1 **Fuzzy-Based Fusion Model for β-Thalassemia Carriers Prediction**

2 **Using Machine Learning Technique**

3 Muhammad Ibrahim[1], Sagheer Abbas,[2] Areej Fatima,[3] Taher M. Ghazal,[4,5] Muhammad Saleem,[6]
4 Meshal Alharbi,[7] Fahad Mazaed Alotaibi,[8] Muhammad Adnan Khan,[9,10] Muhammad Waqas[1]
5 and Nouh Elmitwally[11,12]

6 [1]School of Computer Science, National College of Business Administration & Economics,
7 Lahore, 54000, Pakistan.
8 [2]Department of Computer Science, Bahria University Lahore Campus, Lahore, 54000,
9 Pakistan
10 [3]Department of Computer Science, Lahore Garrison University, Lahore, Pakistan
11 [4]Centre for Cyber Physical Systems, Computer Science Department, Khalifa University
12 [5]Center for Cyber Security, Faculty of Information Science and Technology, UKM, 43600
13 Bangi, Selangor, Malaysia
14 [6]School of Computer Science, Minhaj University Lahore, Pakistan
15 [7]Department of Computer Science, College of Computer Engineering and Sciences, Prince
16 Sattam Bin Abdulaziz University, Alkharj 11942, Saudi Arabia.
17 [8]Faculty of Computing and Information Technology in Rabigh (FCITR), King Abdulaziz
18 University, Jeddah, Saudi Arabia
19 [9]School of Information Technology, Skyline University College, University City Sharjah,
20 Sharjah, 1797, UAE
21 [10]Riphah School of Computing and Innovation, Faculty of Computing, Riphah International
22 University, Lahore Campus, Lahore 54000, Pakistan
23 [11]School of Computing and Digital Technology, Birmingham City University,
24 Birmingham B4 7XG, UK
25 [12]Department of Computer Science, Faculty of Computers and Artificial Intelligence, Cairo
26 University, Giza 12613, Egypt

27 Correspondence should be addressed to Sagheer Abbas; jamsagheer@gmail.com

28 **Abstract**

29 The abnormality of haemoglobin in the human body is the fundamental cause of thalassemia
30 disease. Thalassemia is considered a common genetic blood condition that has received
31 extensive investigation in medical research globally. Likely, inherited disorders will be passed
32 down to children from their parents. If both parents are beta Thalassemia carriers, 25% of their
33 children will have intermediate or major beta thalassemia, which is fatal. An efficient method
34 of beta thalassemia is prenatal screening after couples have received counselling. Identifying
35 Thalassemia carriers involves a costly, time-consuming, and specialized test using quantifiable
36 blood features. However, cost-effective and speedy screening methods must be developed to
37 address this issue. The demise rate due to thalassemia development is outstandingly high
38 around the globe. The passing rate due to thalassemia development can be reduced by
39 following the proper procedure early; otherwise, it significantly impacts the body. A machine
40 learning-based late fusion model proposes the detection of beta-thalassemia carriers by
41 analyzing red blood cells. This study applied the late fusion technique to employ four machine
42 learning algorithms. For identifying the beta thalassemia carriers, Logistics Regression, Naïve
43 Bayes, Decision Tree, and Neural Network, they have achieved an accuracy of 94.01%,

93.15%, 97.93%, and 98.07%, respectively, by using the features-based dataset. The late fusion-based ML model achieved an overall accuracy of 96% for detecting beta-thalassemia carriers. The proposed late fusion model performs better than previously published approaches regarding efficiency, reliability, and precision.

**Keywords:** Machine Learning (ML), Logistics Regression (LR), Naïve Bayes (NB), Decision Tree (DT), Neural Network (NN), Fuzzy Logic (FL), Internet of medical things (IoMT), Late Fusion model.

## Introduction

Thalassemia comes from the Greek terms 'Thalassa' and 'Haima.' "Thalassa" means "the ocean," and "Haima" means "the blood". Thalassemia is a genetic blood disorder characterized by insufficient production of haemoglobin [1]. Haemoglobin plays a crucial role in the human body by transporting oxygen from the lungs to the rest of the body and returning carbon dioxide to the lungs [2]. Thalassemia is one of the world's most frequent diseases, particularly in the Mediterranean. Many countries are currently dealing with the high and rising incidence of thalassemia, which has become a primary public health concern—a significant source of disability and mortality around the world. Early detection of thalassemia can aid in the reduction of death rates. As a result, healthcare practitioners are responsible for making the right decisions. When distinguishing between ordinary people and patients, complete the following options. Who are carriers of diseases, especially when it comes to genetic disorders such as a condition known as thalassemia [3].

There are two divisions of thalassemia based on two polypeptide chains in haemoglobin. These are known as alpha-thalassemia (α) and beta-thalassemia (β). An abnormality causes alpha thalassemia in the alpha polypeptide gene of haemoglobin, whereas beta-thalassemia is caused due to disturbance in the beta polypeptide gene. The development of any of the alpha or beta-thalassemia in a person's body leads to low or abnormal haemoglobin creation in the body [4]. The red blood cells are affected due to inadequate haemoglobin [5].

The classification of thalassemia consists of three stages: major, intermediate, and minor thalassemia. Thalassemia major is the most crucial stage of the disease in which the patient needs a continuous blood transfusion to survive. Thalassemia intermediate is the middle stage of the condition in which the patient needs blood transfusion occasionally. It is also known as mild or moderate anaemia. The patient with thalassemia minor looks physically fit and healthy. They don't need a blood transfusion but maintain their diet and healthy lifestyle [6].

World Health Organization (WHO) identifies that beta-thalassemia has 5.1% carriers worldwide [7]. Many tests are required to diagnose the difference between Iron deficiency anaemia and beta-thalassemia. These tests include serum iron level, Complete blood count, High-performance liquid chromatography, the binding capacity of Iron, and the calculation of ferritin and HBA2. However, these tests are expensive and unavailable everywhere [8].

In many other research disciplines, machine learning approaches are very efficient in producing results. They make managing and analyzing other fields easier and play a significant role in the health sector. A computer-based system can be developed to identify thalassemia with improved accuracy, better results, and more affordable cost. Various machine learning algorithms have offered effective treatments for various biomedical problems. Many models have been presented to analyze the data of other diseases [32,33] like brain tumours [9], kidney diseases [10], lung disorders [11], and Iron deficiency anaemia by using machine learning techniques [25, 26 and 30], including support vector machine [12], K-nearest neighbour [13], fuzzy logic [14, 29 and 31], Deep extreme machine learning [27] and deep neural network [15, 24 and 28].

