# Thyroid Cancer Detection Using Machine Learning Framework

**UTKARSH GUPTA[1], ANANT CHAUHAN[2, *] SIMRAN GOEL[3,] SHOBHIT KUMAR ROY[4]**

**ROHIT KUMAR GUPTA [5] LAVEENA SEHGAL[6]**

[1] IIMT COLLEGE OF ENGINEERING GREATER NOIDA
[2] IIMT COLLEGE OF ENGINEERING GREATER NOIDA
[3] IIMT COLLEGE OF ENGINEERING GREATER NOIDA
[4] IIMT COLLEGE OF ENGINEERING GREATER NOIDA
[5] IIMT COLLEGE OF ENGINEERING GREATER NOIDA
[6] IIMT COLLEGE OF ENGINEERING GREATER NOIDA

*

**Abstract:**

Thyroid illnesses are induced because of dissimilarity underway of chemicals - TSH, T4 and T3. The greater part of the patients of thyroid brokenness go untreated because of late identification or no recognition by any means. AI based models for discovery of thyroid illnesses offer a huge help to medical services. The clinical history of the patient supplies the highlights expected by AI based characterization and expectation models for thyroid brokenness. The point of this examination paper is to gain a characterization model in view of AI methods for appraisal of euthyroidism, hyperthyroidism, and hypothyroidism among guys, females, and youngsters. Different AI arrangement calculations, for example, innocent bayes, choice tree, arbitrary woodland and strategic relapse are utilized for order of genuine information. The exactness of every one of the strategies has laid out utilizing measurements like accuracy, review, particularity and responsiveness. A thyroid dataset has been recovered from two emergency clinics in Haryana from January 2020 to July 2020 to prepare the proposed model. The dataset involves clinical history of 539 thyroid patients including youngsters, men, and ladies of different ages. Out of 539 patients screened, 163 have sporadic TSH, 138 have commonness of raised TSH with 376 having negligible TSH height. Little youngsters prevalently experience more when contrasted with different layers of patients from Thyroid problems. This paper orders thyroid infection and looks at AI calculation among one another for prescient thyroid sickness and tracks down the best precision among them.

**Watchwords: Machine learning, Thyroid Dysfunction, Classification, Hyperthyroidism, Hypothyroidism.**

## 1 Introduction

The thyroid organ is a basic organ of humanoid which is situated in the throat under distinction of thyroid ligament additionally named the Adam's apple. The organ of thyroid deliveries two fundamental chemicals created from the iodine of the food ate - Triiodothyroinine (T3) and Thyroxine (T4) which partakes in managing body's digestion [12]. Too high or too low creation of T3 and T4 chemicals might cause serious problems. Thyroid Simulating Hormone (TSH) is created by nerve center, an organ in mind, which signals pituitary organ of the cerebrum to regulate the development of T4,T3 chemicals. Whenever T4, T3 levels are least and creation should be expanded, endocrine organ delivers more TSH. At the point when T4,T3 levels are most elevated the endocrine organ dials back the creation by providing less measure of TSH to thyroid organ.

Anomaly in thyroid organ's capacities or construction brings about clinical consideration. The side effects trait to expanded or diminished centralizations of plasma of chemicals, is known as hyperthyroidism and hypothyroidism individually. Cardiovascular issues include a super thyroid condition, hypertension, hyper cholesterol levels, clinical despondency, and barrenness. [6] The second-biggest illness in the endocrine world that might prompt a patient is endocrine thyroid organ issue situated on the cervix [4].

Symptomatic disclosures in regards to thyroid infection are of central issue for clinical science. Alongside the operations to decide the thyroid infection, AI likewise assumes a critical part. AI empowers probabilistic models to upgrade the exactness of illness assurance. Productivity of AI based arrangement models is straightforwardly relative to the meticulousness of its preparation utilizing preparing informational collection. Clinical history of the patient including actual details (Age, Weight and so forth), indicative record (tension, hustling heart, weight reduction, swelling eyes and so on) and symptomatic research center results (blood test report, ultrasound report and so on) are usually utilized as elements for preparing the model. [5] Trained an AI based grouping model has following factors as elements - Serum, Age, Name, Total thyroxine (T4),Total triiodothyronine (T3), TSH (fourth Generation) from a dataset of 539 harmless thyroid knobs.

