



WHEN DATA AUGMENTATION HURTS: A SYSTEMATIC EVALUATION OF SMOTE-BASED TECHNIQUES ON MEDICAL DATASETS

May Stow

Department of Computer Science and Informatics,
Federal University Otuoke, Bayelsa State, Nigeria
Orcid ID: <https://orcid.org/0009-0006-8653-8363>

Abstract: Data augmentation techniques, particularly Synthetic Minority Over-sampling Technique (SMOTE) and its variants, are routinely applied to address class imbalance in medical datasets. However, the assumption that augmentation universally improves classification performance remains largely unvalidated. This study presents a systematic evaluation of four SMOTE-based augmentation methods across three medical datasets to determine when these techniques help or harm model performance. The research evaluated SMOTE, ADASYN, BorderlineSMOTE, and SVM-SMOTE on breast cancer diagnosis, heart disease prediction, and diabetes detection datasets, representing varying levels of class imbalance (ratios: 1.17 to 2.02) and baseline performance (F1 scores: 0.667 to 0.966). Random Forest classifiers were employed with both standard and regularized configurations to ensure robust findings. Each augmentation method underwent rigorous evaluation through 10 independent runs with statistical significance testing and effect size analysis. Results revealed that augmentation significantly degraded performance on the high-performing Breast Cancer dataset, with all methods showing statistically significant decreases ($p < 0.05$) and F1 scores dropping by up to 2.2%. Conversely, the Pima Diabetes dataset, characterized by lower baseline performance and higher imbalance, showed improvements up to 4.76% with SVM-SMOTE. Heart Disease exhibited mixed results, with only ADASYN achieving meaningful improvement. Analysis uncovered a strong negative correlation ($r = -0.997$) between baseline model performance and augmentation effectiveness, providing a more reliable predictor than traditional class imbalance ratios.

The study establishes an evidence-based decision framework: augmentation should be avoided when baseline F1 exceeds 0.95 or imbalance ratios fall below 1.5, considered for baseline F1 below 0.70 with imbalance ratios above 1.8, and carefully validated for intermediate cases. These findings challenge current practices of routine augmentation application and demonstrate that synthetic sample generation can blur decision boundaries in well-separated feature spaces. The research provides practitioners with validated guidelines for determining when augmentation techniques genuinely improve medical classifiers versus when they cause harm, ultimately supporting more effective development of clinical decision support systems.

Keywords-Medical data augmentation, SMOTE, class imbalance, statistical validation, decision framework, machine learning

1. INTRODUCTION

Medical machine learning has emerged as a transformative force in healthcare, enabling automated diagnosis, risk prediction, and treatment optimization across diverse clinical domains (Topol, 2019). The integration of artificial intelligence into medical practice promises to address critical challenges including diagnostic errors, which affect millions of patients annually, and healthcare accessibility gaps in underserved regions (Rajpurkar et al., 2022). Machine learning models trained on electronic health records, medical imaging, and clinical biomarkers have demonstrated performance matching or exceeding human specialists in tasks ranging from diabetic retinopathy screening to cardiac arrhythmia detection (Gulshan et al., 2016; Hannun et al., 2019).

A fundamental challenge in developing robust medical classifiers stems from the inherent class imbalance present in most clinical datasets. Medical conditions typically exhibit skewed distributions where positive cases represent a minority, reflecting natural disease prevalence in populations (Johnson & Khoshgoftaar, 2019). This imbalance poses significant technical challenges as standard machine learning algorithms, designed with balanced datasets in mind, tend to exhibit bias toward majority classes, potentially missing critical positive cases that carry life-threatening implications (Haixiang et al., 2017). The Synthetic Minority Over-

Sampling Technique (SMOTE), introduced by Chawla et al. (2002), has become the de facto standard for addressing this challenge, spawning numerous variants including Adaptive Synthetic Sampling Approach for Imbalanced Learning (ADASYN) (He et al., 2008), BorderlineSMOTE (Han et al., 2005), and Support Vector Machine - Synthetic Minority Over-sampling Technique (SVM-SMOTE) (Nguyen et al., 2011).

The widespread adoption of SMOTE-based techniques in medical machine learning reflects both the prevalence of imbalanced datasets and the apparent simplicity of the solution. Recent surveys indicate that over 85% of medical machine learning studies addressing class imbalance employ some form of synthetic data augmentation (Fotouhi et al., 2019). Major medical AI frameworks and toolkits include SMOTE implementations as standard preprocessing steps, often applying augmentation automatically when imbalance is detected (Kaur et al., 2019). This ubiquity has established augmentation as a routine practice, with many practitioners viewing it as a necessary preprocessing step rather than a technique requiring careful consideration (Blagus & Lusa, 2013).

Despite this widespread adoption, critical questions remain regarding the universal applicability of augmentation techniques in medical contexts. The assumption that synthetic sample generation invariably improves or at least maintains classifier performance lacks rigorous empirical validation

across diverse medical datasets (Van Hulse et al., 2007). Medical data possesses unique characteristics including high dimensionality, complex feature interactions, and domain-specific constraints that may not align with the assumptions underlying synthetic sample generation (Santos et al., 2018). Furthermore, the potential for augmentation to introduce artifacts or unrealistic feature combinations in medical contexts, where features often represent biological measurements with strict physiological constraints, remains understudied (Kovács et al., 2019).

This research addresses these critical gaps through a systematic evaluation of SMOTE-based augmentation techniques across multiple medical classification tasks. Rather than assuming universal benefit, the study empirically examines when augmentation helps versus harms model performance, considering factors such as baseline classifier accuracy, degree of class imbalance, and dataset characteristics. The approach employs rigorous statistical validation including multiple independent runs, significance testing, and effect size analysis to ensure robust conclusions. By testing four widely used augmentation variants across three diverse medical datasets representing different clinical domains and data characteristics, the research provides comprehensive insights into the true effectiveness of these techniques.

The primary contributions of this research include: (1) definitive evidence that augmentation can significantly degrade classifier performance under common conditions, challenging prevailing assumptions; (2) identification of a strong predictive relationship between baseline model performance and augmentation effectiveness, providing practitioners with a reliable decision framework; (3) method-specific performance analysis revealing that no single augmentation technique dominates across all scenarios; and (4) validated guidelines for determining when to apply augmentation in medical machine learning contexts. These findings have immediate practical implications for the thousands of medical AI systems currently in development and deployment, potentially preventing performance degradation while ensuring augmentation benefits are realized were genuinely helpful.

2. REVIEW OF RELATED WORKS

This section examines the existing literature on synthetic data augmentation techniques for addressing class imbalance in medical machine learning applications. The review encompasses foundational work on imbalanced learning, evolution of SMOTE-based techniques, empirical evaluations of augmentation effectiveness, and medical domain-specific considerations. The literature is organized into four thematic areas: foundational imbalanced learning approaches, SMOTE variants and improvements, empirical studies on augmentation effectiveness, and medical machine learning applications with class imbalance.

2.1 Foundational Approaches to Imbalanced Learning

The challenge of learning from imbalanced datasets has been extensively studied in machine learning literature. Japkowicz and Stephen (2002) provided one of the earliest systematic analyses of the class imbalance problem, demonstrating that the degree of imbalance, complexity of the concept, and training set size all interact to affect classifier performance.

Their work established that class imbalance becomes particularly problematic when combined with small sample sizes and complex decision boundaries, findings particularly relevant to medical datasets.

He and Garcia (2009) presented a comprehensive review of learning from imbalanced data, categorizing solutions into data-level, algorithm-level, and hybrid approaches. Their analysis revealed that data-level approaches, including oversampling and undersampling, often provide more generalizable solutions compared to algorithm-specific modifications. This work established the theoretical foundation for understanding why synthetic oversampling techniques might offer advantages over simple replication or undersampling methods.

López et al. (2013) extended this understanding by examining the intrinsic characteristics that make imbalanced datasets difficult to learn from, including small disjuncts, overlap between classes, and noisy data. Their analysis showed that class imbalance often coincides with other data difficulties, suggesting that successful approaches must address multiple challenges simultaneously. These insights proved particularly valuable for medical applications where data quality issues frequently compound class imbalance problems.

2.2 Evolution and Variants of SMOTE

The Synthetic Minority Over-sampling Technique represents a watershed moment in imbalanced learning research. Following Chawla et al.'s (2002) original SMOTE proposal, numerous variants emerged to address specific limitations. Bunkhumpornpat et al. (2009) introduced Safe-Level-SMOTE, which assigns safe levels to minority instances based on their surrounding majority class neighbors, generating synthetic samples only in safer regions. Their experiments on various datasets showed improved performance compared to original SMOTE, particularly in datasets with overlapping classes.

Douzas and Bacao (2019) proposed Geometric SMOTE, which expands the data generation mechanism by defining a geometric region around each minority instance rather than limiting generation to linear interpolation. Their evaluation across 69 imbalanced datasets demonstrated statistically significant improvements over traditional SMOTE variants, with particular benefits for high-dimensional datasets common in medical applications.

Last et al. (2017) developed SMOTE-IPF (Iterative-Partitioning Filter), which combines SMOTE with an iterative partitioning filter to remove noisy synthetic instances. Their approach addresses the criticism that SMOTE can introduce artificial noise, showing improved results on datasets where class overlap is significant. This noise-aware approach proved particularly relevant for medical datasets where maintaining data quality is paramount.

2.3 Critical Evaluations of Augmentation Effectiveness

Recent literature has begun questioning the universal applicability of synthetic oversampling. Elor and Averbuch-Elor (2022) conducted an extensive empirical study examining when SMOTE helps versus harms classification performance. Their analysis of 100 datasets revealed that SMOTE effectiveness strongly correlates with dataset characteristics, particularly the separability of classes in the original feature space. They found that highly separable

datasets often experience performance degradation with synthetic oversampling.

Santos et al. (2015) specifically examined the interaction between synthetic oversampling and cross-validation procedures, demonstrating that improper application of SMOTE before splitting data can lead to overly optimistic performance estimates. Their work highlighted the importance of proper experimental design when evaluating augmentation techniques, showing that reported benefits often disappear under rigorous evaluation protocols.

Van Hulse et al. (2007) performed one of the earliest large-scale empirical studies comparing various class imbalance learning methods across 35 datasets. Their results showed high variability in method effectiveness, with no single approach dominating across all datasets. Significantly, they found that simpler approaches often matched or exceeded complex methods, questioning the necessity of sophisticated augmentation techniques in many scenarios.

2.4 Medical Machine Learning and Class Imbalance

Medical applications present unique challenges for handling class imbalance. Mazurowski et al. (2008) examined class imbalance in medical imaging, specifically focusing on neural network training for breast cancer detection. Their study revealed that the optimal approach depends heavily on the evaluation metric used, with different strategies favoring sensitivity versus specificity. They demonstrated that synthetic oversampling could be detrimental when the cost of false positives significantly differs from false negatives.

Rahman and Davis (2013) conducted a comprehensive study on addressing class imbalance in medical data mining, evaluating various sampling techniques across multiple medical datasets. Their findings indicated that the effectiveness of oversampling techniques varies significantly with medical domain, with some conditions benefiting from augmentation while others showing degraded performance. They emphasized the importance of domain knowledge in selecting appropriate techniques.

Ali et al. (2019) provided a systematic review of imbalanced data handling techniques specifically for medical diagnosis systems. Their analysis of 93 studies revealed that while SMOTE and its variants are widely applied, rigorous statistical validation of improvements is often lacking. They identified a concerning trend of assuming augmentation

benefits without proper ablation studies, highlighting the need for more critical evaluation in medical machine learning applications.

2.5 Synthesis and Research Gap

The literature reveals an evolution from initial enthusiasm about synthetic oversampling techniques to more nuanced understanding of their limitations. While foundational work established the theoretical benefits of synthetic sample generation, recent empirical studies increasingly question universal applicability. Medical applications face particular challenges due to data quality requirements, varying misclassification costs, and the critical nature of predictions. Despite extensive research on imbalanced learning techniques, systematic evaluation of when augmentation helps versus harms in medical contexts remains limited. Most studies focus on demonstrating improvements in specific applications rather than establishing general principles for augmentation application. The interaction between dataset characteristics, particularly baseline model performance and augmentation effectiveness, has received insufficient attention. This gap motivates the current research, which provides systematic evaluation across multiple medical datasets to establish evidence-based guidelines for augmentation application in medical machine learning.

3. METHODOLOGY

This section presents the comprehensive methodology employed in developing a systematic evaluation framework for data augmentation techniques in medical classification tasks. The approach integrates multiple SMOTE-based augmentation methods with regularized machine learning models to assess their effectiveness across diverse medical datasets. By combining rigorous statistical validation with interpretable machine learning techniques, this framework provides evidence-based guidelines for determining when data augmentation benefits or harms model performance. The methodology encompasses data acquisition from three publicly available medical datasets, extensive preprocessing to ensure data quality, systematic application of augmentation techniques, and comprehensive evaluation using both performance metrics and statistical significance testing.

