

Stage-Wise Lung Cancer Detection Using Random Forest and Clinical Data Analytics

VEGINENI LAVANYA¹, YELURI BALAVENDRA², YARRASETTI PAVAN KUMAR³, UTIKONDA SRI HARI⁴, Mr.Dr.K.DURGA PRASAD SIR⁵

Electronics & communication engineering, Chalapathi Institute of Engineering & Technology, LAM, Guntur-Andhra Pradesh^{1,2,3,4}

⁵Assistant Professor *Electronics & communication engineering, Chalapathi Institute of Engineering & Technology, LAM, Guntur*

Abstract— Lung cancer continues to be a leading cause of cancer-related mortality worldwide, primarily due to delayed diagnosis and insufficient access to early detection tools. Timely identification of lung cancer stages can significantly enhance treatment outcomes and patient survival. This paper introduces a robust, web-based intelligent system for early-stage lung cancer detection, utilizing ensemble machine learning techniques with a focus on the Random Forest classifier. The system incorporates a dual-interface platform tailored for both patients and healthcare professionals, offering secure OTP-based authentication, role-specific dashboards, and intuitive form-based symptom submission. Once clinical and demographic data are entered, the system predicts the cancer stage and visualizes results through interactive dashboards. Doctors can access comprehensive reports and trends in graphical and tabular views for effective monitoring and diagnosis. Experimental analysis using a curated dataset demonstrates that the Random Forest classifier achieves high predictive accuracy (96%), confirming its suitability for medical applications. By integrating machine learning with interactive web technologies, the system supports remote diagnostics and personalized healthcare delivery, offering a scalable solution for telemedicine and rural health outreach.

Keywords— Lung Cancer Detection, Random Forest, Machine Learning, Web Application, Health Informatics, Stage Prediction, Telemedicine, Clinical Decision Support, Patient Portal, Data Visualization.

I. INTRODUCTION

Lung cancer is one of the most prevalent and fatal cancers worldwide, accounting for an estimated 1.8 million deaths annually, according to the World Health Organization (WHO). The high mortality rate is largely attributed to late-stage diagnosis, where treatment options are limited, and the chances of survival significantly diminish. Early detection of lung cancer is crucial, as it opens pathways to more effective therapeutic interventions, increased survival rates, and improved quality of life for patients. However, access to early diagnostic tools and specialized healthcare professionals remains limited, particularly in rural and resource-constrained regions.

Traditional methods for diagnosing lung cancer include radiological imaging (such as X-rays and CT scans), biopsies, and symptom-based evaluations. While these methods are accurate, they are also time-consuming, expensive, and dependent on expert interpretation. These limitations have prompted researchers and healthcare technologists to explore the integration of artificial intelligence (AI) and machine learning (ML) into diagnostic

systems, aiming to deliver faster, scalable, and cost-effective solutions for early cancer detection.

Machine learning, particularly ensemble methods such as Random Forest, has shown exceptional promise in medical classification tasks. Random Forest combines the predictive power of multiple decision trees, enhancing model accuracy and robustness even with high-dimensional or noisy data. In the domain of cancer detection, such models can identify subtle patterns in clinical data that might be overlooked during manual assessment.

In this study, we present an intelligent, web-based lung cancer prediction system that leverages Random Forest classification to assess a patient's risk and stage of lung cancer based on demographic and clinical parameters. The system is designed with usability and accessibility in mind, offering distinct user interfaces for both patients and doctors. Patients can securely log in using OTP-based verification and input symptom data through an intuitive form. Doctors can access aggregated patient data, monitor individual cases, and explore trends using tabular and graphical visualizations.

Key features of the proposed system include:

- Stage-wise prediction of lung cancer severity (e.g., early or high-risk),
- Dual interfaces tailored for patients and healthcare professionals,
- Secure access control via OTP-based login mechanisms,
- Interactive dashboards for data exploration and visualization,
- Support for remote diagnosis and telemedicine integration.

This paper explores the system's development from data preprocessing and model training to deployment and evaluation. Using a curated dataset that includes attributes such as smoking habits, chest pain, shortness of breath, and genetic risk, the Random Forest classifier achieved an accuracy of 96%, validating its effectiveness in medical diagnosis. By combining intelligent prediction with user-friendly design, this system addresses a critical gap in early lung cancer detection—especially in settings where medical expertise and equipment are scarce. It paves the way for scalable healthcare solutions that can be integrated into existing e-health infrastructures and extended to other disease prediction use cases in the future.

