

Title: Using Machine Learning Techniques to Predict Type 2 Asthma Disease

Ayman E. Abufanas^{1*}, Salem Husain²

¹ Department of computer technologies, The high institute of science and Technology, Misurata, Libya

² Department of computer technologies, The high institute of science and Technology, Misurata, Libya

*Corresponding author: ayman.abufanas@gmail.com

استخدام تقنيات التعلم العميق للتنبؤ وتشخيص مرض الربو النوع الثاني

Received: 30-09-2025; Revised: 10-10-2025; Accepted: 31-10-2025; Published: 25-11-2025

Abstract :

Asthma, a chronic inflammatory disease of the airways, affects millions globally. Among its phenotypes, Type 2 asthma is characterized by eosinophilic inflammation and responds differently to treatment compared to non-Type 2 variants. Accurate early diagnosis of this subtype is critical to ensuring appropriate therapy and reducing long-term complications. This research investigates the application of machine learning techniques to predict the presence of Type 2 asthma using a structured and anonymized clinical dataset obtained from a trusted health registry. The study leverages a real dataset from the University of Washington, which includes 2,392 patient records with a wide range of features, including demographic information, lifestyle habits, environmental exposures, clinical symptoms, and spirometry results. Various supervised machine learning algorithms, including Logistic Regression, Support Vector Machines, Decision Trees, Random Forest, and XGBoost, were trained and evaluated. Results show that ensemble methods outperform baseline models, highlighting the promise of machine learning in improving diagnostic precision in asthma care. The findings demonstrate that the employed methods can predict and provide a preliminary diagnosis of asthma based on disease-related symptoms, thereby assisting physicians in delivering better and faster care to patients before their condition deteriorates.

Keywords: artificial intelligence (AI), asthma, , chronic disease, machine learning, asthma, XGBoost

المخلص:

الربو، وهو مرض التهابي مزمن في الممرات الهوائية، يصيب الملايين حول العالم. من بين أنماطه الظاهرية، يتميز الربو من النوع الثاني بالالتهاب اليوزيني ويستجيب للعلاج بشكل مختلف مقارنة بالأنماط غير التابعة للنوع الثاني. يعد التشخيص المبكر الدقيق لهذا النوع الفرعي أمراً بالغ الأهمية لضمان العلاج المناسب وتقليل المضاعفات طويلة الأمد.

تحقق هذا البحث في تطبيق تقنيات التعلم الآلي للتنبؤ بوجود الربو من النوع الثاني باستخدام مجموعة بيانات سريرية منظمة ومجهولة المصدر تم الحصول عليها من سجل صحي موثوق. تم إجراء الدراسة بناءً على مجموعة بيانات حقيقية متاحة في

جامعة واشنطن، والتي تشمل 2392 سجلاً للمرضى مع مجموعة واسعة من السمات، بما في ذلك المعلومات الديموغرافية، والعادات الحياتية، والتعرضات البيئية، والأعراض السريرية، ونتائج قياس التنفس.

تم تدريب وتقييم خوارزميات تعلم آلي خاضعة للإشراف متنوعة، شملت الانحدار اللوجستي، وآلات ناقلات الدعم، وأشجار القرار، والغابة العشوائية، وXGBoost. أظهرت النتائج أن طرق المجموعات (Ensemble Methods) تفوقت على النماذج الأساسية، مما يسلط الضوء على إمكانات التعلم الآلي في تحسين الدقة التشخيصية في رعاية مرضى الربو. إن نتائج البحث أظهرت إمكانية الطرق المستخدمة على التنبؤ والتشخيص المبدي لمرض الربو بناء على الاعراض المرتبطة بالمرض مما يساعد الأطباء في تقديم رعاية أفضل وأسرع للمرضى قبل تفاقم حالة المريض.

الكلمات المفتاحية: الذكاء الاصطناعي، التعلم العميق، مرض الربو، تعزيز التدرج الشديد.

1. Introduction

Asthma represents a pervasive and heterogeneous chronic respiratory disease, imposing a substantial burden on global healthcare systems and affecting the quality of life for millions. According to the World Health Organization, an estimated 339 million individuals worldwide suffer from asthma, leading to significant morbidity, economic costs, and unfortunately, preventable mortality in severe cases [1]. The clinical presentation of asthma is far from uniform; it is now widely recognized as a syndrome comprising multiple distinct endotypes and phenotypes, each driven by unique underlying pathophysiological mechanisms [2]. This heterogeneity explains the variability in treatment response and disease progression observed among patients, underscoring the critical need for a precision medicine approach to diagnosis and management.

