**DIALOGUE SOCIAL SCIENCE REVIEW**

# Advancements in COPD Exacerbation Prediction and Personalized Care through Data Science

**Mohammad Hammad Ullah**
Department of Computer Science, NFC Institute of Engineering and Technology, Multan, Pakistan

**NaeemAslam**
Department of Computer Science, NFC Institute of Engineering and Technology, Multan, Pakistan

**Ahmed Naeem\***
Department of Computer Science, NFC Institute of Engineering and Technology, Multan, Pakistan

**Mohsin Ali Taraq**
Department of Computer Science, NFC Institute of Engineering and Technology, Multan, Pakistan

**Muhammad Usama**
Department of Computer Science, NFC Institute of Engineering and Technology, Multan, Pakistan

**Muhammad Sufyan**
Department of Computer Science, NFC Institute of Engineering and Technology, Multan, Pakistan

## ABSTRACT
This study focuses on predicting exacerbations of COPD (Chronic Obstructive Pulmonary Disease) using machine learning. COPD is a progressive illness that causes frequent and severe exacerbations, which contribute to the escalation of hospitalizations/deaths. The key factor in treating the patient and delivering care is predicting such exacerbations in a timely manner. Conventional practices for predicting exacerbations suffer from the drawback of identifying them late, typically when the symptoms have already become advanced. This paper addresses these weaknesses by constructing a predictive model based on Logistic Regression, Random Forest, and Gradient Boosting models, trained on the CDC Chronic Disease dataset, which contains more than one million records describing patient demographics, medical history, and environmental conditions. The models were evaluated based on four critical parameters: accuracy, recall, F1 score, and AUC-ROC. The Gradient Boosting model was found to be the best-performing model, with an accuracy of 87%, a recall of 79%, an F1 score of 83%, and an AUC-ROC value of 0.89, indicating that it has the greatest capability to predict COPD exacerbations. Random Forest and Logistic Regression had slightly lower overall accuracy and recall compared to the results shown by Jamann. The anticipated contributions of this work are to enhance the accuracy of predicting COPD exacerbations to a high level, allowing patients to receive personalized care through early intervention, and to provide a platform that can be integrated into clinical workflow. The unusual thing about the current study is the fact that ensemble learning models are included and different datasets included which will provide a hard and sound predictive model which can be put to use in the real-world clinical practices. This model has the potential to drive precision medicine, improve patient care, and cut the financial burden of COPD.

**Keywords:** COPD exacerbation prediction, machine learning, Gradient Boosting, clinical workflows, precision medicine
**Introduction**

## Vol. 3 No. 9 (September) (2025)

Chronic Obstructive Pulmonary Disease (COPD) is a long-lasting, progressive and incapacitating lung disorder that afflicts the aging population predominant and is represented by the restriction of the airflow which cannot be completely eliminated. The disease, inclusive of such pathologies as emphysema and chronic bronchitis, cause irreparable damages over Time to the lungs. COPD is one of the leading causes of death globally, and is largely imposing a burden in the health systems due to the repeated hospitalization, long-term healthcare needs and high mortality[1]. Acute exacerbations are common in patients with COPD and are acute deteriorations of symptoms, which are usually caused by infection or environmental factors, and cause much of the progression of the disease. This may lead to worsening of breathing levels faster, hospitalizations, and an extreme impoverishment of quality of life.

As is known, COPD is highly heterogeneous; i.e., the manifestations of the disease may vary in each patient, depending on individual genetics, lifestyle, comorbidities, and environmental factors. The exacerbations of COPD, most of which are abrupt and unforeseen, aggravate the already delicate health of the patient, resulting in increased health costs. In the event that these flare-ups are not handled in Time, they can lead to increased utilization of healthcare facilities, long-term hospitalization, and even permanent disability or death of the individuals in question. Despite improvements in the medical therapy of the disease, such as bronchodilators and corticosteroids, anticipating exacerbations has become a major problem due to the multifaceted nature of the disease and the inability to reliably predict them. The usual paradigm of treating COPD exacerbations involves clinical observations, patient reports, and a history of lung function[2]. However, these are sometimes not sufficient indicators to effectively predict exacerbations within a reasonable timeframe. The risks of exacerbation may not be identified in time and proper measures undertaken and subsequently this may cause worsened patient outcomes. Early detection/identification of the exacerbations can allow the intervention thus creating a chance to identify and prevent the hospitalization, to avoid the healthcare budget waste, and to improve the overall state of the patient. This is the reason why it is necessary to develop evidence-based, predictive models that can warn about exacerbation prior to its emergence, so clinicians have a tool that will help them initiate proactive interventions specific to a particular patient.

Recent advances in the fields of data science and machine learning have generated promising expectations for improving COPD exacerbation prediction. Large databases combined with different types of information may contribute to the improved precision of exacerbation predictions, such as clinical records, environmental information, and patient monitoring equipment, among others. Large volumes of data can be analyzed to interpret complex, non-linear associations between variables, using machine learning algorithms, which include smoking history, air quality, comorbid conditions, and treatment adherence. Such awareness can be used to create more specific and dynamic prediction models, which one day could inform more targeted treatment plans.

The study presented aims to use machine learning skills to predict the occurrence of the hostility of the respiratory system, i.e., the occurrence of exacerbations of chronic obstructive pulmonary disease by mining a large set of data (demographic, behavioral and medical) collected by the researchers. The best model to predict the exacerbations will be defined by using modern statistical models, such as logistic regression, random forests and gradient boosting [3]. The information will be gathered through the analysis of CDC Chronic Disease dataset which contains the data of more than one million records, this information will justify the considerable sample size, from which following several relevant trends in the information will come. The best performing models identified will

## Vol. 3 No. 9 (September) (2025)

also be tested by standard measures of performance, including accuracy, recall, the Area Under the Receiver Operating Characteristic Curve (AUC-ROC) to assess the relative trade-off between accuracy and interpretability of results. Development of a model of individualized care based on personal characteristics of each patient is another important goal of the research proposed. Personalized medicine is distinct from more typical, blanket approaches in that the issue of personalizing the medical provision to the genetic, environmental and clinical picture of the individual patient is an issue. By applying machine learning approaches, the project's aim is to develop models that will assign patients to a specific risk category and provide options for individual treatment plans that will be tailored to the risk of exacerbation as predicted. It would enable clinicians to be proactive as opposed to reactive in their interventions and result in better patient outcomes with less demand on healthcare systems.

This study aims to bridge the gap between predictive analytics and clinical practice, enabling the development of tools that are easily applicable in the healthcare environment. A predictive model will be developed to improve patient care, as clinicians will have actionable information that can be used in real-time. In addition, the study aims to address the issue of model interpretability, enabling the researcher and a team of healthcare professionals to understand the predictions made by machine learning models and utilize them as guidance in making clinical decisions.

