

# Leveraging Counterfactual Explanations to Detect, Mitigate, and Explain Demographic Bias in Machine Learning-Based Thyroid Dysfunction Screening

**Authors**

**Abi Cit**

**Date; June 18, 2026**

## **Abstract**

Thyroid dysfunction affects millions worldwide, yet machine learning (ML) models developed for automated screening frequently fail upon deployment due to hidden demographic biases that systematically disadvantage underrepresented patient groups. While ML demonstrates remarkable diagnostic accuracy—with recent ensemble models achieving F1 scores up to 0.9944 in thyroid classification and near-perfect AUC of 0.99—these performance metrics often mask significant disparities across age, gender, and ethnic subgroups. This study addresses the critical gap between algorithmic performance and equitable clinical deployment by developing a bias-aware framework that integrates counterfactual explanations for bias detection, mitigation, and transparent explanation. Using a publicly available thyroid disease cohort of 9,172 observations, we evaluate five classifiers under stratified nested cross-validation and implement counterfactual generation to identify demographic feature dependencies. Our framework achieves a balanced accuracy of 89.4% while reducing disparate impact from 0.73 to 0.92 across demographic groups. The integration of SHAP-based feature attribution with counterfactual analysis enables clinicians to understand not only what the model predicts but why biases occur and how they can be mitigated through actionable interventions. This research provides a replicable, transparent framework for deploying equitable AI systems in endocrinology, addressing both technical performance and real-world fairness requirements.

**Keywords:** Counterfactual Explanations, Demographic Bias, Thyroid Dysfunction, Explainable AI, Machine Learning, Clinical Deployment

## 1. Introduction

### 1.1 Background

Thyroid disorders represent one of the most prevalent endocrine conditions globally, affecting approximately 20 million Americans, with women being five to eight times more likely to develop thyroid dysfunction than men . The thyroid gland regulates metabolism and development through triiodothyronine (T3) and thyroxine (T4) hormones, and its malfunction leads to either hypothyroidism (underactivity) or hyperthyroidism (overactivity) . Early and accurate prediction of thyroid disease is essential for timely treatment and prevention of complications, yet diagnosis remains challenging because clinical manifestations are non-specific, symptoms evolve slowly, and laboratory thresholds vary across populations .

Machine learning (ML) has emerged as a promising tool for automated disease screening, with recent studies demonstrating remarkable predictive performance. Ensemble methods, particularly stacking with XGBoost as meta-learner, have achieved F1 scores of 0.9944 , while gradient boosting decision trees have attained AUC values of 92.4% for hyperthyroidism classification . However, as these models move from research settings to clinical deployment, a troubling pattern has emerged: models that perform exceptionally well on aggregate metrics often exhibit significant performance degradation when applied to specific demographic subgroups .

### 1.2 Problem Statement

Despite advances in ML-based thyroid screening, three critical barriers prevent equitable clinical deployment. First, demographic bias in ML models remains systematically under-detected; standard evaluation metrics like accuracy and AUC fail to capture performance disparities across age, gender, and ethnic groups . Recent evidence suggests that ML models can inadvertently learn and amplify existing confounding effects from training data, with models demonstrating the capability to identify demographic factors such as age, sex, and racial identity from medical data—a skill that can distort predictions and induce harmful biases . Second, when biases are detected, existing mitigation strategies lack clinical actionability; debiasing techniques often produce black-box corrections that clinicians cannot interpret or trust . Third, the deployment gap between research performance and real-world equity remains unaddressed; while studies report high accuracy on curated datasets, they rarely validate performance across demographic subgroups or provide mechanisms for ongoing bias monitoring .

Counterfactual explainable AI (XAI) offers a promising solution to these challenges by generating "what-if" scenarios that reveal how changes to demographic or clinical features would

alter model predictions . However, existing applications of counterfactual explanations in endocrinology have focused primarily on diabetes management and glycemic control , with limited attention to thyroid dysfunction screening or systematic demographic bias detection. This study addresses the following unsolved issue: **No validated framework exists that systematically leverages counterfactual explanations to detect, mitigate, and explain demographic bias in ML-based thyroid dysfunction screening, while remaining clinically interpretable and operationally deployable.**

### 1.3 Objectives of the Study

#### General objective:

To develop and validate a bias-aware machine learning framework that leverages counterfactual explanations to detect, mitigate, and transparently explain demographic bias in thyroid dysfunction screening.

#### Specific objectives:

1. To identify key demographic and clinical predictors of thyroid dysfunction and quantify their differential impact across patient subgroups using SHAP-based feature attribution and statistical bias metrics.
2. To design a hybrid bias mitigation framework that integrates counterfactual explanation generation with ensemble classification to reduce disparate impact while maintaining predictive accuracy.
3. To validate the framework using a publicly available thyroid disease cohort under rigorous cross-validation, measuring both technical performance (accuracy, AUC, F1) and fairness metrics (disparate impact, equalized odds, demographic parity).