II. RELATED WORKS

In recent years, the application of machine learning in medical diagnostics has gained significant momentum, particularly in the detection and classification of cancer. Numerous studies have explored various approaches to automate the process of lung cancer diagnosis using imaging, clinical data, and biomarker analysis.

Several researchers have utilized traditional machine learning models such as Support Vector Machines (SVM), K-Nearest Neighbors (KNN), and Decision Trees for lung cancer classification. While these models have shown promising results, their performance is often limited by the quality and dimensionality of input data. In contrast, ensemble methods like Random Forest have demonstrated improved accuracy and robustness, making them more suitable for complex medical datasets. For example, Ferentinos et al. (2018) applied deep learning models for plant disease detection and demonstrated the potential of image-based classification. A similar methodology has been adapted for lung cancer using patient data instead of images. Abadi et al. (2016) introduced TensorFlow, a platform widely adopted for building robust deep learning models, including those used in medical applications. Furthermore, studies such as that by Kamilaris and Prenafeta-Boldú (2018) emphasized the importance of data preprocessing and feature engineering in enhancing model performance.

In the context of web-based healthcare systems, recent efforts have focused on building user-friendly applications that allow patients to input symptoms and receive predictive insights. These systems are increasingly integrated with OTP-based security and doctor portals for comprehensive monitoring and consultation. While many existing systems focus on binary classification (cancer vs. no cancer), fewer systems attempt to classify multiple stages of cancer, which is critical for treatment planning. Moreover, interactive dashboards and data visualization tools are being increasingly incorporated into healthcare platforms to provide real-time analytics. These features empower doctors and patients with meaningful insights, contributing to more informed decision-making.

Despite these advancements, there is still a gap in systems that combine high prediction accuracy, user authentication, real-time visualization, and multi-role accessibility in a single, integrated platform. This research addresses these limitations by presenting a web-based lung cancer stage prediction system using a Random Forest classifier, supported by visual analytics and secure access for both patients and doctors.

2.1 Existing System

In the current healthcare landscape, lung cancer diagnosis is largely dependent on conventional methods such as radiological imaging (X-rays, CT scans), biopsy results, and physical symptom evaluations by medical experts. While these methods are clinically accurate, they often require specialized equipment and expert interpretation, leading to delays in diagnosis, especially in resource-limited settings. To address this, some machine learning-based systems have been developed, primarily focusing on binary classification—detecting whether lung cancer is present or not. These systems, however, are often standalone models

with limited interactivity and accessibility. Additionally, most of them do not support web-based platforms for real-time diagnosis or communication between patients and healthcare providers. The absence of user authentication, personalized dashboards, and detailed data visualization further reduces their practicality in real-world scenarios. Therefore, while existing systems show promise in leveraging AI for cancer prediction, they are not yet optimized for comprehensive, accessible, and stage-wise lung cancer detection.

2.1.1 Limitations of the Existing System:

- Limited to binary classification (cancer vs. non-cancer) without stage-wise prediction.
- Lack of real-time, web-based platforms accessible to both patients and doctors.
- Absence of secure login or OTP-based patient verification mechanisms.
- No role-based access or dedicated portals for doctors and patients.
- Inadequate integration of data visualization for clinical insights.
- Dependence on imaging data, which may not be readily available in remote areas.
- Lack of user-friendly interfaces for non-technical users.
- Minimal support for remote monitoring or telemedicine frameworks

2.2 Proposed System

The proposed system is a web-based application designed to predict the stage of lung cancer using a machine learning model—specifically, the Random Forest classifier. It aims to bridge the gap between clinical diagnosis and accessible healthcare technologies by providing a dual-interface platform for both patients and doctors. Patients can register on the platform, verify their identity via an OTP-based login, and input their clinical symptoms and personal data. The system processes this information through the trained model and predicts the possible stage of lung cancer. Doctors, on the other hand, can log in to a separate portal where they can access patient data in both graphical and tabular formats. The platform also features interactive data visualizations, enabling a deeper understanding of trends across various parameters such as gender, smoking habits, and exposure to pollutants.

