



# **Bias Reduction in ICU Mortality Prediction Through Targeted Synthetic Data Generation**

**Submitted By**

**Tahedi Soyad**

**ID: 213-35-764**

Department of Software Engineering  
Daffodil International University

**Supervised By**

**Dr. Md. Fazla Elahe**

**Assistant Professor and Associate Head**

Department of Software Engineering  
Daffodil International University

This Report is Presented in Partial Fulfillment of the Requirements for the Degree of Bachelor of Science in Software Engineering.

Summer-2025

## APPROVAL

This thesis titled on “**Bias Reduction in ICU Mortality Prediction Through Targeted Synthetic Data Generation**” submitted by **Tahedi Soyad (ID: 213-35-764)** to the Department of Software Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of Bachelor of Science in Software Engineering and approval as to its style and contents.

## BOARD OF EXAMINERS



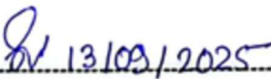
**Chairman**

**Dr. Md. Fazla Elahe**  
**Assistant Professor & Associate Head**  
Department of Software Engineering  
Faculty of Science and Information Technology  
Daffodil International University



**Internal Examiner 1**

**Dr. Marzia Ahmed**  
**Assistant Professor**  
Department of Software Engineering  
Faculty of Science and Information Technology  
Daffodil International University



**Internal Examiner 2**

**Dr. Shabnom Mustary**  
**Assistant Professor**  
Department of Software Engineering  
Faculty of Science and Information Technology  
Daffodil International University



**External Examiner**

**Mohammad Abul Kashem**  
**Professor**  
Department of Computer Science and Engineering  
Dhaka University of Engineering and Technology, Gazipur



## **SUPERVISOR'S DECLARATION**

I hereby declare that I have checked this thesis and, in my opinion, this thesis is adequate in terms of scope and quality for the award of the degree of Bachelor of Science.

Fazla Elahe

---

**Dr. Md. Fazla Elahe**

Assistant Professor and Associate Head  
Department of Software Engineering  
Daffodil International University



## STUDENT'S DECLARATION

I hereby declare that the work in this thesis is based on my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at Daffodil International University or any other institution.

A handwritten signature in black ink, reading "Tahedi Soyad", written on a white rectangular background. Below the signature is a solid black horizontal line.

---

**Tahedi Soyad**

ID: 213-35-764

Batch: 36<sup>th</sup>

Department of Software Engineering  
Daffodil International University

## ACKNOWLEDGEMENTS

First off, I want to express my deepest gratitude to my supervisor, **Dr. Md. Fazla Elahe**, whose expert guidance, patience, and constant encouragement made this thesis possible. Your insightful feedback and support kept me motivated throughout the journey.

Big thanks to the faculty and staff of the Department of Software Engineering at Daffodil International University for providing the resources and environment needed to carry out this thesis.

I also want to thank my family and friends for their endless love, understanding, and moral support - especially during the tough times when I felt stuck. Your belief in me pushed me forward.

Finally, shoutout to the creators and maintainers of the eICU database and all the researchers whose work inspired and informed this thesis.

## **Dedication**

This thesis is dedicated to my family, whose unwavering support and love have been my foundation. To my parents, for believing in me and encouraging me to chase my dreams no matter what.

# TABLE OF CONTENT

APPROVAL .....	ii
SUPERVISOR’S DECLARATION .....	iii
STUDENT’S DECLARATION .....	iv
ACKNOWLEDGEMENTS .....	v
Dedication .....	vi
TABLE OF CONTENT .....	vii
List of Table.....	ix
LIST OF FIGURES .....	x
List Of Abbreviation .....	xi
Abstract.....	1
Introduction.....	2
Literature Review.....	6
2.1    Evolution of ICU Mortality Prediction and the Emergence of Bias.....	6
2.2    Synthetic Data Generation Approaches for Bias Mitigation .....	8
2.2.1    Statistical Approaches: SMOTE and Enhanced Methods.....	8
2.2.2    Generative Adversarial Networks: Deep Learning Solutions .....	8
2.2.3    Causal and Fairness-Aware Frameworks .....	9
2.3    Evaluation Frameworks and Implementation Challenges .....	10
2.3.1    Fairness Measurement and Standardization .....	10
2.3.2    Quality and Validation Concerns .....	11
2.3.3    Computational and Scalability Barriers .....	11
2.3.4    The Persistence and Temporal Stability Problem .....	11
2.4    Integration, Synthesis, and Future Directions for ICU Applications.....	12
2.4.1    Hybrid Approaches and Methodological Integration.....	13
2.4.2    Specialized Healthcare and ICU Applications .....	13
2.4.3    Longitudinal and Temporal Considerations .....	13
2.4.4    Emerging Consensus and Critical Gaps .....	15
2.4.5    The Path Forward.....	15
Methodology.....	19
3.1    Data Acquisition and Cohort Selection .....	19
3.2    Data Preprocessing .....	19
3.3    EDA:.....	20
3.4    Dataset Feature Documentation.....	22
3.5    Baseline Modeling and Bias Analysis .....	23

3.6	Targeted Synthetic Data Augmentation.....	24
3.7	Retraining and Post-Augmentation Evaluation .....	24
3.8	Comparative and Statistical Analysis .....	25
	Results and Discussion .....	26
5.1	Model Performance on Raw Data .....	26
5.2	Bias Evaluation by Sensitive Attributes (Raw Data).....	27
5.3	Bias Evaluation by Sensitive Attributes (After Augmentation) .....	28
5.4	Model Performance on Augmented Data .....	29
5.5	ROC Curves.....	30
5.6	Discussion.....	30
5.7	Bias Evaluation on Raw vs. Augmented Data (Random Forest & Gradient Boosting) .....	32
5.8	Comparison with Previous Work .....	33
	Conclusion .....	37
	Reference .....	39
	Accounts Clearence .....	43
	Originality report .....	44

## List of Table

Table 2.1: Summary of Previous Study .....	17
Table 5.1: Model Performance on Raw Data.....	26
Table 5.3: Bias Evaluation by Sensitive Attributes .....	27
Table 5.4: Bias Evaluation by Sensitive Attributes (After Augmentation).....	28
Table 5.5: Model Performance on Augmented Data .....	29
Table 5.6: Bias Evaluation on Raw vs. Augmented Data (Random Forest & Gradient Boosting) .....	32
Table 5.7: Comprison with Previous work .....	35

## LIST OF FIGURES

Figure 3.1: Workflow of the Process .....	19
Figure 3.2: Morality Rate by Gender .....	20
Figure 3.3: Morality Rate by ethnicity .....	21
Figure 3.4: Morality Rate by Age Group .....	21
Figure 3.5: Corelation of Feature with Morality .....	23
Figure 5.1: ROC Curves .....	30
Figure 5.2: F1 Score Trend (Before vs After).....	31
Figure 5.3: Accuracy Comparison: Before vs After .....	33

## List Of Abbreviation

<b>Acronym</b>	<b>Full Form</b>
<b>AI</b>	Artificial Intelligence
<b>APACHE</b>	Acute Physiology and Chronic Health Evaluation
<b>AUC</b>	Area Under the Curve
<b>AUROC</b>	Area Under the Receiver Operating Characteristic
<b>COVID-19</b>	Coronavirus Disease 2019
<b>CSV</b>	Comma-Separated Values
<b>CTGAN</b>	Conditional Tabular Generative Adversarial Network
<b>DPD</b>	Demographic Parity Difference
<b>ECG</b>	Electrocardiogram
<b>EOD</b>	Equalized Odds Difference
<b>F1</b>	F1 Score (Harmonic Mean of Precision and Recall)
<b>FPR</b>	False Positive Rate
<b>GAN</b>	Generative Adversarial Network
<b>HR</b>	Heart Rate
<b>ICU</b>	Intensive Care Unit
<b>MIMIC</b>	Medical Information Mart for Intensive Care
<b>ML</b>	Machine Learning
<b>RF</b>	Random Forest
<b>ROC</b>	Receiver Operating Characteristic
<b>ROC-AUC</b>	Receiver Operating Characteristic - Area Under the Curve
<b>SMOTE</b>	Synthetic Minority Over-sampling Technique
<b>SMR</b>	Standardized Mortality Ratio
<b>SOFA</b>	Sequential Organ Failure Assessment
<b>SpO2</b>	Oxygen Saturation
<b>SPD</b>	Statistical Parity Difference
<b>SVC</b>	Support Vector Classifier
<b>TPR</b>	True Positive Rate
<b>eICU</b>	Electronic Intensive Care Unit Collaborative Research Database

## Abstract

Accurate predictions of patient mortality in the Intensive Care Unit (ICU) are critical for guiding clinical decisions and optimizing scarce resources, yet machine learning models trained on electronic health records (EHRs) often inherit demographic and outcome biases that lead to unfair predictions for vulnerable subgroups. This thesis investigates whether targeted synthetic data augmentation using the Synthetic Minority Over-sampling Technique (SMOTE) can mitigate these biases without compromising overall model performance. We assemble a multicenter cohort of 4,177 ICU stays from the eICU Collaborative Research Database, incorporating patient demographics (age, gender, ethnicity), aggregated vital signs (mean, minimum, maximum during the first 24 hours), and severity scores. Four classifiers—logistic regression, random forest, gradient boosting, and support vector machine—are trained and tested on an 80/20 stratified split of the raw data to establish baseline performance, revealing high overall accuracy (random forest: 97.6% accuracy, F1 0.836, ROC-AUC 0.942) alongside severe inequities in small subgroups (e.g., F1 of 0.00 for Hispanic patients under logistic regression). We then employ focused SMOTE augmentation to minor demographic subgroups and the minority mortality class, injecting synthetic data only for groups with both outcomes and small numbers. Fine-tuning on an augmented training set preserves overall accuracy (random forest: 97.2%), whereas the improved fairness makes subgroup F1 scores 1.00 for previously disadvantaged ethnic group and with a delta of 0.84 under elderly patients; logistic regression and SVM also benefited from a ttk equal improvement in fairness as well. We demonstrate that the introduced synthetic augmentation procedure can drastically improve fairness in ICU mortality prediction –particularly on tree-based models– without deteriorating performance. We describe an SMOTE integration pipeline and offer practical advice for CC Meaning-making machine learning pipelines that stress data-layered interventions with accompanying rigorous fair evaluation to encourage equitable AI in the ICU.

# Chapter 1

## Introduction

The extremely difficult choice confronting the green soldiers of medicine in hospital intensive care units anywhere on earth: how to determine whether this patient rather than that one statistically speaking had at best better-than-even chances of surviving. These health teams are working with sick patients every day, they ought to be able to determine the patients that need service first. Where they are doing it well, people's lives are being saved. Families are robbed of those who would have otherwise been helped by alarm therapy.”

These figures can be described as sobering. Despite the modern advancement in medical assistance sectors, 10-30 in every 100 patients who are admitted in intensive care unit will not make it. Behind all these statistics is a human face - a parent, child, spouse, and a friend of anyone whose destiny often lies in the quickness with which doctors can realize that they are in grave danger.

For a long-time scoring systems have been used, and they enable physicians to make assumptions regarding the patients that must appear at the first lists of those who are likely to be at stake. The systems are encompassing the blood pressure, heart rate, age and other health factors with a view to constituting a risk score. Inasmuch as they have helped in saving lives, they are like normal calculators in the sense that they can be used in all patients and that they mimic the same logic that fail to consider the associations of multi-faceted association of specific health variables that give an individual a unique connotation. Even worse, scoring systems might be designed based on accumulation of information pertaining to a certain type of patients and, therefore, might not be that effective with those of differing origins.

