

AI model for predicting asthma prognosis in children



Elham Sagheb, MS,^a Chung-Il Wi, MD,^b Katherine S. King, MS,^c Bhavani Singh Agnikula Kshatriya, MS,^d Euijung Ryu, PhD,^c Hongfang Liu, PhD,^{a,e} Miguel A. Park, MD,^f Hee Yun Seol, MD,^g Shauna M. Overgaard, PhD,^d Deepak K. Sharma, PhD,^d Young J. Juhn, MD,^b and Sunghwan Sohn, PhD^a
Rochester, Minn; Houston, Tex; and Yangsan-si, Korea

Background: Childhood asthma often continues into adulthood, but some children experience remission. Utilizing electronic health records (EHRs) to predict asthma prognosis can aid health care providers and patients in developing effective prioritized care plans.

Objective: We aimed to develop artificial intelligence (AI) models using various clinical variables extracted from EHRs to predict childhood asthma prognosis (remission vs no remission) in different age groups.

Methods: We developed AI models utilizing patients' EHRs during the first 6, 9, or 12 years of their lives to predict their asthma prognosis status at ages 6 to 9, 9 to 12, or 12 to 15 years, respectively. We first developed the models based on a manually annotated birth cohort (n = 900). We then leveraged a larger birth cohort (n = 29,594) labeled automatically (with weak labels) by a previously validated natural language processing algorithm for asthma prognosis. Different models (logistic regression, random forest, and XGBoost [eXtreme Gradient Boosting]) were tested with diverse clinical variables from structured and unstructured EHRs.

Results: The best AI models of each age group produced a prediction performance with areas under the receiver operating characteristic curve ranging from 0.85 to 0.93. The prediction model at age 12 showed the highest performance. Most of the AI models with weak labels showed enhanced performance, and models using the top 10 variables performed similarly to those using all of the variables.

Conclusions: The AI models effectively predicted asthma prognosis for children by using EHRs with a relatively small number of variables. This approach demonstrates the potential to enhance prioritized care plans and patient education,

improving disease management and quality of life for asthmatic patients. (J Allergy Clin Immunol Global 2025;4:100429.)

Key words: Asthma, asthma prognosis, natural language processing, artificial intelligence, machine learning, electronic health records, dynamic variables

Asthma is often a lifelong disease, affecting 1 in 12 people in the United States,¹ and the most common chronic disease in children.² A long-term study following individuals up to age 50 years found evidence linking asthma severity during childhood to clinical and lung function outcomes in adulthood.³ Although several studies tried to predict who will have asthma at school age (eg, age 6-10 years) among preschool children with a history of asthma-like symptoms (eg, wheezing for the first 3 years of life),⁴ it is understudied to predict future asthma prognosis (eg, remission vs persistent asthma) among children with asthma, which is informative when clinicians determine continuing versus changing their asthma management based on short-term prediction. A review paper⁴ highlighted little success of several asthma prediction rules in clinical practice, which may in part be accounted for by the later development of asthma, noninformative rules (eg, low likelihood ratio of a positive test result), and/or heterogeneity of definitions of risk factors and/or outcomes (asthma). These challenges may be overcome by using electronic health records (EHRs) to reduce methodologic heterogeneity, including unstructured data (eg, free text), which is estimated to account for more than 80% of currently available health care data.⁵

As EHRs have been increasingly implemented in the US health care system, development of artificial intelligence (AI) models based on EHRs has been explored in many clinical areas. AI offers transformative resources, including diagnostics, predictive analytics, virtual health assistance, drug discovery, and personalized medicine, thereby assisting in treatment planning through analysis of patient EHR data. AI also aids research by finding patterns in data and through medical education providing training scenarios for health care professionals. In particular, use of AI for asthma has been studied largely for predicting the development of asthma at school age or identifying undiagnosed asthma,⁶⁻¹¹ poor asthma control status defined by asthma exacerbation or hospitalization for asthma,¹²⁻¹⁷ or other indicators (eg, response to treatment, asthma management, and phenotype).¹⁸⁻²² A recent review paper on AI approaches using natural language processing (NLP) to advance EHR-based clinical research showed the underuse of AI and NLP in the fields of asthma, allergy, and immunology.²³ Our research group developed several rule-based NLP algorithms for asthma ascertainment and extraction of relevant concepts from EHRs.²⁴⁻²⁸ We also developed an NLP algorithm to retrospectively determine a patient's asthma prognosis status by using

From ^athe Department of Artificial Intelligence and Informatics, ^bthe Department of Pediatric and Adolescent Medicine, ^cthe Department of Quantitative Health Sciences, ^dthe Center for Digital Health, and ^ethe Department of Allergy and Immunology, Mayo Clinic, Rochester; ^fthe UTHealth Houston; and ^gthe Department of Internal Medicine, Pusan National University School of Medicine, Pusan National University Yangsan Hospital, Yangsan-si.

