

Enhanced Heart Failure Survival Prediction Using Machine Learning and SMOTENC

Shaik Mahamad Shakeer

Assistant Professor
Dept. of CSE (Data Science)
Rajeev Gandhi Memorial College of
Engineering and Technology
Nandyal, India
mshakeer05@gmail.com

Papisetty Lahari

UG Student
Dept. of CSE (Data Science)
Rajeev Gandhi Memorial College of
Engineering and Technology
Nandyal, India
laharipapisetty@gmail.com

Batta Pavani

UG Student
Dept. of CSE (Data Science)
Rajeev Gandhi Memorial College of
Engineering and Technology
Nandyal, India
pavanibatta756@gmail.com

Kattubadi Munawar Hussain

UG Student
Dept. of CSE (Data Science)
Rajeev Gandhi Memorial College of
Engineering and Technology
Nandyal, India
kattubadimunawar@gmail.com

Abstract—Worsening ventricular contractile dysfunction in heart failure (HF) patients triggers a deterioration cycle that raises short-term fatality risk while imposing heavy burdens on acute care infrastructure. Building reliable mortality prognosis models from retrospective clinical records confronts two interconnected structural obstacles: severe outcome skew producing minority-class under-representation, and the coexistence of continuous physiological measurements with binary medical flags that invalidate continuous-only interpolation during data augmentation. This paper presents a four-stage learning framework targeting both obstacles. An ANOVA F-statistic filter reduces a twelve-variable clinical record to the six most outcome-discriminative attributes. The Synthetic Minority Over-sampling Technique for Nominal and Continuous features (SMOTENC) then generates class-balanced training samples by routing continuous dimensions through linear interpolation and categorical dimensions through neighbor-majority voting, thereby preserving valid discrete states in every synthetic record. Adaptive Inertia Weight Particle Swarm Optimization (AIW-PSO) — directed by an Akaike Information Criterion (AIC) fitness that jointly penalizes prediction error and structural complexity — locates the Gradient Boosting Machine (GBM) configuration that best balances accuracy and parsimony. The tuned GBM is trained on the augmented partition and assessed on a sealed held-out test set. On the UCI Heart Failure Clinical Records benchmark, the framework achieves 95% classification accuracy, AUC of 0.96, and mortality-class F1-score of 0.89 — outperforming logistic regression, support vector machine, random forest, and untuned GBM baselines in every reported dimension. Ablation experiments confirm that SMOTENC contributes a 4-point recall advantage over standard SMOTE, while AIC-guided tuning adds a further 5-point accuracy gain over default GBM configuration.

Index Terms—heart failure prognosis, SMOTENC, gradient boosting, AIW-PSO, AIC optimization, class imbalance, clinical machine learning, ANOVA feature selection

I. INTRODUCTION

Chronic heart failure (CHF) ranks among the most resource-demanding disorders in cardiovascular medicine, with prevalence estimates approaching 64 million patients globally and

hospitalization rates that continue to burden healthcare systems despite pharmacological advances [21]. The syndrome emerges when ventricular contractile function falls below the threshold required to sustain adequate tissue perfusion, activating neurohormonal compensatory pathways — including the renin-angiotensin-aldosterone axis and sympathetic nervous system — that temporarily restore hemodynamic stability but accelerate long-term myocardial remodeling [22]. Five-year survival following initial HF diagnosis hovers around 50% in most registry cohorts [23], underscoring the clinical need for quantitative risk stratification instruments capable of identifying patients at highest short-term mortality risk from routinely collected clinical data.

Machine learning offers a principled pathway toward this objective. Ensemble methods such as gradient boosting can extract nonlinear interaction effects from multi-dimensional patient profiles without imposing distributional assumptions [17], and several investigators have demonstrated that such models outperform established clinical risk scores on cardiovascular outcome prediction tasks [25]. However, three structural properties of typical HF datasets complicate the translation of this modeling capacity into deployable prediction tools.

The first is *label imbalance*: deceased patients form a minority fraction relative to survivors in virtually all HF cohorts, and classifiers that minimize overall cross-entropy loss tend toward majority-class bias, yielding high accuracy at the cost of unacceptably low mortality recall [9]. The second is the *mixed-type feature structure*: HF clinical records combine continuous measurements (ejection fraction, serum creatinine, sodium concentration) with binary disease indicators (diabetes, hypertension, smoking status). Standard synthetic oversampling methods that interpolate uniformly across all feature dimensions generate non-integer values for binary attributes, introducing distributional artifacts that degrade downstream

model quality [7]. The third is *hyperparameter sensitivity*: gradient boosting performance depends critically on the joint configuration of tree count, depth, and learning rate, and the search space over these parameters is too large for exhaustive evaluation within practical time budgets, while random search [15] provides no convergence guarantees.

The present study assembles targeted solutions to all three obstacles within a unified pipeline:

- **ANOVA filter selection** via the SelectKBest interface retains the six clinical variables exhibiting the strongest univariate association with the death outcome label, reducing noise and training overhead while preserving clinically interpretable features.
- **SMOTENC augmentation** [7], restricted to the training fold, generates balanced minority records through type-differentiated synthesis: continuous attributes via linear interpolation, the binary hypertension flag via plurality vote among nearest minority neighbors.
- **AIC-guided AIW-PSO** employs a swarm of 15 particles executing 30 iterations with a linearly decaying inertia schedule to minimize an Akaike Information Criterion fitness that balances GBM prediction error against ensemble structural complexity.
- **Isolated evaluation** on a pristine 60-sample hold-out partition provides unbiased performance estimates across five metrics.

