

Machine Learning-Based Approaches for the Prognosis and Prediction of Multiple Diseases

Priya Bhardwaj

Department of CSE, Tula's Institute,
Dhoolkot, Dehradun, India
priyabharadwaj_21@yahoo.co.in

Yogesh Kumar

Department of CSE, School of Technology,
Pandit Deendayal Energy University
Gandhinagar, Gujarat, India
yogesh.kumar@sot.pdpu.ac.in

Shakti Mishra

Department of CSE, School of Technology,
Pandit Deendayal Energy University
Gandhinagar, Gujarat, India
Shakti.mishra@sot.pdpu.ac.in

Abstract: The rapid progress in machine learning techniques has significantly transformed healthcare which enables the simultaneous and accurate detection of multiple diseases. This paper delves into the application of diverse machine learning algorithms for multi-disease detection by using a comprehensive dataset which focuses on three diseases i.e. diabetes, gonorrhoea, and typhoid. The multi-disease dataset has been meticulously pre-processed and graphically visualized to discern patterns and represent diseases against emotional states/urges and critical feelings. Subsequently, a range of machine learning classifiers which includes logistic regression, Adaboost, random forest, support vector machine, CatBoost, Light Gradient Boosting Classifier, Naïve Bayes, XGBoost, KNN, and Decision Tree, are trained on this dataset. Their performance across these different classes is rigorously evaluated using various parameters such as accuracy, F1 score, recall, and precision. During execution, Adaboost emerged as the top performer, by achieving an impressive accuracy of 94.37% and maintaining a precision, recall, and F1 score of 0.94, which indicates its robustness in multi-disease detection.

Index Terms: Healthcare, Multidisease, Machine learning classifiers, Adaboost

I. INTRODUCTION

In these recent years, integrating healthcare with technology has revolutionized the way diseases are being diagnosed as well as treated. The most significant advancements that have been seen so far in this domain are the role of machine learning techniques for detecting diseases [1]. Traditional methods to diagnose diseases mostly rely on manual analysis of symptoms as well as medical history and it turns out to be subjective, time-consuming, as well as prone to errors. With the arrival of machine learning, healthcare professionals now are being blessed with the powerful tools at their disposal to enhance the efficiency and accuracy of disease detection [2].

Detecting multi-diseases using machine learning refers to the process of using advanced computational algorithms for the identification of multiple diseases or health conditions in patients. These algorithms influence large amount of medical data which includes patient history, symptoms, diagnostic tests, and genetic information, for learning complex patterns as well as associations [3]. Later, by processing and analyzing this data, machine learning models helps to provide valuable insights into the presence as well as severity of various diseases simultaneously [4].

Such transformative approach to healthcare has the ability to significantly improvise patient outcomes by enabling earlier as well as more accurate diagnoses. Detecting diseases

promptly is crucial, as it allows for timely medical intervention and customized treatment plans which ultimately enhances the quality of patient care and reduces the healthcare costs. Moreover, multi-disease detection using machine learning aids to identify risk factors and predicts the likelihood of developing certain conditions [5].

Based on this, the manuscript explores the multi-disease detection using multiple machine learning techniques and through a deep examination of our research findings, this exploration seeks to highlight the importance of machine learning techniques to shape the future of healthcare by making it more personalized, precise, and efficient.

II. RELATED WORK

Bansal et al. [6] used MobileNetV2, ResNet50, VGG16, VGG19, and DenseNet in deep transfer learning to analyze oral cancer and non-malignant images from histopathologic and real-time datasets. Gaussian blur preprocessing and optimization techniques like SGD, ADAM, and RMSprop achieved impressive accuracies: 95.41% for oral cancer, 92.41% for real-time data, and 92.41% for non-cancerous histopathologic images, demonstrating the efficacy of hybrid optimization methods. Bhardwaj et al. [7] predicted chest disorders using a dataset of 112,120 X-ray images from the National Institute of Health. They employed six transfer-learning approaches such as MobileNet V2, VGG-16, DenseNet-161, VGG-19, ResNet-50, and Inception V3. VGG-16 achieved the highest accuracy of 81%, with a recall rate of 90% and 85% as F1 Score. Kakkar et al. [8] delved into AI-based techniques for cancer detection and treatment, highlighting their potential in processing extensive datasets efficiently. Automated and computer-assisted approaches demonstrated accurate diagnoses and effective solutions. Swapna et al. [9] presented a novel approach employing LSTM and CNN in classifying diabetic and normal HRV signals, combined with SVM for precise categorization. This method outperformed prior work, enhancing CNN and CNN-LSTM models by 0.03% and 0.06%, achieving an accuracy of 95.7%. Bhuiyan et al. [10] showcased machine learning and deep learning's transformative role in health by introducing typhoid fever prediction model pre-clinical trials. Using ten algorithms, XGBoost emerged with 97.87% accuracy, highlighting its effectiveness in early disease prediction. Bao et al. [11] examined 21,273 Australian MSM (2011-2017), employing Gradient Boosting Machine (GBM) for HIV and STI prediction. GBM outperformed traditional methods,

