Does Cohort Selection Affect Machine Learning from Clinical Data?

Atefehsadat Haghighathoseini¹, Janusz Wojtusiak¹, Hua Min¹, Timothy Leslie¹, Cara Frankenfeld², Nirup M Menon^{1,3} (On behalf of the N3C Consortium¹)

¹George Mason University, Fairfax, VA, USA; ²MaineHealth Institute for Research, Scarborough, ME, USA ³On behalf of the N3C Consortium

Abstract

This study investigates cohort selection and its effects on the quality of machine learning (ML) models trained on clinical data, focusing on measurements taken within the first 48 hours of hospital admission. It discusses the potential repercussions of making arbitrary decisions during data processing prior to applying ML methods. Experiments are performed within the framework of the National COVID Cohort Collaborative (N3C) dataset. The research aims to unravel biases and assess the fairness of machine learning models used to predict outcomes for hospitalized patients. Detailed discussions cover the data, decision-making processes, and the resulting impact on model predictions regarding patient outcomes. An experiment is conducted in which four arbitrary decisions are made, resulting in 16 distinct datasets characterized by varying sizes and properties. The findings demonstrate significant differences in the obtained datasets and indicate a high potential for bias based on inclusion or exclusion decisions. The results also confirm significant differences in the performance of models constructed on different cohorts, especially when cross-compared between ones based on different inclusion criteria. The study specifically chose to analyze gender, race, and ethnicity as these social determinants of health played a significant role in COVID-19 outcomes.

Keywords— Prediction, Selection Bias, Data Processing, Machine Learning, National COVID Cohort Collaborative (N3C)

Introduction

In health data processing, decisions regarding cohort and variable selection often lack standardization or reproducibility, relying on the intuition of researchers and existing literature. This common practice creates a significant gap in the assessment of the impact of these choices on data analysis models. Healthcare stakeholders must be vigilant in conducting unbiased data analyses and recognize the potential limitations stemming from arbitrary decisions. Overlooking seemingly minor choices during data processing can result in selection bias and alterations to the dataset. Strict inclusion/exclusion criteria may inadvertently exclude important subpopulations, affecting generalizability. Conversely, lenient inclusion/exclusion criteria may increase the likelihood of introducing noise into the data. A thorough evaluation of these decisions is crucial during data processing.

Selection bias arises when the chosen cohort inadequately represents the broader population under study, resulting in diminished external validity or generalizability. This occurs when individuals, groups, or data are selected for analysis without proper randomization, leading to a sample that may not accurately reflect the population ¹. James Heckman's "Varieties of Selection Bias" in econometrics explores the consequences of selection bias on estimating the effects of specific variables ². Furthermore, selection bias can inflate effect sizes and yield inaccurate results, undermining the reliability of statistical tests and the validity of research findings ³.

Electronic medical records (EMR) are valuable data sources for assessing health outcomes. However, using EMRs and other health data in clinical settings introduces several biases that can significantly impact data analysis results, subsequent healthcare policies, and patient outcomes. These biases may arise from arbitrary decisions in cohort criteria and variable selection, patient selection, and data completeness ⁷. Significant efforts have focused on examining and mitigating selection bias in conventional statistics and experimental research. For example, a study on participants' arbitrary choices between two equally favored options demonstrated that such choices could significantly affect their preferences ⁴. This highlights the substantial impact that arbitrary decisions can have on study outcomes. The consequences of bias extend across various domains, including clinical practice, research, and decision-making. In healthcare, implicit bias has been recognized as a significant contributor to disparities, ultimately influencing patient

¹ Authorship was determined using ICMJE recommendations. The analyses described in this publication were conducted with data or tools accessed through the NCATS N3C Data Enclave covid.cd2h.org/enclave and supported by CD2H - The National COVID Cohort Collaborative (N3C) IDeA CTR Collaboration 3U24TR002306-04S2 NCATS U24 TR002306. This research was possible because of the patients whose information is included within the data from participating organizations (covid.cd2h.org/dtas) and the organizations and scientists (covid.cd2h.org/duas) who have contributed to the on-going development of this community resource. See Hanendel et al 2021 for details.

outcomes ¹³. In outcome prediction, selection bias can introduce systematic errors into predictions generated by machine learning (ML) models. For example, models trained on nonrepresentative data may yield inaccurate predictions when applied to a broader population ⁷. Suri ⁸ and MacNamee ⁹ highlight the critical role of bias in ML systems for outcome prediction, especially in cardiovascular disease risk assessment and anticoagulant drug therapy. Researchers have utilized various data sources, including electronic health records (EHRs), claims data, and registries, to analyze and evaluate biases. This comprehensive approach enhances the understanding of factors affecting the quality and reliability of clinical data. Recognizing and addressing these biases is essential for improving healthcare policies and patient outcomes.

This study delves into the influence of cohort selection on the efficacy of ML models trained using clinical data. It explores how arbitrary decisions made during data processing can affect the quality and fairness of ML models within the National COVID Cohort Collaborative (N3C) dataset. By examining biases and evaluating the fairness of ML models in predicting outcomes for hospitalized patients, particularly in the context of N3C, this research sheds light on the importance of thoughtful cohort selection. The study specifically chose to analyze gender, race, and ethnicity as these social determinants of health played a significant role in COVID-19 outcomes. The study conducted an experiment involving four arbitrary decisions, resulting in 16 distinct datasets with varying sizes and characteristics. Detailed discussions cover the data, decision-making processes, and the resulting impact on model predictions regarding patient outcomes.

Design of the Experiment

This study explores the impact of data preprocessing decisions, particularly cohort selection, on the quality of machine learning models. Traditionally, machine learning projects for predicting patient outcomes follow a standard pipeline (Figure 1 top). After cohort selection, input and outcome variables are generated, and the data is split into training and testing sets. The training set is used to develop and fine-tune models to improve the AUC, while the testing set is reserved for the final evaluation of the models. However, this approach often involves making data processing decisions, such as cohort selection, before splitting the data. This practice limits the ability to assess how well the model generalizes beyond the specific training and testing sets.

This study generates multiple versions of the final dataset to investigate model generalizability, analyzing each version from the point of data splitting into training and testing sets. A matched training/testing split was used to ensure that the same patient did not appear in the training and testing sets across cohorts. Models are trained and tested uniformly, and results are compared to highlight cross-testing performance across cohorts (Figure 1 bottom).



Figure 1. Standard (top) and proposed (bottom) design of experiments.

