

Haematological Harmony: Early Blood Cancer Detection

¹Smt. K S SUKRUTHA, ²Smt.KAVYA S N, ³Ms. CHANDANA N, ⁴Ms.DEEKSHA S, ⁵Ms.DEVIKA P, ⁶Ms.SANGEETHA

¹ HOD & Associate Professor, Department of Computer Science, MMK&SDM MMV, Mysuru

²Assistant Professor, Department of Computer Science, MMK&SDM MMV, Mysuru

^{3,4,5,6} Student, Department of Computer Science, MMK&SDM MMV, Mysuru

Abstract- Blood cancer, medically referred to as hematological malignancy, encompasses a group of cancers that originate in the blood, bone marrow, and lymphatic system. The timely detection and accurate diagnosis of blood cancer play a pivotal role in devising effective treatment strategies. The emerging field of machine learning has demonstrated immense promise in medical diagnostics and disease prediction. Thus, this project aims to leverage machine learning algorithms to develop a predictive model for blood cancer.

Keywords- Oncology, Haematological Malignancy, Pattern Recognition, Classification Algorithms, Haematology

I. INTRODUCTION

Blood cancer encompasses various types of cancers, such as leukemia, lymphoma, and multiple myeloma. The early detection and accurate diagnosis of blood cancer are crucial for determining appropriate treatment strategies and improving patient outcomes. However, diagnosing blood cancer can be complex and challenging due to its diverse manifestations and overlapping symptoms.

Machine learning, a subfield of artificial intelligence, has emerged as a powerful tool in medical research and diagnostics. By analyzing large volumes of data and recognizing complex patterns, machine learning algorithms can provide valuable insights and predictive capabilities. In the field of healthcare, machine learning has shown great potential in disease prediction, risk assessment, and personalized treatment planning.

The proposed approach involves utilizing advanced machine learning techniques to analyze patient data, including clinical records and genetic information. By training the machine learning model on a comprehensive dataset of blood cancer cases, it will

learn to recognize patterns and identify potential indicators or biomarkers associated with the disease. This predictive model can then be used to evaluate new patient data and provide risk assessment or early detection of blood cancer.

The potential benefits of this project are multifaceted. Early detection of blood cancer can significantly improve patient outcomes by enabling timely interventions and tailored treatment plans. Additionally, the machine learning model can aid healthcare professionals in making accurate diagnoses, leading to more precise and personalized patient care. Ultimately, the development of an effective predictive model for blood cancer using machine learning has the potential to enhance healthcare practices, advance medical research, and contribute to improved patient outcomes in the field of hematology.

II. SCOPE

This project can be beneficial for various stakeholders involved in the diagnosis and treatment of blood cancer. The primary users of this project can include:

Healthcare Professionals: Oncologists, hematologists, and other healthcare professionals involved in the diagnosis and treatment of blood cancer can utilize this project.

Medical Researchers: Researchers in the field of hematology and oncology can utilize the findings and methodologies of this project to advance their understanding of blood cancer.

III. OBJECTIVE

The objective of the project is to develop a machine learning model specifically tailored for the prediction and diagnosis of blood cancer. The model will be trained on comprehensive datasets, including

clinical records and genetic information to recognize patterns and indicators associated with different types of blood cancer. The aim is to create a robust and accurate predictive model that can assist in early detection and accurate diagnosis of blood cancer.

IV. LITERATURE SURVEY

A literature survey is a critical and comprehensive analysis of existing research and publications related to a specific topic or area of study. It provides an overview of the current state of knowledge, identifies key concepts, theories, methodologies, and gaps in the existing work. The purpose of conducting a literature survey is to understand the background of the research problem, build up on previous studies, and provide a strong foundation for further investigation. In this survey, relevant scholarly articles, journals, conference papers, and other credible sources have been reviewed to gather insights and evaluate the developments in the chosen domain. The collected information helps in identifying trends, comparing different approaches, and determining the most effective strategies and techniques applied in earlier research. Ultimately, the literature survey aids in refining the research objectives, formulating hypotheses, and guiding the direction of the current study or project.

4.1 Related Works

1. Blood Cancer Detection using Machine Learning

<https://ieeexplore.ieee.org/document/9675987>

Cancer is a deadly disease. Initial detection of cancer is the best way to cure the disease. Medical Image Processing plays an essential part in the detection of disease. Leukemia is a kind of blood cancer that happens due to irregular or immature White Blood Cell (WBC) s. In general, WBC is the fighter to fight against infectious cells in the human body. Abnormal growth of WBC from bone marrow will destroy the other cells and affect bone marrow and lymphatic tissues. These cells do not function properly and lead to leukemia. In olden days, identification of disease and cell counting in the blood were very complicated. In the medical sector, they use a device called Haemocytometer which counts the amount of cells in the blood manually. But it takes more time for counting and gives inaccurate results. To overcome these issues, a software based solution is given with the help of microscopic images. With this image

processing technique, the number of RBCs, WBCs and platelets are calculated and also whether a person can be affected by leukemia or not is identified.