Received for publication May 24, 2024; revised November 25, 2024; accepted for publication November 27, 2024.

Available online January 31, 2025.

Corresponding author: Sunghwan Sohn, PhD, Department of Artificial Intelligence and Informatics, Mayo Clinic, 200 First St SW, Rochester, MN 55905. E-mail: sohn.sunghwan@mayo.edu.

The CrossMark symbol notifies online readers when updates have been made to the article such as errata or minor corrections

2772-8293

© 2025 The Author(s). Published by Elsevier Inc. on behalf of the American Academy of Allergy, Asthma & Immunology. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

<https://doi.org/10.1016/j.jacig.2025.100429>

Abbreviations used

a-GPS:	Asthma guidance and predication system
AI:	Artificial intelligence
AUC:	Area under the (receiver operating characteristic) curve
EHR:	Electronic health record
NLP:	Natural language processing
PAC:	Predetermined asthma criterion
XGBoost:	eXtreme Gradient Boosting

EHR data related to asthma history and risk factors.²⁹ Knowing the likelihood of patients' asthma remission (vs persistence) will give opportunities for providers and patients to come up with a more patient-tailored and prioritized care plan. Yet, whether an EHR-based machine learning algorithm can predict who will have persistent asthma rather than achieving remission in the near future (eg, in the subsequent 3 years) has not been studied.

In this study, we developed machine learning models using patients' EHRs from the first 6, 9, or 12 years of life, incorporating both structured and unstructured data (eg, free text), to predict asthma prognosis (remission vs no remission [ie, persistent asthma]) at ages 6 to 9, 9 to 12, and 12 to 15 years, respectively. We used our previous NLP algorithms²⁴⁻²⁸ to extract relevant variables from unstructured EHR data for developing the prediction model. Our focus was not on long-term prediction, as the clinical prognosis of asthma is likely to be contingent on the time-dependent (eg, puberty) natural course of asthma, current asthma management, and patient compliance.

METHODS**Study design and setting**

This retrospective study used a designed birth cohort to develop AI models to predict asthma prognosis (remission vs no remission). We used a sample of 900 children aged 5 to 18 years from the 1997-2007 Mayo Clinic Birth Cohort (see the Study subjects section) who were included in an ongoing birth cohort study.²⁴ We utilized their EHRs during their first 18 years of life. We conducted a manual chart review of the entire EHRs of these subjects to determine asthma prognosis status using our previously reported predetermined criteria.²⁹

We defined remission as the absence of asthma events for at least 3 consecutive years after the last asthmatic event. Asthma events or active asthma included any of the following: (1) clinic visit for asthma with a physician diagnosis of asthma; (2) asthma symptoms, such as cough plus wheezing, prolonged exhaling, exercise-induced symptoms, chest tightness/pain, night cough, and dyspnea; and (3) current use of asthma medication. We focused on children who had asthma onset based on predetermined asthma criteria (PACs)^{24,27} during the first 6, 9, or 12 years of life (the observation window). We used their 6, 9, or 12 years of EHRs (sixth, ninth, and 12th birthdays as a decision date) to predict remission status within the next 3 years (prediction window). Fig 1 explains the analytic design for model development and shows our prediction model's observation window, decision date, and prediction window at age 6. Because of the small size of the sample of patients in our cohort who met the inclusion

criteria, we could not build a robust prediction model at age 15 years or further.

Study subjects

This study initially involved the EHRs of 900 patients who had been randomly selected from the Mayo Clinic Birth Cohort (1997-2006). This cohort (n = 900) was manually annotated, and their prognosis status was labeled by a clinical expert (H.Y.S.). A supervised machine learning algorithm often requires large data to perform well, but manually labeled data are costly. To alleviate manual labeling costs and enhance the AI prediction performance, we incrementally added computer-generated prognosis status ("weak labels"), as determined by our previously validated NLP prognosis algorithm.²⁹ The NLP prognosis algorithm was applied to the rest of the Mayo Clinic Birth Cohort (1997-2016) (n = 29,594 [ie, excluding the aforementioned 900 patients]) and labeled prognosis status (the weakly labeled cohort). Asthma status was determined by PACs by applying our validated NLP algorithm.^{24,27} Subjects were eligible for each time point only when a PACs index date (ie, first date of asthma onset based on the PACs) was within the observation window and the subject had EHR data at Mayo Clinic until the end of the prediction window (3 years after the observation window).

To evaluate the models, we utilized only the manually annotated cohort (criterion standard) for testing and applied the 5-fold cross-validation technique for training. Initially, both the training and testing data were sourced exclusively from the standard data set. Later, we incorporated augmented data into the training folds, and finally, we tried training the model by using only weakly labeled data. To ensure rigor in testing in all scenarios, we maintained the criterion standard data set as the only source for testing purposes.

