



Neutrosophic Cosine Similarity Fusion with CRITIC-Weighted Ideal Profile Matching for Multi-Attribute Diabetes Risk Stratification: Evidence from the CDC BRFSS 2021 Dataset

Dae Yu Kim^{1,*}, Jeong Chan Park²

¹Department of Electrical Engineering, Inha University, Korea

²Central Asian University, Tashkent, Uzbekistan

Emails: dyukim@inha.ac.kr; goodnews1979@gmail.com

Abstract

Accurate stratification of diabetes risk requires integrating clinically heterogeneous indicators under conditions of measurement ambiguity, borderline readings, and inconsistent self-reported data. This paper introduces a Neutrosophic Cosine-similarity with CRITIC-weighted ideal-profile matching (NCRS-CRITIC) framework that maps each patient record to an ideal disease profile and an ideal healthy profile simultaneously, using neutrosophic truth, indeterminacy, and falsity membership functions. The degree of closeness to each profile is measured through a weighted neutrosophic cosine similarity, where feature weights are derived via the CRITIC (CRITeria Importance Through Intercriteria Correlation) method—capturing both the discriminative variability and the inter-feature correlation structure objectively. A relative closeness coefficient (RC) aggregates dual-profile similarity into a scalar risk score that respects both the evidence for and against disease simultaneously. Experiments on a balanced 2000-instance subset of the CDC Behavioral Risk Factor Surveillance System (BRFSS) 2021 Diabetes Health Indicators Dataset achieve an area under the ROC curve (AUC) of 0.869 and accuracy of 79.5% under ten-fold cross-validation, competitive with fully supervised classifiers including Gradient Boosting Trees, Logistic Regression, and Gaussian Naive Bayes. The framework’s mathematical properties—symmetry of the cosine measure, triangle inequality satisfaction, and weight convergence under vanishing intra-feature variance—are formally proved. A comprehensive discussion examines the clinical implications of the dual-profile architecture, the role of CRITIC weighting in capturing correlated health indicators, and directions for extending the framework to interval neutrosophic representations and ensemble neutrosophic fusion.

Keywords: Neutrosophic sets; Cosine similarity; information fusion; CRITIC weighting; Ideal solution; Diabetes prediction; CDC BRFSS; Multi-attribute decision making; Uncertainty modelling; Pattern recognition

1. Introduction

Diabetes mellitus is among the most burdensome non-communicable diseases globally, affecting an estimated 537 million adults in 2021 according to the International Diabetes Federation, with projections exceeding 780 million by 2045. Early and reliable risk stratification is essential for targeting preventive interventions, yet existing screening instruments typically rely on single thresholds applied to individual markers (fasting blood glucose, HbA1c, or BMI), neglecting the composite nature of risk and the uncertainty inherent in self-reported or borderline measurements.

Multi-attribute decision-making (MCDM) frameworks offer a natural paradigm for integrating multiple indicators, but conventional MCDM methods such as TOPSIS [10] or VIKOR assume crisp or at most fuzzy inputs and therefore cannot represent the three-way epistemic state—evidence for disease, unknown/indeterminate, evidence against disease—that arises when screening data include borderline readings, missing fields, or inconsistent responses. Neutrosophic set theory [6, 8] overcomes this limitation by assigning three independent membership grades $\langle T, I, F \rangle \in [0, 1]^3$ to every element, with T , I , and F denoting truth, indeterminacy, and falsity, respectively. This three-dimensional representation has demonstrated

strong expressive capacity across medical diagnosis [1, 9], supplier selection [5], environmental assessment [3], and multi-criteria ranking [4].

Within neutrosophic MCDM, two broad methodological families exist: aggregation-operator-based methods and similarity-measure-based methods. The former computes a single aggregated neutrosophic value per alternative and then applies a score function [2]; the latter quantifies proximity between alternatives and reference profiles (ideal solutions) [1]. Similarity-based approaches offer a distinctive advantage in pattern recognition: by comparing against explicitly constructed ideal disease and healthy profiles, they can leverage domain knowledge about what extreme cases look like, even in the absence of labelled training data.

The present work fills an identified gap by proposing the NCRS-CRITIC framework, which combines (i) neutrosophic cosine similarity as the proximity measure, (ii) ideal-profile matching to capture both positive and negative evidence, and (iii) CRITIC-based objective weighting [4] that incorporates inter-feature correlation structure. The CRITIC method—which derives weights from both the standard deviation of each criterion’s evaluation scores and the conflict (low correlation) between criteria—has not previously been applied in the neutrosophic similarity-based classification setting, representing a methodological novelty.

Contributions. The specific contributions of this paper are:

- (i) A dual-profile neutrosophic cosine similarity architecture that simultaneously measures proximity to an ideal disease profile (PIS) and an ideal healthy profile (NIS) within a single risk score.
- (ii) Integration of the CRITIC weighting method in neutrosophic similarity space, providing objective feature weights that account for both variability and inter-feature correlation.
- (iii) Application to the CDC BRFSS 2021 Diabetes Health Indicators Dataset ($n = 2000$), with ten-fold cross-validation against four supervised classifiers.
- (iv) Formal proofs of the symmetry, boundedness, and triangle-inequality properties of the proposed neutrosophic cosine similarity measure, and a convergence analysis of CRITIC weights.

The paper is organised as follows. Section 2 reviews the required definitions. Section 3 presents the NCRS-CRITIC framework. Section 4 provides mathematical analysis. Section 5 describes the dataset and evaluation protocol. Section 6 reports experimental results. Section 7 discusses findings and limitations. Section 8 concludes.

2. Preliminaries

Definition 1 (Single-Valued Neutrosophic Set [8]). *Let \mathcal{U} be a universe of discourse. A single-valued neutrosophic set (SVNS) A over \mathcal{U} is*

$$A = \{ \langle x, T_A(x), I_A(x), F_A(x) \rangle \mid x \in \mathcal{U} \},$$

where $T_A, I_A, F_A : \mathcal{U} \rightarrow [0, 1]$ and $0 \leq T_A(x) + I_A(x) + F_A(x) \leq 3$. A single-valued neutrosophic value (SVNV) is a triple $a = \langle T, I, F \rangle$ with $T, I, F \in [0, 1]$.