Data Source:

The National COVID Cohort Collaborative (N3C) is an indispensable tool in the study of COVID-19, offering a rich dataset for detailed analysis through its specialized N3C Data Enclave ¹⁰. This initiative capitalizes on integrating diverse health records from multiple medical institutions, establishing a unified structure that aids researchers in accessing reliable and comprehensive information. More specifically, the presented study used a Limited Data Set ("LDS"), providing access to individual patient-level data with removed patient identifiers.

The focus of this study revolves around predicting outcomes for individuals admitted to hospitals with COVID-19. The primary examination unit is the patient's hospitalization, and the goal is to forecast the associated outcome. The N3C data, structured across various tables within a relational database, adhere to the Observational Medical Outcomes Partnership (OMOP) Common Data Model (CDM) standard. The necessary transformation involves converting this data into a format where each hospitalization corresponds to a single record (flat table, sometimes referred to as an analytic file). This format is standard for applying both ML methods and conventional statistical tools.

Data Preprocessing Decisions

The initial and critical step in data analysis for constructing machine learning models is cohort creation. This study employed a systematic approach to build cohorts from approximately 21 million N3C patients. As shown in Figure 2, this approach involved identifying COVID-19 cases, recognizing inpatient hospitalization records, determining COVID-19-related hospitalizations, and excluding records missing specific admission time data. While other preprocessing decisions were made, they are not explored in detail in this work for simplicity. For example, the data were filtered to include only records with a first COVID-19 diagnosis between August 1, 2020, and December 31, 2021. This timeframe was selected as the study target period, with data prior to August 2020 excluded due to unreliable early pandemic reporting and data post-2021 removed to maintain a fixed cohort timeframe. While reasonable, the choice of these specific dates does not have a precise clinical justification. Hospitalization records were merged with other data using a unique identifier and analyzed for features such as gender, date of birth, race, ethnicity, and age at death. Non-adult patients were excluded, with the study focusing on individuals aged 18 and older, though there is no specific clinical rationale for choosing 18 over other ages like 21. The application of these four preprocessing decisions results in 16 potential datasets. The optimal choice among these datasets for final analysis remains uncertain, as different options may be argued for at each decision level. Detailed descriptions of these decisions are available in $[1^1]$. Although seemingly arbitrary, these choices lead to 16 distinct datasets characterized by varying sizes and properties. This study analyzes the selection of four cohorts (A, B, C, and D) but holds relevance for all potential datasets. These decisions, along with corresponding annotations, are depicted in Figure 2.

The subsequent process involves extracting all pertinent measurements from the measurement table for patients within specific cohorts of the N3C database. These measurements are then meticulously categorized based on their concept IDs, which are distinctly selected from the concept set member's table. The concept IDs encompass a wide range of medical measurements, including but not limited to Respiratory Rate, GFR (Glomerular Filtration Rate), Cardiac Troponin T, Sodium, ABG indices, Prothrombin Time, Alanine Transaminase, Calcium, Systolic BP (Blood Pressure), Heart Rate, CD3 CD8 T cells, Albumin, Erythrocyte Sedimentation Rate, Venous Lactate, BMI (Body Mass Index), CBC with PLT (Complete Blood Count with Platelets), CD3 CD4 T cells, Cardiac Troponin I, Interleukin 6, FiO2, Blood Urea Nitrogen, Temperature, Bilirubin, D-Dimer, Alkaline Phosphatase, fibrinogen, Creatinine, Chloride,



Figure 2. Partial inclusion/exclusion tree outlining extreme choices for the four decisions. PC - patient counts.

Glucose, C-Reactive Protein, Urine Output, NT-ProBNP, Diastolic BP, Weight, Ferritin, and Interleukin 10. This comprehensive extraction and categorization are focused on the critical initial 48 hours of hospital admission, targeting the acute phase of the patient's stay. This approach is designed to facilitate clinical decision-making and support research efforts by providing a structured dataset of key patient measurements during this pivotal period of hospitalization.

Cohort Characteristics

There is a clear difference in the size of the cohorts, ranging from 124,985 to 328,726 patients, indicating an almost three-fold increase in population size. The most significant impact on the cohort size results from the last decision, which is purely data-driven (see Figure 2) and accounts for approximately 60% of the reduction in cohorts. The inclusive cohorts, determined by the timestamp decision (Cohorts B & D), show higher proportions of males, Asians, and Hispanics/Latinos or individuals with unknown ethnicity than their counterparts. Descriptive statistics for gender, race, and ethnicity across the four cohorts reveal notable disparities. For instance, Cohort A has the lowest number of males (64,953), while Cohort D has the highest number of females (161,188). Examining the racial distribution across cohorts, we observe Cohort D generally has the highest counts across all racial categories, indicating a larger overall patient population in that cohort. Cohort A consistently has the lowest counts across most racial categories. The distribution of races significantly differs among the cohorts, as indicated by the low p-values from the Chi-squared test. This suggests that the racial composition varies considerably between the cohorts. Cohort D also has the highest