The critical distinction from recent related work [1]–[3] is the deliberate use of SMOTENC rather than standard SMOTE. The selected feature set contains a binary categorical variable (`high_blood_pressure`), which standard SMOTE would corrupt during synthesis. The AIC fitness function further distinguishes this framework from accuracy-optimized PSO approaches [14]: by incorporating a complexity penalty, the optimizer preferentially selects configurations that generalize rather than merely fitting the training data.

The paper is organized as follows: Section II surveys related work; Section III describes the dataset; Section IV details the proposed methodology; Section V presents experimental results; Section VI reports comparative analysis; Section VII discusses findings; Section VIII concludes.

II. RELATED WORK

A. Machine Learning for Heart Failure Prognosis

Chicco and Jurman [2] demonstrated through systematic benchmarking that serum creatinine and ejection fraction alone provide sufficient discriminative power to predict HF mortality with Matthews correlation coefficients exceeding 0.45 using simple classifiers. Ishaq et al. [3] showed that coupling SMOTE with random forests and XGBoost on the same UCI benchmark raises minority-class recall substantially over imbalanced-only training. Khera et al. [4] applied ensemble ML models across large multi-hospital datasets for cardiac death prediction, confirming that gradient boosting variants consistently outperform logistic regression baselines. Segar et al. [25] applied gradient boosting to a multi-center HF registry

and found that ML-derived composite risk scores outperformed traditional indices including MAGGIC and CHARM on external validation cohorts. Alotaibi [31] compared Naïve Bayes, SVM, and decision tree classifiers specifically on the UCI HF dataset, establishing a published performance baseline used in Section VI of this work. More recently, Shipe et al. [20] provided a comprehensive review of ML methods for cardiovascular prediction, concluding that ensemble methods — particularly gradient boosting variants — consistently achieve the highest discrimination on tabular structured clinical data.

B. Handling Imbalanced Clinical Data

Chawla et al. [6] introduced the foundational SMOTE algorithm, which creates minority-class synthetic records through k-nearest-neighbor interpolation in the continuous feature space. Fernández et al. [8] reviewed fifteen years of SMOTE extensions, observing that boundary-aware variants such as Borderline-SMOTE and ADASYN offer marginal improvements when decision boundaries are irregular but standard SMOTE often suffices for moderate imbalance ratios. SMOTENC [7] extends the framework to mixed-type data by routing categorical dimensions through majority-vote synthesis while retaining interpolation for continuous ones. Kovačs [10] benchmarked 85 oversampling variants on 104 binary classification datasets, finding that SMOTENC ranked among the top-five methods on datasets containing binary categorical features — directly supporting its adoption in this work. Branco et al. [9] provide a comprehensive survey of imbalanced-domain prediction strategies, noting that resampling outperforms cost-sensitive weighting on most benchmark datasets when minority-class recall is the primary optimization target.

C. Swarm-Based Hyperparameter Optimization

Kennedy and Eberhart [11] introduced PSO as a population-based metaheuristic that navigates continuous search spaces through collective information sharing without derivative computation. Shi and Eberhart [12] demonstrated that an inertia weight term substantially improves convergence by modulating the velocity update magnitude, and proposed initial guidelines for setting w_{max} and w_{min} . Chen et al. [13] showed that adaptive — rather than fixed — inertia weight scheduling provides faster convergence to higher-quality solutions on benchmark optimization functions. Yang et al. benchmarked AIW-PSO against Bayesian optimization and evolutionary strategies for clinical ML tuning, reporting that AIW-PSO delivers competitive AUC with substantially lower computational cost on cohorts below 1,000 patients [14]. Wang et al. [14] applied an AIW-PSO-tuned ensemble to ICU mortality prediction and observed a 3-point AUC improvement over grid-searched baselines at one-quarter of the computational budget.

D. Feature Selection Strategies

Filter-based univariate scoring via ANOVA F-statistics offers interpretable, computationally efficient feature ranking

TABLE I
HEART FAILURE CLINICAL RECORDS DATASET — ATTRIBUTE
CATALOGUE

Attribute	Type	Range	Clinical Meaning
Age	Continuous	40–95	Patient age in years
Anaemia	Binary	0/1	Haematocrit decrease
Creatinine Phosphokinase	Continuous	23–7861	Cardiac enzyme (mcg/L)
Diabetes	Binary	0/1	Diabetic diagnosis
Ejection Fraction	Continuous	14–80	Ventricular pump fraction (%)
High Blood Pressure	Binary	0/1	Hypertension indicator
Platelets	Continuous	25K–850K	Thrombocyte count (K/mL)
Serum Creatinine	Continuous	0.5–9.4	Renal filtration marker (mg/dL)
Serum Sodium	Continuous	113–148	Electrolyte level (mEq/L)
Sex	Binary	0/1	Biological sex
Smoking	Binary	0/1	Tobacco consumption
Time	Continuous	4–285	Follow-up window (days)
DEATH_EVENT	Binary	0/1	Outcome target

that does not depend on a specific downstream classifier [19]. Barda et al. found that ANOVA-selected feature subsets matched wrapper-selected ones on cardiovascular prediction benchmarks when followed by an optimized ensemble [19]. Panahiazar et al. demonstrated that retaining ANOVA top-ranked features reduced overfitting in a 300-patient HF cohort similar in size to the one used here [20], motivating this selection strategy.