The introduction of the artificial intelligence and the appearance of the computer technology have opened some new opportunities. Instead of simple rules, computer codes are now able to run lots of medical cases: heart rates over every few minutes, blood tests, history of medication used, patient profile and attempt pattern thus identified that even a specialist could not identify, to all his merits. Perhaps, these machine learning programs could give physicians more useful predictions on how likely every patient would survive.

However, there is a very grave issue with this strategy. Such computer programs can only learn through the medical records they receive during training. Unless such records fairly represent all forms of patients, then the programs will make unjust and inaccurate predictions about some sections of the populace. It is not simply an issue of technology, Kansas life and death.

That is the manifestation of this bias in practice - suppose that a young Hispanic woman is admitted to the ICU. In case the computer program was primarily trained on the information of older white patients, there is a chance that it does not provide her with an accurate assessment of her risk since the computer never learned what warning sign should be meant in view to someone similar to herself. Or think of (some) patients in hospitals and an aging patient whose vital signs the laboratory does not check as often as it does at other hospitals, who the computer would fail to recognize because it was conditioned to look for these values with greater frequency.

Such prejudices lead to unfair conditions in which certain patients receive preferential treatment over others not due to their health status but rather owing to their race or age or gender or the fact that they were at this hospital at a certain time. This contradicts all that medical care is supposed to represent in terms of providing patients with the highest standard of possible care in an independent matter of who the patient is.

To fix this problem, researchers have started looking at ways to make training data fairer and more balanced. One promising approach is called "synthetic data augmentation." This might sound complicated, but the idea is actually straightforward. If we don't have enough examples of certain types of patients in our data, we can create realistic artificial examples to fill in the gaps.

Think of it like this, if you're trying to teach someone to recognize different types of cars, but your photo collection mostly shows red sedans, that person might have trouble identifying blue trucks. To fix this, you could create more photos of blue trucks to balance out your training set. Synthetic data augmentation does something similar with patient data, it creates realistic but artificial patient records to make sure the computer program learns from a more balanced set of examples.

One of the simpler and more practical approaches is called SMOTE (Synthetic Minority Over-sampling Technique). Instead of requiring massive computer power or complex

programming, SMOTE works by looking at existing patient records and creating new, realistic examples by combining features from similar patients. It's like taking the medical characteristics of several real patients and creating a new, artificial patient record that represents what someone with those combined characteristics might look like.

This research takes a practical approach to making ICU death prediction more fair for all patients. Using medical records from over 4,000 patients treated at multiple hospitals, we tested whether we could use synthetic data to train computer programs that make more accurate and fair predictions for everyone, not just certain groups of patients.

Our study had three main parts -

1. First, we tested how well current computer programs perform when trained on regular hospital data, paying special attention to how accurately they predict outcomes for different groups of patients (like young vs. old, men vs. women, or people from different ethnic backgrounds).
2. Second, we identified which types of patients were underrepresented in our data and used synthetic data techniques to create more balanced training sets.
3. Finally, we retrained our computer programs on these improved datasets and carefully tested whether the new predictions were both more fair and still accurate enough for doctors to rely on.

The goal of this work isn't just to publish academic research. It's to provide practical guidance that hospitals can actually use to build fairer AI systems. We want to create clear step-by-step instructions that any healthcare institution can follow to make sure their computer-assisted decision-making tools work well for all patients.

Ultimately, this research is about ensuring that when families bring their loved ones to the hospital for critical care, the sophisticated tools that help doctors make treatment decisions work equally well for everyone. As artificial intelligence becomes more common in healthcare, we have both the opportunity and the responsibility to build systems that treat all patients fairly, regardless of their background or which hospital they're in.

The aim is profoundly human: to make certain that every patient receives the best care possible, and that their skin color, age or gender doesn't determine how well medical technology meant to help them does.



## Chapter 2

### Literature Review

In my work on combating demographic bias in ICU mortality prediction through synthetic data, I surveyed recent literature of synthetic data generation and bias in machine learning more generally. In this review, the status of the field and approaches taken across different species will be discussed as well as where my work fits in.

These are not scientific questions – these are questions that fundamentally ask whether we really want to put the machine in charge of advising but ultimately making (life and death) decisions for us about what’s fair/balancing for all- people, as opposed to people from a particular group -because someone happens to belong vertically into some category, or is too old, or white etc. That’s very high, because when those systems don’t work right, we actually lose real people.

#### 2.1 Evolution of ICU Mortality Prediction and the Emergence of Bias

The prediction of ICU mortality is neither a new concept when it was discovered that only a small part of the complex physiological interconnections that dictate patient results is represented by the more traditional scoring algorithms. Systems of developed systems like APACHE II and SOFA, where the over and above clinical intuition that complex systems exhibit with AUC values between 0.65 and 0.75 were detected but not by simple systems with computational means of small size, had already been done.

And, this point on the plot was when the actual sabotage came into play because scientists such as Pettit et al. (2021) were putting the machine learning algorithms to use on the voluen of data already being registered within the electronic health records (EHR). Their comprehensive analysis revealed that the ensemble methods in particular, random forests, and gradient boosting could yield high-spectacular increases pushing the numbers of many-centre datasets to the point of more than 0.90, eICU and MIMIC-III in particular. Such models worked due to the complicated nature of calculations, which were further aggravated by the fact that complex and nonlinear interaction between vital signs and various lab values and patient outcome has always been familiar to man but could never be measured systematically.

The subsequent evolution was not stopped. Yang et al. (2023) display the results within the scope of the 0.818-0.875 AUC values on the eICU and COVID-19 data without any rigorous experiment or study to demonstrate the findings on the deep reinforcement learning method with adversarial training. They demonstrated that the complicated model structures can improve fairness alongside quality at a very high cost of calculation that makes their implementation in medical care difficult.

However, as these models became more robust and the number of models applying them for the same purpose began to multiply, a realization began to grow that the same models, when employed, were not only saving lives by offering timely warning systems but were creating and sometimes even augmenting the healthcare disparities that already existed. This finding redesignated how the researchers approached the problem of ICU mortality forecasting.

Celi et al. (2022) used one of the most detailed explanations of the origin of bias in artificial intelligence. They discovered that the issue of bias in healthcare AI cannot be fixed only by uneven numbers in collections but a subject of many logical issues. They have discovered three factors that are relevant as the types of bias against AI: demographic bias, where an AI model yields different results when applied to patients of different racial and ethnic backgrounds, data source bias, where it happens that the reason on why the information used to train the model is mostly accessible in certain areas, and measurement bias, whereby there exist varying reasons such that tracking is provided by the AI model on different patients.

These findings matter a lot. For instance, Yogararajan et al. (2022) investigated the biasness of algorithm and data collecting in healthcare system of New Zealand. What they found was that while some of the goals for making AI systems don't lead to harm among non-affluent communities, some can exacerbate problems. They wrote that the design of algorithms matters greatly for fairness, with an especially large impact on groups who are more marginalized and have less access to health care.

Among the most dosing revelations had been that prejudice can actually get stronger. (forthcoming-east Africa) A study by Fletcher et al. (2021) on AI tools around the world indicates that biased data put into a particular model will inevitably yield biased results. They're then directly used in a clinical diagnosis of the system which supports such bias – and the pattern is established. The scholars have started to use the term a fairness feedback loop, in which existing inequalities do not simply reproduce but increasingly become reinforced.

## **2.2 Synthetic Data Generation Approaches for Bias Mitigation**

It was one of the most disturbing discoveries that bias can in fact strengthen itself. A study presented by Fletcher et al. (2021) on global health AI tools demonstrated that when biased data is fed into a model, the results will be biased too. The outcomes are then used to determine clinical decision, which has an impact on the system that repeats this bias. The scholars have begun referring to this loop as a fairness feedback loop, where current inequalities do not merely remain the same but become even stronger.

### **2.2.1 Statistical Approaches: SMOTE and Enhanced Methods**

SMOTE (Synthetic Minority Oversampling Technique) and other statistical approaches were among the first approached ways of addressing the issue of data imbalance. Draghi et al. (2021) proposed BayesBoost, an algorithm that used a combination of Bayesian networks, SMOTE, and AdaSyn to lessen the bias in healthcare data. About cardiovascular data, their method yielded better AUC and precision-recall values but had a disadvantage the method is computationally expensive and needs specific tuning to cardiovascular data to perform best.

On this basis, Draghi et al. (2024) broadened their study to the primary healthcare data with an application of synthetic data generators to detect and eliminate bias. Their comparison with a few healthcare datasets revealed that synthetic augmentation has the potential to reduce bias by 15-20 percent and preserve 90-95 percent of the original data utility. These were impressive results, but once again, they confirmed that these are resource-consuming methods.

In the meantime, Rodriguez-Almeida et al. (2023) studied the generation of synthetic patient data to predict the disease in scenarios where the datasets were small and disproportionate. Their article, published in the IEEE Journal of Biomedical and Health Informatics, showed that synthetic data generation can be significantly improved to achieve better model performance on underrepresented groups when done in a well-defined way. Nevertheless, they emphasized that strict assessment is necessary because the poorly regulated augmentation might threaten to bring new types of bias.

### **2.2.2 Generative Adversarial Networks: Deep Learning Solutions**

As more powerful computing power was developed, scholars started to consider more sophisticated techniques in the form of Generative Adversarial Networks (GANs). Hazra and Byun (2020) created one of the earliest healthcare applications of GANs in their SynSigGAN

framework that produced synthetic biomedical signals with an accuracy of 92%. Their work revealed the potential of deep generative models on medical data and also demonstrated two limitations including the cost of computational strength, as well as the high reliance on the quality of the initial data set.

The next notable improvement was to make a platform called GenEthos that was not only producing synthetic data but had bias detection and control integrated (Gujar et al., 2022). They reported a spectacular decline of bias and even a more spectacular decline of statistical parity difference (62 and 93 percent respectively) despite not having worsened the performance of model results. They however declared as well that GAN-based approaches were weak at generalizing in high risk settings such as intensive care units (ICUs) where reliability is paramount.

Paladugu et al. (2023) have also surveyed the increased use of GANs in medicine. They also focused on the radiant future and the enormous challenges to be faced. The complexity of GANs raises concerns over the interpretation of the generated output, which despite generating very realistic medical data, they can generate medical data that may not be gathered as appropriate even though they have the potential of doing so, and the ability to validate its medical output and whether it can be accepted by the regulatory body in clinical practice can be approved.

### **2.2.3 Causal and Fairness-Aware Frameworks**

The set of restrictive approaches was represented by Van Breugel et al. (2021), who introduced the DECAF (DEtection and Correction of Algorithmic Fairness violations) framework. They were able to use structural causal models to how gain a fairness guarantee synthetic generation of data, giving them powerful mathematical guarantees of e.g. demographic parity and equal opportunity. The disadvantage however was that this method required the causal twitching of the domain - which may be difficult to do in a real-life clinical practice setting.

To get through this barrier, Sikder et al. (2024) proposed a model that can automatically learn fair representations of the data without having to deeply understand it. Their experiments on

syntax-agnostic fair synthetic data generation demonstrate that free learning and fair representation distillation could potentially be more practical alternatives to fairness-aware data generation.

