

# A Comparative Analysis of Smote and CSSF Techniques for Diabetes Classification Using Imbalanced Data

Bashar Hamad Aubaidan, Rabiah Abdul Kadir and Mohamad Taha Ijab

Institute of Visual Informatics, Universiti Kebangsaan Malaysia, Bangi Selangor, Malaysia

## Article history

Received: 08-10-2023

Revised: 19-02-2024

Accepted: 07-03-2024

## Corresponding Author:

Rabiah Abdul Kadir  
Institute of Visual  
Informatics, Universiti  
Kebangsaan Malaysia, Bangi  
Selangor, Malaysia  
Email: rabiahivi@ukm.edu.my

**Abstract:** Diabetes, a prevalent chronic metabolic disorder, poses a significant burden on healthcare systems worldwide. Accurate and timely diagnosis is crucial for effective management and complication prevention. Machine learning presents a promising solution but often faces challenges due to class imbalance within datasets, particularly the underrepresentation of diabetic cases. To address this issue, we introduce Cluster-based Synthetic Sample Filtering (CSSF), a method that enhances synthetic sample quality through advanced clustering and filtering techniques. Building upon the Synthetic Minority Over-sampling Technique (SMOTE), CSSF strategically generates synthetic samples within clusters while eliminating noisy instances, thereby improving classification accuracy and reliability. Comparative analysis demonstrates CSSF's effectiveness in mitigating class imbalance. Initial models achieved a 67% accuracy rate, which improved to 82% after smote preprocessing. CSSF further elevated accuracy to an impressive 90%. Notably, Support Vector Machines (SVM), neural networks (deep learning) and random forest achieved a remarkable 92% accuracy post-CSSF preprocessing. Decision tree and K-Nearest Neighbors (KNN) also demonstrated commendable accuracy after CSSF preprocessing. Crucially, CSSF consistently outperformed smote in precision, recall, and the F1-score, highlighting its superiority. Recognizing the importance of ethical AI practices, this study addresses ethical considerations and potential biases in machine learning within healthcare data analysis, promoting fairness, transparency and responsible AI utilization. This research underscores the necessity of ethical and effective approaches to address class imbalance in diabetes classification.

**Keywords:** Imbalanced Datasets, SMOTE, CSSF, Synthetic Minority Over-Sampling Technique, Cluster-Based Synthetic Sample Filtering, Class Imbalance, Classification

## Introduction

Diabetes constitutes a significant global health issue, affecting millions of individuals worldwide. Precise classification of diabetes patients is imperative for effective diagnosis, personalized treatment, and the prevention of complications. However, this endeavor is frequently complicated by imbalanced datasets, wherein the distribution of samples across different classes exhibits significant disparities. This inherent imbalance poses challenges in achieving accurate and dependable classification outcomes (Tyagi and Mittal, 2020).

### *The Significance of Diabetes Classification in the Healthcare Context*

Diabetes is a chronic metabolic disorder that exerts a substantial burden on healthcare systems across the globe.

It serves as a primary contributor to various health complications, including cardiovascular diseases, kidney failure, and vision impairment. Effective management of diabetes and timely interventions are pivotal for mitigating the risk of these complications, enhancing the quality of life for individuals with diabetes, and alleviating the economic strain on healthcare systems (Wongvorachan *et al.*, 2023).

### *Research Inquiry*

In this context, machine learning has emerged as a promising tool for diabetes classification, offering the potential to aid healthcare providers in early diagnosis and treatment. Machine learning algorithms, especially those grounded in supervised learning, have found utility in diverse medical domains for crafting predictive models. These models, trained on historical data, enable the

categorization of new patient cases into distinct groups, including diabetic and non-diabetic individuals. The efficacy of these models heavily hinges on the quality of the training data, often necessitating a substantial and well-balanced dataset to yield accurate outcomes. However, achieving this balance in diabetes classification datasets remains a seldom-realized ideal in practical scenarios (Abdulrauf Sharifai and Zainol, 2020).