## 2 Related Work

Jayatilake et al surveyed and looked at AI instruments for infection forecast in medical care. The article covers various AI procedures for dynamic which fall in one of various classes of AI for example support learning, solo learning and managed learning. Paper further stresses and organizations the reality with huge factual confirmations that discovery of sickness at beginning phase reduces the therapy upward and enhances the recuperation rate in different basic infections like bosom malignant growth, cellular breakdown in the lungs, heart illnesses, diabetes and so on. Indeed, even in COVID-19 pandemic circumstance, man-made consciousness assumed an indispensable part at different fronts of COVID the board.

Kashyap et al dominatingly concentrated on AI based expectation and characterization models for thyroid illnesses. Creators address a similar understanding of exactness of Artificial Neural Network (ANN) based models and choice tree based models utilizing disarray lattice. ANN gives preferred exactness over choice tree based classifiers.

Thyroid illness aftermaths in enormous impact on wellbeing in the event that not distinguished ahead of time [Min Hu et al]. In Japan just 450,000 patients experiencing thyroid are seeking treatment against 2.4 million all out patients. The grouping model in view of AI proposed by Min Hu et al has fair execution with hypothyroidism when contrasted with hyperthyroidism due to going through levothyroxine treatment for patients with hypothyroidism. Levothyroxine treatment might influence the normal research center qualities.

Priyanka et al directed a review among youthful females from country as well as from metropolitan regions in Bangalore to concentrate on thyroid brokenness. The clinical history of ladies matured 18 to 30 years was procured from two medical clinics - Boring and Lady Cruzon Hospital. Creators utilized IBM SPSS programming to group the information gathered from medical clinics and presumed that young ladies are more inclined to experience the ill effects of thyroid brokenness.

Table 1 notices huge work done by a few specialists in the field of ML based infection forecast for thyroid brokenness.

### 2.1 Problem Descriptions and Justification

In light of the investigation of writing in the field of AI in thyroid infection forecast, conceivable ends can be -

(1) The degree of work done toward this path is fundamentally less.

(2) Several patients experiencing thyroid dysfunctions might go untreated because of late or no identification. In this manner, early identification of thyroid illness lightens the weight of treatment alongside secured recuperation.

(3) Machine learning based discovery models for thyroid illness are unmistakable additional items in medical care.

The above focal points from the investigation of past work in a similar field drives to propose a characterization model for thyroid brokenness in view of AI calculations - arbitrary woods, strategic relapse, innocent bays, support vector machine, k-closest neighbor, and choice tree.

### 3 Thyroid Detection Using Machine Learning

The AI conveyed in medical care for the most part experiences predispositions and blunders. The goal is to stay away from it. AI models are not customized to do a specific assignment but rather it is intended to figure out how to make it happen. Accordingly, a versatile dataset for preparing and testing is required. Contingent upon the preparation information, the AI calculation conveyed in the model will do grouping. The total course of learning based model for grouping at unique level is as per the following:

1. Clinical history information of the patient is maintained.

2. The information is handled to adjust it into usable organization and unsought things are pruned.

3. Significant elements are chosen which will partake in arrangement.

4. The model is tried on preparing information.

5. Prepared model is tried on testing information and precision is assessed.

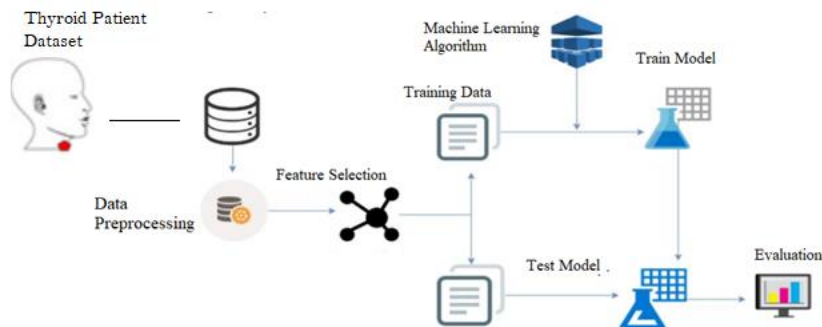6. Until the precision meets predefined edge, stage 4 and 5 are rehashed.