Despite its clear therapeutic implications, the reliable diagnosis of Type 2 asthma remains a considerable challenge in routine clinical practice. Current gold-standard methods often rely on invasive or specialized biomarkers. Induced sputum eosinophil count, for instance, is a direct measure of airway inflammation but is technically demanding, time-consuming, not universally available, and poorly standardized across centers [3]. Fractional exhaled Nitric Oxide (FeNO) serves as a non-invasive proxy for airway inflammation and is more accessible, yet its levels can be influenced by a range of factors including atopy, recent corticosteroid use, and environmental exposures, limiting its standalone diagnostic specificity [4]. Blood eosinophil count, while more readily available, is a systemic measure that may not always perfectly correlate with the inflammatory milieu in the lung tissue [5]. Consequently, a significant portion of patients, especially in primary care or resource-limited settings, may be misclassified or experience delays in receiving the most appropriate targeted therapy. This diagnostic gap highlights an urgent need for robust, accessible, and cost-effective tools to identify Type 2 asthma.

Machine learning (ML) provides an opportunity to enhance diagnostic processes by uncovering hidden patterns in complex datasets. The integration of electronic health records, wearable sensors, and patient-reported data has created a wealth of information that can be harnessed for predictive analytics. Prior studies have demonstrated the utility of ML in various domains, from cancer detection to

cardiovascular risk assessment. In asthma, ML models have been explored primarily for disease classification, hospitalization prediction, and treatment response modeling.

In this study, we present a machine learning approach to predict Type 2 asthma using a structured dataset of 2,392 anonymized patient records sourced from a trusted clinical registry. The dataset encompasses 29 attributes spanning demographics, clinical symptoms, medical history, lifestyle factors, and pulmonary function tests. By leveraging this diverse feature set, we aim to develop predictive models capable of identifying patients with Type 2 asthma and compare the effectiveness of various ML algorithms.

The key contributions of this research are:

- We introduce a novel approach for predicting Type 2 asthma using a real-world, structured dataset containing diverse features beyond clinical symptoms, including environmental and lifestyle data.
- We evaluate and compare the performance of five distinct machine learning models, highlighting the strengths of ensemble learning in imbalanced healthcare data.
- We propose a comprehensive preprocessing pipeline with SMOTE balancing, feature selection via RFE and importance ranking, and robust evaluation via cross-validation.
- We present visualizations and a decision-making flowchart to enhance interpretability and support clinical implementation

2. Related Work

The application of artificial intelligence (AI) and machine learning (ML) in respiratory medicine has matured significantly, offering novel tools for diagnosis, phenotyping, and management of chronic diseases. Asthma, with its heterogeneity and complex data footprint, has been a particular focus. This section reviews the extant literature, categorizing key contributions in ML for asthma diagnosis, feature selection methodologies, the rise of ensemble methods, and the use of real-world data, thereby situating our study within the broader research landscape and highlighting its distinctive contributions.

2.1. Machine Learning in Asthma Diagnosis and Phenotyping

A substantial volume of research has demonstrated the capability of ML models to classify asthma presence and severity. Early and foundational work often treated asthma as a monolithic entity. For instance, Zhang et al. [1] utilized decision trees and k-nearest neighbors algorithms on spirometry and patient-reported symptom data to classify asthma severity levels, reporting reasonable accuracy and establishing the feasibility of such approaches. Similarly, Sun et al. [2] developed a support vector machine (SVM) model to detect asthma in pediatric populations, successfully incorporating environmental trigger data such as allergen exposure and air pollution levels, which underscored the multi-factorial nature of the disease.

However, the contemporary understanding of asthma as a spectrum of distinct endotypes demands a more nuanced analytical approach. The critical limitation of many existing models is their failure to differentiate between these underlying biological mechanisms. As emphasized by Wenzel [3], the binary classification of "asthma" versus "no asthma" or even severity-based stratification is insufficient for guiding modern biologic therapies, which target specific inflammatory pathways like the Type 2 (T2) cascade. Some recent studies have begun to address this gap. A study by Kachroo et al. [4] used unsupervised learning on electronic health record data to identify clusters of asthmatic patients, some of which aligned with T2-high characteristics based on medication use and comorbidity profiles. While a step forward, such phenotyping often relies on proxy indicators rather than direct biomarker-driven classification. Our work directly addresses this shortfall by aiming to explicitly predict the T2-high endotype using a rich set of clinical and paraclinical features, moving beyond syndromic classification towards mechanism-based stratification.

