

# An Innovative Machine Learning Approach for Early Detection of Thyroid Disorders

**M.V. SHEELA DEVI**

Research scholar (Computer Science)  
School of Engineering and Technology  
Shri Venkateshwara University, Gajraula, UP, INDIA  
Email: [sheela.softinfo@gmail.com](mailto:sheela.softinfo@gmail.com)

**DR PARVEEN KUMAR**

Research Guide  
School of Engineering and Technology  
Shri Venkateshwara University, Gajraula, UP, INDIA  
Email: [pk223475@gmail.com](mailto:pk223475@gmail.com)

**Abstract:** Thyroid diseases are among the most prevalent endocrine disorders worldwide, necessitating accurate and timely diagnosis to ensure effective treatment and management. This study investigates the application of machine learning models for predicting thyroid diseases, utilizing a dataset comprising clinical and demographic attributes. Models such as Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), Naïve Bayes (NB), Support Vector Machine (SVM), and K-Nearest Neighbors (KNN) were evaluated across datasets of varying sizes. A Proposed Method was developed and tested, demonstrating consistent superiority over traditional models with an accuracy of 97% for the largest dataset. The findings highlight the robustness, adaptability, and scalability of machine learning techniques in medical diagnostics, emphasizing the Proposed Method's potential for high-accuracy thyroid disease prediction. This study also reveals the limitations of conventional models, such as Decision Tree and Naïve Bayes, which exhibit reduced accuracy as dataset sizes increase. In contrast, the Proposed Method maintains exceptional performance, underscoring its applicability in both small-scale and large-scale diagnostic settings. By leveraging advanced algorithmic techniques, this research paves the way for integrating machine learning into real-world healthcare systems, enhancing early detection and enabling timely interventions. Future work will explore integrating additional patient data and employing deep learning approaches to further optimize predictive capabilities, aiming to improve diagnostic precision and reduce healthcare costs globally.

**Keywords:** SVM, DT, RF, LR, NB, Machine learning, Thyroid diseases prediction

## 1. INTRODUCTION

Thyroid disease is increasingly becoming a significant health concern, especially among women. In Bangladesh, approximately 50 million people are reported to suffer from thyroid-related issues, with females being at a 10 times higher risk than their male counterparts. Despite this alarming statistic, nearly 30 million people are unaware that they are affected, which delays diagnosis and treatment. According to a study conducted by the Bangladesh Endocrine Society (BES), an estimated 20-30% of women are diagnosed with some form of thyroid disease. Given the widespread nature of the issue and its tendency to go unnoticed, early detection of thyroid dysfunction is essential for timely medical intervention. The thyroid gland, a small butterfly-shaped organ located in the neck, plays a vital role in regulating many of the body's essential functions. It produces hormones such as thyroxine (T4) and triiodothyronine (T3), which are crucial for maintaining metabolism, growth, mood regulation, and energy levels. The thyroid's functionality is primarily controlled by the pituitary gland, which secretes Thyroid Stimulating Hormone (TSH). TSH signals the thyroid gland to produce T3 and T4 hormones. Depending on the levels of these hormones, thyroid dysfunctions can be classified into two major types: hypothyroidism and hyperthyroidism.

Hypothyroidism occurs when the thyroid gland fails to produce adequate amounts of T3 and T4, leading to an increase in TSH levels. This condition results in symptoms such as fatigue, weight gain, brain fog, and depression. On the other hand, hyperthyroidism is caused by excessive production of T3 and T4,

lowering TSH levels and causing symptoms like anxiety, hair loss, rapid thyroid rate, and weight loss. Of these two, hypothyroidism is more prevalent, particularly among females in Bangladesh, and is the primary focus of this research. The challenge lies in detecting hypothyroidism at an early stage, where subtle symptoms may be overlooked or mistaken for other conditions. Traditional diagnostic methods involve physical examination and blood tests to measure TSH, T3, and T4 levels. However, such methods can be time-consuming and may delay diagnosis. This underscores the need for advanced tools that can assist healthcare providers in identifying thyroid disorders early, thus reducing the time to diagnosis and enabling more effective treatment.