Data preparation

We used 53 variables clinically relevant to asthma prognosis, including both structured (eg, demographics, laboratory test results) and unstructured data extracted from clinical notes by using our validated NLP algorithms (eg, asthma symptoms, history of viral infections).^{24,25,30,31} Table E1 (available in the Online Repository at www.jaci-global.org) contains the variables that were used in our models along with the EHR sources. We embedded temporality in the time-varying variables to reflect patients' health trajectories. For example, for "flu" and for the model that used 6 years of EHR, 6 variables were generated, namely, "flu at age 1 year", "flu at age 2 years", ..., and "flu at age 6 years," representing flu episodes over the whole observation period at yet-different time points. To handle missing demographic data, we categorized the missing information as "unknown." For temporal variables, we used the frequency during a particular year of a patient's life. For the static variables, we assigned either a numeric value or a category.

Model development and evaluation

We predicted patients' asthma prognosis status for the next 3 years, following an observation period at age 6, 9, or 12 years. To identify eligible patients for the development, the 3 prediction

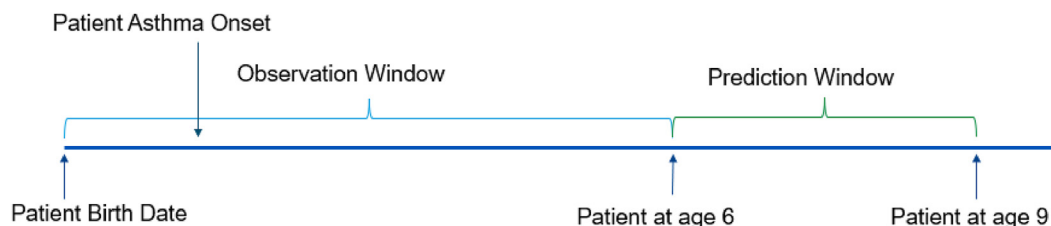


FIG 1. Timeline of asthma prognosis prediction at age 6 years.

models were independent of each other. We included patients with EHR data for the observation window plus for the next 3 years (ie, the prediction model at age 12 required 15 years' worth of EHR data for model development). We used random forest,³² logistic regression,³³ and XGBoost (eXtreme Gradient Boosting)³⁴ algorithms. These algorithms are known for their efficacy in predictive modeling. Random forest was chosen because of its robustness to overfitting, ability to handle large data sets with high dimensionality, and ability to capture complex nonlinear relationships in the data. Logistic regression, a classical statistical method, was selected for its interpretability, simplicity, and efficiency, particularly when dealing with binary classification tasks such as ours. Additionally, XGBoost was included for its state-of-the-art performance in gradient boosting, which often leads to superior predictive accuracy and scalability, especially in structured data settings. By leveraging these 3 algorithms, each with its unique strengths and characteristics, we aimed to provide a comprehensive and reliable prediction framework for asthma prognosis, ensuring robustness and accuracy across different modeling approaches.

A stratified 5-fold cross-validation was established to ensure consistency in training and test sets across all 3 algorithms for a fair performance comparison. We initially partitioned our criterion standard data set into 5 stratified folds based on the labels, preserving consistency throughout. Each iteration involved training models on 4 folds while reserving 1 fold for testing. Weakly labeled data were integrated with the training data from the 4 folds, augmenting the data set's diversity without compromising cross-validation integrity. We also tried training the models using only weak labels and evaluating them with criterion standard data. Notably, we maintained uniformity by using the same test set for all models, drawn from the criterion standard, to ensure a fair evaluation. This approach ensured that models were exposed to a diverse range of data while upholding cross-validation methodology standards and testing with the most accurate labels available. We used the area under the (receiver operating characteristic) curve (AUC)³⁵ for comparing the models' performance, thus providing a comprehensive evaluation of a model's discriminative ability across various thresholds and thereby offering insights into their performance across the entire range of possible classification thresholds.

We tested different scenarios to train the models using (1) the manually labeled data (considered the criterion standard), (2) the weakly labeled data, (3) augmented data (both manually and weakly labeled data), and (4) the top 10 most important variables. We used the Gini Importance scores³⁶ incorporated into random forest, which gave us the best performance while using criterion standard data and compared its performance with the parsimonious model by selecting the top 10 important variables.

To fine-tune our model parameters, we used the RandomizedSearchCV library from scikit-learn (2023).³⁷ In selecting RandomizedSearchCV, we leveraged its efficient approach to hyperparameter tuning. RandomizedSearchCV randomly samples from a predefined hyperparameter space, facilitating the exploration of a broader range of configurations. This randomness enhanced our ability to discover optimal hyperparameter combinations, which is particularly advantageous in high-dimensional feature spaces. By leveraging RandomizedSearchCV, we struck a balance between exploring diverse configurations and computational efficiency, enhancing the robustness and effectiveness of our models in handling complex data sets with numerous variables.

