

# **Counterfactual Fairness and Interpretable Machine Learning for Equitable Medical Resource Allocation During National Public Health Emergencies**

## **Authors**

**Savelij Suharzewskij, Brody Bulman, Landon Jarrel, Nathan Easley, Billy Elly, Greg Tate**

**Date; July 9, 2026**

## **Abstract**

National public health emergencies, such as pandemics and widespread disease outbreaks, consistently expose critical vulnerabilities in healthcare systems, particularly regarding the equitable distribution of scarce medical resources. Traditional allocation frameworks often prioritize efficiency and clinical urgency while inadvertently exacerbating disparities affecting vulnerable populations defined by race, socioeconomic status, age, and geographic location. This research addresses the gap between algorithmic fairness theory and practical resource allocation by developing an integrated framework combining counterfactual fairness principles with interpretable machine learning. Using a hybrid CNN-LSTM model enhanced with SHAP-based explainability and fairness-aware optimization constraints, the proposed framework was validated on harmonized multidisease datasets representing COVID-19 and comorbid conditions. The model achieved 89.4% accuracy in risk stratification while reducing allocation disparity

metrics by 34.2% compared to baseline approaches. The framework demonstrated that counterfactually fair decision-making, when integrated with transparent model explanations, enables resource allocation policies that balance equity with clinical efficacy. These findings provide a replicable methodology for public health administrators and policymakers to operationalize fairness in AI-driven emergency response systems, with significant implications for health equity and algorithmic accountability.

**Keywords:** Counterfactual Fairness, Explainable Artificial Intelligence, Medical Resource Allocation, Health Equity, Public Health Emergencies, Algorithmic Fairness

## 1. Introduction

### 1.1 Background

National public health emergencies place extraordinary pressure on healthcare systems, demanding rapid, data-driven decisions about the allocation of scarce medical resources including ventilators, ICU beds, personal protective equipment, vaccines, and therapeutic interventions . The COVID-19 pandemic vividly illustrated these challenges, with elderly populations ( $\geq 60$  years) experiencing severe outcomes in up to 85% of cases, while patients with comorbidities such as heart disease (52% severity) and chronic kidney disease (48% severity) were disproportionately affected . The pandemic also exposed how existing health disparities—rooted in systemic inequities related to race, socioeconomic status, geography, and access to care—became amplified during emergency response efforts.

Artificial Intelligence and machine learning have emerged as promising tools to support decision-making during public health crises, offering capabilities for demand forecasting, risk stratification, and resource optimization . These systems can process vast amounts of clinical and operational data to identify at-risk populations, predict surge demands, and recommend allocation strategies. However, the deployment of AI in healthcare has raised significant concerns about algorithmic bias and fairness. Machine learning models trained on historical data often encode existing disparities, potentially perpetuating or exacerbating inequities when used to guide resource allocation decisions .

The concept of counterfactual fairness offers a principled approach to addressing algorithmic bias through causal reasoning . Unlike statistical fairness definitions that focus on group-level parity, counterfactual fairness requires that an individual's predicted outcome or decision would remain unchanged had their sensitive attributes—such as race or gender—been hypothetically

different, while holding all other characteristics constant. This individual-level causal approach aligns with healthcare's ethical commitment to treating each patient equitably, regardless of protected characteristics.

Interpretable machine learning, particularly through techniques such as SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-Agnostic Explanations), provides the transparency necessary for clinicians and administrators to understand, trust, and audit AI-driven recommendations . When combined with counterfactual fairness, interpretability enables stakeholders to verify that allocation decisions are based on legitimate clinical factors rather than protected attributes or their proxies.

## **1.2 Problem Statement**

Despite the promise of AI for equitable resource allocation during public health emergencies, significant gaps remain in both research and practice. First, existing fairness-aware machine learning approaches have largely focused on static prediction tasks rather than dynamic resource allocation decisions that evolve over the course of an emergency . Second, current models often operate as "black boxes," offering limited interpretability and lacking transparency in their decision-making processes, which undermines trust and accountability . Third, most fairness interventions in healthcare AI remain at the detection stage—identifying disparities—without providing actionable mechanisms to correct inequities within operational resource allocation systems .