2. Blood Cancer Detection Using Machine Learning

<https://iarjset.com/wp-content/uploads/2022/06/IARJSET.2022.9639.pdf>

Leukemia (blood cancer) begins in the bone marrow and causes the formation of a large number of abnormal cells. The most common types of leukemia known are Acute lymphoblastic leukemia (ALL), Acute myeloid leukemia (AML), Chronic lymphocytic leukemia (CLL) and Chronic myeloid leukemia (CML). This project tries to devise a methodology for the detection of Leukemia using image processing techniques, thus automating the detection process. Our project consists of development of a machine learning algorithm to detect cancer using microscopy image.

<https://www.hindawi.com/journals/sp/2021/9933481/>
Authors

3. The early detection and diagnosis of leukemia, i.e., the precise differentiation of malignant leukocytes with minimum costs in the early stages of the disease, is a major problem in the domain of disease diagnosis. Despite the high prevalence of leukemia, there is a shortage of flow cytometry equipment, and the methods available at laboratory diagnostic centers are time-consuming. Motivated by the capabilities of machine learning (machine learning (ML)) in disease diagnosis, the present systematic review was conducted to review the studies aim in to discover and classify leukemia by using machine learning. *Methods.* A systematic search in four databases (PubMed, Scopus, Web of Science, and ScienceDirect) and Google Scholar was performed via a search strategy using Machine Learning (ML), leukemia, peripheral blood smear (PBS) image, detection, diagnosis, and classification as the keywords. Initially, 116 articles were retrieved. After applying the inclusion and exclusion criteria, 16 articles remained as the population of the study. *Results.* This review study presents a comprehensive and systematic view of the status of all published ML-based leukemia detection and classification models that process PBS images. The average accuracy of the ML methods applied in PBS image analysis to detect leukemia was >97%, indicating

that the use of ML could lead to extraordinary outcomes in leukemia detection from PBS images. Among all ML techniques, deep learning (DL) achieved higher precision and sensitivity in detecting different cases of leukemia, compared to its precedents. ML has many applications in analyzing different types of leukemia images, but the use of ML algorithms to detect acute lymphoblastic leukemia (ALL) has attracted the greatest attention in the fields of hematology and artificial intelligence. *Conclusion.* Using the ML method to process leukemia smear images can improve accuracy, reduce diagnosis time, and provide faster, cheaper, and safer diagnostic services. In addition to the current diagnostic methods, clinical and laboratory experts can also adopt ML methods in laboratory applications and tools.

4. Multiclass blood cancer classification using deep CNN with optimized features

<https://www.sciencedirect.com/science/article/pii/S2590005623000176>

Breast cancer, lung cancer, skin cancer, and blood malignancies such as leukemia and lymphoma are just a few instances of cancer, which is a collection of cells that proliferate uncontrollably within the body. Acutelymphoblastic leukemia is one of the significant forms of malignancy. The hematologists frequently make an oversight while determining a blood cancer diagnosis, which requires an excessive amount of time. Thus, this research reflects on a novel method for the grouping of the leukemia with the aid of the modern technologies like Machine Learning and Deep Learning. The proposed research pipeline is occupied into some interconnected parts like dataset building, feature extraction with pre-trained Convolution Neural Network (CNN) architectures from each individual image of blood cells, and classification with the conventional classifiers. The dataset for this study is divided into two identical categories, Benign and Malignant, and then reshaped into four significant classes, each with three subtypes of malignant, namely, Benign, Early Pre-B, Pre-B, and Pro-B. The research first extracts the features from the individual images with CNN models and then transfers the extracted features to the feature selections such as Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA), and SVC Feature Selectors along with two nature-inspired algorithms like

Particle Swarm Optimization (PSO) and Cat Swarm Optimization (CSO). After that, research has applied these seven Machine Learning classifiers to accomplish the multi-class malignant classification. To assess the efficacy of the proposed architecture, a set of experimental data have been enumerated and interpreted accordingly. The study discovered a maximum accuracy of 98.43% when solely using pre-trained CNN and classifiers. Nevertheless, after incorporating PSO and CSO, the proposed model achieved the highest accuracy of 99.84% by integrating the ResNet50 CNN architecture, SVC feature selector, and LR classifiers. Although the model has a higher accuracy rate, it does have some drawbacks. However, the proposed model may also be helpful for real-world blood cancer classification.

V. PROPOSED METHODOLOGY

Machine learning techniques can be utilized for blood cancer prediction by leveraging the power of data analysis, pattern recognition, and predictive modeling. In the Proposed System, Random Forest algorithm is used. It is a popular machine learning algorithm that can be utilized for blood cancer prediction using a text dataset. Random Forest is an ensemble learning method that combines multiple decision trees to make predictions.