RESULTS

Study cohort characteristics

After applying our inclusion criteria, which required that subjects have their PACs index date (ie, first date of asthma onset based on PACs) fall within the observation window and also have EHR data at Mayo Clinic until at least the end of the prediction window (3 years after observation window), we included 357 subjects for a manually annotated cohort (considered the criterion standard) and 2532 subjects for a weakly labeled cohort data (computer-generated labels).

Table I summarizes the basic characteristics of the 2 experimental cohorts (criterion standard and criterion standard + weakly labeled data), showing similar sociodemographic and asthma-related clinical status. As Table I illustrates, the majority of the cohort was non-Hispanic and White, with a significant male prevalence. Socioeconomic status, represented through quartiles of the HOUSES (Housing-Based Socioeconomic Status) index,³⁸ shows a fairly even distribution across all levels, albeit a slight majority is in the lower SES quartiles.

The proportions of asthma remission between the 2 cohorts are similar, reflecting that the weakly labeled data effectively represents the true nature of the patient cohort (criterion standard). Notably, the remission percentages among patients aged 6 to 9, 9 to 12, and 12 to 15 years show close alignment between the 2 cohorts. Regarding remission of asthma, it has been observed that approximately half of the individuals in each age group achieved remission within the 3-year prediction window.

The AI models' performance

Table II shows the performance of asthma prognosis prediction of the 3 models (logistic regression, random forest, and XGBoost) trained on different sets of data (ie, criterion standard [manually annotated], weakly labeled [computer generated], both criterion standard and weakly labeled, and the 10 most important

TABLE I. Basic characteristics of the 2 cohorts

Characteristic	Criterion standard (manually annotated) cohort (n = 357)	Criterion standard + weakly labeled cohort (n = 2889)
Age at the last follow-up date (y), median (interquartile range)	13.1 (10.9-15.3)	16.4 (13.4-19.2)
Male patient, no. (%)	222 (62.18%)	1787 (61.86%)
Non-Hispanic patient, no. (%)	328 (91.88%)	2699 (93.42%)
White patient, no. (%)	290 (81.23%)	2400 (83.07%)
HOUSES index* in quartile at birth, no. (%)		
N-Miss	37 (10.36%)	259 (8.93%)
Q1 (lowest SES)	83 (23.25%)	781 (27.03%)
Q2	78 (21.85%)	642 (22.22%)
Q3	83 (23.25%)	641 (22.19%)
Q4 (highest SES)	76 (21.29%)	567 (19.63%)
Patient category, no. (%)		
With asthma remission between ages 6 and 9 y	167 (46.78%)	1260 (61.00%)
Without a remission between ages 6 and 9 y	190 (53.22%)	1629 (56.39%)
With asthma remission between ages 9 and 12 y	119 (33.33%)	1014 (35.10%)
Without a remission between ages 9 and 12 y	134 (37.54%)	1190 (41.19%)
With asthma remission between ages 12 and 15 y	58 (16.25%)	656 (22.71%)
Without a remission between ages 12 and 15	65 (18.21%)	700 (23.39%)

HOUSES, Housing-Based Socioeconomic Status; Q, quartile.

*An individual-level housing-based socioeconomic measure.³⁸

variables). The best AI models of each age group produced an AUC ranging from 0.85 to 0.93. The prediction model at age 12 years produced the highest performance (AUC = 0.93 by random forest with criterion standard + weak labels). Overall, the random forest algorithm achieved the best performance for all of the different age groups, albeit not significantly higher than that of the other models. The data augmentation using weak labels enhanced the prediction performance in most cases. The models with the top 10 most important variables performed similarly to the models using all variables. The models trained on weakly labeled data performed similarly to the models trained on criterion standard plus weakly labeled data. The most influential variables include the history of asthma prognosis within the observation window, frequency of rescue medication, age at asthma index date (ie, date of the first asthma diagnosis) of pre-determined asthma criteria, and age at index date of Asthma Predictive Index²⁵ (Table III). The prediction performance tended to be increased as we expanded the patient's medical history used in the models from 6-year EHR data to 9- and 12-year EHR data (ie, a more longitudinal EHR history).

DISCUSSION

To our knowledge, this is the first study to assess the performance of AI models using patients' EHRs, including both structured and unstructured data, to predict asthma prognosis (remission vs no remission) beyond early childhood. The best AI models of each age group produced AUCs ranging from 0.85 to 0.93.