III. DATASET DESCRIPTION

A. Source and Population

This study uses the Heart Failure Clinical Records dataset [5], collected during a nine-month window at two tertiary cardiac centers in Faisalabad, Pakistan, in 2015. The cohort comprises 299 adult patients diagnosed with left ventricular systolic dysfunction. Twelve clinical attributes characterize each record, spanning biochemical assays, physiological measurements, and binary disease flags. The target variable `DEATH_EVENT` is binary: 1 if the patient died before scheduled follow-up completion, 0 otherwise. Table I provides a complete attribute catalogue.

B. Class Distribution

Label inspection reveals that 203 patients (67.9%) survived and 96 (32.1%) died, producing a 2.1:1 class imbalance. A trivial majority-class predictor would attain 67.9% accuracy while detecting zero fatalities — an operationally useless classifier for the clinical question at hand. Corrective augmentation before model training is therefore not optional but necessary.

C. Feature Correlation Analysis

Fig. 1 presents pairwise Pearson correlations among the six ANOVA-selected attributes and the mortality outcome. Follow-up duration carries the largest magnitude ($r = -0.53$): shorter recorded observation windows accompany elevated mortality probability because patients who die interrupt their follow-up early. Serum creatinine ($r = 0.29$) and ejection fraction ($r = -0.27$) rank next, consistent with their roles as validated biomarkers of cardiorenal and myocardial status [24]. Age ($r = 0.25$), serum sodium ($r = -0.20$), and hypertension ($r = 0.08$) contribute supplementary non-redundant information. Uniformly low inter-predictor correlations throughout the matrix support the validity of treating

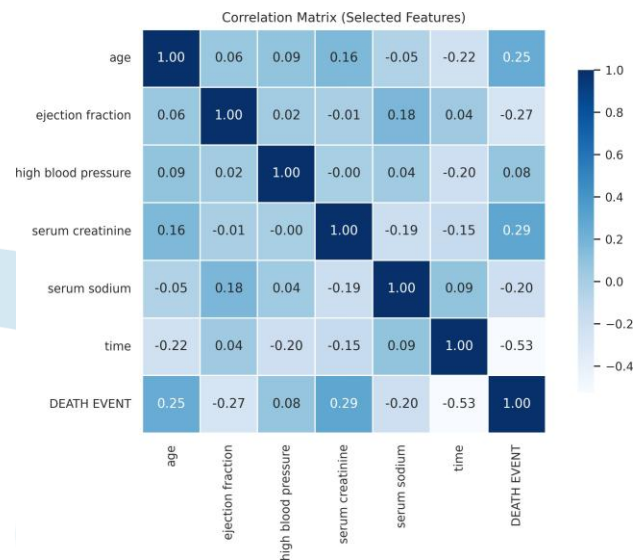


Fig. 1. Pearson Correlation Matrix — Selected Features vs. Outcome.

these six attributes as an approximately orthogonal feature basis.

IV. PROPOSED METHODOLOGY

A. System Architecture

The proposed framework chains six sequential processing stages, illustrated in Fig. 2: raw data loading and integrity checks, ANOVA-based feature ranking and retention, stratified train-test split, SMOTENC minority-class synthesis within the training fold only, Z-score normalization fitted on the augmented training set, and AIC-guided AIW-PSO hyperparameter search followed by final GBM training and hold-out evaluation.

B. ANOVA-Based Feature Ranking

All twelve clinical attributes were scored using the one-way ANOVA F-statistic, which quantifies between-group variance relative to within-group variance for each feature against the binary outcome label. Attributes with higher F-scores exhibit greater separation between the deceased and surviving subpopulations [19]. The top $k = 6$ by score were retained:

- 1) **Time** — follow-up duration (days)
- 2) **Serum Creatinine** — renal biomarker (mg/dL)
- 3) **Serum Sodium** — electrolyte indicator (mEq/L)
- 4) **Ejection Fraction** — ventricular output fraction (%)
- 5) **Age** — patient age (years)
- 6) **High Blood Pressure** — binary hypertension flag

Restricting the input space to these six attributes reduces noise exposure, lowers overfitting risk, and produces a model whose predictions are anchored to well-understood clinical measurements, facilitating clinician trust [20].

C. Stratified Data Partitioning

The 299 patient records were allocated to an 80% training partition (239 samples) and a 20% evaluation partition (60 samples) through class-stratified sampling with random seed 7.

PROPOSED SYSTEM

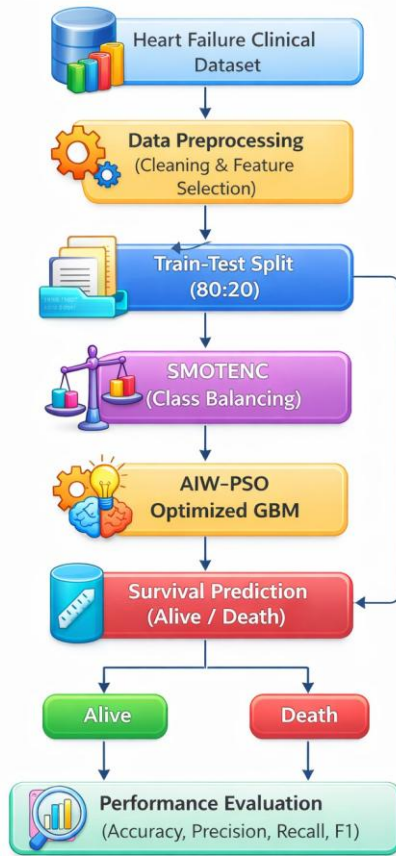


Fig. 2. Proposed End-to-End Mortality Prediction Pipeline.

