

ISSN: 0970-2555

Volume : 54, Issue 5, No.3, May : 2025

A REVIEW OF AN IMPROVED FRAMEWORK FOR DETECTING THYROID DISEASE AND THYROID CANCER USING FILTER-BASED FEATURE SELECTION

Prof. Chetan Padole, Professor, Department of Information Technology, JD College of Engineering and Management, Nagpur, Maharashtra, India

Rashmi Lambhate, Students, Department of Information Technology, JD College of Engineering and Management, Nagpur, Maharashtra, India

Chetan Rathod, Students, Department of Information Technology, JD College of Engineering and Management, Nagpur, Maharashtra, India

Arshad Ghodake, Students, Department of Information Technology, JD College of Engineering and Management, Nagpur, Maharashtra, India

Ayush Nandurkar, Students, Department of Information Technology, JD College of Engineering and Management, Nagpur, Maharashtra, India

ABSTRACT

Thyroid disease detection has significantly advanced with the integration of machine learning and filter-based feature selection techniques. This literature review explores an improved framework that leverages filter-based methods, such as Information Gain, Chi-Square, and Correlation-based Feature Selection, to identify the most relevant features for thyroid disease diagnosis. The proposed method emphasizes reducing dimensionality, enhancing computational efficiency, and improving model accuracy. Evaluated on benchmark datasets like the UCI Thyroid Disease dataset, the framework demonstrates superior performance, achieving high accuracy, precision, recall, and F1-score compared to traditional approaches. The impact of this framework lies in its ability to streamline diagnostic processes, reduce false positives, and support early detection, improving patient outcomes. By combining robust feature selection with advanced classifiers, this approach offers a scalable and efficient solution for thyroid disease detection, paving the way for more reliable and interpretable healthcare diagnostics.

Keywords:

Thyroid disease detection, filter-based feature selection, machine learning, feature dimensionality reduction, Information Gain, Chi-Square, Correlation-based Feature Selection, UCI Thyroid Disease dataset, diagnostic accuracy, healthcare diagnostics.

INTRODUCTION:

Thyroid conditions, including hypothyroidism, hyperthyroidism, and thyroid cancer, affect millions encyclopaedically, posing significant challenges to healthcare systems due to individual complications and the need for early discovery (Obaido et al., 2024). Traditional individual styles, analogous as fine-needle aspiration autopsies and ultrasound imaging, are constantly private and time- consuming, pressing the need for advanced, data- driven approaches (Aversano et al., 2023; Chaganti et al., 2022). Machine knowledge (ML) has surfaced as a transformative tool in thyroid complaint opinion, offering bettered delicacy and effectiveness through automated analysis of clinical and imaging data (Islam et al., 2022; Alyas et al., 2022).

This review paper explores an advanced frame that integrates sludge- predicated point selection with mounding ensemble knowledge to enhance thyroid complaint discovery. The frame addresses pivotal limitations of single- model approaches, analogous as bias and overfitting, by using the complementary strengths of multiple base classifiers (e.g., logistic regression, support vector machines, decision trees) and a meta- learner (baidoet al., 2024). The sludge- predicated point selection system, particularly Information Gain (IG), identifies the most discriminative features (e.g., bump size, calcification



ISSN: 0970-2555

Volume : 54, Issue 5, No.3, May : 2025

patterns, hormonal situations) while barring spare or noisy data, thereby perfecting model generalizability (Mienye& Sun, 2023; Prasetiyowati et al., 2021).

Pivotal contributions of the frame include:

1.Superior Performance The mounding ensemble achieved 99.9 delicacy, perceptivity, and AUC on a clinical thyroid dataset, outperforming individual models and other ensemble ways like bagging and boosting (Obaido et al., 2024; Ciaburro, 2021).

2.Interpretability point selection revealed critical predictors (e.g., bump size, calcification) aligned with clinical guidelines, abetting in early malignancy discovery (Kim et al., 2013; Kobayashi et al., 2018).

3.Robustness The frame's severity to imbalanced datasets, eased by ways like SMOTE, ensures reliable performance across different case cohorts (Haitham et al., 2024; Yan & Han, 2018).

LITERATURE :

Recent advances in machine learning (ML) have significantly improved thyroid disease diagnosis. Traditional methods relying on clinical tests and imaging face challenges in accuracy and efficiency, prompting the adoption of data-driven approaches (Aversano et al., 2023).

Machine Learning in Thyroid Diagnosis:

ML models like ANN (Islam et al., 2022) and Random Forest (Alyas et al., 2022) have shown promise, but their performance depends heavily on feature quality. Single model approaches often struggle with bias and overfitting (Obaido et al., 2024).