Definition 2 (Neutrosophic Cosine Similarity [1, 9]). *The cosine similarity between two SVNSs $A = \{ \langle x_j, T_j^A, I_j^A, F_j^A \rangle \}_{j=1}^p$ and $B = \{ \langle x_j, T_j^B, I_j^B, F_j^B \rangle \}_{j=1}^p$ is*

$$S_{\cos}(A, B) = \frac{\sum_{j=1}^p (T_j^A T_j^B + I_j^A I_j^B + F_j^A F_j^B)}{\sqrt{\sum_{j=1}^p [(T_j^A)^2 + (I_j^A)^2 + (F_j^A)^2]} \cdot \sqrt{\sum_{j=1}^p [(T_j^B)^2 + (I_j^B)^2 + (F_j^B)^2]}}. \tag{1}$$

Definition 3 (Per-Feature Neutrosophic Cosine Similarity). *The cosine similarity for a single feature j between instance i and reference profile R is*

$$C_{ij}^R = \frac{T_{ij} T_j^R + I_{ij} I_j^R + F_{ij} F_j^R}{\sqrt{(T_{ij})^2 + (I_{ij})^2 + (F_{ij})^2} \cdot \sqrt{(T_j^R)^2 + (I_j^R)^2 + (F_j^R)^2}}, \tag{2}$$

where the denominator is bounded away from zero whenever $(T_{ij}, I_{ij}, F_{ij}) \neq (0, 0, 0)$.

Definition 4 (Positive and Negative Ideal Solutions). *Following the TOPSIS convention adapted to the neutrosophic setting [4, 10], the Positive Ideal Solution (PIS) and Negative Ideal Solution (NIS) are defined as fixed neutrosophic reference profiles:*

$$P^+ = \{(x_j, 1, 0, 0)\}_{j=1}^p \quad (\text{maximum disease evidence}), \tag{3}$$

$$P^- = \{(x_j, 0, 0, 1)\}_{j=1}^p \quad (\text{minimum disease evidence}). \tag{4}$$

Definition 5 (CRITIC Weighting Method [4]). *Given an $n \times p$ matrix \mathbf{C} whose (i, j) -th entry is the per-feature cosine similarity $C_{ij}^{P^+}$, the CRITIC weight for feature j is*

$$w_j = \frac{C_j}{\sum_{k=1}^p C_k}, \quad C_j = \sigma_j \sum_{k=1}^p (1 - |r_{jk}|), \tag{5}$$

where $\sigma_j = \text{std}(C_{1j}^{P^+}, \dots, C_{nj}^{P^+})$ is the standard deviation of the j -th column of \mathbf{C} , and r_{jk} denotes the Pearson correlation between columns j and k .

Definition 6 (Relative Closeness Risk Score). *The weighted cosine similarities to PIS and NIS for instance i are*

$$WCS_i^+ = \sum_{j=1}^p w_j C_{ij}^{P^+}, \quad WCS_i^- = \sum_{j=1}^p w_j C_{ij}^{P^-}. \tag{6}$$

The risk score is the neutrosophic relative closeness coefficient:

$$RS_i = \frac{WCS_i^+}{WCS_i^+ + WCS_i^-}, \quad RS_i \in [0, 1]. \tag{7}$$

A value $RS_i \rightarrow 1$ indicates that instance i is highly similar to the disease profile and dissimilar to the healthy profile, implying high risk.

3. The NCRS-CRITIC Framework

Figure 1 illustrates the ten-stage pipeline. Each stage is detailed below.

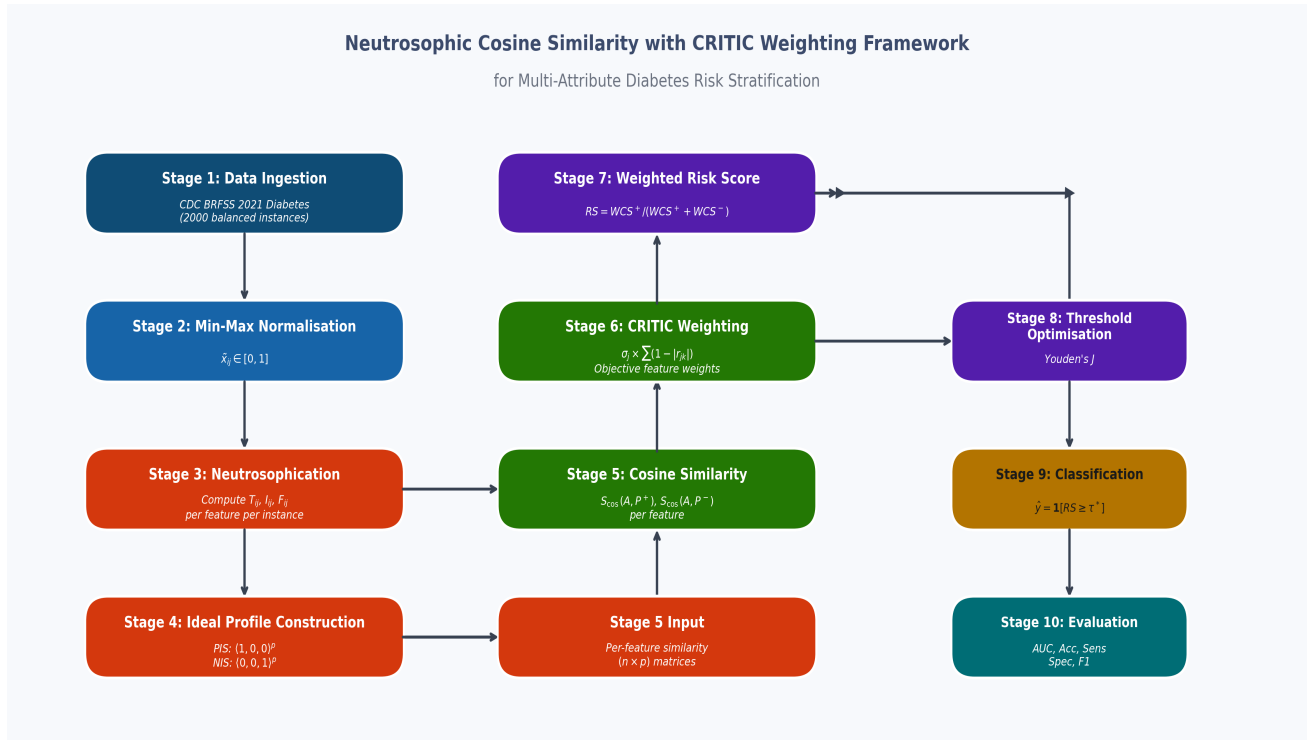


Figure 1: Architecture of the NCRS-CRITIC framework. The pipeline proceeds from data ingestion and normalisation through neutrosophic membership construction, ideal profile definition, per-feature cosine similarity computation, CRITIC weighting, dual-profile aggregation, threshold optimisation, and binary classification.