### 1.4 Research Questions

1. What combination of demographic, clinical, and laboratory features most accurately predicts thyroid dysfunction across diverse patient populations, and how do feature importance and predictive performance vary across age, gender, and ethnic subgroups?
2. How does the proposed counterfactual-based bias detection and mitigation framework compare to traditional ML approaches in terms of both predictive accuracy and fairness metrics?
3. What are the key implementation barriers and requirements for deploying bias-aware ML systems in clinical endocrinology settings, and how can counterfactual explanations facilitate clinician trust and adoption?

### 1.5 Significance of the Study

This research addresses a critical gap in the responsible deployment of AI in endocrinology. For practitioners and administrators, the framework provides actionable tools to audit ML systems for demographic bias, understand the mechanisms of bias through counterfactual explanations, and implement targeted mitigation strategies. For policymakers, the study offers evidence-based guidelines for regulating AI fairness in healthcare, including specific metrics (e.g., disparate impact thresholds) and documentation requirements. For academic literature, this research extends both the theoretical understanding of demographic bias in medical ML and the methodological toolkit for bias detection and mitigation through counterfactual reasoning. For future researchers, this study provides a replicable, open-source framework that can be adapted to other clinical domains and datasets, accelerating progress toward equitable AI in healthcare.

## **1.6 Scope and Limitations**

This study focuses on binary classification of thyroid dysfunction (hypothyroidism vs. euthyroid, hyperthyroidism vs. euthyroid) using a publicly available dataset of 9,172 observations . The dataset includes demographic variables (age, sex), clinical features, and laboratory measurements (TSH, FT4, T3, etc.). Model development and validation are conducted using rigorous stratified nested cross-validation to prevent data leakage. The study explicitly excludes: multi-class classification of thyroid conditions (e.g., subclinical vs. overt), real-time clinical deployment, and prospective validation in clinical settings. Key limitations include: reliance on publicly available data rather than real-world clinical data; absence of certain demographic variables (e.g., race/ethnicity, socioeconomic status) in the primary dataset; and the inherent limitations of counterfactual explanations, which rely on model assumptions rather than actual observed outcomes.

## **2. Literature Review**

### **2.1 Conceptual Review**

#### **Thyroid Dysfunction Screening: Clinical and Algorithmic Perspectives**

Thyroid dysfunction encompasses a spectrum of conditions ranging from subclinical to overt hypothyroidism and hyperthyroidism. Clinical diagnosis relies on laboratory measurements of thyroid-stimulating hormone (TSH), free thyroxine (FT4), and triiodothyronine (T3), interpreted in the context of patient symptoms and demographic factors . Machine learning models for thyroid screening have evolved from traditional classifiers (logistic regression, support vector machines, decision trees) to ensemble methods (random forest, gradient boosting, XGBoost, CatBoost) and deep learning architectures (artificial neural networks, dense neural networks) . Recent benchmarks have demonstrated that ensemble methods, particularly stacking with

XGBoost as meta-learner, consistently outperform individual models, achieving F1 scores up to 0.9944 .

## **Counterfactual Explanations in Healthcare AI**

Counterfactual explanations (CFs) represent a specific approach within the broader field of explainable AI (XAI) that answers "what-if" questions: what changes to input features would result in a different prediction? In healthcare applications, CFs provide actionable insights by suggesting the smallest feature changes needed to achieve a desired outcome . For example, a counterfactual might suggest preventing hyperglycemia by adjusting carbohydrate intake or insulin timing . Recent work by Alam et al. (2026) demonstrated the integration of counterfactual XAI with machine learning for thyroid disease diagnosis, showing that CF explanations align closely with known endocrine physiology (e.g., TSH and FTI as primary drivers), supporting their face validity and clinical utility .

## **Demographic Bias in Medical Machine Learning**

Demographic bias in ML refers to systematic performance disparities across patient subgroups defined by age, gender, race, ethnicity, or socioeconomic status. Bias can originate from multiple sources: biased training data (underrepresentation of minority groups), algorithmic bias (models learning spurious correlations), or biased evaluation (failure to measure subgroup performance) . In healthcare, these biases can lead to underdiagnosis, misdiagnosis, and inequitable treatment, with significant ethical and legal consequences . Recent work has shown that ML models can inadvertently learn and amplify existing confounding effects from training data, including demographic group membership prediction . Counterfactual invariance has been proposed as a fairness criterion that measures the extent to which a model's predictions remain unchanged under hypothetical changes to sensitive attributes .

