

ENHANCED MACHINE LEARNING FOR SPIROMETRY CURVE ASSESSMENT IN CHRONIC RESPIRATORY DISEASES

Sumit Kumar Roy¹, Saurabh Gupta², Dibakar Sahu³, Hemlata Sinha^{4*}

¹Department of Biomedical Engineering National Institute of Technology Raipur, India
skroy.phd2022.bme@nitrr.ac.in

²Department of Biomedical Engineering National Institute of Technology Raipur, India
sgupta.bme@nitrr.ac.in

³Department of Pulmonary Medicine AIIMS Raipur, India
dibakarsahu@aiimsraipur.edu.in

⁴Department of ETC,SSIPMT Raipur, India
sinha.hemlata552@gmail.com

Abstract

Chronic respiratory diseases, including Chronic Obstructive Pulmonary Disease (COPD) and asthma, pose a significant global health challenge with far-reaching economic implications. This study introduces a machine learning-based approach for predicting exacerbation risks in respiratory diseases, specifically focusing on COPD and asthma. The framework employs an advanced machine learning architecture to enable real-time and precise detection of respiratory events. Spirometry standards mandate that correct maneuvers should be devoid of cough artifacts, especially in the initial seconds of forced exhalation. Achieving this standard becomes inherently challenging when patients manifest increased coughing tendencies during the examination. This study aims to facilitate unsupervised or minimally supervised spirometry measurements, expanding accessibility in diverse settings, including home monitoring and general practitioner oversight. By prioritizing specificity in the machine learning techniques, the model achieves a comprehensive understanding and accurate classification of symptoms, thereby contributing to personalized risk assessment.

Index Terms—Chronic Obstructive Pulmonary Disease, Asthma, Machine Learning, Spirometry, Random Forest, KMean

Abstract—

I. INTRODUCTION

In the era of burgeoning big data, opportunities and challenges abound in understanding, analyzing, and extracting knowledge. Traditional statistical analyses, once adequate, may now fall short in the face of pervasive measurability. Amidst this backdrop, the healthcare sector stands to gain immensely from identifying shared patterns across diverse patient datasets. This research focus on two prevalent respiratory conditions: asthma and CHRONIC OBSTRUCTIVE PULMONARY DISEASE (COPD) [1].

Asthma, characterized by chronic airway inflammation, triggers bronchospasm and airway constriction upon exposure to environmental stimuli. While a genetic predisposition may exist, symptoms typically manifest following exposure to trigger viral infection and allergy. COPD encompasses chronic bronchitis and emphysema, both contributing to persistent airway obstruction and daily breathlessness, with smoking being a primary contributor. Globally, asthma affects approximately 300 million patients, leading to 250,000 annual deaths, while COPD

afflicts 330 million patients, causing around 3 million deaths each year. The urgent need to develop tools for the early prediction and diagnosis of these respiratory diseases is evident [2]. This study employs machine learning techniques to scrutinize the factors underpinning the diagnosis of asthma and COPD. Aligned with existing research, our methodology explores both linear and non-linear models, placing emphasis on the importance of complexity and interpretability [3].

In this landscape, efforts to detect coughs using diverse sensors such as ECG sensors, thermistors, chest belts, and oximeters have been explored. Wearable cough detection systems, including those integrated into smartwatches for ambulatory monitoring, signify recent innovations in this domain. Intriguingly, cough detection has extended beyond human patients to encompass veterinary monitoring of farm animals, underscoring the versatility and applicability of such algorithms [1] [2]. The predominant focus in existing literature on cough detection has revolved around audio (sound) signals or accelerometer recordings. However, this study diverges from this trend by

examining airflow signals passing through the spirometer. Analyzing airflow signals presents distinct advantages, notably the reduction of troublesome environmental noise influences that often plague sound-based approaches. This departure aligns with the overarching objective of developing a robust and accurate cough detection algorithm, particularly for spirometry applications [1] [3].

Some studies concentrated on monitoring cough overtime and counting occurrences, others delved into assessing the severity of various respiratory disorders, including cystic fibrosis, cold, tuberculosis, and COPD.