89 Logistic regression models a discrete outcome given an input variable. The most popular logistic
90 regression models a binary result (true/false, yes/no, etc.). When analyzing a classification problem,
91 logistic regression is a helpful analysis tool.
92 Nave Bayes is a superficial learning algorithm that uses the Bayes rule and assumes attributes are
93 class-dependent. Due to its processing efficiency and other benefits, nave Bayes is commonly used
94 in practice [21].
95 A tree has numerous analogies in real life and has inspired machine learning, classification, and
96 regression. A decision tree can represent the decision analysis process visually and explicitly.
97 Feature-based data can be handled very effectively by neural network algorithms. Neural networks
98 are computing systems inspired by human brain neural networks [2, 10].
99 Although machine learning algorithms are currently helpful for identifying illnesses, earlier research
100 models were less accurate because they mainly concentrated on preprocessing methods, data
101 balancing, and other supervised and semi-supervised learning models. A late fusion technique is
102 needed to fuse the accuracy of many machine learning algorithms while maintaining high sickness
103 detection accuracy. This study proposed a late fusion-based ML model that implements Logistics
104 Regression, Naïve Bayes, Decision Tree, and Neural Network for data analysis. The system will use
105 a featured-based dataset of thalassemia reports.
106 It highlights the importance of accurately predicting β-thalassemia carriers to enable early
107 intervention and genetic counselling. The limitations of existing prediction models and the need for
108 an improved approach are discussed. The objectives of the paper are clearly stated as follows:
109     1. To develop a fuzzy-based fusion model that combines multiple machine learning algorithms
110        for β-thalassemia carrier prediction.
111     2. To evaluate the performance of the proposed model using relevant performance metrics and
112        compare it with existing approaches.
113     3. To analyze the effectiveness of fuzzy logic in improving the accuracy and reliability of β-
114        thalassemia carrier prediction.

115 The results of four different machine learning algorithms were combined through fuzzification
116 to decide on beta-thalassemia carrier identification. The outcomes demonstrate that the
117 proposed approach is more precise and effective than existing solutions.

## Related Work

119 The goal of the research is to identify thalassemia sickness early. Hirimutugoda and Wijayarathna
120 [2] implemented a three-layer artificial neural network to detect and differentiate malaria and
121 thalassemia. Both diseases are life-threatening and global health issues. Visual inspection of the
122 images of blood analysis taken with a light microscope is a well-known technique for determining
123 malaria and thalassemia. This technique takes much time and is more consuming and expensive. The
124 model used three and four layers of ANN that merged with methods of image analysis to find the
125 accuracy and effectiveness of the classification for identifying the images related to morphological
126 features of the blood erythrocytes. The study claimed that the three-layered ANN approach generated
127 results with an accuracy of 84.54%.
128 Ayyıldız and Tuncer [5] performed a decision-based diagnosis to identify and discriminate the Iron
129 deficiency anemia (IDA) and beta-thalassemia (β). They implemented red blood cell indices and two
130 effective techniques of machine learning: support vector machine and k-nearest neighbour. Various
131 parameters of Complete blood count were used to differentiate between IDA and β–Thalassemia.
132 Implementation of RBC indices improved the efficiency of the diagnostic model. But if the number
133 of features increases, the system becomes complicated.

Das et al. [16] employed a decision-based system that used decision trees, ANN, and a Naïve Bayes classifier to discriminate β-thalassemia carriers from ordinary people. The Postgraduate Institute of Medical Education and Research in the Indian city of Chandigarh is where the dataset was gathered. Both ratings were determined to be completely sensitive. The screening score for thalassemia characteristics (BTT) was determined to be 79.25 percent and 91.74 percent, respectively, for the combined score of BTT and HbE. Although the mechanism differentiates two main variants related to haemoglobin, it still requires validation with datasets collected from different countries for implementation and unification.

Egejuru et al. [17] implemented a prediction model for identifying the risk of thalassemia disease. The model used supervised machine learning approaches for analyzing the data collected through questionnaires and medical persons. The Waikato Environment for Knowledge Analysis (WEKA) tool was used for data simulation. Identification variables included demographics (age, marital status, gender, social class, and ethnicity) and clinical variables (spleen enlargement, family history, urine colour, diabetes, and inherited disease status). The dataset consisted of 57% disease carriers and 43% non-carriers. The models implemented in the study are multi-layer perception (MLP) and the Naïve Bayes classifier. The study results show that the MLP is a more effective and reliable mechanism for identifying the risk of thalassemia in patients in Nigeria.

Sadiq et al. [18] constructed an ensemble classifier model using a random forest support vector machine and a Gradient boosting machine to identify patients with thalassemia from the Complete blood count (CBC) test data. The model was implemented on the dataset of CBC reports of 5066 patients collected from the Punjab thalassemia prevention program (PTPP). Input parameters used for this study are red blood cells, Haemoglobin, Hematocrit, Mean cell volume, Mean cell haemoglobin concentration, Mean cell haemoglobin, RBC distribution width, platelet count, and white blood cells. The study achieved an accuracy of 93% in identifying β-thalassemia carriers.

Akhtar et al. [19] implemented a Linear discrimination analysis (LDA) classifier to classify the patients with thalassemia using various parameters of a Complete blood count report. The parameters used in the study are Ferritin, HB, RBC, WBC, HCT, and Platelets. The study also used mathematical formulas to discriminate the patients with thalassemia and iron deficiency anaemia. The accuracy achieved 78% results for females and 75% for males.

The fuzzy-based model was developed to classify thalassemia diseases by Susanto et al. [20]. The haemoglobin, MCV, and MCH levels were obtained following the CBC examination to determine the type of thalassemia. Major, Intermedia, Minor, and Not Thalassemia are four output models. The doctor's perspectives on thalassemia were contrasted with the model prediction results against four datasets. Additional data must be used to understand to further test the model's accuracy.

Jahan et al. [21] investigated the research on red cell indices utilizing machine learning techniques, such as an artificial neural network (ANN), to detect Beta-thalassemia traits (BTT) in pregnant women. The optimal cutoff for each index and the BTT detection test characteristics was determined using a Receiver operating characteristic (ROC) curve analysis. The C4.5 and Naive Bayes (NB) classifiers and a back-propagation type ANN were constructed and tested over 3947 patients using the red cell indices. The study emphasizes that none of the red cell features alone helps detect BTT. However, ANN, with a mixture of all red cell indices, exhibited good sensitivity and specificity for this use. Further neural network development might produce a valuable tool for thalassemia screening in remote areas.

Mohammed and Al-Tuwaijari [22] presented various artificial intelligence-based methods and machine learning techniques for classifying and detecting thalassemia utilizing CBC test variables such as RBC, HGB, MCV, HTC, and HB. This system was developed to identify patients with minor thalassemia alpha and major thalassemia beta. The classification methods are decision tree, Naive Bayes, and support vector machine.

4

182 Tyas et al. [23] examined multi-layer perceptron to classify erythrocytes present in thalassemia cases.
183 It combined morphological features with texture and colour features to increase the accuracy of
184 erythrocyte classification. The experimental results of 7108 erythrocytes indicated an accuracy of
185 98.11% for training and 93.77% for testing based on the combination of features. The system's
186 effectiveness was assessed using images captured at various magnifications and on different scanning
187 platforms. The least number of red cells to image for analysis was determined using Poisson
188 modelling, and the results were validated using image sets. Table 1 is showing the comparative
189 analysis with respect to the accuracy of past works that were about anomaly detection in network
190 security.