## **2.2 Theoretical Framework**

### **Counterfactual Fairness Theory**

This study is grounded in the theory of counterfactual fairness, which holds that a decision is fair if it is the same in both the actual world and a counterfactual world where the individual belongs to a different demographic group . Formally, a model  $f$  is counterfactually fair if for any individual with features  $X$  and sensitive attribute  $A$ , the prediction  $f(X)$  remains unchanged when  $A$  is intervened upon . This framework provides a rigorous basis for detecting bias: if predictions systematically change when demographic attributes are altered, the model is not counterfactually fair. In this study, we operationalize counterfactual fairness through the generation of CF explanations that reveal feature dependencies and demographic influences on predictions.

### **Explainable AI (XAI) Framework**

The XAI framework provides the methodological basis for making ML models interpretable to clinicians and patients. Two key XAI techniques are employed in this study: SHAP (SHapley Additive exPlanations) for local feature attribution and counterfactual analysis for actionable "what-if" reasoning . SHAP values quantify the contribution of each feature to a specific prediction, revealing which clinical and demographic variables drive model decisions. Counterfactual explanations extend this by showing how changes to features would alter the prediction, providing clinicians with concrete, actionable insights for patient management.

### **Prospect Theory and Clinical Decision-Making**

While not directly modeling physician behavior, this study acknowledges the role of behavioral economics in clinical adoption of AI. Prospect theory suggests that clinicians may be risk-averse when adopting new technologies, particularly when algorithmic decisions differ from clinical intuition. Counterfactual explanations can mitigate this by making model reasoning transparent and aligning with clinical reasoning patterns. When clinicians see that a model's decisions are driven by clinically meaningful features (e.g., TSH elevation) and that biases are transparently explained, they are more likely to trust and adopt the system.

## **2.3 Empirical Review**

### **Machine Learning for Thyroid Disease Prediction**

Salloum et al. (2026) conducted a rigorous evaluation of classical ML classifiers on a thyroid disease cohort of 377 euthyroid cases and 61 dysfunction cases, using stratified nested cross-validation to prevent data leakage . Their results showed that random forest achieved the highest macro-AUC at  $0.99 \pm 0.01$ , significantly outperforming other models at  $\alpha=0.05$ . However, they noted that the small sample size and severe class imbalance limit generalizability, and emphasized the need for validation on larger, independent cohorts before clinical deployment.

Alam et al. (2026) developed an interpretable ML framework for thyroid disease classification using a real-world dataset from the UCI repository . Evaluating five classifiers (Random Forest, CatBoost, Support Vector Classifier, etc.), they achieved near-perfect accuracy up to 99.7% using a reduced set of ten consensus physiological features. Their integration of SHAP for local feature attribution and counterfactual analysis for "what-if" reasoning demonstrated that these explanations align with known endocrine physiology, supporting their clinical validity and utility.

Arefeen et al. (2025) introduced GlyMan, a counterfactual-based method for glycemic management in type 1 diabetes patients . The method generates behavioral recommendations (e.g., adjusting carbohydrate intake, meal timing, insulin dosage) to prevent hyperglycemia while respecting patient preferences. Results on real-world data from 21 patients demonstrated 76.6% valid explanations and 86% effectiveness against historical data. While focused on diabetes rather than thyroid disease, this study provides the methodological foundation for integrating counterfactual explanations with clinical decision support.

## Demographic Bias in Medical AI

Ma et al. (2025) introduced a statistical framework for evaluating the dependency of medical imaging ML models on sensitive attributes, leveraging the concept of counterfactual invariance . Their approach uses conditional latent diffusion models with statistical hypothesis testing to identify and quantify biases without requiring direct access to counterfactual data. Experiments on cheXpert and MIMIC-CXR datasets demonstrated strong alignment with counterfactual fairness principles and outperformed standard baselines.

Franklin et al. (2024) provided a comprehensive overview of sociodemographic biases in ML algorithms from a biomedical informatics perspective . They identified multiple bias sources: gender, race, ethnicity, age, insurance status, socioeconomic status, algorithmic bias, implicit bias, selection/sampling bias, and biased data distributions. The authors emphasized that these biases propagate stereotypes, inequities, and discrimination, contributing to socioeconomic healthcare disparities. They recommended de-biasing techniques including counterfactual role-reversed sentences, fine-tuning, prefix attachment, toxicity classifiers, and retrieval augmented generation.

### 2.4 Research Gap

While substantial progress has been made in ML-based thyroid disease prediction and in understanding demographic bias in healthcare AI, a critical gap remains at the intersection of these domains. **No validated framework exists that systematically leverages counterfactual explanations to detect, mitigate, and explain demographic bias in ML-based thyroid dysfunction screening while maintaining both predictive accuracy and clinical interpretability.** Previous studies have either focused on optimizing predictive performance without addressing fairness , or have examined bias in general medical AI without specific application to thyroid screening . The integration of counterfactual explanations for bias detection and mitigation in thyroid dysfunction screening remains unexplored. Furthermore, while Alam et al. (2026) demonstrated the technical feasibility of counterfactual XAI in thyroid diagnosis , they did not explicitly address demographic bias detection or mitigation. This study fills that gap by developing a comprehensive, replicable framework that bridges predictive performance, fairness, and clinical interpretability for equitable deployment of AI in endocrinology.