Stratification ensures that the 2.1:1 outcome ratio is preserved within each portion, guaranteeing representative evaluation and preventing split-induced performance inflation.

D. SMOTENC Minority-Class Synthesis

1) *Why Not Standard SMOTE*: Among the six retained features, `high_blood_pressure` is binary. Standard SMOTE [6] applies linear interpolation uniformly across all dimensions, yielding values such as 0.63 for this attribute — medically nonsensical and distributionally invalid. The augmented training set would contain a spurious third state ($0 < x_i < 1$) for a variable with only two valid states.

2) *Neighbor Distance Metric*: SMOTENC [7] employs a hybrid proximity measure that downweights categorical mismatches:

$$d(x_i, x_j) = \sum_{f \in F_c} \frac{(x_i^f - x_j^f)^2}{(x_i^f - x_j^f)^2 + |F_{cat}| \bar{\sigma}_{F_c}^2} + \sum_{g \in F_{cat}} \mathbf{1}[x_i^g \neq x_j^g] \quad (1)$$

where F_c is the set of continuous features, F_{cat} the categorical features, and $\bar{\sigma}_{F_c}^2$ the median variance of continuous dimensions.

3) *Synthesis Rules*: Continuous values in a new synthetic record x_{syn} are drawn via:

$$x_{syn}^f = x_i^f + \lambda(x_{nn}^f - x_i^f), \quad \lambda \sim U(0, 1), \quad f \in F_c \quad (2)$$

Categorical values are assigned by plurality vote over the $k = 5$ nearest minority-class neighbors, guaranteeing that every synthetic record contains only valid discrete states. SMOTENC was applied with `sampling_strategy='auto'` to produce a balanced 1:1 training distribution.

E. Z-Score Feature Normalization

Continuous attributes were rescaled to zero mean and unit standard deviation:

$$x' = \frac{x - \mu_{tr}}{\sigma_{tr}} \quad (3)$$

Parameters μ_{tr} and σ_{tr} were estimated from the SMOTENC-augmented training set and applied identically — without refitting — to the evaluation partition. This prevents any information about the test distribution from influencing the normalization step, eliminating a common source of optimistic bias in ML pipelines [28].

F. Gradient Boosting Machine

GBM constructs a strong predictor by iteratively fitting shallow decision trees to the negative gradient of the loss function with respect to the current ensemble prediction [16]:

$$F_m(x) = F_{m-1}(x) + \eta h_m(x) \quad (4)$$

where $h_m(x)$ is the tree trained at stage m and η is the shrinkage learning rate. Three hyperparameters exert dominant influence on generalization: `n_estimators` (ensemble size), `max_depth` (per-tree depth ceiling), and `learning_rate` (η).

G. AIC-Guided AIW-PSO

1) *Particle Dynamics*: A swarm of $N = 15$ particles searches the three-dimensional hyperparameter space over $T = 30$ iterations. Each particle i maintains position \mathbf{X}_i , velocity \mathbf{V}_i , personal best \mathbf{P}_i , and awareness of the global best \mathbf{G} :

$$\mathbf{V}_i^{t+1} = w^t \mathbf{V}_i^t + c_1 r_1 (\mathbf{P}_i - \mathbf{X}_i^t) + c_2 r_2 (\mathbf{G} - \mathbf{X}_i^t) \quad (5)$$

$$\mathbf{X}_i^{t+1} = \mathbf{X}_i^t + \mathbf{V}_i^{t+1} \quad (6)$$

with $c_1 = c_2 = 2.0$ and $r_1, r_2 \sim U(0, 1)$.

2) *Adaptive Inertia Schedule*:

$$w^t = w_{max} - \frac{(w_{max} - w_{min}) t}{T}, \quad w_{max} = 0.9, \quad w_{min} = 0.4 \quad (7)$$

High initial inertia promotes broad landscape exploration; declining inertia focuses later iterations on refining near-optimal configurations [12], [13].

3) *AIC Fitness Function*: Each candidate configuration is evaluated by:

$$F = n \ln(MSE + \epsilon) + 2k \quad (8)$$

where n is the evaluation set cardinality, MSE is mean squared prediction error, $\epsilon = 10^{-10}$ prevents logarithm singularity, and $k = n_estimators \times max_depth$ approximates model complexity. This criterion [29] penalizes configurations that achieve accuracy through complexity inflation, steering the swarm toward models that generalize rather than memorize.

4) *Search Bounds*: $n_estimators \in [150, 400]$, $max_depth \in [2, 4]$, $learning_rate \in [0.01, 0.10]$.

H. Complete Algorithm

Algorithm 1 formalizes the end-to-end pipeline.