**Figure 1:** Abstract Machine Learning Based Model

The model is unique since it doesn't uncover the calculation utilized for characterization. The following are a few ML based calculations for arrangement, for example,

3.1 Logistic Regression

Strategic relapse is important for Machine Learning arrangement calculations used to characterize infection utilizing likelihood expectation of all out target variable. The likelihood is determined as follows:

$$P(y_i=C_j \mid X) = 1/(1+e^{-(\beta_0+\beta_1 x)})$$

Where, $P(y_i=C_j \mid X)$ is the likelihood of I-th perception having a place with target class $C_j$, $y_i$ is the objective worth of I-th perception, X is the framework, everything being equal, $\beta_0$ and $\beta_1$ are learning boundaries and e is Euler's number. $\beta_0+\beta_1 x$ is otherwise called Sigmoid Function.

3.2 Decision Tree Algorithm

Choice tree is a characterization strategy which is essentially addressed as a parallel tree. The cycle begins at the root hub. Each level compares to one specific element. The leaf hubs contain the choice (Class) [25]. The reason for a Decision Tree is to construct a preparation model that can use to anticipate the worth of the objective variable or class by gaining basic choice principles surmised from earlier data(training data)[26].
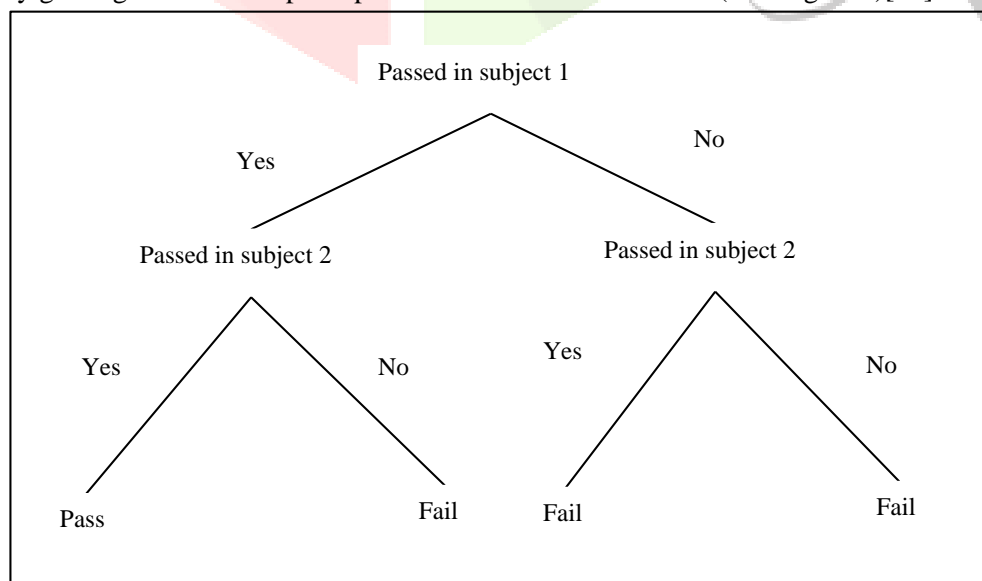


**Figure 2:** Example of decision tree

### 3.3 Naïve Bayes

Innocent Bayes classifier is reasonable for gigantic dataset with high dimensionality [S Vijayarani]. It is a probabilistic model for order.

The likelihood of information record x having the class name C_i is:

P (C_i | x) = (P (x | C_i) P (C_i))/(P(x))

In the given information dataset the class mark C_i with biggest contingent likelihood esteem decides the class of the arrangement. [27]

### 3.4 Random Forest Algorithm

An arbitrary timberland calculation based classifier which aggregates a few choice trees worked over sub-tests of the dataset and the choice lays on the normal of every choice tree. The exactness is improved altogether and over-fitting of information is controlled. It is a directed learning strategy and woodland alludes to assortment of trees (choice trees). Irregular backwoods has huge use in the field of medical care. [2].

### 4 Proposed Algorithm

Grouping process is completed on the highlights separated from clinical history of patient put in a dataset. The dataset is for the most part not unfit for direct use in grouping model. A few qualities in the dataset can be loud, deficient, questionable or missing for certain highlights. Thusly, information cleaning is an essential advance in AI calculations. Result of information cleaning stage is cleaned dataset of its highlights bear amazing qualities yet every one of them may not be needed for arrangement process. Just most huge and adequate highlights are quit from the dataset which take an interest in characterization process. The dataset is additionally parted into two classifications - with treatment and without treatment. The information of patients who went through treatment is placed in with treatment setprepared over it. The model is tried and prepared over without treatment dataset moreover.