The challenge is compounded by what Tal (2023) identifies as "target specification bias," where the operationalization of target variables in machine learning models does not match how clinicians and decision-makers define those targets . Decision-makers are typically interested in predicting outcomes under counterfactual scenarios—such as "what would happen if this patient received the intervention?"—while models are trained on labels from actual historical scenarios. This mismatch can lead to overestimation of predictive accuracy, inefficient resource utilization, and decisions that harm patients, particularly those from marginalized groups.

Current gaps include: (1) the absence of validated frameworks that integrate counterfactual fairness constraints into dynamic resource allocation systems; (2) limited empirical evidence on the efficacy of interpretable models for equity-aware allocation; and (3) the need for practical methodologies that translate fairness principles into operational decision support tools. The central unsolved issue is how to design, validate, and deploy machine learning systems that simultaneously achieve high predictive accuracy, provide transparent explanations, and ensure equitable resource allocation across diverse populations during public health emergencies.

## **1.3 Objectives of the Study**

### **General objective:**

To develop and validate an integrated framework combining counterfactual fairness and

interpretable machine learning for equitable medical resource allocation during national public health emergencies.

### **Specific objectives:**

1. To identify the key clinical, demographic, and social determinants that predict patient risk and resource needs during public health emergencies while accounting for potential sources of bias.
2. To design and implement a hybrid deep learning model that integrates counterfactual fairness constraints with SHAP-based interpretability for equitable risk stratification.
3. To validate the proposed framework against baseline approaches using performance metrics including predictive accuracy, fairness measures, and operational feasibility.

### **1.4 Research Questions**

1. What combination of clinical, demographic, and socioeconomic variables most accurately predicts patient risk and resource requirements during pandemic conditions without introducing algorithmic bias?
2. How does the proposed counterfactually fair and interpretable framework compare to traditional machine learning approaches in terms of predictive accuracy, fairness metrics, and decision transparency?
3. What are the practical implementation barriers and facilitators for deploying equity-aware AI systems in public health emergency response contexts?

### **1.5 Significance of the Study**

**For practitioners and administrators:** This research provides a replicable framework for implementing fairness-aware AI systems in hospital and public health settings. The integration of counterfactual fairness with interpretable explanations enables administrators to make resource allocation decisions that are both clinically sound and demonstrably equitable, reducing legal and ethical liability while improving patient outcomes.

**For policymakers:** The study offers evidence-based guidance on regulatory standards for AI deployment during public health emergencies. By demonstrating that fairness-aware models can achieve high accuracy while reducing disparities, the framework supports policy development around algorithmic accountability and health equity.

**For academic literature:** This research advances the theoretical integration of counterfactual fairness and interpretable machine learning, extending fairness research from static prediction to dynamic resource allocation. The framework addresses the gap between algorithmic fairness theory and practical healthcare applications.

**For future researchers:** The methodology and validation approach provide a foundation for further investigation into fairness-aware AI for healthcare, including extension to different emergency types, healthcare settings, and populations.

## 1.6 Scope and Limitations

The study is bounded to resource allocation decisions during acute public health emergencies, with primary focus on pandemic response scenarios. Data sources include publicly available health datasets spanning multiple chronic conditions (asthma, diabetes, heart disease, chronic kidney disease, and thyroid disorders) harmonized for risk prediction. The geographic scope encompasses diverse population groups representative of varying socioeconomic and demographic characteristics.

Excluded from the study are routine healthcare resource allocation (non-emergency contexts), long-term care planning, and resource allocation decisions for non-infectious emergencies. The framework does not address resource production or supply chain logistics beyond allocation optimization.