191 **Table 1:** Previous Work Accuracy and Dataset Status

| Author | Method | Dataset | Accuracy |
|---|---|---|---|
| Hirimutugoda and Wijayarathna [2] | Artificial Neural Network | Public | 86.54% |
| Sadiq et al. [18] | Random Forest | Private | 91% |
| Sadiq et al. [18] | Support Vector Machine | Private | 90% |
| Sadiq et al. [18] | Gradient Boosting Machine | Private | 91% |
| Susanto et al. [20] | Fuzzy Inference System | Public | 89.26% |
| Jahan et al. [21] | Artificial Neural Network | Private | 85.95% |
| Tyas et al. [23] | Convolutional Neural Network | Public | 93.77% |

192
193 The Aims and objectives of the paper are:
194 To highlight the importance of identifying beta-thalassemia carriers and their impact on reducing the
195 mortality rate associated with the disease.
196 To identify the limitations of current screening methods and propose developing cost-effective and
197 speedy screening methods.
198 To develop a machine learning-based late fusion model for detecting beta-thalassemia carriers by
199 analyzing red blood cells.
200 To compare the proposed late fusion model's accuracy, efficiency, reliability, and precision with
201 previously published approaches.
202 To explore the use of machine learning in medical research to detect beta-thalassemia carriers.
203 To evaluate the performance of four machine learning algorithms, including Logistics Regression,
204 Naïve Bayes, Decision Tree, and Neural Network.
205 To use a features-based dataset for the development of the late fusion model.

## Proposed β-Thalassemia Prediction Model

207 The Late fusion model based on machine learning is proposed for predicting β-thalassemia carriers.
208 The system used a features-based dataset of thalassemia reports obtained from the Internet of Medical
209 Things (IoMT) enabled devices. The novel features dataset was collected from the Punjab
210 Thalassemia Prevention Program (PTPP) database. Table 2 presents a complete overview of the
211 features.

212 **Table 2:** Dataset structure

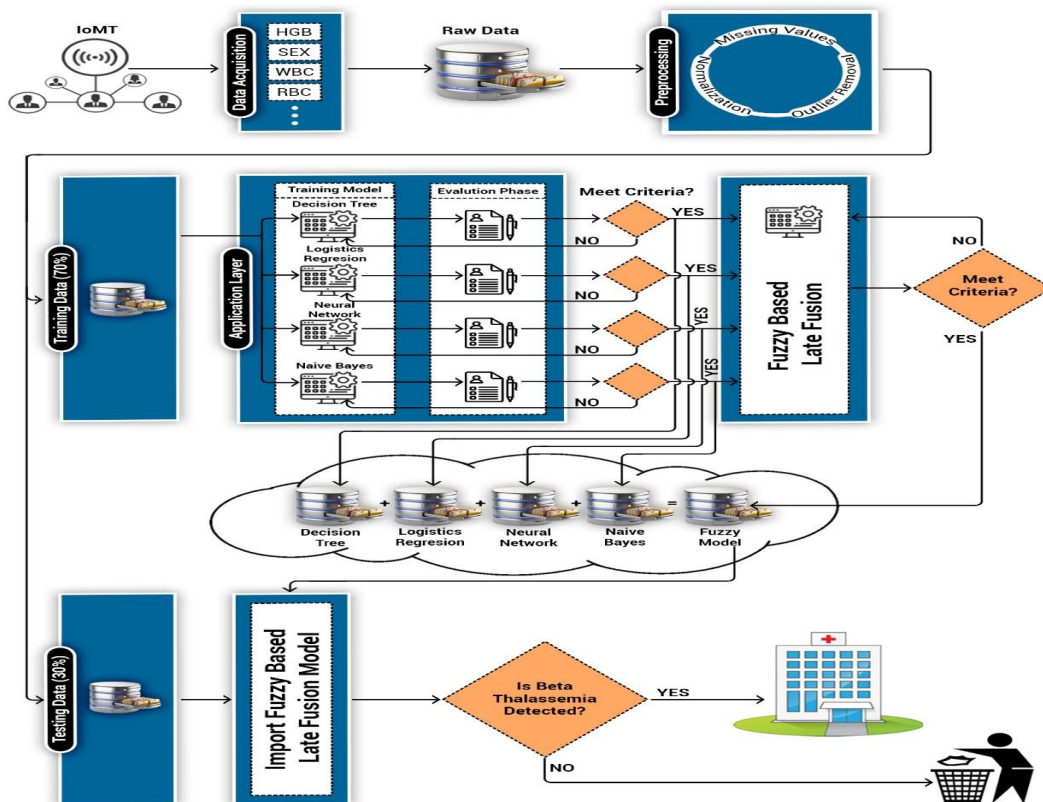| Sr. | Features | Datatype | Sr. | Features | Datatype |
|---|---|---|---|---|---|
| 1. | CS/PS | Integer | 5. | MCH | Integer |
| 2. | ETC | Nominal | 6. | Hb | Integer |
| 3. | Age | Integer | 7. | Hct | Integer |
| 4. | Sex | Integer | 8. | MCV | Integer |

213

214
215 PTPP is the initiative of the Punjab government of Pakistan to protect the people from thalassemia
216 disease. This platform provides support to thalassemia patients in β-Thalassemia carrier screening.
217 Initially, the dataset is divided into training and testing phases. 70% of records were fixed for training
218 and 30% for testing.
219 The Proposed model consists of training and validation phases. The proposed model consists of
220 various layers that help to diagnose beta thalassemia disease. These layers are data acquisition,
221 preprocessing, and application. The proposed model's first layer is the data acquisition layer, which
222 collects the dataset from PTPP based on IoMT devices [18]. It consisted of twelve variables and a
223 total number of 5066 instants. Output is classified into two categories. The first is β-Thalassemia
224 Non-Carriers, which contains 3051 records, and the second is β-Thalassemia Carriers, with 2015
225 patient records. The sex distribution ratio is 53% for males and 47% for females.
226 This unprocessed data may have some missing or noisy values. Normalization of the data and
227 treatment of missing values is accomplished in the preprocessing layer. The normalizing method is
228 used to handle noisy data. In contrast, missing values are driven by calculating existing values' mean
229 and moving averages.
230 In the training phase, the third layer of the model is the application layer, which predicts thalassemia
231 sickness using four different machine learning algorithms: Logistics Regression (LR), Naïve Bayes
232 (NB), Decision Tree (DT), and Neural Network (NN).
233 The LR, NB, DT, and NN results are given to the evaluation phase, which calculates the accuracy. It
234 misrates in the targeted class represented by [0, 1], where 0 is for β-Thalassemia non-carrier, and 1
235 is for β-Thalassemia Carrier investigated. The data is sent to the cloud if the learning criteria are
236 satisfied. Otherwise, it needs to be retrained, as shown in Figure 1.



237

**Figure 1:** Proposed Late Fusion Model for Thalassemia Disease Prediction
239 The results of four different techniques are combined in the following stage using a fuzzy inference

240     system to increase the performance of the suggested beta thalassemia carrier's model.