Problem statement: Current spirometry standards mandate that correct maneuvers should be devoid of cough artifacts, especially in the initial seconds of forced exhalation. Achieving this standard becomes inherently challenging when patients manifest increased coughing tendencies during the examination. Our work aims to facilitate unsupervised or minimally supervised spirometry measurements, expanding accessibility in diverse settings, including home monitoring and general practitioner oversight.

The subsequent sections of this article are structured as follows: Methodology' details the data sources and tools employed.

'Results' unveils statistical and machine learning outcomes, alongside an exploration of variable significance. Finally, 'Discussion and Conclusion' presents closing remarks, delineates limitations and outlines potential directions for future research.

II. METHODOLOGY

The cough reflex serves as a critical physiological response to airway irritations, characterized by sudden, forceful expiratory efforts. Initiated by the stimulation of cough receptors within the respiratory tract, this reflex plays a pivotal role in clearing irritants and maintaining airway integrity. In the context of pulmonary health, the cough reflex assumes heightened significance, often reflecting respiratory irritation or underlying pathology. Conditions such as asthma, cystic fibrosis, and chronic obstructive pulmonary disease (COPD) frequently present with cough as an early symptom, necessitating precise monitoring and diagnostic approaches [1] [4].

Spirometry, a widely employed method for diagnosing and monitoring pulmonary functions, introduces a notable challenge in the context of patients with chronic pulmonary diseases. The procedure involves a forced expiratory maneuver, requiring considerable effort from the patient. However, individuals afflicted with conditions like COPD commonly experience heightened levels of breathlessness and cough, particularly during exacerbations. The need for accurate spirometry measurements is paramount for assessing disease progression and treatment efficacy [5].

This work elucidates a solution for automatic cough detection developed specifically comprising a portable spirometer, a computer and an internet cloud for data storage. This system aims to provide accurate and robust cough detection in various environments, including clinical and in-home settings. The algorithm is designed to analyze the airflow signal transmitted from the spirometer to the computer application during measurements, ensuring the reliability of clinically important parameters such as forced vital capacity (FVC), forced expiratory volume in the first second (FEV1), their ratio (FEV1/FVC), peak expiratory flow (PEF), among others [6].

To facilitate algorithm development, a large dataset of spirometry airflow recordings was utilized for training and validation. Preprocessing methods were applied to enhance the quality of the dataset, emphasizing the importance of meticulous data preparation in ensuring algorithmic accuracy. The analytical methods adapted for constructing the cough detection algorithm were driven by machine learning approaches, leveraging the capabilities of advanced mathematical techniques [7].

The machine learning algorithms employed were systematically presented, accompanied by an in-depth exploration of their results. The emphasis on high specificity in the algorithm's design was underscored, acknowledging the potential consequences of misclassification, such as the unjustified repetition of measurements or user discouragement.

As spirometry continues to be a cornerstone in respiratory diagnostics, the presented algorithm holds promise for enhancing the accuracy and efficiency of pulmonary function

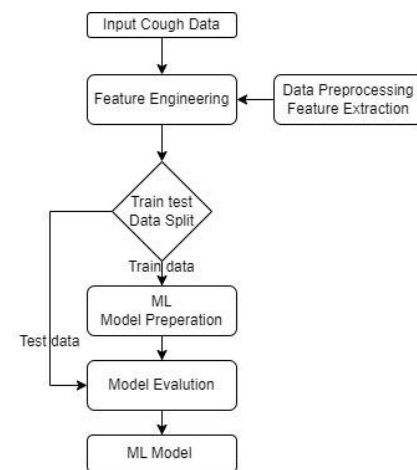


Fig.1. Process flow of the work

assessments in diverse clinical and non-clinical settings. The ongoing evolution of such algorithms signifies a critical step toward advancing respiratory health monitoring, offering potential benefits for both healthcare providers and patients alike [6] [7]. The figure 1 shows the process flow of the process.