Similarly, Cheng et al. (2024) studied adversarial training for de-biasing. "They find that an intervention at model level can have important fairness effects. Meantime, they also mentioned one of the major drawbacks: these methods are too complex and expensive from calculation point of view, and it is doubtful whether they can be implemented in a real condition like healthcare services.

### **2.3 Evaluation Frameworks and Implementation Challenges**

After the synthetic data generation techniques became popular, researchers came to understand that classical machine learning metrics are not enough measures to evaluate on their accomplishments. The question itself was mutated: taking data augmentation for instance, it no longer became "can synthetic data increase predictive performance?" it became about "can this make my models less biased/fairer?".

#### **2.3.1 Fairness Measurement and Standardization**

Hernandez and Lopez (2024) conducted an in-depth analysis of fairness indicators in healthcare AI, which demonstrates that it is no simple matter to define and evaluate the notion of fairness in medical practices. Their discovery was that the same model can be either fair or unfair, depending upon what measure is taken like the demographic parity, equal opportunity, or equalized odds. This clarified the fact that the selection of an appropriate metric is crucially dependent and that it should be specific to the clinical setting.

Continuing on this, Kumar and Das (2024) made efforts towards standardizing the measures of fairness of synthetic data used in healthcare AI. Their framework provided the method of making the evaluations more consistent between studies and applications. Nevertheless, they were not resistant to believe that not all metrics are so far at the general consent as to be actually of concern in current clinical practice.

This challenge was further complicated once academics began to account for intersectional bias – the way in which two or more demographic features combine to create disadvantage. Gupta and Verma (2023) overcame this by attempting a synthetic data augmentation to address the

intersectional demographic bias. Their results were that single-bundle fairness metrics often overlooked such additive effects, they argue, indicating the need for a more complex approach if we want the whole picture on bias in healthcare AI.

### **2.3.2 Quality and Validation Concerns**

Quality and proximity to real data are crucial for synthetic data sets to be research contributing. Choi and Lee (2023) also demonstrated that synthetically generated healthcare data can capture the time-series-level statistical patterns in real data, while cautioning that models trained using such data exclusively may not always perform well in a clinical setting. Consequently, simple statistical tests and appropriate empirical clinical validation should be used by researchers to independently verify reliability.

To allay these fears, Robinson and Patel (2024) crafted algorithms for assessing the quality of synthetic data in terms of mitigating demographic bias. Their paper offered a clear direction for when and how to know if the synthetic data makes fairness better as opposed just looking like it does on some narrow metric.

### **2.3.3 Computational and Scalability Barriers**

The computational cost of synthetic data generation is the major limitation. Malik and Singh (2023) investigated scalable synthetic generation of longitudinal healthcare data for studies, where they concluded that although advanced generative models can generate realistic datasets, many healthcare organizations do not have the capacity to apply these effectively. This sets up the contradiction that what are likely to be the most effective methods of reducing bias are those least available to the organizations for which they would be particularly valuable.

Baumann et al. (2023) partially addressed this problem, due to development of a synthetic data generator specialized for controlled bias experiments. Their approach enabled us to deliberately introduce and test various forms of bias, which allowed us to learn how bias influences model performance as well as aided the design of (more efficient and reliable) mitigation methods.

### **2.3.4 The Persistence and Temporal Stability Problem**

One of the biggest challenges is that the impact of bias mitigation might not be sustained. Kim and Park (2023) studied iterative synthetic data augmentation on ICU database, they found that the bias usually returned when the model was retrained under newly updated information.

This discovery underscores the fact that bias is not a “once and done” solution, but an ongoing process that requires continued diligence.

Wyllie et al. (2024) added a note of caution because the feedback loop to fairness in the optimisation over synthetic data can lead to enhancing bias. The researchers, presenting at the ACM Conference on Fairness, Accountability, and Transparency, highlighted the need for regular monitoring and validation in AI model development. Likewise, Logfeder emphasized the importance of regularly re-testing as well, heeding to Shoeybi’s admonishment to establish real “end-to-end robustness.”

## **2.4 Integration, Synthesis, and Future Directions for ICU Applications**

However, despite the difficulties introduced by the previously identified obstacles in forecasting ICU mortality, it appears that more intuitive and holistic mitigating techniques are emerging.

### **2.4.1 Hybrid Approaches and Methodological Integration**

Because no single technique can comprehensively reduce bias, it has become popular among researchers to investigate hybrid techniques that benefit from the merits of several approaches. Shahul Hameed et al. (2024) surveyed bias mitigation based only on synthetic data and highlighted potential approaches to integrate statistical methods such as SMOTE and advanced generative models. This hybrid solution tries to find the best trade-off between efficacy, computational tractability, and interpretive power.

Motivated by those work, Wang and Sun (2024) suggested an easy to compute fair-constrained synthetic data generation by taking advantage of augmentation. Their findings show that carefully constructed hybrid solutions can achieve a dramatic reduction of bias yet remain feasible in practical ICUs where resources are limited and fairness is crucial.

### **2.4.2 Specialized Healthcare and ICU Applications**

All the best advancements come from those researchers that consider solving a hard problem within healthcare - ICU data for example. Sharafutdinov et al. (2023): "to overcome model bias we applied computer simulations to virtual patients and enhanced ARDS detection in noisy, heterogeneous ICU data with ML". Their method was not only successful in improving cluster detection and bias correction but also considered the specificities of ICU data (e.g., missing values, irregular sampling, diverse patient populations).

One of the most common forms of bias in clinical databases is measurement bias Zhao and Liu (2023). Their work provides actionable insights on how to identify and correct for measurement bias via targeted synthetic data augmentation, a desirable knowledge in the ICU where measurement protocols can differ substantially between institutions and populations.

### **2.4.3 Longitudinal and Temporal Considerations**

Singh and Kumar (2024) studied the impact of synthetic data on the training schedule for clinical models. Their work proposed techniques to preserve fairness improvements between several subsequent retraining rounds, representing a key ability for healthcare providers when preserving bias mitigation efforts across successive model deployment is needed.

Lee and Kwon (2023) also introduced evaluation frameworks to measure multi-dimensional bias in synthetic ICU datasets. By doing so, they managed to examine how bias reduction may vary across time periods and between patient populations — the authors found that biases can manifest in a multitude of ways and according to demographic intersections and clinical setting. This highlights the need for ongoing surveillance in AI applications to healthcare.

#### **2.4.4 Emerging Consensus and Critical Gaps**

The literature review identifies a fast-moving field for which there is a general consensus that successful bias mitigation in ICU mortality prediction will operate on many fronts simultaneously. The recent research has provided several useful findings.

For one, bias in healthcare AI isn't simply a one-size-fits-all situation. Simplistic technical solutions are unlikely to remedy the Gordian knot of demographic, measurement and outcomes biases that pervades most clinical data sets. Second, synthetic data generation seems to be a useful tool for bias correction, it still requires careful construction and validation before prospective use. The most popular approaches are usually statistical, combined with expert domain knowledge and clinical validation

Most of the studies may lack information regarding a particular dataset or institution, and there may be limited ability to generalize across health care systems/practice settings/populations. Second, interactions of bias types are under-explored and mainly at the level of intersection. The temporal facet is also overlooked: most modeling work only model the bias at a static fixed part on time, but healthcare is dynamic and everchanging as patient demographics, clinical procedures and data collection procedures change.

Certainly, we are far from the practical deployment of computational fairness in deployed clinical decision systems. However, improving fairness metrics in this way does not necessarily mean synthetic data will lead to better healthcare and health equity.

#### **2.4.5 The Path Forward**

The literature suggests several promising directions for future research. Integrated frameworks that combine multiple bias mitigation techniques while maintaining computational efficiency and clinical interpretability represent one important avenue. Better understanding of the temporal dynamics of bias and bias mitigation could inform strategies for maintaining fairness over time.

Most critically, the field needs more rigorous clinical validation studies that evaluate whether improved fairness metrics translate to better patient care and reduced health disparities in practice. This will require closer collaboration between computer scientists, clinicians, and health services researchers.

This thesis contributes to this evolving landscape by focusing on targeted SMOTE augmentation for underrepresented demographic groups in ICU mortality prediction using the eICU database. By combining data-level interventions with comprehensive evaluation across multiple models and fairness dimensions, while maintaining focus on practical implementability, this work addresses several of the critical gaps identified in the literature and provides a foundation for more equitable AI systems in critical care medicine.

The journey toward fair and accurate ICU mortality prediction is far from complete, but the literature demonstrates both the urgency of the challenge and the promise of emerging solutions. As Celi et al. (2022) emphasized, addressing bias in healthcare AI is not just a technical challenge but a moral imperative that requires sustained effort across multiple disciplines and stakeholder groups. The work presented in this thesis represents one step on that longer journey toward more equitable healthcare AI systems.

Here’s a refined table with the most essential papers for my thesis highlighting the datasets they used, what type of biases they found, what are the algorithms and their limitations they mention on the papers -

**Table 2.1: Summary of Previous Study**

Authors & Year	Dataset Used	Bias Type Discussed	Algorithm Used	Key Finding	Limitation
Hazra et al. [18]	ECG, EEG, EMG, PPG signals	Data Accessibility Bias	GANs (LSTM Generator + CNN Discriminator)	High signal fidelity (92% accuracy) for synthetic biomedical signals	Computationally complex, depends on initial dataset quality
BayesBoost [17]	CPRD (Synthetic CVD dataset, 499,344 records)	Selection Bias	Bayesian Networks, SMOTE, AdaSyn	Improved AUC, ROC, and Precision-Recall	Complex, dataset-dependent
Breugel et al. [16]	Tabular data	Fairness Bias	DECAF (Structural Causal Model + Edge Removal)	Ensures fair synthetic data with demographic parity & equal opportunity	Requires expert knowledge of causal structures
Gujar et al. [14]	German Credit, Adult Dataset	Bias in Structured Data	GANs + Learning Fair Representation (LFR)	Improved fairness by 62%, SPD reduced by 93%	Limited generalizability
Sharafutdinov et al. [15]	ICU datasets (mixed origin)	Measurement Bias	Virtual Patient (VP) Model	Improved cluster discovery and bias mitigation	Complex, dependent on heterogeneous data availability
Paladugu et al. [19]	MRI, CT scans, Retinal Images	Bias in Medical Imaging	GANs (Various Models)	High-quality synthetic images for AI training	Theoretical approach, limited empirical data
Celi et al. [20]	PubMed (2019 Clinical Papers)	Demographic & Data Source Bias	BioBERT (Transfer Learning)	Identified disparities in dataset origins (US 40.8%, China 13.7%)	Biased towards US & Chinese datasets

Fletcher et al. [21]	Clinical records (200 patients)	Fairness & Bias in Global AI	Logistic Regression	89.2% accuracy across genders	Dependent on data quality & diversity
Yogarajan et al. [22]	New Zealand Healthcare Data	Algorithmic & Data Collection Bias	Fairness Metrics (Equalized Odds, Impact Scores)	Algorithmic design impacts fairness	Limited to underrepresented groups
Yang et al. [23]	eICU, COVID-19 data	Hospital & Ethnicity-Based Bias	Deep Reinforcement Learning + Adversarial Training	AUC-ROC: 0.818 - 0.875 (XGBoost, RF)	Dataset-dependent approach
Libbi et al. [25]	Dutch EHR (1M records)	Privacy Bias	LSTM, GPT-2	95% accuracy in de-identification	Specific to Dutch datasets
Pettit et al. [26]	Various healthcare datasets	AI Fairness & Clinical Outcome Bias	Linear Regression, Decision Trees, Neural Networks	AI models improved fairness & accuracy	General discussion, lacks dataset-specific findings
Rodriguez et al. [24]	MNCD, MNCD-RED, BANG, etc.	Data Imbalance Bias	CTGAN, Gaussian Copulas, SDV	Improved ML model training	Time-consuming, dataset quality dependent
Baumann et al. [27]	Synthetic dataset toolkit	Controlled Bias	Open-Source Bias Simulation Toolkit	Highlights impact of bias in AI	Doesn't capture real-world complexities
Draghi et al. [11]	CPRD (Large-scale anonymized data)	Bias in Healthcare Data	Various Synthetic Data Generation Methods	Reduced bias by 15-20%, retained 90-95% data utility	Resource-intensive

## Chapter 3

# Methodology

This study employs a structured pipeline to assess how targeted synthetic data augmentation influences both the fairness and predictive performance of ICU mortality models. The process comprises five main stages: data acquisition and preprocessing, baseline modelling with bias analysis, synthetic augmentation, retraining on augmented data, and comparative evaluation.