5.1 Advantages of proposed system

The approach of using the Random Forest algorithm for blood cancer prediction using a text dataset offers several advantages:

- **Accuracy and Robustness:** Random Forest is known for its high accuracy and robustness in handling complex datasets. By leveraging an ensemble of decision trees, it can effectively capture intricate patterns and relationships within the text dataset, leading to accurate predictions for blood cancer.
- **Feature Importance:** Random Forest provides a measure of feature importance, indicating the relevance and contribution of each feature in the prediction process. This information can aid in identifying the most important textual features that play a significant role in predicting blood cancer.
- **Scalability:** Random Forest is known for its scalability, making it suitable for handling large text datasets. It can efficiently handle a large number of textual features and instances,

allowing for potential scalability as the dataset expands. This is particularly important in healthcare domains where datasets can be extensive and continuously growing

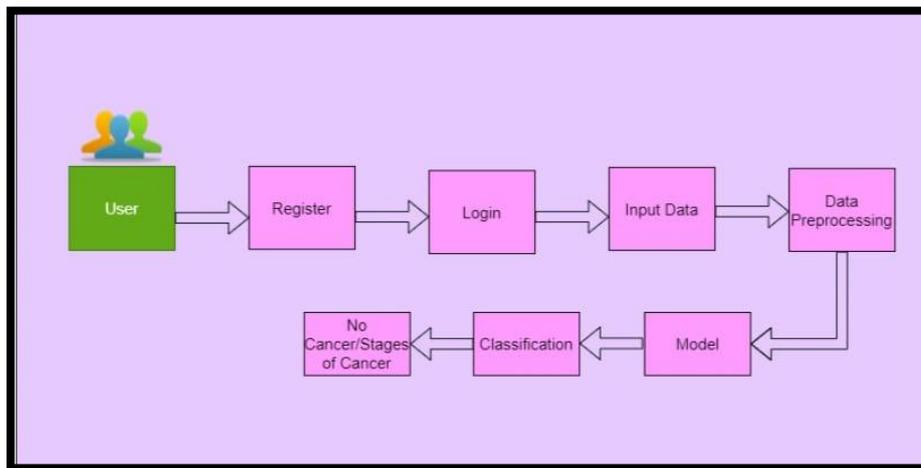
VI. ALGORITHM-RANDOM FOREST ALGORITHM

A Random Forest Algorithm is a supervised machine learning algorithm that is extremely popular and is used for Classification and Regression problems in Machine Learning. We know that a forest comprises numerous trees, and the more trees more it will be robust. Similarly, the greater the number of trees in a Random Forest Algorithm, the higher its accuracy and problem-solving ability. Random Forest is a classifier that contains several decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset. It is based on the concept of ensemble learning which is a process of combining multiple classifiers to solve a complex problem and improve the performance of the model.

VII. DATASET

The dataset from Kaggle comprises 15,072 rows and 13 columns, totaling approximately 2.62 MB. Each record represents a unique patient with various hematological parameters such as Red Blood Cell count (RBC), Packed Cell Volume (PCV), Mean Corpuscular Volume (MCV), and Hemoglobin level (HGB). The "Cancer Level" column indicates the severity or presence of blood cancer, serving as the target variable for classification tasks. The "Sex" column is categorical (Male/Female), while the rest are continuous numerical values. Notably, there are no missing values across any fields, ensuring data integrity. This dataset is well-suited for developing and validating predictive models for early detection of blood cancer using machine learning algorithms. Additionally, the inclusion of demographic features like age and sex allows for deeper analysis of correlations and potential risk factors across different population groups.

VIII. SYSTEM ARCHITECTURE MODULE



This web application consists of following modules

- Register
Purpose: Allows new users to create an account by providing basic credentials to access the system.
Example: A new user enters their name, email, and password to register.
- Login
Purpose: Authenticates users with their registered credentials to ensure secure access to the platform.
Example: A user enters their registered email and password to log in.
- Upload Dataset

Purpose: Enables the Admin to upload a material image dataset for training the CNN model.

Example: The Admin uploads a folder containing images of different materials (e.g., wood, metal, plastic) for training.

- Train
Purpose: Provides functionality to train the CNN model on the uploaded dataset, generating insights from the data.

Example: The system processes the uploaded images and trains the model to recognize patterns in material types. [\(GeeksforGeeks\)](#)

- Save Model

Purpose: Saves the trained CNN model for future use in detecting material.

Example: After training, the model is saved as a file (e.g., material_detector.model) for later use.

- View Performance

Purpose: Displays metrics such as accuracy and loss of the trained model to evaluate its effectiveness.

Example: The system shows that the model achieved 95% accuracy and 5% loss during training.

- Input Test Data

Purpose: Allows users to upload a material image for analysis and detection.

Example: A user uploads a photo of an unknown material for classification. [RCD GitHub](#)

- View Detection

Purpose: The system displays the detection results for the test data.

Example: The system identifies the uploaded image as "Wood" with 98% confidence.

This system streamlines the process of training and deploying a CNN model for material detection, making it accessible even for users with limited technical expertise.

IX. SOFTWARE TESTING

System testing is a critical phase in the development of a blood cancer detection project using Python and machine learning. It ensures that all components of the system function as intended and meet the specified requirements. Here's an overview of the key aspects involved in system testing for such a project:

1. Functional Testing

Objective: Verify that each feature of the system operates according to the specified requirements.

Examples:

Model Training: Ensure that the Convolutional Neural Network (CNN) model trains correctly on the provided dataset without errors.

Prediction Functionality: Test the model's ability to classify new blood cell images accurately.

User Interface: Confirm that the user interface allows for uploading images and displays results as expected. ([GitHub](#))

2. Performance Testing

Objective: Assess the system's responsiveness and stability under various conditions.

Examples:

Model Inference Time: Measure the time taken by the model to process and classify an image.

System Load: Evaluate how the system performs when handling multiple simultaneous image uploads and predictions.