Asthma stands as a pervasive and enduring health challenge.^{1,2,39} Although previous studies attempted to predict asthma onset in school-age children,^{4,6-11} little is known about whether forecasting the future prognosis of asthma among children diagnosed with the condition is feasible. This knowledge gap is particularly significant when considering the distinction between remission and persistence of asthma, which is a crucial factor for clinicians in deciding whether to maintain or modify asthma management strategies based on clinical judgment for short-term prognosis or predictions. Previous attempts at predicting asthma outcomes have encountered challenges owing to the heterogeneity of definitions of risk factors and/or outcomes for asthma.

Our research group has pioneered the development of rule-based NLP algorithms²⁴⁻²⁹ tailored for asthma ascertainment, concept extraction, and notably, determining asthma prognosis retrospectively by using EHR data, which could alleviate an issue of heterogeneity in asthma and related concept definitions. In addition, our group developed the NLP algorithm to identify clinicians' adherence to National Asthma Education and Prevention Program (NAEPP) guidelines including the documentation of asthma medication compliance.³⁰

This study represents an extension of our efforts using machine learning models including diverse EHR variables from distinct age periods (ages 6, 9, and 12 years) to predict asthma prognosis (remission vs persistent asthma) at subsequent ages (6-9, 9-12, and 12-15 years). The emphasis on short-term predictions aligns with the understanding that clinical prognosis is intricately linked to current asthma management practices and patient compliance, presenting an opportunity to optimize asthma care plans for improved patient outcomes.

In informatics research, many studies have used a relatively small set of manually annotated data, but lately, the literature has introduced and applied distant supervision⁴⁰ to augment the training data with weakly labeled (computer-generated labels) data to improve the model prediction performance. Instead of relying on manual annotations, distance supervision leverages existing knowledge bases or heuristics (ie, NLP prognosis algorithms) to assign labels to training data, allowing diversity in data and efficiency/scalability in model development. Our previously developed NLP-PACs algorithm^{24,27} and NLP prognosis algorithm²⁹ allow for identifying consistent asthma onset and weak labels for this purpose when large unlabeled EHR data sets are available. By utilizing distance supervision, most models showed improved performance. We also observed that the models trained by using only weakly labeled data performed similarly to the models trained by using criterion standard plus weakly labeled data (Table II). This indicates that our weak labels are of high quality, potentially thanks to the larger portion of weak labels compared with the criterion standard labels.

The analysis revealed a similarity between manual annotations and weakly labeled data regarding asthma remission rates across

TABLE II. The AI models' performance regarding asthma prognosis prediction

Observation window	Data	LR, median (\pm SD)	RF, median (\pm SD)	XGB, median (\pm SD)
Date of birth to age 6 y	Criterion standard	0.82 (\pm 0.05)	0.84 (\pm 0.04)	0.78 (\pm 0.06)
	Weak labels	0.84 (\pm 0.04)	0.83 (\pm 0.03)	0.82 (\pm 0.02)
	Criterion standard + weak labels	0.84 (\pm 0.04)	0.83 (\pm 0.03)	0.82 (\pm 0.02)
	Criterion standard + weak labels (top 10 variables)	0.84 (\pm 0.04)	0.85 (\pm 0.03)	0.84 (\pm 0.03)
Date of birth to age 9 y	Criterion standard	0.91 (\pm 0.04)	0.91 (\pm 0.05)	0.91 (\pm 0.04)
	Weak labels	0.89 (\pm 0.04)	0.88 (\pm 0.04)	0.85 (\pm 0.04)
	Criterion standard + weak labels	0.90 (\pm 0.04)	0.88 (\pm 0.03)	0.86 (\pm 0.04)
	Criterion standard + weak labels (top 10 variables)	0.88 (\pm 0.02)	0.88 (\pm 0.02)	0.89 (\pm 0.04)
Date of birth to age 12 y	Criterion standard	0.88 (\pm 0.07)	0.90 (\pm 0.10)	0.85 (\pm 0.08)
	Weak labels	0.87 (\pm 0.11)	0.93 (\pm 0.06)	0.92 (\pm 0.07)
	Criterion standard + weak labels	0.87 (\pm 0.12)	0.93 (\pm 0.10)	0.91 (\pm 0.10)
	Criterion standard + weak labels (top 10 variables)	0.90 (\pm 0.09)	0.90 (\pm 0.08)	0.90 (\pm 0.08)

AUC averaged over 5-fold \pm cross-validation (SD).

LR, Logistic regression; RF, random forest; XGB, XGBoost (eXtreme Gradient Boosting).