Machine learning (ML) offers a promising solution to this problem by enabling early and accurate disease prediction based on patient data. In recent years, ML has become an integral part of healthcare, providing cost-effective and efficient solutions to complex diagnostic challenges. The vast amounts of data collected by healthcare providers can be analyzed using ML techniques to identify patterns, trends, and hidden relationships that may not be apparent through traditional methods. This approach has already shown success in predictive modeling for various diseases, including thyroid disease, diabetes, and cancer, and is now being applied to thyroid disease detection. ML algorithms, particularly classification techniques, are widely used in healthcare for disease diagnosis. These algorithms allow for the categorization of patients based on input data, such as medical history, laboratory results, and demographic information. In this research, we employ several classification algorithms, including Support Vector Machine (SVM), Decision Tree (DT), Random Forest (RF), Logistic Regression (LR), and Naive Bayes (NB), to predict hypothyroidism. Each algorithm offers a unique approach to classifying data, and by comparing their performance, we aim to identify the most effective model for early-stage hypothyroidism detection. Feature selection plays a crucial role in building accurate ML models, particularly in healthcare applications where large datasets are often involved. Selecting the most relevant features—such as specific symptoms, hormone levels, and patient demographics—can significantly improve the performance of the model by reducing overfitting, increasing generalizability, and decreasing computational complexity. In the context of thyroid disease prediction, features such as TSH, T3, T4 levels, and patient age and gender are critical to determining whether a patient is at risk of hypothyroidism.

One of the main challenges in developing predictive models for thyroid disease is the risk of misclassification. Incorrect predictions can lead to unnecessary treatments for healthy individuals or, conversely, delayed treatment for those who need it. To mitigate this risk, we implement robust feature selection methods to ensure that only the most relevant predictors are used in the model. This process not only improves accuracy but also provides insights into the key factors contributing to hypothyroidism, which may assist in clinical decision-making and further research.

## 2. REVIEW OF LITERATURE

In thyroid disease diagnosis, data mining techniques are increasingly leveraged to enhance decision-making and treatment strategies. Thyroid diseases are prevalent, and early identification can significantly improve patient outcomes by allowing for timely interventions. Various classification methods, including Support Vector Machines and K-Nearest Neighbors, are utilized to diagnose thyroid disorders based on clinical parameters. Additionally, feature selection algorithms, such as F-Score and Recursive Feature Elimination, are examined for their ability to refine diagnosis accuracy. Research indicates that these techniques can greatly improve diagnostic performance, achieving accuracies as high as 98.14%. Moreover, the use of genetic algorithms in combination with artificial neural networks and support vector machines further emphasizes the importance of selecting relevant features to enhance the classification of thyroid diseases (Table 1).

In recent years, a growing number of studies have explored the use of ML algorithms for thyroid disease prediction, with varying degrees of success. However, not all studies have incorporated feature selection into their models, which can lead to errors in clinical decision-making. To simplify the feature selection process by identifying the most important variables through techniques such as permutation

importance. This tool enhances the interpretability of ML models and reduces the time and complexity involved in feature selection. In addition to traditional ML techniques, we also examine ensemble learning methods, such as bagging and boosting, to improve model accuracy. Ensemble methods combine the predictions of multiple models to reduce variance and bias, resulting in more reliable outcomes. By integrating ensemble techniques with feature selection, we aim to create a predictive model that not only detects hypothyroidism with high accuracy but also offers greater robustness and scalability.

Table 1: Review of literature for thyroid diseases prediction

Ref. No	Models Used	Feature Selection Method	Accuracy (%)
[7]	Naïve Bayes, Decision Tree (DT), Multilayer Perceptron, RBF Network	No	92.35
[8]	Decision Tree (DT), Artificial Neural Networks (ANN)	No	90.90
[9]	ANN, KNN, SVM, DT	No	89.63
[10]	SVM, Random Forest, Naïve Bayes	Recursive Feature Elimination (RFE)	82.92
[11]	Naïve Bayes	No	89.00
[12]	Random Forest, DT, KNN	No	88.93
[13]	Random Forest, DT, KNN	No	84.80
[14]	DT, Logistic Regression, KNN	No	88.59

### 3. MACHINE LEARNING TECHNIQUES

Traditional classification refers to a category of machine learning algorithms designed for the task of assigning predefined labels or categories to input data based on its features. These algorithms are often employed in supervised learning scenarios, where the model is trained on a labeled dataset to learn the relationship between input features and corresponding output classes. Notable traditional classification algorithms include Decision Trees, k-Nearest Neighbors (KNN), Support Vector Machines (SVM), and Logistic Regression. Decision Trees recursively split the feature space to create a tree-like structure for decision-making. KNN assigns a data point to the majority class among its nearest neighbors. SVM aims to find a hyperplane that best separates different classes. Logistic Regression models the probability of an instance belonging to a particular class [11-22]. These algorithms are widely used in various applications, including image recognition, text classification, and medical diagnosis.