**Figure 3:** Framework of machine learning

**Proposed Classifier Algorithm:**

➢ Load the required library files.
➢ Load Dataset
➢ Clean Dataset by removing missing values
➢ Select best features for classification
➢ The dataset is split into two partitions viz. patients with treatment and patient without treatment.
➢ Effect of Treatment is calculated as follows: EOT=E[Y_i (1)-Y_i (0)]
➢ Where, Y_i (1) represents patients with treatment,
➢ Y_i (0) Represents patients without treatment.
➢ Fit() Function is used for training data into model.
➢ Entropy is calculated as follows: $E = -\sum_{i=1} N p_i log_2 p_i$

Where, p_i represents proportion of patients which belong to i-th class.

The above algorithm is developed as a machine learning based model. In the experiments, it is established that the proposed algorithm has training accuracy near to 94%.

### 5 Material And Methods

### 5.1 Description Of Dataset

Rather than utilizing manufactured dataset accessible at different internet based sources like Kaggle, in this paper genuine dataset thyroid turmoil is procured from two emergency clinics in Haryana from January 2020 to July 2020 to prepare the proposed model. The dataset involves clinical history of 539 thyroid patients included youngsters, men, and ladies of different ages. Out of 539 patients screened, 163 are having unpredictable TSH, 138 are having predominance of raised TSH with 376 having negligible TSH rise. TSH levels were low in 25 of the members in the review. Table 2 gives Meta data about the dataset and subordinate/autonomous factors. Age, TSH, T3, TT4, and T4U were consistent autonomous factors in the thyroid dataset, which had 8 free factors. Table 3 address chosen highlights for our model.**Table 2:** Representation of thyroid dataset variable

**Table 3:** Features and domain of their values

| Feature | Values |
|---------|--------|
| Name | Nominal |
| Age | Continuous |
| Sex | 0,1 |
| Married | Continuous |
| TSH | Continuous |
| T3 | Continuous |
| T4 | Continuous |
| Type | 1,2,3 |
| Class | 3,2,1 |

```
th.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 417 entries, 0 to 416
Data columns (total 9 columns):
Name       417 non-null object
Married    417 non-null int64
Age        417 non-null int64
Sex        417 non-null int64
T3         417 non-null float64
T4         417 non-null float64
TSH        417 non-null float64
Type       417 non-null object
Class      417 non-null int64
dtypes: float64(3), int64(4), object(2)
memory usage: 29.4+ KB
```

**Figure 4:** Dataset attributes

## 5.2 Inclusion and Exclusion Criteria

Men, women, and children who were able to participate without medical complications were included in the analysis. Men, women, and children who were severely sick were not permitted to partake.

## 5.3 Performance Measure

The anticipated class of an item fall into different classifications - False Positive (FP), Negative (FN), True

Positive (TP) and True Negative (TN). Table 4 Represents a disarray grid address the comprehension of TP, TN, FP

and FN. Bogus positive suggests that the patient doesn't have thyroid infection, yet has thyroid disease. A misleading

negative implies that the patient is probably not going to have thyroid illness and doesn't have thyroid sickness.

Genuine positive suggests that the patient has

**Table 4:** Confusion matrix

| | | Actual Class | |
|---|---|---|---|
| | | Has thyroid disease | Does not have thyroid disease |
| Predicted Class | Has thyroid disease | *True Positive* | *False Positive* |
| | Does not have thyroid disease | *False Negative* | *True Negative* |

*Precision: Accuracy of characterization is the principal assumption for from the model. Precision just indicates the degree of right evaluations out of complete appraisals. Condition ___ indicates exactness.*

*Accuracy=(Correct appraisals)/(Total Assessments)=(TP+TN)/(TP+TN+FP+FN) (__)*

*Accuracy: Precision is the extent of positive remarks anticipated to be precisely anticipated is the proportion of positive remarks to add up to anticipated. Accuracy (a.k.a. positive anticipated esteem) is an amount of the extent of patients recognized by classifier to have the illness, really have had the infection. Definitively, 'the likelihood that a patient being really wiped out whenever analyzed as debilitated by the classifier'.*

*Accuracy=TP/(TP+FP)For example, if precision=1, it means that all patients diagnosed by the classifier really had the disease.*

*Review (Sensitivity): Recall is the extent of accurately anticipated positive perceptions out of complete positive perceptions. As a rule, it focuses the finger towards the positive cases delegated negative (False Negative) cases by the classifier. Unequivocally, 'Likelihood of a wiped out individual being recognized as debilitated by the classifier'.*