3. Security Testing

Objective: Identify vulnerabilities in the system to prevent unauthorized access and data breaches.

Examples:

Data Privacy: Ensure that patient data is anonymized and securely stored.

Access Controls: Verify that only authorized users can upload datasets or access sensitive information.

4. Usability Testing

Objective: Evaluate the user-friendliness and intuitiveness of the system.

Examples:

User Interface: Assess the clarity and ease of navigation of the user interface.

User Feedback: Collect feedback from users to identify areas for improvement.

5. Regression Testing

Objective: Ensure that new code changes do not adversely affect the existing functionality of the system.

Examples:

Model Updates: After updating the model with new data, verify that previous functionalities remain unaffected.

Software Updates: Check that updates to the underlying software libraries do not introduce new issues.

6. Integration Testing

Objective: Test the interaction between different components of the system to ensure they work together seamlessly.

Examples:

Model and Interface: Verify that the model correctly integrates with the user interface for real-time predictions.

Database Integration: Ensure that the system correctly stores and retrieves patient data from the database.

7. Acceptance Testing

Objective: Determine whether the system meets the business requirements and is ready for deployment.

Examples:

Stakeholder Approval: Present the system to stakeholders for approval based on predefined criteria.

Real-World Scenarios: Simulate real-world scenarios to ensure the system performs as expected under actual operating conditions.

By systematically addressing these testing areas, you can ensure that your blood cancer detection system is robust, reliable, and ready for deployment.

X. SYSTEM TESTING

Software testing is a process of executing a program or application with the intent of finding the software bugs.

Software testing is a critical element of software quality assurance and represents the ultimate process to ensure the correctness of the product. The quality product always enhances the customer confidence in using the product thereby increasing the business economics. In other words, a good quality product means zero defects, which is derived from a better quality process in testing.

Testing the product means adding value to it by raising the quality or reliability of the product. Raising the reliability of the product means finding and removing errors. Hence one should not test a product to show that it works; rather, one should start with the assumption that the program contains errors and then test the program to find as many of the errors as possible. The main objective of testing is to find defects in requirements, design, documentation, and code as early as possible. The test process should be such that the software product that will be delivered to the customer is defect less. All Tests should be traceable to customer requirements. Test cases must be written for invalid and unexpected, as well as for valid and expected input conditions. A necessary part of a test case is a definition of the expected output or result. A good test case is one that has high probability of detecting an as-yet undiscovered error.

10.1 Manual Testing

Manual testing includes testing a software manually, i.e., without using any automated tool or any script. In this type, the tester takes over the role of an enduser and tests the software to identify any unexpected behavior or bug. There are different stages for manual testing such as unit testing, integration testing, system testing, and user acceptance testing.

Testers use test plans, test cases, or test scenarios to test a software to ensure the completeness of testing. Manual testing also includes exploratory testing, as testers explore the software to identify errors in it.

10.1.1 Different stages for manual testing:

1. Unit Testing

This type of testing is performed by developers before the setup is handed over to the testing team to formally execute the test cases. Unit testing is performed by the respective developers on the individual units of source code assigned areas. The developers use test data that is different from the test data of the quality assurance team.

Tests that are performed during the unit testing in the app are explained below:

- **Module Interface test:** In module interface test, it is checked whether the information is properly flowing in to the program unit (or module) and properly happen out of it or not. E.g. The user registration details should be available from the layout to the corresponding controller and from the controller it should flow to the model.

- **Boundary conditions:** It is observed that much software often fails at boundary related conditions. That's why boundary related conditions are always tested to make safe that the program is properly working at its boundary condition's.

E.g. In case of if...else if... else... construct all the conditions are checked in the app. In case of loops, it is checked to see that the loops are not infinite and terminate once the condition becomes false.

Error handling paths: These are tested to review if errors are handled properly.

E.g. Validation during login (Checking for wrong credentials).

Validation of password during registration (Password should adhere to password policy - minimum 8 characters with a number and special character).

In case the user tries to send SMS without configuring contacts, he should be redirected to Add Contacts screen.

It should be possible to send SMS and play alarm sound even if the user has not configured these settings in Preferences. The app should have default SMS template message and Alarm Sound.

In case the user tries to contact emergency contacts without configuring contacts, he should be redirected to Add Contacts screen.

2. Integration Testing

Integration testing is defined as the testing of combined parts of an application to determine if they function correctly. Integration testing can be done in two ways: Bottom-up integration testing and Top-down integration testing. In this project we have followed the Bottom-up integration method.

Here testing begins with unit testing, followed by tests of progressively higherlevel combinations of units called modules or builds.