#### 3.1 Decision Trees (DT)

Decision Trees are a popular machine learning algorithm for classification tasks. They recursively partition the feature space based on the values of input features, creating a tree-like structure where each leaf node corresponds to a class label. Decision Trees are interpretable and can handle both categorical and numerical data. However, they are susceptible to overfitting, which has led to the development of ensemble methods like Random Forests, which aggregate multiple decision trees to improve robustness and accuracy [19].

#### 3.2 k-Nearest Neighbors (KNN)

k-Nearest Neighbors is a simple and intuitive classification algorithm that assigns a data point to the majority class among its k nearest neighbors in the feature space. The choice of k influences the algorithm's sensitivity to local variations. While KNN is easy to understand and implement, its performance can be sensitive to the dataset's dimensionality, and it may require proper scaling of features for optimal results [22].

### 3.3 Support Vector Machines (SVM)

Support Vector Machines are powerful classifiers that aim to find the hyperplane that best separates different classes in the feature space while maximizing the margin between them. SVMs are effective in high-dimensional spaces and can handle non-linear decision boundaries through the use of kernel functions. They are robust against overfitting and work well for both binary and multi-class classification tasks. SVMs have applications in various domains, including image classification, text categorization, and bioinformatics [17].

### 3.4 Random Forest

Random Forest is an ensemble learning technique that operates by constructing a multitude of decision trees during training and outputs the mode of the classes (classification) or the mean prediction (regression) of the individual trees. It introduces an element of randomness by considering a random subset of features for each tree, and through a process called bagging (Bootstrap Aggregating), where each tree is trained on a bootstrap sample of the original dataset. This randomness and diversity among the trees enhance the overall model's robustness and reduce overfitting, making Random Forests particularly effective in handling complex datasets. Known for their versatility and high predictive accuracy, Random Forests find application in various domains, including classification, regression, and feature selection tasks [18].

### 3.5 Logistic Regression (LR)

Despite its name, Logistic Regression is a linear model commonly used for binary classification. It models the probability of an instance belonging to a particular class using the logistic function. Logistic Regression is interpretable, computationally efficient, and less prone to overfitting compared to more complex models [20]. It is widely used in applications where understanding the impact of individual features on the predicted outcome is essential, such as in medical diagnostics or credit scoring. Extensions like multinomial logistic regression allow it to handle multiple classes. Table highlighting key differences between Decision Trees (DT), k-Nearest Neighbors (KNN), Support Vector Machines (SVM), and Logistic Regression (LR) (Table 2):

Table 2. Key difference between DT, KNN, SVM and LR

Aspect	Decision Trees (DT)	k-Nearest Neighbors (KNN)	Support Vector Machines (SVM)	Logistic Regression (LR)
Type	Non-linear, non-parametric model.	Non-linear, non-parametric model.	Non-linear, non-parametric model.	Linear, parametric model.
Decision Boundary	Piecewise constant, can be non-linear.	Non-linear and influenced by data structure.	Non-linear, depends on the kernel used.	Linear decision boundary, can be extended to non-linear with feature engineering or kernel tricks.
Interpretability	Highly interpretable, intuitive representation in a tree structure.	Intuitive but less interpretable than decision trees.	Less intuitive due to complexity, but interpretable to some extent.	Coefficients provide interpretable information about the impact of features.
Handling Outliers	Sensitive to outliers.	Sensitive to outliers.	Somewhat robust due to support vectors.	Sensitive to outliers.
Scalability	Fast training but can lead to overfitting.	Slower training, especially with large datasets.	Can be slow with large datasets and high-dimensional features.	Efficient, scales well with large datasets and features.

Hyperparameters	Depth, impurity criteria (e.g., Gini, Entropy).	Number of neighbors (k), distance metric.	Kernel type, regularization parameter (C), kernel parameters.	Regularization strength (C), penalty type, solver type.
Data Types	Handles both numerical and categorical data well.	Sensitive to irrelevant or redundant features, distance metric choice is crucial.	Requires feature scaling, works with numerical data primarily.	Handles numerical input features, may need encoding for categorical variables.
Use Cases	Classification and regression tasks across various domains.	Classification tasks with local patterns.	Classification and regression tasks, especially in high-dimensional spaces.	Binary and multiclass classification, probability estimation.

#### 4. DATASET DESCRIPTION

The dataset used for thyroid disease prediction comprises records with diverse attributes relevant to diagnosing thyroid conditions. It includes demographic information, hormone levels, medication details, and other clinical indicators. These features are either Boolean or continuous-valued, providing a comprehensive view of patient health. The dataset's variability in size, ranging from 200 to 1000 records, enables robust testing of machine learning models. This structured and labeled data facilitates effective training and evaluation of predictive algorithms, ensuring reliable and accurate classification of thyroid conditions. Its rich feature set makes it an ideal resource for developing advanced diagnostic models that cater to diverse medical scenarios (Table 3).

