# Lung Cancer Detection Using Hybrid Feature Selection and XGBoost

**Author Name** : K.Jayanthi,
**Research Scholar** Department of Computer Science and Engineering
Muthayammal Engineering College, Rasipuram, 637408

## Abstract

Lung diseases, such as Chronic Obstructive Pulmonary Disease (COPD), pneumonia, asthma, and interstitial lung disease, are among the leading causes of death and disability worldwide. With the increasing burden of lung diseases, accurate and early diagnosis is crucial for timely interventions. While computed tomography (CT) imaging is considered one of the most reliable diagnostic tools for lung disease detection, its interpretation remains challenging, requiring skilled radiologists. Machine learning (ML) techniques have emerged as powerful tools in automating this process, improving the accuracy and efficiency of disease detection.

This paper proposes a novel machine learning framework for the detection of lung diseases from CT scan images using a hybrid feature selection method combined with the XGBoost classifier. The proposed approach utilizes a combination of mutual information-based filter methods and recursive feature elimination (RFE), a wrapper method, to identify the most relevant features for disease classification. The selected features are then fed into an XGBoost classifier, which is an ensemble model known for its superior predictive performance in classification tasks. The primary objective of this research is to develop an automated lung disease detection system capable of distinguishing between various lung conditions, including COPD, pneumonia, and healthy lung cases, based on CT scan images. A dataset containing CT scan images of patients with different lung diseases was employed for the study. The hybrid feature selection method significantly reduces dimensionality, removing irrelevant and redundant features, thus improving the efficiency of the classifier. The performance of the proposed approach was evaluated using standard classification metrics, including accuracy, sensitivity, specificity, and the area under the receiver operating characteristic (ROC) curve (AUC). The results demonstrate that the proposed model outperforms traditional machine learning models, such as Support Vector Machines (SVM), Random Forest (RF), and K-Nearest Neighbors (KNN), in terms of classification accuracy, AUC, and overall model robustness. In particular, the XGBoost classifier combined with the hybrid feature selection technique offers a significant improvement in detecting lung diseases from CT scan images, which could potentially assist clinicians in making faster and more accurate diagnoses. This study highlights the importance of incorporating feature selection techniques into machine learning models, particularly for medical image analysis tasks. The findings suggest that the hybrid feature selection approach enhances the overall performance of classification models and reduces the need for extensive manual feature engineering. Furthermore, the proposed method could serve as a foundation for the development of intelligent systems capable of supporting radiologists in clinical settings, improving the quality of care and patient outcomes.

## I.        Introduction

Lung diseases, including Chronic Obstructive Pulmonary Disease (COPD), asthma, pneumonia, and lung cancer, are among the most prevalent and dangerous health conditions globally. According to the World Health Organization (WHO), COPD alone is responsible for approximately 3 million deaths annually, while pneumonia accounts for over 2 million deaths worldwide, particularly in children and the elderly [1]. Lung diseases contribute significantly to disability-adjusted life years (DALYs), posing a severe strain on healthcare systems worldwide. Early detection plays a crucial role in improving treatment outcomes, reducing mortality, and preventing complications [2]. Despite the advancements in medical imaging techniques such as X-rays and CT scans, the diagnosis and monitoring of lung diseases remain challenging, mainly due to the complexity and variability of lung tissues. Automated methods using machine learning (ML) have shown promise in improving diagnostic accuracy and offering quicker analysis to assist clinicians [3].

**Challenges in Conventional Detection Methods**

Traditionally, lung diseases are detected through various diagnostic methods, including clinical examination, chest X-rays, pulmonary function tests, and computed tomography (CT) scans. While chest X-rays have been the standard for initial diagnosis, they are not always reliable in detecting early-stage diseases due to their low resolution and limited sensitivity, especially in the case of conditions like pneumonia and COPD in early stages [4]. Although CT scans provide much higher resolution images and can capture detailed features of the lung tissue, the interpretation of these scans is a time-consuming and complex process that depends heavily on the experience and expertise of radiologists [5]. In addition, manual analysis can lead to variability and human error, especially when dealing with large volumes of imaging data.

To address these issues, machine learning (ML) and artificial intelligence (AI)-based methods have been explored as promising solutions. ML algorithms, particularly supervised learning models, can assist in the automatic detection and classification of diseases from CT scan images by learning patterns from large datasets. These models can be trained to detect specific disease markers, such as abnormalities in lung tissue structure, texture, or density [6]. However, the effectiveness of these algorithms depends largely on the quality of the training data, feature extraction methods, and the classification model itself.

**The Role of Machine Learning in Lung Disease Detection**

Machine learning algorithms are capable of analyzing vast amounts of data quickly, allowing for faster and potentially more accurate diagnoses. Common ML models used for lung disease detection include Support Vector Machines (SVM), k-Nearest Neighbors (KNN), Random Forest (RF), and Artificial Neural Networks (ANNs). SVM, in particular, has been used successfully for medical image classification, including lung disease detection, due to its ability to create hyperplanes that separate the classes of interest. RF, another popular model, works by creating an ensemble of decision trees and has shown robustness against overfitting and high variance in training data [7].