*Recall=TP/(TP+FN)*

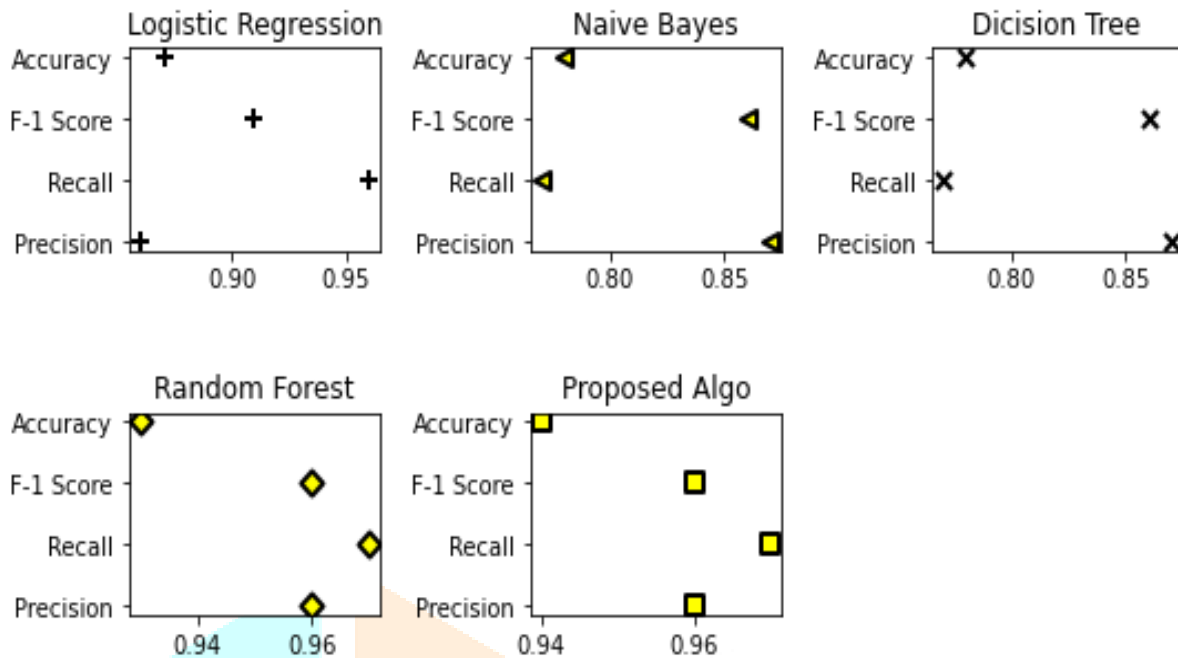*Particularity: Similar to review, 'Likelihood of a solid individual being recognized as sound by the classifier'. .*

*Specificity=TN/(TN+FP)*

*F1 score: The consonant mean of Precision and Recall is F1-score.*

*F1 Score=2/(1/(Precision (P))+1/(Recall (R)))=2\*(P\*R)/(P+R)6.*

**Result and Discussion**

The proposed computation is executed using python 3 (libraries used Skylearn, pandas, Numpy). Survey, expressness, exactness ,accuracy, and F1 scores are used as execution estimations to differentiate capability of proposed computation and other minor request estimations referred to in region 4. This assessment organizes these ailments into appropriate classes like hypothyroidism, hyperthyroidism, and usually established on the value of T4, T3, and TSH. The degree of dataset used for testing and getting ready are 20% and 80% independently. This degree applies to the two characterizations of dataset.