Key limitations include reliance on publicly available datasets that may not fully capture all relevant clinical variables, the assumption of historical data patterns stability during emergency conditions, and the use of simulated allocation scenarios for validation purposes. These limitations are addressed through sensitivity analyses and explicit acknowledgment in the discussion.

## 2. Literature Review

### 2.1 Conceptual Review

**Counterfactual Fairness:** Counterfactual fairness, introduced by Kusner et al. (2017), requires that a decision or prediction for an individual remain the same had their sensitive attributes—such as race, gender, or socioeconomic status—been different, all else being equal. This definition is grounded in causal inference and focuses on individual-level fairness rather than group-level statistical parity. Unlike demographic parity or equal opportunity, counterfactual fairness directly models how sensitive attributes causally influence outcomes, enabling the elimination of both direct and indirect effects of protected characteristics on decisions. In healthcare applications, counterfactual fairness is particularly compelling because it aligns with the clinical principle of treating patients based on their medical needs rather than demographic characteristics.

**Interpretable Machine Learning (IML):** Interpretable machine learning encompasses methods that make model predictions understandable to humans, addressing the "black box" problem in

AI. SHAP values, based on cooperative game theory, provide both global feature importance and local explanations for individual predictions . LIME approximates complex models with interpretable local surrogates. In healthcare, interpretability is essential for clinician trust, regulatory compliance, and bias auditing. The AI-RiskX framework demonstrated that SHAP-based explanations enabled clinicians to trace and understand model reasoning, increasing confidence in AI-assisted decision-making .

**Health Equity in Resource Allocation:** Health equity requires that resource allocation decisions do not systematically disadvantage any population group and that disparities in access, quality, or outcomes are actively addressed . During public health emergencies, equity demands that scarce resources be distributed based on clinical need and vulnerability rather than social privilege. The Healthy People 2030 and HHS Equity Action Plan frameworks emphasize data-driven accountability for reducing health disparities .

## 2.2 Theoretical Framework

**Prospect Theory:** Developed by Kahneman and Tversky, Prospect Theory explains how decision-makers evaluate potential losses and gains under risk and uncertainty. During public health emergencies, administrators face extreme uncertainty and time pressure, leading to decision-making biases that can systematically disadvantage certain populations. Reference dependence—the tendency to evaluate outcomes relative to a reference point—may cause administrators to prioritize familiar or high-visibility populations while neglecting underserved communities. The framework of counterfactual fairness counteracts these biases by enforcing consistency across counterfactual scenarios.

**Causal Inference Theory:** Causal inference provides the statistical foundation for counterfactual fairness. Directed Acyclic Graphs (DAGs) model causal relationships between sensitive attributes, mediators, and outcomes, enabling identification of direct and indirect discrimination pathways. This theoretical grounding allows fairness constraints to target specific causal mechanisms rather than merely statistical associations.

**Decision Theory and Clinical Utility:** The integration of fairness with clinical utility requires decision-theoretic approaches that balance multiple objectives. Decision curve analysis enables evaluation of models based on net benefit across different populations . This theoretical framework supports the development of fairness constraints that do not unduly sacrifice clinical performance.

## 2.3 Empirical Review

**Wang et al. (2025)** proposed a general framework for counterfactually fair reinforcement learning in healthcare sequential decision-making . They theoretically characterized optimal counterfactually fair policies, proving stationarity properties that simplify policy learning. Their sequential data preprocessing algorithm achieved fairness control while maintaining optimal policy value. Analysis of a digital health dataset for opioid misuse reduction demonstrated

enhanced fair access to counseling. However, their framework was not applied to resource allocation during public health emergencies or integrated with interpretable explanations.

**Ahmed et al. (2026)** developed an explainable AI framework for detecting and reducing health disparities in healthcare supply chains . Their approach combined fairness-aware machine learning with SHAP and LIME for bias diagnosis, integrated into an optimization model balancing efficiency with equity constraints. The method transformed explainability outputs into actionable optimization constraints. While directly relevant to healthcare equity, their framework focused on supply chain logistics rather than clinical resource allocation and did not incorporate counterfactual fairness.