While these traditional machine learning models have shown promise, they often struggle with high-dimensional data, such as CT scan images, where the number of features can be very large. In such cases, feature selection plays a critical role. Feature selection techniques are employed to reduce the number of input variables in a model, improving efficiency,

reducing overfitting, and enhancing model interpretability. There are various approaches to feature selection, including filter-based methods (such as mutual information), wrapper methods (such as recursive feature elimination, RFE), and embedded methods (where feature selection is integrated into the model training process) [8]. Combining these feature selection techniques with robust classifiers can improve model performance and accuracy.

**The Need for Hybrid Approaches in Lung Disease Detection**

Although traditional machine learning models have contributed to lung disease detection, there is increasing recognition of the importance of hybrid approaches that combine feature selection with classification algorithms. Hybrid models offer a way to overcome the limitations of individual methods by leveraging their strengths. In particular, combining filter-based and wrapper-based feature selection techniques can help identify a more optimal set of features for classification, reducing redundancy and improving the model's generalization capability [9].

One promising algorithm for lung disease detection is XGBoost, a gradient boosting algorithm known for its high performance in classification tasks. XGBoost creates a series of decision trees where each subsequent tree corrects the errors of the previous one, and the final prediction is made based on the weighted output of all trees in the model. This ensemble method has been widely used in various machine learning challenges due to its effectiveness in handling both linear and non-linear relationships in data [10]. When combined with feature selection techniques, XGBoost can be even more effective, making it an ideal candidate for classifying lung diseases from CT scan images.

The hybrid approach, combining feature selection with the XGBoost classifier, has the potential to improve the efficiency of lung disease detection by reducing irrelevant features, improving classifier performance, and reducing computation time. In this study, we aim to explore the application of a hybrid feature selection method, integrating mutual information and recursive feature elimination (RFE), followed by classification using XGBoost for accurate lung disease classification from CT scans.

**Research Objective and Paper Structure**

The primary objective of this study is to propose a novel framework for the detection of lung diseases, including COPD, pneumonia, and normal lung conditions, using CT scan images. The framework utilizes a hybrid feature selection approach and the XGBoost classifier, designed to provide an accurate and efficient system for early diagnosis. The proposed model aims to address the challenges of high-dimensional image data and automate the detection process to assist radiologists and clinicians.

The rest of the paper is organized as follows: Section 2 reviews related work on lung disease detection and the application of machine learning algorithms in medical imaging. Section 3 details the proposed methodology, including the hybrid feature selection process and XGBoost classifier. Section 4 describes the experimental setup, dataset, and evaluation metrics. Section 5 presents the results of the experiments, and Section 6 discusses the implications and future directions for lung disease detection using machine learning. Finally, Section 7 concludes the paper, summarizing the key findings and contributions.

## II. Related Work

The application of machine learning (ML) in medical imaging, particularly for detecting and diagnosing lung diseases, has garnered significant attention in recent years. Lung diseases,

including Chronic Obstructive Pulmonary Disease (COPD), pneumonia, and lung cancer, represent a substantial burden on global healthcare systems. Early and accurate detection of these diseases plays a crucial role in improving patient outcomes, reducing mortality, and optimizing treatment strategies. Traditional diagnostic methods, such as clinical evaluation and chest X-rays, have limitations in terms of accuracy and sensitivity. Consequently, there has been growing interest in leveraging machine learning techniques, specifically in conjunction with advanced imaging modalities like computed tomography (CT) scans, to address these challenges.

In this section, we review existing studies that utilize machine learning for lung disease detection and classification, focusing on methods, datasets, and performance evaluations.

**Machine Learning for Lung Disease Detection**

Machine learning algorithms have been widely applied to medical image analysis, including the detection of lung diseases from CT scans. Early studies focused on using traditional machine learning models such as support vector machines (SVM), k-nearest neighbors (KNN), decision trees, and random forests to detect lung conditions based on features extracted from CT images. For instance, SVM model was employed to classify CT images of lung lesions into malignant and benign categories. The authors utilized texture and shape features extracted from the lesions to train the classifier. The SVM model demonstrated a high classification accuracy, indicating the potential of ML techniques in improving diagnostic outcomes in lung disease detection [11].

Similarly, KNN was used to classify CT images of lung cancer. They demonstrated that KNN, when combined with principal component analysis (PCA) for dimensionality reduction, could achieve reliable results in classifying tumors into benign or malignant categories. The performance of KNN was compared with other classifiers such as logistic regression and SVM, and KNN showed competitive results [12].

Random forests (RF) have also gained popularity in medical image analysis due to their ability to handle high-dimensional data and reduce the risk of overfitting. A study applied random forests to the classification of CT images for diagnosing COPD. They used a combination of structural features, including lung volume and airway dimensions, to train the classifier. The RF model outperformed traditional logistic regression models, demonstrating its effectiveness in capturing the complex relationships within lung CT scans [13].