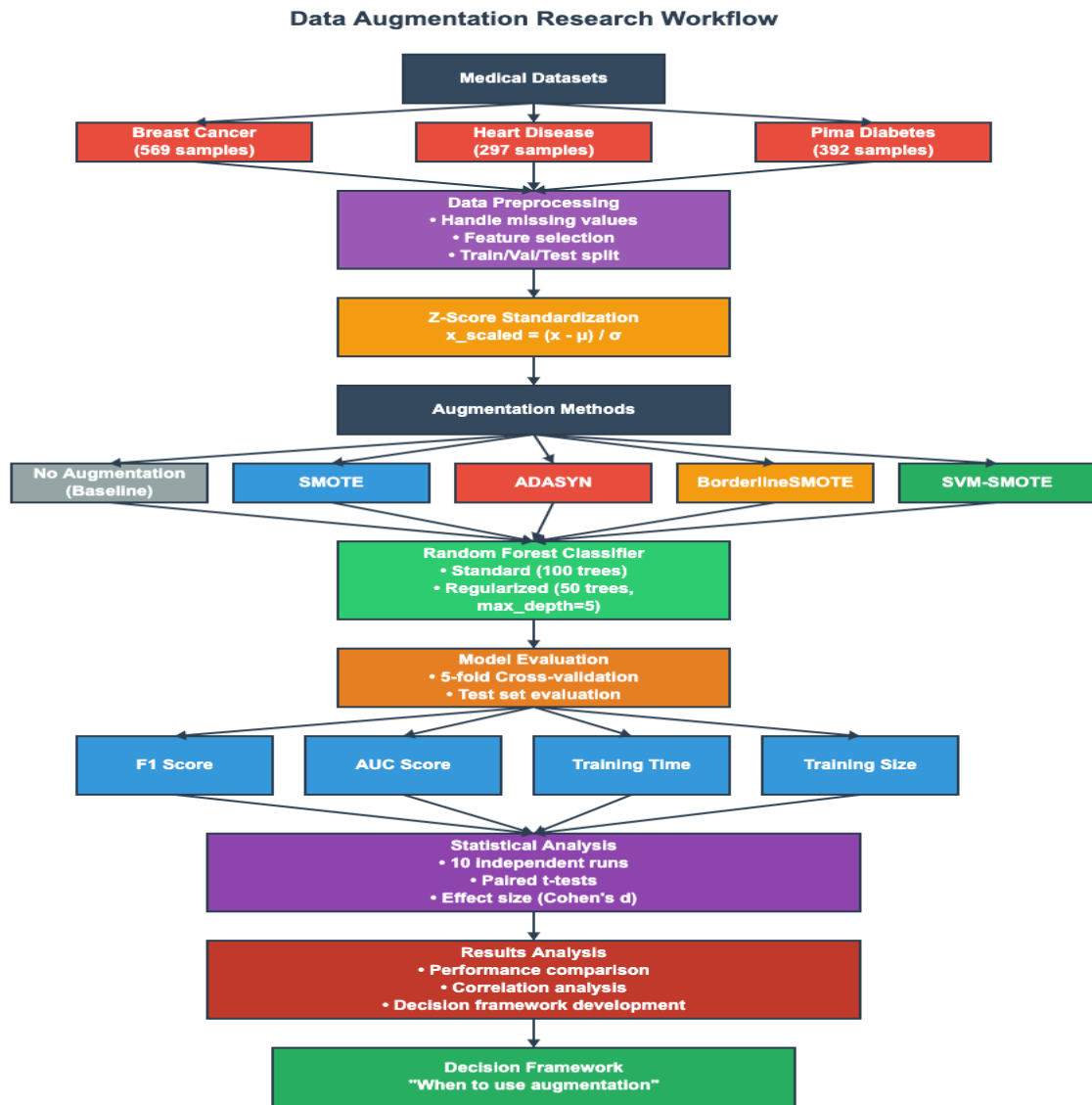


Figure 1: Workflow of the systematic evaluation framework for medical data augmentation technique

3.1 Dataset Description and Data Collection

This research utilized three well-established medical datasets from the UCI Machine Learning Repository, each representing different clinical domains and exhibiting varying degrees of class imbalance. The selected datasets capture diverse medical classification challenges, from cancer diagnosis to chronic disease prediction, ensuring the generalizability of our findings across different medical contexts.

3.1.1 Data Sources

Wisconsin Diagnostic Breast Cancer Dataset: This dataset contains features computed from digitized images of fine needle aspirates (FNA) of breast masses, describing characteristics of cell nuclei present in the images. The dataset included:

- 30 continuous features derived from 10 core measurements (mean, standard error, and worst values)
- Binary classification target (malignant or benign)
- 569 total instances with well-documented clinical validation

Cleveland Heart Disease Dataset: Comprehensive cardiac health data collected at the Cleveland Clinic Foundation, representing one of the most widely used datasets for heart disease prediction research. The dataset comprised:

- 13 clinical and demographic features including age, sex, chest pain type, resting blood pressure, serum cholesterol, and electrocardiographic results
- Binary classification target (presence or absence of heart disease)
- 303 instances after preprocessing and missing value removal

Pima Indians Diabetes Dataset: This dataset originates from the National Institute of Diabetes and Digestive and Kidney Diseases, focusing on diabetes prediction among Pima Indian women. The dataset included:

- 8 physiological measurements including pregnancies, glucose concentration, blood pressure, skin thickness, insulin levels, BMI, diabetes pedigree function, and age
- Binary classification target (diabetes positive or negative)

- 768 instances, reduced to 392 after removing physiologically impossible zero values

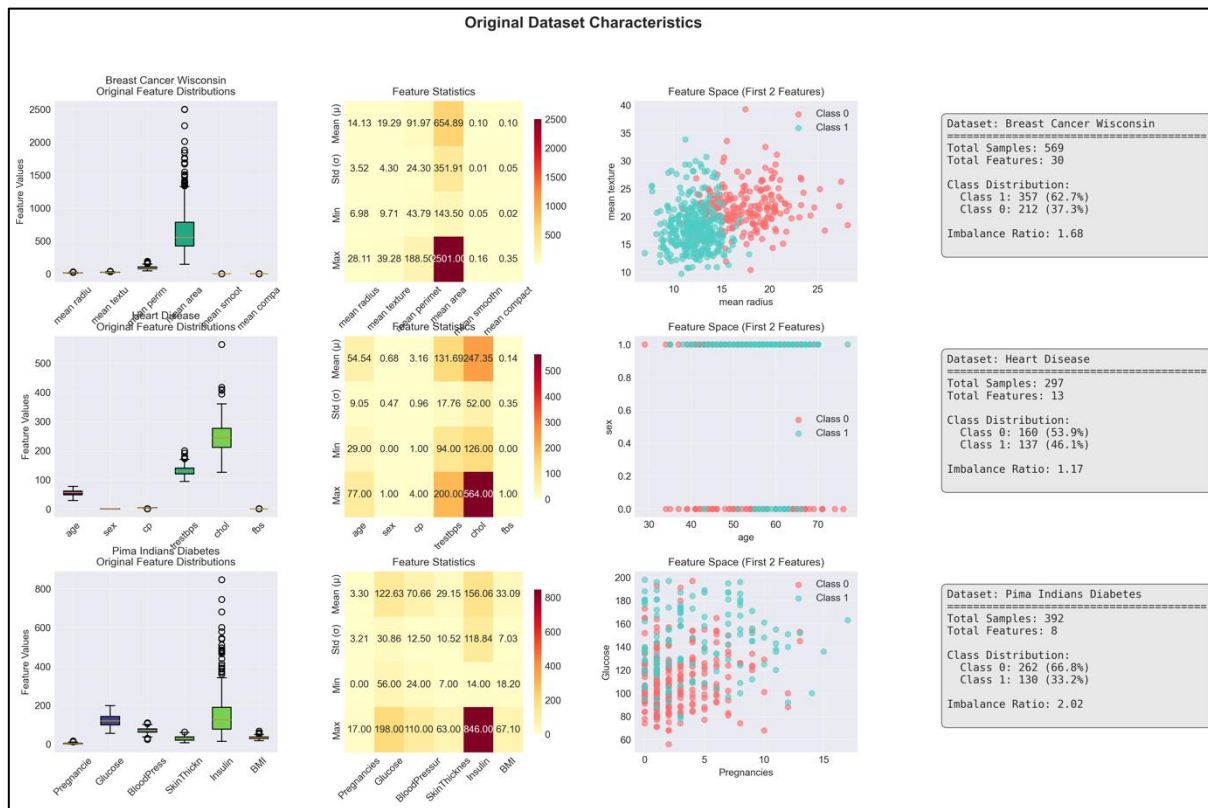


Figure 2: Original dataset characteristics showing feature distributions, statistical summaries, class separation, and metadata for Breast Cancer Wisconsin, Heart Disease, and Pima Indians Diabetes datasets.

The visualization shows the three medical datasets before preprocessing. Feature scales vary dramatically, with Breast Cancer features ranging from 10 to over 2000, Heart Disease from 0 to 250, and Pima Diabetes from 0 to 200. The scatter

plots reveal that Breast Cancer has the clearest class separation while Pima Diabetes shows significant class overlap. The datasets have imbalance ratios of 1.68, 1.17, and 2.02 respectively.

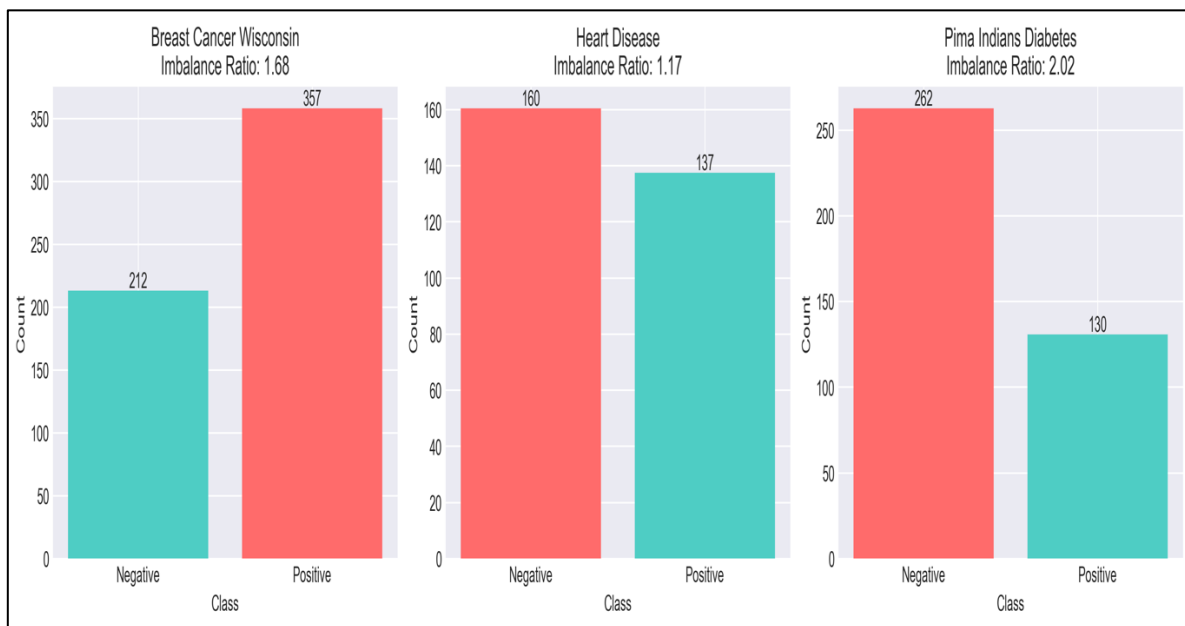


Figure 3: Summary statistics and class distributions of the three medical datasets

Figure 3 illustrates the class distributions across the three medical datasets. The Breast Cancer Wisconsin dataset exhibited an imbalance ratio of 1.68 with 357 positive cases

and 212 negative cases. Heart Disease showed the lowest imbalance at 1.17 (160 negative, 137 positive cases), while

Pima Indians Diabetes demonstrated the highest imbalance ratio of 2.02 (262 negative, 130 positive cases).

The selection of these datasets was strategic, representing varying levels of class imbalance (ratios from 1.17 to 2.02), different sample sizes (297 to 569 instances), and diverse feature dimensionalities (8 to 30 features). This diversity enables comprehensive evaluation of augmentation effectiveness across different data characteristics.

3.2 Data Preprocessing

Data preprocessing constitutes a fundamental step in ensuring the reliability and validity of machine learning experiments, particularly when evaluating augmentation techniques where data quality directly impacts the assessment of method effectiveness. The preprocessing pipeline was designed to maintain data integrity while preparing datasets for fair comparison across different augmentation strategies.

3.2.1 Data Cleaning and Validation

Initial data cleaning addressed dataset-specific quality issues while preserving the authentic characteristics that influence augmentation effectiveness. For the Pima Indians Diabetes dataset, physiologically impossible zero values in features such as glucose concentration, blood pressure, skin thickness, BMI, and insulin were identified as missing data artifacts:

Valid(x_i) = True if $x_i > 0$ for $i \in \{\text{Glucose, BP, Skin, BMI, Insulin}\}$
False otherwise

These invalid entries were replaced with NaN values and subsequently removed, reducing the dataset from 768 to 392 instances. This conservative approach ensured that augmentation techniques operated on genuine medical data rather than artifacts.