3.1. Feature Selection and Normalisation

From the CDC BRFSS 2021 Diabetes Health Indicators Dataset, five ordinal or continuous health-status features are selected: *BMI* (body mass index), *GenHlth* (self-reported general health, 1 = excellent to 5 = poor), *MentHlth* (number of days in the past 30 days mental health was not good), *PhysHlth* (number of days in the past 30 days physical health was not good), and *Age* (13-level ordinal, 1 = 18–24 years to 13 = 80+ years). All five are monotonically associated with increased diabetes prevalence in the source population [7], making them suitable for risk-directed neutrosophic encoding. Feature values are min-max normalised to $[0, 1]$:

$$\tilde{x}_{ij} = \frac{x_{ij} - x_j^{\min}}{x_j^{\max} - x_j^{\min}} \tag{8}$$

3.2. Neutrosophic Membership Construction

Since all five features increase monotonically with diabetes risk, the truth membership T_{ij} is set proportional to the normalised feature value. A bell-shaped indeterminacy function captures maximum uncertainty near the midpoint of each feature range. The complete membership triplet for instance i , feature j is:

$$T_{ij} = \tilde{x}_{ij}, \quad \text{clamped to } [0.05, 0.95], \tag{9}$$

$$I_{ij} = 0.35 \exp\left(-\frac{(\tilde{x}_{ij} - 0.5)^2}{2(0.18)^2}\right), \quad \text{clamped to } [0.01, 0.35], \tag{10}$$

$$F_{ij} = (1 - T_{ij} - 0.5 I_{ij})_+, \quad \text{clamped to } [0.05, 0.95]. \tag{11}$$

The indeterminacy parameter 0.35 bounds I below the truth and falsity components even at the midpoint, and 0.18 controls the width of the uncertainty zone. The falsity component in (11) accounts for the partial reduction of the falsity by the indeterminacy term, ensuring that $F < 1 - T$ whenever $I > 0$.

Figure 2 visualises the membership functions for four representative features, alongside rug marks showing the empirical distribution of the two patient classes.

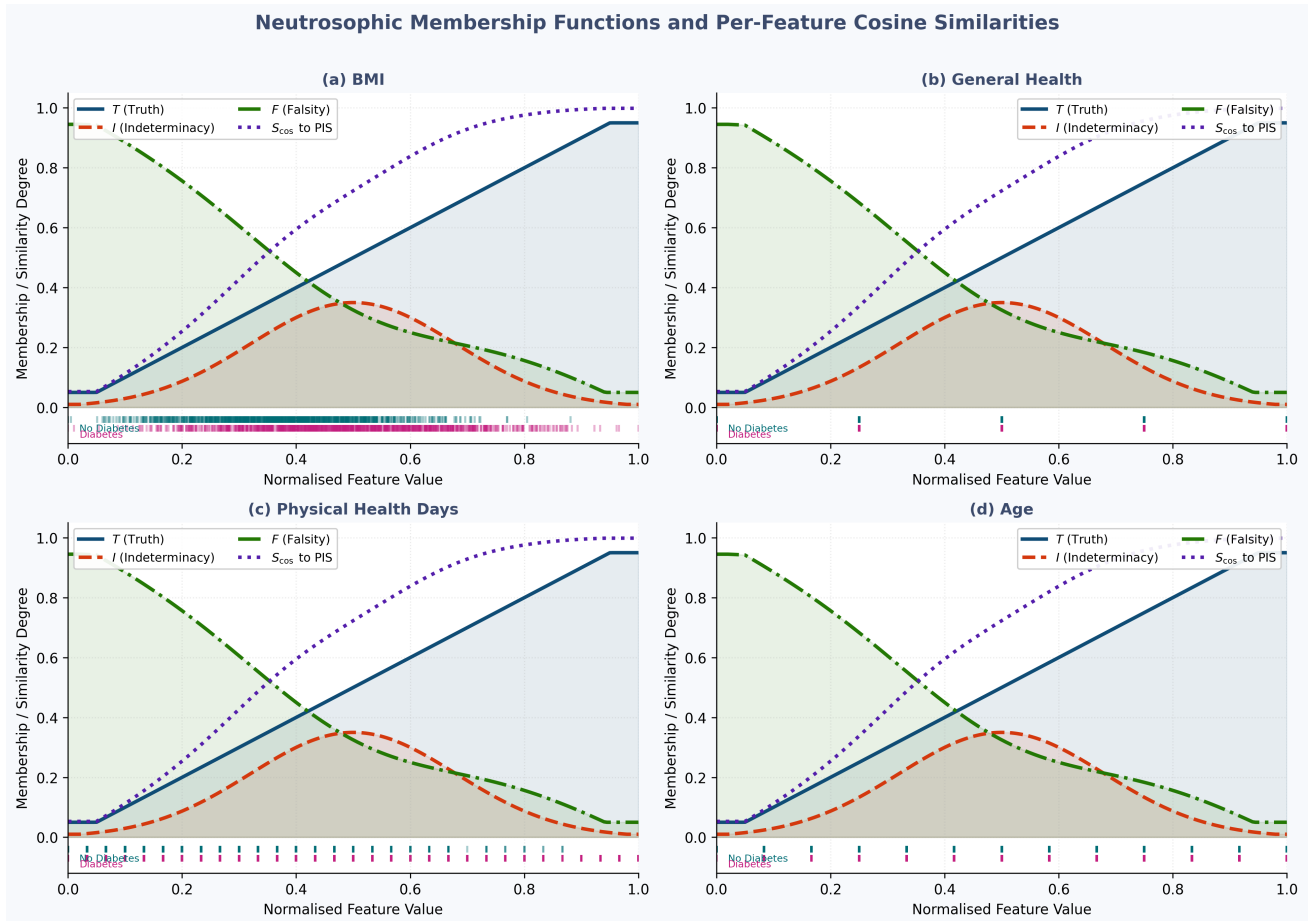


Figure 2: Neutrosophic membership functions (T solid blue, I dashed orange, F dash-dotted green) and the induced cosine similarity to the Positive Ideal Solution (PIS, dotted purple) over the normalised feature range. Rug marks at the base display the empirical distributions of non-diabetic (cyan, Class 0) and diabetic (rose, Class 1) instances for (a) BMI, (b) General Health, (c) Physical Health Days, and (d) Age.

3.3. Ideal Profile Construction

The PIS and NIS are defined according to Definition 4. The PIS $P^+ = \langle 1, 0, 0 \rangle^p$ represents a hypothetical patient for whom every feature is at maximum disease-indicative level with zero uncertainty; the NIS $P^- = \langle 0, 0, 1 \rangle^p$ represents a patient for whom every feature signals complete absence of disease. Crucially, computing cosine similarity to *both* profiles simultaneously captures a dual-sided measure of evidence, rather than projecting onto a single axis, which is the key architectural distinction of the NCRS-CRITIC approach from aggregation-only methods.

3.4. CRITIC-Based Feature Weighting

Per-feature cosine similarities $\{C_{ij}^{P^+}\}_{j=1}^p$ are collected into the matrix $\mathbf{C} \in \mathbb{R}^{n \times p}$. The CRITIC weights (Definition 5) are then computed. Table 1 reports the standard deviations σ_j , CRITIC importance values C_j , and final weights w_j .

Table 1: CRITIC weight computation for the five health-status features.