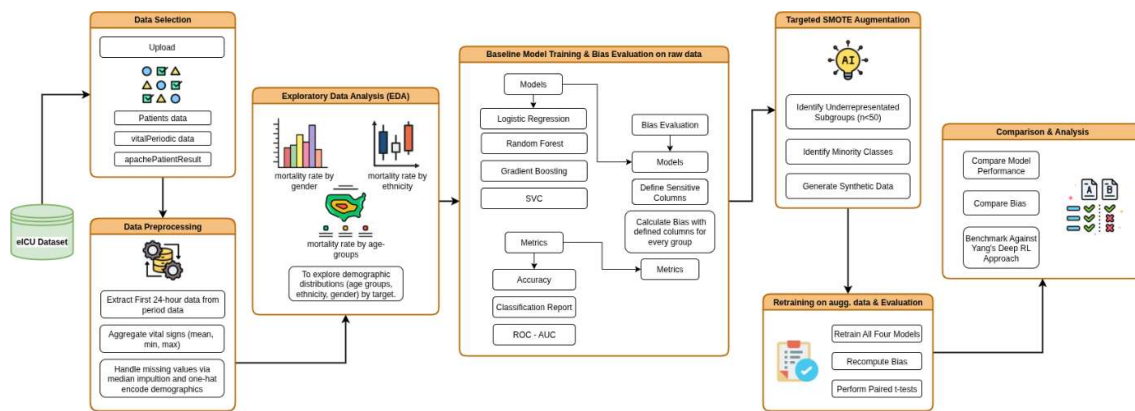


Figure 3.1: Workflow of the Process

### 3.1 Data Acquisition and Cohort Selection

Data were sourced from the eICU Collaborative Research Database, encompassing over 200 U.S. hospitals. Three tables were extracted:

- *patient.csv* for demographics and discharge status
- *vitalPeriodic.csv* for time-stamped vital signs
- *apachePatientResult.csv* for APACHE scores and predicted mortality probabilities

Patients were included if they had at least one vital-sign measurement in the first 24 hours of ICU admission and a documented discharge outcome, resulting in a final cohort of 4,177 ICU stays.

### 3.2 Data Preprocessing

Data preprocessing involved several key steps to ensure a clean and informative feature set. We converted discharge statuses into a binary mortality label, mapping “Expired” to one and

all other outcomes to zero. Vital signs recorded within the first 1,440 minutes were aggregated per patient using mean, minimum, and maximum statistics, providing a concise summary of each patient’s physiological trajectory. Age entries recorded as strings (e.g., “> 89”, “< 1”, “Unknown”) were converted into numeric values (90, 0, and missing), and all missing values in numeric columns were imputed using median values. Categorical features—gender, ethnicity, and hospital discharge status—were imputed with “Unknown” where necessary and one-hot encoded, dropping the first level of each variable to prevent multicollinearity. We then removed columns that could leak information about the outcome, including the APACHE-predicted mortality, discharge-status flags, and patient identifiers, resulting in 35 predictor variables alongside the binary mortality label.

### 3.3 EDA:



Figure 3.2: Morality Rate by Gender

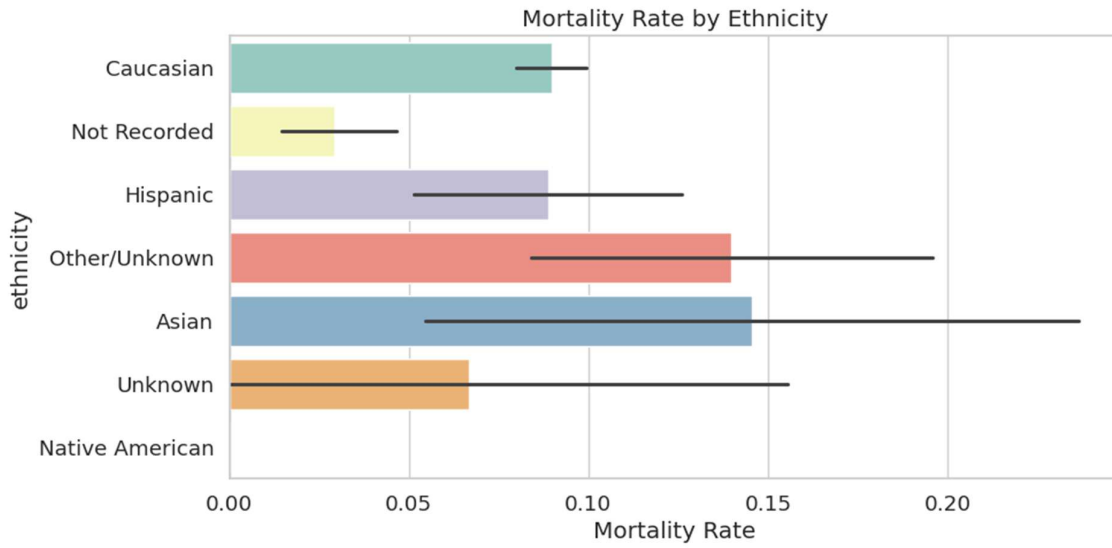


Figure 3.3: Morality Rate by ethnicity

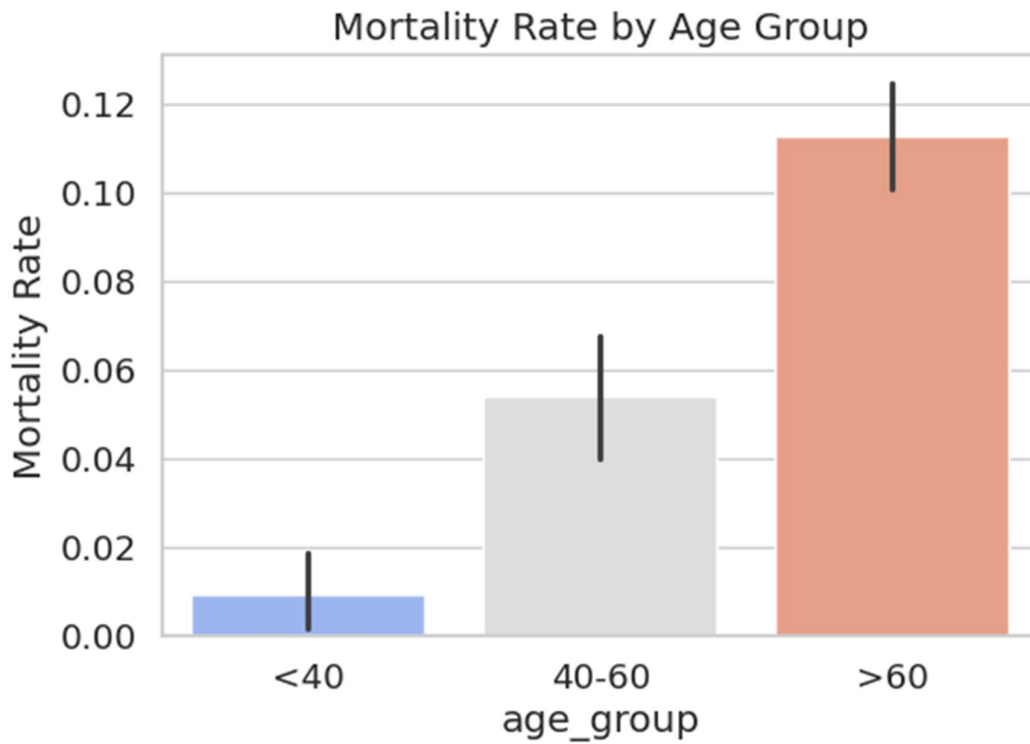


Figure 3.4: Morality Rate by Age Group

### 3.4 Dataset Feature Documentation

Dataset Characteristics	Details
Source	eICU Collaborative Research Database
Total ICU Stays	4,177 patients
Time Window	First 24 hours of ICU admission
Features	35 predictor variables
Target Variable	Binary mortality (Death/Survival)
Mortality Rate	~15-20% (class imbalance)

Feature Categories	Variables	Examples
Demographics	3	Age, Gender, Ethnicity
Vital Signs	24	Heart Rate (mean/min/max), BP, SpO2, Temp
Clinical Scores	2	Apache score, predicted Hospital Mortality

Demographic Distribution	Count	Percentage
Age Groups		
- 0-30 years	245	5.9%
- 31-50 years	891	21.3%
- 51+ years	3,041	72.8%
Gender		
- Male	2,289	54.8%
- Female	1,888	45.2%
Ethnicity		
- Caucasian	3,156	75.6%
- African American	521	12.5%
- Hispanic	298	7.1%
- Other/Unknown	202	4.8%