		Gender* Einnicuy*							Kace*				
Cohort ID	Patients count	Male	Female	Unknown	Hispanic or Latino	Not Hispanic or Latino	Unknown	White	Black	Asian	Native Hawaiian	Multiple races	Unknown
		64,953	60,020	<20	22,425	97,227	5,333	73,539	22,081	4,729	311	130	24,195
	104.005	(52%)	(48%)	(0%)	(18%)	(78%)	(4%)	(59%)	(18%)	(4%)	(0%)	(0%)	(19%)
А	124,985	10,544**	9,709	<20	3,751	15,645	858	11,628	3,651	740	43	26	4,166
		(52%)	(48%)	(0%)	(19%)	(77%)	(4%)	(57%)	(18%)	(4%)	(0%)	(0%)	(21%)
	202 419	155,354	147,800	264	38,549	244,810	20,059	191,356	55,018	8,714	1,452	3,276	43,602
р		(51%)	(49%)	(0%)	(13%)	(81%)	(7%)	(63%)	(18%)	(3%)	(0%)	(1%)	(14%)
в	505,418	23,374	21,997	<20	5,811	37,737	1,827	28,479	8,888	1,120	68	612	6,208
		(52%)	(48%)	(0%)	(13%)	(83%)	(4%)	(63%)	(20%)	(2%)	(0%)	(1%)	(14%)
С	133,992	69,114	64,866	<20	24,425	103,891	5,676	78,357	23,324	5,584	329	137	26,261
		(52%)	(48%)	(0%)	(18%)	(78%)	(4%)	(58%)	(17%)	(4%)	(0%)	(0%)	(20%)
		11,082	10,439	<20	4,018	16,604	900	12,362	3,851	782	44	27	4,456
		(52%)	(48%)	(0%)	(19%)	(77%)	(4%)	(57%)	(18%)	(4%)	(0%)	(0%)	(21%)
	328,726	167,259	161,188	279	42,364	264,819	21,543	207,695	58,792	9,869	1,508	3,603	47,259
D		(51%)	(49%)	(0%)	(13%)	(81%)	(7%)	(63%)	(18%)	(3%)	(0%)	(1%)	(14%)
D		25,017	23,864	<20	6,241	40,676	1,968	30,924	9,362	1,211	69	692	6,627
		(51%)	(49%)	(0%)	(13%)	(83%)	(4%)	(63%)	(19%)	(2%)	(0%)	(1%)	(14%)
BC	25.014	12,762	12,249	<20	2,036	22,013	965	16,798	5,227	379	24	586	2,000
D-C	25,014	(51%)	(49%)	(0%)	(8%)	(88%)	(4%)	(67%)	(21%)	(2%)	(0%)	(2%)	(8%)
ВA	24 084	12,744	12,237	<20	2,032	21,987	965	16,786	5,222	373	24	586	1993
D-A	24,904	(51%)	(49%)	(0%)	(8%)	(88%)	(4%)	(67%)	(21%)	(1%)	(0%)	(2%)	(8%)
DА	28 207	14,291	14,003	<20	2,445	24,751	1,101	19,100	5,650	460	25	666	2396
D-A	20,297	(51%)	(49%)	(0%)	(9%)	(87%)	(4%)	(68%)	(20%)	(2%)	(0%)	(2%)	(8%)
DC	27 220	13,846	13,371	<20	2,195	23,961	1,064	18,494	5,494	421	24	665	2,122
D-C	27.220	(51%)	(49%)	(0%)	(8%)	(88%)	(4%)	(68%)	(20%)	(2%)	(0%)	(2%)	(8%)
DB	3 020	1,385	1,635	<20	390	2,502	128	2,113	374	82	<20	65	385
р-р	5,020	(46%)	(54%)	(0%)	(13%)	(83%)	(4%)	(70%)	(12%)	(3%)	(0%)	(2%)	(13%)
* All	P-values	< 0.001 us	ing Pearso	on's Chi	-squared	test							
** Ea	ch cell sha	aded in gre	y, contair	inform	nation fro	om the tes	t dataset						

Table 1. Cohorts and combinations of specialized test datasets, along with their gender, race, and ethnicity subgroups

 $Gender^*$ $Fthnicity^*$ $Pace^*$

number of patients across all categories of ethnicity, which is consistent with it being the largest cohort. Cohort A has the lowest counts of Hispanic/Latino and NOT Hispanic/Latino patients and the lowest count of patients with unknown ethnicity.

Understanding and addressing potential selection bias based on gender, race, and ethnicity is crucial for developing ML-based models within the N3C. Adjusting for these demographic factors using appropriate techniques is essential to ensure that predictive models are not biased and can be effectively generalized across diverse cohorts. The analysis

indicated that the cohorts are statistically different (p<0.001) in the composition of age, gender, race, ethnicity, and length of stay, as tested with the Mann-Whitney *U* Test and Chi-squared test. The different parts of the cohorts are of most importance, denoted as B-C, D-A, etc. These patients are present in larger cohorts but not in smaller ones, and ones for which the generalizability of models needs to be carefully assessed. Table 1 presents the numbers of patients in these sets split by gender, race, and ethnicity.

For each cohort, the results focus on the measurements taken within the first 48 hours of hospital admission, categorized by the survival status, which indicates whether the patient survived or died. In Cohort A, 86.95% of patients survived, while 13.05% did not. Cohort B shows a slightly different outcome, with 84.5% of patients surviving and 15.5% not surviving. Cohort C's results are like Cohort A's, with 87.12% of patients alive and 12.88% having expired. Lastly, Cohort D's data reveals that 84.77% of patients survived, and 15.23% did not. These figures provide insights into the survival rates of patients within each cohort during the acute phase of their hospital stay.

Machine Learning Models

Model Construction Methods

Standard supervised machine models were trained on the training set, with the predicted outcome being survival at the end of hospitalization. Random Forests (RF)¹², Gradient Boosting (GB)¹², and L2-regularized Logistic Regression (LR)¹³ models were trained. RF, GB, and LR are all supervised learning algorithms for classification and regression tasks. RF and GB are ensemble methods that combine multiple models to improve performance, while LR uses regularization to prevent overfitting. RF is known for its robustness and ability to handle diverse datasets, GB for optimizing arbitrary loss functions, and LR for its regularization capabilities to handle multicollinearity and non-independent features.

This experiment used a holdout training set to construct three ML models - LR, RF, and GB- using their default parameters. The primary goal was to compare the performance of these models across different cohorts rather than training the best possible models. To account for potential variations in the data, the dataset was divided into four distinct cohorts: A, B, C, and D. GridSearchCV was then employed to optimize each cohort's hyperparameters of the LR, RF, and GB models separately. This approach allowed us to fine-tune the hyperparameters specifically for each cohort, potentially improving the overall performance and accuracy of the predictions while maintaining a consistent baseline for comparison across the four cohorts.

For Cohort A, the Logistic Regression model was configured with a maximum iteration of 1000 and balanced class weights. The Random Forest Classifier utilized 1000 estimators, a maximum depth of 12, and balanced class weights. The Gradient Boosting Classifier had 150 estimators, a learning rate of 0.05, and a maximum depth of 5. In Cohort B, the Logistic Regression model maintained the same settings as in Cohort A. The Random Forest Classifier retained the same configuration as well. However, the Gradient Boosting Classifier was slightly modified, with 160 estimators and a maximum depth of 6, while keeping the learning rate at 0.05. For Cohort C, the Logistic Regression and Random Forest Classifier models were identical to those used in Cohort A. The Gradient Boosting Classifier had 150 estimators, a learning rate of 0.05, and a maximum depth of 7, varying from the previous cohorts. Lastly, in Cohort D, the Logistic Regression and Random Forest Classifier models remained unchanged from Cohorts A, B, and C. The Gradient Boosting Classifier had 160 estimators, a learning rate of 0.07 (slightly higher than the other cohorts), and a maximum depth of 6. This approach allowed us to fine-tune the hyperparameters of the models for each cohort separately, potentially improving the overall performance and accuracy of the predictions across the four cohorts.