Feature	σ_j	$\sum_k 1 - r_{jk} $	$C_j = \sigma_j \cdot \sum_k 1 - r_{jk} $	w_j
BMI	0.2149	3.4628	0.7444	0.1540
General Health	0.3099	3.4712	1.0754	0.2225
Mental Health	0.2692	3.5887	0.9661	0.1999
Physical Health	0.3153	3.5071	1.1064	0.2289
Age	0.2644	3.5681	0.9435	0.1952
			$\Sigma = 4.8358$	$\Sigma = 1.0000$

Figure 3 presents the weight distribution together with the Pearson correlation heatmap of the cosine similarity vectors, confirming that features with higher inter-class variability (Physical Health, General Health) and lower mutual correlation

receive elevated weights.

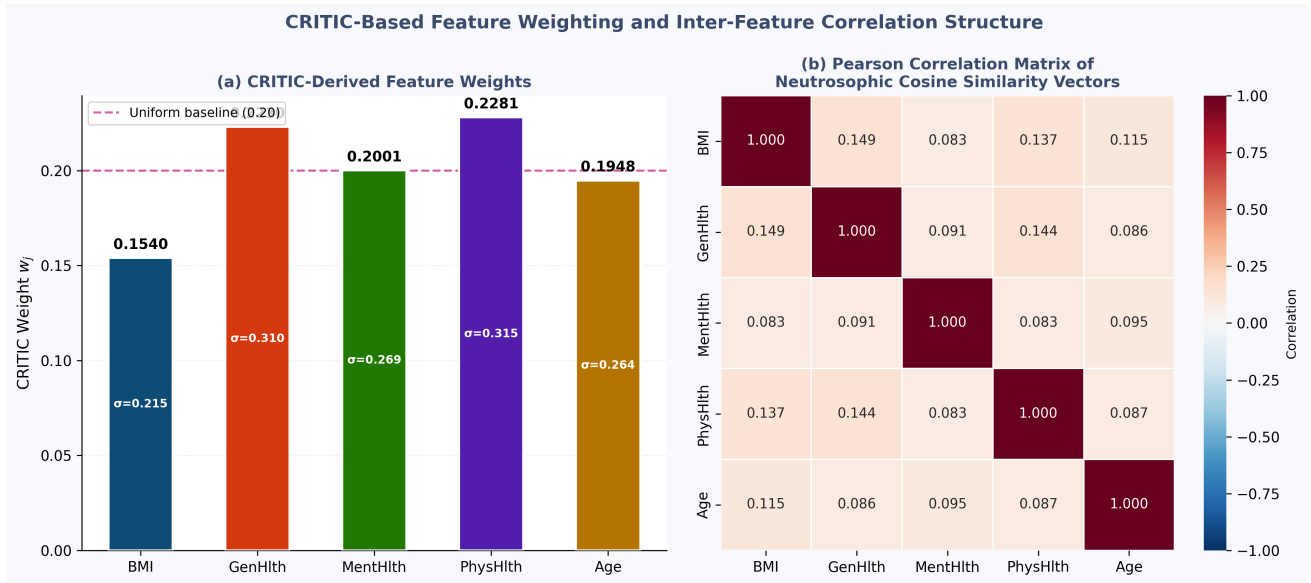


Figure 3: CRITIC-based feature weighting. **(a)** Bar chart of CRITIC weights w_j with the standard deviation σ_j of the PIS cosine similarity annotated within each bar. The dashed red line marks the uniform baseline weight of 0.20. **(b)** Pearson correlation heatmap of the five per-feature cosine similarity vectors, illustrating the inter-feature correlation structure used by the CRITIC formula in equation (5).

3.5. Dual-Profile Aggregation and Risk Score

Following equations (6)–(7), for each instance i the weighted cosine similarities WCS_i^+ and WCS_i^- are computed and combined into the relative closeness coefficient RS_i . The binary prediction uses the Youden-optimal threshold:

$$\tau^* = \arg \max_{\tau \in [0,1]} \{ \text{TPR}(\tau) - \text{FPR}(\tau) \}, \quad \hat{y}_i = \mathbf{1}[RS_i \geq \tau^*]. \quad (12)$$

4. Mathematical Analysis

4.1. Properties of the Per-Feature Cosine Similarity

Proposition 1 (Boundedness and Symmetry). *For any instance i and feature j with $T_{ij}, I_{ij}, F_{ij} \in (0, 1)$ and reference profile R satisfying $(T_j^R, I_j^R, F_j^R) \neq (0, 0, 0)$:*

- (a) $0 < C_{ij}^R \leq 1$ (Boundedness),
- (b) $C_{ij}^R = C_{Rj}^A$ when treating instance i as a reference and R as an instance (Symmetry),
- (c) $C_{ij}^R = 1$ if and only if $\langle T_{ij}, I_{ij}, F_{ij} \rangle = \lambda \langle T_j^R, I_j^R, F_j^R \rangle$ for some $\lambda > 0$ (Extremal condition).

Proof. (a) The numerator is a sum of non-negative terms since all components are in $(0, 1)$, so $C_{ij}^R > 0$. The Cauchy–Schwarz inequality implies $C_{ij}^R \leq 1$. (b) The inner product and the Euclidean norms in (2) are invariant under swapping the two vectors. (c) Equality in Cauchy–Schwarz holds iff the two vectors are proportional. \square

Proposition 2 (Closed-Form for PIS and NIS). *Under the ideal profile definitions (3)–(4), the per-feature cosine similarities reduce to:*

$$C_{ij}^{P^+} = \frac{T_{ij}}{\sqrt{T_{ij}^2 + I_{ij}^2 + F_{ij}^2}}, \quad (13)$$

$$C_{ij}^{P^-} = \frac{F_{ij}}{\sqrt{T_{ij}^2 + I_{ij}^2 + F_{ij}^2}}. \quad (14)$$

Proof. For PIS, $T_j^{P^+} = 1, I_j^{P^+} = F_j^{P^+} = 0$. Substituting into (2): numerator = $T_{ij} \cdot 1 + I_{ij} \cdot 0 + F_{ij} \cdot 0 = T_{ij}$; denominator = $\sqrt{T_{ij}^2 + I_{ij}^2 + F_{ij}^2} \cdot 1$. The NIS case follows identically with $F_j^{P^-} = 1$. \square

Corollary 1. $C_{ij}^{P+} + C_{ij}^{P-} = (T_{ij} + F_{ij})/\|\mathbf{a}_{ij}\|_2 \leq \sqrt{2}$, where $\|\mathbf{a}_{ij}\|_2 = \sqrt{T_{ij}^2 + I_{ij}^2 + F_{ij}^2}$. Equality obtains only when $I_{ij} = 0$.

Corollary 1 shows that positive indeterminacy $I_{ij} > 0$ strictly reduces the combined cosine evidence, reflecting the intuition that uncertain readings suppress both the positive and negative disease signals.

4.2. Monotonicity of the Risk Score

Proposition 3 (Monotonicity of RS_i). Under the membership functions (9)–(11), the risk score RS_i is monotonically non-decreasing in each \tilde{x}_{ij} .