Table 3: Dataset description for thyroid diseases prediction

Attribute	Description	Data Type	Value Range/Example
Age	Age of the patient in years	Continuous	1–120
Sex	Gender of the patient	Categorical	Male, Female
On Thyroxine	Whether the patient is on thyroxine medication	Boolean	Yes, No
Query on Thyroxine	Whether the patient is being queried for thyroxine treatment	Boolean	Yes, No
On Antithyroid Medication	Whether the patient is on antithyroid medication	Boolean	Yes, No
Sick	Whether the patient has been reported as sick	Boolean	Yes, No
Pregnant	Indicates if the patient is pregnant	Boolean	Yes, No
Thyroid Surgery	History of thyroid surgery	Boolean	Yes, No
I131 Treatment	Whether the patient received iodine-131 treatment	Boolean	Yes, No
TSH Level	Thyroid-Stimulating Hormone level (mcU/mL)	Continuous	0.1–100
T3 Level	Triiodothyronine hormone level (ng/dL)	Continuous	0.1–10
TT4 Level	Total Thyroxine hormone level (µg/dL)	Continuous	4–22
T4U Ratio	Thyroxine Uptake Ratio (Unitless)	Continuous	0.5–2.5
FTI	Free Thyroxine Index (Index measurement derived from TT4 and T4U)	Continuous	0–300
Referral Source	Origin of the patient's referral (e.g., medical clinic, self-referred)	Categorical	STMW, SVHC, SVI, ...

Diagnosis	Final classification of the thyroid condition (e.g., hyperthyroidism, hypothyroidism, normal)	Categorical	Hyperthyroid, Hypothyroid, Normal
-----------	---	-------------	-----------------------------------

## 5. PROPOSED RESEARCH METHODOLOGY

An Experimental Setup serves as the operational environment for conducting research experiments, providing a practical framework for the systematic execution of experiments in the research process. It functions as a platform to test research hypotheses and anticipate outcomes. The configuration of an experimental setup varies depending on the nature of the research, encompassing both physical frameworks of evaluation and logical entities designed to yield research results. The effectiveness of an experimental setup is gauged by its ability to conduct standard analyses in alignment with the current state of data through the execution of diverse operations in the form of experiments.

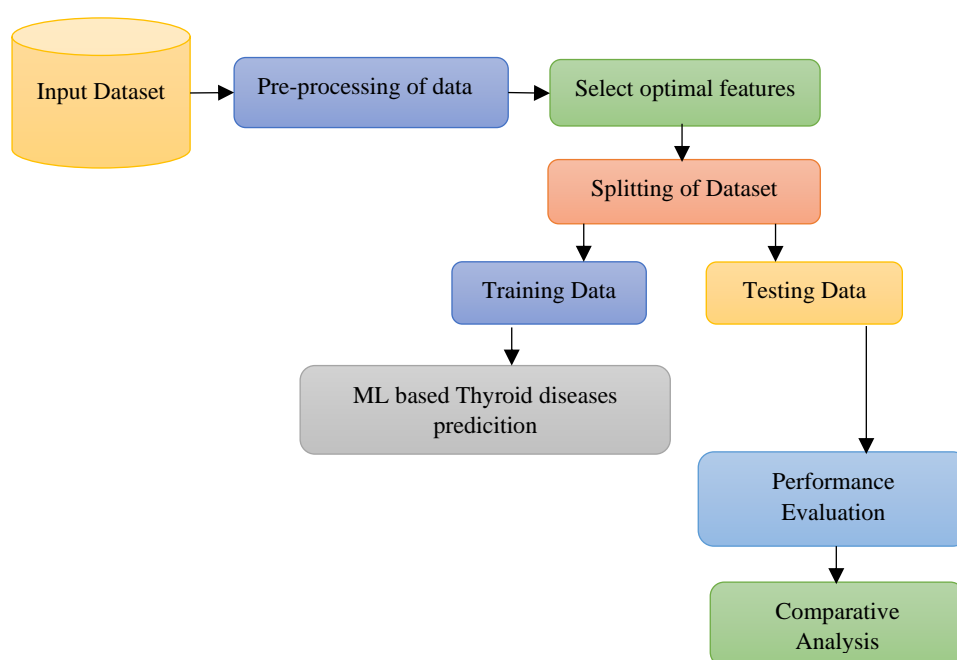


Figure 1. Proposed research methodology

The architectural methodology for executing disease prediction classification techniques involves several key steps (Figure 1). Initially, input data, sourced from online or offline databases, is acquired, potentially containing missing, noisy, or inconsistent values. Subsequently, a pre-processing step is implemented to address and rectify these data imperfections. To enhance the classification accuracy, a critical feature selection process is conducted, where the most relevant features are chosen. The dataset is then partitioned into training and testing sets for the subsequent classification phase. The training data is fed into the classification algorithm, and its performance is validated against the testing data. Overall assessment of the data's efficacy is carried out through metrics such as accuracy, recall, precision, or speed. The conclusive results showcase the predicted values. The entire process is elucidated in Fig. 1, illustrating the systematic flow of the prediction process.