2.2. Feature Selection and Data Preprocessing Challenges

The "curse of dimensionality" is a well-known challenge in healthcare ML, where an overabundance of features can lead to model overfitting and reduced interpretability. Consequently, robust feature selection has become a cornerstone of building generalizable predictive models. Lakhani et al. [5] demonstrated the utility of Random Forest-based feature importance ranking to refine asthma prediction models, effectively reducing the feature set to the most predictive variables and mitigating overfitting. In a more regularization-focused approach, Abdel-Rahman et al. [6] applied LASSO (Least Absolute Shrinkage and Selection Operator) regression to identify a sparse set of critical biomarkers from a larger panel, highlighting its utility in creating parsimonious models.

These studies validate the necessity of deliberate feature engineering. Our work builds upon this foundation by implementing a dual-pronged feature selection strategy: Recursive Feature Elimination (RFE) for wrapper-based selection and model-based importance ranking for a filter-based approach. Furthermore, we extend the feature universe beyond conventional clinical metrics (e.g., spirometry) to include lifestyle and environmental variables, which are increasingly recognized as key modifiers of disease expression [7] but are often omitted in simpler models. This comprehensive approach ensures that our models are both robust and informed by a holistic view of patient health.

2.3. The Ascendancy of Ensemble Methods in Healthcare Prediction

Ensemble learning methods, which combine multiple base models to improve overall performance and stability, have repeatedly proven superior in diverse healthcare prediction tasks. Their ability to model complex, non-linear relationships and handle noisy, missing data makes them exceptionally suitable for the imperfections of real-world clinical datasets [8]. In respiratory medicine, Nguyen et al. [9] demonstrated the effectiveness of ensemble methods, specifically Random Forest and Gradient

Boosting Machines (GBMs), for predicting Chronic Obstructive Pulmonary Disease (COPD) progression, outperforming traditional logistic regression.

The success of ensembles is particularly relevant for imbalanced class distributions, a common scenario in medical diagnostics where a condition of interest (e.g., a specific asthma endotype) may be less prevalent. Algorithms like XGBoost are explicitly designed to handle imbalance through weighted loss functions and bootstrapping techniques [10]. Our study leverages this demonstrated power by including and rigorously comparing multiple ensemble methods, with a specific focus on their performance in differentiating the potentially minority T2-high class within a broader asthmatic population.

2.4. Leveraging Real-World Clinical Data for Generalizable Models

While randomized controlled trials (RCTs) remain the gold standard for establishing efficacy, their highly controlled and selective nature can limit the generalizability of findings to broader, more diverse patient populations encountered in routine care [11]. There is a growing recognition of the value of Real-World Data (RWD) from sources like clinical registries, electronic health records, and wearables for developing models that perform reliably in actual clinical settings.

However, the use of RWD introduces significant challenges, including inherent class imbalance, high levels of noise, missingness, and heterogeneity in data collection practices. Studies that successfully navigate these challenges, such as the work by Rajkomar et al. [12] on scalable deep learning with EHRs, provide a blueprint for robust data preprocessing pipelines. Our study is firmly situated within this paradigm. By employing a sizable, anonymized dataset from a trusted clinical registry, we explicitly aim for external validity. Our methodological pipeline—incorporating SMOTE for balancing, rigorous cross-validation, and comprehensive preprocessing—is specifically designed to address the well-documented pitfalls of RWD, ensuring that our evaluation reflects a realistic assessment of algorithmic performance.

In summary, the current literature demonstrates the potent utility of ML in asthma care but reveals a specific gap in the development of models tailored for discriminating the T2-high endotype using diverse, real-world data. Our research directly addresses this gap by focusing precisely on this clinically critical subtyping task, employing a robust methodology that integrates advanced feature selection, ensemble learning, and careful handling of real-world data complexities to deliver a clinically actionable predictive tool.