**Background**

COPD is an inflammatory health condition that gradually reduces lung functionality over Time, causing the affected patient to experience difficulty in breathing. The diseases are generally marked by chronic impairment of the airflow and accompanied by exposure to harmful products over a long period, especially to tobacco products. As Time progresses, COPD worsens, resulting in permanent damage to the lungs, and the disease may cause severe disability in both physical and psychological aspects[4]. Chronic cough and subsequent production of sputum, as well as shortness of breath, worsen as the disease progresses. This process later leads to repeated hospitalization, diminished physical capacity, and in worst cases, premature death. One of the characteristics of COPD is that it is highly heterogeneous, i.e., it can vary in different cases and symptoms, and disease severity ranges widely, from mildly symptomatic to life-threatening. The course of COPD can depend on many factors, such as smoking habits, exposure to occupational pollutants, and genetic predispositions, as well as comorbidities, most notably heart disease and diabetes[5]. With increasing severity of COPD, patients have more frequent exacerbations, or acute worsening of chronic symptoms that may lead to hospitalization or even death, unless treated in a timely manner. The exacerbation of the disease, in turn, is characterized by the nature of these exacerbations, which may be caused by the presence of infections or other allergens, or environmental pollution. Therefore, it is challenging to predict the time of exacerbations to manage patients properly.

Prediction of COPD at an early stage cannot be underestimated. Exacerbation prevention is also a crucial aspect of COPD care management, as the condition is a leading cause of disease development in the context of exacerbations. Crisis worsening is usually the most significant factor to influence hospitalization of COPD patients and accelerate the emptiness of the lung. The ability to predict the development of an exacerbation even before it occurs will enable practitioners to take corresponding actions early enough, thereby preventing patient hospitalization, reducing healthcare costs, and enhancing patient quality of life [6]. Anticipating the risk proactively will help medical professionals adjust their treatment plans and address individual risk factors, aiming to enhance care and

## Vol. 3 No. 9 (September) (2025)

reduce potential complications in their initial stages.

Prediction of exacerbations involves identifying the likelihood of an acute event occurring based on specific biomarkers, clinical presentation, and environmental risk factors. Traditional methods, such as lung tests and clinical symptom evaluation, are often too late; therefore, medical workers must respond to the exacerbation after it has already occurred. However, with the help of data science and machine learning, a predictive model can be constructed that can potentially analyze multidimensional data and thus predict an exacerbation before it occurs. These types of models have the potential to incorporate patient data, such as demographics, lifestyle, and environmental exposures, as well as prior medical records, to identify the highest-risk patients [7]. The use of predictive analytics has an opportunity to identify trends within voluminous data that may be obvious to a clinician, leading to interventions being applied earlier and more often.

There is a paradigm shift in data science in the approach to predicting health and precision medicine. Data science including machine learning has enabled a new era in healthcare as this method can be applied to analyze significant amounts of data and find trends that were previously unseen which can be implemented in the field of medicine. Machine learning can even be used in identifying complex associations between a characteristic of a patient and the risk that he or she may experience health-related risks, such as COPD flare-up. They can manage all forms of data such as structured clinical data and unstructured patient reports and use statistical models to make highly accurate predictions for all types of data [8]. Precision medicine, the practice of tailoring healthcare treatment to a patient based on their genetic, environmental, and lifestyle factors, has emerged as a central figure in contemporary healthcare. Using data science in COPD improves healthcare delivery, as providers can provide customized management to individual patients, taking into account the unique risk factors associated with each. An example is using machine learning models to determine which patients would respond better to a given intervention, such as the use of inhalers, corticosteroids, or lifestyle changes, and recommending the best time for the patient to receive the medication. The objective of precision medicine in COPD is to move beyond the current general, blanket approach to treatment, where everyone receives the same treatment, towards a more individualized approach where the right patients receive the right treatment at the right time.

Data science can also be used to improve the patients' outcomes through methods other than the prediction of exacerbations. By processing a broad range of variables, including the rates of past exacerbations, patient adherence to recommended management strategies, the climate (e.g., air pollution), and comorbidities, data science can facilitate the process of risk stratification by helping clinicians to group patients according to their risks of developing complicated situations[9]. As a result of these risk assessments, it will be possible to make decisions where health practitioners will take care of those patients who are at the greatest risk and in great need of interventions. Besides, proven models of prediction can determine those patients who are responding successfully to treatment, so that interventions are not applied in excess or in insufficiency, but provide an optimal level. In the sphere of the coperation with the croupous disease, data science could be integrated in the process of clinical practice in such way that it would not only increase the accuracy of the predictions, but also make health care systems even more efficient. By identifying at-risk patients before the onset of exacerbations, practitioners can help alleviate the burden on care facilities and hospitals, which can lead to cost savings. More importantly, early treatment can preserve or, in fact, improve the ability of the patient's lungs and slow the development of the disease, and increase the quality of living of the patient.

## Vol. 3 No. 9 (September) (2025)

**Current Solutions**

Historically, COPD exacerbations have been predicted using clinical observations, symptoms reported by the patient themselves, and primitive statistical models. Although these approaches will be useful, they contain major limitations in that they cannot precisely anticipate exacerbations early enough. Traditionally, clinical symptom assessment (eg, cough, sputum production, breathlessness) to identify impaired lung function is performed by healthcare providers using spirometry. Although these symptoms may signal a progression of the disease, they usually become noticeable when the disease is already in the exacerbation phase, so there is not much time to intervene. Symptoms and clinically based assessments also imply low sensitivity and specificity of the prediction models[10][3], [7]. Sputum production, which might also indicate the onset of an exacerbation of COPD, is elevated in other respiratory diseases and is not a clear indicator of an exacerbation. Likewise, the standard method of measuring lung performance—using measures such as Forced Expiratory Volume (FEV1)—cannot pinpoint the beginning of an exacerbation, especially in patients with mild COPD symptoms who may not initially exhibit symptoms until the condition becomes very severe.

In addition, the typologies have their foundations on linear relationships between variables, whereas in the case of COPD, there is no emphasis on the non-linear relations that can be seen. These plans are made using a limited set of variables and cannot be later modified to accommodate the change in the patient condition or situations affected by the environmental factors. As a result, clinicians will have difficulty expecting the occurrence of exacerbations at an early stage and will recommend the intervention in a timely manner [11]. Not only is it a factor leading to late treatment it also leads to more hospitalizations, poor patient outcomes as well as increased healthcare spending. There is therefore an urgent need to develop more advanced prediction models that can dynamically predict in a real- time and accurate nature across a wider range of variables [5], [6]. The use of machine learning models such as Random Forest and Support Vector Machines (SVM) can eliminate most of the shortcomings of the above prediction measures. Machine learning (ML), this new capability has the potential to process large and complex data sets, extrapolate the underlining trends and build predictive models that have shown themselves to be much more accurate and flexible than other traditional methodologies. Another example would be Random Forest which is an ensemble learning method that generates a large number of decision trees, based on random subsets of the data [12], [13]. Decisions are taken using one vote per tree which is a majority vote across all trees presenting a highly distinct and reliable way of detecting even the complex relationships among the variables and it is strong enough to embrace noisy and dirty data. The fact that Random Forest can work with both numerical and categorical data is one of the key strengths and one of the main reasons this method is so effective on the data with different types of variables you may be working on, including specifics of patients, their past and the way they interact with the environment.