241     The validation phase utilized the 30% records of the thalassemia dataset to validate the model. The

242     trained fusion-based model is imported from the cloud to predict thalassemia. The model discards

243     the value if a beta-thalassemia non-carrier is found. If a beta-thalassemia carrier is found, the patient

244     is referred to the hospital for additional treatment, as shown in Figure 1.

245     The following conditions (if-then rules) are employed in the fuzzy logic of the suggested late fusion

246     model, which is written as follows:

247     The late fusion-based rules identify beta-thalassemia carriers.

248
$$\mu_{LR \cap NB \cap DT \cap NN}(l, n, d, n\,) = \min[\mu_{LR}(l), \mu_{NB}(n), \mu_{DT}(d), \mu_{NN}(n)]$$

249     $\boldsymbol{Rule_{bt}^{1}}$= IF (LR is carrier and NB is carrier and DT is carrier and NN is carrier) THEN (Thalassemia is Beta
250     Carrier)

251     $\boldsymbol{Rule_{bt}^{2}}$= IF (LR is carrier and NB is carrier and DT is carrier and NN is Non-carrier) THEN (Thalassemia is
252     Beta Carrier)

253     $\boldsymbol{Rule_{bt}^{3}}$= IF (LR is carrier and NB is carrier and DT is Non-carrier and NN is carrier) THEN (Thalassemia is
254     Beta Carrier)

255     $\boldsymbol{Rule_{bt}^{4}}$= IF (LR is carrier and NB is carrier and DT is Non-carrier and NN is Non-carrier) THEN (Thalassemia
256     is Beta Carrier)

257     $\boldsymbol{Rule_{bt}^{5}}$= IF (LR is carrier and NB is Non-carrier and DT is carrier and NN is carrier) THEN (Thalassemia is
258     Beta Carrier)

259     $\boldsymbol{Rule_{bt}^{6}}$= IF (LR is carrier and NB is Non-carrier and DT is carrier and NN is Non-carrier) THEN (Thalassemia
260     is Beta Carrier)

261     $\boldsymbol{Rule_{bt}^{7}}$= IF (LR is carrier and NB is Non-carrier and DT is Non-carrier and NN is carrier) THEN (Thalassemia
262     is Beta Carrier)

263     $\boldsymbol{Rule_{bt}^{8}}$= IF (LR is carrier and NB is Non-carrier and DT is Non-carrier and NN is Non-carrier) THEN
264     (Thalassemia is Beta Non-Carrier)

265     $\boldsymbol{Rule_{bt}^{9}}$= IF (LR is Non-carrier and NB is carrier and DT is carrier and NN is carrier) THEN (Thalassemia is
266     Beta Carrier)

267     $\boldsymbol{Rule_{bt}^{10}}$= IF (LR is Non-carrier and NB is carrier and DT is carrier and NN is Non-carrier) THEN (Thalassemia
268     is Beta Non-Carrier)

269     $\boldsymbol{Rule_{bt}^{11}}$= IF (LR is Non-carrier and NB is carrier and DT is Non-carrier and NN is carrier) THEN (Thalassemia
270     is Beta Non-Carrier)

271     $\boldsymbol{Rule_{bt}^{12}}$= IF (LR is Non-carrier and NB is carrier and DT is Non-carrier and NN is Non-carrier) THEN
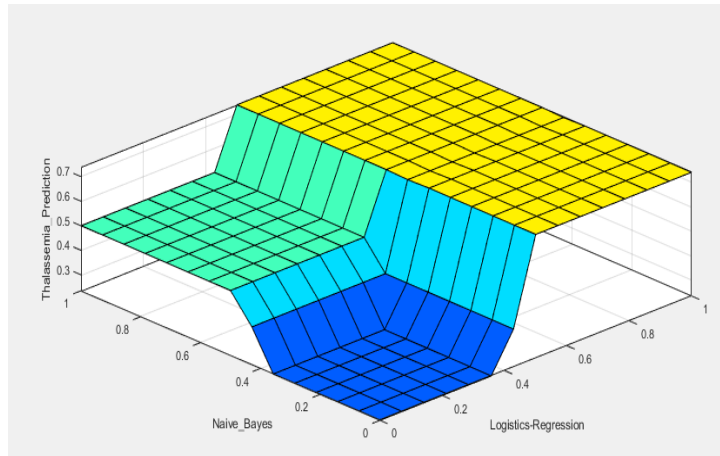272     (Thalassemia is Beta Non-Carrier)

273     $\boldsymbol{Rule_{bt}^{13}}$= IF (LR is Non-carrier and NB is Non-carrier and DT is carrier and NN is carrier) THEN (Thalassemia
274     is Beta Non-Carrier)

275     $\boldsymbol{Rule_{bt}^{14}}$= IF (LR is Non-carrier and NB is Non-carrier and DT is carrier and NN is Non-carrier) THEN
276     (Thalassemia is Beta Non-Carrier)

277     $\boldsymbol{Rule_{bt}^{15}}$= IF (LR is Non-carrier and NB is Non-carrier and DT is Non-carrier and NN is carrier) THEN
278     (Thalassemia is Beta Non-Carrier)

279     $\boldsymbol{Rule_{bt}^{16}}$= IF (LR is Non-carrier and NB is Non-carrier and DT is Non-carrier and NN is Non-carrier) THEN
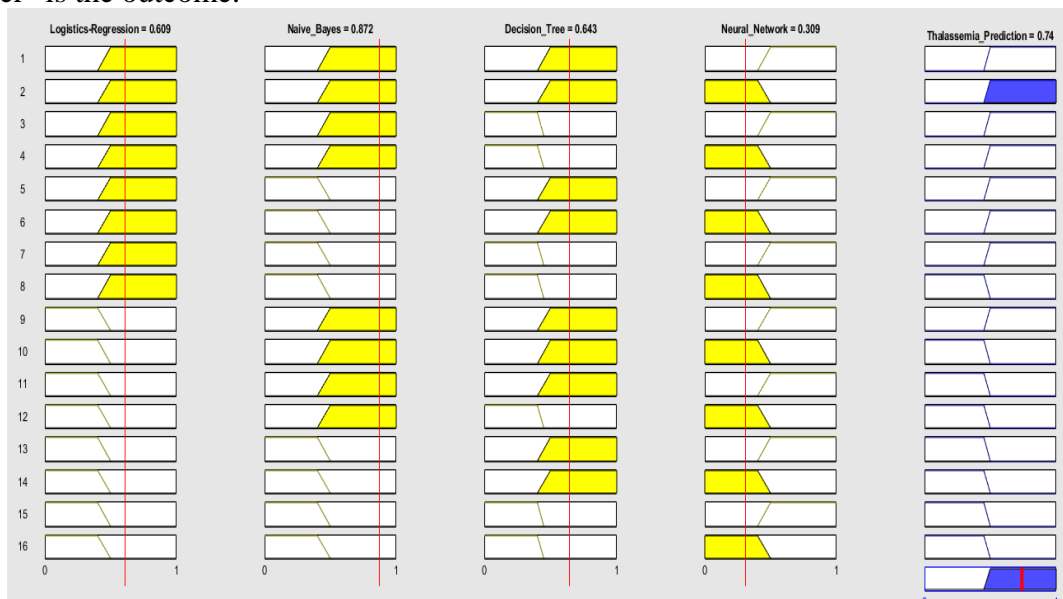280     (Thalassemia is Beta Non-Carrier)

281 The generated fuzzy rules show that the suggested late fusion-based technique will predict the

282 optimal result based on at least three classification strategies (either the beta-thalassemia carrier

283 or beta-thalassemia non-carrier).