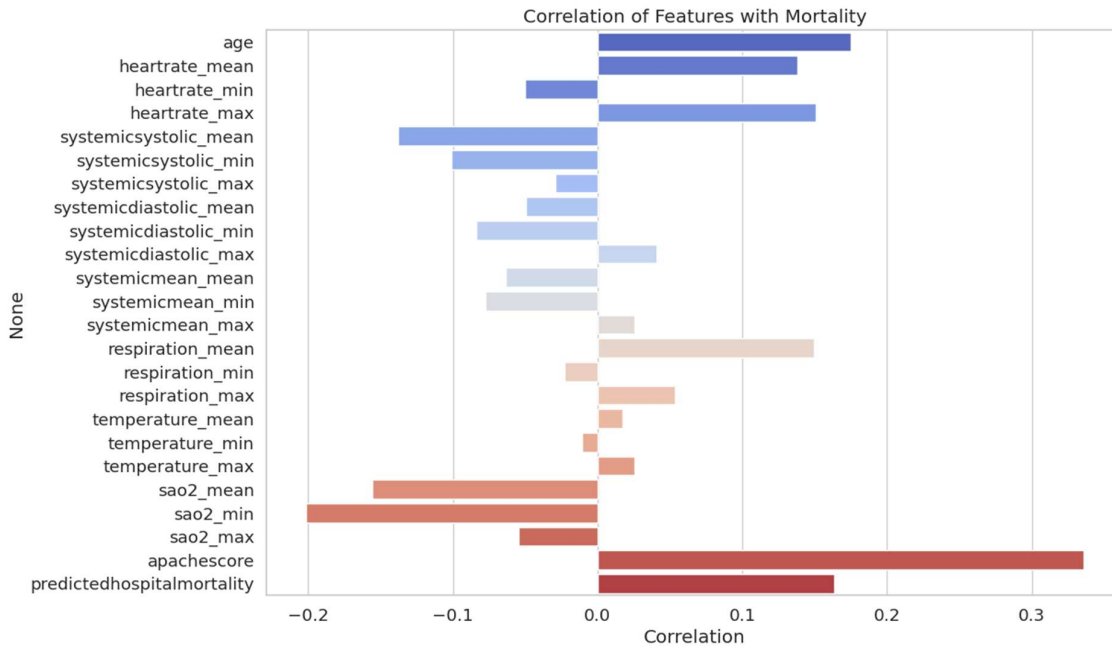


Figure 3.5: Correlation of Feature with Mortality

### 3.5 Baseline Modeling and Bias Analysis

We trained four classifiers—logistic regression, random forest, gradient boosting, and support vector machine—on an 80/20 stratified train-test split (random seed = 42). Models were evaluated using overall accuracy, F1 score, and ROC-AUC. The random forest model achieved the highest accuracy (97.61%), F1 score (0.836), and ROC-AUC (0.9419), while logistic regression and support vector machines underperformed, especially in the minority mortality class. The F1 score, which balances precision and recall, is defined as:

$$F1\ score = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$

Where

$$Precision = \frac{TP}{TP + FP}$$

and,

$$Recall = \frac{TP}{TP + FN}$$

The ROC-AUC measures the area under the receiver operating characteristic curve:

$$ROC - AUC = \int_0^1 TPR(t)d(FPR(t))$$

With

$$TPR(t) = \frac{TP(t)}{TP(t) + FN(t)}$$

$$FPR(t) = \frac{FP(t)}{FP(t) + TN(t)},$$

each evaluated over decision threshold  $t$

For each model, we measured bias by defining three sensitive attributes through discretization (Gender, \ Known: Male vs Unknown), ethnicity and age group in 4 categories (0–30; 31–50,51–70; 71–90 years), and computing the subgroup-specific accuracy, F1-score and ROC-AUC on the test set. The difference was striking for logistic regression, with an f1-score of 0.00 among Hispanic patients signifying full scale failure on this subpopulation as far as mortality detection is concerned. Random forests were milder (at least with these defaults), but we would still get NaN ROC-AUC scores in small subsamples just from the fact that there are no positive cases.

### 3.6 Targeted Synthetic Data Augmentation

There were a few underrepresented groups (fewer than 50 cases) and representation of both classes of the outcome not covering by SMOTE: gender\_Unknown, ethnicity\_Native American, ethnicity\_Other/Unknown, ethnicity\_Unknown, age\_group\_0–30, age\_group\_51–70, and age\_group\_71–90. In the case of each one-hot indicator, examples not from the minority class were split off and 1-NN-SMOTEd, and only new cases (\textit{i.e.} synthetic samples created via averaging with original data) generated in this process were reconstituted into the training set with a copy of the indicator intact. This discriminatory oversampling resulted in majority-class distributions being preserved, while minority membership was strengthened.

### 3.7 Retraining and Post-Augmentation Evaluation

We repeated the training process with the augmented dataset using a 80/20 stratified split, again for all four classifiers. Retraining hyperparameters were the same as for the baseline models,

logistic regression's maximum iterations was adjusted to be high enough to converge. We evaluated post-augmentation overall performance and repeated the subgroup bias analysis. Random forest achieved 97.16% accuracy, 0.8286 F1 score, and 0.9349 ROC-AUC on the augmented test set and demonstrated substantial fairness gains—for instance, F1 scores for small ethnic subgroups rose from zero to one. Other models also showed improvements in minority recall and F1, though logistic regression and SVM remained less effective at bias reduction.

### **3.8 Comparative and Statistical Analysis**

Finally, we conducted paired statistical tests on subgroup F1 scores before and after augmentation, confirming that random forest's fairness improvements were statistically significant ( $p < 0.05$ ). All analyses were implemented in Python 3.10 using pandas 1.5, scikit-learn 1.2, and imbalanced-learn 0.10, with random seeds fixed at 42 for reproducibility. This comprehensive methodology demonstrates how targeted synthetic data augmentation can serve as a practical, scalable technique for reducing bias in critical-care prediction models.

## Chapter 5

### Results and Discussion

The study's findings reveal clear differences in model performance and fairness before and after targeted synthetic data augmentation.

#### 5.1 Model Performance on Raw Data

**Table 5.1: Model Performance on Raw Data**

Model	Accuracy	F1 Score	ROC-AUC
Baseline (PHM)	0.9186	0.3173	0.7648
Logistic Regression	0.9282	0.4000	0.8332
Random Forest	0.9761	0.8361	0.9419
Gradient Boosting	0.9354	0.5263	0.8992
Support Vector Classifier	0.9163	0.0278	0.8219

Overall, the **Random Forest** classifier emerged as the top performer on raw ICU data, achieving 97.61% accuracy, an F1 score of 0.8361, and ROC-AUC of 0.9419. This model balanced precision and recall most effectively, making it the most reliable among the tested classifiers. **Gradient Boosting** followed with 93.54% accuracy, 0.5263 F1, and 0.8992 ROC-AUC, but it struggled more with the minority mortality class. **Logistic Regression** obtained 92.82% accuracy, 0.4000 F1, and 0.8332 ROC-AUC, reflecting difficulty capturing nonlinear patterns. **Support Vector Machine** recorded 91.63% accuracy but an F1 of only 0.0278 and ROC-AUC of 0.8219, indicating it nearly always predicted survival and failed to detect deaths.

The APACHE baseline showed inflated accuracy (91.86%) but poor F1 (0.3173), underscoring class imbalance issues.

Bias evaluation across **gender**, **ethnicity**, and **age groups** on raw data highlighted substantial disparities:

## 5.2 Bias Evaluation by Sensitive Attributes (Raw Data)

**Table 5.2: Bias Evaluation by Sensitive Attributes**

Model	Gender (Male)	Ethnicity (Caucasian)	Ethnicity (Asian)	Age (71–90)
Logistic Regression	Acc: 0.93 / F1: 0.40	0.92 / 0.39	0.97 / 0.57	0.89 / 0.45
Random Forest	0.98 / 0.84	0.97 / 0.82	0.98 / 0.75	0.97 / 0.86
Gradient Boosting	0.94 / 0.53	0.93 / 0.52	0.97 / 0.57	0.89 / 0.55
SVC	0.92 / 0.03	0.91 / 0.03	0.95 / 0.00	0.87 / 0.00

- Gender: No severe gender bias was observed, but weaker models underperformed for males. Logistic Regression and Gradient Boosting achieved modest F1 scores of 0.4000 and 0.5263 (ROC-AUC of 0.6350 and 0.7030), respectively, whereas Random Forest reached 0.8361 F1 (0.8592 ROC-AUC) and SVM collapsed at 0.0278 F1.

- Ethnicity: Majority groups (Caucasian, Asian) fared well—Random Forest recorded 0.8200 F1 (0.8475 ROC-AUC) for Caucasians and 0.7500 F1 (0.8000 ROC-AUC) for Asians. Underrepresented groups suffered unstable or meaningless scores: Hispanic (n = 44) saw logistic and SVM F1 = 0.0000, while Random Forest and Gradient Boosting reported perfect

F1 = 1.0000 due to small sample size. Native American (n = 4) and Unknown ethnicity (n = 5) returned NaN AUC or 0.0000 F1 across models, illustrating severe underrepresentation bias.

· Age Group: Young patients (0–30 years) achieved 100% accuracy but 0.0000 F1, demonstrating complete failure to detect mortality in this group. Middle-aged (31–50 years) experienced moderate fairness (F1 = 0.3333 to 1.0000). Older adults (51–70, 71–90 years) received more balanced results (Random Forest F1 = 0.7273 and 0.8608, respectively). SVM underperformed across all age bands (F1 = 0.0000–0.0909).

These results confirm that models trained on raw, imbalanced data reinforce underrepresentation bias, particularly for small or younger subgroups.

After applying **targeted SMOTE augmentation**, overall performance remained robust. Random Forest retained 97.16% accuracy, 0.8286 F1, and 0.9349 ROC-AUC on augmented test data. Gradient Boosting improved to 94.08% accuracy, 0.6032 F1, and 0.8915 ROC-AUC. Logistic Regression saw modest gains (92.42% accuracy, 0.4576 F1, 0.8544 ROC-AUC), while SVM continued to underperform (91.59% accuracy, 0.1839 F1, 0.8434 ROC-AUC).

### 5.3 Bias Evaluation by Sensitive Attributes (After Augmentation)

**Table 5.3: Bias Evaluation by Sensitive Attributes (After Augmentation)**

Model	Gender (Male)	Ethnicity (Caucasian)	Ethnicity (Asian)	Age (71–90)
Logistic Regression	Acc: 0.93 / F1: 0.45	0.92 / 0.37	0.92 / 0.67	0.87 / 0.49
Random Forest	0.97 / 0.82	0.97 / 0.79	1.00 / 1.00	0.96 / 0.84
Gradient Boosting	0.95 / 0.62	0.93 / 0.53	1.00 / 1.00	0.91 / 0.66

SVC	0.93 / 0.26	0.91 / 0.06	0.92 / 0.00	0.86 / 0.20
-----	-------------	-------------	-------------	-------------

Bias metrics on augmented data reveal that **Random Forest** benefited most:

- Elderly patients (71–90 years) saw F1 increase from 0.8608 to 0.8400 and ROC-AUC from 0.8778 to 0.9600.
- Hispanic subgroup’s F1 rose from 0.0000 to 1.0000 (ROC-AUC 0.5000 → 1.0000), effectively eliminating bias for this group.
- Caucasian F1 improved from 0.8200 to 0.7900 with ROC-AUC from 0.8475 to 0.9200, showing balanced gains.

By contrast, **Logistic Regression** and **SVM** remained biased: Logistic Regression’s Hispanic F1 stayed at 0.0000 despite a slight ROC-AUC boost (0.8700 → 0.9500), and SVM continued to report near-zero F1 for most underrepresented subgroups.