accurately forecasting infections (HIV, syphilis, gonorrhea, chlamydia) in MSM.

III. PROPOSED METHODOLOGY

This section covers the flow, as shown in Figure 1 that has been proposed to develop the model for the detection of multi-diseases such as diabetes, gonorrhoea, and typhoid.

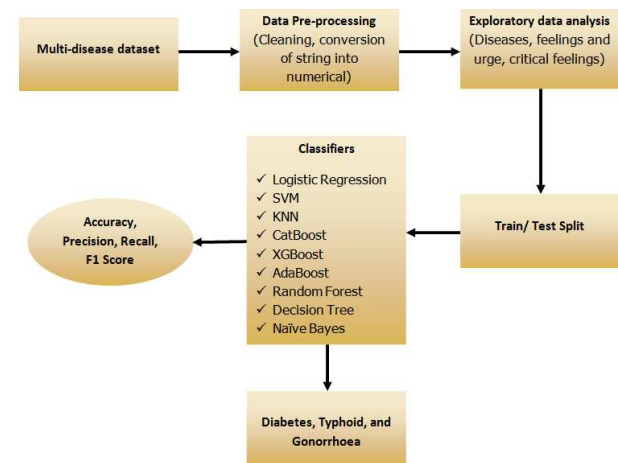


Fig 1: Proposed system to detect and classify multi-disease

A. Dataset: The collected dataset is particularly well-suited for multi-class disease detection, representing a valuable resource for supervised learning tasks. In this dataset, diseases serve as the labels to be predicted, with symptoms provided as input features. Comprising 7 columns and 4000 rows, the dataset offers a substantial volume of data for analysis. The 'disease' column stands out as the intuitive target variable, making it the focal point of prediction [12]. The dataset attributes are displayed in Table 1.

Table 1: Attributes of the multi-disease dataset

Attribute	Values
Discharge	“Frequent urination”, “Yellow, green discharge”, “difficulties in stooling”, “painful urination”
Feelings and Urge	“Tired”, “Hunger”, “urge to urinate”, “Fever”
Pains and Infection	“Blurred vision”, “Frequent infection”, “Swollen genitals”, “ Anal Itching and Pain”, “Abdominal pain”, “Muscle Aches”, “head ache”
Physical Conditions	“Swollen Testicles”, “Rose spot”, “Excessive Urination and thrust”, “Rashes”, “Bloody Diarrhea”
Critical Feelings	“Severe Pelvic Pain”, “Slow heart rate and pulse”, “Disorientation”, “Confusion”
Critical Feature	“Critical” “not Critical”
Diseases	“Typhoid”, “Diabetes”, “Gonorrhea”

B. Preprocessing: The data has undergone thorough cleaning and now requires the application of the LabelEncoder class to transform string-based labels into numeric representations. This step, constituting a vital aspect of data preprocessing, ensures that the dataset is in a suitable format for machine learning algorithms.

C. Exploratory data analysis: In this exploratory data analysis, with attention to 'feelings,' 'urges,' and 'critical feelings' columns box plots were used to illustrate the

distribution patterns across various diseases. This approach allowed for a comprehensive understanding of the emotional states associated with different conditions, shedding light on potential correlations and outliers. Figure 2 depicts the analysis of three diseases i.e. typhoid, diabetes, and gonorrhoea with respect to their feelings and urge.