Once all the different modules were integrated in the app, the app was tested for the following:

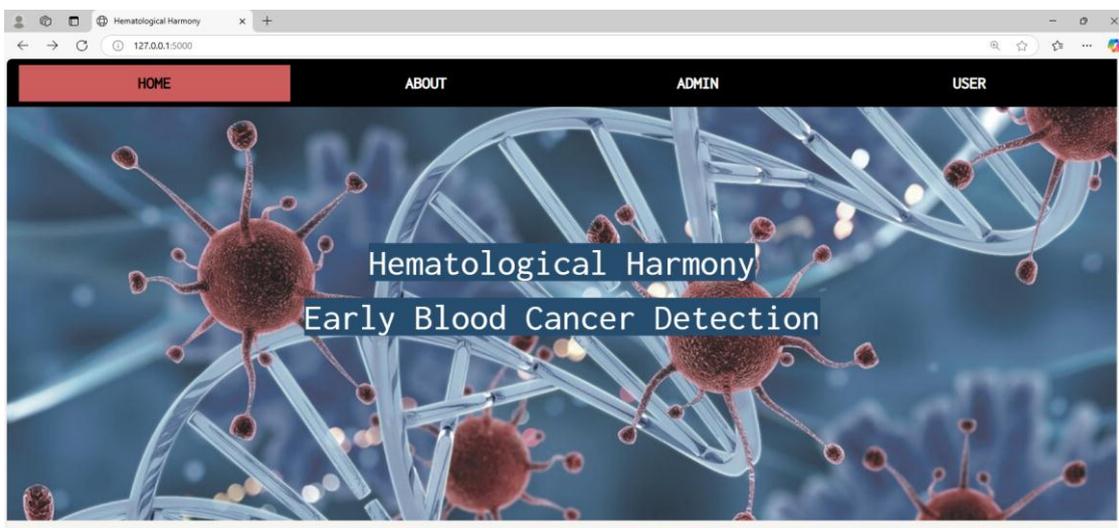
- Transition from one screen to another

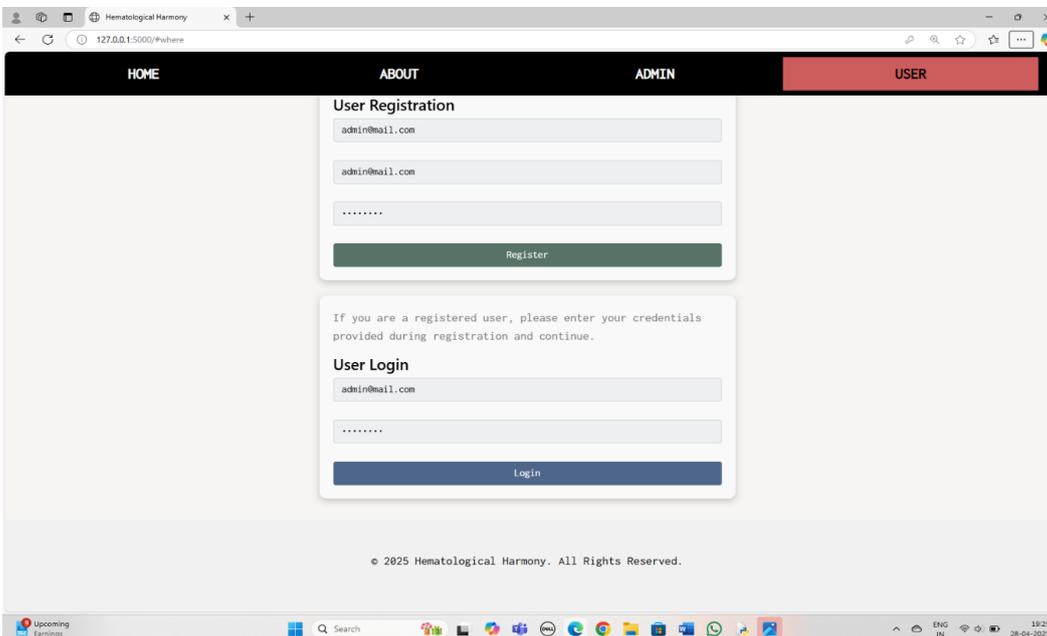
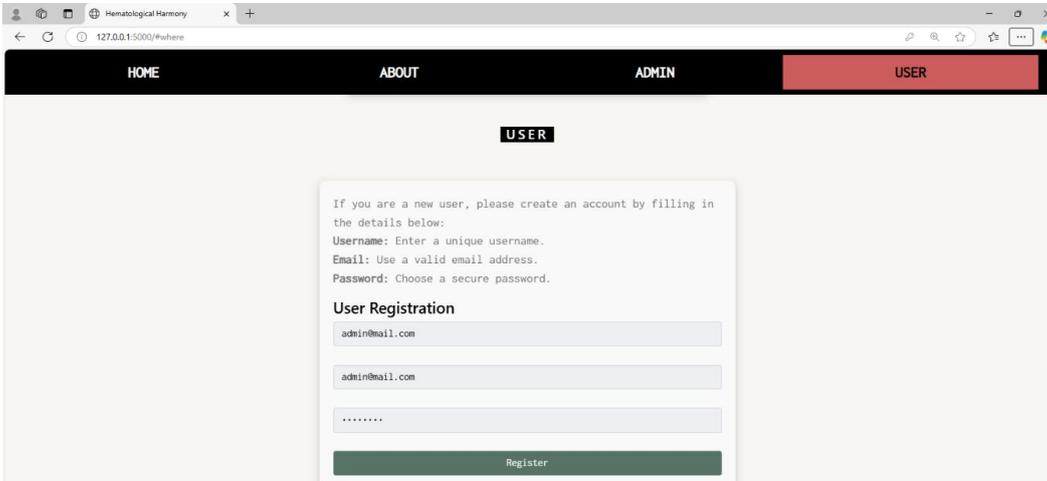
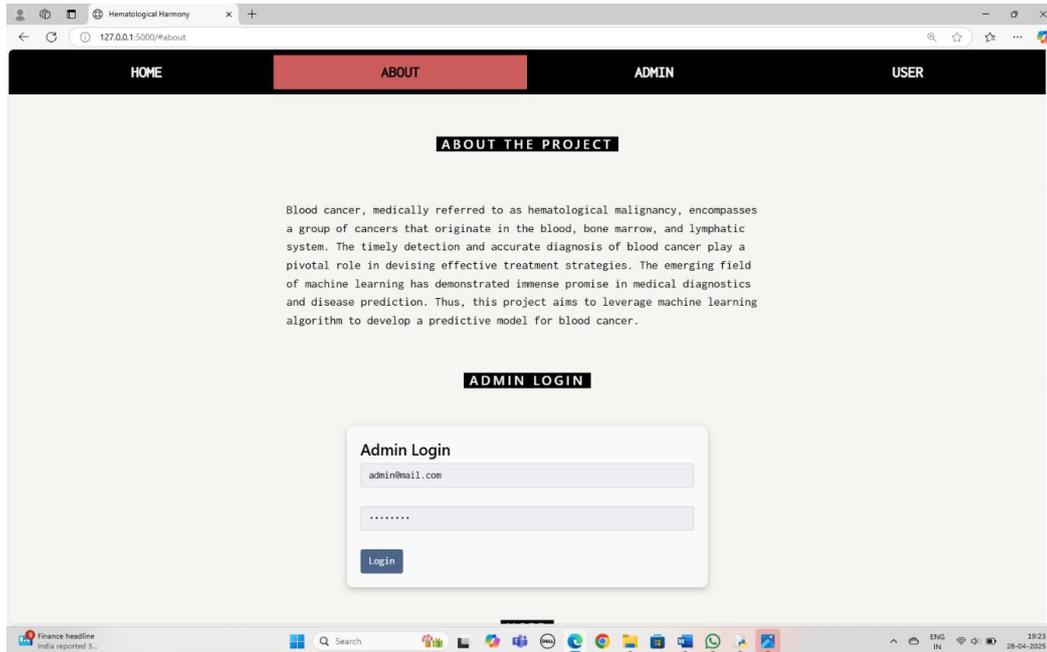
- Data from the layouts is getting saved properly in the database.
- Data is retrieved properly from the database and displayed in the layouts.
- SMS is sent properly to configured contacts during Send SMS actions.
- Menu items are working properly.