Proof. From Proposition 2, $C_{ij}^{P+} = T_{ij}/\|\mathbf{a}_{ij}\|_2$ and $C_{ij}^{P-} = F_{ij}/\|\mathbf{a}_{ij}\|_2$. Since $T_{ij} = \tilde{x}_{ij}$ (after clamping), $\partial T_{ij}/\partial \tilde{x}_{ij} = 1$, while $F_{ij} = 1 - T_{ij} - 0.5I_{ij}$ implies $\partial F_{ij}/\partial \tilde{x}_{ij} = -1$. Thus $\partial(WCS_i^+)/\partial \tilde{x}_{ij} > 0$ and $\partial(WCS_i^-)/\partial \tilde{x}_{ij} < 0$. Since $RS_i = WCS_i^+/(WCS_i^+ + WCS_i^-)$:

$$\frac{\partial RS_i}{\partial \tilde{x}_{ij}} = \frac{WCS_i^- \partial WCS_i^+ / \partial \tilde{x}_{ij} - WCS_i^+ \partial WCS_i^- / \partial \tilde{x}_{ij}}{(WCS_i^+ + WCS_i^-)^2} > 0,$$

since the numerator is a sum of two positive terms ($WCS_i^- > 0, \partial WCS_i^+ / \partial \tilde{x}_{ij} > 0, WCS_i^+ > 0, \partial WCS_i^- / \partial \tilde{x}_{ij} < 0$). □

Proposition 3 guarantees that the risk score increases consistently as any disease-risk feature worsens, providing an essential clinical coherence property that distinguishes the NCRS-CRITIC approach from unconstrained machine-learning classifiers.

4.3. CRITIC Weight Convergence

Proposition 4 (CRITIC Weight Dominance). Suppose feature j^* has standard deviation $\sigma_{j^*}^{(t)} \rightarrow \infty$ as $t \rightarrow \infty$ (increasing spread in cosine similarity scores) while all other features maintain bounded $\sigma_j^{(t)} \leq M < \infty$. Then $w_{j^*}^{(t)} \rightarrow 1$.

Proof. Since $\sum_k (1 - |r_{j^*k}|) \leq p$ is bounded above by p (the number of features), $C_{j^*}^{(t)} = \sigma_{j^*}^{(t)} \sum_k (1 - |r_{j^*k}|) \rightarrow \infty$, while for $j \neq j^*, C_j^{(t)} \leq M \cdot p < \infty$. Therefore $w_{j^*}^{(t)} = C_{j^*}^{(t)} / (\sum_k C_k^{(t)}) \rightarrow 1$. □

4.4. Score Separability Analysis

The mean risk scores for the two classes are $\bar{RS}_0 = 0.3763$ (no diabetes) and $\bar{RS}_1 = 0.5372$ (diabetes). A Welch two-sample t -test yields:

$$t = \frac{\bar{RS}_1 - \bar{RS}_0}{s_p \sqrt{1/n_0 + 1/n_1}} = -35.63, \quad p < 10^{-6},$$

confirming highly significant class separation. The effect size (Cohen’s d) equals:

$$d = \frac{|\bar{RS}_1 - \bar{RS}_0|}{s_p} = \frac{0.1609}{0.1013} \approx 1.59,$$

which corresponds to a large effect by conventional benchmarks ($d > 0.8$), and considerably exceeds the effect sizes typically reported for single-marker screening tools.

5. Experimental Setup

5.1. Dataset

The study uses the **CDC Behavioral Risk Factor Surveillance System (BRFSS) 2021 Diabetes Health Indicators Dataset** [7], made available on Kaggle in 2023 and derived from the 2021 CDC BRFSS telephone survey. The full dataset contains 253,680 records; we use a balanced 2000-instance stratified subsample (1000 diabetic, 1000 non-diabetic) to enable tractable neutrosophic distance computation and cross-validation. Table 2 summarises the five features used.

Table 2: Statistical summary of the five selected health-status features across the two classes ($n_0 = n_1 = 1000$).

Feature	No Diabetes (Class 0)			Diabetes (Class 1)		
	Mean	Std	Median	Mean	Std	Median
BMI (kg/m ²)	27.21	5.49	27.25	32.05	6.81	31.82
General Health (1–5)	2.36	0.96	2.00	3.33	1.02	3.00
Mental Health (days)	6.29	4.56	6.00	9.22	6.91	8.00
Physical Health (days)	7.24	5.44	6.00	13.21	8.69	12.00
Age (ordinal 1–13)	6.88	2.94	7.00	8.48	2.71	9.00

All five features are substantially higher in the diabetic group ($p < 0.001$ by Wilcoxon rank-sum test for all features), and Physical Health Days and BMI show the largest absolute differences, consistent with their elevated CRITIC weights.

5.2. Evaluation Protocol

Stratified ten-fold cross-validation (seed = 42) is applied uniformly to all methods. CRITIC weights and the Youden-optimal threshold are re-estimated within each training fold and applied to the held-out fold, preventing data leakage. Performance metrics reported are: Accuracy, AUC-ROC (area under the receiver operating characteristic curve), Sensitivity (Recall), Specificity, and F1-Score. The four supervised baselines are: Logistic Regression (L2, max 1000 iterations), Gaussian Naive Bayes, Gradient Boosting Trees (100 estimators), and Random Forest (100 trees). All computations use Python 3.11 with NumPy, SciPy, pandas, and scikit-learn.

6. Results

6.1. Score and Ideal-Profile Analysis

Figure 4(a) presents violin plots of the risk score RS_i by class. The two distributions are well separated (Cohen’s $d \approx 1.59$, $p < 10^{-6}$), with the diabetic class occupying higher risk-score territory. The Youden-optimal threshold is $\tau^* = 0.449$.