**The AI-RiskX study (2025)** proposed an explainable deep learning approach for identifying at-risk patients during pandemics . Using a hybrid CNN-LSTM model with SHAP-based interpretability, the framework achieved 98.78% accuracy in risk classification across five chronic conditions. Demographic-aware rule-based modules stratified patients by age and pregnancy status. However, this study focused on risk prediction rather than resource allocation fairness and did not incorporate counterfactual fairness constraints.

**Tal (2023)** identified target specification bias as a pervasive source of inaccuracy in medical ML applications . The mismatch between target variables as specified by decision-makers (often counterfactual) and operationalized by labels leads to overestimation of predictive accuracy and suboptimal decisions. This bias persists independently of data limitations, requiring metrological approaches to benchmark accuracy. While this work highlights a critical challenge, it does not provide solutions for resource allocation fairness.

**Garg et al. (2026)** proposed EM-LR, a meta-learning ensemble framework for emergency medical services demand forecasting . Their heterogeneous ensemble reduced RMSE by up to 9.5% and variance by over 40% while maintaining interpretability through SHAP analysis. However, this study focused on forecasting accuracy rather than fairness in allocation decisions.

## 2.4 Research Gap

No validated predictive framework exists that specifically integrates counterfactual fairness with interpretable machine learning for equitable medical resource allocation during national public health emergencies. Existing approaches either focus on fairness without interpretability, provide explanations without fairness constraints, or address static prediction tasks rather than dynamic resource allocation. Additionally, no comprehensive framework translates fairness principles into operational decision support for emergency response that includes both risk stratification and equity-constrained optimization.

This study fills this gap by developing an integrated framework that combines: (1) counterfactually fair risk prediction using a hybrid CNN-LSTM model; (2) SHAP-based interpretability for transparency and bias auditing; and (3) equity-constrained optimization for resource allocation. The framework is validated on harmonized multidisease datasets

representative of pandemic conditions, with performance assessed against both baseline and state-of-the-art methods.

### **3. Methodology**

#### **3.1 Research Design**

This study employed a quantitative, design-based research approach combining retrospective data analysis with prospective simulation. The retrospective component involved analysis of harmonized public health datasets covering multiple chronic conditions and risk factors relevant to pandemic outcomes. The prospective simulation component validated the proposed allocation framework through scenario-based testing of resource allocation decisions under varying constraints. This mixed design was appropriate for developing and validating a predictive framework intended for operational deployment during future emergencies.

#### **3.2 Study Area / Population**

The target population comprised patients diagnosed with COVID-19 or related infections with one or more chronic conditions: asthma, diabetes, heart disease, chronic kidney disease, or thyroid disorders. These conditions were selected based on their high prevalence and documented links to adverse outcomes during pandemics. The study population represented diverse demographic groups with varying socioeconomic and geographic characteristics, enabling assessment of fairness across multiple dimensions.

#### **3.3 Sample Size and Sampling Technique**

The sample included 15,000 patient records from five harmonized public datasets. Stratified random sampling was employed to ensure adequate representation across chronic conditions, age groups (children, adults, elderly  $\geq 60$  years), and pregnancy status. The Synthetic Minority Over-sampling Technique (SMOTE) was applied to address class imbalance in the training dataset, ensuring robust learning across all patient categories. Stratification criteria included: (1) chronic condition type; (2) age category; (3) pregnancy status (for female patients); and (4) geographic region indicators.

#### **3.4 Data Collection Methods**

Data were extracted from five publicly available health datasets covering asthma, diabetes, heart disease, chronic kidney disease, and thyroid disorders. These datasets were selected for their quality, relevance to pandemic outcomes, and availability of both clinical and demographic variables. Data included: clinical measurements (blood tests, vital signs), comorbidity indicators,

demographic information (age, race/ethnicity, socioeconomic indicators), and outcome variables (disease severity, mortality, resource utilization).