TABLE III. Top 10 variables for prognosis prediction for 6 years, 9 years, and 12 years of EHR, respectively

Variable name
Observation with 6 years' worth of EHRs
Asthma prognosis history until age 6 y
No. of visits with an indication of taking a rescue medication at age 6 y
No. of visits with an indication of checking for asthma at age 6 y
No. of visits with an indication of taking an asthma medication at age 6 y
No. of visits with an indication of taking a rescue medication at age 5 y
Age at PACs index date
Age at API index date
No. of visits with an indication of having a cough at age 6 y
No. of visits with an indication of checking for asthma at age 4 y
No. of visits with an indication of being on rescue medication at age 4 y
Observation with 9 years' worth of EHRs
Asthma prognosis history until age 9 y
No. of visits with an indication of taking an asthma medication at age 9 y
Age at PACs index date
No. of visits with an indication of taking an asthma medication at age 7 y
No. of visits with an indication of taking a rescue medication at age 9 y
No. of visits with an indication of taking a rescue medication at age 8 y
No. of visits with a diagnosis of asthma at age 9 y
No. of visits with an indication of checking for asthma at age 9 y
No. of visits with an indication of having nighttime symptoms at age 2 y
Age at API index date
Observation with 12 years of EHRs
No. of visits with an indication of taking a rescue medication at age 12 y
No. of visits with an indication of checking for asthma at age 11 y
No. of visits with an indication of being on rescue medication at age 10 y
No. of visits with an indication of being on rescue medication at age 11 y
No. of visits with an indication of being on asthma medication at age 12 y
Asthma prognosis history until age 12 y
Age at API index date
No. of visits with an indication of being on asthma medication at age 11 y
No. of visits with an indication of checking for asthma at age 9 y
No. of visits with a diagnosis of asthma at age 12 y

API, Asthma Predictive Index.

different age categories. By expanding the cohort, we also observed that nearly half of the asthmatic patients in our age groups achieved remission within the 3-year prediction window.

We used diverse clinical variables in EHRs relevant to asthma prognosis, extracted from both structured and unstructured data and embedded temporality in time-varying variables to better represent patients' health trajectories. Temporality accounting for time-varying variables, such as patients' histories relevant to asthma conditions and management, can provide valuable insights into their future asthma prognosis status. When we used the top 10 variables, most of the models performed similarly to when all 53 clinical variables were used. This suggests that we could develop a parsimonious model focusing on certain clinical variables to reasonably predict asthma prognosis. The accuracy of the model's predictions improved with the incorporation of more longitudinal EHR data. Specifically, the model using 12 years of EHRs outperformed the model using 9 years of EHRs, and the model using 9 years of data performed better than the model using 6 years of data. This may suggest that the model can be enhanced through learning capability via knowledge accumulation.

Finding a cohort with a complete EHR history in the health system and pulling, processing, and validating the data for a large number of variables is a time-consuming process and often difficult in real-world clinical settings, especially when we are extracting information from unstructured data. Over the years, changes in EHR platforms and data structures may have also led to the loss of certain variables in certain periods of time. Considering these facts, a parsimonious model would be a more appropriate choice for model implementation in real-world clinical settings, leading to easier understanding, interpretability, and generalization.

The optimization of an AI model for predicting asthma prognosis can be achieved through the collective utilization of various asthma-related data sources. A prime example of this approach is the Asthma-Guidance and Prediction System (a-GPS), a clinical decision support system powered by AI.⁴¹ The a-GPS approach significantly enhances asthma management by automatically aggregating and presenting clinicians with the most pertinent information from EHRs, including the risk prediction of asthma exacerbation, thereby improving asthma outcomes. Moreover, it alleviates the burden on clinicians involved in asthma care, as exemplified by a reduction in EHR review

time. This not only leads to improved efficiency but also has the potential to reduce health care costs.⁴¹ The efficacy of this novel technology-driven approach to asthma care is currently being evaluated through a randomized controlled trial. Additionally, our machine learning algorithm for predicting asthma prognosis is poised to further enhance the capabilities of a-GPS when integrated into the system. All of these algorithms, including NLP-PACs and prognosis, are part of the AI landscape to leverage the EHR data to improve diagnostic accuracy to better manage and care asthma patients.

Our AI models have been developed and tested in a single clinical setting and warrant validation in other clinical settings. The model used the past 6, 9, or 12 years of EHR data to predict asthma prognosis status within the following 3 years. We intend to use this framework to extend the study by expanding the observation window and modifying the prediction window by leveraging our approaches' strengths such as its foundation in a birth cohort study and the comprehensive capture of clinical services for asthma within medical records.

Conclusion

Asthma prognosis is a crucial outcome in pediatric asthma management, and its prediction supports better clinical decisions in asthma care. Our AI models have demonstrated the potential to provide discriminative insight into predicting asthma prognosis for the next 3 years in children, leveraging their EHR data. The models use the time of diagnosis and relevant events in their predictions, enabling health care providers to gain a more accurate understanding of patients' asthma prognosis and develop more effective treatment plans for them.