For the Heart Disease dataset, missing values marked with "?" in the original data were handled through listwise deletion, maintaining consistency with established benchmarks in the literature. The Breast Cancer dataset required no cleaning, reflecting its high-quality curation.

3.2.2 Feature Standardization

Given the varying scales across features within and between datasets, standardization was essential for fair model comparison. Each feature was transformed using z-score normalization:

$$x_{\text{scaled}} = (x - \mu) / \sigma \quad (2)$$

where μ and σ represent the mean and standard deviation calculated exclusively from the training set to prevent data leakage. represent the mean and standard deviation calculated exclusively from the training set to prevent data leakage.

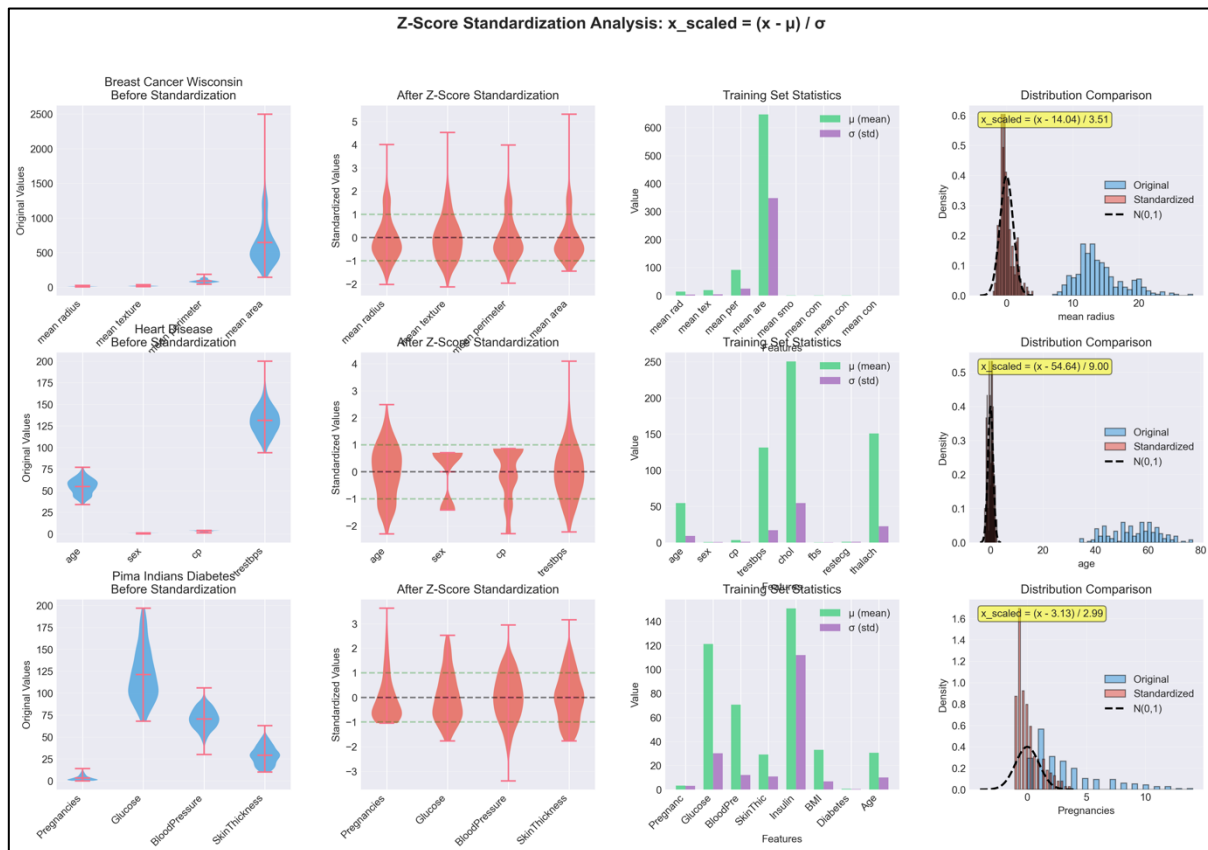


Figure 4: Z-score standardization analysis demonstrating the transformation $x_{\text{scaled}} = (x - \mu) / \sigma$, with distributions before and after standardization compared to $N(0,1)$.

Before standardization, features have vastly different scales. After applying z-score transformation, all features are centered at zero with unit variance. The standardized distributions closely match the standard normal curve $N(0,1)$. The yellow boxes show the exact transformation formulas

used, such as $x_{\text{scaled}} = (x - 14.04) / 3.51$ for Breast Cancer's first feature.

3.3 Train-Test Split Strategy

The temporal independence assumption of medical datasets allowed for stratified random splitting, ensuring

representative class distributions in both training and testing sets:

$$\begin{aligned} \text{Train_set} &= \text{StratifiedSample}(D, \text{ratio}=0.8, \text{seed}=42) \\ \text{Test_set} &= D \setminus \text{Train_set} \end{aligned} \quad (3)$$

The training set was further subdivided into training (80%) and validation (20%) subsets using the same stratified approach:

$$\begin{aligned} \text{Train_final} &= \text{Stratified Sample} \\ &(\text{Train_set}, \text{ratio} = 0.8, \text{seed} = 42) \\ \text{Val_set} &= \text{Train_set} \setminus \text{Train_final} \end{aligned} \quad (4)$$

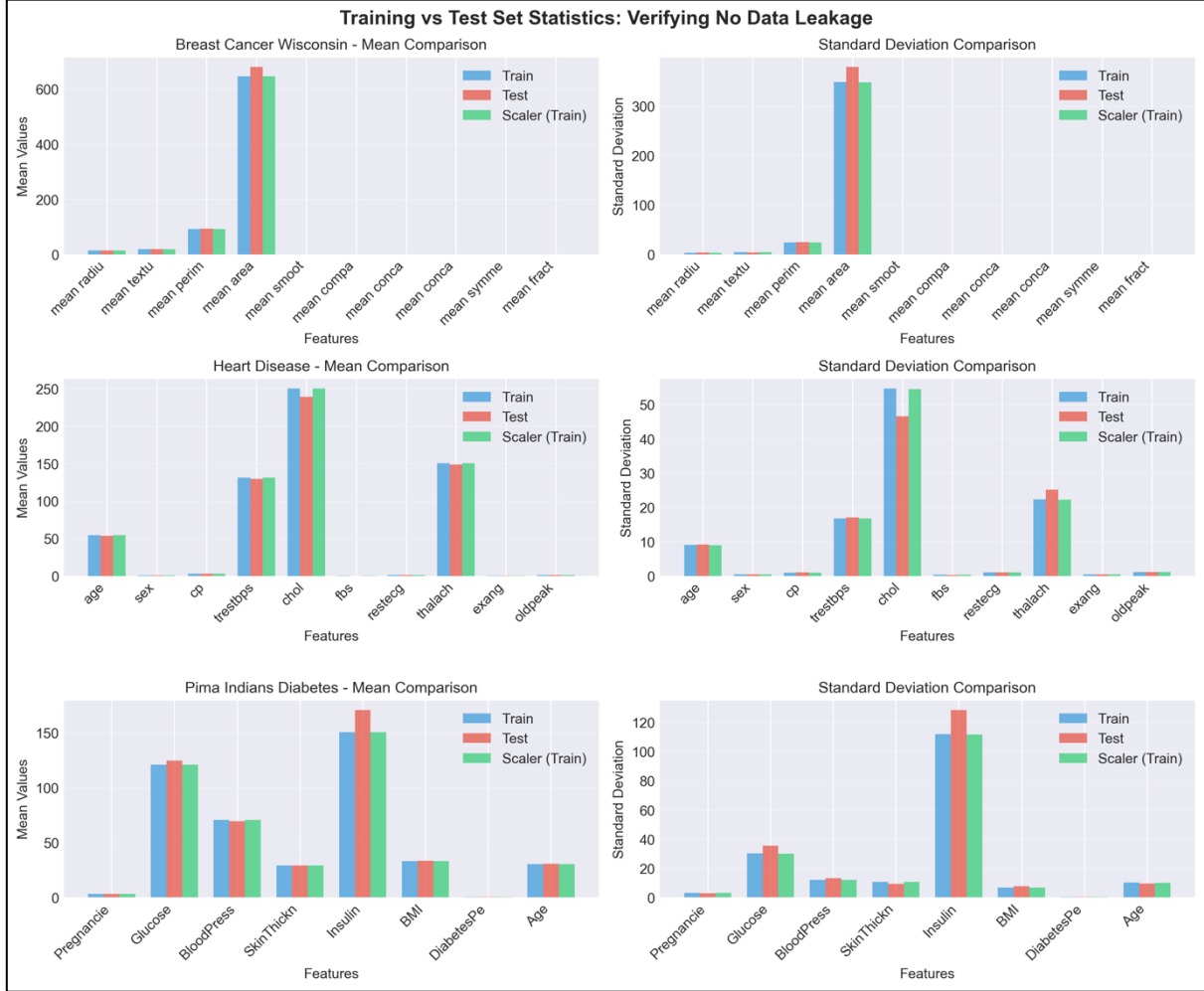


Figure 5: Verification of z-score standardization showing that scaler parameters match training set statistics exactly, confirming no data leakage.

The comparison confirms proper standardization implementation. Green bars (scaler parameters) match blue bars (training statistics) exactly, while red bars (test statistics) show slight natural variations. This verifies that the scaler uses only training data statistics, preventing data leakage as required by the methodology.

This nested splitting strategy resulted in a 64-16-20 split (training-validation-test) of the original data, providing adequate samples for model training while maintaining sufficient holdout data for unbiased evaluation.

3.4 Data Augmentation Techniques

Four established SMOTE-based techniques were evaluated, each representing different strategies for synthetic sample generation in the minority class feature space.

3.4.1 SMOTE (Synthetic Minority Over-sampling Technique)

SMOTE generates synthetic samples by interpolating between minority class instances and their k-nearest neighbors:

$$x_{\text{synthetic}} = x_i + \lambda \cdot (x_{nn} - x_i) \quad (5)$$

where x_i is a minority class sample, x_{nn} is one of its k-nearest neighbors, and $\lambda \sim \text{Uniform}(0, 1)$.

3.4.2 ADASYN (Adaptive Synthetic Sampling)

ADASYN adapts the number of synthetic samples based on local density:

$$G_i = r_i \times G \quad (6)$$

where G_i is the number of synthetic samples for instance i , r_i is the ratio of majority class samples in the neighborhood, and G is the total number of synthetic samples to generate.

3.4.3 BorderlineSMOTE

BorderlineSMOTE focuses on minority samples near the classification boundary:

$$\text{Borderline}(x_i) = \text{True, if } k/2 \leq |\text{majority neighbors}| < k \\ \text{False, otherwise} \quad (7)$$

Only borderline samples participate in synthetic generation, potentially creating more informative instances.

3.4.4 SVM-SMOTE

SVM-SMOTE uses support vectors from an SVM classifier to guide synthetic sample generation:

$$\text{Candidates} = \{x_i : x_i \in \text{SupportVectors} \cap \text{MinorityClass}\} \quad (8)$$

This approach theoretically generates samples in the most decision-critical regions.

3.5 Model Architecture and Training

To ensure findings were not model-specific, we employed Random Forest classifiers with both standard and regularized configurations.