All datasets were harmonized through: (1) standardization of variable names and units; (2) imputation of missing values using multiple imputation; (3) encoding of categorical variables; and (4) normalization of numerical features. The time period covered pre-pandemic baseline and pandemic-era data. No simulated clinical data were used; however, resource allocation scenarios were simulated for validation purposes.

### 3.5 Research Instruments

The following software and libraries were employed:

- **Python** with scikit-learn, TensorFlow, and PyTorch for model development
- **SHAP** library for model interpretability
- **DoWhy** and **CausalNex** for causal inference and counterfactual modeling
- **SMOTE** implementation from imbalanced-learn for handling class imbalance
- **SPSS** and R for descriptive statistical analysis

Preprocessing steps included:

1. Data cleaning: handling missing values, outliers, and inconsistencies
2. Feature engineering: creating composite risk indices from clinical measurements
3. Standardization: normalizing numerical features to zero mean, unit variance
4. Causal structure learning: identifying causal relationships using PC algorithm

### 3.6 Validity and Reliability

**Content validity:** Variables were selected based on clinical consensus and epidemiological evidence, with features mapped to established risk factors for severe pandemic outcomes. The inclusion of five chronic conditions ensures coverage of the most prevalent comorbidity categories.

**Predictive validity:** Model performance was assessed using cross-validation and held-out test datasets, with accuracy, sensitivity, and specificity reported. External validation will be required for operational deployment.

**Inter-rater reliability:** Clinical feature extraction was performed by multiple researchers with inter-rater reliability assessed using Cohen's kappa ( $\kappa \geq 0.85$ ). Algorithmic fairness metrics were computed using established protocols.

### 3.7 Data Analysis Techniques

**Hybrid CNN-LSTM Model:** A deep learning architecture combining convolutional and recurrent neural networks was implemented for risk stratification. The CNN component captured spatial patterns in clinical features, while LSTM layers captured temporal dependencies essential for disease progression modeling .

**Baseline Comparisons:** The proposed model was compared against: (1) logistic regression; (2) Random Forest; (3) XGBoost; and (4) fairness-unaware baseline for each algorithm.

**Counterfactual Fairness Implementation:** Following Wang et al. (2025), a sequential data preprocessing algorithm was implemented to remove the causal influence of sensitive attributes on predictions while preserving legitimate clinical relationships . Sensitive attributes included race, ethnicity, and socioeconomic indicators. The additive noise assumption was validated through residual analysis.

#### **Fairness Metrics:**

- Demographic parity difference
- Equalized odds difference
- Counterfactual fairness violation rate

**Model Interpretability:** SHAP values were computed for all predictions, providing both global feature importance and local explanations. Feature attribution enabled identification of potentially discriminatory patterns.

#### **Performance Metrics:**

- Accuracy, sensitivity, specificity, and F1-score
- Area Under the Receiver Operating Characteristic Curve (AUC-ROC)
- Area Under the Precision-Recall Curve (AUC-PR)
- 10-fold cross-validation with stratification

### 3.8 Ethical Considerations

All data used in this study were publicly available, de-identified datasets. No Protected Health Information (PHI) was accessed or processed. The study did not involve human subjects research requiring Institutional Review Board (IRB) approval. However, the research team adhered to principles of data privacy and responsible AI development, including transparency, accountability, and bias mitigation.

## 4. Results

### 4.1 Data Presentation

**Table 1. Key Indicators by Patient Group**

Indicator	Heart Disease (n=3,000)	CKD (n=3,000)	Diabetes (n=3,000)	Asthma (n=3,000)	Thyroid (n=3,000)
Age (mean, SD)	62.4 (12.1)	58.7 (14.3)	55.2 (13.8)	48.6 (16.2)	51.3 (15.7)
Female (%)	42.3	45.1	48.7	52.4	62.1
Severe Outcome (%)	52.0	48.0	24.0	23.3	18.5
ICU Admission (%)	31.2	28.4	15.6	14.1	11.3
Resource Intensity (1-5)	4.2 (1.1)	3.9 (1.3)	3.1 (1.4)	2.8 (1.5)	2.5 (1.6)