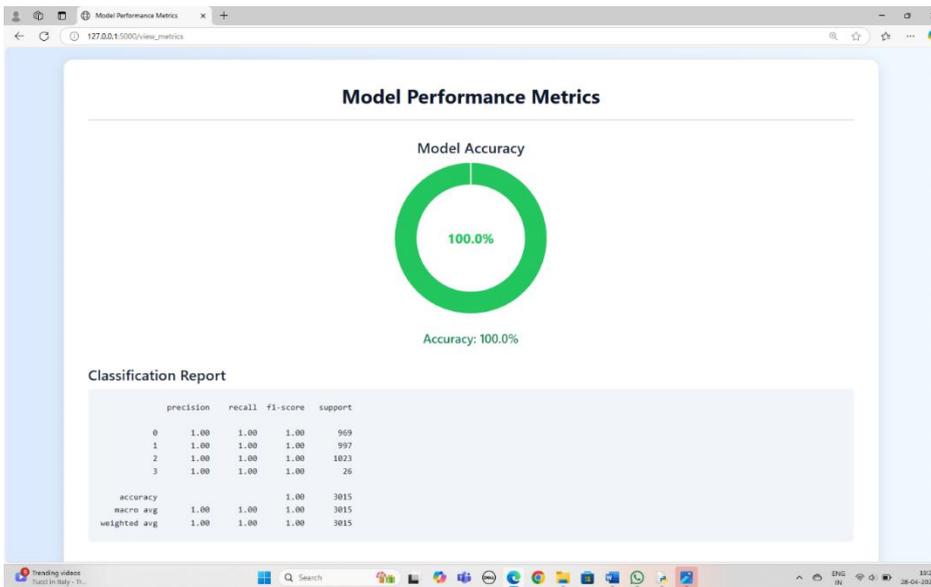
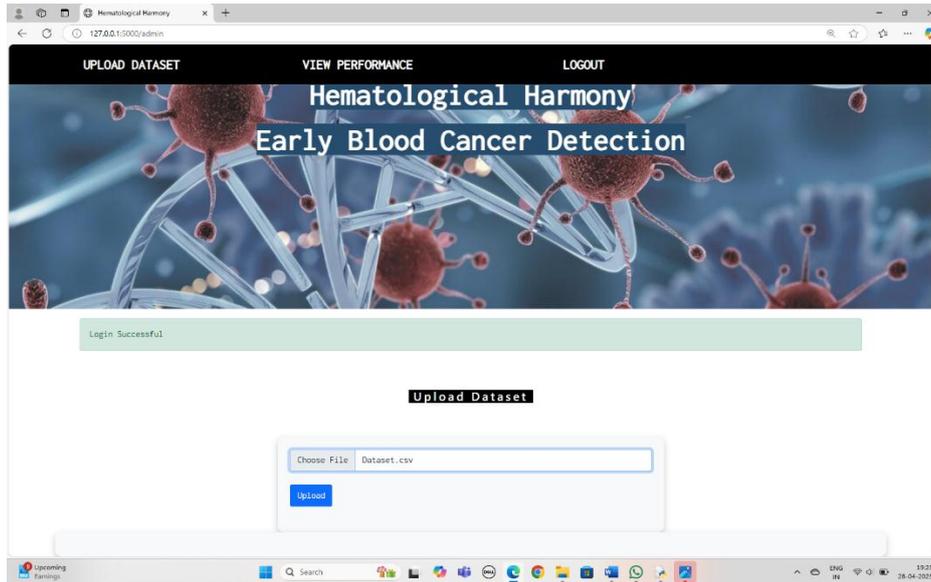
TC 01	Data Collection	To import data from location	Import from file location	Valid location and file name	import data successful	Import data successful	Pass
	Data Collection	To import data from location	Import from file location	Invalid location and file name	Import data unsuccessful	Import data unsuccessful	Pass
TC 02	Data preprocessing	To verify that data has preprocessed	Enter the data file imported	Valid detail	preprocessing successfully	preprocessing successfully	Pass
	Data preprocessing	To verify that data has preprocessed	If invalid data file imported	Invalid detail	Preprocessing unsuccessfully	Preprocessing unsuccessfully	Pass
TC 03	Train the data	To train data	Implement algorithms	Success output	Training successful	Training successful	Pass
TC 04	Detection	Test the trained algorithm	Implement best algorithm to dataset	Success output	Prediction successful	Prediction successful	Pass