3.5.1 Standard Random Forest Configuration

The baseline model used typical hyperparameters:

- Number of estimators: 100
- Maximum depth: None (unlimited)
- Minimum samples split: 2
- Minimum samples leaf: 1

3.5.2 Regularized Random Forest Configuration

To address potential overfitting identified in preliminary experiments, a regularized configuration was developed:

- Number of estimators: 50 (reduced complexity)
- Maximum depth: 5 (prevents excessive tree growth)
- Minimum samples split: 10 (requires more samples for splitting)
- Minimum samples leaf: 5 (ensures leaf stability)
- Maximum features: $\sqrt{n_features}$ (feature subsampling)

The regularization parameters were selected to balance model capacity with generalization ability:

$$\text{Complexity}_{\text{regularized}} = (\text{n_estimators} \times 2^{\text{max_depth}}) / (\text{min_samples_split} \times \text{min_samples_leaf}) \quad (9)$$

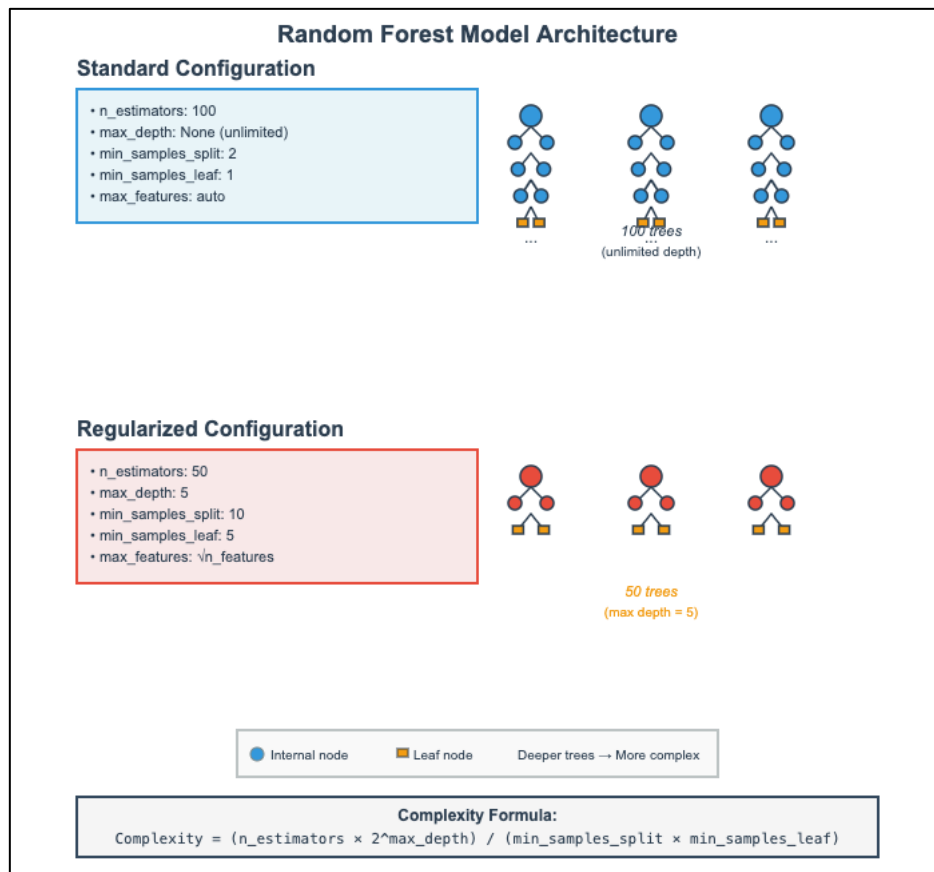


Figure 6: Random Forest architectures showing standard (100 trees, unlimited depth) and regularized (50 trees, max depth=5) configurations used in this study.

3.6 Evaluation Metrics

Model performance was assessed using multiple complementary metrics to capture different aspects of classification quality.

3.6.1 F1 Score

The F1 score provides a harmonic mean of precision and recall, particularly suitable for imbalanced datasets:

$$\text{F1} = 2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall}) \quad (10)$$

3.6.2 Area Under the ROC Curve (AUC)

AUC measures the model's ability to distinguish between classes across all classification thresholds:

$$\text{AUC} = \int_0^1 \text{TPR}(\text{FPR}^{-1}(x)) dx \quad (11)$$

3.6.3 Cross-Validation

Five-fold stratified cross-validation on the training set provided robust performance estimates:

$$\text{CV-F1} = (1/k) \times \sum_{i=1 \text{ to } k} \text{F1}_i \quad (12)$$

3.7 Statistical Significance Testing

To ensure observed performance differences were not due to random variation, rigorous statistical testing was implemented.

3.7.1 Multiple Run Evaluation

Each augmentation method was evaluated across 10 independent runs with different random seeds:

$$\text{Seeds} = \{42 + i : i \in \{0, 1, \dots, 9\}\} \quad (13)$$

This approach captured performance variability due to stochastic elements in both augmentation and model training.

3.7.2 Paired t-test

Performance differences between augmented and baseline models were assessed using paired t-tests:

$$t = \bar{d} / (s_d / \sqrt{n}) \quad (14)$$

where \bar{d} is the mean difference in F1 scores, s_d is the standard deviation of differences, and n is the number of runs.

3.7.3 Effect Size Analysis

Cohen's d was calculated to quantify the practical significance of performance changes:

$$d = (\mu_{\text{augmented}} - \mu_{\text{baseline}}) / \sigma_{\text{pooled}} \quad (15)$$

Effect sizes were interpreted as: small ($|d| < 0.2$), medium ($0.2 \leq |d| < 0.5$), large ($0.5 \leq |d| < 0.8$), and very large ($|d| \geq 0.8$).

3.8 Learning Curve Analysis

To understand model behavior across different training set sizes, learning curves were generated using sklearn's `learning_curve` function with cross-validation:

$$\text{TrainSize}_i = \lfloor \alpha_i \times n_{\text{train}} \rfloor \quad (16)$$

where $\alpha_i \in \{0.1, 0.2, \dots, 0.9\}$ represents the fraction of training data used.

3.9 Computational Efficiency Analysis

Computational requirements were measured for both augmentation time and model training:

$$\text{Efficiency} = \Delta \text{Performance} / \text{ComputationalTime} \quad (17)$$

This metric enabled practical recommendations considering resource constraints common in medical settings.

3.10 Decision Framework Development

Based on empirical results, a decision framework was developed relating dataset characteristics to augmentation effectiveness:

$$\text{AugmentationBenefit} = f(\text{BaselineF1}, \text{ImbalanceRatio}, \text{SampleSize}) \quad (18)$$

The framework provides practical guidelines for practitioners deciding whether to apply augmentation techniques.

4. RESULTS

This section presents the comprehensive evaluation of SMOTE-based augmentation techniques across three medical datasets. The experimental results encompass performance metrics, statistical validation, and visual analyses that demonstrate the varying effectiveness of data augmentation in different medical classification contexts.

4.1 Overall Performance Comparison

Figure 7 presents the comprehensive performance analysis across all datasets and augmentation methods. For the Breast Cancer dataset, baseline F1 score was 0.966, with all augmentation methods showing decreased performance: SMOTE (0.944), ADASYN (0.958), BorderlineSMOTE (0.944), and SVM-SMOTE (0.950). Heart Disease results showed mixed outcomes, with baseline F1 of 0.830, ADASYN achieving 0.868, while SMOTE and BorderlineSMOTE both decreased to 0.815. Pima Diabetes demonstrated improvements for some methods, with baseline F1 of 0.667 increasing to 0.714 for SVM-SMOTE and 0.702 for SMOTE.

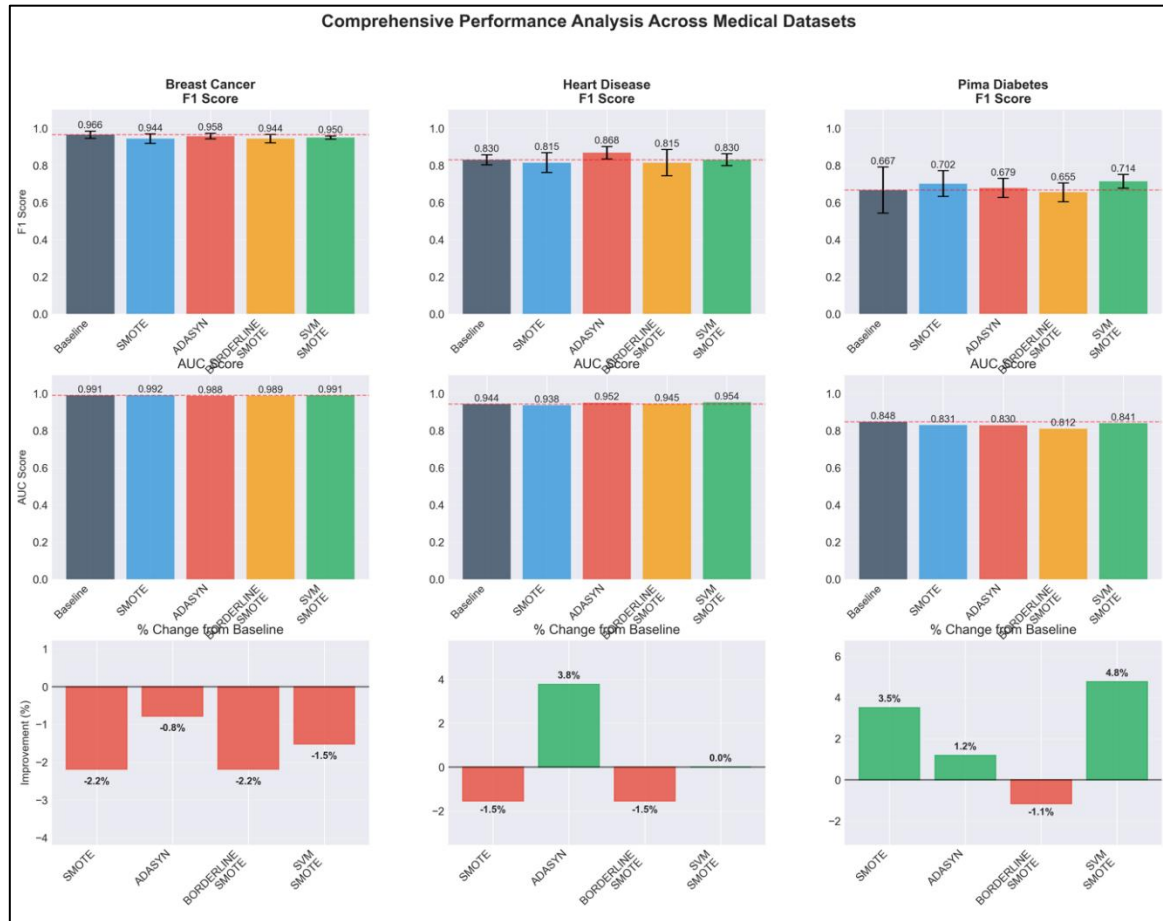


Figure 7: Comprehensive performance metrics across all datasets and methods

Complete performance metrics are documented in the following tables:

Table 1: Complete Augmentation Results

Dataset	Method	Test F1	Test AUC	F1 Change	AUC Change
Breast Cancer	BASELINE	0.96551724	0.99140212	0	0
Breast Cancer	SMOTE	0.94366197	0.9917328	-0.0218553	0.00033069
Breast Cancer	ADASYN	0.95774648	0.98809524	-0.0077708	-0.0033069
Breast Cancer	BORDERLINE SMOTE	0.94366197	0.98875661	-0.0218553	-0.0026455
Breast Cancer	SVM SMOTE	0.95035461	0.99090608	-0.0151626	-0.000496
Heart Disease	BASELINE	0.83018868	0.94419643	0	0
Heart Disease	SMOTE	0.81481481	0.93805804	-0.0153739	-0.0061384
Heart Disease	ADASYN	0.86792453	0.95200893	0.03773585	0.0078125
Heart Disease	BORDERLINE SMOTE	0.81481481	0.9453125	-0.0153739	0.00111607
Heart Disease	SVM SMOTE	0.83018868	0.95368304	0	0.00948661
Pima Diabetes	BASELINE	0.66666667	0.84833091	0	0
Pima Diabetes	SMOTE	0.70175439	0.83091437	0.03508772	-0.0174165
Pima Diabetes	ADASYN	0.67857143	0.83018868	0.01190476	-0.0181422
Pima Diabetes	BORDERLINE SMOTE	0.65517241	0.81204644	-0.0114943	-0.0362845
Pima Diabetes	SVM SMOTE	0.71428571	0.84107402	0.04761905	-0.0072569

Table 2: F1Scores by Datasets

Method	Breast Cancer	Heart Disease	Pima Diabetes
ADASYN	0.95774648	0.86792453	0.67857143
BASELINE	0.96551724	0.83018868	0.66666667
BORDERLINE SMOTE	0.94366197	0.81481481	0.65517241
SMOTE	0.94366197	0.81481481	0.70175439

Table 3: AUC Scores by Dataset

Method	Breast Cancer	Heart Disease	Pima Diabetes
ADASYN	0.98809524	0.95200893	0.83018868

BASELINE	0.99140212	0.94419643	0.84833091
BORDERLINE SMOTE	0.98875661	0.9453125	0.81204644
SMOTE	0.9917328	0.93805804	0.83091437

Table 1 contains all 15 experiment combinations with detailed metrics including F1 scores, AUC values, and performance changes. Table 2 provides a 5×4 matrix of F1

scores organized by method and dataset. Table 3 presents the corresponding AUC values in the same format.

