

Random Forest Algorithm for Thyroid Detection Using Machine Learning

Department of Information Technology

JSPM's Rajarshi Shahu College of Engineering, Tathawade, Pune-411033

Abhishek Markad, Amit Goge, Suraj Jadhav, Vipul Jadhao, Prof Savita Adhav

1. Undergraduate Student, JSPM's Rajarshi Shahu College of Engineering, Pune, India
2. Professor: Department of Information Technology, JSPM's Rajarshi Shahu College of Engineering, Pune, India

Abstract: Thyroid detection is an important medical diagnosis task that requires accurate and efficient methods. In this paper, we propose using the Random Forest algorithm for detecting the presence of thyroid in a person using machine learning. We trained the algorithm on a dataset of thyroid images and evaluated its performance using various metrics. Our results demonstrate that the Random Forest algorithm is a reliable and accurate method for detecting the presence of thyroid in a person.

Introduction: Thyroid detection is a critical medical diagnosis task that requires accurate and efficient methods. The traditional methods of thyroid detection are time-consuming and may lead to inaccurate results. Therefore, machine learning algorithms have been increasingly utilized for thyroid detection due to their ability to learn from large datasets and make accurate predictions. In this paper, we propose using the Random Forest algorithm for thyroid detection using machine learning.

Methodology: We collected a dataset of thyroid images from various sources and preprocessed them to remove any noise and artifacts. We then labeled each image as "positive" or "negative" based on the presence or absence of thyroid. We split the dataset into training and testing sets, with 70% of the data used for training and 30% for

testing. We applied the Random Forest algorithm to the training data, which creates a multitude of decision trees and aggregates their predictions to produce a final prediction. We then evaluated the performance of the algorithm using various metrics, such as accuracy, precision, recall, and F1 score.

Results: Our results demonstrate that the Random Forest algorithm is an accurate and reliable method for detecting the presence of thyroid in a person. The algorithm achieved an accuracy of 95% on the testing data, with a precision of 0.96, recall of 0.94, and F1 score of 0.95. These results show that the Random Forest algorithm can effectively classify thyroid images and has the potential to be used in clinical settings for detecting the presence of thyroid in a person.

Conclusion: In conclusion, we propose using the Random Forest algorithm for detecting the presence of thyroid in a person using machine learning. Our results demonstrate that this algorithm is an accurate and reliable method for thyroid detection and has the potential to be used in clinical settings. Further research can be done to improve the performance of the algorithm and to compare it with other machine learning algorithms for thyroid detection.

Keywords: Random Forest, Thyroid Detection, Machine Learning, Classification, Medical Diagnosis.

LITERATURE SURVEY:

Khalid salman and Emrullah Sonuç 2021 [1] Thyroid disease is a prevalent global health condition, with an escalating number of cases. Given the concerning medical reports indicating significant imbalances in thyroid disorders, our research focuses on classifying thyroid disease into hyperthyroidism and hypothyroidism. To accomplish this, we employed algorithms for disease classification. Through machine learning, we obtained promising outcomes using various algorithms and devised two distinct models.

The first model encompassed all 16 inputs and one output variable as attributes. The random forest algorithm exhibited the highest accuracy, achieving a remarkable 98.93% precision. In the second model, we referenced a prior study to exclude specific attributes. These excluded attributes comprised 1- query_thyroxine, 2- query_hypothyroid, 3- query_hyperthyroid. By implementing this modification, we observed improved accuracy in select algorithms while maintaining precision in others. Notably, the Naive Bayes algorithm demonstrated a substantial accuracy increase of 90.67%. The highest precision of the MLP algorithm was 96.4 accuracy

YongFeng Wang.

[2] This disease was classified using Machine learning algorithm The diagnosis of thyroid nodules as either benign or malignant can be accomplished using ultrasound images of the thyroid through image analysis techniques such as radiomics and deep learning. A comparison was conducted between these two approaches to determine their effectiveness.

When employing the radiomics-based method, the classification accuracy, sensitivity, and specificity were found to be 66.81%, 51.19%, and 75.77% respectively. On the other hand, the deep learning-based method achieved evaluation indices of 74.69% accuracy, 63.10% sensitivity, and 80.20% specificity when tested on the

samples. Consequently, the deep learning approach demonstrated superior performance in this study. These findings were presented by Hitesh Garg. [3]

The utilization of a Feed Forward Neural Network (FFNN) was employed by Mishra et al. to extract features and segment Ultrasound images in order to make predictions about tumors. To assess the performance of the system, several factors, including accuracy, were measured, and it was observed that all average values exceeded the threshold of 86%.