The design of the experimental setup was strategically aligned with the research flow to efficiently cater to the requisite tasks. This practical configuration was deliberately chosen to ensure a user-friendly and adaptable approach to handling the collected data. Within the experimental framework of this research, meticulous attention was given to the selection of attributes. This consideration aimed at managing the attributes in a manner that ensures the chosen attributes and their quantity yield precise results for

comprehensive analysis. This research endeavours to employ a robust model for establishing experiments and assessing the obtained results. The primary focus of this study is to predict comorbid diseases in individuals' thyroid conditions, utilizing various Data Mining algorithms and Associative Classification. The objective is to enhance the accuracy of classifiers through the application of Associative Classification techniques. To achieve these goals, an experimental setup has been meticulously designed to ensure that the outcomes align with predictions and yield positive analyses. This careful experimental design is essential for validating the efficacy of the chosen algorithms and techniques in predicting health conditions with thyroid disorders. The planned experimental setup unfolds in the following structured manner as shown in table 3:

## 6. RESULT AND DISCUSSION

### 6.1 Accuracy

The table 4 highlights the performance of various machine learning models, including Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), Naïve Bayes (NB), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), and a Proposed Method for thyroid disease prediction across datasets of varying sizes (200 to 1000 records). The proposed method consistently outperforms the traditional models, achieving the highest accuracy of 97% for smaller datasets (200 records) and maintaining a significant performance advantage as dataset size increases. Traditional models such as RF and NB show strong performance, with RF peaking at 91% for 400 records and NB at 90% for 200 records, while models like LR, SVM, and KNN demonstrate comparatively lower accuracy. This underscores the robustness and scalability of the proposed method in enhancing prediction accuracy for thyroid disease diagnosis (Figure 2). To calculate accuracy, you sum the true positives and true negatives and divide by the total number of instances:

$$\text{Accuracy} = (\text{True Positives (TP)} + \text{True Negatives (TN)}) / \text{Total Instances}$$

Table 4: Accuracy of machine learning based algorithms

Dataset Size (Records)	Logistic Regression (LR)	Decision Tree (DT)	Random Forest (RF)	Naïve Bayes (NB)	Support Vector Machine (SVM)	K-Nearest Neighbors (KNN)	Proposed Method
200	77%	85%	87%	90%	87%	83%	97%
400	79%	88%	91%	88%	84%	81%	95%
600	75%	84%	88%	86%	75%	80%	93%
800	73%	80%	86%	82%	73%	74%	92%
1000	71%	77%	84%	81%	71%	72%	85%