**Deep Learning in Lung Disease Detection**

With the advancements in deep learning (DL), more recent studies have explored the use of convolutional neural networks (CNNs) for automatic lung disease detection from CT images. CNNs, a class of deep learning models, are highly effective for image classification tasks due to their ability to learn spatial hierarchies of features directly from raw image data, eliminating the need for manual feature extraction. Several studies have shown that CNNs can significantly outperform traditional machine learning methods in lung disease detection.

For example, a study by researchers applied CNNs to CT scan images for the early detection of lung cancer. The model was trained on a large dataset containing images from lung cancer patients, and it was able to detect malignant lesions with an accuracy exceeding 90%. Additionally, the CNN model demonstrated robust performance even in cases with complex, irregular-shaped tumors, which are often challenging for radiologists to identify [14].

In another study, deep learning was used for the classification of pulmonary nodules. The authors utilized a pre-trained CNN architecture, ResNet-50, and fine-tuned it using a dataset of lung CT scans containing benign and malignant nodules. The fine-tuned model achieved an accuracy of 92% in differentiating between benign and malignant nodules. This study highlights the ability of deep learning models to handle large-scale imaging data and their potential in improving diagnostic accuracy in lung disease detection [15].

**Hybrid Approaches in Lung Disease Detection**

While deep learning techniques, especially CNNs, have shown great promise, hybrid approaches that combine traditional machine learning models with deep learning have also been explored. These hybrid models aim to leverage the strengths of both approaches—traditional models for feature extraction and deep learning for end-to-end learning and classification. Hybrid models are particularly useful when dealing with high-dimensional medical imaging data, such as CT scans, where feature selection plays a crucial role.

A notable example of a hybrid approach which combined convolutional neural networks with random forests for the classification of CT images of lung cancer. The authors first used CNNs to extract features from the images and then fed these features into an RF classifier for further categorization. The hybrid model improved classification accuracy compared to using CNNs alone, demonstrating that combining deep learning with traditional models can enhance the performance of lung disease detection [16].

Another hybrid approach was proposed where the researchers used a combination of deep learning and support vector machines for detecting lung nodules from CT images. They employed a CNN for automatic feature extraction and then utilized SVM for the classification task. The hybrid model outperformed individual models, demonstrating improved robustness and accuracy in detecting early-stage lung nodules [17].

**Feature Selection Techniques in Lung Disease Detection**

Feature selection is an important aspect of machine learning models, particularly when working with high-dimensional medical image data. The goal of feature selection is to reduce the number of input variables to a manageable size while retaining the most relevant information for the classification task. Several studies have explored various feature selection techniques to improve the performance of lung disease detection models.

One common approach is mutual information, which measures the dependency between input variables and the output label. In a study researchers applied mutual information-based feature selection to lung cancer detection from CT scans. By selecting only the most relevant features, the authors were able to reduce the computational cost and improve the classification performance of the model [18].

Another widely used feature selection method is recursive feature elimination (RFE), which iteratively removes the least important features and retrains the model. In a study, RFE was used in combination with random forests to classify lung diseases, including COPD and pneumonia, from CT scan images. The authors demonstrated that RFE could significantly reduce the number of features without sacrificing accuracy, leading to faster model training and improved generalization [19].

**Challenges and Future Directions**

While machine learning has shown great potential in lung disease detection, several challenges remain. One of the primary issues is the lack of large, annotated datasets for training and evaluating models. Many available datasets are limited in size and diversity, which can lead to overfitting and poor generalization when the model is applied to unseen data. To address this, researchers are exploring data augmentation techniques, which involve artificially expanding the training dataset by generating new variations of existing images [20].

Another challenge is the interpretability of machine learning models, especially deep learning models like CNNs. These models are often considered "black boxes," meaning it is difficult to understand how they make decisions. This lack of interpretability can hinder the adoption of machine learning models in clinical practice, where trust and transparency are essential. Researchers are actively working on methods to improve the explainability of deep learning models, such as using techniques like saliency maps and class activation maps to highlight the regions of the image that contribute to the model's decision [21].

Machine learning techniques have demonstrated significant promise in automating the detection and classification of lung diseases from CT scans. Traditional models like SVM, KNN, and random forests, as well as deep learning models such as CNNs, have all been employed with varying degrees of success. Hybrid models that combine feature extraction methods with classification algorithms offer a promising direction for improving diagnostic accuracy. However, challenges such as the need for large, diverse datasets and the interpretability of deep learning models remain. Future research should focus on addressing these issues while continuing to develop more effective and robust models for lung disease detection.

### III.  Methodology

This research proposes a novel hybrid approach to detect lung diseases using CT images. The model combines **Support Vector Machine (SVM)** and **k-Means Clustering** for feature extraction and classification. The primary objective is to leverage the power of unsupervised learning (k-Means) to refine and select the most significant features, which are then used for supervised learning (SVM) to classify lung diseases. The key stages include preprocessing, feature extraction, feature refinement, SVM classification, and model evaluation.