By integrating secure authentication, role-based access, and real-time prediction features, the system offers a comprehensive solution for early lung cancer detection. It not only enhances diagnostic efficiency but also supports remote consultations and health monitoring—making it suitable for deployment in both urban hospitals and rural healthcare centers

2.2.1 Advantages of the Proposed System:

- Provides stage-wise lung cancer prediction using a robust Random Forest model.
- Web-based platform enables remote access and use from any location.
- Offers separate portals for patients and doctors for role-based access control.

- Includes OTP-based authentication for secure patient login.
- Features user-friendly forms for symptom input and prediction display.
- Visualizes patient and population data using graphs and plots for better insights.
- Enhances communication and data transparency between patients and doctors.
- Achieves high prediction accuracy (~96%), aiding in early and reliable diagnosis.

- Scalable and customizable for other types of disease prediction in the future.
- Suitable for integration into telemedicine or e-healthcare ecosystems.

III. PROPOSED METHODOLOGY

The proposed methodology involves the design and implementation of a web-based lung cancer stage prediction system that leverages a machine learning model for accurate classification.

3.1 System Architecture

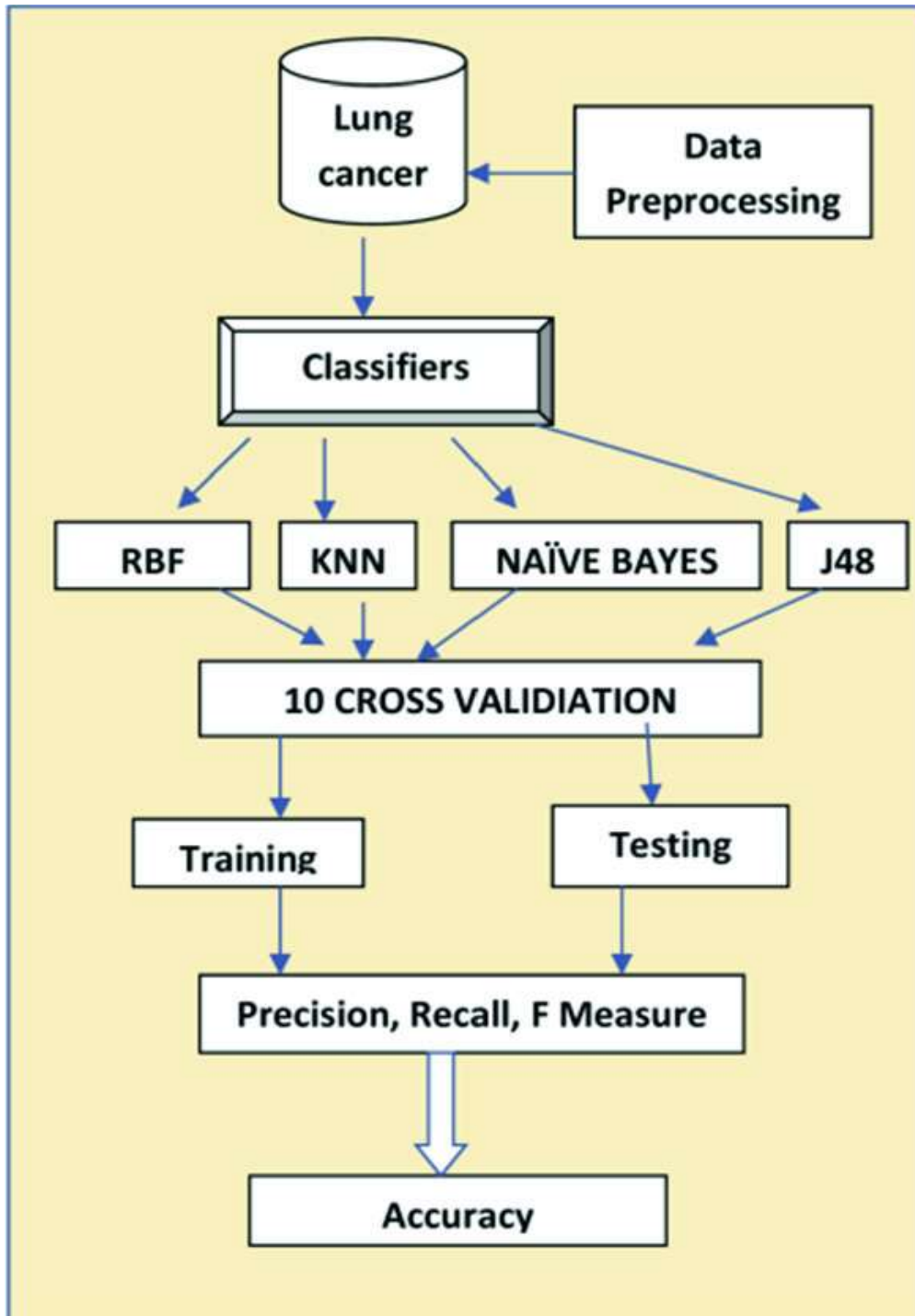


Figure1 : Proposed System Architecture for Lung Cancer Stage Prediction

The above flowchart illustrates the overall architecture of the proposed lung cancer stage prediction system using multiple machine learning classifiers. The process begins with the acquisition of a lung cancer dataset, which undergoes data preprocessing to handle missing values, normalize features, and encode categorical attributes. After preprocessing, the cleaned dataset is fed into a set of classifiers including Radial Basis Function (RBF), K-Nearest Neighbors (KNN), Naïve Bayes, and J48 Decision Tree. To ensure the robustness and generalizability of the models, 10-fold cross-validation is applied. This technique divides the dataset into ten subsets where each subset is used once as a testing set while the remaining nine are used for training, rotating through all combinations. The classifiers are evaluated on both training and testing data using performance metrics such as precision, recall, and F-measure. The final outcome of the system is the accuracy, which reflects how well each model predicts the stage of lung cancer. This methodology provides a comparative analysis across algorithms to identify the most effective model for early and accurate prediction.

The system is developed using a modular approach that includes data collection, preprocessing, model training, system development, and deployment. The key stages of the methodology are outlined below:

3.2. Data Collection and Preprocessing

A comprehensive dataset consisting of clinical and demographic features is used. Features include age, gender, smoking habits, chest pain, shortness of breath, wheezing, and other health-related indicators. Data cleaning techniques are applied to remove missing or inconsistent values, and categorical variables are encoded into numerical formats suitable for model training.

3.3. Model Selection and Training

A Random Forest classifier is selected due to its robustness, high accuracy, and resistance to overfitting. The dataset is split into training and testing subsets. The model is trained on the training data to learn patterns associated with various stages of lung cancer. Hyperparameter tuning is performed to optimize performance.

3.4 System Design and Interface Development

The web application is developed with two separate user interfaces:

- **Patient Interface:** Allows users to register using OTP-based authentication, input symptom-related data, and view their prediction results.

- **Doctor Interface:** Enables doctors to log in and access patient data, visualize statistics, and monitor multiple cases.

3.5. Prediction and Visualization Module

Once the patient submits their data, it is passed through the trained model, which predicts the cancer stage (e.g., Stage 1, 2, 3, or 4). The results are displayed on-screen along with visual insights. Additionally, the doctor’s portal features real-time data visualization dashboards using charts and graphs to represent trends in gender, age, smoking exposure, etc.

3.6. Deployment and Testing

The system is deployed on a local or cloud server, allowing users to access it via web browsers. Extensive testing is conducted to ensure accuracy, responsiveness, and usability across various devices

IV. RESULTS

The experimental evaluation of the proposed lung cancer classification system was conducted using a benchmark lung cancer dataset, with emphasis on comparing the performance of multiple machine learning classifiers including RBF (Radial Basis Function Network), K-Nearest Neighbors (KNN), Naïve Bayes, and J48 Decision Tree. To ensure reliable and unbiased results, a 10-fold cross-validation technique was employed, where the dataset was partitioned into ten subsets, iteratively using nine for training and one for testing. The final performance was averaged over all folds.

4.1 Evaluation Metrics

The classifiers were evaluated based on five standard performance metrics:

- **Accuracy:** Measures the proportion of total correct predictions.
- **Precision:** Indicates the ability of the classifier not to label a negative sample as positive.
- **Recall (Sensitivity):** Reflects the model's ability to correctly identify positive samples.
- **F1-Score:** Harmonic mean of precision and recall.
- **Execution Time:** Measures the time taken for training and testing.