Tree, Naïve Bayes, and Logistic Reg

**Figure 5:** Performance metrics of various classifiers
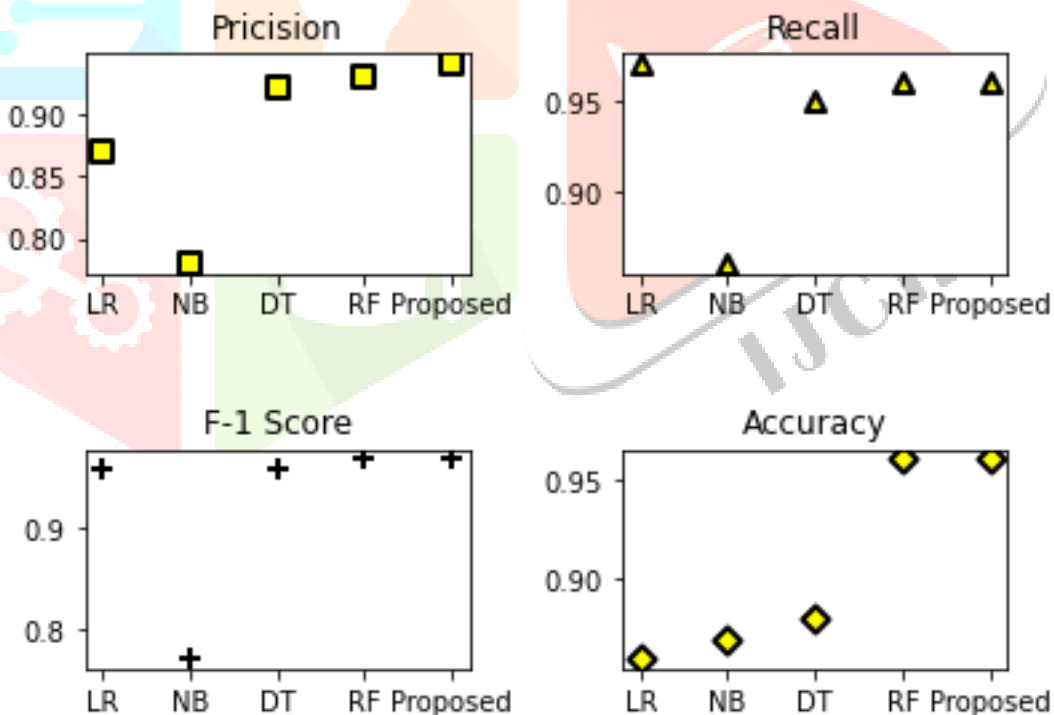


**Figure 6:** presents a performance of various classifiers projected through performance metrics.

It is clear from the above figure that the proposed calculation has accuracy practically 86% which is superior to Decision Tree, Logistic Regression and Naïve bayes and is equivalent to Random Forest.

Review addresses the likelihood of an item being placed in the legitimate class. It is clear from the figure 6 that the proposed calculation is superior to Random Forest Classifier, Naïve Bayes, Logistic Regression, and equivalents Decision Tree.

F-1 score of the proposed calculation is 0.96 which is superior to Random Forest Classifier, Naïve Bayes, Logistic Regression, and is equivalent to that of Decision Tree. The proposed calculation beats each of the four customary classifiers as far as accuracy. By the above investigation it is laid out that the proposed calculation effectively

arranges the patients to their right classes with the precision of 94%.

## Conclusion

Thyroid sickness analyze fundamentally spins around three chemicals T3, T4 and TSH. Scattered equilibrium of these chemicals brings about different sorts of thyroid sicknesses. Alongside indicative information comparing to T3, T4 and TSH, the clinical history of the patient which contains suggestive information is additionally utilized by proposed calculation for demonstrative grouping. The exhibition of proposed calculation is contrasted and deeply grounded customary classifiers like Random Forest Classifier, Logistic Regression, Naïve Bayes and Decision Tree. The calculation has been tried on the genuine information gathered from different clinics and through overviews. Less yet critical elements are picked for grouping by the proposed calculation so the patient doesn't need to go for too many check-ups. Accuracy, review, F-1 score and Accuracy are utilized as execution measurements. The proposed calculation beats LR, NB and DT classifier in accuracy (0.96). Review (0.97) of proposed calculation is superior to LR, NB and DT classifier. F-1 score (0.96) of the proposed calculation is superior to LR, NB and DT classifier. The proposed calculation exactness is 94% which is improved than each of the four customary calculations.

**Conflicts of Interest:** The authors declare that they have no conflicts of interest to report regarding the present study.