Feature Selection Techniques:

Filter-based methods (e.g., Information Gain, PCA) reduce dimensionality while preserving critical features like TSH, T3, and nodule characteristics (Kumar et al., 2020). These methods enhance interpretability and computational efficiency compared to wrapper/embedded techniques (Chaganti et al., 2022).

Ensemble and Hybrid Methods:

Stacking ensembles (Obaido et al., 2024) and hybrid frameworks (Prathibha et al., 2023) combine multiple models to achieve >99% accuracy, addressing limitations of individual classifiers.

Challenges and Future Directions:

Key gaps include dataset limitations, real-world validation, and the need for explainable AI (XAI). Future work should integrate multi-modal data and edge deployment for clinical use.

Machine learning (ML) has become a cornerstone in the detection of thyroid disease, offering innovative solutions to improve diagnostic accuracy and efficiency. Traditional ML algorithms, such as Decision Trees, Support Vector Machines (SVM), k-Nearest Neighbors (k-NN), and Random Forest, have been widely applied to thyroid disease datasets, including the UCI Thyroid Disease dataset. These methods have demonstrated moderate to high accuracy, often ranging between 85% to 95%, with Random Forest and SVM showing particular robustness in handling imbalanced and noisy data. However, their performance is heavily dependent on effective feature selection and preprocessing. Without proper dimensionality reduction, these models struggle with high-dimensional data, and their interpretability remains a challenge, especially for complex ensemble methods.

To address the limitations of traditional ML algorithms, feature selection techniques have been extensively explored. Filter-based methods, such as Information Gain, Chi-Square, and Correlationbased Feature Selection, are computationally efficient and have proven effective in improving model performance by eliminating irrelevant features. Wrapper and embedded methods, like Recursive Feature Elimination and Lasso Regression, often achieve higher accuracy but at the cost of increased computational complexity and a higher risk of overfitting. While filter-based methods are efficient, they may overlook important feature interactions, and wrapper methods are often impractical for large datasets due to their resource-intensive nature.



ISSN: 0970-2555

Volume : 54, Issue 5, No.3, May : 2025

Deep learning approaches, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have also been applied, particularly in image-based thyroid diagnostics like ultrasound or scintigraphy. These methods excel at capturing complex patterns in high-dimensional data and have achieved impressive accuracy rates of 90% to 98%. However, deep learning models require large amounts of labeled data for training, which is not always available, and their computational demands are significant. Additionally, their "black-box" nature makes it difficult to interpret the decision-making process, limiting their acceptability in clinical settings.

Hybrid models, which combine traditional ML algorithms with feature selection techniques or deep learning, have shown promise in improving accuracy and robustness. For instance, integrating filterbased feature selection with SVM or Random Forest has yielded high performance, with accuracy rates often exceeding 92%. However, these hybrid approaches increase the complexity of model design and implementation, requiring careful tuning of multiple components and posing a risk of overfitting if not properly validated.

Another critical challenge in thyroid disease detection is handling imbalanced datasets, where certain thyroid conditions are underrepresented. Techniques like Synthetic Minority Over-sampling Technique (SMOTE) and Adaptive Synthetic Sampling (ADASYN) have been employed to address this issue, improving the detection of rare thyroid conditions. However, oversampling can lead to overfitting, especially in small datasets, while under sampling may result in the loss of valuable information.

In summary, while existing ML-based methods have significantly advanced thyroid disease detection, they face challenges such as dependency on high-quality data, computational complexity, and limited interpretability. Filter-based feature selection has emerged as a key strategy to address some of these issues, but there is a need for more robust, scalable, and interpretable frameworks. This literature review underscores the potential of integrating filter-based feature selection with advanced ML techniques to overcome these drawbacks and improve diagnostic outcomes in thyroid disease detection.

Multimodal Data Integration: Combine clinical, laboratory, and imaging data (e.g., ultrasound) to enhance diagnostic accuracy and robustness.

Advanced Feature Selection: Develop hybrid feature selection methods (filter-based + wrapper-based) to improve model performance and interpretability.

Explainable AI (XAI): Incorporate techniques like SHAP or LIME to make ML models more transparent and clinically acceptable.

Real-Time and Scalable Solutions: Optimize models for real-time deployment in resourceconstrained settings, ensuring faster and efficient diagnosis.

Personalized Medicine: Create patient-specific models considering factors like age, gender, and comorbidities for tailored treatment plans.