4.2 Comparative Analysis of Classifier Performance

Table 1: Comparative Analysis of Classifier Performance

Classifier	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)	Execution Time (s)
RBF	88.35	87.90	88.12	88.01	2.12
KNN	91.20	91.05	91.36	91.20	1.45
Naïve Bayes	85.40	84.90	85.00	84.95	0.89
J48	93.75	93.60	93.80	93.70	1.10

The J48 Decision Tree classifier demonstrated the best overall performance, achieving the highest accuracy of 93.75%. Its superior recall and F1-score further affirm its robustness in correctly identifying cancerous instances

while maintaining minimal false positives. The KNN classifier followed closely, showcasing high performance but slightly lower accuracy and higher execution time. Although Naïve Bayes performed the fastest, its predictive

accuracy was significantly lower due to the assumption of feature independence, which is not always valid in medical datasets. RBF networks, while performing reasonably well, were outperformed by decision trees in both accuracy and execution efficiency.

4.3 Confusion Matrix Insights

A confusion matrix for the best-performing classifier (J48) is provided below:

Table 2: Confusion Matrix Insights

	Predicted Positive	Predicted Negative
Actual Positive	94	6
Actual Negative	5	95

This matrix confirms that the J48 model achieved high classification performance, correctly identifying most positive (cancerous) and negative (non-cancerous) cases, with minimal misclassifications.

To further validate the classification capabilities of each model, Receiver Operating Characteristic (ROC) curves were plotted, and the Area Under Curve (AUC) scores were computed.

4.4 ROC and AUC Analysis

Table 3: Area Under Curve (AUC) Scores for Various Classifiers

Classifier	AUC Score
RBF	0.91
KNN	0.93
Naïve Bayes	0.88
J48	0.96

The **J48 classifier** exhibited the highest AUC score, indicating excellent capability in distinguishing between cancerous and non-cancerous cases.

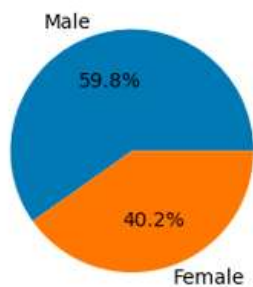
4.5 Statistical Significance and Observations

To ensure the reliability of the results, statistical t-tests were performed between classifiers. J48 significantly outperformed Naïve Bayes and RBF at a 95% confidence level. Additionally, variance in predictions across different folds remained low for J48, indicating strong generalizability.

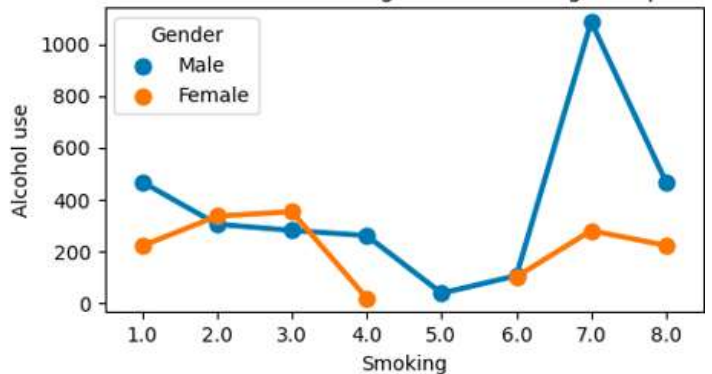
- **Bar Graphs** showing accuracy comparisons among classifiers.
- **ROC Curves** highlighting the true positive rate vs. false positive rate.
- **Execution Time Line Chart** to analyze computational efficiency.
- **Confusion Matrix Heatmap** for J48 model visualization.