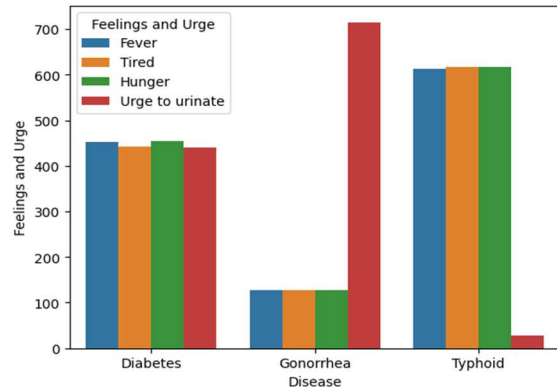


Fig 2: Analysis of feelings and Urge vs Disease

The data depicted in the figure reveals significant trends among individuals with different health conditions. In the case of diabetes, over 450 people exhibit symptoms of both fever and hunger, and a similar number experience tiredness and a frequent urge to urinate. Meanwhile, individuals with Gonorrhea primarily display a strong urge to urinate, affecting more than 700 people, and around 110 individuals also report symptoms of fever, tiredness, and hunger. Interestingly, those diagnosed with typhoid predominantly suffer from fever, tiredness, and hunger. However, only a small fraction, fewer than 50 people, experiences an urge to urinate.

Similarly, the analysis of the data for the people suffering from diabetes, gonorrhoea, and typhoid having the symptoms of critical feelings has been also presented in Figure 3. Many individuals with diabetes often experience no noticeable symptoms, making it challenging to detect. However, a small percentage of around 250 people may exhibit symptoms such as a slow heart rate, seizures, disorientation, confusion, and cough, along with severe pelvic pain and breathing difficulties.

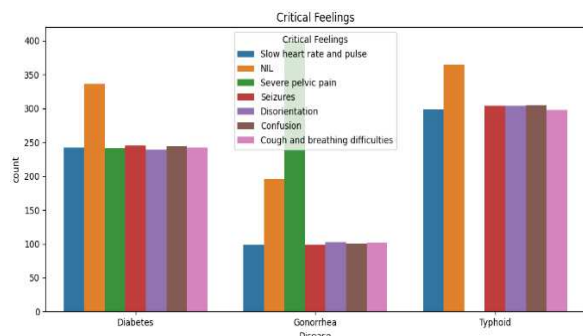


Fig 3: Analysis of critical feelings vs diseases

On the other hand, those affected by gonorrhoea commonly present with pelvic pain, a prominent symptom observed in approximately 400 cases. Additionally, about 100 individuals may show signs such as a slow heart rate, seizures, disorientation, confusion, and cough. In the case of typhoid, more than 300 people indicate the presence of symptoms such

as seizures, slow heart rate, confusion, disorientation, and cough. However, unlike diabetes, severe pelvic pain is not a symptom associated with typhoid.

D. Applied Machine Learning Classifiers: This section explores the role of various machine learning classifiers to detect and classify typhoid, diabetes, and Gonorrhoea disease. The methods that have been used are explained in brief.

Logistic Regression: This algorithm on dealing with multiclass data such as multi-disease detection, transforms into Multinomial Logistic Regression. This technique calculates the probabilities of each disease class based on the given set of input features. On including the softmax function, Multinomial Logistic Regression ensures that the sum of probabilities equals to 1 across all disease classes [13]. The general equation for Multinomial Logistic Regression to predict the probability of an instance which belongs to each class k , for X set of input features is presented in eq (i)

$$P(Y = k|X) = \frac{e^{X \cdot W_k + b_k}}{\sum_{j=1}^K e^{X \cdot W_j + b_j}} \quad (i)$$

Here, $P(Y = k|X)$ is the probability that the instance belongs to class k , X represents the input features, W_k and b_k are the weights and bias specific to class k , K is the total number of class, $\sum_{j=1}^K e^{X \cdot W_j + b_j}$ is the sum of exponential scores for all classes, e represents the Euler's number.

Random Forest: In the realm of multi-disease detection, Random Forest stands out as a powerful machine learning approach. Each tree independently evaluates input data, providing specific disease predictions [14]. The equation for a Random Forest algorithm can be represented in eq (ii):

$$\hat{y} = \frac{1}{N} \sum_{i=1}^N f_i(X) \quad (ii)$$

Here, N is the number of decision trees, \hat{y} is the predicted output, $f_i(X)$ represents the prediction of the i -th decision tree for the input features X .