Our findings illustrate the efficacy of weakly labeled data in mirroring the precision of manual annotations, highlighting its potential as a reliable alternative for large-scale medical studies, in which manual annotation may be impractical owing to time or resource constraints. Our approach also determined that nearly half of children with asthma can achieve remission within 3 years of the prediction window.

DISCLOSURE STATEMENT

Supported by the National Institutes of Health (grants R01 HL126667, R21 AI142702, and R21 AG065639).

Disclosure of potential conflict of interest: Y. J. Juhn is principal investigator of the Respiratory Syncytial Virus Incidence Study, which was supported by GlaxoSmithKline, and principal investigator of the Artificial Intelligence Development Study for Asthma, which was supported by Genentech. The rest of the authors declare that they have relevant conflicts of interest.

We thank Mrs Kelly Okeson for her administrative assistance.

REFERENCES

- Center for Disease Control and Prevention. Asthma in the US (May 2011). Available at: <https://www.cdc.gov/vitalsigns/asthma/index.html#:~:text=About%201%20in%2012%20people,numbers%20are%20increasing%20every%20year>. Accessed June 12, 2023.
- Center for Disease Control and Prevention. Asthma in the US (May 2011). Available at: <https://www.cdc.gov/healthyschools/asthma/index.htm>. Accessed September 12, 2023.
- Tai A, Tran H, Roberts M, et al. Outcomes of childhood asthma to the age of 50 years. *J Allergy Clin Immunol* 2014;133:1572-8.e1573.
- Savenije OE, Kerkhof M, Koppelman GH, Postma DS. Predicting who will have asthma at school age among preschool children. *J Allergy Clin Immunol* 2012;130:325-31.
- Martin-Sanchez F, Verspoor K. Big data in medicine is driving big changes. *Yearb Med Inform* 2014;9:14-20.
- Kothalawala DM, Murray CS, Simpson A, Custovic A, Tapper WJ, Arshad SH, et al. Development of childhood asthma prediction models using machine learning approaches. *Clin Transl Allergy* 2021;11:e12076.
- Jeddi Z, Gryech I, Ghogho M, El Hammoumi M, Mahraoui C. Machine learning for predicting the risk for childhood asthma using prenatal, perinatal, postnatal and environmental factors. *Healthcare (Basel)* 2021;9:1464.
- Yu G, Li Z, Li S, Liu J, Sun M, Liu X, et al. The role of artificial intelligence in identifying asthma in pediatric inpatient setting. *Ann Transl Med* 2020;8:1367.
- Patel D, Hall GL, Broadhurst D, Smith A, Schultz A, Foong RE. Does machine learning have a role in the prediction of asthma in children? *Paediatr Respir Rev* 2022;41:51-60.
- Spathis D, Vlamos P. Diagnosing asthma and chronic obstructive pulmonary disease with machine learning. *Health Informatics J* 2019;25:811-27.
- Chatzimichail E, Paraskakis E, Sitzimi M, Rigas A. An intelligent system approach for asthma prediction in symptomatic preschool children. *Comput Math Methods Med* 2013;2013:240182.
- Lisspers K, Stållberg B, Larsson K, Janson C, Müller M, Łuczko M, et al. Developing a short-term prediction model for asthma exacerbations from Swedish primary care patients' data using machine learning - based on the ARCTIC study. *Respir Med* 2021;185:106483.
- Sills MR, Ozkaynak M, Jang H. Predicting hospitalization of pediatric asthma patients in emergency departments using machine learning. *Int J Med Inform* 2021;151:104468.
- Zein JG, Wu CP, Attaway AH, Zhang P, Nazha A. Novel machine learning can predict acute asthma exacerbation. *Chest* 2021;159:1747-57.
- Hussain Z, Shah SA, Mukherjee M, Sheikh A. Predicting the risk of asthma attacks in children, adolescents and adults: protocol for a machine learning algorithm derived from a primary care-based retrospective cohort. *BMJ Open* 2020;10:e036099.
- Patel SJ, Chamberlain DB, Chamberlain JM. A machine learning approach to predicting need for hospitalization for pediatric asthma exacerbation at the time of emergency department triage. *Acad Emerg Med* 2018;25:1463-70.
- Farion KJ, Wilk S, Michalowski W, O'Sullivan D, Sayyad-Shirabad J. Comparing predictions made by a prediction model, clinical score, and physicians: pediatric asthma exacerbations in the emergency department. *Appl Clin Inform* 2013;4:376-91.
- Lović M, Banić I, Lacić E, Pavlović K, Kern R, Turkalj M. Predicting treatment outcomes using explainable machine learning in children with asthma. *Children (Basel)* 2021;8:376.
- Qin Y, Wang J, Han Y, Lu L. Deep learning algorithms-based CT images in glucocorticoid therapy in asthma children with small airway obstruction. *J Healthc Eng* 2021;2021:5317403.
- Kercsmar CM, Sorkness CA, Calatroni A, Gergen PJ, Bloomberg GR, Gruchalla RS, et al. A computerized decision support tool to implement asthma guidelines for children and adolescents. *J Allergy Clin Immunol* 2019;143:1760-8.
- Bhardwaj P, Tyagi A, Tyagi S, Antão J, Deng Q. Machine learning model for classification of predominantly allergic and non-allergic asthma among preschool children with asthma hospitalization. *J Asthma* 2023;60:487-95.
- Ross MK, Yoon J, van der Schaar A, van der Schaar M. Discovering pediatric asthma phenotypes on the basis of response to controller medication using machine learning. *Ann Am Thorac Soc* 2018;15:49-58.
- Juhn Y, Liu H. Artificial intelligence approaches using natural language processing to advance EHR-based clinical research. *J Allergy Clin Immunol* 2020;145:463-9.
- Seol HY, Rolfes MC, Chung W, Sohn S, Ryu E, Park MA, et al. Expert artificial intelligence-based natural language processing characterises childhood asthma. *BMJ Open Respir Res* 2020;7:e000524.
- Kaur H, Sohn S, Wi CI, Ryu E, Park MA, Bachman K, et al. Automated chart review utilizing natural language processing algorithm for asthma predictive index. *BMC Pulm Med* 2018;18:34.
- Wi CI, Sohn S, Rolfes MC, Seabright A, Ryu E, Voge G, et al. Application of a natural language processing algorithm to asthma ascertainment. An automated chart review. *Am J Respir Crit Care Med* 2017;196:430-7.
- Wi CI, Sohn S, Ali M, Krusemark E, Ryu E, Liu H, et al. Natural language processing for asthma ascertainment in different practice settings. *J Allergy Clin Immunol Pract* 2018;6:126-31.
- Wu ST, Sohn S, Ravikumar KE, Waghlikar K, Jonnalagadda SR, Liu H, et al. Automated chart review for asthma cohort identification using natural language