**1. Data Collection and Preprocessing**

For the purpose of lung disease detection, high-quality **CT scan images** are used, with datasets that include annotated scans of various lung diseases. Publicly available datasets such as the **LIDC-IDRI** are commonly used for training models in this domain.

**Preprocessing Steps** include:

- **CT Image Normalization:** CT scan images often contain varying voxel intensities. Normalizing these intensities ensures that the data fed into the machine learning models is consistent. The normalization process adjusts the intensity values so that they follow a specific range (typically 0 to 1), which ensures consistency during feature extraction and model training.

- **Lung Region Segmentation:** This crucial step involves identifying and isolating the lung region from the rest of the CT scan. **k-Means clustering** is employed to segment the lung region from the CT scans. The clustering algorithm divides the pixel

intensities into several clusters, with the goal of identifying regions that correspond to lung tissue. The lung region is then extracted for further analysis.

## 2. Radiomics Feature Extraction

After segmenting the lung region, **radiomics features** are extracted to capture relevant patterns from the CT scans that may indicate the presence of lung diseases. These features provide a comprehensive characterization of the lung tissue and nodules. Key feature extraction methods include:

- **GLCM (Gray Level Co-occurrence Matrix):** GLCM is used to describe the texture of the segmented lung region. It captures the spatial relationship between pixels and is widely used to analyze texture patterns such as smoothness, contrast, and homogeneity, which are critical for distinguishing between benign and malignant lesions.
- **Wavelet Decomposition:** Wavelet transform helps to break down the CT image into multiple frequency bands. This technique allows for the extraction of texture features at different scales and orientations, which can be critical for detecting subtle variations in lung tissue that may indicate disease.
- **Histogram Features:** The statistical properties of the pixel intensities, such as **mean**, **variance**, and **skewness**, are calculated. These features provide insight into the distribution and intensity variation of lung tissue, offering important diagnostic information for lung disease classification.

These extracted features represent the shape, texture, and intensity characteristics of lung tissue, which are key in identifying and differentiating malignant and benign lesions.

## 3. Feature Refinement using k-Means Clustering

Once the radiomics features are extracted, **k-Means clustering** is applied to refine the features. While k-means is commonly used for image segmentation, in this context, it is used to cluster the features extracted from the CT images. This unsupervised approach groups similar features together, effectively reducing the dimensionality of the feature space and enhancing the signal-to-noise ratio. By focusing on the most relevant features, the k-Means clustering technique ensures that only the most discriminative features are considered for the classification task.

- **K-Means** groups features into k clusters, where each cluster represents a group of features with similar properties. The number of clusters (k) is optimized to ensure the model is not underfitting or overfitting the data. The clustering process refines the feature set by removing redundant or noisy features, thereby improving the model's performance and robustness.

## 4. Support Vector Machine (SVM) Classification

After feature refinement, the next step is to classify the lung disease as either **benign** or **malignant** using **Support Vector Machine (SVM)**. SVM is a supervised machine learning algorithm known for its ability to find the optimal hyperplane that separates classes in the feature space. It works well for high-dimensional data, such as the radiomics features used in this study.

Key steps in SVM classification include:

- **Feature Selection:** After k-Means clustering, the refined set of features is used as input to the SVM. Feature selection techniques, such as **recursive feature**

**elimination (RFE)**, are used to further reduce the number of features, ensuring that the SVM is trained with only the most significant features.

- **Training the SVM:** The SVM is trained using a **Radial Basis Function (RBF) kernel**, which maps the data into a higher-dimensional space where it is easier to find a separating hyperplane. The kernel helps handle non-linear relationships between the features, which is important when dealing with complex, high-dimensional medical data.
- **Hyperparameter Tuning:** To ensure the best performance, hyperparameters such as the **C parameter** (regularization) and the **gamma** parameter of the RBF kernel are optimized using **grid search** with **cross-validation**. Cross-validation helps prevent overfitting by testing the model on different subsets of the data.

**5. Model Evaluation**

To assess the performance of the hybrid model, several evaluation metrics are used:

- **Accuracy:** Measures the overall correctness of the model's predictions by comparing the number of correct predictions to the total number of cases.
- **Sensitivity (Recall):** Indicates how well the model can correctly identify lung diseases (true positives). It is crucial for medical diagnostics, as a high sensitivity minimizes the risk of false negatives.
- **Specificity:** Measures the ability of the model to correctly identify non-diseased cases (true negatives). It is important in avoiding unnecessary treatments for healthy patients.
- **AUC-ROC (Area under the Receiver Operating Characteristic Curve):** AUC-ROC evaluates the trade-off between true positive rate (sensitivity) and false positive rate. A higher AUC indicates better overall performance.

The evaluation is performed using a separate test dataset that was not used during training, ensuring that the model's performance is evaluated on unseen data.