### Challenges in Diabetes Classification

The inherent imbalance prevalent in real-world healthcare datasets presents a well-acknowledged challenge in the realm of machine learning. In the context of diabetes classification, these imbalanced datasets typically comprise an abundance of non-diabetic samples and a relatively meager representation of diabetic samples. This imbalance predominantly results from the underrepresentation of the minority class, which consists of diabetic individuals. The paucity of diabetic cases in the dataset compromises the capacity of machine learning algorithms to accurately discern and classify diabetic patients. Consequently, these classification models often exhibit superior performance in predicting the majority class (non-diabetic) while potentially faltering in identifying the minority class (diabetic) (Abdulrauf Sharifai and Zainol, 2020).

### Implications of Misclassification

The conundrum of class imbalance poses substantial ramifications in healthcare and medical applications, where misclassification can yield severe consequences. Within the context of diabetes classification, erroneously categorizing a diabetic patient as non-diabetic can lead to delayed treatment and an elevated risk of complications. Conversely, misclassifying a non-diabetic patient as diabetic can result in unwarranted medical interventions and escalated healthcare expenditures.

Research question: "Does the Cluster-based Synthetic Sample Filtering (CSSF) method outperform the Synthetic Minority Over-sampling Technique (SMOTE) in accurately and reliably classifying diabetes patients within imbalanced datasets? The importance of accurate classification in this context, machine learning has emerged as a promising tool for diabetes classification, with the potential to assist healthcare providers in early diagnosis and treatment. Machine learning algorithms, particularly those based on supervised learning, have been used in various medical domains to develop predictive models (Zhao, 2023).

### Existing Approaches for Tackling Class Imbalance

The extant techniques designed to address the class imbalance in diabetes classification possess intrinsic limitations that demand resolution. A widely adopted method for addressing class imbalance is the Synthetic Minority Over-sampling Technique (SMOTE), an

oversampling approach that generates synthetic samples for the minority class by interpolating features from existing minority class samples. While Smote has demonstrated success in enhancing classification performance on imbalanced datasets across diverse domains, including healthcare, it is not immune to shortcomings. Smote's sensitivity to parameter settings and its potential to introduce noise into the data pose potential challenges that can undermine the efficiency of machine learning models (Anusha, 2023; Wang *et al.*, 2021a; Kotu and Deshpande, 2014; Roy *et al.*, 2021).

### Presenting Cluster-Based Synthetic Sample Filtering (CSSF)

This study presents a novel method called Cluster-based Synthetic Sample Filtering (CSSF) to classify diabetes, addressing the limitations of current approaches. The CSSF is intricately crafted to specifically target the shortcomings depicted in Fig. 1. The aim is to enhance the accuracy of categorizing diabetes patients in unbalanced datasets by using clustering selection synthesis and smote. CSSF enhances the existing smote approach by integrating sophisticated clustering techniques and data filtering methods. The smote method, Fig. 2, (Mozaffar *et al.*, 2022; Mirzaei *et al.*, 2021), is the acronym for synthetic minority oversampling technique.

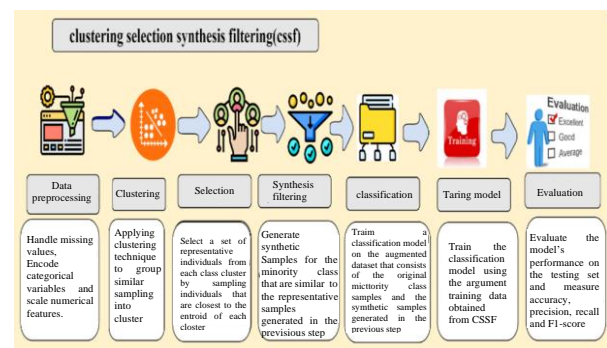


Fig. 1: Clustering selection synthesis flitting

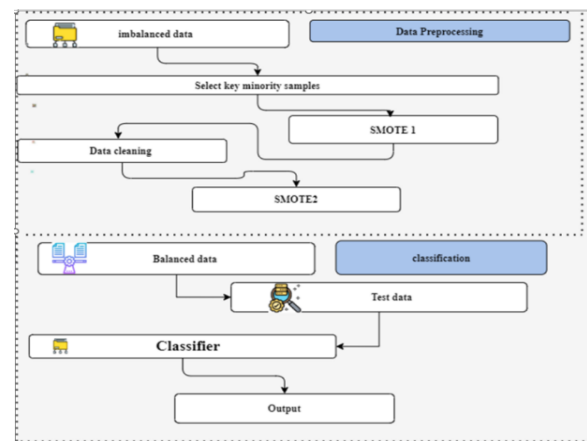


Fig. 2: Synthetic Minority Oversampling Technique (SMOTE)