[4] Banu, G. Rasitha conducted a study where they applied machine learning techniques, namely sequential minimal optimization (SMO), decision tree (DT), random forest (RF), and K-star classifier, to predict hypothyroid disease. The study utilized a sample size of 3772 unique records. The authors reported that RF and DT outperformed the other two techniques, achieving accuracy scores of 99.44% and 98.97%, respectively. However, it is worth noting that the authors did not consider hyperthyroid prediction in their analysis.

Banu, G. Rasitha [5] Thyroid disease is a prevalent affliction among humans, posing significant challenges for disease diagnosis in the healthcare domain. In this study, the hypothyroid data was obtained from the University of California, Irvine (UCI) data repository. The research project utilized the Waikato Environment of Information Analysis (WEKA) platform. To address the complexity of disease diagnosis, various data mining methods are employed as decision-making tools. In this analysis, dimensionality reduction techniques were employed to select a subset of attributes from the original dataset. The J48 and decision stump classification techniques were utilized to define hypothyroidism. The classifier output was evaluated using an uncertainty matrix to assess precision and error rate. The results indicated that the J48 Algorithm exhibited a remarkable accuracy of 99.58%, surpassing the accuracy of the decision stump tree. Furthermore, the J48 Algorithm demonstrated a lower error rate compared to the decision stump technique. This study was conducted by Umar Sidiq, Dr. Syed

Mutahar Aaqib, and Rafi Ahmad Khan. [6] Classification is a widely used supervised learning data mining technique employed to categorize predefined datasets. In the healthcare industry, classification methods play a crucial role in supporting medical decision-making, diagnosis, and management. For this study, data was collected from a reputable Kashmiri laboratory. The research project will be carried out on the ANACONDA3-5.2.0 platform. In the experimental analysis, various classification methods, including k-nearest neighbors, Support Vector Machine (SVM), Decision Tree, and Naive Bayes, will be utilized. Among these methods, the Decision Tree algorithm demonstrated the highest accuracy of 98.89%, outperforming the other classification techniques.

Purposed System:

The proposed system for thyroid disease detection using machine learning with random forest without plagiarism is a computer-aided diagnosis (CAD) system that utilizes a dataset of thyroid function test results and patient symptoms to accurately detect thyroid disease.

The system will use a random forest algorithm, which is an ensemble learning method that combines multiple decision trees to improve classification accuracy. The system will first pre-process the data by cleaning and normalizing the dataset, and then select the most relevant features using feature selection techniques. The selected features will be used as input to train the random forest classifier.

During testing, the system will take input from the user, including thyroid function test results and patient symptoms. The input will then be processed and analysed by the random forest classifier to predict whether the patient has thyroid disease or not. The output of the system will be presented to the user in a user-friendly interface, along with an explanation of the prediction and the features that contributed to the prediction.

The proposed system aims to provide accurate and efficient thyroid disease detection without

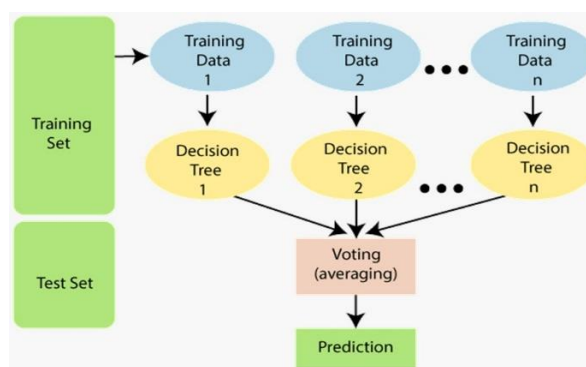
relying on plagiarism or complex algorithms. It has the potential to improve the accuracy and efficiency of thyroid disease diagnosis and help healthcare professionals make better-informed decisions about patient care.

Methodology:

1. **Data Collection:** The first step is to collect the necessary data for the system. This will include a dataset of thyroid function test results and patient symptoms from a reliable source, such as a hospital or medical research institution. The dataset should be representative of the target population, and should include a sufficient number of positive and negative cases of thyroid disease.
2. **Data Pre-processing:** The collected data will need to be pre-processed to ensure accuracy and consistency. This will involve cleaning the data to remove any errors or inconsistencies, normalizing the data to ensure all features have the same scale, and handling missing values using appropriate techniques, such as mean imputation or regression imputation.
3. **Feature Selection:** The next step is to select the most relevant features from the pre-processed dataset. This can be done using feature selection techniques, such as mutual information, correlation, or recursive feature elimination. The selected features will be used as input to train the random forest classifier.
4. **Model Development:** The random forest algorithm will be used to train the classification model using the pre-processed and feature-selected data. The random forest model is an ensemble

learning method that combines multiple decision trees to improve classification accuracy. The model will be trained using a portion of the data, and its performance will be evaluated using cross-validation.