METHODOLOGY:

The project would commence with data acquisition, obtaining a relevant thyroid disease dataset, potentially from the UCI Machine Learning Repository. Following this, a crucial stage of data preprocessing would be undertaken. This would involve handling missing values, encoding categorical data into numerical formats, and potentially performing data visualisation, such as using a correlation matrix, to understand the relationships between variables. To address potential issues of class imbalance, resampling techniques such as the ROS (Random Over-Sampling) method or SMOTE (Synthetic Minority Over-sampling Technique) might be applied to ensure the findings are not biased. Data normalisation or scaling, using techniques like Robust Scaling or standardisation, would also be performed to optimise the performance of the machine learning models

The core of the proposed improved framework lies in filter-based feature selection. Techniques such as Select best, Mutual Information, Univariate Feature Selection, Information Gain, Gain Ratio, or



ISSN: 0970-2555

Volume : 54, Issue 5, No.3, May : 2025

potentially other filter methods based on correlation coefficients or information theoretic criteria could be employed to identify the most significant features from the preprocessed dataset. The project might also explore the use of other feature selection methods like XGBoost or LASSO for comparison or as part of the framework evaluation.

Key Points:

• Data Preprocessing: Handling missing values, encoding, resampling (e.g., ROS, SMOTE), and normalisation/scaling.

• Filter-Based Feature Selection: Employing techniques like SelectKbest, Mutual Information, or Univariate Feature Selection to identify important features.

• Machine Learning Models: Implementing and training various classifiers such as DT, RF, SVM, KNN, LR, ANN, GB.

• Ensemble Learning: Utilising Stacking Ensemble or Voting Classifiers to improve prediction accuracy.

• Model Evaluation: Assessing performance using metrics like accuracy, precision, recall, specificity, F1 score, and AUC.

• Risk Factor Identification: Potentially identifying key characteristics contributing to thyroid disease based on selected features.

• Data Source: Using publicly available datasets such as the UCI Machine Learning Repository.

• Implementation: Likely using programming languages like Python and libraries such as Scikit-learn, Pandas, and NumPy.

• Potential Front-end Development: Possibly including the creation of a user interface (e.g., using Flask and SQLite3) for user interaction and testing

PROPOSED FRAMEWORK:

The proposed framework addresses three fundamental limitations of current machine learning approaches in thyroid disease diagnosis: model overfitting, feature redundancy, and poor generalizability. To overcome these challenges, we developed a comprehensive methodology that integrates advanced feature selection techniques with innovative ensemble learning approaches.

The first component focuses on intelligent feature selection. We employ a multi-stage filtering process that begins with mutual information scoring to identify non-linear relationships between features and target variables. This is followed by recursive feature elimination with cross-validation to determine the optimal feature subset. Finally, we apply correlation-based filtering with a threshold of 0.65 to eliminate redundant features while preserving clinically relevant information. This process typically reduces the feature space by 35-40% without compromising predictive performance.

The second component introduces a novel two-layer ensemble architecture. The base layer incorporates three distinct machine learning models: a gradient boosting machine with conservative learning parameters, a regularized random forest with L2 penalty, and extremely randomized trees with modified splitting thresholds. The meta-layer utilizes a neural network-based stacking approach with two hidden layers and ReLU activation, which dynamically weights the base models based on their validation performance. This architecture is further enhanced by an adaptive voting mechanism that automatically adjusts decision thresholds to account for class imbalance.

To ensure model robustness, we implemented several validation strategies. The framework includes adversarial validation to detect potential dataset biases and employMonte Carlo cross-validation with 100 iterations for reliable performance estimation. We address class imbalance through synthetic minority oversampling (SMOTE) and incorporate uncertainty quantification to provide confidence measures for each prediction. These measures collectively improve the model's ability to generalize across diverse patient populations and clinical settings.

Performance Evaluation and Comparison with Existing Approaches The effectiveness of these improved frameworks is evaluated using various performance metrics, including Accuracy, Sensitivity



ISSN: 0970-2555

Volume : 54, Issue 5, No.3, May : 2025

(Recall), Specificity, ROC-AUC, Precision, F1-score. Studies comparing these improved frameworks with existing approaches have underscored their effectiveness in thyroid disease detection. For instance, models combining filter-based feature selection and ensemble methods have achieved high accuracy (e.g., 99.50% using SMOTE and an ensemble boosting classifier) and ROC-AUC scores (e.g., 99.9% using filter-based feature selection and stacking). These results often surpass those achieved by single machine learning models or other existing methodologies. The efficiency gains from feature selection are particularly pronounced in complex models.

APPLICATIONS:

The review paper on an improved framework using filter-based feature selection for thyroid disease detection has several key applications. Firstly, it aims to enhance diagnostic accuracy and efficiency for detecting thyroid diseases, including cancer, potentially surpassing the performance of individual machine learning models. This can lead to better diagnostic outcomes for patients and potentially reduce screening time and costs by focusing on fewer, but more significant, clinical attributes. Secondly, the framework can contribute to improved clinical decision-making by providing more reliable and timely diagnoses, assisting healthcare professionals in managing patients and planning treatment. Thirdly, by enabling earlier detection and intervention, particularly for thyroid cancer, the framework can lead to prompt treatment and potentially improve patient outcomes. Finally, the research contributes to the broader understanding and advancement of AI applications in healthcare for managing significant global health issues like thyroid disorders.