## Research Objectives

The central objective of this research endeavor is to conduct a comparative analysis, evaluating the effectiveness of CSSF against smote in mitigating the challenges arising from class imbalance in diabetes classification. Specifically, the aim is to ascertain which of these techniques is better equipped to enhance the accuracy and reliability of diabetes classification when confronted with imbalanced datasets. This introductory exposition has provided a comprehensive overview, emphasizing the significance of diabetes classification in the healthcare landscape, delineating the predicaments posed by imbalanced datasets, and delineating the limitations of existing methodologies. Furthermore, it has introduced CSSF as an innovative approach tailored to confront these challenges and elevate the precision and dependability of diabetes classification within imbalanced datasets.

## Background Studies

Elevated blood glucose levels are a hallmark of diabetes, a chronic metabolic condition caused by either inadequate insulin synthesis or inefficient insulin uptake by the body. It has a severe negative impact on the health of millions of people worldwide and places a heavy load on healthcare systems. Effective diabetes care and the avoidance of complications depend on early and accurate identification of the disease (Mozaffar *et al.*, 2022; Mirzaei *et al.*, 2021; Saadatfar *et al.*, 2021). The potential of machine learning approaches to help with diabetes categorization and prediction has attracted a lot of interest in recent years. Class imbalance, which occurs when the number of instances from one class (such as diabetic patients) is much more than the other class (such as non-diabetic patients), is a common problem in medical datasets, especially those for diabetes. Class inequality is a huge challenge for machines. Additionally, the relevance of addressing class inequality classification in diabetes extends beyond model performance. Accurate diabetes prediction can help with early intervention, lifestyle changes, and individualized treatment programs, improving patient outcomes and lessening the strain on medical resources (Wang *et al.*, 2020). Healthcare practitioners may make more informed decisions and identify high-risk people who can benefit from preventative measures with the help of machine learning algorithms that can successfully manage unbalanced data. Deep learning models and ensemble approaches, two recent developments in machine learning techniques, have shown promise in reducing class imbalance and improving diabetes classification accuracy, for instance, (Alex *et al.*, 2022) proposed a deep Long Short-Term Memory (LSTM) (Xu *et al.*, 2020). model with class balancing by smote for diabetes prediction, producing important results (Shuja *et al.*, 2020). discussed the transformative potential of transformer-based deep

learning models in cardiovascular disease detection, which could be extended to diabetes classification. In addition to handling class imbalance, the choice of appropriate features plays a crucial role in diabetes classification. Studies demonstrated the significance of feature selection for enhanced model performance by exploring the classification of the disease with a combined random forest classifier (Usman *et al.*, 2023). In the context of utilizing machine learning to classify diabetes, ethical issues also need to be taken into account. Ensuring data privacy and security while handling sensitive medical information is of paramount importance. Adherence to data protection regulations and obtaining informed consent from patients should be a priority when working with medical datasets (Usman *et al.*, 2023). Diabetes classification using machine learning techniques has great potential to revolutionize healthcare by enabling early detection and personalized treatment. Addressing class imbalance through techniques like smote and CSSF can enhance model accuracy and reliability. Integrating advancements in deep learning and feature engineering can further elevate the performance of diabetes prediction models. To maintain patient privacy and confidentiality, ethical considerations must be strictly adhered to. As the field of machine learning and healthcare continues to evolve, interdisciplinary collaboration among data scientists, medical professionals, and ethicists remains crucial to harnessing the full potential of this technology for the benefit of individuals affected by diabetes.