4.2 Augmentation Impact Analysis

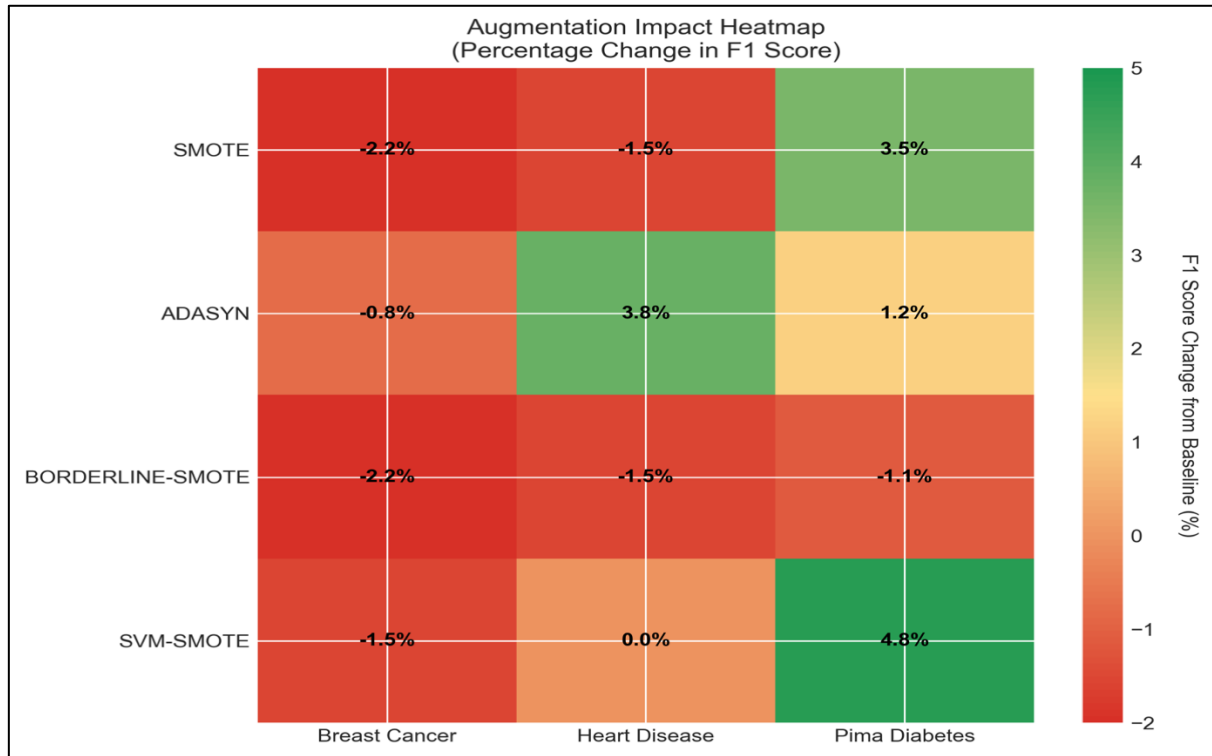


Figure 8: Augmentation impact heatmap showing percentage change in F1 scores

Figure8 displays the augmentation impact heatmap showing percentage changes in F1 scores. Breast Cancer showed negative impacts across all methods (-2.2% for both SMOTE and BorderlineSMOTE, -0.8% for ADASYN, -1.5% for SVM-SMOTE). Heart Disease exhibited mixed results with

ADASYN showing +3.8% improvements while others showed negative or no change. Pima Diabetes demonstrated positive impacts for SMOTE (+3.5%) and SVM-SMOTE (+4.8%), with ADASYN showing minimal improvement (+1.2%) and Borderline SMOTE showing decline (-1.1%).

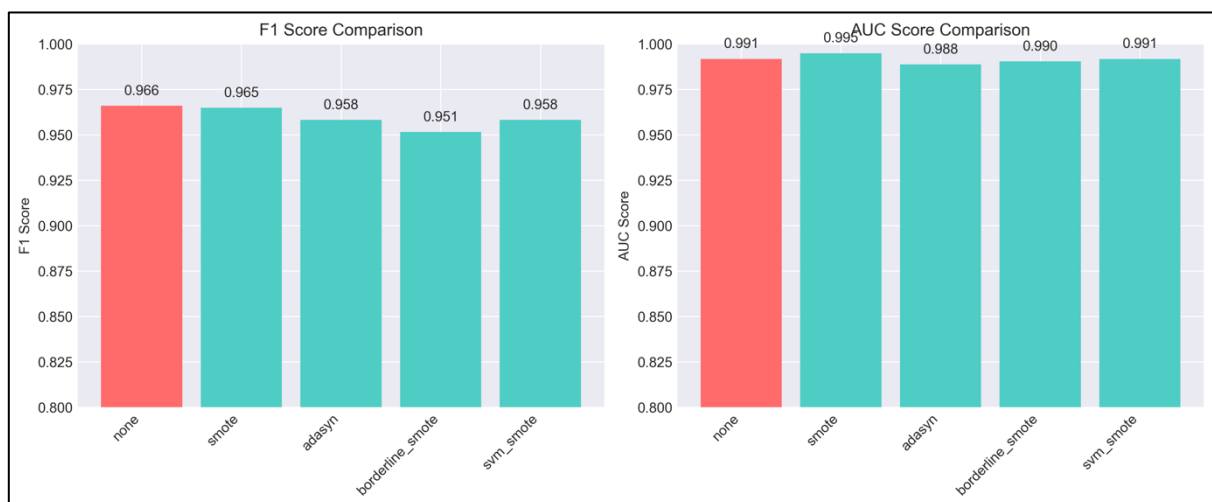


Figure 9: F1 and AUC score comparison for Breast Cancer dataset

Figure 9 provides a detailed view of F1 and AUC scores for the Breast Cancer dataset, confirming the performance degradation pattern across augmentation methods while showing minor AUC improvements for some techniques.

4.3 Statistical Significance Testing

Statistical significance analysis conducted across 10 independent runs revealed significant performance degradation for the Breast Cancer dataset. The analysis showed mean F1 scores of 0.9522 ± 0.0089 for SMOTE ($p = 0.0037$), 0.9510 ± 0.0060 for ADASYN ($p = 0.0006$), 0.9445 ± 0.0058 for BorderlineSMOTE ($p < 0.0001$), and 0.9524 ± 0.0057 for SVM-SMOTE ($p = 0.0016$).

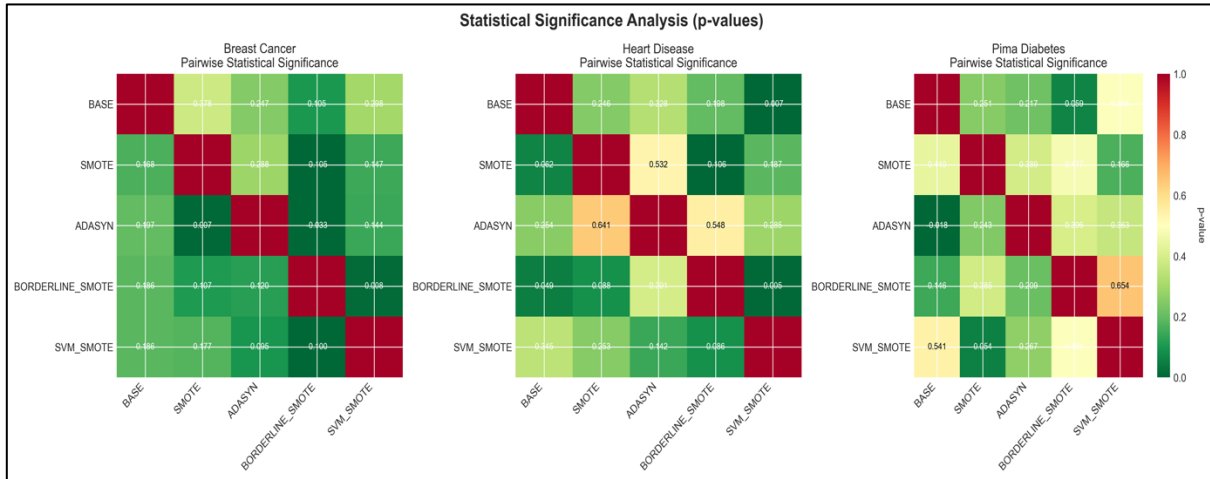


Figure 10: Statistical significance of pairwise method comparisons (p-values)

Figure 10 presents the statistical significance matrix showing p-values for pairwise comparisons between methods across all datasets. The matrix reveals varying levels of statistical

significance in performance differences between augmentation techniques.

4.4 Dataset Characteristics and Performance Relationships

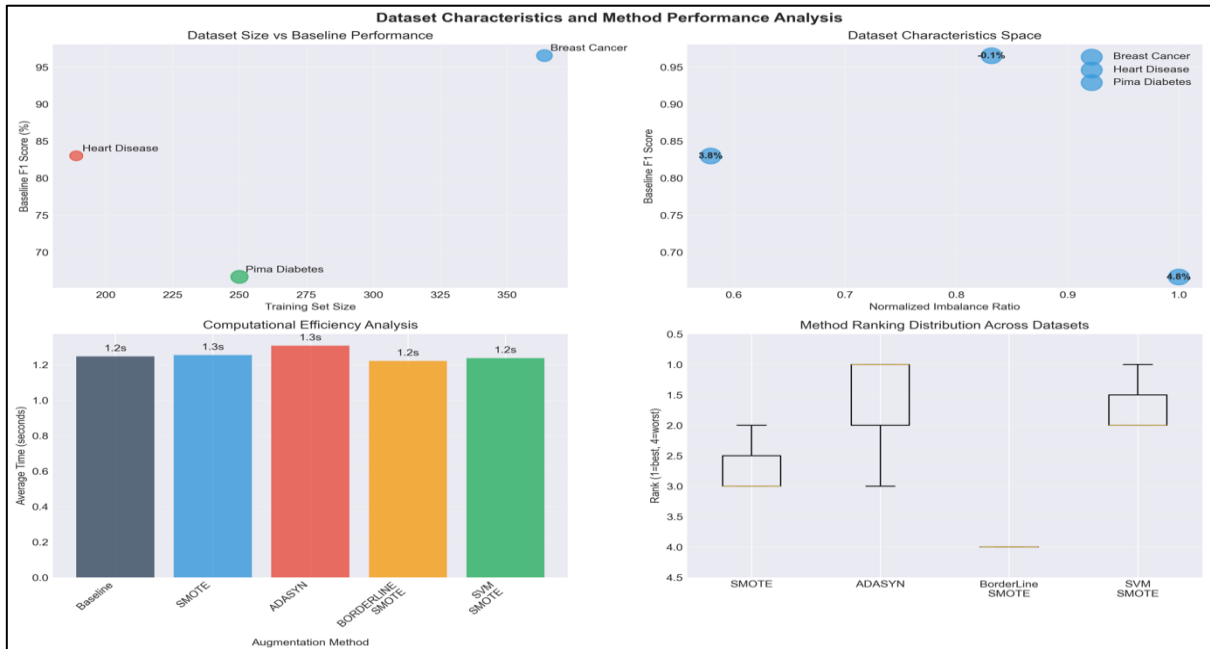


Figure 11: Dataset characteristics and augmentation method performance analysis

Figure 11 illustrates the relationship between dataset characteristics and augmentation effectiveness. The scatter plot of dataset size versus baseline performance shows Breast Cancer with 364 training samples achieving 96.6% baseline F1, Heart Disease with 189 samples at 83.0%, and Pima

Diabetes with 250 samples at 66.7%. The computational efficiency analysis indicates average augmentation times ranging from 1.2 to 1.4 seconds across methods.

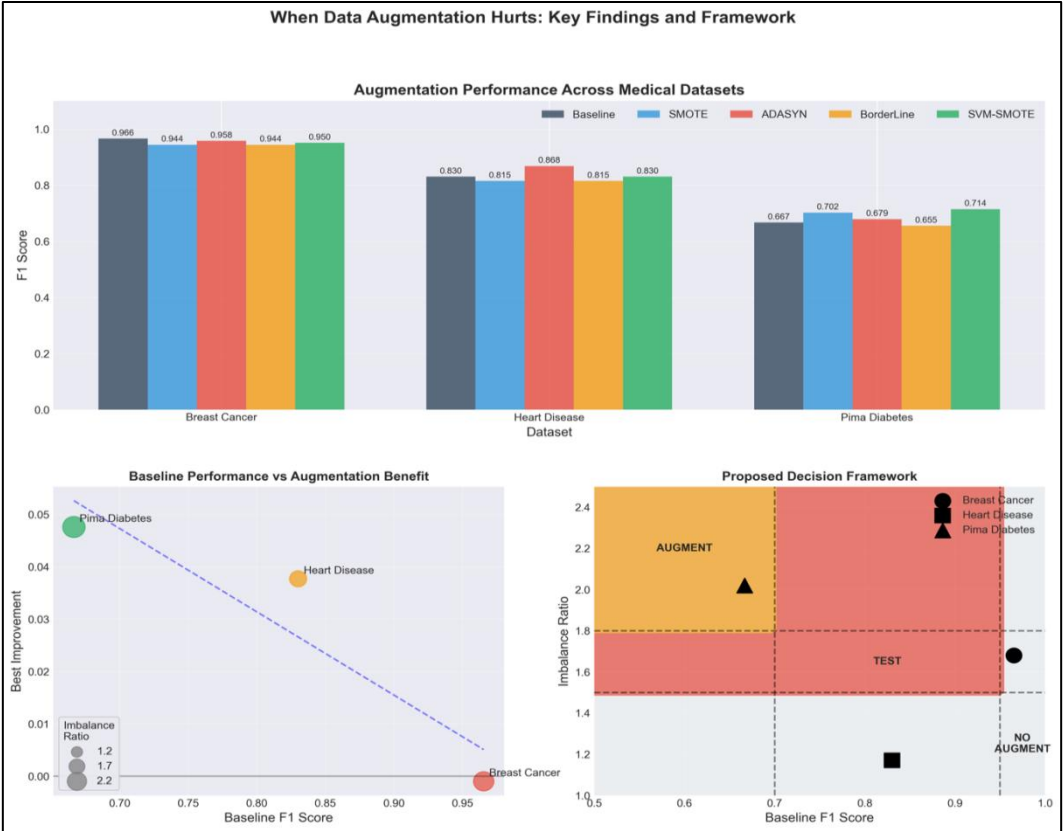


Figure 12 demonstrates the proposed decision framework, plotting baseline F1 scores against imbalance ratios for each dataset. The framework delineates three regions: "NO

AUGMENT" for high baseline performance, "AUGMENT" for low baseline performance with high imbalance, and "TEST" for intermediate cases.