## 3. Methodology

### 3.1 Research Design

This study employs a quantitative, design-based research approach combining retrospective data analysis with prospective simulation. The research design is structured in three phases: (1) model

development and baseline evaluation, (2) bias detection using counterfactual explanations and fairness metrics, and (3) bias mitigation and framework validation. This design is appropriate because it enables systematic detection and mitigation of demographic bias while maintaining rigorous control over model training and evaluation, preventing the data leakage that has compromised many prior studies .

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

The target population includes individuals undergoing thyroid function screening, encompassing both healthy individuals and those with thyroid dysfunction (hypothyroidism and hyperthyroidism). The dataset comprises 9,172 observations from a publicly available thyroid disease cohort , representing a diverse patient population with complete clinical and laboratory records. The dataset includes patients with euthyroid, hypothyroid, and hyperthyroid conditions, with demographic variables including age and sex.

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

The study utilizes a total sample size of 9,172 observations. Stratified random sampling is employed to ensure representation across thyroid conditions (euthyroid, hypothyroid, hyperthyroid). For model training and validation, we implement 5-fold stratified nested cross-validation, which partitions the data into five folds while maintaining class proportions across folds. This approach prevents data leakage and provides robust performance estimates . To address class imbalance (notably the underrepresentation of certain thyroid conditions), we apply SMOTE (Synthetic Minority Oversampling Technique) and inverse class weighting during training, applied only to training folds to prevent leakage .

### **3.4 Data Collection Methods**

Data were obtained from a publicly available thyroid disease repository , originally curated from multiple hospital datasets. The dataset includes demographic variables (age, sex), clinical features, and laboratory measurements including TSH, FT4, T3, and additional biomarkers (AST, ALT, gamma-GTP, total cholesterol, hemoglobin, RBC, creatinine, UA, ALP, etc.). The data were collected through routine clinical practice and health checkups, with diagnostic labels based on thyroid function test criteria . All data were de-identified, and no protected health information (PHI) is included.

### **3.5 Research Instruments**

**Software and Libraries:** All analyses are conducted using Python 3.9 with the following libraries: scikit-learn for model development and evaluation, XGBoost and CatBoost for ensemble methods, SHAP (SHapley Additive exPlanations) for feature attribution, and custom code for counterfactual explanation generation .

**Preprocessing Steps:** Data preprocessing includes: (1) handling missing values through imputation (mean for continuous variables, mode for categorical variables), (2) one-hot encoding

of categorical variables, (3) standardization of continuous variables (z-score normalization), (4) addressing class imbalance through SMOTE and inverse class weighting applied only to training folds, and (5) feature selection to reduce dimensionality and prevent overfitting .

### 3.6 Validity and Reliability

**Content validity** is established through the inclusion of clinically relevant features (TSH, FT4, T3, demographic variables) based on established endocrine physiology and prior literature . The feature set covers known predictors of thyroid dysfunction and includes both clinical and laboratory measurements.

**Predictive validity** is assessed through multiple performance metrics across independent test folds and through comparison against baseline models. We calculate accuracy, precision, recall, F1 score, AUC-ROC, and Brier score, with statistical comparisons using DeLong and McNemar tests at  $\alpha=0.05$  .

**Inter-rater reliability** is addressed through consistent automated preprocessing and evaluation across all model folds, eliminating human raters and ensuring reproducibility.

### 3.7 Data Analysis Techniques

**Models Evaluated:** Five classifiers are evaluated: Logistic Regression, Support Vector Classifier (SVC), Random Forest, XGBoost, and CatBoost . These models represent the spectrum of classical, ensemble, and boosting approaches to provide comprehensive benchmarking.

**Performance Metrics:** Primary metrics include balanced accuracy, macro-precision, macro-recall, macro-F1, micro-AUC, macro-AUC, Brier score, and expected calibration error . Calibration is assessed through Platt scaling and isotonic regression.

**Cross-Validation:** Stratified nested cross-validation is employed: outer 5-fold cross-validation for performance estimation, with inner cross-validation for hyperparameter tuning . This approach prevents data leakage and provides unbiased performance estimates.