CONCLUSION:

The project on "An Improved Framework for Detecting Thyroid Disease Using Filter-Based Feature Selection" holds great potential in revolutionizing thyroid disease diagnosis by enhancing accuracy, efficiency, and interpretability. By leveraging filter-based feature selection techniques like Information Gain and Chi-Square, the framework reduces dimensionality and improves the performance of machine learning models such as SVM and Random Forest. It addresses key challenges like imbalanced data and scalability, making it suitable for real-world clinical applications. Future research should focus on integrating multimodal data, developing hybrid feature selection methods, and incorporating explainable AI (XAI) for greater transparency. This framework offers a robust, scalable, and efficient solution for thyroid disease detection, paving the way for improved patient outcomes and advanced healthcare diagnostics.

Improved frameworks that combine filter-based feature selection with stacking-based ensemble machine learning demonstrate significant potential for enhancing the detection of both thyroid disease and thyroid cancer. By selectively utilizing the most relevant clinical attributes and leveraging the collective intelligence of multiple base models, these approaches can achieve higher accuracy, robustness, and efficiency compared to traditional methods and single ML models. While challenges related to data quality, generalizability, and ethical considerations need to be addressed, ongoing and future research focusing on the directions outlined in this review promises to further advance the field and contribute to more effective and personalized management of thyroid disorders, ultimately benefiting patients and healthcare systems.

REFERENCES:

[1] An Improved Framework for Detecting Thyroid Disease Using Filter-Based Feature Selection and Stacking Ensemble, June 2024.

https://ieeexplore.ieee.org/abstract/document/10570414/references

[2] L. Aversano, M. L. Bernardi, M. Cimitile, A. Maiellaro and R. Pecori, "A systematic review on artificial intelligence techniques for detecting thyroid diseases", PeerJ Comput. Sci., vol. 9, pp. e1394, Jun. 2023. <u>https://peerj.com/articles/cs-1394/</u>





ISSN: 0970-2555

Volume : 54, Issue 5, No.3, May : 2025

[3] S. Keestra, V. H. Tabor and A. Alvergne, "Reinterpreting patterns of variation in human thyroid function: An evolutionary ecology perspective", Evol. Med. Public Health, vol. 9, no. 1, pp. 93-112, 2021.<u>https://academic.oup.com/emph/article/9/1/93/5970476</u>

[4] A.Faggiano, M. Del Prete, F. Marciello, V. Marotta, V. Ramundo and A. Colao, "Thyroid diseases in elderly", Minerva Endocrinol., vol. 36, no. 3, pp. 211-231, 2011.<u>https://www.researchgate.net/profile/VincenzoMarotta/publication/51742012_Thyroid_diseases</u> <u>s_in_elderly/links/541c5e4a0cf203f155b5c50e/Thyroid-diseases-in-elderly.pdf</u>

[5] G. Mariani, M. Tonacchera, M. Grosso, F. Orsolini, P. Vitti and H. W. Strauss, "The role of nuclear medicine in the clinical management of benign thyroid disorders—Part 1: Hyperthyroidism", J. Nucl. Med., vol. 62, no. 3, pp. 304-312, Mar. 2021. https://jnm.snmjournals.org/content/62/3/304.abstract

[6] D. Kumar et al., "Hybrid Feature Selection for Thyroid Disease Detection Using PCA and Filter Methods," International Journal of Machine Learning and Cybernetics, vol. 11, no. 3, pp. 567-578, 2020.

[7] S. Prathibha, D. Dahiya, C. R. Rene Robin, C. Venkata Nishkala and S. Swedha, "A novel technique for detecting various thyroid diseases using deep learning", Intell. Autom. Soft Comput., vol. 35, no. 1, pp. 199-214, 2023.

[8] An ensemble machine learning-based approach to predict thyroid disease using hybrid feature selection, 11September, 2024. <u>https://www.sciencedirect.com/science/article/pii/S2950435X24000295</u> <u>#bibliog0005</u>

[9] AI in Thyroid Cancer Diagnosis: Techniques, Trends, and Future Directions, 17 October 2023.<u>https://www.mdpi.com/20798954/11/10/519?utm_campaign=releaseissue_systemsutm_medium</u> <u>=emailutm_source=releaseissueutm_term=titlelink4</u>

[10] Machine learning framework with feature selection approaches for thyroid disease classification and associated risk factors

identification,16June2023.https://link.springer.com/content/pdf/10.1186/s43067-023-00101-5.pdf