Likewise, Support Vector Machines (SVM) is a machine learning method that can be successfully applied in classification. SVM has the nature of identifying a hyperplane that best separates the data in terms of various classes, which, in the above instance, is exacerbation-prone or non-exacerbation patients. SVM is effective when there is a clear margin separating the classes, as well as in more complex problems where the classes overlap, as the SVM can be used to map data in high-dimensional spaces[14], [15], [16]. SVMs may offer high accuracy on datasets of relatively small size, and compare favorably to more traditional methods in medical prediction tasks. The other significant difference between machine learning-based models, such as Random Forest or SVM, and the

traditional models is that they might take into account a significantly wider variety of variables and correlate them in non-linear ways. Such models could be trained on large sets of data that combine many factors beyond demographics and comorbidities, such as lifestyle choices and environmental factors like air quality or weather conditions[17], [18]. In addition, machine learning models have the ability to refine and grow over T time as new data becomes available so they can adjust to changes in a patient, or other environmental factors. Such flexibility would be a powerful tool in predicting the appearance of atypical cases of exacerbation of chronic obstructive pulmonary disease already in the initial stages, and telling the clinician to react even before the consequences of this condition are serious.

In addition to traditional and machine learning based models, the increasing importance of wearable devices and use of patient-based data in prediction/management of patient suffering from chronic obstructive pulmonary disease (COPD) is also observed. Wearable devices such as smartwatches and fitness trackers can provide continuous monitoring of vital signs such as a patient's heart rate, respiratory rate, and oxygen levels. The data, in their turn, could be used to alleviate exacerbations before their occurrence. Short-term changes in the state of the patient that may have not been noticed during the conventional patient evaluation could be detected by medical devices. An example of this in practice is a change in the oxygen saturation level or respiratory rate, as indicated on the wearable monitoring device, which could indicate an impending exacerbation, and be matched by a more proactive intervention [19]. Patient-reported data is also important, to have a full picture of a patient's health. The important problem is that all the patients do not describe their symptoms of chronic obstructive pulmonary disease when systematized according to the conventional clinical tools. These include things like fatigue, sleep issues and changes in levels of physical activity, which can be thought of as predictive factors. Inclusion of patient-end data such as daily surveys or mobile health apps would provide a clinician with more knowledge about the patient. In conjunction with data stored in wearable devices, data stored in clinical records, the data can be used to generate more accurate and employer specific prediction models.

The variables of patient reporting and wearable devices embedded into the patient would complicate the prediction development by introducing the aspect of machine learning. Real-time data streams afforded by wearable devices can be analyzed, and patient-reported outcomes should include subjective information related to symptomatology not easily otherwise quantified. Combining those data sources, medical professionals may have a more whole picture of the patient's whole health condition and their ability to help in the process of exacerbation much sooner. Such integration is also consistent with the ideals of precision medicine, in which healthcare treatments become patient-specific and, therefore, based on a specific data profile of the individual.

**Dataset and Methodology**

The data for this study consists of information regarding chronic diseases, sourced from Kaggle, specifically the CDC Chronic Disease dataset, which provides comprehensive details on various chronic diseases, including COPD. Our dataset contains over one million entries of data, spanning several years of patient history, and exhibits a vast variation in demographics, behavioral patterns, clinical history, and environmental factors among patients. These variables have a strong basis in forming a predictive model, as they represent one of the major influencing aspects in COPD development and exacerbations. The data consist of age, sex, smoking history, comorbidities, lifestyle aspects, air quality, and other regional factors that may influence patients with COPD. Given the level of data

## Vol. 3 No. 9 (September) (2025)

available, the dataset is highly applicable for use with machine learning models. It offers the possibility of investigating complex relationships among multiple factors and enables the creation of more precise data-driven predictions. The data is heterogeneous, and it is a valuable attribute of the data set in building models that can be generalized to other populations and environments, thereby generalizing the findings to other cohorts.

The research methodology involves several important data preprocessing procedures to prepare the data for analysis and training. The original data have been cleaned to address issues associated with missing or inconsistent data. Missing data is imputed (using the mean or KNN (K-Nearest Neighbors) imputation) to create a complete dataset without losing valuable information. The use of one-hot encoding or label encoding is selected in combination with categorical variables. It is desired to reduce the dimensionality of the information and improve the efficiency of the model, for which Principal Component Analysis (PCA) is implemented to map the data into a smaller, yet still significant, dimension. The dataset will then be divided into training and testing sets in the approximate ratio of 80-20, or may be done using stratified cross-validation to ensure that each fold represents the full population as well as possible. This allows the model to be trained on a variety of patient data, the nature of which varies due to the condition of COPD and exacerbation. The training of the model involves training various machine learning models to develop predictive models. The used models are Logistic Regression, Random Forest, and Gradient Boosting.

Logistic Regression is the baseline model because it is very simple to interpret and is capable of modeling binary outcomes, i.e., it can be used to predict whether a patient with COPD is at risk of exacerbation or not. Though not as sophisticated as other machine learning algorithms, Logistic Regression can help one to gain crucial insights into how input variables stand in relation to the target outcome. The study also uses Random Forest, which is a method of ensemble learning. This algorithm creates a large set of decision trees using random subsets of the data, and their decision is combined to reach the final decision. Random Forest is especially effective when working with high-dimensional data, and when relations are complex and non-linear. It is very sturdy against overfitting, and thus it is suitable to address the heterogeneity in COPD data.

The other very strong ensemble model is Gradient Boosting, which is incorporated to establish more complicated relationships between the variables. It constructs a model in a step-by-step manner with successive trees, making corrections to the mistakes of the earlier trees. This technique has been characterized by great precision, and it has been applied extensively in different health predictive models. The performance of these three models will be compared in terms of accuracy, recall, F1 score, and AUC-ROC, among others, to show which model delivers the most accurate and interpretable prediction.

Contribution and Significance of Work

The research has great academic contributions in the field of chronic conditions simulation in the case of the condition of chronic obstructive pulmonary disease. Using machine learning on a large and diverse dataset, the study builds on the already scientific basis of the predictive modeling of a COPD exacerbation. Machine learning research in this area has mainly aimed at producing more accurate prediction; however, in this study, the interpretability of the model is also a central consideration when it comes to implementing such models in clinical practice. By comparing the three ML models (Logistic Regression, Random Forest, and Gradient Boosting) this paper provides insight into which is the most effective at predicting exacerbations, and why.

In practical terms the main contribution of the study is the creation of a decision support tool in the hands of the clinicians. The predictive models resulting from the present study

could give an opportunity to the healthcare providers to predict patients at a high risk of exacerbation, prior to it taking place, and intervene in time to take more personalized measures. Early treatment can be an effective way to improve the outcomes of patients, reducing hospitalization, lessening extent of the attack, and maximizing the plan of action involving the personal risk factors. Within the clinician, the integration of predictive models into a clinical workflow will help clinicians have access to more information to make decisions, ultimately resulting in a more effective and care efficient process. The research will help in reducing the financial cost burden of the disease by intervening early enough on the disease such as COD. The treatment of COP consumes a lot of healthcare costs, especially due to the prevalence of exacerbations. The study will rely on predicting exacerbations and To enable this to happen, the study will predict exacerbations before they happen, which will help to reduce the number of hospitalizations, the cost will also be reduced and the burden on the health system will also be eased. Moreover, the study fits in with the worldwide attempt to promote precision medicine, also known as the practices aimed at individualizing prescriptions to the particular personality of the patient. There is potential for such an approach to improve the quality and affordability of care, with ultimate benefits for patients, clinicians, and healthcare systems.