Model Evaluation Methods

The presented work followed the general evaluation framework described by Wojtusiak^{14 15}. Evaluation metrics are crucial for assessing the performance of ML models. They provide insights into how well a model's predictions match the data. This study utilizes several metrics, including AUC (Area Under the Receiver-Operator Curve), Precision, Recall, and F1 Score, as well as measures such as model calibration. Assessing machine learning models across subpopulations delineated by gender, race, and ethnicity is imperative to test for fairness and reliability. Subpopulations include categories such as male and female for gender and white, black, Asian, Native Hawaiian, and multiple races for racial classification. Ethnicity is distinguished by Hispanic and non-Hispanic categories. Proper evaluation across these diverse groups is critical to mitigating biases and ensuring equitable outcomes. To ensure algorithmic fairness, we assess the performance of machine learning models to predict the outcome of survival at the end of hospitalization by evaluating the AUC. This evaluation allows us to audit fairness across the entire dataset and within subgroups.

Model Comparison

For each model— LR, RF, and GB—there is a standard train-test evaluation where the model is trained on a cohort's data (e.g., cohort A) and then tested on that same cohort's test data (A). This provides a baseline understanding of the model's performance on data similar to what it was trained on. However, the evaluation goes further by testing each model on data from other cohorts. For example, a model trained on cohort A's data is also tested on cohorts B, C, and D's test data.

This is done by creating specialized test datasets for each combination, ensuring that individuals from the training set are not included in the test sets to avoid data leakage and to ensure a fair assessment of the model's predictive power on entirely new data. This cross-comparison is repeated for all possible combinations, such as training on cohort B, testing on A, C, and D, and so on for the other cohorts.

Moreover, we have specialized test datasets. These specialized datasets are created by combining test data from different cohorts while ensuring no individual is included in more than one test dataset, leading to combinations such as B-A, B-C, D-A, D-B, and D-C. For the D-A test dataset, as an example, all individuals from cohort D's test data are included except for those who also appear in cohort A's test data.

Comparison of Feature Importance

In ML, identifying which input variables most significantly predict the outcome is crucial, and this is where feature importance becomes essential. Methods for determining feature importance include statistical validation, Shapley value analysis, and notably, the Random Forest Feature Importance (RFFI) method ¹⁶.

RFFI, embedded within the Random Forest algorithm, assesses feature significance by measuring each feature's contribution to reducing the model's impurity, with common metrics such as Gini impurity or mean squared error being used for this impurity-based importance ¹⁷. In this study, we have implemented the RFFI method to evaluate the significance of features across cohorts A, B, C, and D.

In this method, we have observed that 'age' occupies a larger portion of the bar plot. Therefore, we have decided to exclude it from our calculation. Consequently, we implemented the RFFI without taking 'age' into consideration. Each value represents the calculated importance score for a specific feature within a cohort, with higher scores indicating a greater impact on the predictive model.

Results

Accuracy comparison

The accuracy comparison is centered around the AUC metric. This metric is used to evaluate the performance of classification models, providing an aggregate measure of their ability to classify positive and negative cases across different thresholds 18 correctly. Based on the AUC scores, which measure the model's ability to distinguish between classes, we determine the best model by comparing the AUC values across three models. At first, this comparison is conducted across all cohorts (A, B, C, and D) using the entire dataset (see Table 2).

Table 2 illustrates the consistency and variation in model performance across different cohort combinations, showing that models tend to perform similarly within the same cohort combinations (e.g., Train = A and Test = A vs. Train = C and Test = C).

In contrast, performance can vary when models are trained and tested on different cohorts. This evaluation helps assess the models' robustness across different data distributions.

The table effectively assesses how well the models generalize to unseen data from different cohorts by comparing AUC values across various cohort combinations. These results indicate that the similarity of the training and testing

						Test				
		А	В	С	D	B-A	B-C	D-A	D-B	D-C
	Α	0.75, 0.75, 0.76*	0.71, 0.71, 0.72	0.75, 0.75, 0.76	0.71, 0.71, 0.73	0.69, 0.68, 0.70		0.70, 0.69, 0.71		
Tr	В	0.74, 0.74, 0.74	0.73, 0.75, 0.75	0.74, 0.74, 0.74	0.74, 0.75, 0.75	0.73, 0.74, 0.74	0.73, 0.74, 0.74		0.79, 0.78, 0.79	
ain	С	0.75, 0.75, 0.76	0.71, 0.71, 0.72	0.75, 0.76, 0.76	0.71, 0.71, 0.73		0.69, 0.68, 0.70			0.69, 0.68, 0.70
_	D	0.74, 0.74, 0.73	0.73, 0.75, 0.75	0.74, 0.74, 0.73	0.74, 0.75, 0.75			0.73, 0.75, 0.75	0.79, 0.78, 0.79	0.73, 0.74, 0.75
*Th	e seau	ence of numbers i	in each cell, listed	from left to right.	represents the AU	C values for LR.	RF, and GB mode	ls.		

Table 2. The AUC values for three models were calculated across all cohorts using the entire dataset

data affects the model's ability to generalize, with the best performance seen when both are closely aligned. This underscores the importance of training on diverse and representative data for better model generalization.