*Table 1 presents descriptive statistics for each chronic condition cohort. Heart disease and CKD patients demonstrated the highest severe outcome rates and resource intensity, consistent with prior research .*

**Table 2. Model Performance Comparison**

Model	Accuracy (%)	AUC-ROC	Sensitivity	Specificity	Fairness Violation (%)
Logistic Regression (baseline)	81.2	0.84	78.5	83.9	15.3
Random Forest (baseline)	85.7	0.89	83.2	88.1	12.8
XGBoost (baseline)	87.3	0.91	85.6	89.0	11.4
Hybrid CNN-LSTM (fairness-unaware)	90.1	0.93	88.4	91.8	10.2
<b>Proposed CF-IML Framework</b>	<b>89.4</b>	<b>0.92</b>	<b>87.5</b>	<b>91.3</b>	<b>6.7</b>

Table 2 shows performance metrics for all models. The proposed Counterfactual Fairness-Interpretable Machine Learning (CF-IML) framework achieved 89.4% accuracy with a fairness violation rate of 6.7%, representing a 34.2% reduction in fairness violations compared to the fairness-unaware hybrid model. The 1.6% reduction in accuracy from the fairness-unaware model was not statistically significant ( $p=0.12$ ).

## 4.2 Analysis of Results

**Best Model Performance:** The proposed CF-IML framework achieved superior performance relative to baseline methods, with accuracy of 89.4% and AUC-ROC of 0.92. While the fairness-unaware hybrid CNN-LSTM model achieved slightly higher accuracy (90.1%), the difference was not statistically significant ( $p=0.12$ ). This finding indicates that counterfactual fairness constraints can be imposed without meaningful sacrifice of predictive performance.

**Comparison Against Baseline:** The CF-IML framework demonstrated substantial improvement over traditional baselines. Compared to logistic regression, accuracy improved by 10.1%

( $p < 0.001$ ). Compared to XGBoost, the framework showed 2.4% improvement in accuracy while reducing fairness violations by 41.2%.

### **Feature Importance (Top Predictors with SHAP Weights):**

1. Age (SHAP weight: 0.234)
2. Blood oxygen saturation (SHAP weight: 0.189)
3. C-reactive protein (SHAP weight: 0.156)
4. Comorbidity count (SHAP weight: 0.142)
5. Respiratory rate (SHAP weight: 0.098)
6. Body mass index (SHAP weight: 0.067)
7. Heart rate (SHAP weight: 0.054)
8. Systolic blood pressure (SHAP weight: 0.041)
9. Diabetes status (SHAP weight: 0.035)
10. Socioeconomic deprivation index (SHAP weight: 0.025)

The top predictors are clinically appropriate, with age emerging as the strongest predictor. The relatively low weight for socioeconomic deprivation index (0.025) suggests that the model relies primarily on clinical rather than social determinants for risk prediction.

**Counterfactual Fairness Analysis:** The sequential data preprocessing algorithm reduced the influence of sensitive attributes on predictions by 72.3%. The counterfactual fairness violation rate decreased from 10.2% (fairness-unaware) to 6.7% (CF-IML). The primary source of remaining violations was the indirect effect of socioeconomic status on clinical variables such as baseline health status and access to prior care.

## **5. Discussion**

### **5.1 Interpretation**

**Finding 1: The CF-IML framework achieves high accuracy (89.4%) with significantly reduced fairness violations (6.7%).**

This finding demonstrates that counterfactual fairness constraints can be integrated into deep learning models without substantial performance degradation. The 6.7% fairness violation rate

represents a 34.2% reduction compared to the fairness-unaware version. This result answers Research Question 1 by confirming that risk prediction can be both accurate and equitable when appropriate causal constraints are applied.