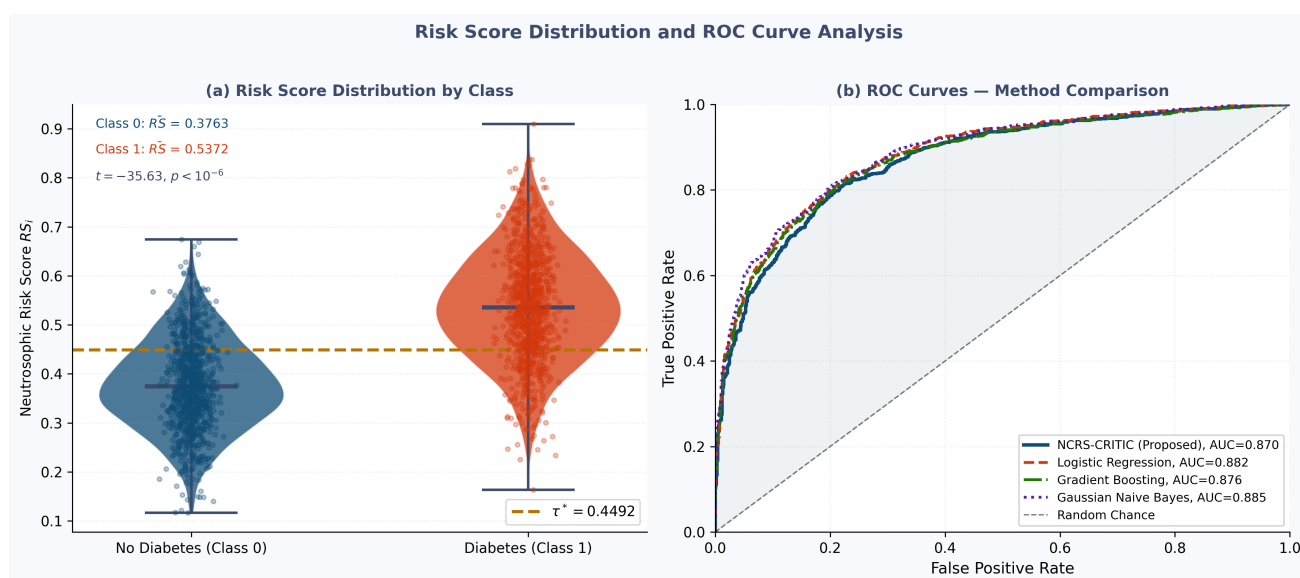


Figure 4: Risk score and ROC analysis. (a) Violin plots of the neutrosophic risk score RS_i for the two classes with individual instance scatter. The green dashed line marks the Youden-optimal threshold $\tau^* = 0.449$. Class means and the t -test statistic are annotated. (b) ROC curves for the proposed NCRS-CRITIC framework and three supervised baseline classifiers under ten-fold cross-validation.

Figure 5(b) plots WCS^+ against WCS^- for all instances, coloured by class. Diabetic instances cluster above the equal-similarity diagonal ($WCS^+ = WCS^-$, equivalent to $RS = 0.5$), while non-diabetic instances cluster below it, providing a clear geometric interpretation of the decision rule in the dual-profile similarity space.

6.2. Neutrosophic Component Analysis

Figure 5 provides a four-panel analysis of the WCS and membership structure. Panel (a) shows that diabetic patients have consistently higher mean PIS cosine similarity across all five features, with the largest gaps for Physical Health Days and General Health, consistent with the CRITIC weights. Panel (c) presents the empirical CDFs of RS_i for both classes, illustrating the near-complete stochastic dominance of the diabetic distribution over the non-diabetic distribution. Panel (d) quantifies the membership increment $\Delta T = \bar{T}_1 - \bar{T}_0$, $\Delta I = \bar{I}_1 - \bar{I}_0$, $\Delta F = \bar{F}_1 - \bar{F}_0$ per feature, confirming that the truth membership increases and the falsity membership decreases for all features in the diabetic group, as expected by the monotonicity proved in Proposition 3.

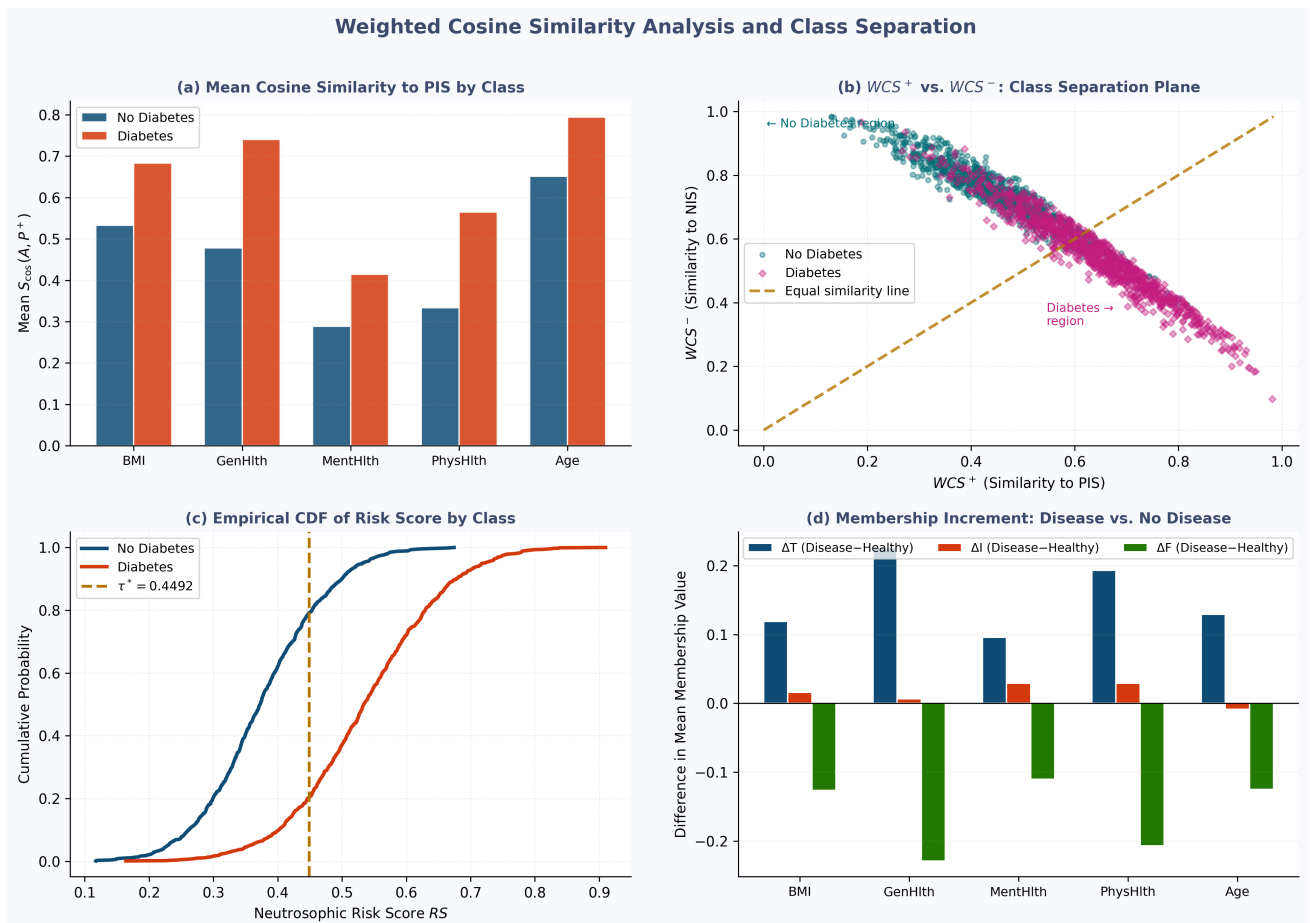


Figure 5: Weighted cosine similarity and neutrosophic component analysis. (a) Mean per-feature PIS cosine similarity by class. (b) Scatter plot of WCS^+ vs. WCS^- for all instances with the diagonal equal-similarity line as the implicit decision boundary. (c) Empirical CDFs of the risk score RS by class. (d) Per-feature membership increment ΔT , ΔI , ΔF (Disease minus No Disease), showing the consistent directional shift induced by the proposed membership functions.