		Test (Gender)											
		А	В	С	D	B-A	B-C	D-A	D-B	D-C			
Trair	Α	0.74, 0.75, 0.76*	0.71, 0.70, 0.72	0.74, 0.75, 0.75	0.71, 0.70, 0.72	0.69, 0.67, 0.70		0.69, 0.67, 0.70					
		0.75, 0.76, 0.76	0.71, 0.71, 0.73	0.75, 0.76, 0.76	0.71, 0.71, 0.73	0.69, 0.69, 0.70		0.70, 0.69, 0.71					
	В	0.74, 0.74, 0.74	0.73, 0.75, 0.75	0.73, 0.73, 0.73	0.73, 0.74, 0.75	0.73, 0.74, 0.74	0.73, 0.74, 0.74		0.73, 0.72, 0.74				
		0.74, 0.74, 0.74	0.73, 0.74, 0.75	0.74, 0.75, 0.75	0.74, 0.75, 0.75	0.72, 0.74, 0.74	0.72, 0.74, 0.74		0.83, 0.83, 0.83				
	С	0.74, 0.75, 0.76	0.71, 0.70, 0.72	0.74, 0.75, 0.75	0.71, 0.70, 0.72		0.69, 0.68, 0.70			0.69, 0.68, 0.70			
-		0.75, 0.76, 0.76	0.71, 0.71, 0.72	0.75, 0.76, 0.76	0.71, 0.72, 0.73		0.68, 0.68, 0.70			0.69, 0.69, 0.70			
	D	0.74, 0.74, 0.73	0.73, 0.75, 0.75	0.73, 0.73, 0.73	0.73, 0.75, 0.75			0.73, 0.74, 0.75	0.73, 0.73, 0.74	0.73, 0.74, 0.75			
	D	0.74, 0.74, 0.74	0.73, 0.74, 0.75	0.74, 0.75, 0.74	0.74, 0.75, 0.75			0.73, 0.75, 0.75	0.83, 0.83, 0.83	0.73, 0.74, 0.75			
* The	e sequ	uence of numbers	in each cell, listed	from left to right,	represents the AU	JC values for LR,	RF, and GB mode	els. Also, the seque	ence of three num	pers represents			
first t	he M	lale and then the F	emale.										

Table 3. The AUC values for three models were calculated across all conorts using the gender data	Table 3.	The AUC	values for three	models were	calculated across	all cohorts	using the	gender datas
--	----------	---------	------------------	-------------	-------------------	-------------	-----------	--------------

The performance analysis across gender subgroups reveals consistent results for the LR, RF, and GB models (Table 3). Minimal differences in AUC values are observed for both training and testing sets within each gender category (male and female), indicating stability in model performance across genders. These findings suggest that there is no significant gender bias present in the AUC metrics, highlighting the models' consistent performance regardless of gender distinctions.

Table 4 highlights the variability in model performance across different training and test cohorts, emphasizing the impact of cohort selection on model generalizability. The GB model generally performs better across different racial groups, indicating its effectiveness in distinguishing between cohorts. However, challenges arise in the Native Hawaiian group, suggesting potential overfitting during training and poor generalization ability. Overall, the table

Table 5. The AUC values for three models were calculated across all cohorts using the race dataset

						Test (Race)				
		А	В	С	D	B-A	B-C	D-A	D-B	D-C
		0.71, 0.72, 0.72*	0.69, 0.69, 0.70	0.71, 0.72, 0.72	0.69, 0.69, 0.71	0.68, 0.67, 0.69		0.69, 0.68, 0.69		
		0.74, 0.75, 0.76	0.70, 0.70, 0.72	0.74, 0.75, 0.75	0.70, 0.70, 0.72	0.68, 0.68, 0.70		0.69, 0.68, 0.70		
	Α	0.76, 0.78, 0.77	0.76, 0.76, 0.77	0.76, 0.78, 0.77	0.76, 0.77, 0.77	0.76, 0.67, 0.78		0.77, 0.77, 0.79		
		0.86, 0.78, 0.82	0.75, 0.62, 0.67	0.86, 0.79, 0.82	0.75, 0.62, 0.67	0.48, 0.30, 0.42		0.48, 0.24, 0.42		
		0.88, 0.92, 0.88	0.72, 0.71, 0.72	0.88, 0.92, 0.88	0.72, 0.71, 0.72	0.71, 0.70, 0.72		0.72, 0.70, 0.71		
Trair		0.70, 0.71, 0.71	0.72, 0.73, 0.73	0.70, 0.71, 0.71	0.72, 0.73, 0.73	0.72, 0.74, 0.74	0.72, 0.74, 0.74		0.75, 0.67, 0.75	
		0.73, 0.73, 0.73	0.72, 0.74, 0.74	0.73, 0.73, 0.73	0.72, 0.74, 0.74	0.71, 0.73, 0.73	0.71, 0.73, 0.73		0.77, 0.76, 0.78	
	В	0.76, 0.75, 0.73	0.78, 0.78, 0.78	0.76, 0.75, 0.73	0.78, 0.79, 0.79	0.81, 0.81, 0.81	0.81, 0.81, 0.81		0.89, 0.81, 0.90	
		0.86, 0.80, 0.78	0.75, 0.65, 0.64	0.86, 0.81, 0.78	0.75, 0.65, 0.64	0.41, 0.41, 0.41	0.41, 0.43, 0.45		Nan, Nan, Nan	
		0.88, 0.88, 0.92	0.75, 0.76, 0.75	0.88, 0.88, 0.92	0.75, 0.75, 0.75	0.74, 0.75, 0.75	0.74, 0.75, 0.75		0.84, 0.77, 0.76	
		0.71, 0.72, 0.72	0.69, 0.69, 0.70	0.71, 0.72, 0.72	0.69, 0.66, 0.70		0.68, 0.67, 0.69			0.68, 0.68, 0.69
_		0.74, 0.75, 0.76	0.70, 0.70, 0.72	0.74, 0.75, 0.76	0.70, 0.70, 0.72		0.68, 0.68, 0.70			0.68, 0.68, 0.70
	С	0.76, 0.78, 0.78	0.76, 0.76, 0.76	0.76, 0.78, 0.78	0.76, 0.77, 0.76		0.76, 0.76, 0.75			0.77, 0.77, 0.77
		0.86, 0.80, 0.79	0.75, 0.62, 0.67	0.86, 0.80, 0.79	0.75, 0.62, 0.67		0.45, 0.30, 0.51			0.45, 0.23, 0.51
		0.88, 0.88, 0.88	0.71, 0.71, 0.71	0.88, 0.88, 0.88	0.72, 0.71, 0.70		0.71, 0.70, 0.71			0.72, 0.71, 0.70
		0.70, 0.71, 0.71	0.72, 0.73, 0.74	0.70, 0.71, 0.71	0.72,0.73, 0.74			0.73, 0.74, 0.75	0.75, 0.75, 0.75	0.73, 0.74, 0.75
		0.73, 0.73, 0.72	0.72, 0.74, 0.74	0.73, 0.73, 0.72	0.72, 0.74, 0.74			0.72, 0.74, 0.74	0.78, 0.77, 0.78	0.72, 0.73, 0.74
	D	0.76,0.75,0.71	0.78, 0.78, 0.78	0.76, 0.75, 0.71	0.78, 0.79, 0.79			0.82, 0.82, 0.81	0.89, 0.83, 0.91	0.82, 0.82, 0.81
		0.86, 0.81, 0.77	0.74, 0.65, 0.63	0.86, 0.82, 0.77	0.75, 0.65, 0.64			0.41, 0.46, 0.43	Nan, Nan, Nan	0.41, 0.39, 0.43
		0.88,0.88, 0.92	0.75, 0.76, 0.75	0.88, 0.88, 0.92	0.75, 0.76, 0.75			0.74, 0.75, 0.75	0.84, 0.77, 0.74	0.74, 0.75, 0.75
* The	e seq	uence of numbers	in each cell, listed	from left to right,	represents the AU	JC values for LR,	RF, and GB mode	ls. Also, the seque	ence of three numb	bers represents
first t	he B	lack and then the	White, Asian, Nati	ive Hawaiian. Mul	lti Race .					-