KNN: In this approach, diseases are classified based on the majority class among their K nearest neighbors in the feature space. When a new case is presented, the algorithm identifies its K nearest neighbors and assigns the majority disease class among these neighbours to the new case [15]. It is mentioned by eq (iii).

$$Distance(X_{new}, X_i) = \sqrt{\sum_{j=1}^n (X_{new,j} - X_{i,j})^2} \quad (iii)$$

Naïve Bayes: This probabilistic machine learning algorithm, is adept at multi-disease detection due to its efficiency and simplicity. Using statistical methods, it estimates the probability of each disease given the symptoms and selects the disease with the highest probability as the prediction [16]. It is presented by eq (iv).

$$\frac{\text{Posterior Prob of class } A = \text{Prior Prob of class } A \times \text{Likelihood prob of class } A}{\text{Prior prob without any specific class}} \quad (iv)$$

SVM: In this context, SVM works by finding a hyperplane that best separates different disease classes within the feature space, maximizing the margin between classes. This margin signifies the distance between the closest data points of different classes, ensuring optimal classification [16]. It is mentioned by eq (v)

$$f(X) = \text{argmax}_i (w_i \cdot X + b_i) \quad (v)$$

Here, $f(X)$ is the prediction function for multiclass SVM, i represent each class from 1 to k , w_i is the weight vector, X is the input feature vector, and b_i bias.

LGB Classifier: LightGBM is a powerful gradient boosting framework that efficiently handles large and diverse datasets, making it ideal for identifying patterns across multiple diseases. Through intelligent feature selection and high predictive accuracy, LGB classifier enhances multi-disease detection, aiding early diagnosis and timely medical interventions for improved patient outcomes [17]. The general equation for the prediction function in LightGBM can be represented as eq (vi):

$$F(x) = \sum_{m=1}^M f_m(x) \quad (vi)$$

Here, $F(x)$ represents final prediction for x , M is total number of trees, and $f_m(x)$ represents the prediction of the m^{th} tree.

CatBoost: CatBoost, a high-performance gradient boosting algorithm, stands out in multi-disease detection and classification tasks. CatBoost employs an innovative technique called ordered boosting, which arranges categorical variables naturally, enhancing predictive accuracy. Moreover, CatBoost automatically selects optimal hyperparameters, reducing manual tuning efforts [17].

AdaBoost: Adaptive Boosting, in multi-disease detection, excels at handling complex and diverse medical data, to enhance the identification of various diseases. Its ability to adapt and correct errors iteratively makes it particularly effective in capturing intricate patterns within the data, leading to precise and reliable predictions [16]. The mathematical equation for AdaBoost's final prediction $F(x)$ in the context of multiclass classification with α as weight and $h_m^k(x)$ as prediction of weak learners, is presented in eq(vii).

$$F(x) = \text{argmax}_k (\sum_{m=1}^M \alpha_m \cdot h_m^k(x)) \quad (vii)$$

XGBoost: Extreme Gradient Boosting captures intricate patterns within medical datasets and makes it valuable to diagnose multiple diseases simultaneously. It has the ability to handle missing data, scales it efficiently to large datasets, as well as prevents overfitting by incorporating regularization techniques [17].

E. Performance evaluation: The parameters to examine the algorithms in this research are as: *Accuracy* is simply the percentage of predictions that are correct. *Precision* measures the percentage of positive predictions that are actually correct. *Recall* measures the percentage of actual positives that are correctly predicted. *F1 score* is calculated as the harmonic mean of precision as well as recall and gives more weight to precision and recall than accuracy [18].

IV. RESULTS

This section presents the performance of the applied machine learning models on the basis of the aforementioned parameters after being trained with the multi-disease dataset. In the evaluation of various machine learning models, several metrics such as accuracy, precision, recall, and F1 score were considered to assess their performance for the whole dataset, as shown in Table 2.