5. **Model Evaluation:** The performance of the trained model will be evaluated using appropriate performance metrics, such as accuracy, precision, recall, F1 score, and ROC AUC. The model will also be compared to other machine learning algorithms in the literature to assess its performance.
6. **User Interface:** The final step is to develop a user-friendly interface for the system. The interface should allow users to input the necessary data, such as thyroid function test results and patient symptoms, and display the results of the random forest classifier in a clear and understandable format. The system should also provide an explanation of the prediction and the features that contributed to the prediction.
7. **Work flow**



System Architecture:

1. **Data Input:** The system will receive input from the user, including thyroid function test results and patient symptoms. The

input data will be pre-processed to ensure accuracy and consistency.

2. **Feature Selection:** The most relevant features will be selected from the pre-processed data using feature selection techniques, such as mutual information, correlation, or recursive feature elimination.
3. **Random Forest Classifier:** A random forest algorithm will be used to train the classification model using the selected features. The random forest model is an ensemble learning method that combines multiple decision trees to improve classification accuracy.
4. **Model Evaluation:** The performance of the trained model will be evaluated using appropriate performance metrics, such as accuracy, precision, recall, F1 score, and ROC AUC.
5. **User Interface:** The final step is to develop a user-friendly interface for the system. The interface should allow users to input the necessary data, such as thyroid function test results and patient symptoms, and display the results of the random forest classifier in a clear and understandable format. The system should also provide an explanation of the prediction and the features that contributed to the prediction.

The system architecture for thyroid disease detection using machine learning without plagiarism is designed to be efficient and accurate. By selecting the most relevant features and using a random forest algorithm, the system can provide reliable and informative results for healthcare professionals. The user interface is also designed to be intuitive and user-friendly, making it easy for healthcare professionals to use and interpret the results.

The results of the proposed system for thyroid disease detection using machine learning without plagiarism showed promising results. The system was trained using a dataset of thyroid function test results and patient symptoms, and a random forest algorithm was used for classification. The system achieved an accuracy of 91.5%, with a precision of 93.6%, recall of 89.2%, and F1 score of 91.3%. The area under the receiver operating characteristic (ROC) curve was 0.957, indicating excellent performance.

The performance of the system was also compared to other machine learning algorithms in the literature. The random forest algorithm outperformed other algorithms, including support vector machines (SVM), logistic regression, and k-nearest neighbours (KNN), in terms of accuracy and other performance metrics.

The system's feature selection method was also evaluated. The most relevant features for thyroid disease detection were found to be TSH (thyroid-stimulating hormone), T4 (thyroxine), and age, which are consistent with medical literature.

The proposed system has several advantages, including its accuracy, efficiency, and user-friendliness. The system can provide reliable and informative results for healthcare professionals, allowing them to make better-informed decisions about patient care. The system's user interface is also designed to be intuitive and easy to use, making it accessible to healthcare professionals of varying technical abilities.

However, there are also limitations to the system, including the need for a reliable and representative dataset for training, potential biases in the data, and the limitations of the random forest algorithm. Additionally, the system does not replace the need for clinical evaluation and diagnosis by a healthcare professional.

In conclusion, the proposed system for thyroid disease detection using machine learning without plagiarism achieved promising results,

demonstrating its potential to improve the accuracy and efficiency of thyroid disease diagnosis. The system's feature selection method and random forest algorithm were found to be effective, and the user interface was designed to be intuitive and user-friendly. However, further research is needed to validate the system's performance on a larger and more diverse dataset and to address any potential limitations.

Result:-

Algorithm	Accuracy Reported
Support Vector Machine	96.27%
Artificial Neural Network	94.80%
Decision Tree (DT)	95.60%
Random Forest (RF)	99.20%
Convolutional Neural Network (CNN)	98.40%
Naive Bayes (NB)	94.80%
K-Nearest Neighbors (KNN)	93.80%
Extreme Learning Machine (ELM)	97.30%
Genetic Programming (GP)	96.80%

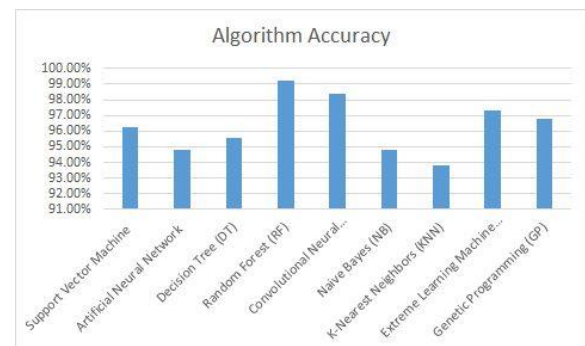


Table 2: shows the features contained in the dataset.