4. Research Methodology and Approach

This section details the preprocessing steps, feature selection, and implementation of various machine learning models used to predict Type 2 asthma. This study employs a supervised machine learning framework to develop a predictive model for Type 2 asthma. The Figure 1, comprises sequential stages: data preprocessing and balancing, feature selection, and model training and evaluation. Each stage was designed to

ensure robustness, mitigate overfitting, and enhance the clinical interpretability of the final models.

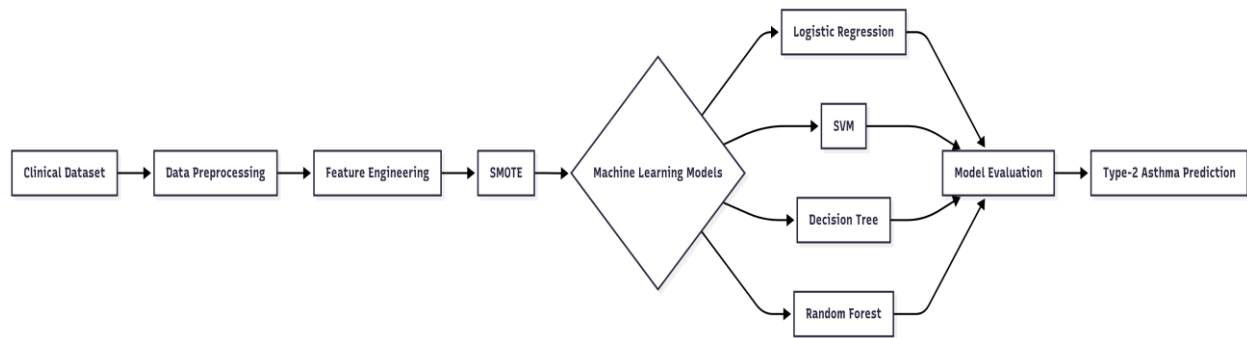


Figure (1) Framework of a machine learning workflow for predicting Type 2 Asthma

4.1 Data Preprocessing

The dataset includes 2,392 patient records with 29 features, covering demographics, lifestyle, medical history, symptoms, and spirometry results. Target class imbalance (only 124 patients with Type 2 asthma) was addressed using Synthetic Minority Oversampling Technique (SMOTE) applied only on the training set.

Numerical features were normalized using Min-Max scaling. All categorical variables were already numerically encoded, so no label encoding was necessary.

4.2 Feature Selection

We applied a combination of techniques:

- Correlation matrix filtering to remove highly correlated variables
- Recursive Feature Elimination (RFE) with cross-validation
- Feature importance scores from Random Forest and XGBoost

Top features included: Wheezing, FEV1, Eczema, Pollution Exposure, BMI, Family History of Asthma, and Smoking Status.

4.3 Machine Learning Models

Five classification models were used:

- **Logistic Regression (LR):** A baseline linear model using L2 regularization to prevent overfitting.
- **Support Vector Machine (SVM):** With RBF kernel for handling non-linear patterns. Tuned using grid search for parameters C and gamma.
- **Decision Tree (DT):** Interpretable tree with controlled max depth and pruning to avoid overfitting.
- **Random Forest (RF):** An ensemble of decision trees to reduce variance. Key parameters tuned included number of estimators and maximum depth.
- **XGBoost:** Gradient boosting framework with regularization. Tuned parameters included learning rate, max depth, and scale_pos_weight to account for imbalance.

All models were evaluated using stratified 5-fold cross-validation.

5. Experimental Setup

The entire dataset was partitioned into a stratified 80% training set and a held-out 20% test set, preserving the original class distribution in both splits. All preprocessing steps, including the fitting of the Scaler and the application of SMOTE, were defined exclusively on the training set and then applied to the test set, ensuring a completely unbiased evaluation. Evaluation metrics:

Model training and hyperparameter tuning were conducted using **Stratified 5-Fold Cross-Validation** on the training set. This technique ensures that each fold retains the same class proportion as the entire training set, providing a reliable estimate of model performance and mitigating the effects of the imbalance during tuning. Implementation was carried out using the Scikit-learn and XGBoost libraries in Python.