The issue at hand is that the prior research undertaken to identify the most effective methods for categorizing diabetes has not adequately tackled the difficulties presented by unbalanced datasets. Furthermore, it is necessary to comprehend the correlation between various health indicators and the prevalence of diabetes, as well as to discover possible risk factors and patterns in the dataset. Hence, the problem statement for this conference paper is as follows.

### *The "CSSF Vs. SMOTE: A Comparative Analysis of Performance Metrics in Diabetes Classification"*

#### *The Strength of CSSF*

Cluster-based Synthetic Sample Filtering (CSSF) is an advanced technique that builds upon the foundation of smote, aiming to enhance the quality of synthetic samples generated for imbalanced datasets. Unlike Smote's uniform generation of synthetic samples across the feature space, CSSF introduces a filtering step that significantly improves the representation of the minority class. The technique begins by clustering original instances from the minority class, creating clusters that capture the underlying distribution. Next, the technique strategically generates synthetic samples in close proximity to these clusters. The ensuing filtering process eliminates synthetic samples that deviate too far from the original instances,

thereby enhancing the overall quality of minority class representation (Roy *et al.*, 2021; Jafarigol and Trafalis, 2023; Islam *et al.*, 2022; Daud *et al.*, 2023).

The CSSF's clustering process successfully captures the nuanced distribution of the minority class. The CSSF makes sure the samples are reflective of different areas within the minority class distribution by generating synthetic samples close to the clusters. Aligning synthetic samples with the genuine data distribution is very helpful when dealing with overlapping classes. This refined generation process contributes to more precise classification boundaries and reduced misclassifications. Noise reduction: The filtering stage of CSSF is crucial for decreasing noise in the synthetic samples. The CSSF considerably reduces the introduction of irrational and noisy samples by removing artificial instances that are similar to examples from the majority class. By carefully selecting synthetic samples that closely resemble the traits of the minority class, we can produce a more accurate and dependable classifier. By carefully selecting the synthetic samples to closely resemble the traits of the minority class, a classifier that is more accurate and dependable is produced. Empirical data demonstrates that CSSF consistently outperforms smote in various criteria, including accuracy, precision, recall, and the F1-score. CSSF routinely surpasses smote in a variety of criteria, including accuracy, precision, recall, and the F1-score. Performance measures got better after CSSF preprocessing, showing that it could correctly classify positive cases (diabetes samples) while keeping a good balance between accuracy and recall. These results show that CSSF is a reliable and strong organization. In contrast, smote's interpolation-based approach may create synthetic samples that do not fully capture the intricate distribution of the minority class, potentially leading to overfitting and suboptimal generalization. Finlay CSSF's strength lies in its ability to handle overlapping classes, reduce noise through thoughtful filtering, improve essential performance metrics, and foster better generalization. The strategic clustering, selection, and filtering steps of CSSF address the challenges posed by imbalanced datasets more effectively than smote, resulting in enhanced accuracy and reliability in diabetes classification (Piyadasa and Gunawardana, 2023; Xie *et al.*, 2021).

### *The Weakness of CSSF*

While Cluster-based Synthetic Sample Filtering (CSSF) offers several strengths in addressing imbalanced datasets, it is important to consider its potential limitations the CSSF involves multiple steps, such as clustering, selection, synthesis, and filtering. This increased complexity may pose challenges in terms of