284
285 **Figure 2:** Proposed Late Fusion Rule Surface for NB and LR
286 The proposed late fusion technique of rule surface for predicting beta-thalassemia carriers
287 based on NB and LR is shown in Figure 2. If both classification methods indicate that "beta
288 thalassemia = carrier" is the outcome, then the suggested technique will also mean that "beta
289 thalassemia = carrier" is the outcome. If both methods indicate that "beta thalassemia = non-
290 carrier" is the outcome, then the proposed technique will suggest that "beta thalassemia = non-
291 carrier" is the outcome.
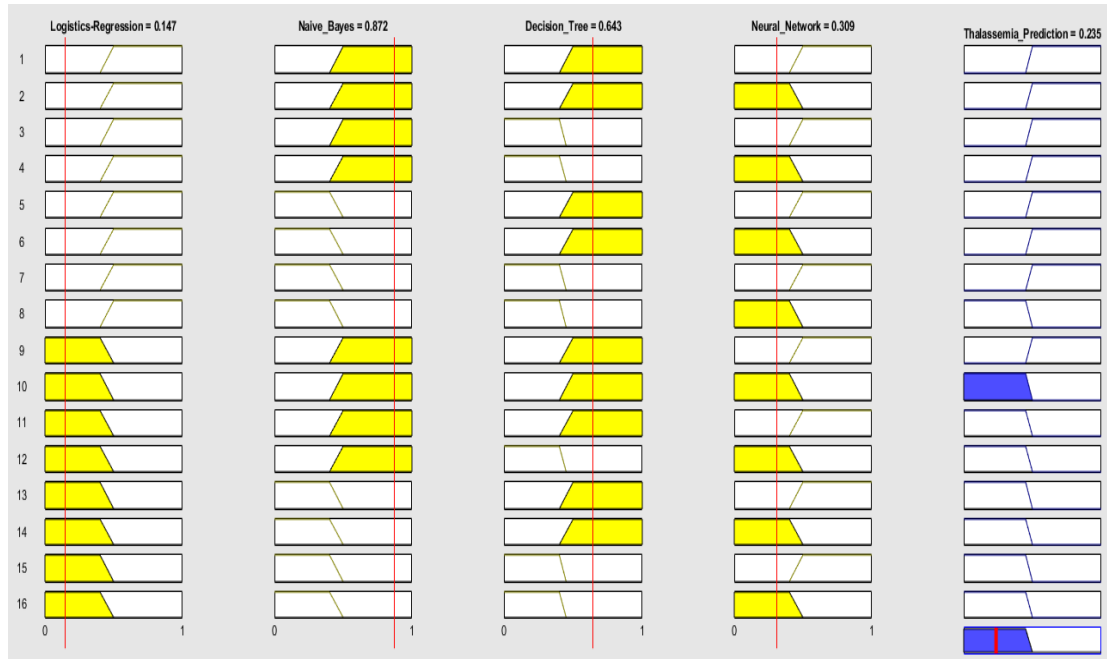


292
293 **Figure 3:** Result of Proposed Late Fusion Model Beta Thalassemia Carrier
294 Figure 3 demonstrates that the suggested late fusion technique will also predict "beta
295 thalassemia = carrier" if NB, DT, and NN make this prediction 'beta thalassemia = carrier'.
296

**Figure 4:** Result of Proposed Late Fusion Model Beta Thalassemia Non-Carrier

Figure 4 shows that if LR and NB show "beta thalassemia = non-carrier," even if DT and NN show "beta thalassemia = carrier," the proposed method will still show "beta thalassemia = non-carrier."

**Table 3:** Fuzzy-based Graphical and Mathematical Membership Function Representation

| Sr No. | Input / Output Variables | Membership Functions (MF) | Graphical Representation of MF |
|---|---|---|---|
| 1 | LR = $\mu_{LR}((lr))$ | $\mu_{LR,c}(lr) = \left\{ max\left( min\left(1, \frac{0.5 - lr}{0.1}\right), 0\right)\right\}$ <br><br> $\mu_{LR,nc}(lr) = \left\{ max\left( min\left(\frac{lr - 0.4}{0.1}, 1\right), 0\right)\right\}$ |  |
| 2 | NB = $\mu_{NB}((nb))$ | $\mu_{NB,c}(nb) = \left\{ max\left( min\left(1, \frac{0.5 - nb}{0.1}\right), 0\right)\right\}$ <br><br> $\mu_{NB,nc}(nb) = \left\{ max\left( min\left(\frac{nb - 0.4}{0.1}, 1\right), 0\right)\right\}$ |  |
| 3 | DT= $\mu_{DT}((dt))$ | $\mu_{DT,c}(dt) = \left\{ max\left( min\left(1, \frac{0.5 - dt}{0.1}\right), 0\right)\right\}$ <br><br> $\mu_{DT,nc}(dt) = \left\{ max\left( min\left(\frac{dt - 0.4}{0.1}, 1\right), 0\right)\right\}$ |  |
| 4 | NN = $\mu_{NN}((nn))$ | $\mu_{NN,c}(nn) = \left\{ max\left( min\left(1, \frac{0.5 - nn}{0.1}\right), 0\right)\right\}$ <br><br> $\mu_{NN,nc}(nn) = \left\{ max\left( min\left(\frac{nn - 0.4}{0.1}, 1\right), 0\right)\right\}$ |  |
| 5 | Beta-Thalassemia = $\mu_{BT}((bt))$ | $\mu_{BT,c}(bt) = \left\{ max\left( min\left(1, \frac{0.5 - bt}{0.05}\right), 0\right)\right\}$ <br><br> $\mu_{BT,nc}(bt) = \left\{ max\left( min\left(\frac{bt - 0.45}{0.05}, 1\right), 0\right)\right\}$ |  |

9

304 Table 3 shows membership functions based on fuzzy rules. The system testing layer predicts
305 beta-thalassemia carriers. A fuzzy-based cloud model is used to achieve an outcome that
306 stores real-time patient data for evaluation.

## Results and Simulation

308 The late fusion-based model is proposed for the earliest prediction of beta-thalassemia carriers. The
309 results are obtained using MATLAB tool 2022. The proposed model comprises four machine learning
310 techniques, LR, NB, DT, and NN are applied to 5066 features. For both methods, 30% of the fused
311 samples were utilized for validation, while the remaining 70% were used for training. The proposed
312 model diagnoses the beta thalassemia carrier and beta thalassemia non-carrier. The statistical metrics
313 used to evaluate the suggested late fusion model's predicted effectiveness and other categorization
314 methods are explained below. βTc represents beta thalassemia true predicted, βTnc represents beta
315 thalassemia false predicted, βFnc represents beta thalassemia non-carrier false expected, and βFc
316 means expected false beta thalassemia carrier.