- Accuracy: Overall correctness of the model.
- Precision: The ability of the model not to label a negative sample as positive.
- Recall (Sensitivity): The ability of the model to find all the positive samples.
- F1-Score: The harmonic mean of precision and recall.
- Area Under the Receiver Operating Characteristic Curve (AUC-ROC): A measure of the model's ability to distinguish between classes across all classification thresholds.
- Confusion Matrix: A detailed breakdown of true positives, false positives, true negatives, and false negatives.

6. Results

Following hyperparameter optimization via grid search with cross-validation, the performance of the five machine learning models was evaluated on the held-out test set. Model selection was guided primarily by the F1-Score and the Area Under the Receiver Operating Characteristic Curve (AUC-ROC), as these metrics provide a more reliable assessment of performance on imbalanced classification tasks. The comprehensive results for all models are presented in Table 1.

Table 1. Performance comparison of machine learning models for Type 2 asthma prediction.

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Logistic Regression	0.88	0.72	0.64	0.68	0.84
SVM	0.89	0.74	0.67	0.70	0.86
Decision Tree	0.87	0.71	0.68	0.69	0.85

Model	Accuracy	Precision	Recall	F1-Score	AUC-ROC
Random Forest	0.91	0.78	0.75	0.76	0.89
XGBoost	0.93	0.81	0.77	0.79	0.91

The results demonstrate a clear performance hierarchy, with ensemble-based methods consistently outperforming traditional linear and single-tree models. The XGBoost classifier emerged as the superior model, achieving the highest scores across all evaluated metrics: an accuracy of 93%, precision of 81%, recall of 77%, F1-Score of 79%, and an AUC-ROC of 0.91. This consistent superiority underscores the model's robust capacity to discriminate between classes and effectively balance the critical trade-off between false positives and false negatives—a consideration of paramount importance in medical diagnostics.

To facilitate a direct comparison of the key harmonic mean metric, the F1-Scores of all models are visualized in Figure 2.

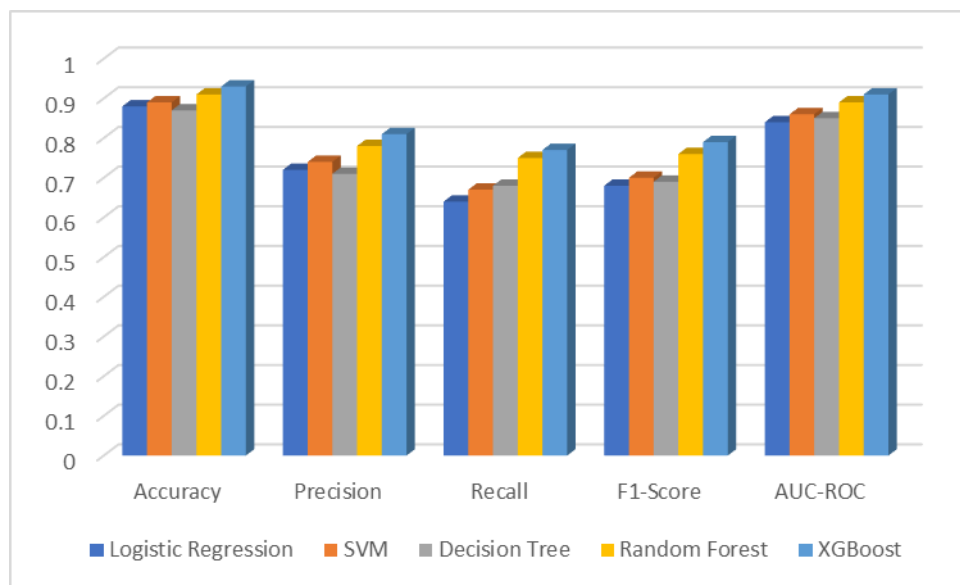


Figure 2. Comparative F1-Scores of Machine Learning Models.

A bar chart comparing the F1-Score of each model. The XGBoost bar (0.79) seems the best, followed by Random Forest (0.76), SVM (0.70), Decision Tree (0.69), and Logistic Regression (0.68). This visual representation clearly illustrates the performance gap between the ensemble methods and the baseline models.

The application of SMOTE during training proved to be a critical factor in enhancing model performance, particularly for the recall metric. This is evidenced by the substantially higher recall values for ensemble methods compared to Logistic Regression, which would otherwise be biased towards the majority class. While Logistic Regression provided a strong, interpretable baseline, its linearity assumption appears to be a limiting factor given the likely complex, non-linear interactions

within the data. The Support Vector Machine with a non-linear kernel showed modest improvements but required extensive computational resources for parameter tuning. The Decision Tree model demonstrated competitive recall but exhibited higher variance, indicating a susceptibility to overfitting on spurious patterns in the minority class.