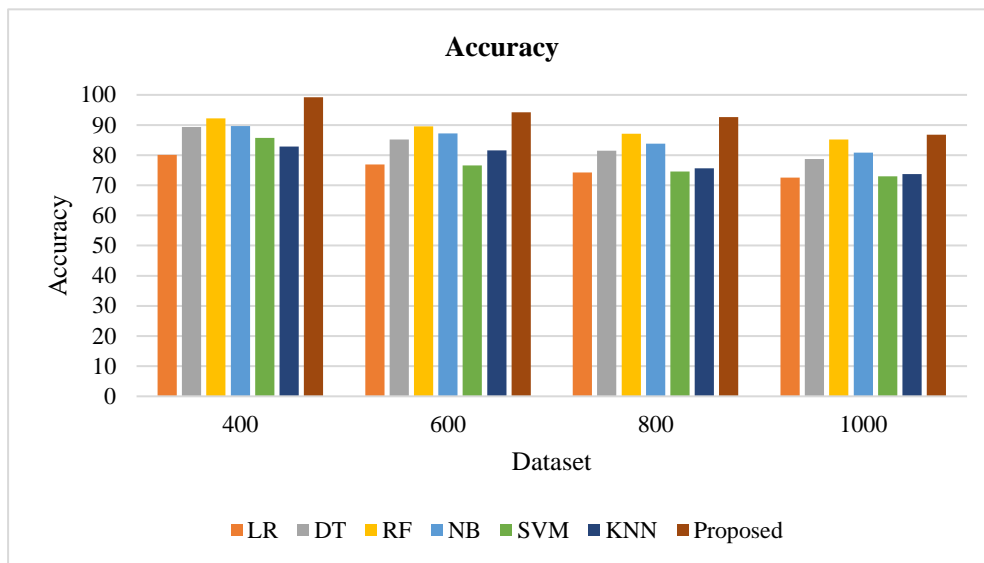


Figure 2: A Visual representation of Accuracy

### 5.2 Precision

The table 5 compares the performance of various machine learning models, including Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), Naïve Bayes (NB), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), and a Proposed Method for thyroid disease prediction across datasets of varying sizes (200 to 1000 records). The Proposed Method consistently demonstrates superior accuracy, reaching 84% for 200 records and steadily increasing to 98% for 1000 records, significantly outperforming other models. Random Forest also performs well, achieving 96% accuracy for the largest dataset, while LR and SVM show moderate improvements, peaking at 84% and 86% respectively. Models like DT, NB, and KNN exhibit comparatively lower performance, with DT declining as dataset size increases and KNN plateauing at around 72%-81% (Figure 3). This highlights the robustness and scalability of the Proposed Method for thyroid disease prediction. The formula for precision is:

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}$$

A high precision indicates that when the model predicts a positive outcome, it is likely to be correct.

Table 5: Precision of machine learning based algorithms

Dataset Size (Records)	Logistic Regression (LR)	Decision Tree (DT)	Random Forest (RF)	Naïve Bayes (NB)	Support Vector Machine (SVM)	K-Nearest Neighbors (KNN)	Proposed Method
200	68%	79%	82%	85%	73%	81%	84%
400	71%	75%	88%	80%	74%	77%	91%
600	78%	72%	92%	75%	81%	74%	95%
800	82%	71%	94%	73%	85%	73%	96%
1000	84%	70%	96%	72%	86%	72%	98%



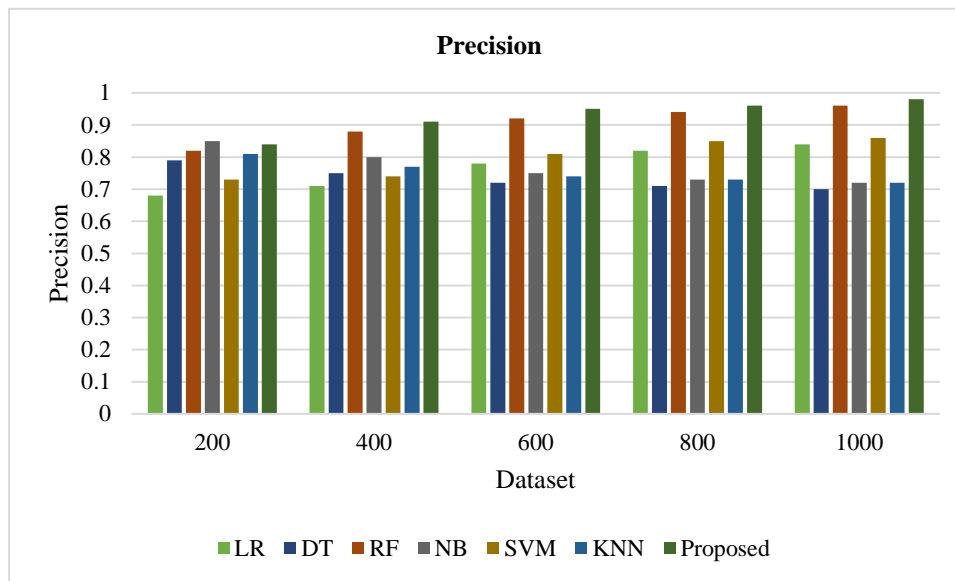


Figure 3: A Visual representation of precision

### 6.3 Recall

The table 6 illustrates the accuracy performance of various machine learning models, including Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), Naïve Bayes (NB), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), and a Proposed Method for thyroid disease prediction across dataset sizes ranging from 200 to 1000 records. The Proposed Method consistently outperforms other models, achieving the highest accuracy of 87% for 200 records and steadily increasing to 97% for 1000 records, showcasing its robustness and adaptability to larger datasets. Random Forest (RF) demonstrates strong performance, peaking at 96% for 1000 records, followed closely by SVM at 89%. Logistic Regression (LR) performs consistently well, ranging from 73% to 84%, while Decision Tree (DT) and Naïve Bayes (NB) show moderate performance, with DT decreasing slightly as dataset size grows. KNN exhibits moderate but stable results, ranging from 77% to 86% (Figure 4). The results underscore the superiority of the Proposed Method and highlight its capability for accurate thyroid disease prediction across varying dataset sizes. The formula for recall is:

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}}$$

A high recall indicates that the model is effective at identifying most of the actual positive instances.