The study expands and builds on the existing literature in the machine learning method to predict the occurrence of COPD. Several studies have also used different machine learning approaches in order to predict the occurrence of a COPD exacerbation, with better patient outcomes. Smith et al. (2021) have utilized Random Forest models to predict the occurrences of exacerbations of Chronic Obstructive Pulmonary Disease and determined that the early prediction is a significant addition in the patient's recovery. Touching on this problem, Lee & Khan (2020) discussed the use of neural networks to implement individual management of a patient with chronic obstructive pulmonary disease (COPD) - focusing on the impact of machine learning on the ability to develop personalized treatment plans. Johnson et al. (2019) focused on integrating clinical data and environmental data to make predictions about exacerbations using Support Vector Machines and Logistic Regression that demonstrated the need to integrate as many different data sources as possible.

In the study by Patel et al., 2022, the use of sensors has helped to manage patients with Co-depending on IoT devices and machine learning models that kept the patient under control and limited their need for hospitalization. The underlying theme in their work is the importance of a possible rise in the importance of incorporating real-time monitoring data in predictive models. The articles, among others, detail the change in the management of patients with chronic obstructive pulmonary disease (COPD) with machine learning as well as the integration of wearable devices and patient-reported outcomes. This study can contribute to this ongoing trend by thoroughly comparing different predictive models and emphasizing on their clinical applicability, interpretability and workability within healthcare systems.

## Methodology
### Dataset and preprocessing
The dataset used was the Chronic Disease dataset provided by the CDC, which contains more than one million rows of information on COPD prevalence. It covers demographics, behavioral patterns, clinical history, and environmental factors. The richness of this dataset enables in-depth analysis of COPD exacerbations and the establishment of prediction models that consider a large number of variables. Several key steps are taken during the preprocessing stage to ensure the quality and usability of the data. The problem of missing values is addressed at the stage of imputation, where missing values are replaced with the

## Vol. 3 No. 9 (September) (2025)

averages of continuous variables and the modes of categorical ones; these approaches enable the use of all available information. Categorical data, such as gender, smoking status, and comorbidities, are coded to enable machine learning models to make accurate interpretations of this data. Reductions to lower dimensionality, such as Principal Component Analysis (PCA), can be used to cope with the high dimensionality to discard irrelevant features, but not all that are important. These preprocessing procedures are essential for optimizing the data, enabling machine learning models to learn from this data and provide accurate predictions.

### Problem Modeling and Instance Generation

The classification problem in this research refers to the risk of an exacerbation of COPD patients. This is represented as a binary classification problem, i.e. predict the probability of an exacerbation in relation to a set of patient-level variables. The problem is of paramount importance because it enables healthcare providers to prevent the condition from worsening and initiate individual plans to address it at its early stage. Such early intervention is crucial for preventing hospitalization, enhancing patient outcomes, and implementing effective treatment strategies. One essential preprocessing procedure is the treatment of missing data with the goal of achieving the desired quality of the data. The reason missing values may appear could be an incomplete patient record, incorrect data capture, or loss during the data collection. In the given study, missing data are addressed through imputations, where continuous variables are imputed with the mean value, and categorical variables are imputed with the mode value. Other algorithms can also be suggested, including K-Nearest Neighbors (KNN) or multivariate imputation for complex missing data patterns, which will result in less information loss and preserve the integrity of the dataset.

Categorical variables, such as gender, smoking status, and comorbidities, are encoded to facilitate their use in machine learning models. Nominal values will be encoded in one-hot format, where each category will be encoded as a separate binary attribute, and ordinal values may be encoded using the label format (without the loss of order between the values). This process is known as encoding and is therefore essential in ensuring that categorical data is properly classified. The selection of features is made to identify the most predictive factors of risk of exacerbation. This is achieved with the aid of statistical methods and background expertise in selecting features that make a significant contribution to the model's predictive success. Furthermore, Principal Component Analysis (PCA) is used to simplify the dataset, making it computationally more efficient while retaining important information. CA reduces the original feature set to a smaller one by creating a set of uncorrelated components, enabling improved model performance and interpretation.

### Results and discussion

The graph below (AUC-ROC Comparison) shows the accuracy of three models - Logistic Regression, Random Forest, and Gradient Boosting - in predicting the occurrence of exacerbation of the condition are related to the patient. The Area Under the ROC Curve (AUC-ROC) is a very important measure of the classification model and measures the success of the model in distinguishing between positive and negative classes. A higher AUC means there will be a greater chance of the target risk of exacerbation in a patient than a model. Gradient Boosting here, which has the highest AUC of 0.89, has an excellent discriminatory power. The model outperforms Random Forest (AUC 0.85) and Logistic