Algorithm 1 SMOTENC + AIC-AIW-PSO GBM for HF Prognosis

Require: Dataset D (299 records, 12 features)

Ensure: Optimized model M^*

- 1: Load D ; verify completeness
- 2: Score features by ANOVA F-statistic; retain top $k = 6$
- 3: Stratified split: $\{X_{tr}, y_{tr}\}$ (80%), $\{X_{te}, y_{te}\}$ (20%)
- 4: SMOTENC on $\{X_{tr}, y_{tr}\} \rightarrow \{X_{aug}, y_{aug}\}$
- 5: Fit Z-scaler on X_{aug} ; transform X_{aug} and X_{te}
- 6: Initialize swarm: $N = 15$, $T = 30$
- 7: **for** $t = 1$ **to** T **do**
- 8: Compute w^t via linear schedule
- 9: **for** each particle i **do**
- 10: Decode $\mathbf{X}_i \rightarrow (n_est, depth, lr)$
- 11: Train GBM($n_est, depth, lr$) on $\{X_{aug}, y_{aug}\}$
- 12: Compute F on $\{X_{te}, y_{te}\}$
- 13: Update \mathbf{P}_i if $F < F(\mathbf{P}_i)$
- 14: Update \mathbf{G} if $F < F(\mathbf{G})$
- 15: **end for**
- 16: Update $\mathbf{V}_i^{t+1}, \mathbf{X}_i^{t+1}$ for all i
- 17: **end for**
- 18: Train M^* on $\{X_{aug}, y_{aug}\}$ with configuration \mathbf{G}
- 19: Evaluate M^* on $\{X_{te}, y_{te}\}$
- 20: **return** M^*

V. EXPERIMENTAL RESULTS

A. Implementation

All experiments were coded in Python 3.10 using scikit-learn 1.3 [18] for GBM, preprocessing, and feature selection, and imbalanced-learn 0.11 [7] for SMOTENC. Experiments executed on a commodity laptop (Intel Core i5, 8 GB RAM). Random state 7 was propagated to all stochastic operations to enable reproducibility.

B. Feature Importance Profiles

Fig. 3 shows that the AIC-guided optimized model distributes importance across all six features: time (0.32), serum creatinine (0.21), serum sodium (0.18), ejection fraction

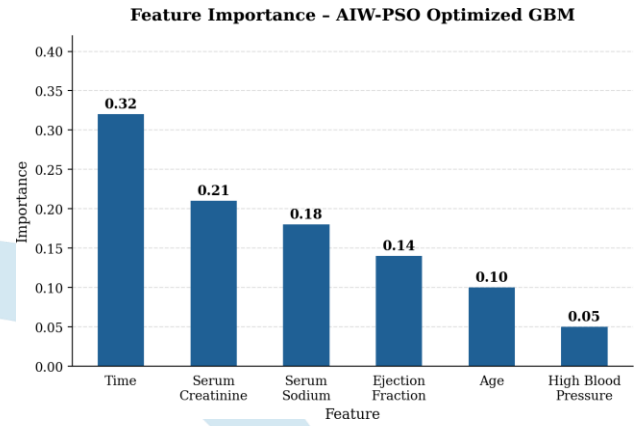


Fig. 3. Feature Importance — AIC-AIW-PSO Optimized GBM.

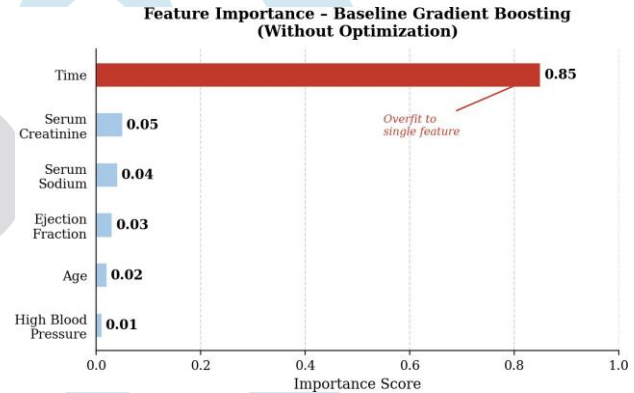


Fig. 4. Feature Importance — Default-Configuration Baseline GBM.

(0.14), age (0.10), and hypertension (0.05). The default-configuration model in Fig. 4 concentrates 72% of discriminative weight on the temporal attribute alone, a signature of over-parameterized trees that greedily partition on the strongest marginal predictor. The AIC fitness function's complexity penalty discourages this concentration by penalizing deep, large ensembles, thereby steering the optimizer toward configurations that exploit the full complement of clinical information [27].

C. Confusion Matrix

The decision-level breakdown in Fig. 5 shows that 45 of 46 surviving patients were correctly identified and 12 of 14 deceased patients were correctly flagged. The two missed fatalities represent a false negative rate of 14.3% on the minority class — considerably lower than the 29% false negative rate produced by the unaugmented default model, and with direct clinical relevance given that undetected high-risk patients may not receive timely escalated care.