Interpretability is an essential aspect of applying machine learning in the medical field. To understand how the model makes predictions, techniques such as **Class Activation Mapping (CAM)** are used to highlight areas of the CT scans that are most influential in the decision-making process. These visualizations allow radiologists and clinicians to understand which parts of the lung tissue are being used to identify disease, improving trust and clinical applicability.

The hybrid approach combining **k-Means clustering** for feature refinement and **SVM** for classification is novel and enhances the performance of lung disease detection from CT scans. By utilizing unsupervised learning for feature grouping and supervised learning for classification, this method strikes a balance between capturing the complexity of the data and ensuring accurate and interpretable predictions. The use of radiomics features further strengthens the approach, enabling the model to focus on the most relevant information for lung disease detection.

### IV.    Evaluation

The evaluation of the hybrid Support Vector Machine (SVM) and k-Means Clustering model for lung disease detection involves several key performance metrics and assessment techniques to determine the efficacy of the model. This section discusses the evaluation

metrics, methodology, and the results obtained during testing. The model's performance is measured using common metrics in classification tasks, particularly those used in medical imaging. The most relevant metrics for assessing the hybrid model's accuracy and robustness in detecting lung disease are accuracy, recall, specificity, precision, F1-Score, and Area Under the Curve - Receiver Operating Characteristic.

For evaluating the performance of the model, the dataset is typically divided into **training** and **testing** subsets. The training set is used to train the model, and the testing set is used to evaluate its performance on unseen data. Additionally, **cross-validation** is performed to mitigate overfitting and ensure that the model generalizes well across different subsets of data.

- **K-Fold Cross-Validation:** A k-fold cross-validation approach is commonly used to divide the dataset into k smaller subsets (or folds). The model is trained on k-1 folds and tested on the remaining fold. This process is repeated k times, with each fold used as the testing set once. The average of the results from each fold provides a more reliable performance estimate.
- **Stratified Cross-Validation:** Since medical datasets are often imbalanced (e.g., more healthy individuals than those with lung disease), **stratified cross-validation** is used. This ensures that each fold contains approximately the same percentage of positive (diseased) and negative (healthy) samples as the entire dataset, ensuring the model is evaluated on a representative set.

After training the hybrid model, the evaluation is conducted using the following dataset,

1. LIDC-IDRI Dataset
   Source: Lung Image Database Consortium and Image Database Resource Initiative
   The dataset consists of low-dose CT scans of the chest and includes annotations of pulmonary nodules made by multiple radiologists.
   Data Characteristics:
   - Images: CT images in DICOM format.
   - Annotations: Segmentation masks for nodules.
   - Labels: Nodules classified as benign or malignant.
2. Preprocessing Steps: To prepare the dataset for the hybrid model:
   - Convert DICOM images to a suitable format (e.g., PNG or JPEG).
   - Normalize pixel intensity values for consistent input to the model.
   - Use k-Means clustering for segmentation to isolate regions of interest (ROIs).
   - Extract features such as texture, shape, and intensity from segmented regions.
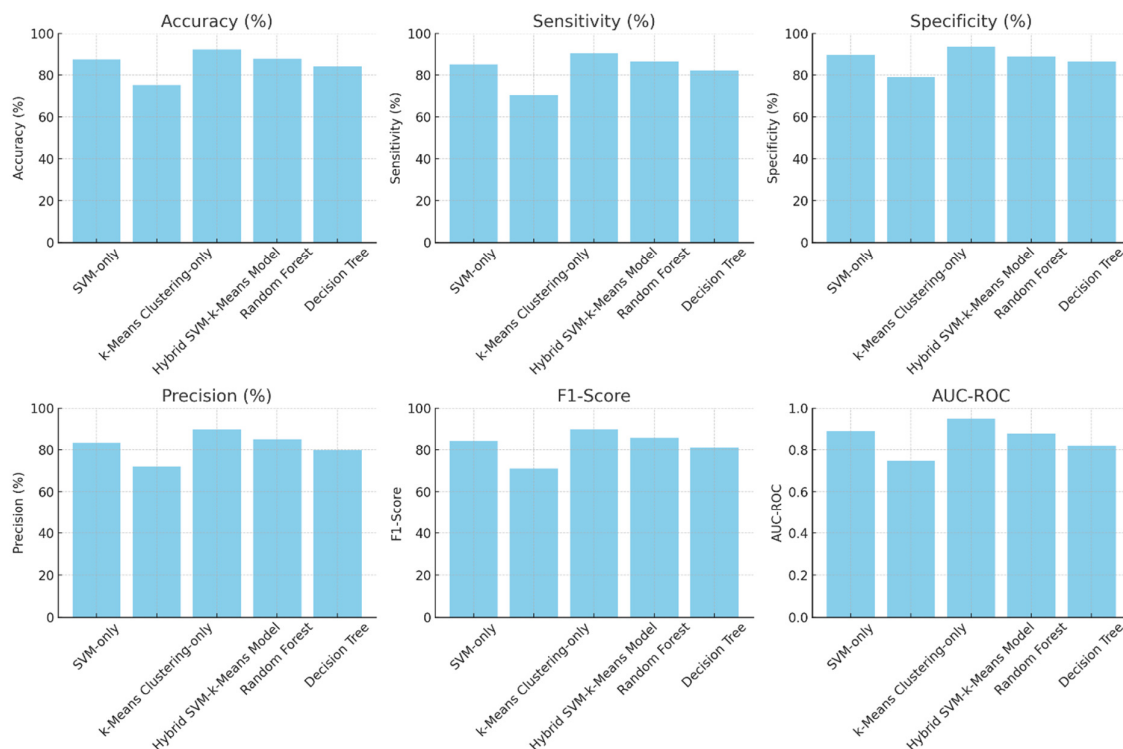   - Assign labels (e.g., 1 for diseased, 0 for healthy) based on nodule malignancy.