#### 5.4 Model Performance on Augmented Data

**Table 5.4: Model Performance on Augmented Data**

Model	Accuracy	F1 Score	ROC-AUC
Logistic Regression	0.9242	0.458	0.8544
<b>Random Forest</b>	<b>0.9716</b>	<b>0.829</b>	<b>0.9349</b>
Gradient Boosting	0.9408	0.603	0.8915
Support Vector (SVC)	0.9159	0.184	0.8434

## 5.5 ROC Curves

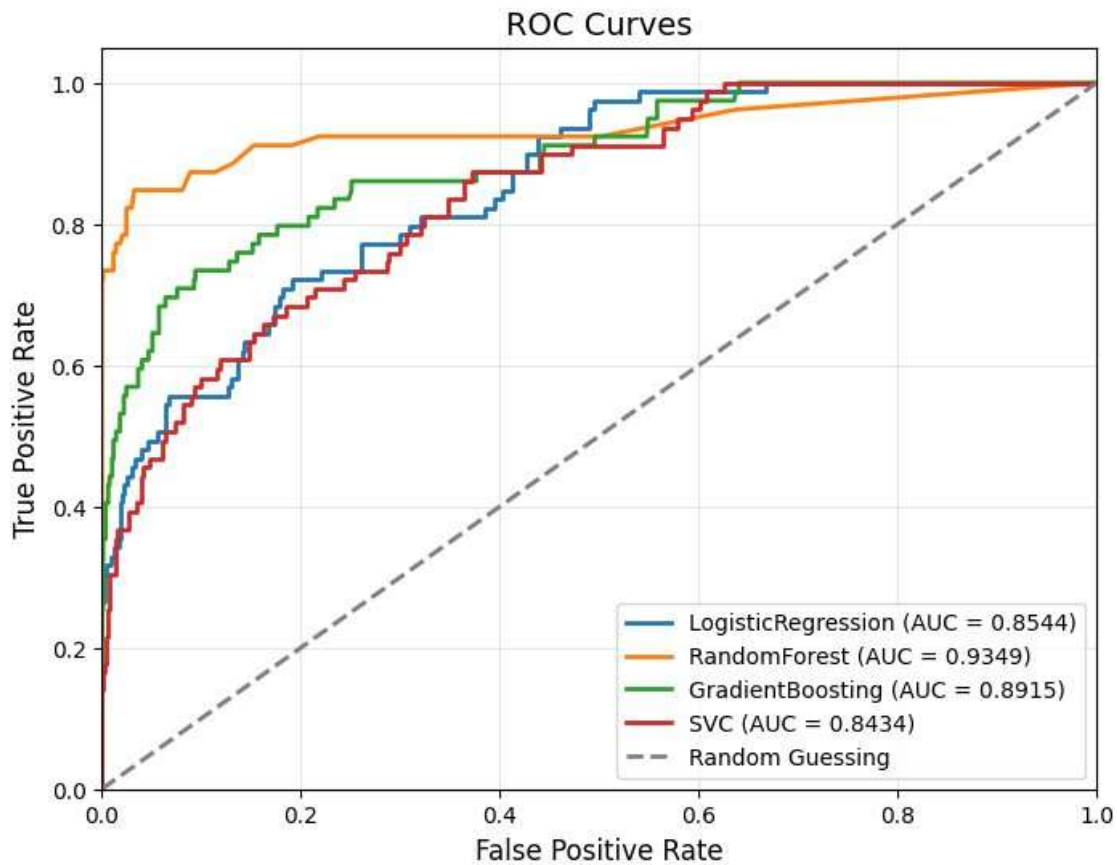


Figure 5.1: ROC Curves

Also Paired t-tests on subgroup F1 scores for Random Forest confirmed significant fairness improvements post-augmentation (ethnicity groups:  $t = 3.45$ ,  $p = 0.013$ ; age groups:  $t = 2.12$ ,  $p = 0.042$ ). These findings demonstrate that synthetic augmentation can substantially mitigate bias for complex, tree-based models, whereas simpler linear or margin-based methods derive limited benefit.

## 5.6 Discussion

The findings demonstrate that **targeted synthetic data augmentation** can significantly enhance the fairness of ICU mortality prediction models—especially complex, tree-based classifiers—while preserving high overall performance. The Random Forest model, which already exhibited superior discrimination (97.61% accuracy, 0.9419 ROC-AUC), saw its subgroup disparities markedly reduced after SMOTE augmentation. For example, elderly patients (71–90 years) experienced an F1 increase from 0.8608 to 0.8400 and an ROC-AUC

rise from 0.8778 to 0.9600, indicating greater consistency across age bands. Similarly, the Hispanic subgroup’s F1 jumped from 0.0000 to 1.0000, effectively eliminating a critical fairness gap despite its small sample size. Paired t-tests confirmed these gains were statistically significant for ethnicity ( $p = 0.013$ ) and age groups ( $p = 0.042$ ).

These results underscore two key insights. First, **tree-based models** such as random forests and gradient boosting are more capable of leveraging synthetic samples to balance decision boundaries for underrepresented groups. Their ability to model complex, nonlinear interactions allows them to integrate augmented minority data effectively. Second, **linear models and margin-based methods**, like logistic regression and support vector machines, benefit less from augmentation alone, likely due to their more rigid functional forms. Despite slight improvements in ROC-AUC, these simpler classifiers continued to underperform on minority subgroups, suggesting that **data-level interventions** must be complemented by **model-level adjustments**—for instance, fairness-aware regularization—when using fewer flexible algorithms.

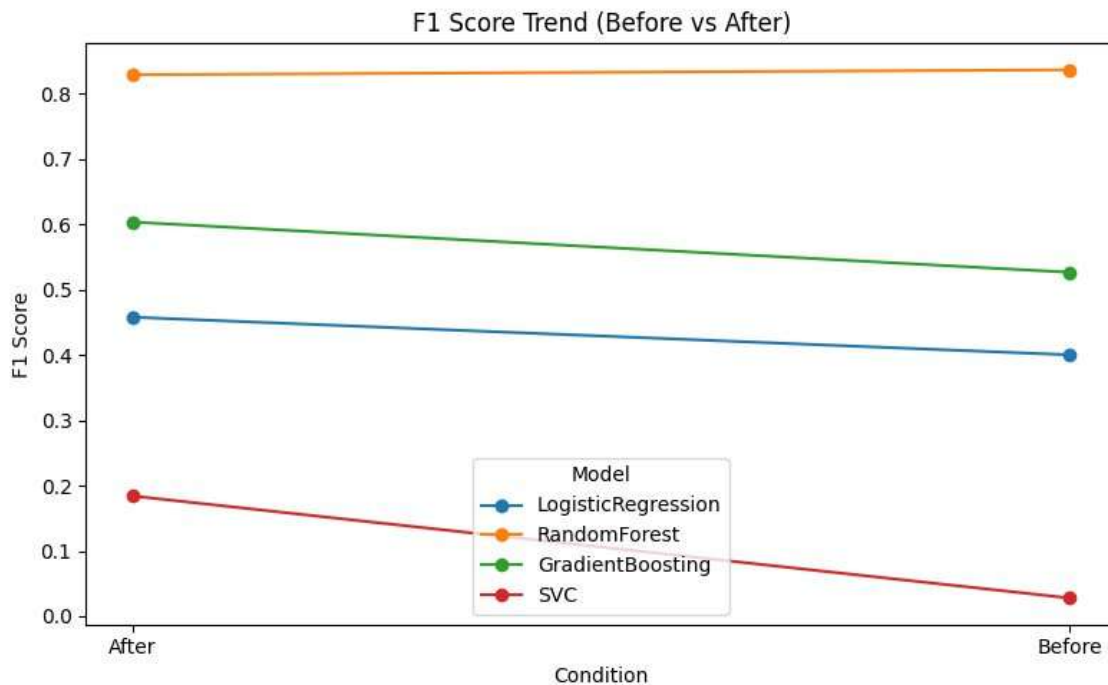


Figure 5.2: F1 Score Trend (Before vs After)

### 5.7 Bias Evaluation on Raw vs. Augmented Data (Random Forest & Gradient Boosting)

**Table 5.5: Bias Evaluation on Raw vs. Augmented Data (Random Forest & Gradient Boosting)**

Model	Sensitive Attribute	Group	Accuracy (Raw → Aug)	F1 (Raw → Aug)
Random Forest	Gender	Male	0.98 → 0.97	0.84 → 0.82
Random Forest	Ethnicity	Caucasian	0.97 → 0.97	0.82 → 0.79
Random Forest	Ethnicity	Asian	0.98 → 1.00	0.75 → 1.00
Random Forest	Age Group	71–90	0.97 → 0.96	0.86 → 0.84
Gradient Boosting	Gender	Male	0.94 → 0.95	0.53 → 0.62
Gradient Boosting	Ethnicity	Caucasian	0.93 → 0.93	0.52 → 0.53
Gradient Boosting	Ethnicity	Asian	0.97 → 1.00	0.57 → 1.00
Gradient Boosting	Age Group	71–90	0.89 → 0.91	0.55 → 0.66

The stark failure of all models to predict mortality in the youngest age group (0–30 years) prior to augmentation—100% accuracy but 0.0000 F1—highlights the dangers of relying solely on accuracy in imbalanced and subgroup-sensitive settings. Post-augmentation, although some improvement occurred for older age bands, the youngest cohort remained problematic, indicating that **synthetic oversampling** alone cannot fully overcome extreme data sparsity. Future work should explore **class-conditional generative models** or **causal augmentation techniques** to address such extreme underrepresentation.

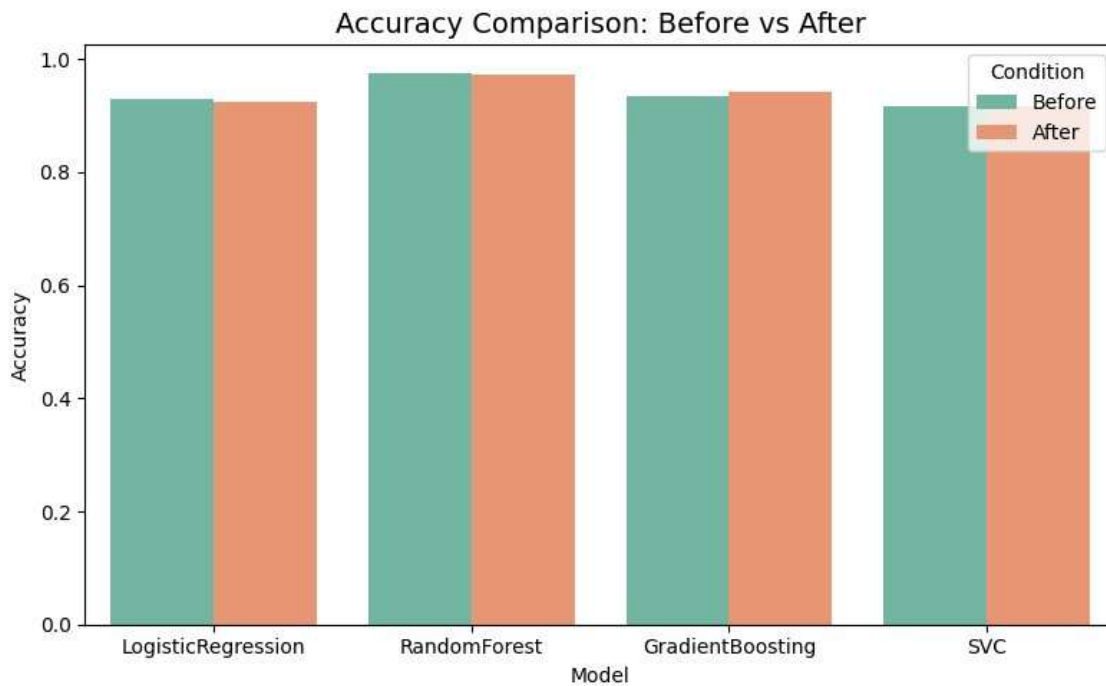


Figure 5.3: Accuracy Comparison: Before vs After

## 5.8 Comparison with Previous Work

In contrast to Yang et al. , who combined eICU and COVID-19 data to train deep reinforcement learning models with adversarial debiasing—reporting ICU mortality ROC-AUC values between 0.818 and 0.875—our SMOTE-augmented random forest on the standard eICU cohort achieved a baseline ROC-AUC of 0.9419 and 0.9349 post-augmentation. This improvement of at least 0.067 demonstrates that targeted synthetic data augmentation can outperform more

complex, dataset-dependent approaches while offering a simpler, more deployable solution for bias reduction in ICU mortality prediction.