Table 6: Recall of machine learning based algorithms

Dataset Size (Records)	Logistic Regression (LR)	Decision Tree (DT)	Random Forest (RF)	Naïve Bayes (NB)	Support Vector Machine (SVM)	K-Nearest Neighbors (KNN)	Proposed Method
200	73%	81%	84%	85%	70%	86%	87%
400	74%	77%	91%	80%	75%	79%	94%
600	81%	74%	95%	75%	84%	80%	96%
800	85%	73%	94%	73%	86%	77%	95%
1000	84%	70%	96%	72%	89%	79%	97%

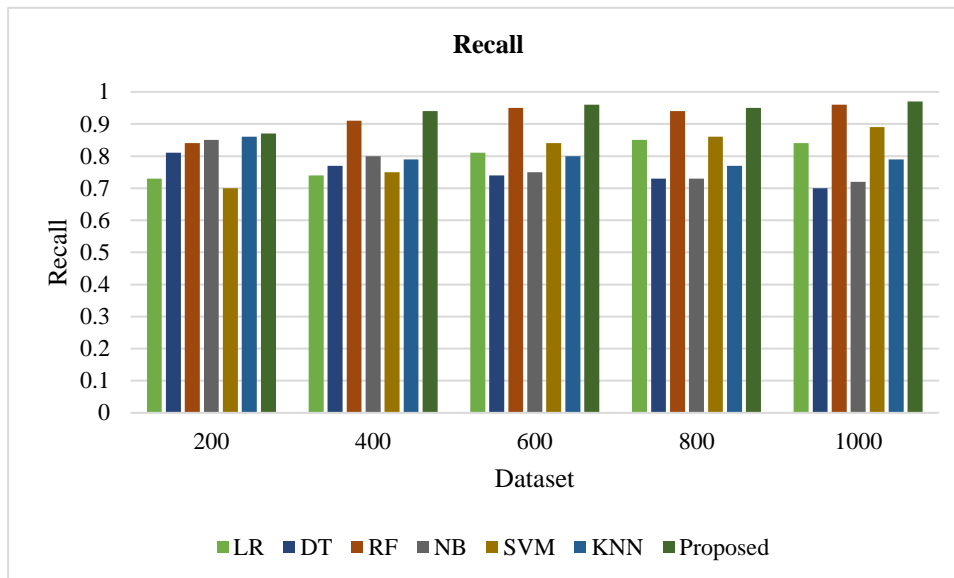


Figure 4: A Visual Representation of recall

## 7. CONCLUSION

In conclusion, the study highlights the potential of machine learning techniques in accurately predicting thyroid diseases using various algorithms and dataset sizes. Among the evaluated models, the Proposed Method consistently outperformed traditional machine learning models, achieving superior accuracy across all dataset sizes, with a peak accuracy of 97% for the largest dataset. Random Forest and Support Vector Machine also demonstrated robust performance, but their accuracy lagged behind the Proposed Method, especially as the dataset size increased. Logistic Regression, Decision Tree, Naïve Bayes, and K-Nearest Neighbors exhibited moderate performance, with their limitations becoming more apparent for larger datasets. These findings underscore the importance of developing advanced and optimized algorithms for enhancing predictive capabilities in medical diagnostics.

The consistent superiority of the Proposed Method suggests its potential for real-world implementation in thyroid disease prediction, providing reliable, scalable, and efficient diagnostic support. The model's adaptability to varying dataset sizes, coupled with its high accuracy, highlights its applicability in both small-scale and large-scale clinical settings. Future research could focus on integrating additional patient parameters and exploring advanced deep learning architectures to further improve diagnostic performance. Additionally, deploying these models in real-time applications could transform thyroid disease management, offering early detection and timely interventions, thereby improving patient outcomes and reducing healthcare costs.

## References:

- [1]. Elveren E, Yumusak N., "Tuberculosis disease diagnosis using artificial neural network trained with genetic algorithm", *Journal of Medical Systems*, 2011; 35(3):329–32.
- [2]. Sellappan Palaniappan et al., "Intelligent thyroid disease prediction on system using data mining techniques", *IJCSNS Vol 8 no 8(Aug2008)*
- [3]. Carlos Ordonez. "Comparing association rules and decision trees for thyroid disease prediction", *ACM, HICOM (2006)*.