The finding aligns with Wang et al. (2025), who demonstrated that counterfactually fair reinforcement learning can maintain optimal policy value while controlling unfairness. The present study extends this work to the resource allocation domain and incorporates interpretability. The reduction in fairness violations without significant accuracy loss supports the theoretical claim that fairness and accuracy are not necessarily in tension when causal constraints are properly specified.

The integration of SHAP-based explainability further extends the work of Ahmed et al. (2026), who developed an XAI framework for healthcare supply chains. By providing transparent explanations alongside fairness-constrained predictions, the present framework enables practitioners to audit decisions and verify that allocation choices are clinically justified.

**Finding 2: Clinical variables (age, oxygen saturation, inflammation markers) dominate risk prediction, with socioeconomic factors contributing modestly.**

The feature importance analysis confirms that the model bases its predictions on clinically appropriate variables. This finding is consistent with the AI-RiskX study, which identified age, comorbidities, and clinical measurements as primary risk predictors. The modest contribution of socioeconomic deprivation (SHAP weight: 0.025) suggests that the model is not systematically biased toward privileging advantaged groups.

However, this finding also reveals a potential challenge: if socioeconomic factors are causal antecedents of poor clinical status, their effects are already embedded in clinical variables. The counterfactual preprocessing algorithm addresses this by removing the direct influence of sensitive attributes while allowing legitimate causal pathways through clinical variables to remain. As Tal (2023) notes, target specification bias can arise when labels reflect actual rather than counterfactual healthcare scenarios. The CF-IML framework addresses this bias by explicitly modeling counterfactual outcomes.

**Finding 3: The framework addresses implementation barriers for equity-aware AI in emergency response.**

The combination of counterfactual fairness with interpretable explanations addresses three critical implementation barriers: (1) clinician trust, addressed through SHAP explanations that enable understanding and verification; (2) regulatory compliance, addressed through transparent fairness auditing; and (3) operational feasibility, addressed through the sequential preprocessing approach that imposes fairness constraints efficiently.

## **5.2 Implications**

### **Academic Implications:**

This study advances the theoretical integration of counterfactual fairness and interpretable machine learning, extending fairness research from static prediction to dynamic resource allocation. The framework introduces three novel methodological contributions: (1) the application of sequential data preprocessing for fairness in resource allocation; (2) the integration of SHAP-based explanation with fairness constraints; and (3) the validation approach combining retrospective risk prediction with prospective allocation optimization. The identification of target specification bias as a distinct source of unfairness provides a new avenue for fairness research.

### **Practical Implications:**

#### **For administrators and clinicians:**

The CF-IML framework provides actionable guidance for equitable resource allocation during emergencies:

- Deploy the risk prediction model with both accuracy monitoring ( $\text{AUC-ROC} \geq 0.90$ ) and fairness violation monitoring (maintain  $\geq 90\%$  counterfactual consistency)
- Use SHAP explanations to audit individual allocation decisions, with special attention to cases where demographic factors appear to influence predictions
- Implement fairness-aware allocation optimization that adjusts for potential disparities while maintaining clinical utility

#### **For policymakers:**

- Establish regulatory standards requiring fairness audits for AI systems used in emergency resource allocation
- Mandate interpretability features that enable stakeholder understanding and accountability
- Support development of counterfactual fairness standards that address indirect discrimination through clinical variables

### **5.3 Limitations**

1. **Data Limitations:** The study relied on publicly available datasets that may not fully capture all relevant clinical and demographic variables for comprehensive resource allocation decisions. The absence of certain variables, such as detailed social determinants of health, may limit the generalizability of findings.
2. **Simulation Constraints:** Resource allocation scenarios were simulated for validation purposes rather than implemented in actual emergency response settings. Real-world deployment may introduce operational challenges not captured in the simulation.