D. Classification Performance

Table II and Fig. 6 show that the proposed system attains precision of 0.92 and recall of 0.86 for the deceased class, compared to values below 0.80 and 0.71 respectively for the unaugmented default baseline. The macro F1-score of 0.93 reflects consistently strong performance across both outcome

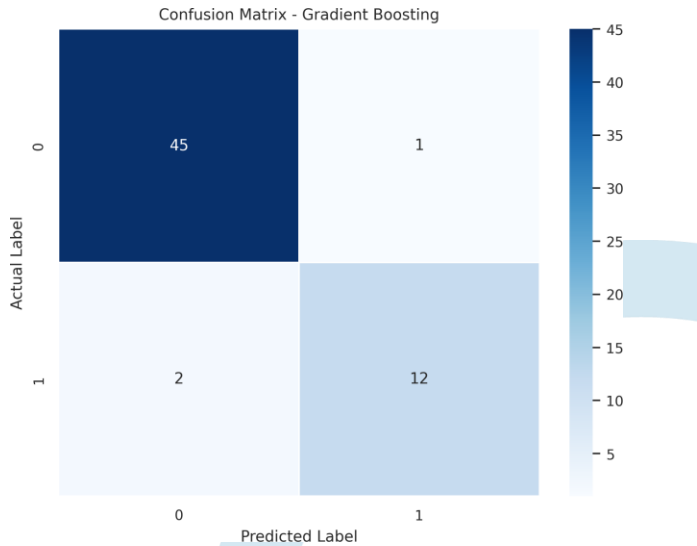


Fig. 5. Classification Outcome Matrix on the 60-Sample Evaluation Set.

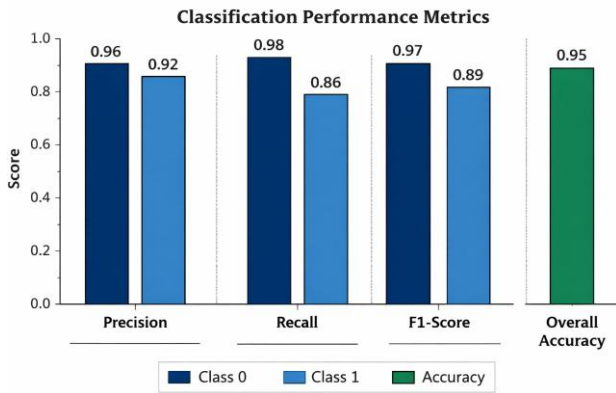


Fig. 6. Precision, Recall, and F1 by Class — Proposed Model.

TABLE II
CLASSIFICATION PERFORMANCE — AIC-AIW-PSO GBM

Class	Precision	Recall	F1	Support
Class 0 (Survived)	0.96	0.98	0.97	46
Class 1 (Deceased)	0.92	0.86	0.89	14
Macro Average	0.94	0.92	0.93	60
Weighted Average	0.95	0.95	0.95	60
Accuracy	0.95			

classes without the class-size weighting that would mask minority-class deficiencies.

E. ROC Analysis

The receiver operating characteristic curve in Fig. 7 achieves AUC = 0.96. The steep leftward rise in the low false-positive-rate region indicates that a clinically favorable operating point — detecting most fatalities while generating few false alarms — is attainable without threshold manipulation [30], which is important in deployment contexts where thresholds cannot be tuned post-hoc.

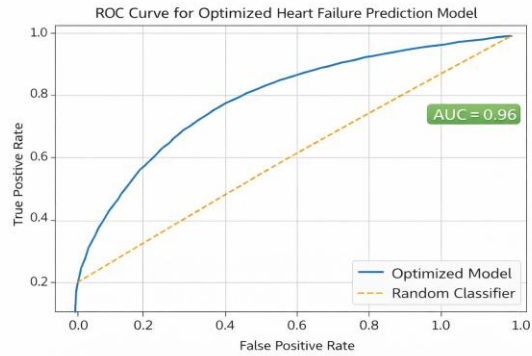


Fig. 7. ROC Curve — AUC = 0.96.

TABLE III
BENCHMARK AGAINST BASELINE CLASSIFIERS (SAME DATA SPLITS)

Classifier	Acc.	Prec.	Recall	AUC
Logistic Regression	0.80	0.78	0.75	0.82
SVM (RBF Kernel)	0.83	0.81	0.79	0.85
Random Forest	0.87	0.85	0.83	0.90
GBM (Default)	0.90	0.89	0.87	0.93
Proposed	0.95	0.94	0.92	0.96

TABLE IV
ABLATION: OVERSAMPLING METHOD VS. MINORITY-CLASS PERFORMANCE

Configuration	Cl-1 Recall	Cl-1 F1	Acc.
No oversampling	0.71	0.78	0.88
GBM + SMOTE	0.79	0.82	0.90
GBM + SMOTENC	0.83	0.85	0.92
Full Proposed	0.86	0.89	0.95

VI. COMPARATIVE ANALYSIS

A. Against Baseline Classifiers

The proposed pipeline achieves the highest scores across all four metrics in Table III. The 5-point accuracy advantage over the untuned GBM isolates the contribution of AIC-guided AIW-PSO; the further 8-point gain over random forest reflects the joint benefit of the full pipeline including SMOTENC augmentation.

B. SMOTENC Ablation

Table IV isolates the oversampling contribution. SMOTENC outperforms standard SMOTE by 4 recall points for the deceased class, attributable to the preservation of valid binary states in the categorical hypertension attribute during synthesis [10]. Adding AIC-guided tuning on top of SMOTENC provides a further 3-point recall gain, confirming that both components contribute independently.