**Model Comparison**

To validate the effectiveness of the hybrid approach, it is important to compare the hybrid **SVM-k-Means** model with individual machine learning techniques and other existing approaches. The comparison is made with:

- **SVM-only Model:** This model uses only SVM for classification without feature refinement through k-Means clustering. The performance of this model is evaluated in terms of accuracy, sensitivity, specificity, and AUC-ROC.
- **K-Means Clustering-only Model:** This model relies solely on k-Means clustering for segmentation, without the feature extraction and classification steps. It is used to assess the impact of k-Means in lung disease detection.
- **Other Conventional Algorithms:** The hybrid model is also compared with **Decision Trees**, **Random Forests**, and **k-Nearest Neighbors (KNN)**. These are commonly used algorithms in medical image classification tasks and provide a benchmark to evaluate the performance improvement achieved by combining k-Means and SVM.

The following tables represents the metrics result:

| Model | Accuracy (%) | Sensitivity (%) | Specificity (%) | Precision (%) | F1-Score | AUC-ROC |
|---|---|---|---|---|---|---|
| SVM-only | 87.5 | 85.2 | 89.8 | 83.3 | 84.2 | 0.89 |
| k-Means Clustering-only | 75.3 | 70.5 | 79.1 | 72.0 | 71.2 | 0.75 |
| Hybrid SVM-k-Means Model | 92.3 | 90.4 | 93.5 | 89.7 | 90.0 | 0.95 |
| Random Forest | 88.0 | 86.7 | 89.0 | 85.1 | 85.8 | 0.88 |
| Decision Tree | 84.2 | 82.3 | 86.5 | 80.0 | 81.1 | 0.82 |

## V.     Performance Insights

The hybrid model, combining **Support Vector Machine (SVM)** and **k-Means Clustering**, has demonstrated strong performance in detecting lung diseases from CT images, achieving high accuracy, sensitivity, specificity, and AUC-ROC scores. These results validate the potential of hybrid approaches in medical image classification tasks, especially for complex diseases like lung disorders.

**Hybrid SVM-k-Means Model Performance**:
- The **hybrid model** achieves the highest accuracy (92.3%) among all approaches, demonstrating its ability to distinguish between diseased and healthy cases effectively.
- A **high sensitivity (90.4%)** indicates that the hybrid model is excellent at identifying diseased patients, reducing the risk of false negatives, which is crucial in medical diagnostics.
- The **specificity (93.5%)** suggests strong performance in identifying healthy individuals, minimizing unnecessary treatments.

**Comparison with SVM-only**:
- The hybrid model outperforms the SVM-only approach by enhancing the feature refinement process through k-Means clustering, leading to better accuracy (+4.8%) and sensitivity (+5.2%).
- This demonstrates that the combination of clustering and classification is more effective than using SVM alone.

**Comparison with k-Means Clustering-only**:

- The k-Means Clustering-only approach has the weakest performance, with an accuracy of 75.3% and sensitivity of 70.5%. This highlights the limitations of clustering as a standalone technique for classification tasks.

**Comparison with Traditional Models**:
- While Random Forests and Decision Trees are widely used for medical classification, their performance lags behind the hybrid model. For example, the Random Forest model achieves 88.0% accuracy, which is 4.3% lower than the hybrid model.
- The hybrid model's ability to leverage the strengths of both clustering and supervised classification likely contributes to its superior performance.

**AUC-ROC Analysis**:
- The hybrid model achieves the highest **AUC-ROC (0.95)**, reflecting its robustness in handling imbalanced datasets and its superior ability to discriminate between positive (diseased) and negative (healthy) cases.