Table 4. The AUC values for three models were calculated across all cohorts using the ethnicity dataset

		Test (Ethnicity)											
		А	В	С	D	B-A	B-C	D-A	D-B	D-C			
Trair	٨	0.77, 0.78, 0.78*	0.73, 0.71, 0.75	0.77, 0.78, 0.79	0.73, 0.71, 0.75	0.66, 0.59, 0.68		0.66, 0.60, 0.69					
	A	0.74, 0.74, 0.75	0.70, 0.70, 0.72	0.73, 0.74, 0.75	0.70, 0.70, 0.72	0.68, 0.65, 0.70		0.69, 0.65, 0.71					
	D	0.77, 0.78, 0.77	0.76, 0.76, 0.77	0.77, 0.78, 0.77	0.76, 0.76, 0.77	0.73, 0.74, 0.74	0.72, 0.73, 0.74		0.71, 0.71, 0.75				
	Б	0.73, 0.73, 0.73	0.72, 0.72, 0.74	0.73, 0.73, 0.73	0.73, 0.74, 0.74	0.72, 0.70, 0.73	0.72, 0.70, 0.73		0.78, 0.75, 0.78				
	C	0.77, 0.78, 0.78	0.73, 0.71, 0.74	0.77, 0.78, 0.79	0.73, 0.71, 0.75		0.65, 0.59, 0.68			0.65, 0.59, 0.68			
-	C	0.73, 0.74, 0.75	0.70, 0.70, 0.72	0.73, 0.74, 0.75	0.70, 0.70, 0.72		0.68, 0.65, 0.70			0.69, 0.65, 0.71			
	D	0.77, 0.78, 0.78	0.76, 0.76 0.77	0.77, 0.78, 0.78	0.75, 0.76, 0.77			0.73, 0.74, 0.75	0.71, 0.71, 0.76	0.72, 0.74, 0.75			
	D	0.73, 0.73, 0.72	0.72, 0.74, 0.74	0.73, 0.73, 0.72	0.73, 0.74, 0.74			0.73, 0.70, 0.73	0.78, 0.75, 0.78	0.72, 0.70, 0.73			
* The	e sequ	uence of numbers	in each cell, listed	from left to right,	represents the AU	JC values for LR,	RF, and GB mode	els. Also, the seque	ence of three num	oers represents			
first	he F	Ispanic and then t	the Not Hispanic.										

underscores the importance of training on diverse and representative data for better model generalization and fairness across different racial groups.

Table 5 emphasizes training on diverse and representative data for better model generalization and fairness across different ethnic groups. The GB model consistently demonstrates superior performance across various ethnic groups, indicating its effectiveness in distinguishing between cohorts. However, it encounters challenges in the Native Hawaiian group, suggesting potential overfitting during training and a diminished ability to generalize effectively.

Comparison of Feature Importance

The RFFI analysis across cohorts A, B, C, and D, excluding 'age', shows 'ABG indices', 'Respiratory Rate', and 'Chloride' as top predictors in Cohort A, with 'Cardiac Troponin I', 'Heart Rate', and 'Creatinine' also notable. In Cohort B, 'ABG indices', 'Cardiac Troponin I', and 'Respiratory Rate' are most influential, alongside 'Not Hispanic or Latino', 'Hispanic or Latino', and 'FiO2'. Cohort C highlights 'ABG indices', 'Respiratory Rate', and 'Chloride', with 'Cardiac Troponin I', 'Heart Rate', and 'Creatinine' also significant. Cohort D emphasizes 'ABG indices', 'Cardiac Troponin I', and 'Not Hispanic or Latino', with 'Hispanic or Latino', 'CBC with PLT', and 'MSRMNTS_FiO2' also important. This analysis underscores the need for cohort-specific feature consideration in healthcare ML models, with 'ABG indices', 'Respiratory Rate', and 'Cardiac Troponin I' consistently important, guiding feature selection and model optimization. The analysis also reveals the importance of cohort-specific features, such as ethnicity indicators in Cohorts B and D, suggesting that tailored models considering these unique characteristics could improve predictive accuracy. This insight underscores the need for a nuanced approach to feature selection and model development, considering both universally important features and those specific to individual cohorts.

Correlation between models trained on the same cohorts

We conducted correlation analyses using both Pearson's and Spearman's methods to assess the performance of three different ML models across multiple cohorts. Each cohort was provided with its own training and testing datasets, and we examined the correlations between these datasets when the models were trained and tested under various combinations. You can observe Pearson's and Spearman's correlation results in Table 6, illustrating the correlations between LR and RF, LR and GB, and RF and GB. Table 6 highlights a strong positive correlation between different ML models, particularly between LR and GB, across all cohorts. This indicates a high degree of agreement in their predictions. However, some variability in correlation strength between other model pairs, such as LR-RF and RF-GB, suggests nuances in data processing. Spearman's correlation consistently yields higher scores than Pearson's, indicating its effectiveness in capturing the relationship between model predictions.