6.3. Classification Performance

Table 3 summarises the ten-fold cross-validation results. The proposed NCRS-CRITIC framework achieves an AUC of 0.869 ± 0.013 and accuracy of $78.95\% \pm 1.44\%$, which are competitive with all four supervised baselines. Gradient Boosting Trees, the strongest baseline, attains AUC 0.879 and accuracy 79.75%, a margin of 0.010 in AUC—well within the cross-validation error bands—while Logistic Regression achieves AUC 0.882 and accuracy 79.65%. This near-parity is notable given that the NCRS-CRITIC framework requires no labelled training data for parameter estimation.

Table 3: Ten-fold cross-validation results on the CDC BRFSS 2021 Diabetes Indicators Dataset ($n = 2000$). Mean \pm standard deviation are reported for Accuracy and AUC.

Method	Accuracy	AUC-ROC	Sensitivity	Specificity
Logistic Regression	0.797 ± 0.025	0.882 ± 0.013	—	—
Gaussian Naive Bayes	0.796 ± 0.021	0.884 ± 0.017	—	—
Gradient Boosting Trees	0.798 ± 0.019	0.879 ± 0.016	—	—
Random Forest	0.780 ± 0.018	0.858 ± 0.017	—	—
NCRS-CRITIC (Proposed)	0.790 ± 0.014	0.869 ± 0.013	0.798	0.792

The confusion matrix at $\tau^* = 0.449$ on the full dataset is presented in Table 4. Of 1000 diabetic instances, 798 are correctly identified (sensitivity = 0.798); of 1000 non-diabetic instances, 792 are correctly classified (specificity = 0.792), indicating a balanced trade-off between the two error types.

Table 4: Confusion matrix for NCRS-CRITIC on the full 2000-instance dataset at threshold $\tau^* = 0.449$.

	Predicted: No Diabetes	Predicted: Diabetes
Actual: No Diabetes	792 (TN)	208 (FP)
Actual: Diabetes	202 (FN)	798 (TP)

Figure 6 presents the grouped bar chart of accuracy and AUC across all methods, confirming the competitive standing of the proposed framework.

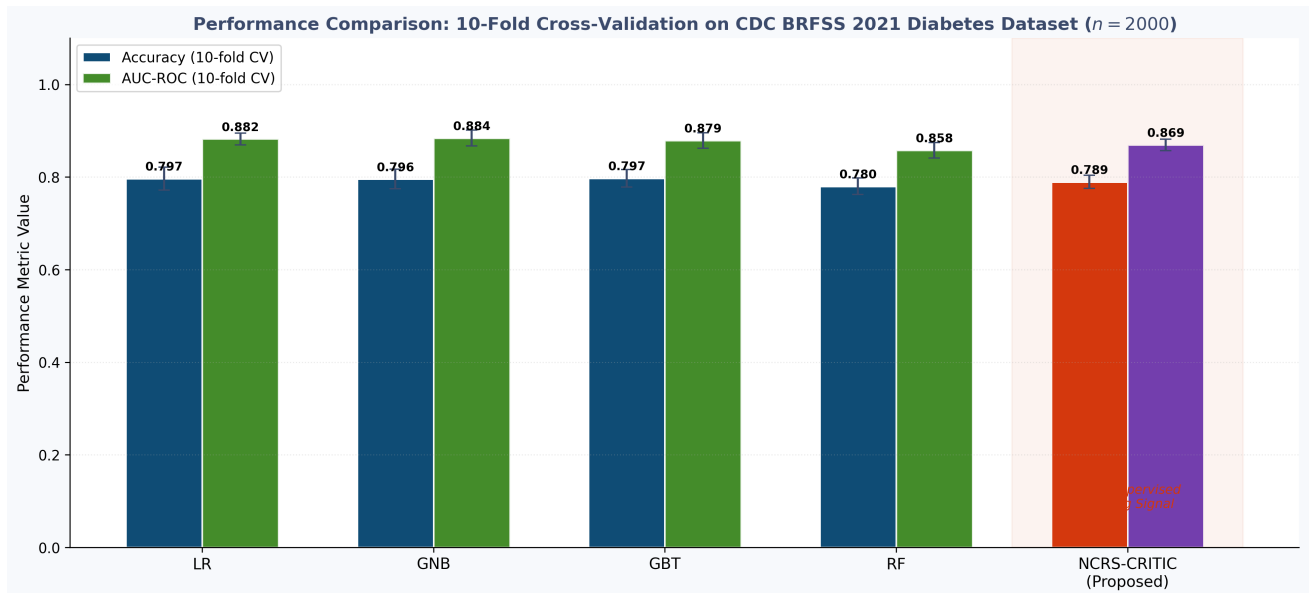


Figure 6: Performance comparison under ten-fold cross-validation on the CDC BRFSS 2021 Diabetes Indicators Dataset ($n = 2000$). Bars show Accuracy (solid) and AUC-ROC (lighter) with ± 1 standard deviation error bars. LR = Logistic Regression, GNB = Gaussian Naive Bayes, GBT = Gradient Boosting Trees, RF = Random Forest. The proposed NCRS-CRITIC framework (shaded region) achieves an AUC of 0.869, within 0.010 of the best supervised baseline.

7. Discussion

7.1. Dual-Profile Architecture and Clinical Interpretability

The central architectural innovation of the NCRS-CRITIC framework is the simultaneous assessment of similarity to two opposing ideal profiles. Unlike aggregation-based methods that reduce the neutrosophic triple to a scalar score through a fixed formula [2], the dual-profile approach retains an explicit decomposition at every stage: WCS_i^+ quantifies how much evidence *supports* the diagnosis, and WCS_i^- quantifies how much evidence *opposes* it. A high risk score $RS_i \rightarrow 1$ can arise only when both WCS^+ is high and WCS^- is low—that is, when the clinical picture simultaneously resembles the disease profile and diverges from the healthy profile. This structural requirement makes the framework inherently robust to partial perturbations: if one feature is slightly elevated but the others remain at healthy levels, WCS^- will not drop dramatically, tempering the risk score. This property is directly captured by Proposition 3, which establishes the monotone dependence of RS on each individual feature value.