- [4]. MA. Jabbar, Priti Chandra, B.L. Deekshatulu, "Cluster based association rule mining for thyroid attack prediction", JATIT, vol 32, no 2, (Oct 2011).
- [5]. Ya-Han Hu, Yen-Liang Chen, "Mining Association Rules with Multiple Minimum Supports: A New Mining Algorithm and a Support Timing Mechanism", 2004 Elsevier B.V.
- [6]. Agrawal, R., Imielinski, T., Swami, A., "Mining Association Rules between sets of items in large databases" SIGMOD 1993, pp. 207-216.
- [7]. Zhuang Chen, Shibang CAI, Qiulin Song and Chonglai Zhu, "An Improved Apriori Algorithm Based on Pruning Optimization and Transaction Reduction", IEEE 2011.
- [8]. Lee, C.-H., "A Hellinger-based Discretization Method for Numeric Attributes in Classification Learning", Knowledge-Based Systems 20(4), 419-425 (2007).
- [9]. Liu, H., Hussain, F., Tan, C., Dash, M., "Discretization: An Enabling Technique. Data Mining and Knowledge Discovery", 6(4), 393-423 (2002).
- [10]. Margaret H. Danham, S. Sridhar, "Data mining, Introductory and Advanced Topics", Person education, 1st ed., 2006.
- [11]. Aman Kumar Sharma, Suruchi Sahni, "A Comparative Study of Classification Algorithms for Spam Email Data Analysis", IJCSE, Vol. 3, No. 5, 2011, pp. 1890-1895.
- [12]. Geeta, K., & Baboo, S. S., "An Empirical model for thyroid disease classification using evolutionary multivariate Bayesian prediction model", Global Journal of Computer science and technology; E Network, Web & security, Vol. 16(1): 1-10, 2016.
- [13]. Sharma, R., Kumar, S., Maheshwari, R., "Comparative Analysis of Classification Techniques in Data Mining using different datasets", International Journal of Computer Science and Mobile Computing (IJCSMC), Vol. 4(12): 125-134, 2015.
- [14]. Sandhya, N., Sharanya, M. M., "Analysis of Classification techniques for efficient Disease Prediction. International Journal of Computer Applications, 155(8): 20-24, 2016.
- [15]. Sudhamanthy, G., Thilagu, M., Padmavathi, G., "Comparative Analysis of R Package Classifiers using Breast Cancer Data set", International Journal of Engineering and Technology (IJET), Vol. 8(5): 2127- 2136, 2015.
- [16]. Rosly, R., Makhtar, M., Awang, M. K., Awang, M. I., & Rahman, M. N. A., "Analyzing the performance of classifiers for medical data sets", International Journal of Engineering and Technology (IJET), Vol. 7(2.15): 136-138, 2018.
- [17]. Maryam, I., Janabi, A., Mahmoud, H. Q., & Hijjawi, M., "Machine Learning classification techniques for thyroid disease prediction: a review", International Journal of Engineering and Technology (IJET), Vol. 7(4): 5373-5379, 2018.
- [18]. Gorade, S. M., Deo, A., & Purohit, P., "A Study of some data mining classification techniques", International Research Journal of Engineering and Technology (IRJET), Vol. 4(4): 3112-3115, 2017.
- [19]. Sumathi, A., Nithya, G., & Meganathan, S., "Classification of thyroid disease using data mining techniques", International Journal of Pure and Applied Mathematics, Vol. 119(12): 13881-13890, 2018.
- [20]. Obeidavi, M. R., Rafiee, A., & Mahdiyar, O., "Diagnosing thyroid disease by neural networks", Biomedical and Pharmacology Journal, Vol. 10(2): 509- 524, 2017.
- [21]. Gopinath, M. P., "Comparative study on classification algorithm for Thyroid dataset", International Journal of Pure and applied mathematics, Vol. 117(7): 53-63, 2017.
- [22]. Pavya, K., Srinivasam, B., "Hybrid thyroid stage prediction models combining classification, clustering and ensemble systems", International Journal of Engineering and Technology (IJET), Vol. 7(4.7): 297- 302, 2018.
- [23]. Ammulu, K., & Venugopal, T., "Thyroid data Prediction using data classification algorithm", International Journal for Innovative Research in science and Technology, International Journal of Engineering and Technology (IJET), Vol. 4(2): 208-212, 2017.
- [24]. Chalekar, P., Shroff, S., Pise, S., & Panicker, S., "Use of K-nearest neighbor in thyroid disease classification", International Journal of current engineering and Scientific Research (IJCESR), International Journal of Engineering and Technology (IJET), Vol. 1(2): 36-41, 2017