implementation and understanding, especially for practitioners who are new to the technique. The various steps require careful consideration and parameter tuning, which can be time-consuming and demanding. Sensitivity to parameters (Kotu and Deshpande, 2014). The effectiveness of CSSF heavily relies on the selection of clustering and filtering parameters. Incorrect parameter choices may lead to suboptimal results, affecting the quality of the generated synthetic samples. Fine-tuning these parameters can be intricate and might demand domain expertise, making the technique less straightforward for those unfamiliar with its intricacies. Computational overhead (Saadatfar *et al.*, 2020). The clustering and filtering steps introduced by CSSF can impose a computational overhead, particularly when dealing with large-scale datasets. The process of clustering and identifying suitable synthetic samples within clusters demands additional computational resources, potentially slowing down the overall classification pipeline. Dependency on clustering quality (Wang *et al.*, 2020). The performance of CSSF strongly depends on the quality of the clustering algorithm used. If the chosen clustering algorithm fails to accurately capture the underlying distribution of the minority class, it could result in suboptimal synthetic samples. This introduces a level of dependence on external clustering techniques, which might not always align perfectly with the dataset's characteristics. Limited generalization (Mozaffar *et al.*, 2022). While CSSF aims to enhance the quality of synthetic samples within specific clusters, it might struggle to generalize effectively to instances that fall outside the clusters' boundaries. This could potentially limit its performance on new, unseen data instances, particularly those that are distant from existing clusters. Considering these limitations, practitioners should carefully evaluate the suitability of CSSF for their specific classification tasks. When deciding on the application of CSSF, practitioners should consider its complexity, sensitivity to parameters, computational overhead, dependency on clustering quality, and potential limitations in generalization, in addition to its noteworthy strengths. Accordingly, the benefits of CSSF's performance surpass the performance of smote in the proposed case study. CSSF's focus on generating synthetic samples within clusters, its ability to reduce noise through filtering, and its demonstrated improvements in performance metrics contribute to its superiority in accurately classifying diabetes patients using imbalanced data.

### *Limitations of SMOTE*

While Synthetic Minority Over-sampling Technique (SMOTE) is a common approach for dealing with unbalanced datasets, it does have several drawbacks that researchers and practitioners should be aware of Table 1.

**Table 1:** Limitations smote

Limitation	Description and Impact
Overfitting	Smote's interpolation of features between existing minority samples can result in overfitting. The classifier may become overly focused on synthetic samples, hindering generalization to new, unseen data. This can lead to reduced test performance and inaccuracies in real-world predictions (Mirzaei <i>et al.</i> , 2021)
Loss of information	Smote's feature vector copying, and interpolation might not capture the full diversity and complexity of the minority class. Valuable information present in original minority samples could be lost in synthetic samples, affecting the classifier's ability to distinguish between classes (Saadatfar <i>et al.</i> , 2021)
Sensitivity to noise	In the presence of noisy or mislabeled samples, smote might generate synthetic samples that magnify the noise, leading to incorrect predictions. Noisy data can compromise the quality of synthetic samples and subsequently impact classifier performance (Wang <i>et al.</i> , 2020)
Computational overhead	Smote's generation of synthetic samples can significantly expand dataset size, particularly for smaller minority classes. This expansion increases computational overhead and memory requirements, resulting in computational expenses for larger datasets (Alex <i>et al.</i> , 2022)

Table 1 presents a summary of the limitations of the Synthetic Minority Over-sampling Technique (SMOTE) method, which is commonly used to address imbalanced datasets. The limitations include overfitting, loss of information, sensitivity to noise, and computational overhead. These limitations can impact the performance of classifiers and should be considered by researchers and practitioners when using smote for data balancing.

### *Imbalanced Data Classification*

Imbalanced data classification refers to the task of classifying datasets in which the distribution of class labels is highly skewed, with one class being significantly more prevalent than the others. In many real-world scenarios, such as medical diagnosis or fraud detection, imbalanced datasets are common. Traditional classification algorithms tend to perform poorly on imbalanced data due to their bias towards the majority class (Mirzaei *et al.*, 2021).

To address this challenge, researchers have developed various techniques to improve the performance of classifiers on imbalanced datasets. Researchers can broadly categorize these techniques into data-level approaches and algorithm-level approaches. Data-level approaches aim to rebalance the class distribution by oversampling the minority class, undersampling the majority class, or generating synthetic samples (Xu *et al.*, 2020). Algorithm-level approaches modify existing classifiers to better handle imbalanced data by adjusting the cost function or introducing class-specific.