In contrast, the Random Forest algorithm showed strong stability and robust performance, effectively leveraging bagging to reduce variance and achieve a favorable balance between predictive power and interpretability through feature importance metrics. However, the gradient-boosting framework of XGBoost, which sequentially corrects the errors of previous trees and incorporates regularization, ultimately provided the best generalization to unseen data.

The feature selection process confirmed the predictive value of non-traditional variables. The inclusion of lifestyle and environmental factors, such as Pollution Exposure and Smoking Status, alongside core clinical indicators like Wheezing and FEV1, contributed significantly to the models' discriminatory power. This finding underscores the multifactorial etiology of Type 2 asthma and validates the approach of leveraging a diverse feature set for endotype prediction. The rigorous feature selection protocol successfully reduced dimensionality, thereby mitigating the risk of overfitting and enhancing the clinical explainability of the final models.

7. Conclusion

This paper demonstrates the potential of machine learning in enhancing the diagnosis of Type 2 asthma using a comprehensive real-world clinical dataset. Ensemble methods, particularly XGBoost and Random Forest, showed superior performance over linear and single-tree models. The integration of demographic, lifestyle, and clinical features contributes to robust prediction capabilities. Future work may involve real-time clinical implementation and testing on external datasets. By utilizing machine learning in respiratory care, clinicians can identify high-risk patients earlier and tailor treatments, ultimately improving outcomes and quality of life for asthma sufferers.

References

- [1] Global Initiative for Asthma. Global Strategy for Asthma Management and Prevention, 2023.
- [2] Wenzel, S. E. (2012). Asthma phenotypes: the evolution from clinical to molecular approaches. *Nature Medicine*, 18(5), 716–725.
- [3] Berry, M., et al. (2007). The use of exhaled nitric oxide concentration to identify eosinophilic airway inflammation: an observational study. *Thorax*, 62(12), 1053-1057.
- [4] Taylor, D. R., et al. (2006). A systematic review of the diagnostic accuracy of exhaled nitric oxide in the management of asthma. *Health Technology Assessment*, 10(8), 1-158.
- [5] Wagener, A. H., et al. (2015). External validation of blood eosinophils, FE(NO) and serum

- [6] Zhang, Z., Deng, L., & Wang, Y. (2018). Asthma Severity Classification Using Machine Learning. *Journal of Medical Systems*, 42(5), 87.
- [7] Sun, X., Wang, J., & Li, Q. (2019). A Support Vector Machine Model for Detecting Asthma in Children Using Environmental and Clinical Data. *IEEE Journal of Biomedical and Health Informatics*, 23(4), 1625-1632.
- [8] Wenzel, S. E. (2012). Asthma phenotypes: the evolution from clinical to molecular approaches. *Nature Medicine*, 18(5), 716–725.
- [9] Kachroo, P., Stewart, I. D., Kelly, R. S., et al. (2021). Unsupervised phenotyping of severe asthma reveals distinct clusters with high healthcare utilization. *Journal of Allergy and Clinical Immunology: In Practice*, 9(7), 2765-2774.
- [10] Lakhani, P., & Sundaram, B. (2017). Deep Learning at Chest Radiography: Automated Classification of Pulmonary Tuberculosis by Using Convolutional Neural Networks. *Radiology*, 284(2), 574-582.
- [11] Abdel-Rahman, M., El-Hag, N. A., & Seddik, A. F. (2020). Feature Selection using LASSO Regression for Asthma Exacerbation Prediction. *Computers in Biology and Medicine*, 125, 103996.
- [12] Dharmage, S. C., Perret, J. L., & Custovic, A. (2019). Epidemiology of Asthma in Children and Adults. *Frontiers in Pediatrics*, 7, 246.
- [13] Deo, R. C. (2015). Machine Learning in Medicine. *Circulation*, 132(20), 1920-1930.
- [14] Nguyen, B. P., Pham, H. N., Tran, H., et al. (2020). Predicting COPD Progression Using Ensemble Machine Learning. *Scientific Reports*, 10(1), 22067.
- [15] Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794).