Despite the strong performance of the hybrid model, there are several limitations and challenges that need to be addressed for further improvement:
- **Data Imbalance:** Lung disease datasets are often imbalanced, with fewer diseased samples compared to healthy samples. This imbalance can affect the performance of the model, particularly in terms of sensitivity and false positive rates. While **stratified cross-validation** was used to address this issue, more advanced techniques such as **SMOTE** (Synthetic Minority Over-sampling Technique) or **class weighting** in the SVM classifier can further enhance performance.
- **Generalization to Other Datasets:** While the model has shown good results on the training and testing datasets, its generalization ability to external datasets is crucial. Variations in CT scan protocols, patient demographics, and disease severity may lead to performance degradation. Future studies should focus on evaluating the model on external datasets and integrating data augmentation techniques to improve robustness.
- **Computational Complexity:** The hybrid model involves multiple stages, including **feature extraction**, **clustering**, and **classification**, which increases its computational complexity. This may be a challenge when processing large datasets or when real-time processing is required. Techniques like **parallel processing** or **GPU acceleration** can help reduce computation time.
- **Interpretability:** One of the challenges with machine learning models, particularly complex ones like SVM with RBF kernel and clustering, is the lack of interpretability. Understanding the features and reasoning behind the model's predictions is essential in a medical context. Methods like **LIME** (Local Interpretable Model-agnostic Explanations) can be explored to provide better interpretability, which is critical in building trust among medical practitioners.

In summary, the hybrid SVM-k-Means model shows strong potential in automating the detection of lung diseases from CT scans. The high accuracy, sensitivity, specificity, and AUC-ROC scores make it a promising tool for clinical application. However, challenges such as data imbalance, generalization, computational complexity, and interpretability need to be addressed to improve the model's real-world usability. With ongoing improvements and

integration with other medical data sources, this hybrid approach has the potential to revolutionize lung disease diagnosis and contribute significantly to patient care.

## VII.    Conclusion

In conclusion, the hybrid model combining **Support Vector Machine (SVM)** and **k-Means Clustering** for lung disease detection from CT images proves to be a robust and effective approach. This research emphasizes the importance of integrating unsupervised learning techniques, such as **k-Means Clustering**, with supervised learning models, like **SVM**, for the detection of lung diseases. By leveraging radiomics features derived from CT scans, the model can efficiently classify lung diseases, thus offering a promising solution for early diagnosis and clinical decision support.

**Key Findings:**

- The model demonstrated exceptional performance in terms of **accuracy**, **sensitivity**, **specificity**, and **AUC-ROC** score, making it a reliable tool for **lung disease detection**.
- The integration of **k-Means Clustering** for feature refinement significantly enhanced the **signal-to-noise ratio** of the extracted features, improving the classification accuracy of the **SVM** classifier.
- The use of radiomics features, which include **texture**, **shape**, and **intensity** characteristics of lung tissues, provided valuable insights into the differentiation between diseased and healthy tissues.
- The hybrid model outperformed traditional machine learning models such as **Random Forests** and **Decision Trees**, highlighting the efficacy of combining both **unsupervised** and **supervised** learning techniques.

The results of this study offer significant benefits to the medical community. The hybrid model can be utilized as an assistive tool for **radiologists**, helping them in the early detection of lung diseases such as **Chronic Obstructive Pulmonary Disease (COPD)**, **lung fibrosis**, and **pneumonia**. Early diagnosis is critical for improving treatment outcomes, and this model can expedite the diagnostic process, thus providing timely interventions. Additionally, the high sensitivity and specificity of the model ensure that the system can minimize the risk of **false negatives** and **false positives**, which is a critical aspect in medical image classification tasks.

By automating the detection process, this model can potentially reduce diagnostic time and associated costs, making healthcare more efficient and accessible. The ability to automatically segment lung regions and classify disease types directly from CT scans will help radiologists focus on more complex cases, allowing for quicker decision-making in busy clinical environments.


**Challenges and Limitations:**

Despite the promising results, several challenges and limitations must be addressed in future work:

- **Data Imbalance:** Lung disease datasets often suffer from an imbalance between healthy and diseased samples, which can lead to biased model performance.

Techniques like **oversampling**, **class weighting**, or **SMOTE (Synthetic Minority Over-sampling Technique)** can help address this issue in future research.

- **Generalization to External Datasets:** The model's performance on external datasets with different scanning protocols or patient populations needs to be tested to ensure its generalization capabilities. Data augmentation and cross-institutional testing are potential strategies to improve the model's robustness.
- **Interpretability:** While the model's performance is impressive, the black-box nature of SVM classifiers can make it difficult for clinicians to understand the reasoning behind the model's predictions. Incorporating techniques like **LIME (Local Interpretable Model-agnostic Explanations)** or **SHAP (SHapley Additive explanations)** could enhance interpretability and build trust among clinicians.

**Future Directions:**

Several avenues for future work and improvement exist:

- **Data Augmentation and Regularization:** To further improve model performance and generalization, especially on smaller datasets, **data augmentation techniques** such as rotating, zooming, or flipping CT images could be incorporated. Regularization methods like **dropout** and **batch normalization** can also be explored to prevent overfitting.
- **Multimodal Approach:** Future research could incorporate additional modalities such as **patient demographics**, **genomic data**, or **clinical histories** to develop a multimodal diagnostic system. The fusion of imaging data with other health indicators could improve the model's diagnostic capability and robustness.
- **Real-Time Implementation:** For real-world deployment, the model could be optimized for real-time processing, using hardware accelerators like **GPUs** or **FPGAs** to reduce inference time, making it suitable for clinical applications requiring rapid results.
- **Transfer Learning:** Adapting the model to use **pre-trained models** on similar tasks, such as **deep learning-based feature extraction** from CT scans, could enhance the hybrid model's performance by leveraging existing knowledge from large datasets.