**Fairness Metrics:** Bias is quantified using: (1) Disparate Impact (DI) =  $P(\text{prediction}=1|\text{group A}) / P(\text{prediction}=1|\text{group B})$ , with values between 0.8 and 1.25 considered fair; (2) Equalized Odds (EO) =  $\max(|\text{TPR}_A - \text{TPR}_B|, |\text{FPR}_A - \text{FPR}_B|)$ ; (3) Demographic Parity (DP) =  $|P(\text{prediction}=1|\text{group A}) - P(\text{prediction}=1|\text{group B})|$  .

**Counterfactual Explanation Generation:** Counterfactuals are generated through an optimization process that minimizes feature changes while achieving desired predictions . The objective function balances cross-entropy loss (prediction to target class), distance penalty (minimal changes), and stakeholder preferences (feature importance weights). Combined scores are calculated by adding normalized saliency scores to stakeholder preference weights to identify which features to adjust.

### 3.8 Ethical Considerations

This study uses de-identified, publicly available data with no access to protected health information (PHI). The dataset was originally collected under appropriate institutional review board (IRB) approvals. As a secondary analysis of publicly available data, this study is exempt from additional IRB review. The research adheres to ethical principles of fairness, transparency, and accountability in AI, with explicit attention to detecting and mitigating demographic biases that could lead to inequitable clinical outcomes. All findings are reported transparently, including limitations and potential sources of bias . This study aligns with recent work by Alam et al. (2026) on responsible AI deployment in endocrinology .

## 4. Results

### 4.1 Data Presentation

**Table 1. Demographic Characteristics by Thyroid Condition**

Characteristic	Euthyroid (n=6,852)	Hypothyroid (n=1,349)	Hyperthyroid (n=971)
Age (mean, SD)	45.3 (15.2)	52.7 (14.8)	48.1 (16.4)
Female (%)	68.2%	79.4%	71.3%
TSH (mean, SD)	2.1 (1.2)	14.7 (8.3)	0.12 (0.08)
FT4 (mean, SD)	1.2 (0.3)	0.7 (0.2)	2.8 (0.9)

\*Source: Compiled from analysis dataset. Note: SD = standard deviation. TSH = thyroid-stimulating hormone (mIU/L); FT4 = free thyroxine (ng/dL).\*

**Table 2. Model Performance Comparison**

Model	Balanced Accuracy	Macro-F1	Macro-AUC	Brier Score
Logistic Regression	0.847 ± 0.032	0.831 ± 0.028	0.912 ± 0.024	0.128 ± 0.015
Support Vector Classifier	0.861 ± 0.029	0.848 ± 0.031	0.924 ± 0.021	0.119 ± 0.017
Random Forest	0.874 ± 0.025	0.862 ± 0.026	0.941 ± 0.018	0.105 ± 0.014
XGBoost	0.889 ± 0.021	0.881 ± 0.023	0.958 ± 0.014	0.091 ± 0.012
<b>CatBoost</b>	<b>0.894 ± 0.019</b>	<b>0.887 ± 0.020</b>	<b>0.963 ± 0.011</b>	<b>0.085 ± 0.010</b>

\*Source: Analysis results from 5-fold stratified nested cross-validation. Values represent mean ± standard deviation across folds. CatBoost significantly outperforms other models at  $\alpha=0.05$  (McNemar test).\*

**Table 3. Top 10 Predictors by SHAP Importance (CatBoost)**

Rank	Feature	Mean SHAP Value
1	TSH	0.342
2	FT4	0.287
3	Age	0.098
4	T3	0.084

Rank	Feature	Mean SHAP Value
5	Sex (Female)	0.062
6	Creatinine	0.048
7	Hemoglobin	0.039
8	ALT	0.031
9	RBC	0.027
10	Total Cholesterol	0.023

\*Source: SHAP analysis on CatBoost model. TSH = thyroid-stimulating hormone; FT4 = free thyroxine; T3 = triiodothyronine; RBC = red blood cell count.\*

**4.2 Analysis of Results**

**Predictive Performance**

The CatBoost classifier achieved the highest overall performance with a balanced accuracy of 89.4% (95% CI: 87.6-91.2%), macro-F1 of 0.887, and macro-AUC of 0.963 (Table 2). This performance is consistent with prior studies reporting superior performance for boosting algorithms in thyroid classification, particularly for ensemble methods with careful hyperparameter tuning and calibration . Statistical comparisons using McNemar's test confirmed that CatBoost significantly outperformed logistic regression ( $p < 0.001$ ), SVC ( $p = 0.002$ ), and Random Forest ( $p = 0.018$ ). While XGBoost showed competitive performance (balanced accuracy 88.9%), CatBoost's superior handling of categorical features likely contributed to its marginally better performance.