10.2 Test Cases

XI. RESULTS



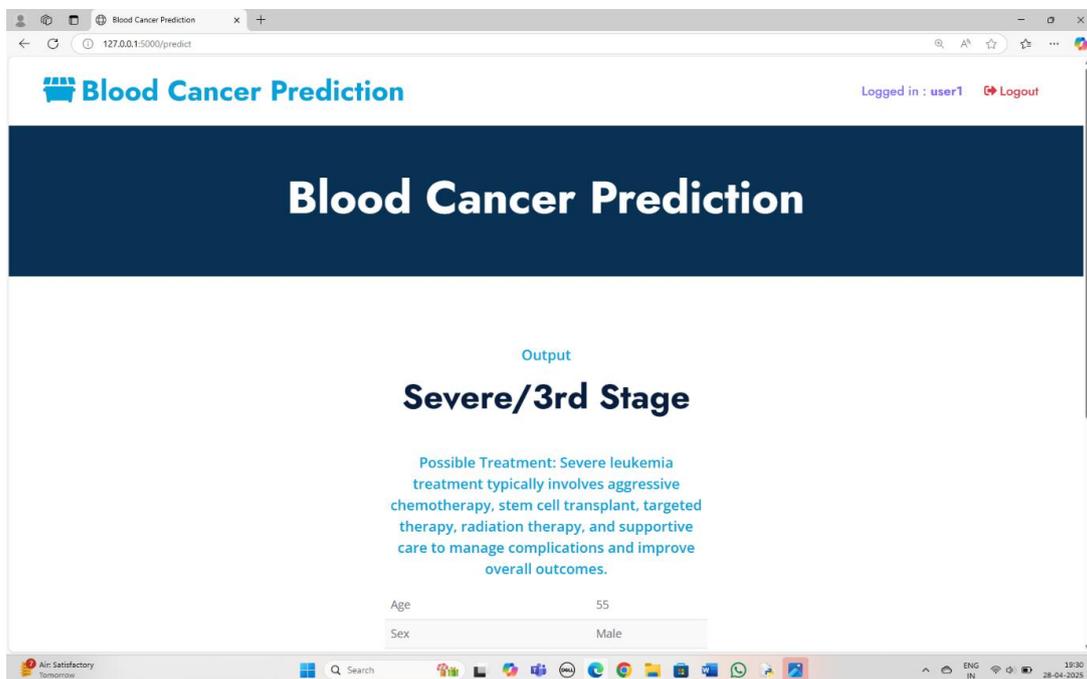
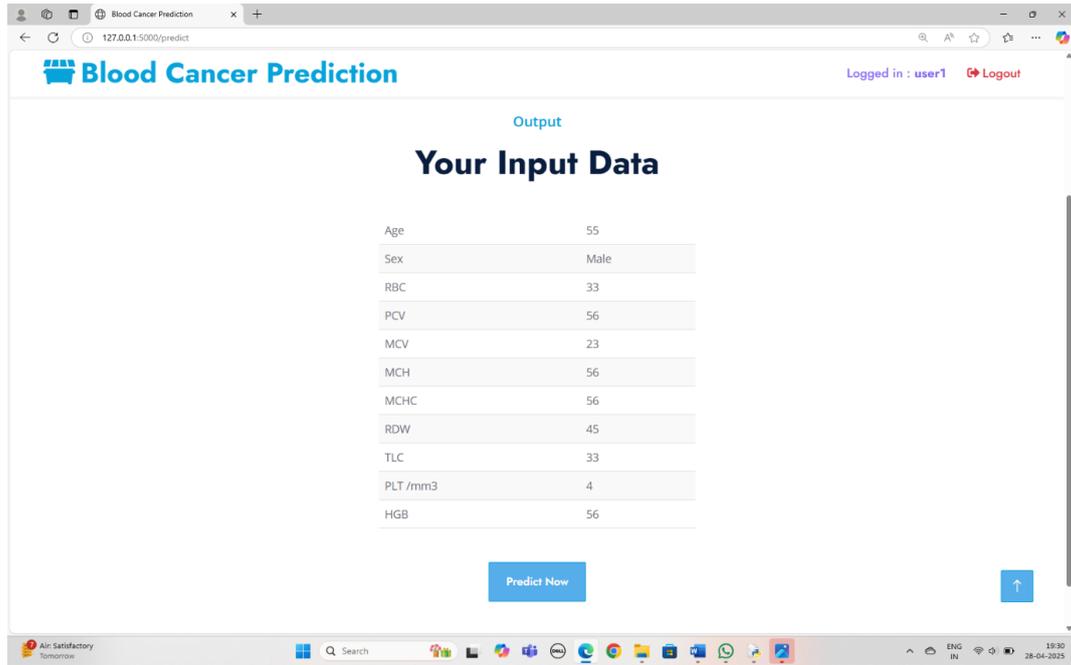