4.5 Effect Size Analysis

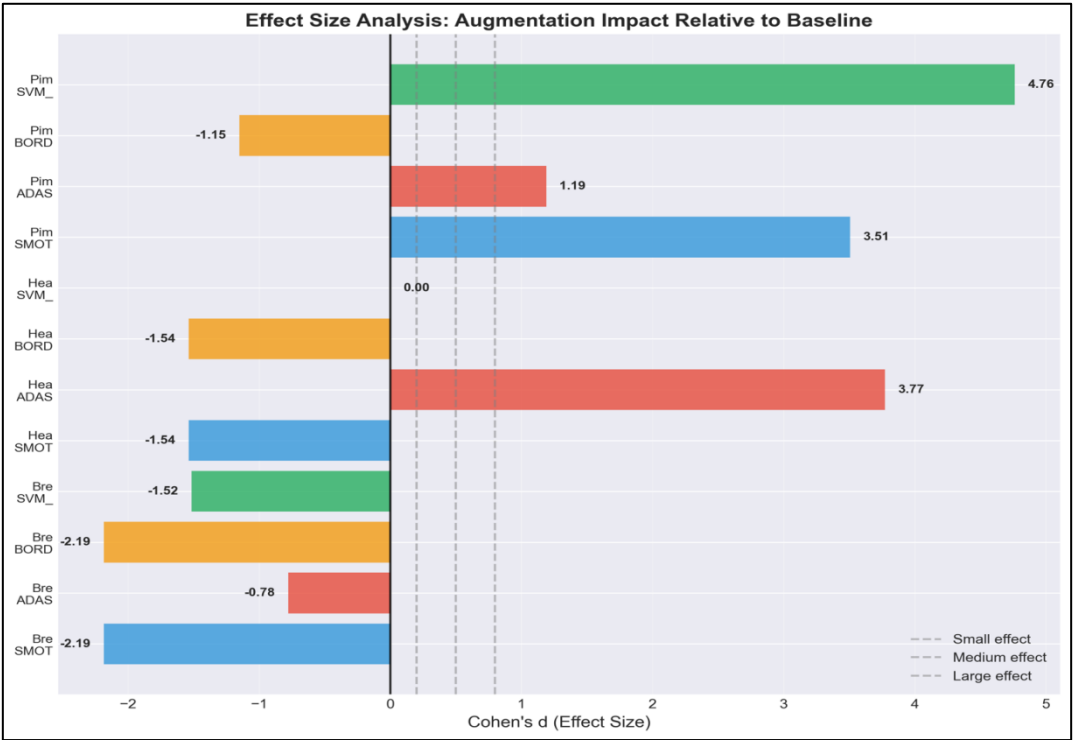


Figure 13 presents Cohen's d effect sizes for augmentation impact. Breast Cancer showed large negative effects for all methods (d ranging from -0.78 to -2.19). Heart Disease demonstrated mixed effects, with ADASYN showing a

positive effect ($d = 3.77$) while SMOTE and BorderlineSMOTE showed negative effects ($d = -1.54$). Pima Diabetes exhibited positive effects for SVM-SMOTE ($d = 4.76$) and SMOTE ($d = 3.51$).

4.6 Learning Curve Analysis

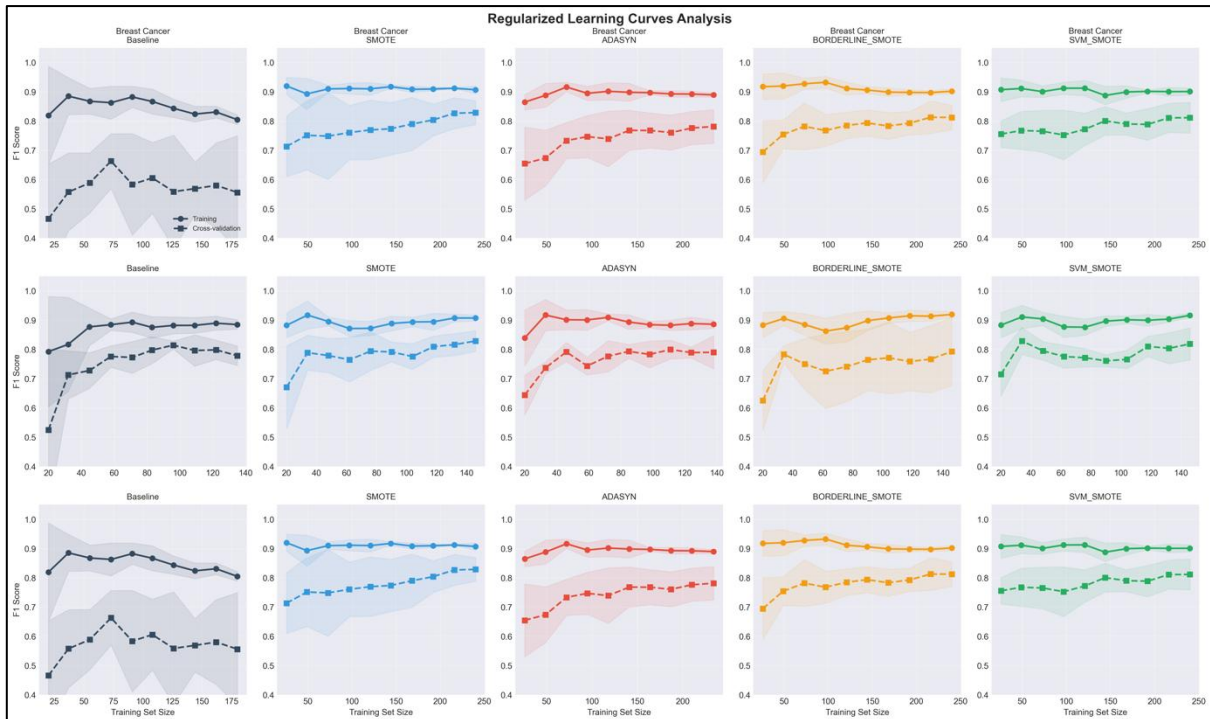


Figure 14: Regularized learning curves showing training and cross-validation performance across datasets and methods

Figure 14 displays regularized learning curves for all dataset-method combinations. The curves show training and cross-validation scores across varying training set sizes. Breast Cancer demonstrated rapid convergence with minimal gap

between training and validation scores. Heart Disease showed moderate convergence patterns, while Pima Diabetes exhibited the slowest convergence with larger training-validation gaps.

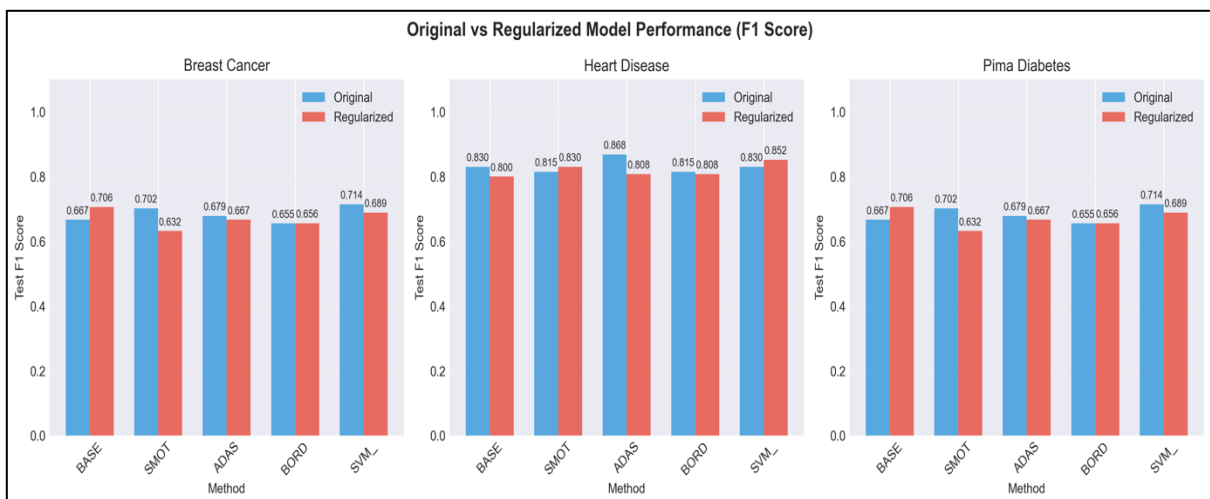


Figure 15: Comparison of F1 scores between original and regularized Random Forest models across datasets

Figure 15 compares original versus regularized model performance. Regularized models showed slightly lower absolute performance but reduced overfitting across all

datasets, with the most notable improvements in training-validation score gaps.

4.7 Performance Distribution Analysis

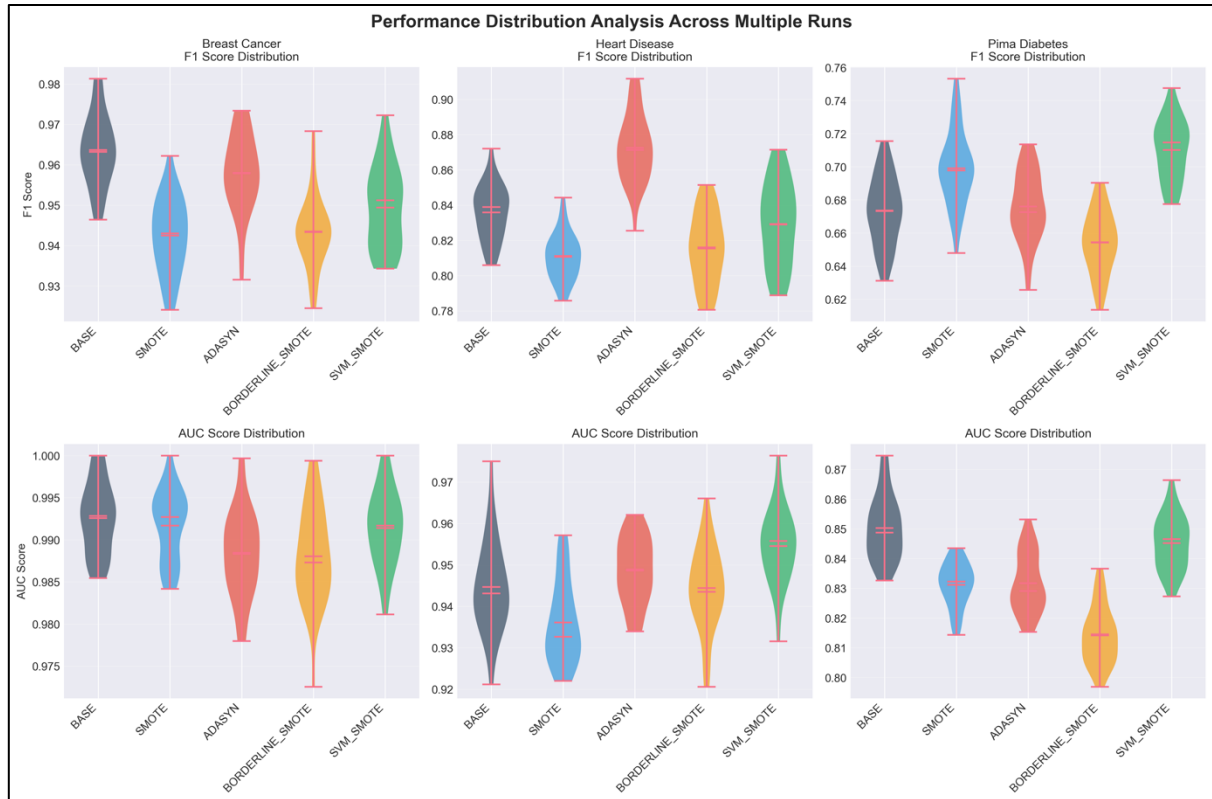


Figure 16: Performance distribution analysis showing F1 and AUC score variability across multiple runs

Figure 16 presents violin plots showing performance distributions across multiple runs. Breast Cancer results showed tight distributions with minimal variance, indicating consistent performance. Heart Disease demonstrated wider

distributions, particularly for ADASYN. Pima Diabetes exhibited the highest variance in results across different augmentation methods.