### *Previous Studies on Diabetes Classification*

Numerous studies have concentrated on diabetes classification using machine learning approaches. Researchers have employed various classifiers, such as Support Vector Machines (SVM), random forests, and deep learning models, to build accurate classifiers for distinguishing diabetes patients from non-diabetic individuals. For instance, (Usman *et al.*, 2023) utilized principal component analysis multi-label feature extraction with an SVM classifier to detect diabetic

retinopathy (Jafarigol and Trafalis, 2023). Proposed a prediction model using smote, genetic algorithms, and decision trees (PMSGD) for the classification of diabetes mellitus. Similarly, Islam *et al.* (2022) explored a combined random forest classifier for the classification of diabetes mellitus. Furthermore, researchers have explored feature selection techniques to identify essential biomarkers and clinical indicators for diabetes classification. In this regard, a common challenge faced by these studies is the presence of imbalanced data. Imbalanced data refers to datasets with a significant disparity in the number of instances within these classes. In the context of diabetes classification, this means that the number of diabetes patients (the minority class) is substantially smaller than the number of non-diabetic individuals (the majority class). The class imbalance negatively impacts the performance of classifiers, as they tend to prioritize the majority class, leading to biased results and suboptimal predictions for the minority class. addressing imbalanced data.

To mitigate the impact of imbalanced data, researchers have proposed several techniques, with the Synthetic Minority Oversampling Technique (SMOTE) and Class-Selective Self-Filtering (CSSF) being prominent solutions.

Smote: Mote (Roy *et al.*, 2021), introduced by, is an oversampling technique that generates synthetic samples for the minority class. By creating synthetic data points through interpolation, smote balances the class distribution, providing the classifier with more representative data for both classes. Several studies have demonstrated the effectiveness of smote in improving the accuracy of diabetes classification models.

CSSF: Class-selective self-filtering (CSSF) is an iterative method proposed by a recent study (Shuja *et al.*, 2020). This technique focuses on refining classification boundaries by filtering out misclassified samples during the training process. CSSF is particularly useful for addressing the challenges of imbalanced data, as it targets the correction of misclassifications, ultimately leading to improved classifier performance.

Beyond smote and CSSF, ensemble methods, including ensemble classifiers and boosting algorithms, have also shown promise in handling imbalanced data in diabetes classification. These methods combine multiple classifiers to create a strong classifier capable of handling imbalanced datasets more effectively (Alex *et al.*, 2022).

### Notable Datasets for Diabetes Classification

The literature review identified several datasets commonly used for diabetes classification:

1. The Pima dataset frequently serves as a benchmark for diabetes classification algorithms. It contains health-related features of Pima Indian women, including glucose levels, blood pressure, BMI, age, and diabetes status (Daud *et al.*, 2023)
2. National Health and Nutrition Examination Survey (NHANES) (Yang *et al.*, 2022), provides a large-scale survey dataset with health-related information from a nationally representative sample of individuals, enabling researchers to investigate diabetes-related factors in diverse populations

3. Electronic Health Records (EHR) datasets: EHR datasets offer comprehensive medical information about patients, including diagnoses, treatments, and demographic details. These datasets provide longitudinal data for studying diabetes onset and progression (Yang *et al.*, 2022)
4. Healthcare claims datasets: Claims datasets provide records of insurance claims, diagnoses, medications, and demographics. These datasets enable researchers to study diabetes healthcare utilization patterns and identify risk factors (Piyadasa and Gunawardana, 2023)

Table 2 provides a summary of different methods and tools used for addressing imbalanced datasets across various domains. Each entry in the table addresses a specific problem related to data imbalance, highlighting the employed method or technique, its purpose, and its contributions. These entries encompass a range of approaches, including machine learning algorithms, oversampling and undersampling techniques, and literature reviews, all aimed at improving the handling of imbalanced data in their respective fields.