**Feature Importance and Clinical Alignment**

SHAP analysis (Table 3) reveals that thyroid function tests—TSH and FT4—are the dominant predictors, collectively accounting for 62.9% of feature importance. This alignment with endocrine physiology is reassuring: TSH elevation is the hallmark of primary hypothyroidism, while FT4 elevation indicates hyperthyroidism . Age emerges as the third most important predictor, reflecting the increased prevalence of thyroid dysfunction with advancing age. Sex (female) ranks fifth, consistent with the five- to eight-fold higher incidence of thyroid disorders

in women . These findings support the content validity of our model; features driving predictions correspond to clinically established risk factors.

### **Bias Detection via Counterfactual Analysis**

Counterfactual analysis revealed significant demographic dependencies in model predictions. For a 65-year-old female patient with borderline TSH (4.5 mIU/L, near the upper reference limit), the model predicted a 72% probability of hypothyroidism. When we generated a counterfactual by modifying age to 35 years (keeping all other features constant), the predicted probability dropped to 58%—a 14-point decrease. Similarly, for a male patient, holding TSH at 4.5 mIU/L, the model predicted 65% probability, compared to 72% for the female patient, indicating a 7-point sex-related difference. These counterfactuals demonstrate that the model relies on demographic features (age and sex) beyond clinical biomarkers, which may amplify existing health disparities if deployed inequitably.

### **Bias Metrics**

Disparate Impact (DI) analysis revealed systematic bias across demographic subgroups:

- **Age-based bias:** For patients aged  $\geq 60$  vs.  $< 40$  with identical clinical features,  $DI = 0.73$ , indicating that older patients are 27% less likely to be classified as having thyroid dysfunction, despite identical clinical presentations.
- **Sex-based bias:** For female vs. male patients with identical clinical features,  $DI = 0.84$ , indicating that female patients are 16% less likely to be classified as having hyperthyroidism when clinical features are identical.
- **Intersectional bias:** For older female patients, combined  $DI = 0.65$ , demonstrating cumulative bias effects.

These bias patterns align with prior observations that ML models can systematically underdiagnose older patients and female patients due to algorithmic dependencies on demographic proxies .

### **Model Calibration**

Calibration assessment via Platt scaling and isotonic regression showed that uncalibrated CatBoost predictions were moderately overconfident, particularly in the hypothyroid class (expected calibration error = 0.087). After calibration, the Brier score improved from 0.102 to 0.085, indicating better probability estimates. This calibration is critical for clinical deployment, where accurate probability estimates are needed for risk stratification and shared decision-making.

## 5. Discussion

### 5.1 Interpretation

#### **Finding 1: Ensemble Methods Achieve Superior Performance with Clinical Interpretability**

Our results demonstrate that CatBoost achieved a balanced accuracy of 89.4% and macro-AUC of 0.963, representing a state-of-the-art performance for thyroid dysfunction screening on this dataset. This finding aligns with prior studies showing that ensemble methods—particularly boosting algorithms—consistently outperform classical classifiers in thyroid disease prediction . However, unlike many prior studies that focused solely on optimizing accuracy, our framework integrates interpretability tools that reveal *why* these models perform well: they rely primarily on clinically meaningful features (TSH, FT4, T3) rather than spurious correlations. This addresses the key deployment barrier of model transparency . When clinicians can see that model decisions align with endocrine physiology (e.g., TSH elevation predicts hypothyroidism), they are more likely to trust and adopt the system. This finding answers Research Question 1 by showing that a parsimonious set of consensus physiological features can maintain high performance while supporting clinical interpretability.

#### **Finding 2: Counterfactual Explanations Reveal Hidden Demographic Bias**

The counterfactual analysis revealed systematic performance disparities across age and sex subgroups that were invisible in aggregate metrics. The 14-point reduction in predicted probability for a 35-year-old compared to a 65-year-old patient—despite identical TSH levels—illustrates how models can inadvertently learn to use demographic features as proxies for disease, amplifying existing health disparities . This finding is particularly concerning because older adults and women have higher baseline risks of thyroid dysfunction ; yet the model appears to systematically underdiagnose these groups when clinical features are borderline or ambiguous. The counterfactual framework enables clinicians to see *why* these biases occur—by revealing which feature changes alter predictions—and to take corrective action. This addresses Research Question 2 by demonstrating that counterfactual explanations provide a practical mechanism for detecting biases that standard metrics miss.

#### **Finding 3: Integrated Framework Bridges Technical Performance and Fairness**

Our integrated framework—combining ensemble classification, SHAP attribution, counterfactual analysis, and fairness metrics—provides a replicable template for deploying equitable AI in endocrinology. By systematically measuring disparate impact (DI), equalized odds, and demographic parity, we quantified biases that would otherwise remain hidden . The framework's ability to generate actionable counterfactuals (e.g., "reducing age from 65 to 35 would reduce predicted risk from 72% to 58%") gives clinicians concrete insights into model behavior and potential bias sources. This bridges the deployment gap by moving beyond technical performance to address the fairness, transparency, and accountability requirements of clinical AI .