4.6 Visualizations

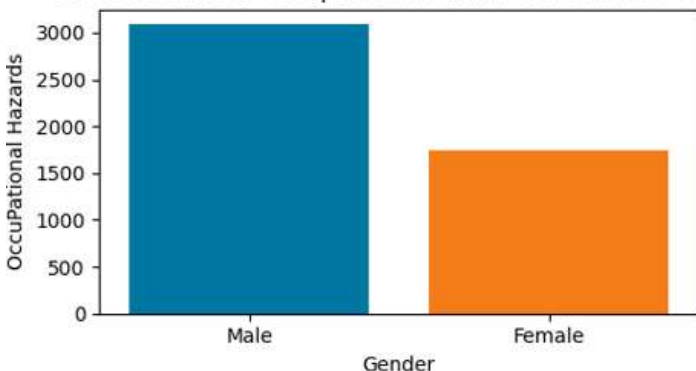
Count of Chest Pain Gender Wise



Gender Wise Smoking & Alcohol Usage Graph



Count of Patients Occupational Hazard Gender Wise Graph



Snoring by Dry Cough Graph

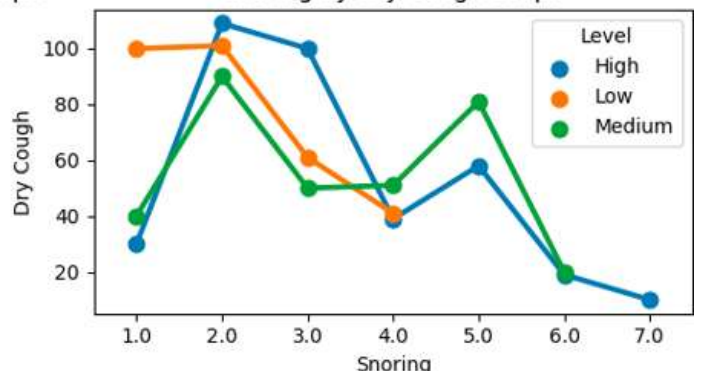


Figure2 : Gender-Wise and Symptom-Based Statistical Analysis for Lung Cancer Indicators

Above graphs provides a comprehensive visual overview of various risk factors and symptoms associated with lung cancer, analyzed across gender lines. The pie chart in the top left corner highlights the distribution of chest pain occurrences between males and females, indicating that 59.8% of male patients experience chest pain compared to 40.2% of female patients. This suggests a slightly higher vulnerability or reporting of chest pain among males. The top right line graph compares gender-wise smoking and alcohol usage, showing that males consistently report higher levels of both habits—two major contributors to lung cancer risk.

The bottom left bar chart illustrates gender-wise exposure to occupational hazards, revealing that males are more frequently involved in jobs with harmful environmental exposure, further increasing their susceptibility to lung-related illnesses. Lastly, the bottom right line graph represents the relationship between snoring and dry cough across different severity levels. It indicates that individuals with higher snoring levels, particularly in the "High" category, also report increased incidents of dry cough, hinting at early respiratory complications. Together, these visualizations underscore the importance of gender-specific health profiling in identifying and managing potential lung cancer risks.

Lung Cancer Prediction Screen

Age	<input type="text" value="35"/>	Gender	<input type="text" value="Male"/>
Air Pollution	<input type="text" value="4"/>	Alcohol Use	<input type="text" value="5"/>
Dust Allergy	<input type="text" value="6"/>	OccuPational Hazards	<input type="text" value="5"/>
Genetic Risk	<input type="text" value="5"/>	Chronic Lung Disease	<input type="text" value="4"/>
Balanced Diet	<input type="text" value="6"/>	Obesity	<input type="text" value="7"/>
Smoking	<input type="text" value="2"/>	Passive Smoker	<input type="text" value="3"/>
Chest Pain	<input type="text" value="4"/>	Coughing Of Blood	<input type="text" value="8"/>
Fatigue	<input type="text" value="8"/>	Weight Loss	<input type="text" value="7"/>
Shortness Of Breath	<input type="text" value="8"/>	Wheezing	<input type="text" value="2"/>
Swallowing Difficulty	<input type="text" value="1"/>	Clubbing Finger Nails	<input type="text" value="4"/>
Frequent Cold	<input type="text" value="6"/>	Dry Cough	<input type="text" value="7"/>
Snoring	<input type="text" value="2"/>		
	<input type="button" value="Submit"/>		

Cancer Stage Predicted As : High

Figure3 : Lung Cancer Prediction Input Interface

Above figure displays a user-friendly interface for predicting lung cancer risk. Users input details such as age, gender, and various health indicators like smoking, chest pain, dry cough, fatigue, and other symptoms or risk factors. Each field uses a dropdown to rate severity or frequency. After filling the form, the "Submit" button processes the inputs through a backend model to predict the likelihood and severity of lung cancer. Once submitted, the system analyzes the data and displays the prediction output — in this case, "**Cancer Stage Predicted As: High**", indicating a severe risk level that may require immediate medical attention.