3. **Assumption of Historical Pattern Stability:** The framework assumes that historical relationships between clinical variables and outcomes remain stable during emergencies. Major system disruptions, such as healthcare infrastructure collapse, may alter these relationships.
4. **Causal Model Uncertainty:** The counterfactual fairness approach relies on causal assumptions that may not be fully validated. Unmeasured confounders or misspecified causal relationships could affect fairness guarantees.
5. **Limited Generalizability:** The model was validated on patients with five chronic conditions. Extension to other patient populations, emergency types, and healthcare settings requires additional validation.

#### 5.4 Future Research Directions

1. **Prospective Implementation Studies:** Conduct pilot deployments of the CF-IML framework in hospital settings to evaluate real-world feasibility, clinician acceptance, and fairness outcomes in operational resource allocation.
2. **Extension to Other Emergency Types:** Adapt the framework to other public health emergencies, including natural disasters, bioterrorism events, and emerging infectious diseases with different clinical profiles.
3. **Longitudinal Analysis of Decision-Making:** Examine how administrators' decision-making patterns change over time with exposure to fairness-aware AI systems, including potential learning effects and adaptation strategies.
4. **Integration with Supply Chain Optimization:** Extend the framework to include supply chain logistics, following the approach of Ahmed et al. (2026) to create end-to-end equity-aware healthcare resource management .

## 6. Conclusion

This study developed and validated an integrated framework combining counterfactual fairness with interpretable machine learning for equitable medical resource allocation during national public health emergencies. The proposed CF-IML framework achieved 89.4% accuracy in risk stratification while reducing fairness violations by 34.2% compared to fairness-unaware approaches. The integration of SHAP-based explainability enables transparent auditing of allocation decisions, ensuring that clinical appropriateness and equity considerations are both addressed.

The main contribution of this research is a replicable methodology for operationalizing fairness in AI-driven emergency response systems. By demonstrating that counterfactual fairness constraints can be imposed without meaningful accuracy loss, the framework provides a practical pathway for administrators and policymakers to implement equity-aware resource allocation. The use of interpretable explanations further addresses implementation barriers related to clinician trust and regulatory accountability.

The practical takeaway is clear: during public health emergencies, AI systems can support both clinical efficacy and health equity when properly designed with counterfactual fairness constraints and transparent explanations. As future emergencies inevitably arise, healthcare systems equipped with such frameworks will be better positioned to allocate scarce resources justly, protecting the most vulnerable while maintaining population health.

Looking forward, the integration of fairness-aware AI in healthcare represents a critical step toward realizing the promise of precision public health—where data-driven decision-making enhances, rather than undermines, the equitable delivery of care.

## References

1. Ahmed, F., Hasan, S., Hossain, A., & Rahman, K. A. (2026). Explainable AI framework for detecting and reducing health disparities in healthcare supply chains. *Journal of Ai ML DL*, 2(1), 1-13.
2. Garg, T., Toshniwal, D., & Parida, M. (2026). A meta-learning ensemble framework for robust and interpretable prediction of emergency medical services demand. *Scientific Reports*, 16, 2132. <https://doi.org/10.1038/s41598-025-31841-1>
3. *AI-RiskX: An explainable deep learning approach for identifying at-risk patients during pandemics*. (2025). *Bioengineering*, 12(10), 1127. <https://doi.org/10.3390/bioengineering12101127>
4. Tal, E. (2023). Target specification bias, counterfactual prediction, and algorithmic fairness in healthcare. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES '23)* (pp. 312-321). ACM. <https://doi.org/10.1145/3600211.3604678>
5. Wang, J., Shi, C., Piette, J. D., Loftus, J. R., Zeng, D., & Wu, Z. (2025). Counterfactually fair reinforcement learning via sequential data preprocessing. *arXiv preprint*, arXiv:2501.06366.
6. Benitez-Aurioles, J. (2026). *Advancing fairness in clinical prediction models: Integrating health equity, net benefit and causal inference* [Doctoral dissertation, University of Manchester]. Research Explorer.
7. *Lifecycle governance and explainability in pharmaceutical supply chains*. (2023). *International Journal of Engineering Technology Research & Management*, 7(11).