317
$$\text{Accuracy} = \frac{\beta Tc + \beta Tnc}{\beta Fc + \beta Fnc + \beta Tc + \beta Tnc} \tag{1}$$

318 Accuracy is the number of correctly labelled cases out of the total number of cases.

319
$$\text{Misrate} = \frac{\beta Fc + \beta Fnc}{\beta Fc + \beta Fnc + \beta Tc + \beta Tnc} \tag{2}$$

320 The percentage of real positives and negatives missed during an experiment is known as the
321 miss rate.

322
$$\text{Sensitivity} = \frac{\beta Tc}{\beta Tc + \beta Fnc} \tag{3}$$

323 Sensitivity measures the capacity of the proposed model to identify positive cases.

324
$$\text{Specificity} = \frac{\beta Tnc}{\beta Tnc + \beta Fc} \tag{4}$$

325
$$\text{Positive Predication Value} = \frac{\beta Tc}{\beta Tc + \beta Fc} \tag{5}$$

326
$$\text{Negative Predication Value} = \frac{\beta Tnc}{\beta Fnc + \beta Tnc} \tag{6}$$

327 Predictive values, positive and negative, are calculated by dividing each set of results by the
328 proportion of actual successes and failures.

329
$$\text{False Postive Ratio} = 1 - \frac{\beta Tnc}{\beta Tnc + \beta Fc} \tag{7}$$

330
$$\text{False Negative Ratio} = 1 - \frac{\beta Tc}{\beta Tc + \beta Fnc} \tag{8}$$

331
$$\text{Likelihood Ratio Positive} = \frac{\beta Tc}{\beta Tc + \beta Fnc} \bigg/ 1 - \frac{\beta Tnc}{\beta Tnc + \beta Fc} \tag{9}$$

332
$$\text{Likelihood Ratio Negative} = 1 - \frac{\beta Tc}{\beta Tc + \beta Fnc} \bigg/ \frac{\beta Tnc}{\beta Tnc + \beta Fc} \tag{10}$$

333 The dataset contains 5066 instances. 70% of the dataset is used for training which consists of
334 3,546 records, while the remaining 30% is used for testing, which consists of 1,520 records.
335 The 3546 records were used for training with the LR approach, in which 1715 were beta-
336 thalassemia non-carriers, and 1831 were beta-thalassemia carriers. When trained with LR, 1623
337 out of 1715 occurrences were non-carriers, while 1717 out of 1831 were found to be carriers.
338 Table 4 displays the results of a comparison between actual and predicted performance
339 throughout training. Results showed an accuracy of 94.2% with a miss rate of 5.8%.

340 **Table 4:** Proposed LR-based Training Confusion Matrix

| Total Samples = 3546 | O$_{BT\text{-Non-Carrier}}$ | O$_{BT\text{-Carrier}}$ |
|---|---|---|
| I$_{BT\text{-Non-Carrier}}$ = 1715 | 1623 | 92 |
| I$_{BT\text{-Carrier}}$ = 1831 | 114 | 1717 |

10

341 In contrast, during the testing of LR, 716 records out of 757 were identified as non-carriers,
342 while 713 records out of 763 were classified as carriers, as shown in Table 5. In LR testing, the
343 attained accuracy was 94.01%, and the miss rate of 5.99%.

**Table 5:** Proposed LR-based Testing Confusion Matrix

| Total Samples = 1520 | $O_{BT-Non-Carrier}$ | $O_{BT-Carrier}$ |
|---|---|---|
| $I_{BT-Non-Carrier = 757}$ | 716 | 41 |
| $I_{BT-Carrier = 763}$ | 50 | 713 |

345 The 3546 records were used for training with the NB approach, in which 1715 were beta-
346 thalassemia non-carrier, and 1831 were beta-thalassemia carriers. When trained with NB, 1618
347 out of 1715 occurrences were found to be non-carriers, while 1658 out of 1831 instances were
348 found to be carriers. Table 6 displays the results of a comparison between actual and predicted
349 performance throughout training. Results showed an accuracy of 92.4% with a miss rate of
350 7.6%.

**Table 6:** Proposed NB-based Training Confusion Matrix

| Total Samples = 3546 | $O_{BT-Non-Carrier}$ | $O_{BT-Carrier}$ |
|---|---|---|
| $I_{BT-Non-Carrier = 1715}$ | 1618 | 97 |
| $I_{BT-Carrier = 1831}$ | 173 | 1658 |

352 In contrast, during the testing of NB, 721 records out of 757 were identified as non-carriers,
353 while 695 records out of 763 were classified as carriers, as shown in Table 7. In NB testing,
354 the attained accuracy was 93.15% and a miss rate of 6.85%

**Table 7:** Proposed NB-based Testing Confusion Matrix

| Total Samples = 1520 | $O_{BT-Non-Carrier}$ | $O_{BT-Carrier}$ |
|---|---|---|
| $I_{BT-Non-Carrier = 757}$ | 721 | 36 |
| $I_{BT-Carrier = 763}$ | 68 | 695 |

356 The 3546 records were used for training with the DT approach, in which 1715 were beta-
357 thalassemia non-carrier, and 1831 were beta-thalassemia carriers. When trained with DT, 1703
358 out of 1715 occurrences were non-carriers, while 1813 out of 1831 were found to be carriers.
359 Table 8 displays the results of a comparison between actual and predicted performance
360 throughout training. Results showed an accuracy of 99.15% with a miss rate of 0.85%.

**Table 8:** Proposed DT-based Training Confusion Matrix

| Total Samples = 3546 | $O_{BT-Non-Carrier}$ | $O_{BT-Carrier}$ |
|---|---|---|
| $I_{BT-Non-Carrier = 1715}$ | 1703 | 12 |
| $I_{BT-Carrier = 1831}$ | 18 | 1813 |

362 In contrast, during the testing of DT, 756 records out of 757 were identified as non-carriers,
363 while 763 records out of 763 were classified as carriers, as shown in Table 9. In DT testing,
364 the attained accuracy was 99.93% , and a miss rate of 0.07%

**Table 9:** Proposed DT-based Testing Confusion Matrix

| Total Samples = 1520 | $O_{BT-Non-Carrier}$ | $O_{BT-Carrier}$ |
|---|---|---|
| $I_{BT-Non-Carrier = 757}$ | 756 | 1 |
| $I_{BT-Carrier = 763}$ | 0 | 763 |

366 The 3546 records were used for training with the NN approach, in which 1715 were beta-
367 thalassemia non-carrier, and 1831 were beta-thalassemia carriers. When trained with NN, 1700
368 out of 1715 occurrences were found to be non-carriers, while 1824 out of 1831 instances were
369 found to be carriers. Table 10 displays the results of a comparison between actual and predicted

370 performance throughout training. Results showed an accuracy of 99.4% with a miss rate of
371 0.6%.