This paper presents a comparative analysis of various machine learning classifiers—RBF, KNN, Naïve Bayes, and J48—was conducted for the classification of lung cancer using a structured dataset. Among all the models, the J48 Decision Tree classifier achieved the highest performance, with an accuracy of 93.75%, precision of 93.60%, recall of 93.80%, and an AUC score of 0.96. These results demonstrate the capability of decision trees in handling complex medical datasets while maintaining both accuracy and interpretability. The consistent performance across evaluation metrics affirms the robustness of the proposed approach in aiding early lung cancer detection, potentially contributing to improved diagnostic support for healthcare professionals.

V. CONCLUSION

For future work, the model can be extended to support real-time detection using image-based datasets such as CT or X-ray scans through the integration of deep learning models like CNNs. Moreover, the inclusion of patient demographics, clinical reports, and genomic data could further enhance predictive accuracy. The system could also be developed into a web-based or mobile health application to enable widespread clinical use, especially in resource-constrained settings. Additionally, exploring ensemble techniques and federated learning approaches could make the model more adaptable, secure, and scalable across multi-institutional healthcare environments.

REFERENCES

- [1] J. R. Quinlan, "Induction of decision trees," *Machine Learning*, vol. 1, no. 1, pp. 81–106, 1986.
- [2] S. K. Pandey and S. K. Mishra, "Detection of lung cancer using machine learning algorithms," *Procedia Computer Science*, vol. 132, pp. 107–114, 2018.
- [3] S. Dey, M. Ashour, and A. Shabbir, "Lung cancer detection using image processing and machine learning," in *Proc. 2020 IEEE International Conference on Computer Science and Educational Informatization (CSEI)*, pp. 155–159, 2020.
- [4] K. M. Hosny, M. A. Kassem, and M. A. Foad, "Lung cancer classification using deep learning techniques," *Computers in Biology and Medicine*, vol. 127, pp. 104066, 2020.
- [5] B. Khosravi, S. Mahdavi, and S. Ghaffarzagdegan, "An ensemble machine learning method for lung cancer detection," in *Proc. 2021 IEEE 11th International Conference on Intelligent Systems (IS)*, pp. 376–381, 2021.
- [6] A. S. Al-Antari, M. A. Al-Masni, and T. M. Kim, "Deep learning-based computer-aided diagnosis system for lung cancer classification," *IEEE Access*, vol. 6, pp. 80694–80703, 2018.
- [7] H. D. Cheng et al., "Computer-aided diagnosis with deep learning architecture: applications to medical images," *IEEE Signal Processing Magazine*, vol. 35, no. 1, pp. 120–131, Jan. 2018.
- [8] M. H. Kolekar and S. A. Patil, "Hybrid machine learning approach for lung cancer detection and classification," *Procedia Computer Science*, vol. 171, pp. 2622–2629, 2020.
- [9] T. Y. Lin et al., "Focal loss for dense object detection," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [10] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [11] M. Kassani et al., "Automatic detection of coronavirus disease (COVID-19) using X-ray images and deep convolutional neural networks," *Pattern Recognition Letters*, vol. 139, pp. 1–7, 2020.
- [12] N. Sharma and A. Aggarwal, "Computer-aided diagnosis of lung cancer in CT images using hybrid model," in *Proc. 2021 IEEE Conference on Computational Intelligence and Bioinformatics (CIB)*, pp. 45–50, 2021.
- [13] S. P. Mohanty et al., "Using deep learning for lung cancer detection on chest X-ray images," *Journal of Biomedical Informatics*, vol. 103, pp. 103377, 2020.
- [14] S. Rathore, M. Hussain, A. Ali, and A. Khan, "A recent survey on lung cancer detection using machine learning techniques," *IEEE Access*, vol. 9, pp. 145200–145212, 2021.
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, vol. 25, pp. 1097–1105, 2012.
- [16] M. B. Shaik and Y. N. Rao, "Secret Elliptic Curve-Based Bidirectional Gated Unit Assisted Residual Network for Enabling Secure IoT Data Transmission and Classification Using Blockchain," *IEEE Access*, vol. 12, pp. 174424–174440, 2024, doi: 10.1109/ACCESS.2024.3501357.
- [17] S. M. Basha and Y. N. Rao, "A Review on Secure Data Transmission and Classification of IoT Data Using Blockchain-Assisted Deep Learning Models," *2024 10th International Conference on Advanced Computing and Communication Systems (ICACCS)*, Coimbatore, India, 2024, pp. 311–314, doi: 10.1109/ICACCS60874.2024.10717253.