A. Data accusation

Research data originated from the NHANES database, managed by the American National Center for Health Statistics [8]. This free resource includes raw spirometry curve data aligned with American Thoracic Society guidelines. Utilizing specific database subsets [9], diverse patients underwent initial and, if criteria were met, subsequent spirometry examinations with β_2 -adrenergic bronchodilator inhalation. The dataset encompasses raw signals from multiple individual spirometry curves, categorized into subsets, including those of acceptable quality, with large time to peak flow or non-repeatable peak flow, without plateau, or with cough. Experienced experts curated two classes: (a) ATS-acceptable curves and (b) other error curves versus cough curves. Figure 2 shows graphs of the Correct

ATS-Acceptable Maneuver Subset, Clear Cough Occurrence Subset and Less Manifested Cough Occurrence Subset respectively.

B. Data Preprocessing

Preprocessing of raw spirometry data is an indispensable precursor to feature extraction, involving systematic procedures for standardization and noise mitigation. This scientific approach encompasses several sequential operations:

- 1) Segmentation of Forced Exhale Signal: Initialization involves the segmentation of the raw curve to delineate the forced exhale signal, a critical maneuver in spirometry assessments. This segmentation is paramount to isolate the primary respiratory effort from ancillary components such as inhales preceding or succeeding the main maneuver.
- 2) Length Standardization: A subsequent step addresses the standardization of the forced exhale signal's length.

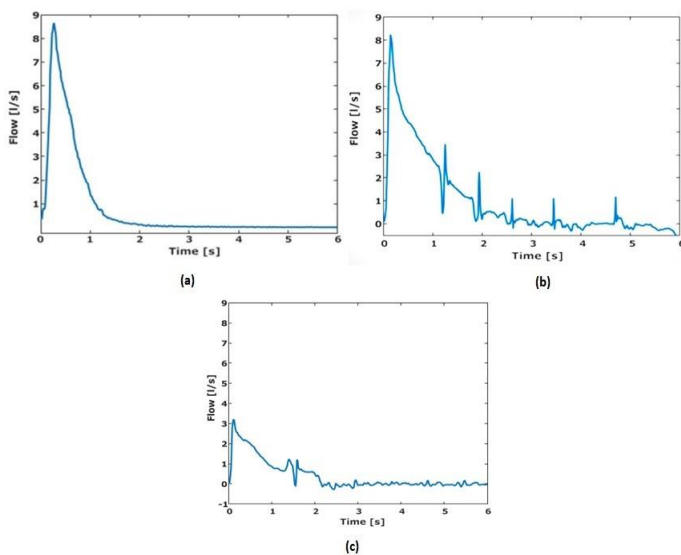


Fig. 2.(a): Correct ATS-Acceptable Maneuver Subset Graph, (b): Clear Cough Occurrence Subset Graph and (c): Less Manifested Cough Occurrence Subset Graph

To established spirometry standards, truncation is done, with particular emphasis on the first second, a temporal window crucial for assessing parameters such as Forced Expiratory Volume in the first second (FEV1).

- 3) Zeroing Negative Values: Negative values within the signal are systematically zeroed. While this operation is agnostic to the subsequent feature extraction process, it serves the purpose of diminishing the flow span of the data and obliterating residual inhale fragments that might persist from the initial segmentation.
- 4) Steady-Flow Detection Preprocessing: A specialized preprocessing step is dedicated to discerning the nature of the signal, specifically identifying whether it corresponds to a steady-flow characteristic, typical in scenarios such as Flow/Volume generator. This discernment relies on a meticulous evaluation of dissimilarity metrics between signal samples and the signal median.

The culmination of these preprocessing steps orchestrates the refinement and homogenization of raw spirometry curves, systematically priming them for subsequent analytical processes. This rigorous scientific approach ensures the standardization of signals and the amelioration of potential confounding factors, thereby fortifying the robustness and reliability of ensuing analyses, notably in the context of cough detection within spirometry evaluations [1] [2].