Table 2: Evaluation of applied techniques for multi-disease dataset

Models	Accuracy (%)	Precision	Recall	F1 score
Logistic Regression	57.59	0.57	0.57	0.57
Random Forest	93.03	0.92	0.93	0.92
Support Vector Machine	76.49	0.79	0.74	0.75
K Nearest Neighbour	80.18	0.79	0.79	0.79
LGBM	76.49	0.79	0.74	0.75
AdaBoost	94.37	0.94	0.94	0.94
CatBoost	93.36	0.93	0.93	0.93
XGBoost	93.11	0.93	0.93	0.93
Naïve Bayes	64.48	0.64	0.64	0.64

Logistic Regression and Naïve Bayes computed the lowest accuracy of 57.79% and 64.48% respectively while as the highest was computed by AdaBoost with 94.37%. Random Forest, CatBoost, and XGBoost also demonstrated great results with accuracy levels around 93-94% and showcase their ability to make highly accurate predictions across classes. Besides this, these techniques also displayed outstanding precision, recall, and F1 scores of 0.92-0.94 which indicate their robustness. Support Vector Machine showed a relatively lower accuracy of 76.49% and slightly lower recall (0.74) but maintained decent precision and F1 score at 0.79 and 0.75 respectively. K Nearest Neighbour and LGBM also performed well with an accuracy of around 80% and 76.49% respectively, demonstrating a good balance between precision, recall, and F1 score.

Besides this, the confusion matrix of the algorithms (Figure 4) have been also generated in a size of 3 x 3 to fetch true label as well as predicted label for each class.