		Test											
		А	В	С	D	B-A	B-C	D-A	D-B	D-C			
Train	Α	0.90, 0.88, 0.87* 0.91, 0.95, 0.95	0.87, 0.88, 0.86 0.87, 0.93, 0.90	0.90, 0.88, 0.87 0.91, 0.96, 0.95	0.87, 0.88, 0.86 0.87, 0.93, 0.90	0.78, 0.91, 0.78 0.86, 0.92, 0.88		0.77, 0.91, 0.77 0.86, 0.92, 0.88					
	В	0.90, 0.87, 0.86 0.89, 0.93, 0.93	0.91, 0.90, 0.89 0.92, 0.95, 0.96	0.90, 0.87, 0.86 0.90, 0.93, 0.94	0.91, 0.90, 0.89 0.92, 0.93, 0.96	0.86, 0.93, 0.82 0.91, 0.95, 0.96	0.91, 0.90, 0.89 0.91, 0.95, 0.96		0.86, 0.93, 0.83 0.93, 0.96, 0.96				
	С	0.91, 0.84, 0.85 0.91, 0.93, 0.95	0.88, 0.81, 0.82 0.88, 0.88, 0.89	0.91, 0.84, 0.85 0.91, 0.93, 0.95	0.88, 0.81, 0.82 0.88, 0.88, 0.89		0.86, 0.78, 0.77 0.86, 0.84, 0.82			0.77, 0.91, 0.78 0.87, 0.85, 0.84			
	D	0.90, 0.84, 0.85 0.89, 0.91, 0.93	0.91, 0.88, 0.89 0.92, 0.94, 0.96	0.90, 0.84, 0.85 0.90, 0.91, 0.93	0.92, 0.88, 0.89 0.92, 0.95, 0.96			0.86, 0.93, 0.83 0.92, 0.94, 0.96	0.87, 0.93, 0.83 0.94, 0.96, 0.97	0.86, 0.93, 0.83 0.92, 0.94, 0.96			
* The three	* The sequence of numbers in each cell, listed from left to right, represents the correlations between LR and RF, LR and GB, and RF and GB. Also, the sequence of three numbers represents first the Pearson's and then the Spearman's correlation.												

Table 6. Correlation Coefficient between two models on the same cohorts

Table 7. Correlation Coefficient between two models on the different cohorts

						Test				
Train	Train	А	В	С	D	B-A	B-C	D-A	D-B	D-C
А	В	0.98, 0.95, 0.87	0.90, 0.86, 0.87	0.98, 0.95, 0.87	0.90, 0.86, 0.87	0.85, 0.79, 0.83				
		0.98,0.95,0.94	0.90,0.85,0.91	0.98,0.95,0.94	0.90,0.85,0.91	0.85,0.79,0.87				
	С	1.00, 1.00, 0.97	1.00, 1.00, 0.94	1.00, 1.00, 0.97	1.00, 1.00, 0.94					
		1.00, 1.00, 0.98	1.00, 1.00, 0.97	1.00, 1.00, 0.98	1.00, 1.00, 0.97					
	D	0.98, 0.95, 0.83	0.90, 0.86, 0.86	0.98, 0.95, 0.84	0.90, 0.86, 0.86			0.86, 0.80, 0.82		
		0.98,0.95,0.93	0.90,0.86,0.91	0.98,0.95,0.93	0.90,0.86, 0.91			0.86,0.80,0.87		
	C	0.98, 0.95, 0.84	0.89, 0.86, 0.83	0.98, 0.95, 0.84	0.89, 0.86, 0.83		0.84, 0.79, 0.78		1.00, 1.00, 0.99	
р	C	0.98,0.95,0.92	0.89,0.85,0.89	0.98,0.95,0.93	0.89,0.85,0.89		0.84,0.79, 0.86		1.00, 1.00, 0.99	
В	D	1.00, 1.00, 0.98	1.00, 1.00, 0.99	1.00, 1.00, 0.98	1.00, 1.00, 0.99					
	D	1.00, 1.00, 0.99	1.00, 1.00, 0.99	1.00, 1.00, 0.99	1.00, 1.00, 0.99					
C	D	0.98, 0.95, 0.81	0.90, 0.86, 0.82	0.98, 0.95, 0.82	0.89, 0.86, 0.82					0.84, 0.80, 0.77
C	D	0.98,0.95,0.91	0.90,0.86,0.89	0.98,0.96,0.91	0.89,0.86,0.89					0.84,0.79,0.85
* The s	sequence	e of numbers in ea	ch cell, listed fron	n left to right, repr	esents the correlat	ions between LR a	und LR, RF and R	F, and GB and GB	on the different c	ohorts. Also, the
sequen	ce of thr	ee numbers repres	sents first the Pear	son's and then the	Spearman's correl	ation.				

Correlation between models trained on the different cohorts

Table 7 underscores the strong positive correlations between ML models trained on different cohorts, indicating a high degree of agreement in their predictions. This is particularly evident in the high correlation coefficients for the LR model, which demonstrates consistent correlations across different cohorts. Despite the overall strong positive correlation, there is some variability in the strength of the correlation between different pairs of models. For instance, the RF model shows moderate variability, while the Gradient Boosting GB model exhibits significant fluctuations when generalizing across cohorts. Table 7 includes Pearson's and Spearman's correlation coefficients, with Spearman's correlation generally yielding higher scores. This indicates that Spearman's correlation captures the relationship between the models' predictions more effectively, underscoring the robust correlation between different ML models across different cohorts. This emphasizes the importance of considering different correlation measures to understand the relationship between models and their generalization ability fully.

Model Calibration

Model calibration indicates how far "scores" outputted from an ML model are from the actual probabilities. As one can see in Figure 3, Gradient Boost is the best calibrated model, yet it does not output any scores above about 0.7. Logistic Regression and Random Forest models significantly overestimate predictions. More importantly, the calibration worsens when models are applied to data from outside the cohort. This is the most prominent model trained on Cohort A, which is applied to test data from Cohort D (right side of the figure).



Figure 3: Example calibration plots for models trained on Cohort A and tested on Cohorts A, B, C and D.

Conclusion

This study meticulously underscores the profound influence that decisions made during data preprocessing exert on the composition and scale of cohorts. Such decisions can introduce selection bias, potentially compromising the integrity of research findings. The inclination towards constructing more exclusive cohorts necessitates a judicious approach, as such exclusions could diminish the representativeness of the data and raise concerns regarding the fairness of the outcomes. Moreover, comparing feature importance across cohorts illuminates the nuanced interplay between input variables and predictive outcomes, emphasizing the need for tailored approaches to model development that account for cohort-specific characteristics. Furthermore, correlation analyses between models trained on the same and different cohorts reveal varying degrees of agreement, underscoring the complexity of model generalization and the importance of cohort-specific considerations. The research presented herein highlights the imperative for ongoing inquiry into the ramifications of preprocessing choices on biases within ML models, advocating for the development of bespoke strategies to mitigate these biases. Our team's research endeavors are dedicated to evaluating the impact of these preprocessing decisions on the quality of ML models derived from the data. Our current findings definitively establish that the cohorts differ, and these differences significantly affect ML outcomes. Future endeavors will concentrate on domains where biases manifest more prominently, scrutinizing the effects of diverse preprocessing methodologies on the generalizability of research. These efforts are anticipated to yield critical insights that will fortify the external validity of studies. Adopting a discerning and informed methodology in data preprocessing is paramount to enhancing the reliability and equity of ML applications across a spectrum of fields. Social determinants of health, such as Gender, Race, and Ethnicity, have been shown to play a significant role in COVID-19 outcomes, further emphasizing the need for careful cohort selection and bias mitigation in ML models to ensure fair and accurate predictions.