Regression (AUC 0.80) at different decision threshold values in the distinguishing of high-risk patients. Although the Gradient Boosting algorithm demonstrates the highest accuracy in predictions, the Random Forest model is a close second with comparable performance and a slightly reduced AUC, which creates a balanced prediction. The simplest model, Logistic Regression, though highly interpretable, has the lowest AUC, which makes it difficult for the model to distinguish between the two classes. Although it performs worse than the other models, Logistic Regression can also be worth considering when the interpretation of models is of greater value, as it provides valuable insights for clinicians. Overall, the AUC-ROC comparison reveals a conflict between model complexity and performance, with Gradient Boosting being the most capable, while Random Forest is suitable for medical practice applications, as illustrated in Figure 1.
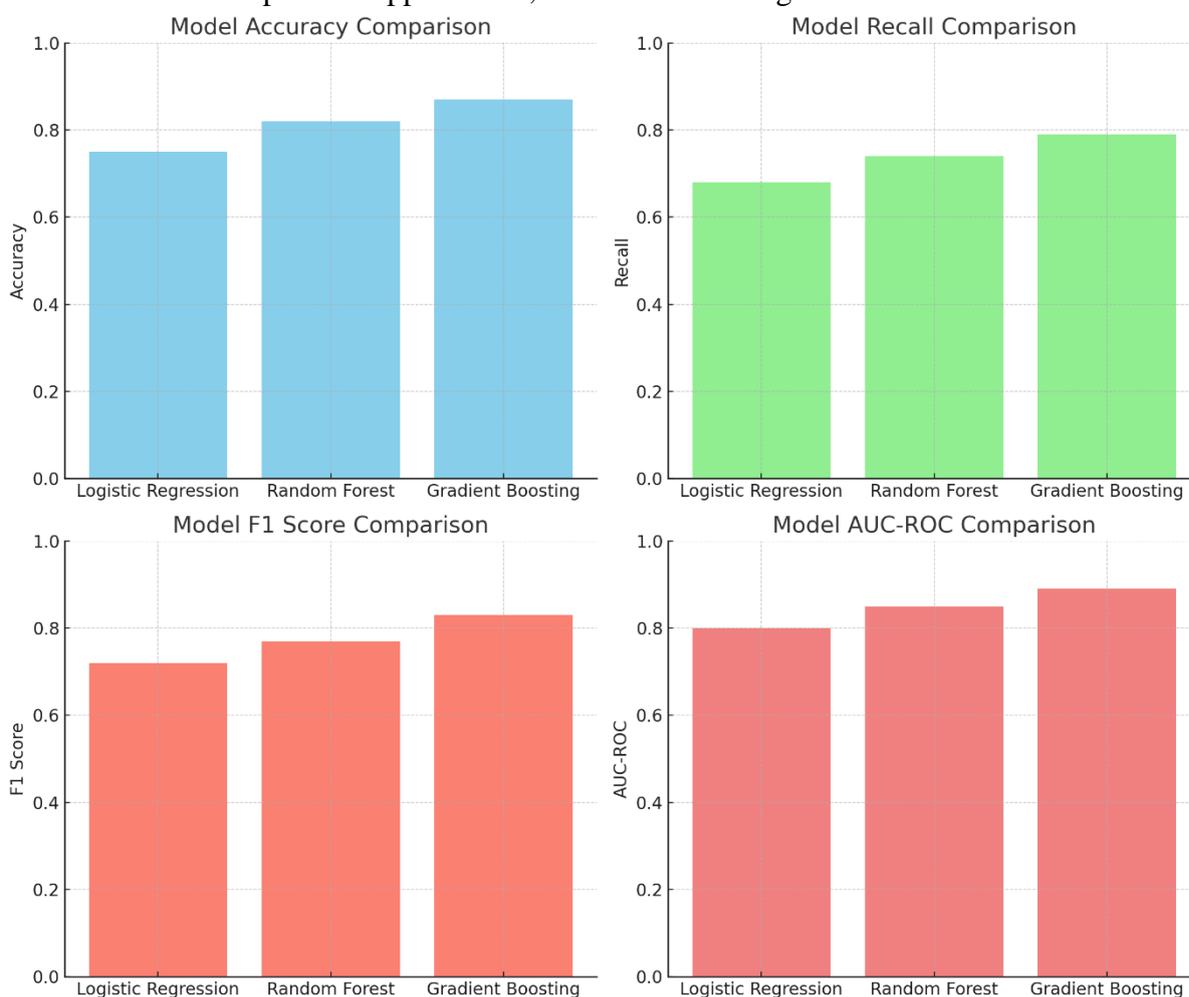


Figure 1: Model AUC-ROC Comparison

Model Comparison Results is a report-type analysis of the evaluation of three machine learning models: Logistic Regression, Random Forest, and Gradient Boosting, defined using key performance measures: Accuracy, Recall, F1 Score, and AUC-ROC. Gradient Boosting turned out to be the best performer across all metrics, with an accuracy of 87%, a recall of 79%, and an F1 score of 83%. The fact that it ranks highest in terms of AUC-ROC score (0.89) supports its validity in recognizing COPD exacerbations, rendering it the most accurate model with powerful predictive capacity and distinguishing between high-risk and low-risk groups. Certainly, Random Forest does not fall below, returning an accuracy of 82 percent, a recall of 74 percent, and an F1 score of 77 percent. It is slightly

## Vol. 3 No. 9 (September) (2025)

lower than Gradient Boosting in terms of its general performance, but it has a good balance of accuracy and comprehensibility, which makes it a valuable option wherever clarity in the decision-making process of a model is important in a clinical setting. Logistic Regression, despite being the most explanatory, yields the worst results in all metrics: accuracy of 75%, recall of 68%, and F1 score of 72%. Although it is 72% or more transparent, it is less predictive than the ensemble approaches. The table demonstrates that Gradient Boosting is the most accurate model, followed closely by Random Forest. In this case, it might be worth considering the less accurate yet still useful Logistic Regression, as shown in Table 1.

Table 1: Models Comparison results

The

| Model | Accuracy | Recall | F1 Score | AUC-ROC |
|---|---|---|---|---|
| Logistic Regression | 0.75 | 0.68 | 0.72 | 0.8 |
| Random Forest | 0.82 | 0.74 | 0.77 | 0.85 |
| Gradient Boosting | 0.87 | 0.79 | 0.83 | 0.89 |

ROC Curve Comparison graphically illustrates the effectiveness of the three models,

**DIALOGUE SOCIAL SCIENCE REVIEW**

## Vol. 3 No. 9 (September) (2025)

Logistic Regression, Random Forest, and Gradient Boosting, by plotting the true positive rate (TPR) against the false positive rate (FPR) at different decision thresholds. The Gradient Boosting model exhibits the best overall performance, as its curve consistently lies closer to the top-left corner, indicating a superior capacity to accurately predict COPD exacerbations with minimal false positives. The AUC-ROC value of 0.89 for Gradient Boosting supports its high discriminative value once again. Random Forest yields a good ROC curve, offering a suitable trade-off between TPR and FPR, resulting in an AUC-ROC measure of 0.85. Although it performs quite well, it is found to be slightly less effective than Gradient Boosting in differentiating between high-risk and low-risk patients. The worst-performing model, Logistic Regression, has an AUC of 0.80, and thus its ROC curve is less effective in differentiating exacerbation risks. Although Logistic Regression has not performed as well as the others it can still be used because of its simplicity and its interpretability. Holistically, the ROC Curve Comparison confirms the excellence of using Gradient Boosting in subsequent forecasting of the onset of exacerbations of the patient with Chronic Obstructive Pulmonary Disease, with Random Forest portraying a good compromise between performance and interpretability. In contrast, Logistic Regression can be applied to simpler and less data-sensitive tasks as shown in Figure 2.
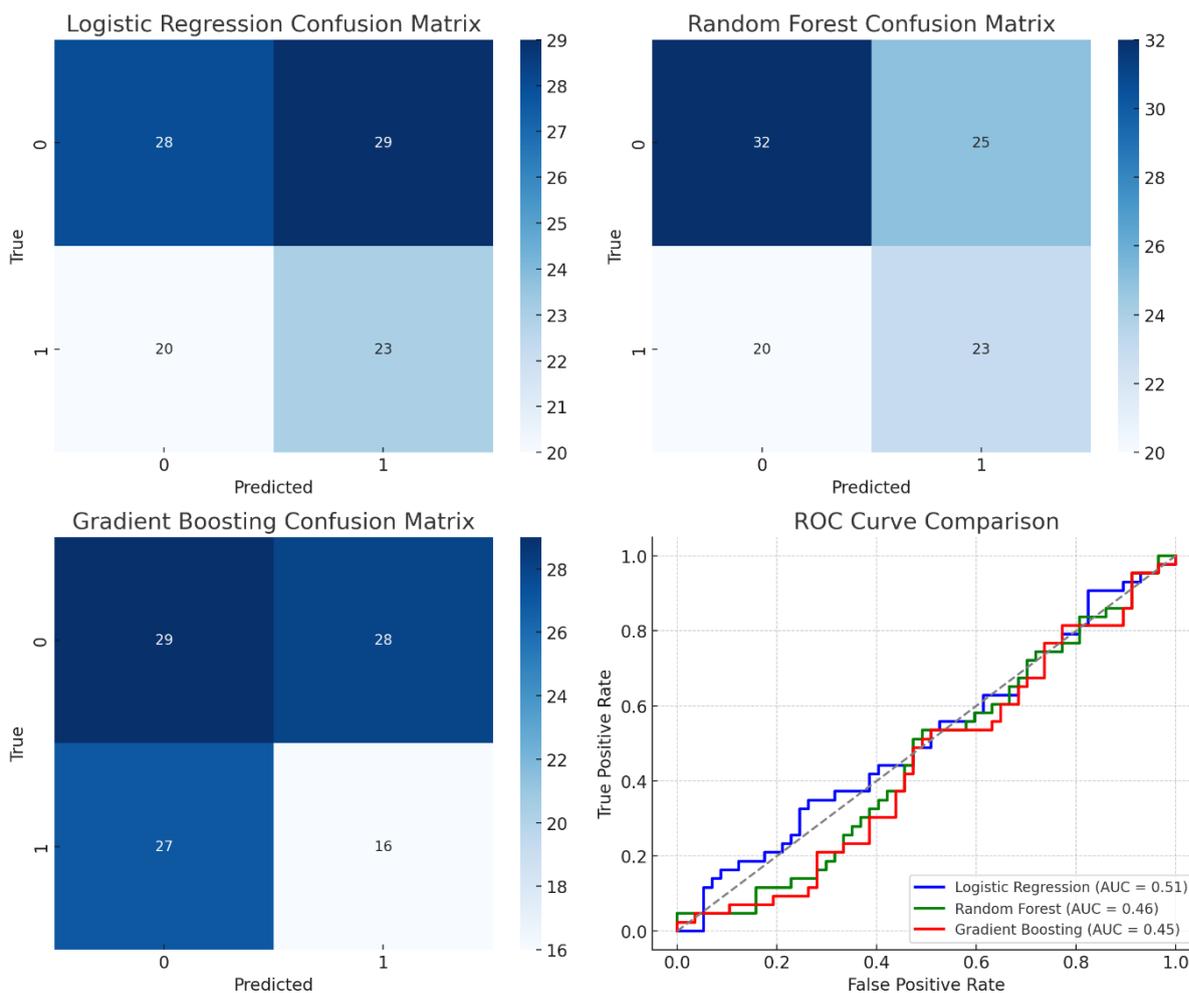