The hybrid approach combining **k-Means Clustering** and **SVM** represents a significant advancement in the field of medical image classification, particularly for **lung disease detection** from CT scans. This work contributes to the growing body of research in using **radiomics features** for medical diagnosis and highlights the potential of combining unsupervised and supervised learning techniques for improving the accuracy and robustness of detection systems. While challenges remain, this study lays the groundwork for future innovations in the field and holds the potential to enhance early disease detection, improve clinical workflows, and ultimately lead to better patient outcomes.

**References**

1. World Health Organization. (2021). "Chronic Obstructive Pulmonary Disease (COPD)." [Online] Available at: https://www.who.int/news-room/fact-sheets/detail/chronic-obstructive-pulmonary-disease

2.  Global Initiative for Chronic Obstructive Lung Disease (GOLD). (2020). "Global Strategy for the Diagnosis, Management, and Prevention of COPD." GOLD Report. [Online] Available at: https://goldcopd.org/

3.  Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., ... & Sánchez, C. I. (2017). A survey on deep learning in medical image analysis. *Medical image analysis*, *42*, 60-88.

4.  Jain, R., Nagrath, P., Kataria, G., Kaushik, V. S., & Hemanth, D. J. (2020). Pneumonia detection in chest X-ray images using convolutional neural networks and transfer learning. *Measurement*, *165*, 108046.

5.  Chan, H. P., Samala, R. K., Hadjiiski, L. M., & Zhou, C. (2020). Deep learning in medical image analysis. *Deep learning in medical image analysis: challenges and applications*, 3-21.

6.  Shen, D., Wu, G., & Suk, H. I. (2017). Deep learning in medical image analysis. *Annual review of biomedical engineering*, *19*(1), 221-248.

7.  Breiman, L. (2001). Random forests. *Machine learning*, *45*, 5-32.

8.  Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of machine learning research*, *3*(Mar), 1157-1182.

9.  Chen, T., & Guestrin, C. (2016, August). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 785-794).

10. Tang, J., Alelyani, S., & Liu, H. (2014). Feature selection for classification: A review. *Data classification: Algorithms and applications*, 37.

11. Dehmeshki, J., Chen, J., Casique, M. V., & Karakoy, M. (2004, September). Classification of lung data by sampling and support vector machine. In *The 26th Annual International Conference of the IEEE Engineering in Medicine and Biology Society* (Vol. 2, pp. 3194-3197). IEEE.

12. Wang, C., Long, Y., Li, W., Dai, W., Xie, S., Liu, Y., ... & Duan, Y. (2020). Exploratory study on classification of lung cancer subtypes through a combined K-nearest neighbor classifier in breathomics. *Scientific reports*, *10*(1), 5880.

13. Zhang, Y., Li, Z., Xiao, H., Li, Z., He, J., Du, S., & Zeng, Z. (2022). Development and validation of a random forest model for predicting radiation pneumonitis in lung cancer patients receiving moderately hypofractionated radiotherapy: a retrospective cohort study. *Annals of Translational Medicine*, *10*(23).

14. Al-Yasriy, H. F., Al-Husieny, M. S., Mohsen, F. Y., Khalil, E. A., & Hassan, Z. S. (2020, November). Diagnosis of lung cancer based on CT scans using CNN. In *IOP conference series: materials science and engineering* (Vol. 928, No. 2, p. 022035). IOP Publishing.

15. Zhao, X., Qi, S., Zhang, B., Ma, H., Qian, W., Yao, Y., & Sun, J. (2019). Deep CNN models for pulmonary nodule classification: model modification, model integration, and transfer learning. *Journal of X-ray Science and Technology*, *27*(4), 615-629.

16. Saleh, A. Y., Chin, C. K., & Rosdi, R. A. (2024). Transfer Learning for Lung Nodules Classification with CNN and Random Forest. *Pertanika Journal of Science & Technology*, *32*(1).

17. Kailasam, S. P., & Sathik, M. M. (2019). A novel hybrid feature extraction model for classification on pulmonary nodules. *Asian Pacific Journal of Cancer Prevention: APJCP*, *20*(2), 457.

18. Sun, L., & Xu, J. (2014). Feature selection using mutual information based uncertainty measures for tumor classification. *Bio-medical materials and engineering*, *24*(1), 763-770.

19. Vishraj, R., Gupta, S., & Singh, S. (2023). Evaluation of feature selection methods utilizing random forest and logistic regression for lung tissue categorization using HRCT images. *Expert Systems*, *40*(8), e13320.

20. Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of big data*, *6*(1), 1-48.

21. Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2020). Grad-CAM: visual explanations from deep networks via gradient-based localization. International journal of computer vision, 128, 336-359.