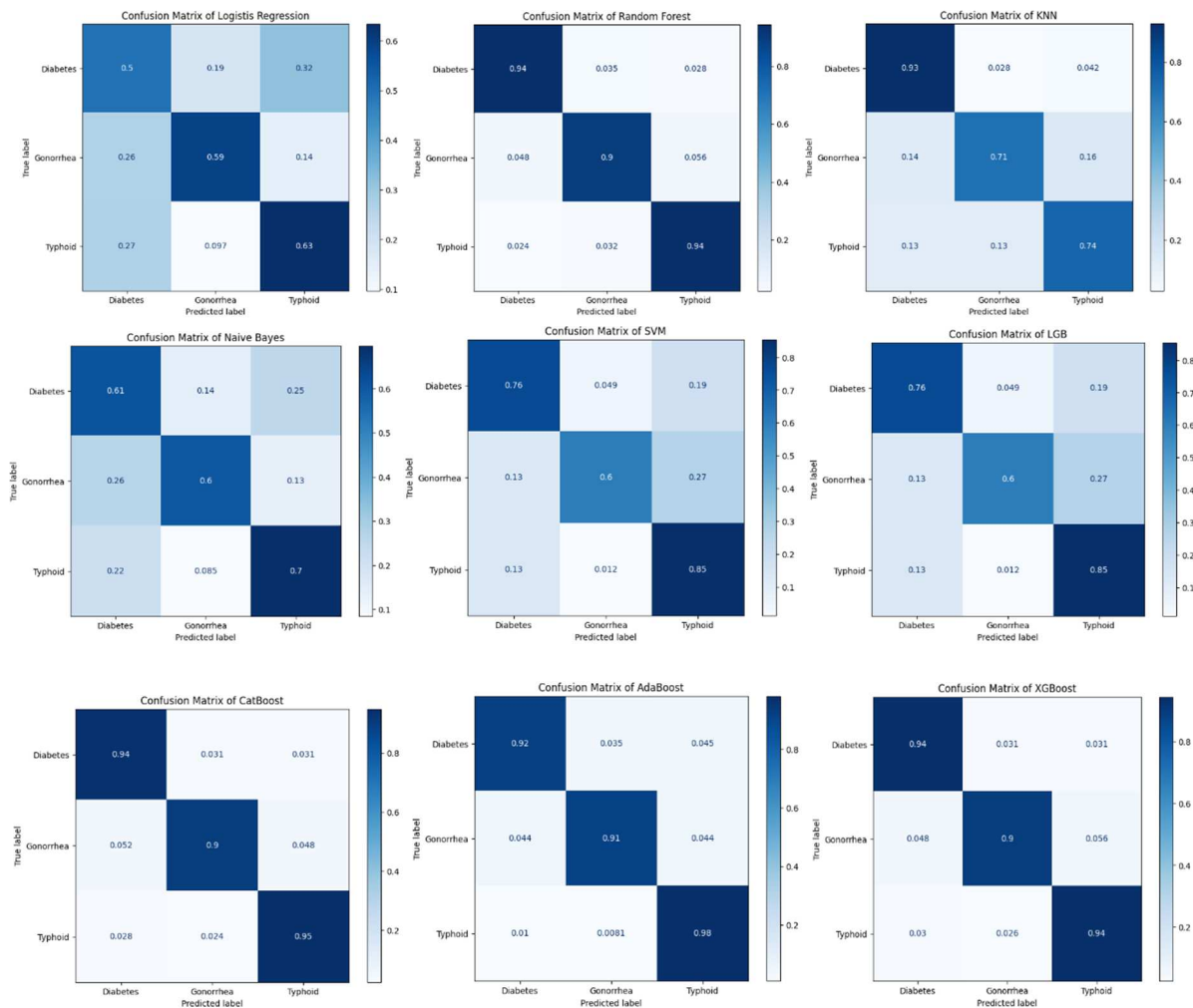


Fig 4: Confusion matrix of applied algorithms

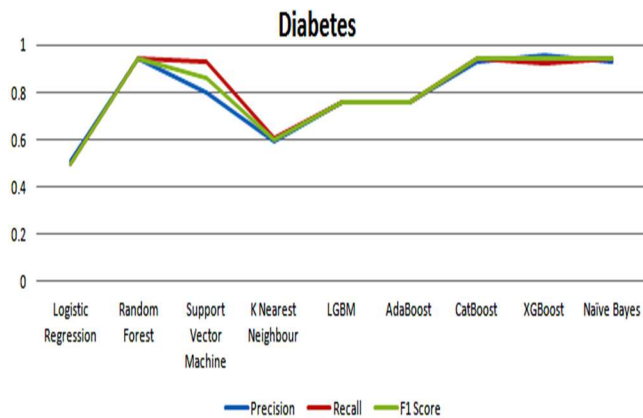
Based on it, the precision, recall, and F1 score value of the applied machine learning models for the Diabetes, Gonorrhoea, and Typhoid class has been also computed and mentioned in Table 3, 4, and 5 respectively. The same has

been also presented graphically in Figure 5, 6, and 7 respectively.

Table 3: Evaluation of techniques for Diabetes class

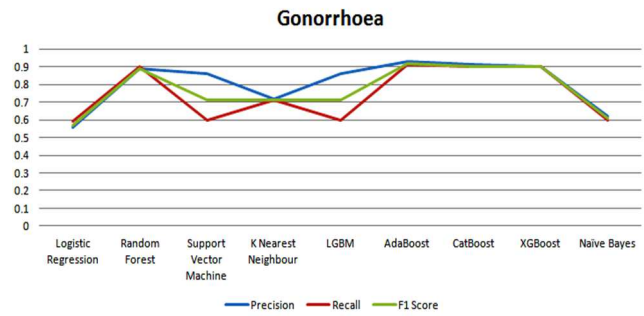
Models	Precision	Recall	F1 Score
Logistic Regression	0.51	0.50	0.50
Random Forest	0.94	0.94	0.94
K Nearest Neighbour	0.80	0.93	0.86
Naïve Bayes	0.59	0.61	0.60
SVM	0.76	0.76	0.76
LGB classifier	0.76	0.76	0.76
CatBoost	0.93	0.94	0.94
AdaBoost	0.96	0.92	0.94
XGBoost	0.93	0.94	0.94

In case of diabetic class, Logistic Regression presented the lowest performance by obtaining precision of 0.51, recall of 0.50, and F1 score of 0.50 which indicates it's below par accuracy in predicting both positive and negative class. K Nearest Neighbour, Support Vector Machine, and LGB classifier had a moderate performance with F1 scores ranging from 0.76 to 0.80 and showcased a decent balance between precision and recall. CatBoost, Random Forest, and XGBoost displayed strong overall performance by boasting high precision, recall, and F1 scores of 0.94 which highlights their ability to make accurate predictions across the classes. Among them all, AdaBoost stood at the top with a high precision of 0.96 but a slightly lower recall of 0.92, resulting in an overall F1 score of 0.94.