Discussion

The study highlights the critical importance of thoughtful cohort selection in health data processing, emphasizing that arbitrary decisions can introduce biases affecting the representativeness, fairness, and generalizability of ML models in healthcare. Standardized and transparent selection criteria are crucial to ensure data representativeness, avoiding exclusion of important subpopulations or introducing noise. Cross-cohort validation is needed to ensure model robustness and reliability across different clinical settings. The insight that certain features are consistently important across cohorts while others vary can guide feature selection and model optimization. Implementing rigorous data preprocessing protocols, including sensitivity analyses and considering ensemble methods or meta-learning, can improve generalization. Ongoing inquiry into the ramifications of preprocessing choices on biases within ML models is essential, particularly in domains where biases are prominent. Given the significant role of social determinants of health in COVID-19 outcomes, it is imperative to consider these factors in ML model development to ensure fair and accurate predictions.

It is important to acknowledge certain limitations within the scope of this study. Primarily, the analysis is confined to data from August 1, 2020, to December 31, 2021, which may not encompass subsequent evolutions in COVID-19 healthcare practices, data documentation, or patient demographics. Additionally, not all preprocessing decisions outlined in the decision tree have been thoroughly examined, i.e., the selection of patients aged 18 and over remains a point of contention, as alternative age thresholds such as 21 could be equally justifiable.

References

- 1. Heckman, J. J. Selection Bias. in *Encyclopedia of Social Measurement* (ed. Kempf-Leonard, K.) 463–468 (Elsevier, New York, 2005). doi:10.1016/B0-12-369398-5/00115-8.
- 2. Heckman, J. Varieties of Selection Bias. The American Economic Review 80, 313-318 (1990).
- 3. Lu, H., Cole, S. R., Howe, C. J. & Westreich, D. Toward a clearer definition of selection bias when estimating causal effects. *Epidemiology* 33, 699–706 (2022).
- 4. Nakamura, K. & Kawabata, H. I Choose, Therefore I Like: Preference for Faces Induced by Arbitrary Choice. *PLoS One* 8, e72071 (2013).
- 5. Quick Safety Issue 23: Implicit bias in health care | The Joint Commission. https://www.jointcommission.org/resources/newsand-multimedia/newsletters/newsletters/quick-safety/quick-safety-issue-23-implicit-bias-in-health-care/.
- 6. Gopal, D. P., Chetty, U., O'Donnell, P., Gajria, C. & Blackadder-Weinstein, J. Implicit bias in healthcare: clinical practice, research and decision making. *Future Healthc J* 8, 40–48 (2021).
- Vokinger, K. N., Feuerriegel, S. & Kesselheim, A. S. Mitigating bias in machine learning for medicine. *Commun Med* 1, 1–3 (2021).
- 8. Suri, J. S. *et al.* Understanding the bias in machine learning systems for cardiovascular disease risk assessment: The first of its kind review. *Comput Biol Med* 142, 105204 (2022).
- 9. Mac Namee, B., Cunningham, P., Byrne, S. & Corrigan, O. I. The problem of bias in training data in regression problems in medical decision support. *Artif Intell Med* 24, 51–70 (2002).
- 10. N3C Home. https://covid.cd2h.org/.
- 11. Haghighathoseini, A., Qodrati, M., Min, H., Leslie, T.F., Frankenfeld, C.L., Menon, N.M. and Wojtusiak, J. Selection Bias from Data Processing in N3C. *Proceedings of the IEEE 12th International Conference on Healthcare Informatics (ICHI), Orlando, FL, 2024.*
- 12. Cha, G.-W., Moon, H.-J. & Kim, Y.-C. Comparison of Random Forest and Gradient Boosting Machine Models for Predicting Demolition Waste Based on Small Datasets and Categorical Variables. *Int J Environ Res Public Health* 18, 8530 (2021).
- 13. Nusinovici, S. *et al.* Logistic regression was as good as machine learning for predicting major chronic diseases. *Journal of Clinical Epidemiology* 122, 56–69 (2020).
- 14. Domanski, P. A., Steven Brown, J., Heo, J., Wojtusiak, J. & McLinden, M. O. A thermodynamic analysis of refrigerants: Performance limits of the vapor compression cycle. *International Journal of Refrigeration* 38, 71–79 (2014).
- 15. Irvin, K. & Wojtusiak, J. Comparison of Classification Learning Methods for Medical Claims Payments. in AMIA (2012).
- Lee, Y. & Seo, J. Suggestion of statistical validation on feature importance of machine learning. Annu Int Conf IEEE Eng Med Biol Soc 2023, 1–4 (2023).
- 17. Algehyne, E. A., Jibril, M. L., Algehainy, N. A., Alamri, O. A. & Alzahrani, A. K. Fuzzy Neural Network Expert System with an Improved Gini Index Random Forest-Based Feature Importance Measure Algorithm for Early Diagnosis of Breast Cancer in Saudi Arabia. *BDCC* 6, 13 (2022).
- Ling, C. X., Huang, J. & Zhang, H. AUC: A Better Measure than Accuracy in Comparing Learning Algorithms. in *Advances in Artificial Intelligence* (eds. Xiang, Y. & Chaib-draa, B.) 329–341 (Springer, Berlin, Heidelberg, 2003). doi:10.1007/3-540-44886-1_25.

Acknowledgement

The analyses described in this publication were conducted with data or tools accessed through the NCATS N3C Data Enclave https://covid.cd2h.org and N3C Attribution & Publication Policy v 1.2-2020-08-25b supported by NCATS Contract No. 75N95023D00001, Axle Information Subcontract: NCATS-P00438-B. This research was possible because of the patients whose information is included within the data and the organizations (https://ncats.nih.gov/n3c/resources/data-contribution/data-transfer-agreement-signatories) and scientists who have contributed to the on-going development of this community resource [https://doi.org/10.1093/jamia/oca196]