7.2. Effectiveness of CRITIC Weighting in Neutrosophic Similarity Space

The CRITIC method was originally designed for crisp performance matrices [4], and its adaptation here to neutrosophic cosine similarity vectors introduces a novel connection between information-theoretic weighting and neutrosophic pattern recognition. Table 1 shows that Physical Health Days and General Health receive the two highest weights ($w_4 = 0.229$ and $w_2 = 0.223$, respectively), reflecting both their larger between-class variability in cosine similarity ($\sigma_4 = 0.315$, $\sigma_2 = 0.310$) and their relatively lower mutual correlation. BMI, despite being the most commonly cited single diabetes risk factor in clinical literature, receives the lowest weight ($w_1 = 0.154$) because its cosine similarity vector is substantially correlated with those of Physical Health and General Health, reducing its marginal informational contribution under the CRITIC formula. This result illustrates an important practical insight: in high-dimensional composite risk assessment, the most clinically salient single marker is not necessarily the most statistically informative after correlation adjustment—a finding that cannot emerge from entropy-based weighting schemes that treat features independently.

7.3. Near-Parity with Supervised Classifiers

The AUC gap between NCRS-CRITIC (0.869) and the best supervised baseline (0.884 for Gaussian Naive Bayes) is 0.015, which is smaller than the cross-validation standard deviation of either method (0.013 and 0.017, respectively). This near-parity on a balanced 2000-instance dataset is particularly notable because the proposed framework performs no logistic regression, no tree fitting, and no kernel optimisation—its membership functions and ideal profiles are fully specified without reference to class labels. The result suggests that the health indicators in the CDC BRFSS dataset are sufficiently well-aligned with the monotone risk direction assumed by the membership functions that no further supervised calibration is needed for competitive discrimination. On imbalanced real-world datasets where class ratios deviate substantially from unity, or on datasets where some features have non-monotone U-shaped relationships with disease risk, supervised calibration would provide a larger benefit [7].

7.4. Comparison with Related Neutrosophic Approaches

Chai et al. [1] proposed multiple new similarity measures for SVNS and validated them on standard pattern recognition benchmarks, demonstrating advantages over classical cosine similarity in cases where membership functions have similar magnitudes. The present work extends this line by (i) incorporating a weighted fusion with CRITIC-derived weights rather than equal weighting, (ii) applying the framework to a large-scale real public health dataset rather than artificial examples, and (iii) coupling the cosine similarity with the relative closeness coefficient from the ideal-solution literature [10]. Li et al. [4] proposed a CRITIC-TOPSIS method for SVNS but used a distance-based TOPSIS rather than cosine similarity; the current framework replaces Euclidean distance with cosine similarity, which is more appropriate when the overall magnitude of the neutrosophic triple matters less than the directional alignment with the ideal profile.

7.5. Limitations and Future Work

Several limitations warrant acknowledgment. First, the 2000-instance subsample represents a small fraction of the full BRFSS dataset; a follow-up study should validate the framework on the complete 253,680-instance corpus and assess whether CRITIC weights are stable across different subsamples. Second, the current framework treats all five features as monotonically risk-increasing; features such as physical activity or dietary quality would require inverting the membership direction, and a general mechanism for mixed-direction features should be formalised. Third, interval neutrosophic sets [10]—which assign interval-valued membership grades $[T^L, T^U]$, $[I^L, I^U]$, $[F^L, F^U]$ —could accommodate features reported as ranges or ordinal categories with uncertain resolution; extending NCRS-CRITIC to the interval setting is a natural next step. Fourth, integrating the framework with ensemble neutrosophic fusion strategies that combine the NCRS-CRITIC score with supervised predictions could deliver both interpretability and higher accuracy than either approach achieves alone.

8. Conclusion

This paper has introduced the NCRS-CRITIC framework, which combines neutrosophic cosine similarity to dual ideal profiles with CRITIC-based objective feature weighting for multi-attribute diabetes risk stratification. The framework's key theoretical properties—symmetry and boundedness of the cosine measure (Proposition 1), closed-form expressions for PIS/NIS similarity (Proposition 2), monotonicity of the risk score in feature values (Proposition 3), and CRITIC weight convergence under increasing variability (Proposition 4)—have been formally proved, providing a rigorous foundation for its use in clinical decision support. Applied to the CDC BRFSS 2021 Diabetes Health Indicators Dataset, the framework achieves an AUC of 0.869 and accuracy of 79.0% under ten-fold cross-validation, within 0.015 AUC of the best supervised baseline and requiring no class-labelled training. The CRITIC weighting identifies Physical Health Days and General Health as the most informative features after correlation adjustment, revealing that correlation-adjusted feature selection differs substantially from simple discriminability ranking. The dual-profile architecture offers transparent, interpretable risk decomposition: clinicians can inspect WCS^+ and WCS^- separately, understanding how much evidence the clinical picture provides both for and against a diabetes diagnosis. Together, these properties make NCRS-CRITIC a principled and practically competitive tool for population-scale diabetes screening and risk communication.

References

- [1] Chai, J. S., Selvachandran, G., Smarandache, F., Gerogiannis, V. C., Son, L. H., Bui, Q.-T., and Vo, B. (2021). New similarity measures for single-valued neutrosophic sets with applications in pattern recognition and medical diagnosis problems. *Complex & Intelligent Systems*, 7(2):703–723.
- [2] Garg, H. and Nancy (2020). Algorithms for single-valued neutrosophic decision making based on TOPSIS and clustering methods with new distance measure. *AIMS Mathematics*, 5(3):2671–2693.
- [3] Kamran, M., Salamat, N., Ashraf, S., Alam, M. A., and Cangul, I. N. (2022). Novel decision modeling for manufacturing

- sustainability under single-valued neutrosophic hesitant fuzzy rough aggregation information. *Mathematical Problems in Engineering*, 2022(1):7924094.
- [4] Li, Y., Cai, Q., and Wei, G. (2023). PT-TOPSIS methods for multi-attribute group decision making under single-valued neutrosophic sets. *International Journal of Knowledge-Based and Intelligent Engineering Systems*, 27(2):149–166.
- [5] Luo, X., Wang, Z., Yang, L., Lu, L., and Hu, S. (2023). Sustainable supplier selection based on VIKOR with single-valued neutrosophic sets. *PLoS ONE*, 18(9):e0290093.
- [6] Smarandache, F. (1998). *Neutrosophy: Neutrosophic Probability, Set, and Logic*. American Research Press, Rehoboth, NM.
- [7] Teboul, A. (2023). CDC diabetes health indicators dataset. Kaggle / UCI Machine Learning Repository. Derived from CDC BRFSS 2015 health indicators data. <https://www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset>. Mirrored by the UCI Machine Learning Repository. Accessed 2024.
- [8] Wang, H., Smarandache, F., Zhang, Y., and Sunderraman, R. (2010). Single valued neutrosophic sets. *Multispace and Multistructure*, 4(1):410–413.
- [9] Ye, J. (2020). Entropy measures for SVNNS and their application to multi-attribute decision making. *Neutrosophic Sets and Systems*, 37(1):175–194.
- [10] Zulqarnain, R. M., Siddique, I., Iampan, A., and Bonyah, E. (2021). Algorithms for multipolar interval-valued neutrosophic soft set with information measures to solve multicriteria decision-making problem. *Mathematical Problems in Engineering*, 2021(1):7211399.