**Table 10:** Proposed NN-based Training Confusion Matrix

| Total Samples = 3546 | $O_{BT\text{-}Non\text{-}Carrier}$ | $O_{BT\text{-}Carrier}$ |
|---|---|---|
| $I_{BT\text{-}Non\text{-}Carrier = 1715}$ | 1700 | 15 |
| $I_{BT\text{-}Carrier = 1831}$ | 7 | 1824 |

373 In contrast, during the testing of NN, 757 records out of 757 were identified as non-carriers,
374 while 763 records out of 763 were classified as carriers, as shown in Table 11. In NN testing,
375 the attained accuracy was 100%.

**Table 11:** Proposed NN-based Testing Confusion Matrix

| Total Samples = 1520 | $O_{BT\text{-}Non\text{-}Carrier}$ | $O_{BT\text{-}Carrier}$ |
|---|---|---|
| $I_{BT\text{-}Non\text{-}Carrier = 757}$ | 757 | 0 |
| $I_{BT\text{-}Carrier = 763}$ | 0 | 763 |

377 Table 12 displays detailed results for validation of all used classification machine learning
378 techniques (LR, NB, DT, and NN). It can be observed that all four machine learning
379 techniques performed well and achieved an average accuracy is 96.77% and misrate of
380 3.23%.

**Table 12:** ML-based Proposed Model Performance (Validation)

| Samples for validation (30% Records) | | |
|---|---|---|
| **Approaches** | **Accuracy** | **Miss Rate** |
| Logistics Regression (LR) | 94.01% | 5.99% |
| Naïve Bayes (NB) | 93.15% | 6.85% |
| Decision Tree (DT) | 99.93% | 0.07% |
| Neural Network (NN) | 100% | 0% |
| **Average Performance Proposed Model** | 96.77% | 3.23% |

382
383 Four machine learning techniques are finally provided to the fuzzy system as input for the final
384 prediction. Input to the fuzzy system consists of LR, NB, DT, and NN classifiers and the output
385 class Beta Thalassemia Carriers classifiers. By employing fuzzy rules, the suggested machine
386 learning late fusion-based fuzzy system attained an accuracy of 96% and a miss rate of 4%.
387 The fuzzy system randomly takes twenty-five input ranges for generating the fusion-based
388 results. Based on the fuzzy rules, 12 outputs show beta-thalassemia carriers, and 12 outcomes
389 non-carriers truly predicted. The remaining one is between the carrier and non-carrier stages
390 that, showed the system's error.

**Table 13:** Comparison of Previous Approaches with Proposed Late Fusion-based ML Model

| Author | Method | Accuracy | Misrate |
|---|---|---|---|
| Sadiq et al. [18] | Random Forest | 91% | 9% |
| Sadiq et al. [18] | Support Vector Machine | 90% | 10% |
| Sadiq et al. [18] | Gradient Boosting Machine | 91% | 9% |
| Susanto et al. [20] | Fuzzy Inference System | 89.26% | 10.74% |
| Jahan et al. [21] | Naïve Bayes | 82.49% | 17.51% |
| Tyas et al. [23] | Convolutional Neural Network | 93.77% | 6.23% |
| **Proposed Model** | **Late fusion-based ML Model** | **96%** | **4%** |

393

394 Table 13 displays the results of a comparison between the suggested fused machine learning
395 model and the various thalassemia illness prediction methods described in the literature. The
396 proposed late fusion model is compared with RF [18], SVM [18], GBM [18], FIS [20], NB
397 [21], and CNN [23]. Advanced methods are contrasted with the proposed late fusion model. In
398 comparison to the other methods, the proposed late fusion model excelled. The proposed fused
399 model outperformed the different approaches. The suggested machine learning fusion-based
400 system can be included in intelligent healthcare systems for early and accurate beta thalassemia
401 carrier prediction. The proposed model has shown the accuracy of beta thalassemia carrier
402 prediction is 96%.

## Conclusions

404 The critical point of this study is to develop a system to analyze beta-thalassemia carrier
405 patients using the late fusion-based ML model. This system is fundamental and more accessible
406 for medical experts and non-experts. Hence, any person can examine the status of thalassemia
407 just by feeding the required input data. The goal of this study is to analyze the various
408 dimensions of thalassemia. The total precision of this proposed late fusion-based ML model is
409 96%. The presented framework can be enhanced in the future by utilizing different methods,
410 including federated learning. The study can also be extended by applying Short-Term Long
411 Memory (LSTM) and other ML algorithms and diagnosing the other stages of Thalassemia like
412 Alpha Max and Min, Beta Max, and Min.

## Data Availability

414 Data will provide on demand.

## Conflicts of Interest

416 The authors declare no conflicts of interest to report regarding the present study.

## Funding Statement

418 The authors received no specific funding for this study.

## Acknowledgements

420 Thanks to Gapico PVT, who provide us simulation platform.

## References

422 1. Hossain, M. S., Hasan, M. M., Petrou, M., Telfer, P., & Al Mosabbir, A. (2021). The parental perspective
423 of thalassaemia in Bangladesh: lack of knowledge, regret, and barriers. Orphanet Journal of Rare
424 Diseases, 16(1), 1-10.
425 2. Hirimutugoda, Y. M., & Wijayarathna, G. (2010). Image analysis system for detection of red cell disorders
426 using artificial neural networks. Sri Lanka Journal of Bio-Medical Informatics, 1(1).
427 3. Q. Zhuang et al., "The value of combined detection of HbA2 and HbF for the screening of thalassemia
428 among individuals of childbearing ages," Zhonghua yi xue yi Chuan xue za zhi= Zhonghua Yixue
429 Yichuanxue Zazhi= Chinese J. Med. Genet., vol. 39, no. 1, pp. 16–20, 2022.
430 4. Z. Rustam, A. Kamalia, R. Hidayat, F. Subroto, and A. Suryansyah, "Comparison of Fuzzy C-Means, Fuzzy
431 Kernel C-Means, and Fuzzy Kernel Robust C-Means to Classify Thalassemia Data," Update, vol. 1, p. 1,
432 2019.

5. Ayyıldız, H., & Tuncer, S. A. (2020). Determination of the effect of red blood cell parameters in the discrimination of iron deficiency anemia and beta thalassemia via Neighborhood Component Analysis Feature Selection-Based machine learning. Chemometrics and Intelligent Laboratory Systems, 196, 103886.

6. Nabi, A. T., Muttu, J., Chhaparwal, A., Mukhopadhyay, A., Pattnaik, S. J., & Choudhary, P. (2022). Implications of β-thalassemia on oral health status in patients: A cross-sectional study. Journal of Family Medicine and Primary Care, 11(3), 1174.

7. Sari, D. P., Wahidiyat, P. A., Setianingsih, I., Timan, I. S., Gatot, D., & Kekalih, A. (2022). Hematological Parameters in Individuals with Beta Thalassemia Trait in South Sumatra, Indonesia. Anemia, 2022.