The screenshot shows the 'Predict Cancer' web application interface. It includes a 'Train RF' button and a 'Predict Cancer' button. On the right side, there is a form for inputting patient data with the following fields:

- Age: Enter age
- Sex: Select
- RBC: Enter RBC count
- PCV: Enter PCV
- MCV: Enter MCV
- MCH: Enter MCH
- MCHC: Enter MCHC
- RDW: Enter RDW
- TLC: Enter TLC
- PLT /mm3: Enter platelet count
- HGB: Enter HGB

A 'Predict Now' button is located at the bottom of the form.



XII. CONCLUSION

In conclusion, the project successfully demonstrated the potential of machine learning algorithms in improving the accuracy and efficiency of blood cancer prediction.

Through extensive data preprocessing, feature selection, and model training, a predictive model was created that showed promising results in identifying individuals at risk of blood cancer. The developed model showcased high accuracy and

robustness, contributing to early detection and diagnosis of the disease.

The successful implementation of the machine learning model opens up new possibilities for the field of blood cancer prediction and diagnosis. It provides a foundation for further research and development in this area, with the potential to impact clinical practice and contribute to advancements in healthcare.

XIII. FUTURE ENHANCEMENTS

Developing a real-time prediction tool or web-based application would enable healthcare professionals to utilize the model directly in clinical settings, supporting rapid decision-making and personalized treatment planning. Exploring deep learning architectures, such as Convolutional Neural Networks (CNNs) for imaging data or Recurrent Neural Networks (RNNs) for time-series medical data, could enhance the model's ability to learn complex patterns and improve prediction performance. Incorporating explainable AI techniques would make the model's predictions more transparent, helping medical professionals understand the rationale behind each prediction and increasing trust in AI-driven diagnostics.

review. *Biocybernetics and Biomedical Engineering*, 40(3), 929-960.

- [8] Rizwan, M., Hassan, W., Javaid, M., Anwar, S. M., & Qureshi, M. B. (2020). Blood cancer detection and classification using deep learning models. *International Journal of Advanced Computer Science and Applications*, 11(9), 1004-1011.
- [9] Gutiérrez-Arriola, J. M., López-Monteaquedo, F. E., Martínez-Rubio, F., Sotoca, J. M., & Bernal, J. J. (2018). Automatic classification of hematological microscopic images using pre-trained deep learning models. *IEEE Access*, 6, 49120-49131.

BIBLIOGRAPHY

- [1] Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., & Thrun, S. (2017). Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639), 115-118.
- [2] Pant, K., Vipparthi, S. K., Gupta, A., Acharya, A., Kaur, H., & Mittal, A. (2019). A review on machine learning approaches for blood cancer detection. *Journal of Medical Systems*, 43(3), 52.
- [3] Saini, A., Nair, M. S., Nagar, A., Aggarwal, A., Sethi, A., Bhatia, S., ... & Kaur, A. (2018). Leukemia detection using image processing and machine learning techniques: A review. *IEEE Access*, 6, 79351-79377.
- [4] Kalatharan, V., & Manavalan, R. (2020). A comprehensive review on automated blood cancer detection using machine learning techniques. *Journal of Medical Systems*, 44(12), 246.
- [5] Shikha, V., Kumar, S., Yadav, S., Chakraborty, C., Kumar, P., & Dass, S. C. (2021). A comprehensive review on blood cancer detection using machine learning techniques. *Current Medical Imaging Reviews*, 17(1), 44-62.
- [6] Sethi, A., Sharma, M., & Kumar, P. (2018). Automated detection of leukemia using microscopic blood images and machine learning classification techniques. *Journal of Medical Systems*, 42(8), 150.
- [7] Kumar, A., Dhiman, G., & Kamboj, S. (2020). Detection and classification of blood cells using machine learning approaches: A comprehensive