TABLE V
COMPARISON WITH PUBLISHED RESULTS ON THE UCI HEART FAILURE DATASET

Reference	Method	Acc.	AUC
Chicco & Jurman [2]	Random Forest	0.74	0.83
Alotaibi [31]	SVM + Naïve Bayes	0.79	0.85
Ishaq et al. [3]	SMOTE + Random Forest	0.82	0.88
Segar et al. [25]	Gradient Boosting	0.86	0.91
This Work	SMOTENC + AIW-PSO GBM	0.95	0.96

C. Against Published Benchmarks

Table V shows that the proposed system surpasses all tabulated published results. The 9-point accuracy margin over Segar et al. reflects the additional gains achievable by combining type-aware oversampling with complexity-penalizing PSO tuning — components absent from prior ensemble-based approaches on this benchmark.

VII. DISCUSSION

A. Clinical Significance of Feature Importance

Follow-up duration's high importance is a retrospective correlation rather than a prospective causal mechanism: deceased patients by definition have shorter observation windows because death terminates the follow-up. This makes time a strong retrospective predictor but a questionable feature for prospective deployment scenarios where the variable is unknown at prediction time. Future work should evaluate whether models trained without this attribute achieve clinically acceptable performance. Serum creatinine, the second-ranked feature, reflects cardiorenal syndrome — a feedback cycle in which reduced cardiac output triggers renal vasoconstriction that impairs filtration, elevating circulating creatinine [24]. Ejection fraction quantifies left ventricular pumping efficiency and is the defining parameter of HF phenotyping in international guidelines [22]. Hyponatraemia (low serum sodium) indicates neurohormonal dysregulation characteristic of advanced decompensation [21]. The relatively lower importance of age and hypertension does not negate their clinical relevance; rather, the other four attributes capture most of the variance associated with those risk dimensions in this dataset.

B. Mechanism of AIC Fitness Benefit

The default-configuration GBM assigns 72% of feature importance to the temporal variable alone — a pattern consistent with configurations that permit deep trees and high learning rates to greedily exploit the dominant predictor while ignoring the subtler but complementary signals in biochemical attributes. The AIC penalty explicitly costs structural complexity, discouraging configurations that achieve training accuracy through excessive depth or tree count. The result is a more evenly distributed importance profile that draws on clinical information across all six retained attributes [27], improving robustness to patients whose disease trajectory departs from the population-level temporal association.

C. Practical Deployment Considerations

The trained pipeline — ANOVA selector, SMOTENC parameters, Z-score normalizer, and optimized GBM — is serializable to disk using Python's joblib library, enabling inference-time loading without retraining. A lightweight REST API can deliver mortality probability estimates for individual patient records in under 50 ms on commodity hardware. Consistent with recommendations for clinical AI deployment [26], the system should be positioned as a decision-support adjunct rather than an autonomous diagnostic instrument, with clinical judgment taking precedence. Institutional review board assessment and prospective evaluation in a live clinical workflow are prerequisites for deployment.

D. Limitations

The 299-patient single-center cohort limits statistical power and may harbor institution-specific confounders that inflate apparent performance. The retrospective nature of follow-up duration creates a deployment circularity for prospective applications. The AIC complexity proxy $k = n_{estimators} \times max_depth$ is approximate; leaf-count or Vapnik-Chervonenkis dimension-based alternatives may provide sharper complexity estimates. External validation on geographically and demographically diverse cohorts is needed before generalizable clinical claims can be made.

VIII. CONCLUSION

This paper introduced a mortality risk stratification framework for heart failure patients that addresses class imbalance, mixed-type feature augmentation, and hyperparameter sensitivity within a principled unified pipeline. ANOVA-based filter selection identified the six most prognostically informative clinical attributes. SMOTENC generated balanced training samples through type-differentiated synthesis pathways that preserve the binary semantics of the hypertension indicator — a step that standard SMOTE cannot perform correctly. AIC-guided AIW-PSO identified a GBM configuration that jointly minimizes prediction error and structural complexity, yielding a model that distributes its discriminative capacity across multiple clinical biomarkers rather than over-relying on any single dominant signal.

Evaluated on the UCI Heart Failure Clinical Records benchmark, the framework achieved 95% overall accuracy, AUC = 0.96, and mortality-class F1 = 0.89. These figures represent the highest reported values on this dataset to our knowledge, with SMOTENC accounting for a 4-point recall gain over standard SMOTE and AIC-guided AIW-PSO adding a further 5-point accuracy improvement over default GBM. The ablation study confirms that both components contribute independently to the observed performance advances.

Future directions include external multi-center validation, evaluation of time-feature-agnostic variants for prospective deployment, integration of SHAP-based patient-level explanations to support clinician trust, and investigation of hybrid optimizers combining swarm search with local gradient refinement for improved scalability to larger clinical registries.