## References

1. Suresh Kumar Kashyap, Dr. Neelam Sahu A Comparative Study Of Machine Learning Based Model For Thyroid Disease Prediction International Journal of Creative Research Thoughts (IJCRT)

2. Mohammad Savargiv ,Behrooz Masoumi ,(2021) A New Random Forest Algorithm Based on Learning Automata Hindawi Computational Intelligence and Neuroscience Volume 2021, Article ID 5572781, 19 pages https://doi.org/10.1155/2021/5572781

3. Min Hu, Chikashi Asami 2021 Thyroid dysfunction diagnosis from routine laboratory tests based on machine learning doi: https://doi.org/10.1101/2021.03.30.21254605;

4. Liyong Ma , 1 Chengkuan Ma,1 Yuejun Liu,2 and Xuguang Wang Thyroid Diagnosis from SPECT Images Using Convolutional Neural Network with Optimization Volume 2019, Article ID 6212759, 11 pages https://doi.org/10.1155/2019/6212759

5. Jordi L. Reverter , 1,2 Irene Rosas-Allende,3 Carlos Puig-Jove,1,2 Carles Zafon,2,4 Ana Megia,2,5 Ignasi Castells,2,6 Eduarda Pizarro,2,7 Manel Puig-Domingo,1,2 and M. Luisa Granada2 Prognostic Significance of Thyroglobulin Antibodies in Differentiated Thyroid Cancer Volume 2020, Article ID 8312628, 6 pages https://doi.org/10.1155/2020/8312628

6. Gyanendra Chaubey1, Dhananjay Bisen1, Siddharth Arjaria1 ,Vibhash Yadav Thyroid Disease Prediction Using Machine Learning Approaches DOI: 10.1007/s40009-020-00979-z

7. Farzad PAKDEL 1 , Roghayeh GHAZAVI 2 (2019) Effect of Selenium on Thyroid Disorders: Scientometric Analysis Iran J Public Health, Vol. 48, No.3, Mar 2019, pp.410-420

8. Barbara L Parry*, and Daniel F Kripke (2020) Antidepressant and Stabilizing Effects of Thyroid Hormone Augmentation in Women's Mood Disorders Volume 5, Issue 3,March 2020

9. Shaik Razia, P. Swathi Prathyusha, N. Vamsi Krishna, N. Sathya Sumana (2018) A Comparative study of machine learning algorithms on thyroid disease prediction International Journal of Engineering & Technolog vol 7 issue 2 .8 (2018) pp 315-319

10. Rahul Sindhwani, Punj Lata Singh (2019) Agile System in Health Care: Literature Review Springer Nature Singapore Pte Ltd. 2019 , https://doi.org/10.1007/978-981-13-6412-9_61 pp 643-652.

11. Saurabh Bilgaiyan, Santwana Sagnika 2017 A Systematic Review on Software Cost Estimation in Agile Software Development, Journal Of Engineering Science And Technology Review, pp 51-64 , DOI: 10.25103/jestr.104.08

12. Dhyan Chandra Yadav1 , Saurabh Pal 2020 Discovery of Hidden Pattern in Thyroid Disease by Machine Learning Algorithms Indian Journal of Public Health Research and Development DOI Number: 10.37506/v11/i1/2020/ijphrd/193785 pp 61-66

13. Vaishnavi Kannan, Mujeeb A Basit 2019 User stories as lightweight requirements for agile clinical decision support development Journal of the American Medical Informatics Association, Volume 26, Issue 11, November 2019, Pages 1344–1354, https://doi.org/10.1093/jamia/ocz123

14. Agile Estimation Techniques: A True Estimation In An Agile Project 2020 https://www.softwaretestinghelp.com/agile-estimation-techniques

15. Murad Ali, Zubair A Shaikh , Eaman Ali 2015 Estimation of Project Size Using User Stories DOI: 10.2991/racs-15.2016.9 International Conference on Recent Advances in Computer Systems (RACS 2015)

16. Andrew L. Beam 2018  Big Data and Machine Learning in Health Care  American Medical Association doi:10.1001/jama.2017.18391

17. Charan R, Akash Yadav M, Aprameya N Katti (2019) Prediction of Thyroid Disease Based on Classification Using Hierarchical Structure Journal of Emerging Technologies and Innovative Research (JETIR) ISSN-2349-5162

18. Michael T. Biochemical Testing of the Thyroid: TSH is the Best and, Oftentimes, Only Test Needed – A Review for Primary Care doi:10.3121/cmr.2016.1309 6 Marshfield Clinic Health System

19. Jamot and Pettersson (2016). Agile challenges within regulated healthcare environments,  Journal of Advanced Nursing, 48(5), 454–462

20. Grantina Modern (2020) Correlates of diarrhea and stunting among under-five children in Ruvuma, Tanzania; a hospital-based cross-sectional study Scientific African https://doi.org/10.1016/j.sciaf.2020.e00430 2