**Fig 5:** Analysing algorithms using diabetes class**Table 4:** Evaluation of techniques for Gonorrhoea class

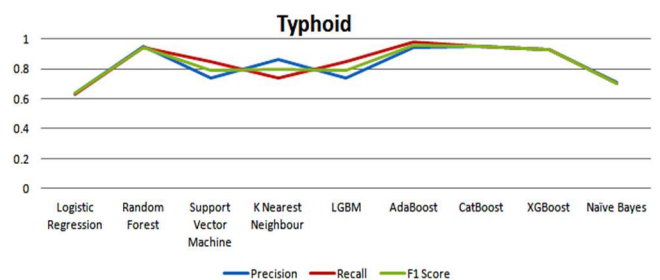
Models	Precision	Recall	F1 Score
Logistic Regression	0.56	0.59	0.57
Random Forest	0.89	0.90	0.89
SVM	0.86	0.60	0.71
K Nearest Neighbour	0.72	0.71	0.71
LGBM	0.86	0.60	0.71
AdaBoost	0.93	0.91	0.92
CatBoost	0.91	0.90	0.90
XGBoost	0.90	0.90	0.90
Naïve Bayes	0.62	0.60	0.61

In case of gonorrhoea class, Logistic Regression and Naïve Bayes computed the least precision of 0.56 and 0.62, recall of 0.59 and 0.60, and F1 score of 0.57 and 0.61 respectively. The highest value has been computed by AdaBoost with 0.93, 0.91, and 0.92 precision, recall, and F1 score respectively followed by CatBoost and XGBoost which demonstrates its effectiveness to handle the given classification task. Support Vector Machine, K Nearest Neighbour, and LGBM displayed similar F1 scores of 0.71, although SVM had a higher precision (0.86) but lower recall (0.60) which indicates a bias towards positive class predictions.

**Fig 6:** Analysing algorithms using gonorrhoea class**Table 5:** Evaluation of techniques for Typhoid class

Models	Precision	Recall	F1 Score
Logistic Regression	0.64	0.63	0.64
Random Forest	0.95	0.94	0.94
SVM	0.74	0.85	0.79
K Nearest Neighbour	0.86	0.74	0.80
LGBM	0.74	0.85	0.79
AdaBoost	0.94	0.98	0.96
CatBoost	0.95	0.95	0.95
XGBoost	0.93	0.93	0.93
Naïve Bayes	0.71	0.70	0.70

In case of typhoid class, Logistic Regression and XGBoost demonstrated moderate performance with a balanced precision, recall, and F1 score of 0.64 and 0.93 respectively which indicates a relatively accurate classification across the classes. Random Forest and CatBoost exhibited strong performance with high precision, recall, and F1 score of around 0.95 and showcase their ability to make precise and accurate predictions. Support Vector Machine, K Nearest Neighbour, and LGBM showed similar F1 scores around 0.79-0.80, with SVM and LGBM having a slightly higher recall which indicates a better ability to capture positive instances. AdaBoost performed remarkably well with F1 scores of 0.96 while as Naïve Bayes, obtained a lower F1 score of 0.70.

**Fig 7:** Analysing algorithms using typhoid class

V. CONCLUSION

The paper explores the potential of machine learning in identifying multiple diseases in the healthcare field. The text delves into the use of different algorithms that have been trained on a diverse dataset of diseases to achieve exceptional accuracy. During execution, it was observed that AdaBoost achieved excellent performance across all classes of the dataset in terms of precision, accuracy, recall, and F1 score. However, the research does face some challenges like restricted diseases were used and certain algorithms perform poorly in terms of their evaluative metrics. This suggests that more diseases should be taken and the hyperparameters of the

models need to be optimized to improve their performance in detecting and diagnosing diseases in future.