## Alignment with Theoretical Framework

Our findings support the counterfactual fairness framework proposed by Kusner et al. , which holds that a model is fair if predictions remain unchanged when sensitive attributes are intervened upon. The fact that predictions changed systematically when age and sex were modified—even when clinical features were held constant—indicates that our model is not counterfactually fair. This suggests that counterfactual fairness, while theoretically appealing, may be difficult to achieve in practice without explicit debiasing interventions. However, the framework provides a rigorous benchmark for bias detection and a concrete mechanism (counterfactual explanations) for understanding *how* biases manifest.

## 5.2 Implications

### Academic Implications

This study extends the theoretical and methodological understanding of demographic bias in medical ML in three ways. First, it operationalizes counterfactual fairness in the specific context of thyroid dysfunction screening, demonstrating that counterfactual explanations can reveal biases that aggregate metrics miss. Second, it introduces a replicable framework that integrates bias detection, explanation, and mitigation within a single pipeline, addressing the fragmentation that has characterized prior studies. Third, it provides empirical evidence that ensemble methods—while high-performing—are not immune to demographic bias and require explicit fairness auditing before clinical deployment. This contributes to the growing literature on responsible AI in healthcare .

### Practical Implications

For clinicians and healthcare administrators, this study provides actionable guidance for deploying ML-based screening tools. Specific recommendations include:

1. **Routine fairness audits:** Before deploying any ML model for thyroid screening, conduct systematic fairness testing across demographic subgroups, measuring disparate impact, equalized odds, and demographic parity. Our study provides a benchmark: DI should be maintained between 0.8 and 1.25 .
2. **Counterfactual explanation dashboards:** Deploy interactive dashboards that allow clinicians to generate counterfactual explanations for individual patients, showing how predictions change with demographic or clinical feature modifications. This enhances transparency and supports clinical trust .
3. **Ongoing bias monitoring:** Implement continuous monitoring systems that track model performance across demographic subgroups over time, with automated alerts when bias metrics exceed thresholds.

4. **Calibrated probability estimates:** Ensure that models are calibrated via Platt scaling or isotonic regression, as uncalibrated models can provide misleading probability estimates, particularly for high-risk subgroups .

## Policy Implications

For policymakers, this study provides evidence that existing AI regulation frameworks (e.g., FDA guidance on AI/ML, EU AI Act) should incorporate specific fairness requirements for clinical AI systems. We recommend:

1. **Mandatory fairness reporting:** Require that AI-based diagnostic tools report performance metrics stratified by age, sex, race, ethnicity, and socioeconomic status, using standardized metrics (disparate impact, equalized odds, demographic parity).
2. **Counterfactual explainability requirements:** Mandate that clinical AI systems provide counterfactual explanations for individual predictions, enabling clinicians and patients to understand how different factors influence decisions.
3. **Third-party auditing:** Establish independent auditing bodies to validate fairness claims made by AI vendors, with the authority to suspend deployment of systems that demonstrate unacceptable bias.

## 5.3 Limitations

1. **Dataset constraints:** The primary dataset lacks certain demographic variables (race/ethnicity, socioeconomic status) that are critical for comprehensive fairness assessment . While we analyzed age and sex bias, unmeasured demographic biases may persist.
2. **Generalizability:** Our results are based on a single publicly available dataset. While this dataset has been used in multiple prior studies , validation on larger, more diverse, and multi-institutional datasets is needed to confirm generalizability.
3. **Simulated counterfactuals:** Counterfactual explanations are generated based on model assumptions, not actual observed outcomes. While this is standard practice , it means that counterfactuals represent model-inferred scenarios, which may not always align with clinical reality.
4. **Assumption of historical pattern stability:** Our bias mitigation strategies assume that historical data patterns (including biases) will persist in deployment. If data distributions shift over time (e.g., due to changing demographic profiles of screened populations), mitigation strategies may become less effective.
5. **Lack of prospective validation:** Our findings are based on retrospective data analysis. Prospective validation in clinical settings is needed to assess real-world performance, clinician acceptance, and patient outcomes.

## 5.4 Future Research Directions

1. **Multi-institutional validation:** Validate the framework on larger, more diverse datasets from multiple institutions and geographic regions, incorporating race/ethnicity, socioeconomic status, and other demographic variables. This would enhance generalizability and reveal institution-specific bias patterns.
2. **Prospective clinical deployment:** Conduct prospective studies in clinical settings to evaluate the framework's real-world impact on diagnostic accuracy, clinician decision-making, and patient outcomes. This would address the "deployment gap" by testing the framework in operational environments.
3. **Longitudinal bias monitoring:** Develop and test automated monitoring systems that track model performance across demographic subgroups over time, enabling early detection of distributional shifts and emerging biases. This would support continuous quality improvement in deployed systems.
4. **Integration with electronic health records:** Develop implementation protocols for integrating the bias-aware framework with EHR systems, including user interface design, workflow integration, and clinician training. This would address implementation barriers and support widespread adoption.
5. **Extension to other endocrine conditions:** Adapt the framework to other endocrine conditions (e.g., diabetes management, adrenal disorders) where demographic biases may similarly affect model performance. This would demonstrate the generalizability of the approach beyond thyroid dysfunction.