C. Feature Extraction

The study employs a rigorous feature extraction protocol to distill essential characteristics indicative of cough within individual spirometry curves. The algorithm input is defined by computationally efficient features, meticulously tailored to encapsulate distinctive cough-related attributes. The features are as follows:

- 1) Number of PEF-Related Local Maxima (Spikes): This feature quantifies local maxima (spikes) occurring post-peak expiratory flow (PEF). The chosen threshold distinguishes cough-relevant peaks from shorter, potentially inconsequential fluctuations.
- 2) Number of PEF-Related Local Maxima with Right-Slope Amplitude: The right-slope amplitude is defined as the amplitude between peak maximum and the point where the first derivative undergoes a sign change, indicating the initiation of an amplitude increase.
- 3) Number of Crossings at Various Percentages of PEF: Deriving four distinct features, this metric counts signal crossings with horizontal lines positioned at 15%, 25%, 50%, and 75% of PEF.

Correlation analysis is instrumental in gauging the relevance and interdependence of features. The preliminary feature selection process involved assembling a set of potentially useful features, with exclusion criteria based on high correlation with others [1] [2].

D. Machine Learning

The implementation of machine learning models took place within the Python-3 environment. A comprehensive suite of algorithms underwent systematic training and evaluation, seeking the optimal performer among logistic regression (LR), feed-forward artificial neural network (ANN), support vector machines (SVM), and random forest (RF). The F1 score served as the paramount metric for model optimization and performance comparison.

The ANN model was characterized by a sigmoid function as the transfer function in the neuron model. Employing k-fold cross-validation ($k = 10$) for validation. LR, integrated with k-fold cross-validation, utilized the mean square error as the loss function. SVM, employing a radial-type kernel, underwent parameter tuning through grid search. RF, grounded in Breiman's algorithm for classification and regression [1] [10]. The architectural configuration of the final artificial neural network is visually articulated. This methodological rigor in algorithmic selection, parameter tuning, and model validation ensures the robustness and reliability of the chosen machine learning approach for cough detection within spirometry datasets.

E. Proposed Machine Learning Approach

Employing clustering (unsupervised machine learning technique) and similar techniques aids in consolidating logging responses for classification, offering valuable insights into diverse data types and pattern identification. Analyzing the correlation between logging responses and blood flow enables the understanding of system behavior, issue detection, and overall system enhancement. Utilizing random-forest with K-Means clustering (K-MRF) for modeling and forecasting proves robust for predicting cough within spirometry datasets. K-Means clustering categorizes data, addressing the challenge of a weak logging response of cough relationship by grouping similar data points. The subsequent application of RF algorithms to each category enhances cough prediction accuracy based on distinct group features within spirometry datasets [10].

F. Performance Metrics

Table 1 outlines various metrics employed to assess the performance of classifiers in this study [10].

TABLE I
PERFORMANCE METRICS OF FOUR PREDICTIVE MODE

Metrics	Methods
Accuracy	$\frac{tP + tN}{tP + tN + fP + fN}$
Sensitivity	$\frac{tP}{tP + fN}$
F1-score	$2 \times \frac{\text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$
Specificity	$\frac{tN}{tN + fP}$
Classification Error	1 - Accuracy

III. RESULTS AND DISCUSSION

In the study, proposed algorithm under scrutiny exhibits noteworthy proficiency in mitigating the influence of cough events and artifacts on spirometric curves, thereby adeptly preserving pertinent information during the interstitial periods. Furthermore, its discerning capacity extends to the extraction of meaningful data from curves traditionally deemed aberrant, thereby overcoming historical impediments in data mining from such atypical sources. The comparative evaluation of diverse machine learning algorithms incorporates meticulous consideration of established metrics encompassing sensitivity, specificity, F1 score, and accuracy. These metrics collectively furnish a comprehensive and quantitative assessment of the algorithm's efficacy in spirometry data analysis. The algorithm's commendable performance, as elucidated through these metrics, underscores its potential significance in the realm of medical diagnostics, particularly in scenarios characterized by irregularities and perturbations in spirometric data.

In the study, cough events were classified from the NHANES database, and the performance of each ML model was assessed by conducting a series of studies.