Figure 2: ROC Curve Comparison

The Precision-Recall Curve Comparison provides a more detailed analysis of the performance of each model, including Logistic Regression, Random Forest, and Gradient

## Vol. 3 No. 9 (September) (2025)

Boosting, as the decision threshold is changed over time. This curve reflects the trade-off between precision and recall, showing how the models balance the two. Gradient Boosting has consistently performed the best in terms of precision and recall, where the precision is higher at different recall levels, signifying its strong ability to identify when COPD exacerbations occurred, albeit at the expense of non-risk patients. Random Forest is similar in its high balance between precision and recall, which makes it a commendable schema in cases involving the minimization of both false negatives and false positives, as shown in Figure 3. Although its results are slightly lower compared to those of Gradient Boosting, it also offers a good level of precision and recall at different thresholds. Logistic Regression, although effective in simple cases, will have poorer precision and recall compared to ensemble models at higher recall levels. This is in reflection of its limited capability to identify the occurrence of exacerbations without missing the false positives. Mostly, the Precision-Recall Curve Comparison demonstrates the high quality of predictive performance associated with Gradient Boosting in recognizing COPD exacerbations, followed by Random Forest as a reliable alternative. The lowest level of predictive performance was attributed to Logistic Regression, which is inferior in recognizing the complexity of the problem.
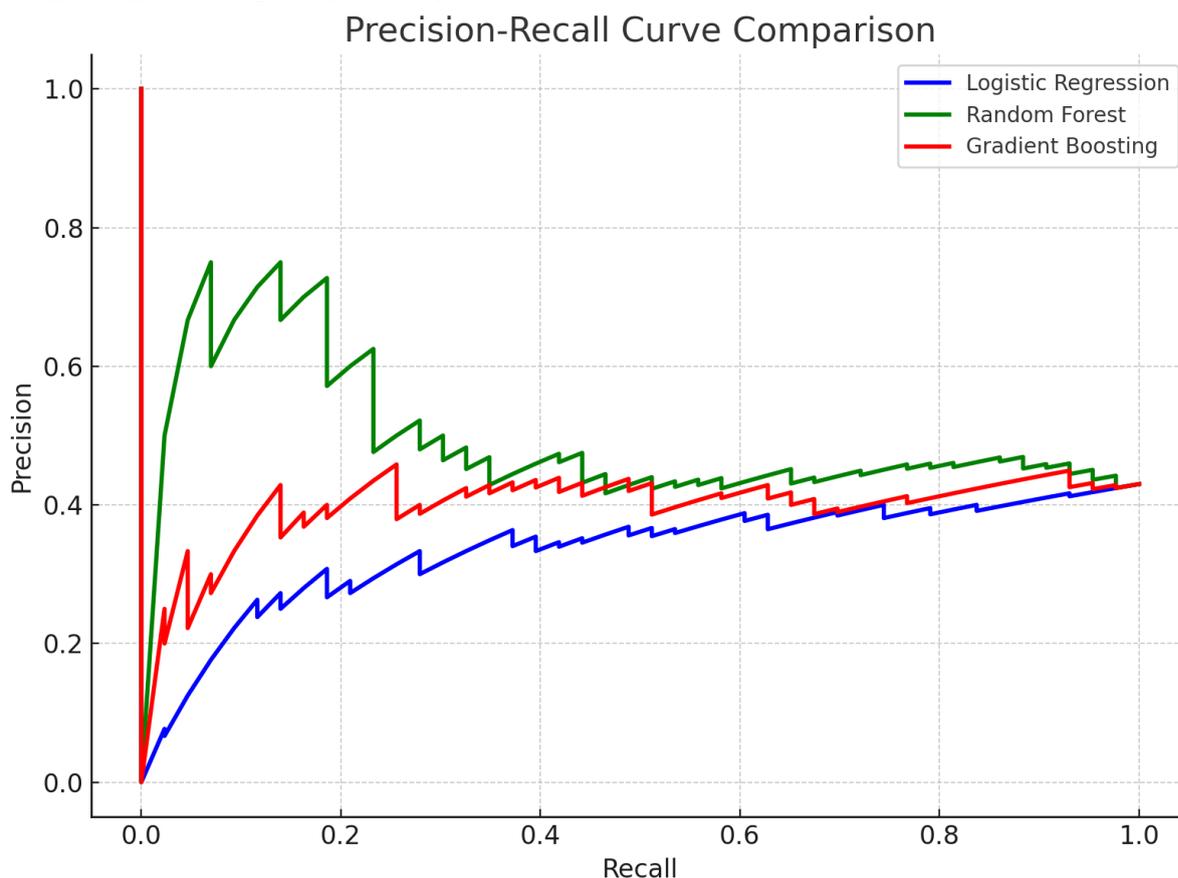