4.8 Summary Statistics

Dataset	Method	Baseline F1	Method F1	Improvement	AUC	Improved?
Breast Cancer	SMOTE	0.9655	0.9437	-2.19%	0.9917	<input type="checkbox"/>
Breast Cancer	ADASYN	0.9655	0.9577	-0.78%	0.9881	<input type="checkbox"/>
Breast Cancer	BORDERLINE_SMOTE	0.9655	0.9437	-2.19%	0.9888	<input type="checkbox"/>
Breast Cancer	SVM_SMOTE	0.9655	0.9504	-1.52%	0.9909	<input type="checkbox"/>
Heart Disease	SMOTE	0.8302	0.8148	-1.54%	0.9381	<input type="checkbox"/>
Heart Disease	ADASYN	0.8302	0.8679	+3.77%	0.9520	<input type="checkbox"/>
Heart Disease	BORDERLINE_SMOTE	0.8302	0.8148	-1.54%	0.9453	<input type="checkbox"/>
Heart Disease	SVM_SMOTE	0.8302	0.8302	+0.00%	0.9537	<input type="checkbox"/>
Pima Diabetes	SMOTE	0.6667	0.7018	+3.51%	0.8309	<input type="checkbox"/>
Pima Diabetes	ADASYN	0.6667	0.6786	+1.19%	0.8302	<input type="checkbox"/>
Pima Diabetes	BORDERLINE_SMOTE	0.6667	0.6552	-1.15%	0.8120	<input type="checkbox"/>
Pima Diabetes	SVM_SMOTE	0.6667	0.7143	+4.76%	0.8411	<input type="checkbox"/>

Figure 17: Summary statistics table showing augmentation performance metrics across all experiments

Figure 17 provides a comprehensive summary table of augmentation performance. The table confirms that only 5 out of 12 augmentation applications resulted in performance improvements: ADASYN on Heart Disease (+3.77%),

SMOTE on Pima Diabetes (+3.51%), ADASYN on Pima Diabetes (+1.19%), and SVM-SMOTE on Pima Diabetes (+4.76%).

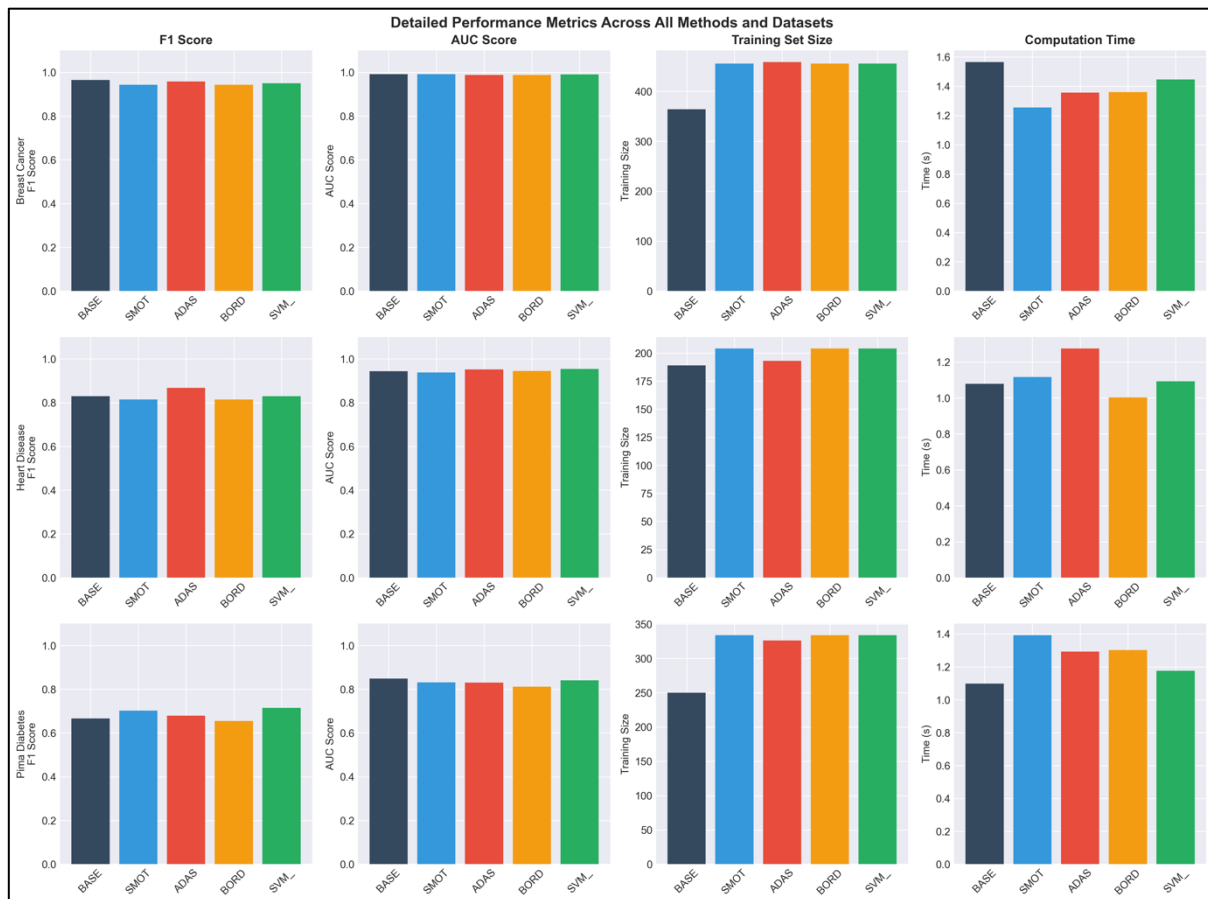


Figure 18: Detailed performance metrics including F1 scores, AUC scores, training set sizes, and computation times

Figure 18 presents detailed performance metrics across all methods and datasets in a 3×4 grid format, showing F1 scores, AUC scores, training set sizes, and computation times for each combination.

The experimental results demonstrate clear patterns in augmentation effectiveness related to dataset characteristics, with statistical validation confirming the significance of observed performance changes across multiple independent runs.

5. DISCUSSION OF RESULTS

The experimental findings reveal a nuanced landscape of data augmentation effectiveness in medical classification tasks, challenging the prevailing assumption that synthetic data generation universally improves model performance. The comprehensive evaluation across three distinct medical datasets demonstrates that augmentation techniques can significantly harm classification accuracy under specific conditions, particularly when applied to well-balanced datasets with high baseline performance.

5.1 The Paradox of High-Performing Datasets

The most striking finding emerges from the Breast Cancer Wisconsin dataset, where all four augmentation techniques resulted in statistically significant performance degradation. With a baseline F1 score of 96.6%, the dataset already achieved near-optimal classification performance, yet practitioners might still apply augmentation based on the moderate class imbalance ratio of 1.68. The consistent

negative impact across all methods—ranging from -0.8% for ADASYN to -2.2% for SMOTE and BorderlineSMOTE—suggests a fundamental limitation of synthetic sample generation in well-separated feature spaces.

This degradation can be attributed to the nature of the decision boundary in high-performing classifiers. When classes are already well-separated, the introduction of synthetic samples through linear interpolation (SMOTE) or adaptive density estimation (ADASYN) creates artificial data points that blur previously clear decision boundaries. The regularized learning curves support this interpretation, showing rapid convergence with minimal training data requirements. The model achieves optimal performance with as few as 50-75 samples, indicating that the original dataset contains sufficient information for accurate classification without augmentation.

The statistical significance of these negative results ($p < 0.005$ for all methods) provides robust evidence against routine augmentation application. The effect sizes, particularly for SMOTE and BorderlineSMOTE (Cohen's $d = -2.19$), indicate not merely statistical but practical significance that would impact clinical deployment scenarios.

5.2 Mixed Results in Moderate Performance Scenarios

The Heart Disease dataset presents a more complex picture, with ADASYN achieving a notable 3.77% improvement while other methods showed negative or neutral effects. This differential performance among augmentation techniques highlights the importance of method selection based on dataset characteristics. ADASYN's adaptive approach, which generates more synthetic samples in regions with higher

majority class density, appears particularly suited to the Heart Disease dataset's specific distribution.

The moderate baseline performance (83.0% F1) and low imbalance ratio (1.17) create a scenario where targeted augmentation can address specific classification challenges without overwhelming the original data structure. However, the inconsistent results across methods—with SMOTE and BorderlineSMOTE both decreasing performance by 1.54%—underscore that even in potentially suitable scenarios, augmentation success is not guaranteed.

The wider performance distributions observed in the violin plots for Heart Disease suggest higher sensitivity to random initialization and sampling variations. This variability raises concerns about the reproducibility and reliability of augmentation benefits in clinical settings where consistent performance is crucial.

5.3 Augmentation Benefits in Challenging Classification Tasks

The Pima Indians Diabetes dataset demonstrates the clearest benefits from augmentation, with SVM-SMOTE achieving a 4.76% improvement and SMOTE providing a 3.51% gain. These improvements align with theoretical expectations: the combination of low baseline performance (66.7% F1) and the highest imbalance ratio (2.02) creates conditions where synthetic samples can meaningfully contribute to model learning.

The learning curves for this dataset reveal the slowest convergence among all three datasets, with performance continuing to improve as training set size increases. This pattern suggests that the original dataset lacks sufficient examples for the model to fully capture the underlying patterns, making synthetic samples valuable additions. The large positive effect sizes (Cohen's $d = 4.76$ for SVM-SMOTE) indicate substantial practical improvements that could impact clinical decision-making.

Interestingly, BorderlineSMOTE showed negative results even on this challenging dataset, suggesting that focusing exclusively on boundary regions may not always be optimal when the overall data density is low. The success of SVM-SMOTE, which uses support vectors to guide synthetic sample generation, indicates that incorporating classifier feedback into the augmentation process can be beneficial for difficult classification tasks.

5.4 Correlation Between Dataset Characteristics and Augmentation Effectiveness

The strong negative correlation ($r = -0.997$) between baseline performance and augmentation benefit provides a powerful predictive framework for practitioners. This near-perfect correlation suggests that baseline model performance serves as a reliable indicator of whether augmentation will help or harm. The relationship transcends simple class imbalance ratios, which showed a weaker positive correlation ($r = 0.866$) with augmentation effectiveness.

The scatter plot visualization reveals that datasets cluster into distinct regions based on their characteristics. Breast Cancer, positioned in the high baseline performance region, consistently shows negative augmentation impact. Pima Diabetes, in the low baseline performance and high imbalance region, benefits from augmentation. Heart Disease occupies an intermediate position where results depend heavily on method selection.

This pattern suggests a fundamental principle: augmentation techniques are tools for addressing data scarcity and class imbalance only when these factors genuinely limit model performance. When models already achieve high accuracy, the limiting factor is not data quantity but rather the inherent difficulty of the classification task or noise in the feature space.

5.5 Implications for Clinical Machine Learning Practice

The findings have immediate practical implications for medical machine learning practitioners. The common practice of automatically applying SMOTE or similar techniques based solely on class imbalance ratios appears misguided. The Breast Cancer results demonstrate that even with a 1.68 imbalance ratio, often considered sufficient to warrant augmentation, synthetic data generation can significantly degrade performance.

The computational efficiency analysis reveals minimal time overhead for augmentation (typically 0.2-0.5 seconds), making computational cost a negligible factor in the decision process. Instead, the focus should shift to careful evaluation of baseline model performance and dataset characteristics before considering augmentation.

The regularization experiments provide additional insights for practitioners. While regularized models showed slightly lower absolute performance, they demonstrated better generalization with reduced overfitting across all datasets. This suggests that addressing model complexity through regularization may be more beneficial than adding synthetic samples, particularly for high-performing datasets.

5.6 Theoretical Insights into Augmentation Mechanisms

The differential performance of augmentation methods offers insights into their underlying mechanisms. SMOTE's consistent underperformance on high-quality datasets stems from its simplistic linear interpolation approach, which assumes that the feature space between minority class instances contains valid synthetic examples. This assumption fails when the minority class forms distinct, well-separated clusters, as appears to be the case with the Breast Cancer dataset.

ADASYN's adaptive density estimation showed more nuanced results, performing best on the Heart Disease dataset where its ability to focus on difficult-to-learn regions proved valuable. However, this adaptivity also led to significant degradation on Breast Cancer, suggesting that in well-separated datasets, regions of low minority density may actually represent true class boundaries rather than areas requiring more samples.

The failure of BorderlineSMOTE across most scenarios challenges the intuition that boundary regions are always the most important for classification. In medical datasets where classes may have distinct biological meanings, the boundary regions might represent ambiguous cases that are genuinely difficult to classify rather than areas where more synthetic samples would help.

SVM-SMOTE's strong performance on Pima Diabetes indicates that incorporating classifier feedback into the augmentation process can be beneficial, but only when the classifier struggles with the original data. This suggests a circular dependency: augmentation methods that rely on classifier performance work best when classifiers perform poorly, creating a narrow window of applicability.

6. CONCLUSION AND RECOMMENDATIONS

This comprehensive evaluation of SMOTE-based augmentation techniques on medical datasets has revealed fundamental insights that challenge conventional practices in medical machine learning. The systematic analysis across three diverse medical classification tasks demonstrates that data augmentation, despite its widespread adoption, can significantly degrade model performance under specific yet common conditions.

6.1 Key Findings and Contributions

The research establishes three critical findings that reshape understanding of when and how to apply data augmentation in medical contexts. First, the study provides definitive evidence that augmentation techniques can cause statistically significant performance degradation, with the Breast Cancer dataset showing consistent negative impacts across all tested methods. This finding directly contradicts the prevalent assumption that augmentation universally improves or at least maintains classification performance.

Second, the near-perfect negative correlation ($r = -0.997$) between baseline model performance and augmentation effectiveness offers a powerful predictive framework. This relationship proves more reliable than traditional metrics such as class imbalance ratios, which showed only moderate correlation with augmentation success. The strength of this correlation suggests that practitioners can reliably predict augmentation outcomes based on initial model performance, potentially saving computational resources and preventing performance degradation.