8. Anari, Shokofeh, Nazanin Tataei Sarshar, Negin Mahjoori, Shadi Dorosti, and Amirali Rezaie. "Review of deep learning approaches for thyroid cancer diagnosis." Mathematical Problems in Engineering 2022 (2022).

9. Ranjbarzadeh, R., Bagherian Kasgari, A., Jafarzadeh Ghoushchi, S., Anari, S., Naseri, M., & Bendechache, M. (2021). Brain tumor segmentation based on deep learning and an attention mechanism using MRI multi-modalities brain images. Scientific Reports, 11(1), 1-17.

10. V. A. Binson, M. Subramoniam, Y. Sunny, and L. Mathew, "Prediction of pulmonary diseases with electronic nose using SVM and XGBoost," IEEE Sens. J., vol. 21, no. 18, pp. 20886–20895, 2021.

11. Asif, M., Abbas, S., Khan, M. A., Fatima, A., Khan, M. A., & Lee, S. W. (2021). MapReduce based intelligent model for intrusion detection using machine learning technique. Journal of King Saud University-Computer and Information Sciences.

12. Islam, M. M., Rahman, M., Roy, D. C., Islam, M., Tawabunnahar, M., Ahmed, N. A. M., & Maniruzzaman, M. (2022). Risk factors identification and prediction of anemia among women in Bangladesh using machine learning techniques. Current Women's Health Reviews, 18(1), 118-133.

13. Haseli, G., Ranjbarzadeh, R., Hajiaghaei-Keshteli, M., Ghoushchi, S. J., Hasani, A., Deveci, M., & Ding, W. (2023). HECON: Weight assessment of the product loyalty criteria considering the customer decision's halo effect using the convolutional neural networks. Information Sciences, 623, 184-205..

14. Kollias, D., Tagaris, A., Stafylopatis, A., Kollias, S., & Tagaris, G. (2018). Deep neural architectures for prediction in healthcare. Complex & Intelligent Systems, 4(2), 119-131.

15. Das, R., Datta, S., Kaviraj, A., Sanyal, S. N., Nielsen, P., Nielsen, I., ... & Saha, S. (2020). A decision support scheme for beta thalassemia and HbE carrier screening. Journal of advanced research, 24, 183-190.

16. Egejuru, N. C., Olusanya, S. O., Asinobi, A. O., Adeyemi, O. J., Adebayo, V. O., & Idowu, P. A. (2019). Using data mining algorithms for thalassemia risk prediction. J Biomed Sci Eng, 7(2), 33-44.

17. Sadiq, S., Khalid, M. U., Ullah, S., Aslam, W., Mehmood, A., Choi, G. S., & On, B. W. (2021). Classification of β-Thalassemia Carriers From Red Blood Cell Indices Using Ensemble Classifier. IEEE access, 9, 45528-45538.

18. Akhtar, F., Shakeel, A., Li, J., Pei, Y., & Dang, Y. (2020, May). Risk Factors Selection for Predicting Thalassemia Patients using Linear Discriminant Analysis. In 2020 Prognostics and Health Management Conference (PHM-Besançon) (pp. 1-7). IEEE.

19. Susanto, E. R., Syarif, A., Muludi, K., Perdani, R. R. W., & Wantoro, A. (2021). Implementation of Fuzzy-based Model for Prediction of Thalassemia Diseases. In Journal of Physics: Conference Series (Vol. 1751, No. 1, p. 012034). IOP Publishing.

20. Jahan, A., Singh, G., Gupta, R., Sarin, N., & Singh, S. (2021). Role of red cell indices in screening for beta thalassemia trait: an assessment of the individual indices and application of machine learning algorithm. Indian Journal of Hematology and Blood Transfusion, 37(3), 453-457.

21. Mohammed, M. Q., & Al-Tuwaijari, J. M. (2021). A Survey on various Machine Learning Approaches for thalassemia detection and classification. Turkish Journal of Computer and Mathematics Education (TURCOMAT), 12(13), 7866-7871.

22. Tyas, D. A., Hartati, S., Harjoko, A., & Ratnaningsih, T. (2020). Morphological, Texture, and Color Feature Analysis for Erythrocyte Classification in Thalassemia Cases. IEEE Access, 8, 69849-69860.

23. Ihnaini, B., Khan, M. A., Khan, T. A., Abbas, S., Daoud, M. S., Ahmad, M., & Khan, M. A. (2021). A smart healthcare recommendation system for multidisciplinary diabetes patients with data fusion based on deep ensemble learning: computational Intelligence and Neuroscience, 2021.

24. Siddiqui, S. Y., Athar, A., Khan, M. A., Abbas, S., Saeed, Y., Khan, M. F., & Hussain, M. (2020). Modelling, simulation and optimization of diagnosis cardiovascular disease using computational intelligence approaches. Journal of Medical Imaging and Health Informatics, 10(5), 1005-1022.

25. Khan, M. A., Abbas, S., Atta, A., Ditta, A., Alquhayz, H., Khan, M. F., & Naqvi, R. A. (2020). Intelligent cloud based heart disease prediction system empowered with supervised machine learning.

26. Rehman, A., Athar, A., Khan, M. A., Abbas, S., Fatima, A., & Saeed, A. (2020). Modelling, simulation, and optimization of diabetes type II prediction using deep extreme learning machine. Journal of Ambient Intelligence and Smart Environments, 12(2), 125-138.

493   27. Ahmad, G., Alanazi, S., Alruwaili, M., Ahmad, F., Khan, M. A., Abbas, S., & Tabassum, N. (2021).
494        Intelligent ammunition detection and classification system using convolutional neural network.
495   28. Fatima, A., Adnan Khan, M., Abbas, S., Waqas, M., Anum, L., & Asif, M. (2019). Evaluation of planet
496        factors of smart city through multi-layer fuzzy logic (MFL). The ISC International Journal of Information
497        Security, 11(3), 51-58.
498   29. Muhammad, M. U. U. A. H., & Saleem, A. M. S. F. M. Intelligent Intrusion Detection System for Apache
499        Web Server Empowered with Machine Learning Approaches.
500   30. Naeem, Z., & Naeem, F. (2022). Predicting the performance of governance factor using fuzzy inference
501        system. International Journal of Computational and Innovative Sciences, 1(2), 35-50.
502   31. Muneer, S., & Rasool, M. A. (2022). AA systematic review: Explainable Artificial Intelligence (XAI) based
503        disease prediction. International Journal of Advanced Sciences and Computing, 1(1), 1-6.
504   32. Appiahene, P., Asare, J. W., Donkoh, E. T., Dimauro, G., & Maglietta, R. (2023). Detection of iron
505        deficiency anemia by medical images: a comparative study of machine learning algorithms. BioData
506        Mining, 16(1), 1-20.
507   33. Appiahene, P., Arthur, E. J., Korankye, S., Afrifa, S., Asare, J. W., & Donkoh, E. T. (2023). Detection of
508        anemia using conjunctiva images: A smartphone application approach. Medicine in Novel Technology and
509        Devices, 18, 100237.