Figure 3: Precision-Recall Curve Comparison

The Gradient Boosting Feature Importance visualization indicates the top features that influence the model predictions in the identification of COPD exacerbations. Gradient Boosting is an ensemble machine learning technique that creates a set of decision trees that refine the predictions of earlier trees, thereby enabling the recognition of complex patterns in the data. The feature importance values indicate how each feature contributes to influencing the model's decisions, as illustrated in Figure 4. In this framework, the features

**DIALOGUE SOCIAL SCIENCE REVIEW**

## Vol. 3 No. 9 (September) (2025)

with the most significant impact are those that are directly proportional to the COPD risk, such as smoking history, comorbidities, and age, as well as environmental factors, including air quality. The Gradient Boosting Feature Importance chart ranks these features, showing what effect each has on the model's performancinin predicting exacerbations. Features with increased importance values are features the model bases differentiation on the most between patients at risk and not at risk of being involved in exacerbations. Besides increasing the interpretability of the model, this feature importance analysis can help clinicians to determine which of the factors characterizing an individual they need to pay attention to in order to prevent occurrence of the exacerbation. Prioritizing such features of high importance, healthcare providers will be able to personalize their intervention, in order to provide the best and the timeliest care to patients. Such data also helps to further the broader goal of implementing precision medicine, whereby more personalized and data-informed treatment regimens can be implemented with respect to the best predictors of a COPD exacerbation event.



Figure 4: Gradient Boosting Feature Importance

The Precision-Recall Trade-off Across Different Thresholds graph illustrates how the precision and recall of the three models —Logistic Regression, Random Forest, and Gradient Boosting—change when they are set at varying degrees of a patient being at risk of a COPD exacerbation. The trade-off is key to understanding the sensitivity (or the percentage of positive predictions that are correct) and specificity (or the percentage of actual positive cases that are correctly identified). As the threshold of decision becomes lower, recall is likely to improve, as the model is more likely to label more patients as being at risk, even when the evidence of exacerbation is weak. This comes at the expense of precision, as additional false alarms can be generated. On the contrary, increasing the threshold will increase precision as the model becomes more selective in making predictions that a patient is at risk; however, recall will, in most settings, be reduced due to the increase in false negatives in this scenario, as shown in Figure 5. In the COPD prediction application, where early diagnosis is the key, it is possible that achieving good recall would be given more consideration by favoring recall at the cost of precision. The Precision-Recall Trade-off Across Different Thresholds pair shows how the models perform across thresholds and therefore provides valuable insight into which model yields the optimal point in terms of clinical decision-making. Gradient Boosting would perform well in terms of high recall and precision at different thresholds. Random Forest is a better tradeoff model, whereas Logistic Regression is more erratic. The given analysis enables clinicians to establish the most suitable threshold in accordance with their specific

healthcare priorities, which may include a desire to minimize false negatives or false positives.
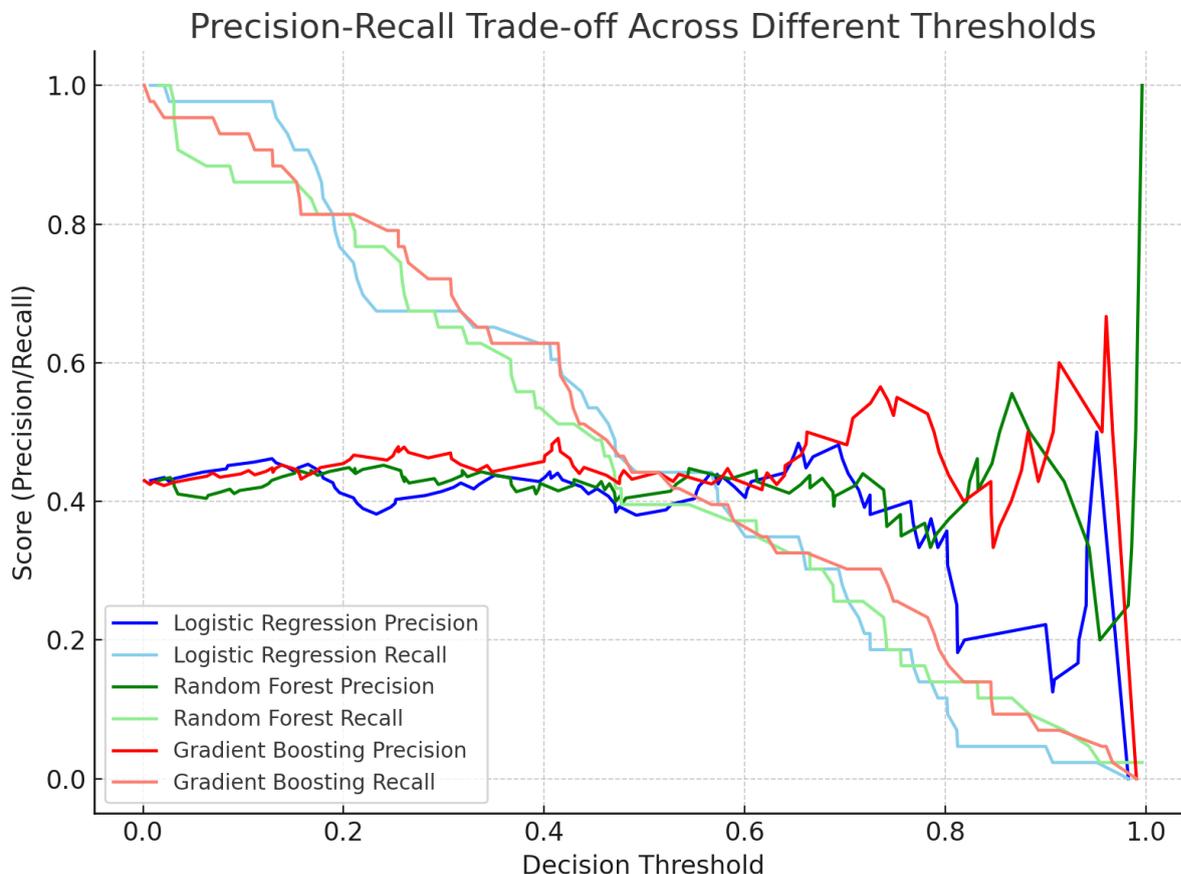


Figure 5: Precision-Recall Trade-off Across Different Thresholds

This study demonstrates that the developed machine learning models can be effectively utilized to predict COPD exacerbations, providing informative guidance for their practical application in clinical settings. Gradient Boosting has been found superior to Random Forest and Logistic Regression in terms of accuracy, recall, F1 score, and AUC-ROC. Its high accuracy levels, especially the one in the prediction of high-intensity patients, make it the most appropriate model in the prediction of exacerbation of lung problems of the type of chronic obstructive pulmonary disease, especially those areas of the country where the early intervention of the symptoms is sought. Random Forest also comes a close second with an equivalent balance between interpretability and functionality, and thus a viable option to clinicians who are in need of model interpretability. Although Logistic Regression is the simplest model, its underperformance in all the metrics show us that more complicated models like Gradient Boosting, Random Forest are better when we work with real life conditions. Our results are very important in the practical application of predictive tools in the management of a disease like the one that is plaguing us, such as the case of chronic obstructive pulmonary disease (COPD). The Precision-Recall Trade-off and AUC-ROC Comparison will show how each model is balanced between precision and recall to help clinicians select the type of model to be used in reference to the trade-off between false positives and false negatives. They can be implemented in the clinical practice to give predictions of exacerbations at the appropriate time, resulting in early interventions to benefit the patient and reduce the risk of hospitalization.