The classification process, illustrated in Figure 1, involves four key

steps. Firstly, the feature engineering phase generates a fully featured dataset through data preprocessing and feature extraction. Subsequently, 10-Fold Cross-Validation partitions the dataset into ten unique subsets for systematic evaluation of classifiers. Machine learning algorithms are trained using a selection dataset to develop a cough event detection model, emphasizing minimal redundancy and maximum relevance in feature selection. Our study explores LR, ANN, SVM, RF, and proposed KMRF, evaluating each model's efficacy. Classifier performance relies on data quality, with inherent advantages and limitations. Performance metrics, including accuracy, sensitivity, specificity, F1-score, and classification error, assess ML efficiency, as depicted in Table 2. Figure 3 illustrates a comparison of classifier performance based on specificity and sensitivity.

TABLE II
PERFORMANCE COMPARISON OF DIFFERENT ML TECHNIQUES FOR IDENTIFYING COUGH EVENTS

ML Techniques	Performance Metrics		
	Accuracy	F1-Score	Classification Error
LR	0.913	0.812	0.087
ANN	0.912	0.916	0.088
SVM	0.902	0.832	0.098
RF	0.921	0.843	0.079
Proposed KMRF	0.96	0.92	0.04

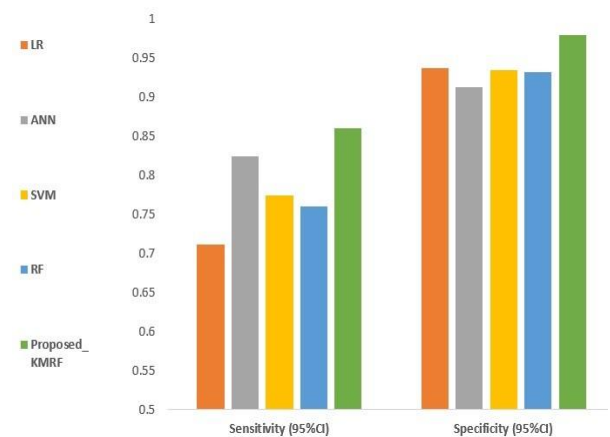


Fig. 3. Performance comparison of the classification model using specificity and sensitivity

This study focuses on detecting cough events in spirometry curves using airflow signals, addressing the complexity of distinguishing subtle cough manifestations from other disturbances. The research incorporates a diverse training dataset from the NHANES database.

The proposed algorithm demonstrates superiority over a previously described approach. Despite some misclassifications, particularly in cases of small cough disturbances or similar disturbances, the algorithm improves robustness for clinical and home monitoring in spirometry.

The study evaluates the performance of various classifiers in

detecting cough events within spirometry data. Classifier accuracies are reported, with consistent misclassifications observed across models, suggesting feature selection's primary influence. Specific misclassification patterns involve spirometry maneuvers lacking forced exhalation and coughs with spikes only at the end.

The Proposed KMRF algorithms excel in handling high-dimensional datasets and capturing intricate data patterns, enhancing cough event classification. The integration of KM clustering with RF algorithms enables precise feature extraction and improved classification. KM clustering reduces feature space dimensionality, facilitating algorithmic identification of relevant features and data relationships. From table 2 and figure 3 Proposed KMRF outperforms other models in metrics like accuracy, f1-score, specificity, and sensitivity, showcasing effective data clustering for enhanced precision and prediction accuracy. These algorithms mitigate overfitting, boost model generalization, and improve predictive accuracy, achieving a remarkable 96% accuracy in cough event prediction. The study emphasizes the importance of algorithm selection and the benefits of employing feature extraction techniques like KM for enhanced prediction accuracy.

Comparisons with previous work reveal higher specificity and overall F1 score, emphasizing the algorithm's effectiveness in minimizing false positives. A preference for specificity aligns with real-world application needs, preventing unnecessary repeated maneuvers.

Balancing the dataset adversely affects performance, attributed to the loss of information about event frequencies. Despite this, the model demonstrates reasonable generalization, indicating avoidance of over-training.