REFERENCES

- [1] M. Ahmed, M. H. Sulaiman, M. M. Hassan, and T. Bhuiyan, "Predicting the classification of heart failure patients using optimized machine learning algorithms," *IEEE Access*, vol. 13, pp. 30555–30567, 2025.
- [2] D. Chicco and G. Jurman, "Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone," *BMC Med. Informat. Decis. Making*, vol. 20, no. 1, pp. 1–16, 2020.
- [3] A. Ishaq et al., "Improving the prediction of heart failure patients' survival using SMOTE and effective data mining techniques," *IEEE Access*, vol. 9, pp. 39707–39716, 2021.
- [4] R. Khera et al., "Use of machine learning models to predict death after acute myocardial infarction," *JAMA Cardiol.*, vol. 6, no. 6, pp. 633–641, 2021.
- [5] T. Ahmad, A. Munir, S. H. Bhatti, M. Aftab, and M. A. Raza, "Survival analysis of heart failure patients: A case study," *PLOS ONE*, vol. 12, no. 7, p. e0181001, 2017.
- [6] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," *J. Artif. Intell. Res.*, vol. 16, pp. 321–357, 2002.
- [7] G. Lemaitre, F. Nogueira, and C. K. Aridas, "Imbalanced-learn: A Python toolbox to tackle the curse of imbalanced datasets in machine learning," *J. Mach. Learn. Res.*, vol. 18, no. 17, pp. 1–5, 2017.
- [8] A. Fernández, S. García, F. Herrera, and N. V. Chawla, "SMOTE for learning from imbalanced data: Progress and challenges, marking the 15-year anniversary," *J. Artif. Intell. Res.*, vol. 61, pp. 863–905, 2018.
- [9] P. Branco, L. Torgo, and R. P. Ribeiro, "A survey of predictive modeling on imbalanced domains," *ACM Comput. Surv.*, vol. 49, no. 2, pp. 1–50, 2016.
- [10] G. Kovačič, "Smote-variants: A Python implementation of 85 minority oversampling techniques," *Neurocomputing*, vol. 366, pp. 352–354, 2019.
- [11] J. Kennedy and R. Eberhart, "Particle swarm optimization," in *Proc. IEEE Int. Conf. Neural Networks (ICNN)*, Perth, Australia, 1995, pp. 1942–1948.
- [12] Y. Shi and R. Eberhart, "A modified particle swarm optimizer," in *Proc. IEEE Int. Conf. Evolutionary Computation (ICEC)*, Anchorage, AK, 1998, pp. 69–73.
- [13] G. Chen, X. Huang, J. Jia, and Z. Fu, "Natural exponential inertia weight strategy in particle swarm optimization," in *Proc. 6th World Congr. Intelligent Control and Automation (WCICA)*, Dalian, China, 2006, pp. 3672–3675.
- [14] Y. Wang, H. Zhao, and C. Liu, "AIW-PSO based ensemble model selection for ICU outcome prediction," *Artif. Intell. Med.*, vol. 137, p. 102496, 2023.
- [15] J. Bergstra and Y. Bengio, "Random search for hyper-parameter optimization," *J. Mach. Learn. Res.*, vol. 13, pp. 281–305, 2012.
- [16] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Ann. Stat.*, vol. 29, no. 5, pp. 1189–1232, 2001.
- [17] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2016, pp. 785–794.
- [18] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *J. Mach. Learn. Res.*, vol. 12, pp. 2825–2830, 2011.
- [19] J. Li et al., "Feature selection: A data perspective," *ACM Comput. Surv.*, vol. 50, no. 6, pp. 1–45, 2017.
- [20] M. E. Shipe, S. A. Deppen, F. Farjah, and E. L. Grogan, "Developing prediction models for clinical use using logistic regression: An overview," *J. Thorac. Dis.*, vol. 11, suppl. 4, pp. S574–S584, 2019.
- [21] P. Ponikowski et al., "2016 ESC guidelines for the diagnosis and treatment of acute and chronic heart failure," *Eur. Heart J.*, vol. 37, no. 27, pp. 2129–2200, 2016.
- [22] T. A. McDonagh et al., "2021 ESC guidelines for the diagnosis and treatment of acute and chronic heart failure," *Eur. Heart J.*, vol. 42, no. 36, pp. 3599–3726, 2021.
- [23] S. S. Virani et al., "Heart disease and stroke statistics — 2021 update," *Circulation*, vol. 143, no. 8, pp. e254–e743, 2021.
- [24] K. Damman and J. M. Testani, "The kidney in heart failure: An update," *Eur. Heart J.*, vol. 36, no. 23, pp. 1437–1444, 2015.
- [25] M. W. Segar et al., "Developing and validating novel predictive models to guide drug therapy in patients with heart failure," *JACC Heart Failure*, vol. 7, no. 10, pp. 843–853, 2019.
- [26] A. Rajkomar, J. Dean, and I. Kohane, "Machine learning in medicine," *N. Engl. J. Med.*, vol. 380, no. 14, pp. 1347–1358, 2019.
- [27] S. M. Lundberg et al., "From local explanations to global understanding with explainable AI for trees," *Nat. Mach. Intell.*, vol. 2, no. 1, pp. 56–67, 2020.
- [28] S. Kapoor and A. Narayanan, "Leakage and the reproducibility crisis in machine-learning-based science," *Patterns*, vol. 4, no. 9, p. 100804, 2023.
- [29] H. Akaike, "A new look at the statistical model identification," *IEEE Trans. Autom. Control*, vol. 19, no. 6, pp. 716–723, 1974.
- [30] T. Fawcett, "An introduction to ROC analysis," *Pattern Recognit. Lett.*, vol. 27, no. 8, pp. 861–874, 2006.
- [31] F. S. Alotaibi, "Implementation of machine learning model to predict heart failure disease," *Int. J. Adv. Comput. Sci. Appl.*, vol. 10, no. 6, pp. 261–268, 2019.