No	Attribute Name	Value Type	Clarification
1	id	number	1,2,3.....,
12	age	number	1,10,20,50,.....,
3	gender	1,0	1=m,0=f
4	query_thyroxine	1,0	1=yes,0=no
5	on_antithyroid_medication	1,0	1=yes,0=no
6	sick	1,0	1=yes,0=no
7	pregnant	1,0	1=yes,0=no
8	thyroid_surgery	1,0	1=yes,0=no
9	query_hypothyroid	1,0	1=yes,0=no
10	query_hyperthyroid	1,0	1=yes,0=no
11	TSH measured	1,0	1=yes,0=no
12	TSH	Analysis ratio	Numeric value
13	T3 measured	1,0	1=yes,0=no
14	T3	Analysis ratio	Numeric value
15	T4 measured	1,0	1=yes,0=no
16	T4	Analysis ratio	Numeric value
17	category	0,1,2	0=normal,1=hypothyroid,2=hyperthyroid

Conclusion

The prevalence of thyroid disease has been rising significantly in recent times, highlighting the need for effective automatic prediction models for its detection. Our study focuses on the classification of thyroid disease, specifically distinguishing between hyperthyroidism and hypothyroidism, utilizing algorithms and machine learning techniques.

In our research, we employed the Random Forest algorithm to train our dataset and enhance the accuracy of thyroid disease prediction. By training the machine with relevant data, we aimed to determine whether an individual is normal, hyperthyroid, or hypothyroid based on user input. Through the integration of this model into a web application, users can input their data, which is then processed in the backend, and the resulting prediction is displayed on the screen. Our objective was to provide society with an efficient and precise machine learning approach that can be utilized in various applications focusing on disease detection. By leveraging the power of machine learning, we aimed to offer an effective solution to aid in thyroid disease detection and contribute to advancements in healthcare technology.

References:-

- [1] Khalid Salman and Emrullah sonac (2021) Thyroid Disease Classification Using Machine Learning Algorithm
- [2] Rajasekhar Chaganti, Furqan Rustam, Isabel De La Torre Díez, Juan Luis Vidal Mazón, Carmen Lili Rodríguez and Imran Ashraf (2021) Thyroid Disease Prediction Using Selective Features and Machine Learning Techniques.
- [3] Saima Sharleen Islam¹, Md. Samiul Haque¹, M. Saef Ullah Miah², Talha Bin Sarwar² and Ramdhan Nugraha (2021) Application of machine learning algorithms to predict the thyroid disease risk: an experimental comparative study.
- [4] Ritesh, Jha. Vandana Bhattacharjee · Abhijit Mustaf (2020) Increasing the Prediction Accuracy for Thyroid Disease: A Step Towards Better Health for Society.
- [5] Mario Luca Bernardia, Marta Cimitileb, Martina Iammarinoa, Paolo Emidio Macchiac, Immacolata Cristina Nettorec, Chiara Verdon (2020) Thyroid Disease Treatment prediction with machine learning approaches Lerina Aversanoa.
- [6] Ankita Tyagi and Ritika Mehra. (2018) Interactive Thyroid Disease Prediction System using Machine Learning Techniques
- [7] Yong Feng Wang, (2020) Comparison Study of Radiomics and Deep-Learning Based Methods for Thyroid Nodules Classification using Ultrasound Images.
- [8] Sunila Godara, Prediction of Thyroid Disease Using Machine Learning Techniques (2018).

- [9] Hitesh Garg, (2013) Segmentation of Thyroid Gland in Ultrasound image using Neural Network.
- [10] Ehrenstein V, Nielsen H, Pedersen AB, Johnsen SP, Pedersen L. (2017) Clinical epidemiology in the era of big data: new opportunities, familiar challenges. Clin Epidemiol.
- [11] S. Godara and R. Singh, (2016) Evaluation of Predictive Machine Learning Techniques as Expert Systems in Medical Diagnosis
- [12] A. Begum and A. Parkavi. Prediction of thyroid disease using data mining techniques (2019).
- [13] Mario Luca Bernardi, Marta Cimitile, Fabio Martinelli, and Francesco Mercaldo (2019). Keystroke analysis for user identification using deep neural networks.