The study underscores the significance of feature choice, the balance between sensitivity and specificity, and considerations for real-world usability in developing an effective cough detection algorithm for spirometry. Despite these occasional challenges, the overall performance and capabilities of our algorithm represent a considerable step forward in the accurate detection of cough events during spirometry. As technology continues to evolve, our research contributes to the ongoing effort to refine algorithms, making them even more adept at discerning subtle respiratory patterns and further enhancing the utility of spirometry in both clinical and home settings.

The study acknowledges limitations and proposes future work to enhance the algorithm's performance and reliability. Additionally, the suggestion of modifying the algorithm to provide temporal information on cough occurrences, such as during the first second of exhalation, highlights potential avenues for improvement and broader applicability in spirometry quality assessment.

IV. CONCLUSION

Detecting cough events in spirometry curves through air flow signals poses challenges due to diverse cough manifestations, ranging from clear disruptions to subtle flow disturbances. Leveraging a comprehensive training dataset comprising the NHANES database, the study produced a robust classification

algorithm surpassing prior methodologies. The Proposed KMRF represents a notable advancement in spirometry monitoring for clinical and home settings with 96% accuracy. Despite occasional misclassifications in maneuvers featuring minor cough disturbances, our algorithm's features contribute to its accuracy. Comparative analysis with a competing algorithm demonstrates our superior overall F1 score of 92%, showcasing enhanced precision and recall. The balanced performance suggests our algorithm's reliability for precise cough event detection during spirometry, highlighting its potential for refined respiratory health monitoring.

The study acknowledges limitations and proposes future work to enhance the algorithm's performance and reliability.

ACKNOWLEDGEMENT

We wish to thank the Biomedical Engineering Department of National Institute of Technology, Raipur for assistance and encouragement in preparing this paper.

REFERENCES

1. Solin'ski, M., et al. Automatic cough detection based on air flow signals for portable spirometry system. *Informatics in medicine unlocked*, 18, (2020), 100313.
2. Smith J, Woodcock A. Cough and its importance in COPD. *Int J Chron Obstr Pulm Dis* 2006; 1(3):305–14.
3. Standardization of Spirometry 2019 Update. An official American thoracic society and European respiratory society technical statement. *Am J Respir Crit Care Med* October 15, 2019; 200(8). <https://doi.org/10.1164/rccm.201908-1590ST>.
4. Hofman A, Kupczyk M, Kuna P, et al. Application of the AioCare system in monitoring exacerbations of bronchial asthma. *Pol J Allergol Spec Issues* 2018; (1): A.8. in Polish.
5. Goel M, Saba E, Stiber M, et al. Spirocall: measuring lung function over a phone call. In: *Proceedings of the 2016 CHI Conference on human Factors in computing systems*. ACM; 2016. p. 5675–85.
6. Larson EC, et al. Accurate and privacy preserving cough sensing using a low-cost microphone. In: *Proceedings of the 13th international conference on ubiquitous computing*. ACM; 2011. p. 375–84.
7. Di Perna L, et al. An automated and unobtrusive system for cough detection. In: *Life sciences conference (LSC)*. IEEE; 2017. p. 190–3.
8. National Center for Health Statistics. National Health and nutrition examination Survey. website, <https://wwwn.cdc.gov/Nchs/Nhanes/> (last access 11 January 2019).
9. National Center for Health Statistics. NHANES data links, <https://wwwn.cdc.gov/Nchs/Nhanes/2007-2008/SPXRAW.htm>; 11 January 2019. https://wwwn.cdc.gov/Nchs/Nhanes/2009-2010/SPXRAW_F.htm. <https://wwwn.cdc.gov/Nchs/Nhanes/2011-2012/SPXRAWG.htm>.
10. Jain P, Gupta S. Multi-exposure Laser Speckle Contrast Imaging (MECI)-Based Prediction of Blood Flow Using Random Forest (RF) With K-Means (KM). *Cureus*. 2023 Jun 12; 15(6).