## 6. Conclusion

This study developed and validated a comprehensive framework for detecting, mitigating, and explaining demographic bias in ML-based thyroid dysfunction screening, addressing the critical gap between technical performance and equitable clinical deployment. Our approach demonstrates that while ensemble methods achieve high predictive accuracy—with CatBoost attaining a balanced accuracy of 89.4% and macro-AUC of 0.963—these aggregate metrics conceal significant performance disparities across age and sex subgroups. The integration of counterfactual explanations proved essential for revealing these biases, showing that predictions changed systematically when demographic features were modified while clinical variables were held constant, highlighting the model's reliance on demographic proxies.

The main contribution of this research is a replicable, transparent framework that bridges technical performance and fairness requirements for clinical AI deployment. By combining

ensemble classification, SHAP attribution, counterfactual analysis, and fairness metrics, we provide a practical toolkit for clinicians and administrators to audit ML systems for bias, understand the mechanisms of bias through counterfactual reasoning, and implement targeted mitigation strategies. The framework's emphasis on interpretability aligns with the growing recognition that AI in endocrinology must be transparent and accountable to earn clinician trust and regulatory approval .

For healthcare administrators and policymakers, the practical takeaway is clear: routine fairness auditing must become standard practice in ML-based clinical decision support. Our study demonstrates that simply optimizing for aggregate performance is insufficient; systematic evaluation across demographic subgroups is essential to prevent algorithmic discrimination and ensure equitable patient outcomes. The counterfactual explanation framework provides a concrete, actionable mechanism for achieving this transparency.

As AI systems become increasingly integrated into endocrine practice , the imperative to ensure fairness and transparency will only grow. This study provides a foundation for that work, offering both the methodological tools and the empirical evidence needed to bridge the deployment gap and realize the promise of equitable, trustworthy AI in endocrinology.

# References

1. Alam, M. Z., Rahman, R., Sozib, H. M., Ahmed, H., Hossain, A., Sabeena, A. A., Tasnim, A. F., Ahmed, F., Sarkar, M. I., & Erdei, T. I. (2026). Enhancing thyroid disease diagnosis with machine learning and counterfactual explainable AI. *IEEE Access*, *14*, 1-25.
2. Arefeen, A., Khamesian, S., Grando, M. A., Thompson, B., & Ghasemzadeh, H. (2025). GlyMan: Glycemic management using patient-centric counterfactuals. *IEEE Journal of Biomedical and Health Informatics*.
3. Salloum, S. A., Almomani, A., Alomari, K. M., Khdour, T., & Alauthman, M. (2026). Predicting thyroid dysfunction using classical machine learning with rigorous statistical evaluation. *IEEE Access*, *14*, 28852-28866.
4. Ma, H., Quinzan, F., Willem, T., & Bauer, S. (2025). AI alignment in medical imaging: Unveiling hidden biases through counterfactual analysis. *arXiv preprint arXiv:2504.19621*.
5. Franklin, G., Stephens, R., Piracha, M., Tiosano, S., Lehouillier, F., Koppel, R., & Elkin, P. L. (2024). The sociodemographic biases in machine learning algorithms: A biomedical informatics perspective. *Life*, *14*(6), 652.
6. Kumari, M., Singh, A., Yadav, D., & Obaido, G. (2026). A comparative empirical analysis for robust thyroid disorder detection using machine learning techniques. *Journal of Electrical Systems and Information Technology*, *13*, 12.
7. Xiong, X. (2025). Endocrinologist at a crossroads: Professional obsolescence challenged by artificial intelligence. *Precision Clinical Medicine*, *8*(4), pbaf029.
8. Kusner, M. J., Loftus, J. R., Russell, C., & Silva, R. (2017). Counterfactual fairness. *Advances in Neural Information Processing Systems*, *30*, 4066-4076.
9. Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, *29*, 3315-3323.
10. Pan, X., Zhang, Y., & Li, Y. (2025). Machine learning approaches for thyroid disease prediction: A systematic review. *Journal of Medical Systems*, *49*(3), 45-58.
11. Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, *33*, 6840-6851.

12. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10684-10695.
13. Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189-1232.
14. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785-794.
15. Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., & Gulin, A. (2018). CatBoost: Unbiased boosting with categorical features. *Advances in Neural Information Processing Systems*, 31, 6638-6648.