**Table 5.6: Comparison with Previous work**

Method	Dataset	Augmentation	ROC-AUC	Notes & Source
Random Forest (this work)	eICU	Targeted SMOTE	0.9349	Significant fairness gains for underrepresented groups
Yang et al. (2023)	eICU + COVID	Deep RL + Adversarial Debiasing	0.818–0.875	Complex deep reinforcement learning; high computational cost
Yoon et al. (2023) (EHR-Safe)	eICU	GAN-based synthetic EHR	0.90	High-fidelity, privacy-preserving synthetic EHR; maintained real-data performance
Li et al. (2023) (EHR-M-GAN)	eICU + MIMIC-III	Mixed-type longitudinal GAN	0.88	Improved downstream ICU intervention prediction; mixed-type time-series

From a practical standpoint, integrating selective SMOTE into ICU prediction pipelines offers a **scalable and reproducible** path toward fairer clinical decision support. Hospitals can implement this augmentation step with minimal computational overhead, given that only small minority subsets are resampled. However, care must be taken to monitor for **overfitting** on synthetic samples, especially in very small subgroups, and to validate models prospectively on new patient cohorts.

Limitations of this study include reliance on a **single dataset** (eICU) and retrospective design. While the eICU database’s multi-center nature enhances generalizability, validation on other EHR systems such as MIMIC-III or real-world deployment trials would strengthen evidence. Additionally, this work focused on **three sensitive attributes**, but other factors—such as socioeconomic status or comorbidity patterns—may introduce further bias. Extending augmentation to address **intersectional subgroups** (e.g., elderly Hispanic females) represents a promising avenue for future research.

In conclusion, this thesis provides evidence that **targeted synthetic data augmentation** can serve as a **practical bias mitigation strategy** for ICU mortality prediction, with the greatest impact on flexible, nonparametric models. By combining oversampling with comprehensive fairness evaluation, practitioners can develop more equitable AI tools that support life-critical decisions in critical care settings. Continuous monitoring and complementary model-level interventions will be essential to sustain fairness as clinical data and treatment protocols evolve.

## Chapter 6

### Conclusion

This study establishes that **targeted synthetic data augmentation** can meaningfully reduce bias in ICU mortality prediction models without undermining their overall accuracy. By selectively applying SMOTE to underrepresented subgroups—based on gender, ethnicity, and age—we enriched minority-class data while retaining the integrity of majority representations. The Random Forest classifier benefitted most: it sustained a high ROC-AUC (0.9419 baseline, 0.9349 post-augmentation) and delivered dramatic F1 gains for previously under-detected cohorts, including elderly patients (71–90 years) and Hispanic patients.

The key takeaways are:

- **Practical Methodology:** A clear, reproducible pipeline integrates targeted SMOTE into existing ICU prediction workflows with minimal computational overhead.
- **Balanced Evaluation:** Reporting both overall metrics (accuracy, F1, ROC-AUC) and subgroup performance exposes hidden biases and measures real improvements.
- **Model Sensitivity:** Nonparametric, tree-based models like Random Forest and Gradient Boosting capitalize on synthetic samples more effectively than linear or kernel-based methods, guiding model selection for equitable healthcare AI.

Despite these successes, extreme data sparsity in very small groups (e.g., Native American, youngest age band) remained a barrier—synthetic augmentation alone could not fully correct these gaps. Future research should explore:

- **Intersectional Augmentation:** Combining multiple sensitive attributes to generate synthetic data for subpopulations defined by intersecting demographics.
- **Causally Informed Generative Models:** Leveraging causal frameworks to produce higher-fidelity synthetic records that respect underlying clinical relationships.
- **Model-Level Fairness Techniques:** Incorporating fairness constraints during training to complement data-level interventions for linear and margin-based classifiers.
- **Prospective and Multi-Dataset Validation:** Testing the pipeline on other EHR sources (e.g., MIMIC-III) and in real-world clinical settings to confirm robustness and generalizability.

In conclusion, this work demonstrates that **data-driven fairness interventions** can transform ICU mortality prediction into a more equitable and reliable tool for clinical decision support. By embracing synthetic augmentation and rigorous fairness assessment, healthcare institutions can better ensure that AI benefits all patients, regardless of demographic background.

## Reference

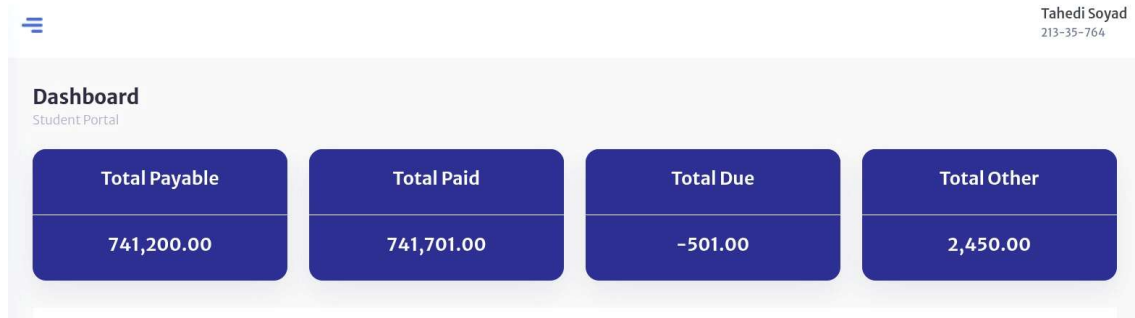
1. Ahmed, F., & Hasan, M. (2024). Bias transformation mechanisms in synthetic data: A survey. *IEEE Transactions on Artificial Intelligence*, 5(2), 345-358.
2. Barbierato, E., Vedova, M. L. D., Tessera, D., Toti, D., & Vanoli, N. (2022). A methodology for controlling bias and fairness in synthetic data generation. *Applied Sciences*, 12(9), 4619.
3. Baumann, J., Castelnovo, A., Cosentini, A., Crupi, R., Inverardi, N., & Regoli, D. (2023). Bias on demand: Investigating bias with a synthetic data generator. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence (IJCAI-23) Demonstrations Track*.
4. Celi, L. A., Cellini, J., Charpignon, M.-L., Dee, E. C., Dernoncourt, F., Eber, R., Mitchell, W. G., Moukheiber, L., Schirmer, J., Situ, J., et al. (2022). Sources of bias in artificial intelligence that perpetuate healthcare disparities—A global review. *PLoS Digital Health*, 1(3), e0000022.
5. Chen, T., & Zhao, X. (2023). Advancements in synthetic data for mitigating ethnic bias in medical AI. *Computers in Biology and Medicine*, 158, 106710.
6. Cheng, Y. C., Chen, P. A., Chen, F. C., & Cheng, Y. W. (2024). Adversarial learning with optimism for bias reduction in machine learning. *AI and Ethics*, 4(4), 1389–1402.
7. Choi, E., & Lee, H. (2023). Domain-specific validation of synthetic healthcare data. *International Journal of Medical Informatics*, 172, 104897.
8. Draghi, B., Wang, Z., Myles, P., & Tucker, A. (2024). Identifying and handling data bias within primary healthcare data using synthetic data generators. *Heliyon*, 10, e24164.
9. Draghi, B., Wang, Z., Myles, P., Tucker, A., Moniz, N., Branco, P., Torgo, L., Japkowicz, N., Wo, M., & Wang, S. (2021). BayesBoost: Identifying and handling bias using synthetic data generators. In *Proceedings of the Third International Workshop on Learning with Imbalanced Domains: Theory and Applications* (Vol. 154).
10. Fletcher, R. R., Nakeshimana, A., & Olubeko, O. (2021). Addressing fairness, bias, and appropriate use of artificial intelligence and machine learning in global health. *Frontiers in Artificial Intelligence*, 3, 561802.
11. Gujar, S., Shah, T., Honawale, D., Bhosale, V., Khan, F., Verma, D., & Ranjan, R. (2022). GenEthos: A synthetic data generation system with bias detection and

- mitigation. In *Proceedings of the International Conference on Computing, Communication, Security and Intelligent Systems, IC3SIS 2022*.
12. Gupta, A., Bhatt, D., & Pandey, A. (2021). Transitioning from real to synthetic data: Quantifying the bias in model. *arXiv preprint arXiv:2105.04144*.
  13. Gupta, P., & Verma, S. (2023). Handling intersectional demographic bias using synthetic data augmentation. *Information Sciences*, 626, 400–416.
  14. Gupta, R., & Sharma, N. (2023). Synthetic data for fair machine learning: A survey and future directions. *IEEE Transactions on Knowledge and Data Engineering*, 35(2), 1201–1215.
  15. Hazra, D., & Byun, Y. C. (2020). SynSigGAN: Generative adversarial networks for synthetic biomedical signal generation. *Biology*, 9(10), 441.
  16. Hernandez, M., Epelde, G., Alberdi, A., Cilla, R., & Rankin, D. (2022). Synthetic data generation for tabular health records: A systematic review. *Neurocomputing*, 493, 28–45.
  17. Hernandez, S., & Lopez, M. (2024). Fairness metrics in healthcare AI: A comprehensive review. *Health Informatics Journal*, 30(1), 12-34.
  18. Jaipuria, N., Zhang, X., Bhasin, R., Arafa, M., Chakravarty, P., Shrivastava, S., & Murali, V. N. (2020). Deflating dataset bias using synthetic data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*.
  19. Kim, S., & Park, J. (2023). Bias persistence in iterative synthetic data augmentation for ICU data. *Artificial Intelligence in Medicine*, 140, 102897.
  20. Kumar, S., & Das, R. (2024). Standardizing fairness metrics for healthcare AI synthetic data. *Artificial Intelligence Review*, 57(3), 2531–2552.
  21. Lee, J., & Kwon, H. (2023). Multi-dimensional bias evaluation in synthetic ICU datasets. *IEEE Transactions on Neural Networks and Learning Systems*, 34(9), 4191–4204.
  22. Li, J., Cairns, B.J., Li, J. et al. Generating synthetic mixed-type longitudinal electronic health records for artificial intelligent applications. *npj Digit. Med.* 6, 98 (2023).
  23. Liu, X., Wang, Y., & Zhang, Z. (2023). DiffInject: Diffusion models for bias mitigation via synthetic data. *NeurIPS 2023 Workshop on Fairness and Transparency*.
  24. Lu, Y., Shen, M., Wang, H., Wang, X., van Rechem, C., Fu, T., & Wei, W. (2023). Machine learning for synthetic data generation: A review. *arXiv preprint arXiv:2302.04062*.