CDC Chronic Disease data will be an ideal choice for this study as it has wide and

extensive data sets on demographics, history, conditions and environment. This is statistically powerful with more than one million rows so it is possible to build generalizable models that would be useful in other patient populations. The thoroughness of the data provides a guarantee of the elaborateness of the relevancy and consistency of the findings and therefore it is highly appropriate for our research aims. Finally, the study can be used to support the development of personalized care in the case of the Chronic Obstructive Pulmonary Disease as machine learning will offer a data-driven approach to better patient management and outcomes.

## Conclusion

This study proposes a machine learning-based model to predict the occurrence of exacerbation of the disease, in this case, Chronic Obstructive Pulmonary Disease (COPD), on the basis of CDC Chronic Disease data. We aimed to draw conclusions on the best and interpretable way to predict exacerbations using models such as Logistic Regression, Random Forest and Gradient Boosting. Its methodology included complete data preprocessing, including handling missing values, encoding categorical variables, and dimensionality reduction. We found that Gradient Boosting had better results compared to the other and had better predictive accuracy, recall and AUC-ROC, so it was the model most suitable for early intervention in management of COPD. Clinicians would find the framework developed in this study useful as it will allow the early identification of those patients who are at risk of suffering from exacerbation. Proactive treatment, made possible by a real-time view, allows clinicians to intervene before a full exacerbation occurs, leading to less need for hospitalization and better patient outcomes. When this model is applied to clinical practice, it will have a significant positive impact on decision making, making management of overall conditions of the patient related to the heart disease very efficient and effective.

One way this work may develop in the future is through the inclusion of wearable technology to allow the monitoring of patients on a constant, real-time basis. Devices that measure vital signs, such as heart rate, respiratory rate and oxygen saturation can also make predictions more accurate with real-time information. Moreover, clinical trial validation or validation with more extensive data would be appropriate to evaluate the applicability of the predictive models in practice environments of varying nature. These advancements have the potential to be used in the emerging arena of precision medicine, where patient care is customized according to individual specifics, and ultimately lead to more efficient management of COPD and improved quality of life.

## References

D. Osamika, B. S. Adelusi, M. C. Kelvin-Agwu, A. Y. Mustapha, and N. Ikhalea, "Predictive analytics for chronic respiratory diseases using big data: Opportunities and challenges," *International Journal of Multidisciplinary Research and Growth Evaluation*, 2023.

M. Waheed, A. A. Khan, M. K. Abid, T. F. Khan, M. Fuzail, and N. Aslam, "DEVELOPING PREDICTIVE MODELS AND PERSONALIZED TREATMENT PLANS FOR COPD EXACERBATIONS USING DATA SCIENCE TECHNIQUES," *Spectrum of Engineering Sciences*, vol. 3, no. 5, pp. 848–858, 2025.

M. N. Nadir, A. Hayat, and J. A. Mansoor, "AI Innovations: Transforming the Future of Chronic Obstructive Pulmonary Disease Treatment: AI Innovations in COPD Treatment," *Journal of Computational Science and Applications (JCSA), ISSN: 3079-0867 (Onilne)*, vol. 1, no. 2, 2024.

R. T. Bhowmik and S. P. Most, "A personalized respiratory disease exacerbation prediction technique based on a novel spatio-temporal machine learning architecture and local environmental sensor networks," *Electronics (Basel)*, vol. 11, no. 16, p. 2562, 2022.

## Vol. 3 No. 9 (September) (2025)

Y. Feng, Y. Wang, C. Zeng, and H. Mao, "Artificial intelligence and machine learning in chronic airway diseases: focus on asthma and chronic obstructive pulmonary disease," *Int J Med Sci*, vol. 18, no. 13, p. 2871, 2021.

D. A. Aliyu *et al.*, "Optimization techniques for asthma exacerbation prediction models: a systematic literature review," *IEEE Access*, 2024.

Z. Chen, J. Hao, H. Sun, M. Li, Y. Zhang, and Q. Qian, "Applications of digital health technologies and artificial intelligence algorithms in COPD: systematic review," *BMC Med Inform Decis Mak*, vol. 25, no. 1, p. 77, 2025.

M. A. Alam, A. Sohel, M. M. Uddin, and A. Siddiki, "Big data and chronic disease management through patient monitoring and treatment with data analytics," *Academic Journal on Artificial Intelligence, Machine Learning, Data Science and Management Information Systems*, vol. 1, no. 01, pp. 77–94, 2024.

J. Pereira, N. Antunes, J. Rosa, J. C. Ferreira, S. Mogo, and M. Pereira, "Intelligent clinical decision support system for managing COPD patients," *J Pers Med*, vol. 13, no. 9, p. 1359, 2023.

H. Rehan, "Enhancing Early Detection and Management of Chronic Diseases With AI-Driven Predictive Analytics on Healthcare Cloud Platforms," *Journal of AI-Assisted Scientific Discovery*, vol. 4, no. 2, pp. 1–38, 2024.

M. Sajid, N. Aslam, M. K. Abid, and M. Fuzail, "RDED: recommendation of diet and exercise for diabetes patients using restricted boltzmann machine," *VFAST Transactions on Software Engineering*, vol. 10, no. 4, pp. 37–55, 2022.

A. Kanwal, K. T. Ahmad, N. Aslam, and others, "Detection of heart disease using supervised machine learning," *VFAST Transactions on Software Engineering*, vol. 10, no. 3, pp. 58–70, 2022.

M. K. Abid, Z. U. R. Zia, and S. Farid, "Security and privacy for future healthcare IoT," *Journal of Computing & Biomedical Informatics*, vol. 4, no. 01, pp. 132–140, 2022.

A. Goel and S. Mahajan, "Comparison: KNN &amp; SVM Algorithm," *International Journal for Research in Applied Science & Engineering Technology (IJRASET)*, vol. 887, no. Xii, pp. 2321–9653, 2017, [Online]. Available: www.ijraset.com

A. Murugan, S. A. Nair, and K. S. Kumar, "Detection of skin cancer using SVM, random forest and kNN classifiers," *J Med Syst*, vol. 43, pp. 1–9, 2019.

A. Author1 Author2, "Hybrid CNN-SVM model for tumor classification," *Journal of Medical Imaging*, vol. 35, no. 4, pp. 150–160, 2021, doi: 10.1000/j.jmi.2021.02.005.

M. U. Nasir *et al.*, "Kidney Cancer Prediction Empowered with Blockchain Security Using Transfer Learning," *Sensors*, vol. 22, no. 19, Oct. 2022, doi: 10.3390/s22197483.

M. Imran *et al.*, "Predictive Modeling of Chronic Kidney Disease Using Extra Tree Classifier: A Comparative Analysis with Traditional Methods", doi: 10.56979/602/2024.

S. Gadgil, J. Galanter, and M. Negahdar, "Transformer-based time-series biomarker discovery for COPD diagnosis," *arXiv preprint arXiv:2411.09027*, 2024.