- processing: an exploratory study. *Ann Allergy Asthma Immunol* 2013;111:364-9.
29. Sohn S, Wi CI, Wu ST, Liu H, Ryu E, Krusemark E, et al. Ascertainment of asthma prognosis using natural language processing from electronic medical records. *J Allergy Clin Immunol* 2018;141:2292-4.e3.
30. Sagheb E, Wi CI, Yoon J, Seol HY, Shrestha P, Ryu E, et al. Artificial intelligence assesses clinicians' adherence to asthma guidelines using electronic health records. *J Allergy Clin Immunol Pract* 2022;10:1047-56.e1.
31. Juhn Y, Moon S, Wi CI, Fu S, Weston J, Porcher J, et al. Automated chart review for identifying pre- and peri-natal risk factors associated with childhood asthma. *Am J Respir Crit Care Med* 2018;197:A2032.
32. Ho TK. Random decision forests. In: *Proceedings of the 3rd international conference on document analysis and recognition*. Vol 1. Montréal, Canada: ICDAR; 1995. pp. 278-82.
33. Cox DR. The regression analysis of binary sequences. *J Roy Stat Soc B* 1958;20:215-32.
34. Chen T, Guestrin C. XGBoost: a scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York, NY: ACM; 785-94.
35. Hanley JA, McNeil BJ. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 1982;143:29-36.
36. Louppe G, Wehenkel L, Suter A, Geurts P. 2013. Understanding variable importances in forests of randomized trees. In: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Vol 1 (NIPS'13)*. Red Hook, NY: Curran Associates; 2013. pp. 431-9.
37. scikit-learn. RandomizedSearchCV. 2023. Available at: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.RandomizedSearchCV.html. Accessed September 23, 2023.
38. Juhn YJ, Beebe TJ, Finnie DM, Sloan J, Wheeler PH, Yawn B, et al. Development and initial testing of a new socioeconomic status measure based on housing data. *J Urban Health* 2011;88:931-44.
39. Zhong W, Finnie DM, Shah ND, Wagie AE, St. Sauver JL, Jacobson DJ, et al. Effect of multiple chronic diseases on health care expenditures in childhood. *J Prim Care Community Health* 2015;6:2-9.
40. In distant supervision, we make use of an already existing database, such as Freebase or a domain-specific database, to collect examples for the relation we want to extract. We then use these examples to automatically generate our training data. Available at http://deepdive.stanford.edu/distant_supervision. Accessed June 12, 2023.
41. Seol HY, Shrestha P, Muth JF, Wi CI, Sohn S, Ryu E, et al. Artificial intelligence-assisted clinical decision support for childhood asthma management: a randomized clinical trial. *PLoS One* 2021;16:e0255261.