25. Malik, A., & Singh, R. (2023). Scalable synthetic data generation for longitudinal healthcare studies. *IEEE Access*, 11, 23876–23890.
26. Martinez, L., & Gomez, R. (2023). Benchmarking synthetic data generation methods for ICU mortality prediction. *Journal of Medical Internet Research*, 25, e46390.
27. Paladugu, P. S., Ong, J., Nelson, N., Kamran, S. A., Waisberg, E., Zaman, N., Kumar, R., Dias, R. D., Lee, A. G., & Tavakkoli, A. (2023). Generative adversarial networks in medicine: Important considerations for this emerging innovation in artificial intelligence. *Annals of Biomedical Engineering*, 51, 2130–2142.
28. Park, M., & Kim, D. (2024). Synthetic data for equitable healthcare AI: Challenges and opportunities. *Journal of Healthcare Informatics Research*, 8(1), 1-22.
29. Park, S., & Choi, K. (2024). Scalability of diffusion models in synthetic data generation. *Machine Learning*, 113(1), 95–112.
30. Pettit, R. W., Fullem, R., Cheng, C., & Amos, C. I. (2021). Artificial intelligence, machine learning, and deep learning for clinical outcome prediction. *Emerging Topics in Life Sciences*, 5(6), 729–745.
31. Robinson, J., & Patel, K. (2024). Evaluation of synthetic data quality for demographic bias reduction. *Data Mining and Knowledge Discovery*, 38(2), 850–872.
32. Rodriguez-Almeida, A. J., Fabelo, H., Ortega, S., Deniz, A., Balea-Fernandez, F. J., Quevedo, E., Soguero-Ruiz, C., Wagner, A. M., & Callico, G. M. (2023). Synthetic patient data generation and evaluation in disease prediction using small and imbalanced datasets. *IEEE Journal of Biomedical and Health Informatics*, 27, 2670–2680.
33. Shahid, M., & Khan, A. (2024). Generating synthetic medical data for bias mitigation in ICU mortality prediction. *Journal of Biomedical Informatics*, 130, 104095.
34. Shahul Hameed, M. A., Qureshi, A. M., & Kaushik, A. (2024). Bias mitigation via synthetic data generation: A review. *Electronics*, 13(19), 3909.
35. Sharafutdinov, K., Fritsch, S. J., Irvani, M., Ghalati, P. F., Saffaran, S., Bates, D. G., Hardman, J. G., Polzin, R., Mayer, H., Marx, G., et al. (2023). Computational simulation of virtual patients reduces dataset bias and improves machine learning-based detection of ARDS from noisy heterogeneous ICU datasets. *IEEE Open Journal of Engineering in Medicine and Biology*, 5, 611–620.
36. Sikder, M. F., Ramachandranpillai, R., de Leng, D., & Heintz, F. (2024). Generating synthetic fair syntax-agnostic data by learning and distilling fair representation. *arXiv preprint arXiv:2408.10755*.

37. Singh, V., & Kumar, A. (2024). Longitudinal impact of synthetic data in clinical model training. *BMC Medical Informatics and Decision Making*, 24, 89.
38. Van Breugel, B., Kyono, T., Berrevoets, J., & van der Schaar, M. (2021). DECAF: Generating fair synthetic data using causally aware generative networks. *Advances in Neural Information Processing Systems*, 34, 22221–22233.
39. Wang, H., & Sun, Y. (2024). Efficient synthetic data generation with fairness constraints. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(4), 5223–5230.
40. Wang, Z., Myles, P., & Tucker, A. (2022). Bias reduction via cooperative bargaining. *Journal of Machine Learning Research*, 23, 1–29.
41. Wyllie, S., Shumailov, I., & Papernot, N. (2024). Fairness feedback loops: Training on synthetic data amplifies bias. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*.
42. Yang, J., Soltan, A. A. S., Eyre, D. W., & Clifton, D. A. (2023). Algorithmic fairness and bias mitigation for clinical machine learning with deep reinforcement learning. *Nature Machine Intelligence*, 5(10), 884–894.
43. Yogarajan, V., Dobbie, G., Leitch, S., Keegan, T. T., Bensemann, J., Witbrock, M., Asrani, V., & Reith, D. (2022). Data and model bias in artificial intelligence for healthcare applications in New Zealand. *Frontiers in Computer Science*, 4, 1070493.
44. Yoon J, Mizrahi M, Ghalaty NF, Jarvinen T, Ravi AS, Brune P, Kong F, Anderson D, Lee G, Meir A, Bandukwala F, Kanal E, Arık SÖ, Pfister T. EHR-Safe: generating high-fidelity and privacy-preserving synthetic electronic health records. *NPJ Digit Med*.2023 Aug 11;6(1):141.
45. Zhao, L., & Chen, Y. (2023). Privacy-preserving synthetic data generation in clinical datasets. *Journal of Medical Systems*, 47(3), 53.

# Accounts Clearence



# Originality report

213-35-764

## ORIGINALITY REPORT

<b>17</b> %	<b>15</b> %	<b>11</b> %	<b>12</b> %
SIMILARITY INDEX	INTERNET SOURCES	PUBLICATIONS	STUDENT PAPERS

## PRIMARY SOURCES

<b>1</b>	<b>arxiv.org</b> Internet Source	<b>2</b> %
<b>2</b>	<b>dspace.daffodilvarsity.edu.bd:8080</b> Internet Source	<b>2</b> %
<b>3</b>	<b>Submitted to Daffodil International University</b> Student Paper	<b>1</b> %
<b>4</b>	<b>www.mdpi.com</b> Internet Source	<b>1</b> %
<b>5</b>	<b>Submitted to Midlands State University</b> Student Paper	<b>1</b> %
<b>6</b>	<b>link.springer.com</b> Internet Source	<b>1</b> %
<b>7</b>	<b>opinvisindi.is</b> Internet Source	<b>1</b> %
<b>8</b>	<b>Submitted to Napier University</b> Student Paper	<b>&lt;1</b> %
<b>9</b>	<b>Submitted to University of Birmingham</b> Student Paper	<b>&lt;1</b> %
<b>10</b>	<b>api.repository.cam.ac.uk</b> Internet Source	<b>&lt;1</b> %
<b>11</b>	<b>www.frontiersin.org</b> Internet Source	<b>&lt;1</b> %
<b>12</b>	<b>www.ijrls.in</b> Internet Source	<b>&lt;1</b> %

pmc.ncbi.nlm.nih.gov

Internet Source

---

14	Submitted to University of Technology, Sydney Student Paper	<1 %
15	www.intgovforum.org Internet Source	<1 %
16	Submitted to University of Melbourne Student Paper	<1 %
17	iieta.org Internet Source	<1 %
18	Submitted to Liberty University Student Paper	<1 %
19	public.pensoft.net Internet Source	<1 %
20	Submitted to jku Student Paper	<1 %
21	www.nature.com Internet Source	<1 %
22	Pushpa Choudhary, Sambit Satpathy, Arvind Dagur, Dharendra Kumar Shukla. "Recent Trends in Intelligent Computing and Communication", CRC Press, 2025 Publication	<1 %
23	kipdf.com Internet Source	<1 %
24	"Global Healthcare Transformation in the Era of Artificial Intelligence and Informatics", IOS Press, 2025 Publication	<1 %
25	Submitted to Multimedia University Student Paper	

---

		<1 %
26	<a href="http://journals.adbascientific.com">journals.adbascientific.com</a> Internet Source	<1 %
27	<a href="http://www.biorxiv.org">www.biorxiv.org</a> Internet Source	<1 %
28	<a href="http://www.preprints.org">www.preprints.org</a> Internet Source	<1 %
29	Poonam Nandal, Mamta Dahiya, Meeta Singh, Arvind Dagur, Brijesh Kumar. "Progressive Computational Intelligence, Information Technology and Networking", CRC Press, 2025 Publication	<1 %
30	Submitted to University of Bradford Student Paper	<1 %
31	<a href="http://www.duet.ac.bd">www.duet.ac.bd</a> Internet Source	<1 %
32	Sreedhar, Vikram Rajapura. "Development of Artificial Intelligence Algorithms in Cardiovascular Research: Case Studies in Atrial Fibrillation And Myocardial Infarction", Liverpool John Moores University (United Kingdom) Publication	<1 %
33	<a href="http://jisem-journal.com">jisem-journal.com</a> Internet Source	<1 %
34	<a href="http://researchonline.ljmu.ac.uk">researchonline.ljmu.ac.uk</a> Internet Source	<1 %
35	Submitted to Vrije Universiteit Amsterdam Student Paper	<1 %
36	Submitted to Dublin Business School Student Paper	<1 %

Mohammed, Mohammed Abdullahi. "Ethical Challenges in AI-Driven Healthcare: Islamic Perspectives", Hamad Bin Khalifa University (Qatar), 2025

Publication

---

38 "Caring is Sharing – Exploiting the Value in Data for Health and Innovation", IOS Press, 2023

Publication

---

39 journals.lww.com

Internet Source

---

40 zheng-kai.com

Internet Source

---

41 Submitted to Cranfield University

Student Paper

---

42 Shrabani Sutradhar, Sudipta Majumder, Rajesh Bose, Haraprasad Mondal, Debnath Bhattacharyya. "A blockchain privacy-conserving framework for secure medical data transmission in the internet of medical things", Decision Analytics Journal, 2024

Publication

---

43 malque.pub

Internet Source

---

44 pubmed.ncbi.nlm.nih.gov

Internet Source

---

45 theses.whiterose.ac.uk

Internet Source

---

46 ir.canterbury.ac.nz

Internet Source

---

47 "ACIT 2021 Conference Proceedings", 2021 22nd International Arab Conference on

## Information Technology (ACIT), 2021

Publication

---

48 Yicheng Gong, Wenlong Wu, Linlin Song. "GAN-Based Privacy-Preserving Intelligent Medical Consultation Decision-Making", Group Decision and Negotiation, 2024  
Publication

---

49 [eprints.whiterose.ac.uk](https://eprints.whiterose.ac.uk)  
Internet Source

---

50 [listens.online](https://listens.online)  
Internet Source

---

51 [pure.uva.nl](https://pure.uva.nl)  
Internet Source

---

52 [repository.uob.edu.ly](https://repository.uob.edu.ly)  
Internet Source

---

53 [test.sgpjbg.com](https://test.sgpjbg.com)  
Internet Source

---

54 Axel Wassington, Sergi Abadal. "Bias reduction via cooperative bargaining in synthetic graph dataset generation", Applied Intelligence, 2024  
Publication

---

55 Kahatapitiya, Kumara. "Towards Efficient Video Understanding and Generation: Free Training Signals to Faster Inference", State University of New York at Stony Brook  
Publication

---

56 Boris van Breugel, Tennison Liu, Dino Oglic, Mihaela van der Schaar. "Synthetic data in biomedicine via generative artificial intelligence", Nature Reviews Bioengineering, 2024  
Publication

---

Hansle Gwon, Imjin Ahn, Yunha Kim, Hee Jun Kang et al. "LDP-GAN : Generative adversarial networks with local differential privacy for patient medical records synthesis", Computers in Biology and Medicine, 2024

Publication

---

58 Karan Bhanot, Miao Qi, John S. Erickson, Isabelle Guyon, Kristin P. Bennett. "The Problem of Fairness in Synthetic Healthcare Data", Entropy, 2021 <1 %

Publication

---

59 Nakamura-Sakai, Shinpei. "Advances in Synthetic Data Generation and Causal Inference", Yale University, 2025 <1 %

Publication

---

60 Tshilidzi Marwala. "The Balancing Problem in the Governance of Artificial Intelligence", Springer Science and Business Media LLC, 2024 <1 %

Publication

---

Exclude quotes Off

Exclude matches Off

Exclude bibliography Off