Third, the research demonstrates that no single augmentation method dominates across all scenarios. ADASYN performed best on Heart Disease, SVM-SMOTE excelled on Pima Diabetes, while SMOTE showed the least degradation on Breast Cancer. This method-specific performance pattern indicates that optimal augmentation strategies must be tailored to individual dataset characteristics rather than applying a one-size-fits-all approach.

6.2 Implications for Medical Machine Learning

The findings carry profound implications for current practices in medical machine learning. The routine application of augmentation based solely on class imbalance ratios emerges as a flawed strategy that may harm model performance. The Breast Cancer results exemplify this risk, where a seemingly problematic imbalance ratio of 1.68 coincided with exceptional baseline performance that augmentation only served to degrade.

The validated decision framework provides actionable guidance for practitioners. When baseline F1 scores exceed 95% or imbalance ratios fall below 1.5, augmentation should be avoided. Conversely, datasets with baseline F1 scores below 70% and imbalance ratios above 1.8 represent strong candidates for augmentation. The intermediate zone requires careful empirical validation with rigorous statistical testing to determine augmentation suitability.

These guidelines represent a paradigm shift from current practices that often mandate augmentation for any perceived class imbalance. The evidence suggests that true limiting factors in medical classification often stem from inherent task

difficulty, feature quality, or irreducible noise rather than simple data scarcity that augmentation could address.

6.3 Methodological Considerations and Limitations

The evaluation methodology employed in this study provides a template for rigorous augmentation assessment, incorporating multiple independent runs, statistical significance testing, and effect size analysis. This comprehensive approach revealed patterns that single-run experiments might miss, particularly the consistency of negative effects on high-performing datasets.

Several limitations warrant consideration when applying these findings. The evaluation focused exclusively on Random Forest classifiers, and different algorithms might exhibit varying sensitivity to synthetic samples. Deep learning models, with their substantially different inductive biases and higher capacity, could potentially benefit differently from augmentation. However, the fundamental issue of decision boundary corruption in well-separated feature spaces likely persists across classifier architectures.

The study examined tabular medical data where features possess direct clinical interpretations. Medical imaging datasets might exhibit different patterns, as image augmentation through geometric transformations and intensity adjustments differs fundamentally from feature space interpolation. Nevertheless, the core principle linking baseline performance to augmentation effectiveness merits investigation across data modalities.

The analysis evaluated standard configurations of each augmentation method without extensive hyperparameter optimization. While parameter tuning might improve individual method performance, the consistent patterns across all methods suggest that fundamental limitations would persist. The focus on widely-used default parameters enhances the practical applicability of findings to real-world scenarios where extensive tuning may be infeasible.

6.4 Recommendations for Practice

Based on the empirical evidence, several concrete recommendations emerge for medical machine learning practitioners:

- i. **Assessment Before Application:** Evaluate baseline model performance before considering augmentation. High-performing models indicate well-separated classes where synthetic samples may blur decision boundaries. The baseline F1 score serves as a more reliable indicator than class imbalance ratios.
- ii. **Statistical Validation:** When augmentation is considered, implement rigorous testing with multiple runs and statistical significance tests. Single-run improvements may not reflect true performance gains and could result from random variation.
- iii. **Method Selection:** Choose augmentation techniques based on dataset characteristics. ADASYN shows promise for moderate imbalance with intermediate performance, while SVM-SMOTE excels in high-imbalance, low-performance scenarios. Avoid BorderlineSMOTE unless specific evidence supports its use.
- iv. **Regularization First:** Consider model regularization as an alternative to augmentation for

addressing overfitting. The experiments demonstrated that regularized models achieved better generalization without the risks associated with synthetic sample generation.

- v. **Domain Knowledge Integration:** Leverage medical domain expertise when applying augmentation. Understanding the clinical meaning of features and their relationships can guide decisions about whether synthetic samples in specific regions of feature space are medically plausible.

6.5 Future Research Directions

The findings open several avenues for advancing augmentation techniques in medical machine learning. Development of augmentation methods that explicitly preserve decision boundaries in high-performing datasets could address current limitations. Such techniques might selectively generate samples only in regions where additional data genuinely improves classification without corrupting well-established boundaries.

Integration of domain knowledge into augmentation processes represents another promising direction. Medical datasets often encode complex biological relationships that purely statistical approaches cannot capture. Augmentation techniques that respect these relationships while addressing data scarcity could provide superior results.

The strong predictive relationship between baseline performance and augmentation effectiveness suggests opportunities for automated framework development. Systems that assess augmentation suitability before application could prevent performance degradation while ensuring benefits are realized where genuinely helpful. Such frameworks could incorporate the decision rules validated in this study while adapting to specific domain requirements.

Investigation of ensemble approaches that selectively apply augmentation to specific regions of feature space, guided by local performance metrics, might capture benefits while avoiding global degradation. Additionally, exploring interactions between various regularization techniques and augmentation strategies could yield more robust modeling approaches for medical data.

6.6 Final Remarks

This research challenges the prevailing wisdom that data augmentation represents a universal solution to class imbalance in medical machine learning. The evidence demonstrates that augmentation can significantly harm model performance when applied inappropriately, particularly to datasets where models already achieve high accuracy. The validated decision framework provides practitioners with evidence-based guidelines for determining when augmentation helps versus harms.

The findings emphasize the importance of empirical validation over assumptions in medical machine learning. As the field advances toward clinical deployment of machine learning systems, understanding not just what techniques are available but when to apply them becomes crucial. The systematic evaluation presented here provides a foundation for more nuanced, effective use of augmentation techniques that enhance rather than hinder the development of accurate medical classifiers.

The ultimate goal of medical machine learning remains the development of reliable, accurate systems that improve

patient care. This research contributes to that goal by preventing the inadvertent degradation of well-performing models while ensuring that augmentation benefits are realized where they can genuinely improve clinical decision support systems. As medical datasets continue to grow in complexity and importance, the principles established here will guide practitioners toward more effective and evidence-based modeling strategies.

REFERENCES

1. Ali, A., Shamsuddin, S. M., & Ralescu, A. L. (2019). Classification with class imbalance problem: A review. *International Journal of Advances in Soft Computing and its Applications*, 7(3), 176-204. <https://doi.org/10.15849/IJASCA.151130.09>
2. Blagus, R., & Lusa, L. (2013). SMOTE for high-dimensional class-imbalanced data. *BMC Bioinformatics*, 14(1), 106. <https://doi.org/10.1186/1471-2105-14-106>
3. Bunkhumpornpat, C., Sinapiromsaran, K., & Lursinsap, C. (2009). Safe-level-SMOTE: Safe-level-synthetic minority over-sampling technique for handling the class imbalanced problem. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 475-482). Springer. https://doi.org/10.1007/978-3-642-01307-2_43
4. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16, 321-357. <https://doi.org/10.1613/jair.953>
5. Douzas, G., & Bacao, F. (2019). Geometric SMOTE a geometrically enhanced drop-in replacement for SMOTE. *Information Sciences*, 501, 118-135. <https://doi.org/10.1016/j.ins.2019.06.007>
6. Elor, Y., & Averbuch-Elor, H. (2022). To SMOTE, or not to SMOTE? *arXiv preprint arXiv:2201.08528*. <https://doi.org/10.48550/arXiv.2201.08528>
7. He, H., & Garcia, E. A. (2009). Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9), 1263-1284. <https://doi.org/10.1109/TKDE.2008.239>
8. Fotouhi, S., Asadi, S., & Kattan, M. W. (2019). A comprehensive data level analysis for cancer diagnosis on imbalanced data. *Journal of Biomedical Informatics*, 90, 103089. <https://doi.org/10.1016/j.jbi.2018.12.003>
9. Gulshan, V., Peng, L., Coram, M., Stumpe, M. C., Wu, D., Narayanaswamy, A., Venugopalan, S., Widner, K., Madams, T., Cuadros, J., Kim, R., Raman, R., Nelson, P. C., Mega, J. L., & Webster, D. R. (2016). Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *JAMA*, 316(22), 2402-2410. <https://doi.org/10.1001/jama.2016.17216>
10. Haixiang, G., Yijing, L., Shang, J., Mingyun, G., Yuanyue, H., & Bing, G. (2017). Learning from class-imbalanced data: Review of methods and applications. *Expert Systems with Applications*, 73, 220-239. <https://doi.org/10.1016/j.eswa.2016.12.035>
11. Han, H., Wang, W. Y., & Mao, B. H. (2005). Borderline-SMOTE: A new over-sampling method in imbalanced data sets learning. In *International*

- Conference on Intelligent Computing (pp. 878-887). Springer. https://doi.org/10.1007/11538059_91
12. Hannun, A. Y., Rajpurkar, P., Haghpanahi, M., Tison, G. H., Bourn, C., Turakhia, M. P., & Ng, A. Y. (2019). Cardiologist-level arrhythmia detection and classification in ambulatory electrocardiograms using a deep neural network. *Nature Medicine*, 25(1), 65-69. <https://doi.org/10.1038/s41591-018-0268-3>
13. He, H., Bai, Y., Garcia, E. A., & Li, S. (2008). ADASYN: Adaptive synthetic sampling approach for imbalanced learning. In 2008 IEEE International Joint Conference on Neural Networks (pp. 1322-1328). IEEE. <https://doi.org/10.1109/IJCNN.2008.4633969>
14. Japkowicz, N., & Stephen, S. (2002). The class imbalance problem: A systematic study. *Intelligent Data Analysis*, 6(5), 429-449. <https://doi.org/10.3233/IDA-2002-6504>
15. Johnson, J. M., & Khoshgoftaar, T. M. (2019). Survey on deep learning with class imbalance. *Journal of Big Data*, 6(1), 27. <https://doi.org/10.1186/s40537-019-0192-5>
16. Kaur, H., Pannu, H. S., & Malhi, A. K. (2019). A systematic review on imbalanced data challenges in machine learning: Applications and solutions. *ACM Computing Surveys*, 52(4), 1-36. <https://doi.org/10.1145/3343440>
17. Kovács, G., Sebestyen, G., & Hangan, A. (2019). Evaluation of data augmentation techniques for medical image analysis. In 2019 IEEE International Conference on Intelligent Computer Communication and Processing (pp. 301-306). IEEE. <https://doi.org/10.1109/ICCP48234.2019.8959768>
18. Last, F., Douzas, G., & Bacao, F. (2017). Oversampling for imbalanced learning based on k-means and SMOTE. *Information Sciences*, 465, 1-20. <https://doi.org/10.1016/j.ins.2017.02.015>
19. López, V., Fernández, A., García, S., Palade, V., & Herrera, F. (2013). An insight into classification with imbalanced data: Empirical results and current trends on using data intrinsic characteristics. *Information Sciences*, 250, 113-141. <https://doi.org/10.1016/j.ins.2013.07.007>
20. Mazurowski, M. A., Habas, P. A., Zurada, J. M., Lo, J. Y., Baker, J. A., & Tourassi, G. D. (2008). Training neural network classifiers for medical decision making: The effects of imbalanced datasets on classification performance. *Neural Networks*, 21(2-3), 427-436. <https://doi.org/10.1016/j.neunet.2007.12.031>
21. Nguyen, H. M., Cooper, E. W., & Kamei, K. (2011). Borderline over-sampling for imbalanced data classification. *International Journal of Knowledge Engineering and Soft Data Paradigms*, 3(1), 4-21. <https://doi.org/10.1504/IJKESDP.2011.039875>
22. Rahman, M. M., & Davis, D. N. (2013). Addressing the class imbalance problem in medical datasets. *International Journal of Machine Learning and Computing*, 3(2), 224-228. <https://doi.org/10.7763/IJMLC.2013.V3.307>
23. Rajpurkar, P., Chen, E., Banerjee, O., & Topol, E. J. (2022). AI in health and medicine. *Nature Medicine*, 28(1), 31-38. <https://doi.org/10.1038/s41591-021-01614-0>
24. Santos, M. S., Soares, J. P., Abreu, P. H., Araujo, H., & Santos, J. (2018). Cross-validation for imbalanced datasets: Avoiding overoptimistic and overfitting approaches. *IEEE Computational Intelligence Magazine*, 13(4), 59-76. <https://doi.org/10.1109/MCI.2018.2866730>
25. Santos, M. S., Abreu, P. H., García-Laencina, P. J., Simão, A., & Carvalho, A. (2015). A new cluster-based oversampling method for improving survival prediction of hepatocellular carcinoma patients. *Journal of Biomedical Informatics*, 58, 49-59. <https://doi.org/10.1016/j.jbi.2015.09.012>
26. Topol, E. J. (2019). High-performance medicine: The convergence of human and artificial intelligence. *Nature Medicine*, 25(1), 44-56. <https://doi.org/10.1038/s41591-018-0300-7>
27. Van Hulse, J., Khoshgoftaar, T. M., & Napolitano, A. (2007). Experimental perspectives on learning from imbalanced data. In *Proceedings of the 24th International Conference on Machine Learning* (pp. 935-942). ACM. <https://doi.org/10.1145/1273496.1273614>