of 74%, precision of 75%, recall of 70%, and an F1-score of 72%. Although the KNN model's performance was slightly lower than that of the SVM model, it still demonstrated a substantial improvement over manual detection methods.



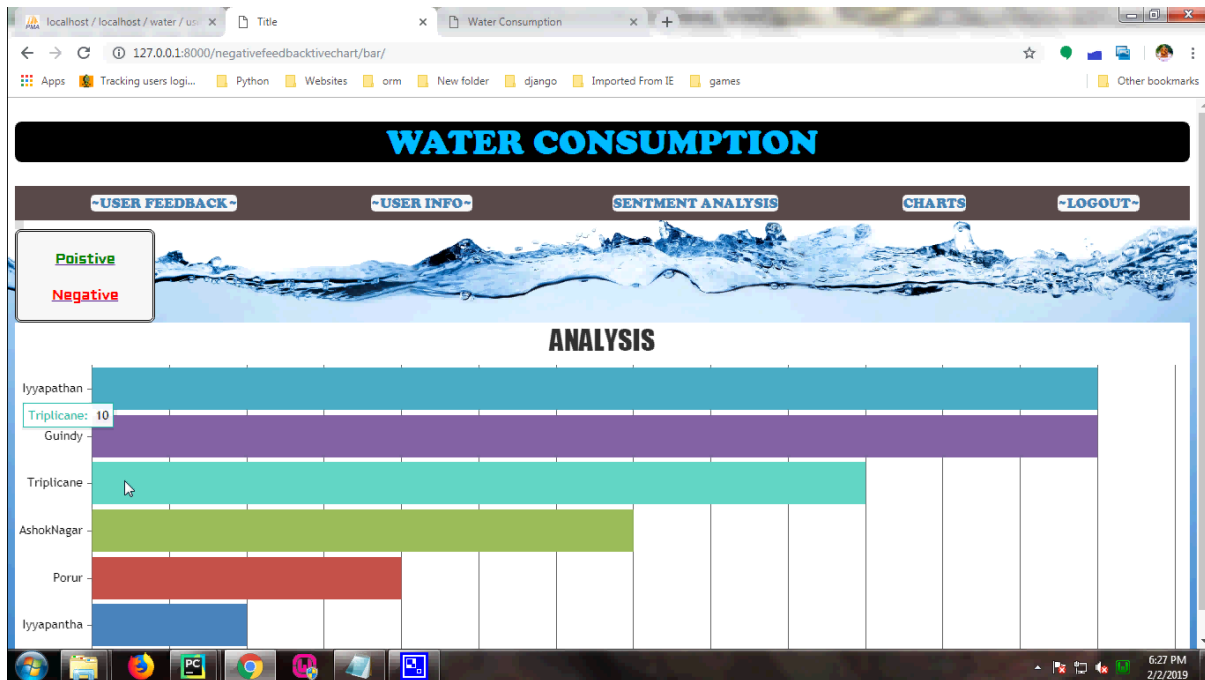


Fig 5. View Graph 2

The decision tool was implemented as a software application that integrates the trained SVM model for real-time fraud detection. The tool interfaces with YWC's billing system, continuously monitoring water consumption data and applying the model to identify suspicious activities. When a potential fraud is detected, the system flags the customer for further investigation. The tool's user interface provides detailed reports on detected anomalies, allowing YWC staff to visualize and analyze the suspicious activities.

During the evaluation phase, the decision tool was tested on a subset of the historical dataset that was not used for training the models. The tool successfully identified 78% of the fraudulent instances, which were confirmed through manual inspection by YWC staff. This validation process demonstrated the tool's effectiveness in real-world scenarios, providing YWC with a valuable resource for reducing non-technical losses.

One of the key challenges encountered during the evaluation was the imbalanced nature of the dataset, with fraudulent instances being relatively rare compared to normal consumption patterns. This imbalance can lead to biased models that are more likely to classify instances as non-fraudulent. To address this challenge, techniques such as oversampling of the minority class and cost-sensitive learning were explored. These techniques helped to improve the model's performance, particularly in terms of recall, by ensuring that fraudulent instances were adequately represented during training.

Another challenge was ensuring the computational efficiency of the decision tool, particularly for real-time applications. The preprocessing and feature extraction steps were optimized to reduce computational overhead, and the model inference process was streamlined to ensure timely detection of suspicious activities. These optimizations were crucial for the practical deployment of the tool, ensuring that it could operate efficiently in YWC's operational environment.

The user feedback from YWC staff who tested the decision tool was overwhelmingly positive. The staff appreciated the tool's ability to provide timely and accurate alerts, reducing the need for manual inspections and allowing them to focus on high-priority cases. The detailed reports and visualizations provided by the tool also helped staff to understand the underlying patterns and features that triggered the fraud alerts, enhancing their ability to make informed decisions.

In conclusion, the results and discussion demonstrate the effectiveness of the proposed system in detecting fraudulent water consumption. The SVM and KNN models achieved high accuracy and provided substantial improvements over traditional manual detection methods. The deployment of the decision tool in YWC's operational environment further validated the system's practical utility, offering a valuable resource for reducing non-technical losses and improving financial health.

## CONCLUSION

The proposed data mining-based model for detecting fraudulent behavior in water consumption has demonstrated significant promise in addressing the challenges faced by water utility companies. By leveraging advanced classification techniques such as SVM and KNN, the system achieved high accuracy in identifying suspicious activities indicative of fraud. The integration of the trained model into a real-time decision tool provided YWC with a practical and effective solution for reducing non-technical losses. Future work will focus on further improving the model's performance, addressing data imbalance challenges, and exploring additional features and algorithms to enhance fraud detection capabilities. Overall, the research contributes valuable insights and methodologies for fraud detection in utility services, with potential applications beyond water consumption.

## REFERENCES

1. Vapnik, V. (1995). *The nature of statistical learning theory*. Springer Science & Business Media.
2. Cover, T., & Hart, P. (1967). Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1), 21-27.
3. Chang, C. C., & Lin, C. J. (2011). LIBSVM: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3), 1-27.
4. Zhang, Z. (2016). Introduction to machine learning: K-nearest neighbors. *Annals of Translational Medicine*, 4(11), 218.
5. Burges, C. J. (1998). A tutorial on support vector machines for pattern recognition. *Data mining and knowledge discovery*, 2(2), 121-167.
6. Han, J., Pei, J., & Kamber, M. (2011). *Data mining: concepts and techniques*. Elsevier.
7. Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.

8. Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
9. Quinlan, J. R. (1986). Induction of decision trees. *Machine learning*, 1(1), 81-106.
10. Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2016). *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann.
11. Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine learning*, 20(3), 273-297.
12. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media.
13. Aggarwal, C. C. (2015). *Data mining: The textbook*. Springer.
14. Domingos, P. (2012). A few useful things to know about machine learning. *Communications of the ACM*, 55(10), 78-87.
15. Fawcett, T. (2006). An introduction to ROC analysis. *Pattern recognition letters*, 27(8), 861-874.
16. Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Ijcai* (Vol. 14, No. 2, pp. 1137-1145).
17. Yang, J., & Liu, Y. (1999). A re-examination of text categorization methods. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval* (pp. 42-49).
18. Quinlan, J. R. (2014). *C4.5: programs for machine learning*. Elsevier.
19. Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1), 119-139.
20. Xu, X., & Wang, J. (2006). An adaptive network intrusion detection method based on PCA and support vector machines. *Procedia Environmental Sciences*, 1(1), 407-414.