REFERENCES

- [1] S. Gupta and Y. Kumar, "Cancer prognosis using artificial intelligence-based techniques," *SN Comput. Sci.*, vol. 3, no. 1, pp. 1–8, 2022.
- [2] Y. Kumar, S. Gupta, R. Singla, and YC Hu, "A systematic review of artificial intelligence techniques in cancer prediction and diagnosis," *Arch. Comput. Methods Eng.*, vol. 29, pp. 2043-2070, 2021.
- [3] A. Koul, RK Bawa, and Y. Kumar, "Artificial intelligence techniques to predict the airway disorders illness: a systematic review," *Arch. Comput. Methods Eng.*, vol. 30, pp. 831-864, 2022.
- [4] P.S. Sisodia, GK Ameta, Y. Kumar, and N. Chaplot, "A Review of Deep Transfer Learning Approaches for Class-Wise Prediction of Alzheimer's Disease Using MRI Images," *Arch. Comput. Methods Eng.*, vol. 30, no. 4, pp. 2409–2429, 2023.
- [5] K. Kaur, C. Singh, and Y. Kumar, "Diagnosis and Detection of Congenital Diseases in New-Borns or Fetuses Using Artificial Intelligence Techniques: A Systematic Review," *Arch. Comput. Methods Eng.*, vol. 30, pp. 3031-3058, 2023.
- [6] K. Bansal, RK Bathla, and Y. Kumar, "Deep transfer learning techniques with hybrid optimization in early prediction and diagnosis of different types of oral cancer," *Soft Comput.*, vol. 26, no. 21, pp. 11153–11184, 2022.
- [7] P. Bhardwaj, Y. Kumar, and G. Bhandari, "AI-Enabled Computational Techniques for Cancer Diagnosis," *IEEE 8th Uttar Pradesh Section International Conference on Electrical, Electronics and Computer Engineering (UPCON)*, pp. 1–7, 2021.
- [8] B. Kakkur, et al., "An IoMT-Based Federated and Deep Transfer Learning Approach to the Detection of Diverse Chest Diseases Using Chest X-Rays," *HUMAN-CENTRIC Comput. Inf. Sci.*, vol. 12, no. 24, 2022.
- [9] G. Swapna, R. Vinayakumar, and K. P. Soman, "Diabetes detection using deep learning algorithms." *ICT express*, vol. 4, no. 4, pp. 243-246, 2018.
- [10] MA. Bhuiyan, et al. "Prediction of Typhoid Using Machine Learning and ANN Prior To Clinical Test." 2023 International Conference on Computer Communication and Informatics (ICCCI). IEEE, 2023.
- [11] Y. Bao, et al. "Predicting the diagnosis of HIV and sexually transmitted infections among men who have sex with men using machine learning approaches." *Journal of Infection*, vol. 82, no. 1, pp. 48-59, 2021.
- [12] M. Sulyma, "Diarrhoeal disease," May 02, 2017. <https://www.who.int/news-room/fact-sheets/detail/diarrhoeal-disease>
- [13] G.P. Kanna, et al. "A Review on Prediction and Prognosis of the Prostate Cancer and Gleason Grading of Prostatic Carcinoma Using Deep Transfer Learning Based Approaches," *Arch. Comput. Methods Eng.*, vol. 30, no. 5, pp. 3113–3132, 2023.
- [14] S. Kaur, Y. Kumar, A. Koul, and S. K. Kamboj, "A Systematic Review on Metaheuristic Optimization Techniques for Feature Selections in Disease Diagnosis: Open Issues and Challenges," *Arch. Comput. Methods Eng.*, vol. 30, pp. 1863-1895, 2022.
- [15] F. Yousaf, S. Iqbal, N. Fatima, T. Kousar, and M. S. M. Rahim, "Multi-class disease detection using deep learning and human brain medical imaging," *Biomedical Signal Processing and Control*, vol. 85, 2023.
- [16] G. Angayarkanni and S. Hemalatha, "Evaluating the Performance of Supervised Machine Learning Algorithms for Predicting Multiple Diseases: A Comparative Study," 2023 9th International Conference on Advanced Computing and Communication Systems (ICACCS). vol. 1. IEEE, 2023.
- [17] JD Huang, et al., "Applying artificial intelligence to wearable sensor data to diagnose and predict cardiovascular disease: a review," *Sensors*, vol. 22, no. 20, pp.1-28, 2022.
- [18] Y. Kumar, A. Koul, S. Kaur, and YC Hu, "Machine Learning and Deep Learning Based Time Series Prediction and Forecasting of Ten Nations' COVID-19 Pandemic," *SN Comput. Sci.*, vol. 4, no. 1, pp. 1–27, 2023.