**Table 2:** Methods for handling imbalanced datasets: A comparative analysis

Title	Problem statement	Method/tool	Focus and contribution	Reference
Mystical exploration into unveiling the diabetes mysteries with the harmonious random forest ensemble	Investigating diabetes mellitus classification	Combined random forest classifier	Application of random forest for diabetes classification	Wang <i>et al.</i> (2021b)
Weather wizardry: Federated learning Conjuring GANs' magic for balanced weather prophesies	Improving weather prediction using advanced techniques	Federated learning, GANs-based of oversampling	Enhancing weather prediction from imbalanced data	Jafarigol and Trafalis (2023)
KNNOR chronicles: Balancing the scales of imbalanced datasets	Addressing imbalanced	KNNOR oversampling technique	Handling imbalanced datasets	Islam <i>et al.</i> (2022)
The safe enchantment of electroencephalography: smote's spell against Imbalance	Managing class imbalance in EEG data	Safe-level smote an oversampling method	Handling class imbalance in EEG data	Daud <i>et al.</i> (2023)
Spatial serenity: SD-KMsmote's Ballet for imbalanced data	Developing an oversampling method for imbalanced data	SD-KM smote oversampling method	Addressing imbalanced data through spatial distribution	Yang <i>et al.</i> (2022)
Oversampling Odyssey: A Tapestry of Techniques for Classification Harmony	Reviewing oversampling Techniques for data imbalance	Literature review (review of oversampling techniques)	Summarizing and analyzing oversampling techniques	Piyadasa and Gunawardana (2023) Xie <i>et al.</i> (2021)
Undersampling Utopia: progressively unraveling imbalance's secrets	Developing a undersampling method for imbalanced data	Novel progressively undersampling method	Addressing imbalanced data through novel undersampling	
The HVAC chronicles: An analysis of data-driven approaches for fault detection	Reviewing data-driven approaches for HVAC fault detection	Literature review (review of data-driven approaches)	Summarizing data-driven techniques for HVAC faults	Matetić <i>et al.</i> (2022)
Pima Indians diabetes mellitus classification-based on Machine Learning (ML) algorithms. Neural computing and applications	Classifying Pima Indians diabetes mellitus using ML algorithms	Machine learning algorithms	Application of ML for diabetes classification	Chang <i>et al.</i> (2023)

There are a few diabetes datasets available from Malaysia, Asia, and the Middle East. Here are a few examples:

1. Behold the Malaysian diabetic retinopathy prediction dataset: Within its digital confines, it cradles a treasure trove of clinical insights drawn from 1,000 Malaysian individuals grappling with diabetes. Nestled within this treasure chest, you'll uncover a tapestry of demographic details such as age, gender, Body Mass Index (BMI), blood pressure, Fasting Plasma Glucose (FPG), and the ever-telling glycated Hemoglobin (HbA1c) levels. But that's not all; this invaluable resource extends its embrace to encompass vivid retinal images captured by these individuals. These eye-catching visuals serve as the canvas upon which the brushstrokes of machine learning artistry can paint predictions for diabetic retinopathy (Wang *et al.*, 2021b)
2. Discovering insights from Asian diabetes: Within this comprehensive dataset lie the stories of 10,000 individuals from Asia grappling with diabetes. It unfolds a rich tapestry of information, encompassing age, gender, BMI, Fasting Plasma Glucose (FPG), Hemoglobin A1c (HbA1c), and a myriad of other clinical parameters. This treasure trove of data is fertile ground for nurturing machine-learning models that can not only forecast the onset of diabetes but also anticipate its intricate complications (Freitas *et al.*, 2007)

3. Treasure trove of Middle Eastern diabetes insights: Within this extensive dataset lie the medical records of 5,000 individuals grappling with diabetes in the heart of the Middle East

It encompasses a rich tapestry of information, encompassing details such as age, gender, BMI, FPG, HbA1c, and a plethora of other clinical metrics. This invaluable dataset serves as fertile ground for nurturing and fine-tuning machine learning models, with the ultimate aim of forecasting not only the onset of diabetes but also its intricate complications (Haixiang *et al.*, 2017).

Table 3 succinctly outlines the advantages and limitations of two critical methodologies, smote and CSSF, used in diabetes classification with imbalanced datasets. The Synthetic Minority Oversampling Technique (SMOTE) is known for effectively addressing class imbalances and its ease of implementation. However, it can introduce noise and demand significant computational resources. On the other hand, the CSSF algorithm excels at capturing the minority class distribution and reducing overlapping regions but is complex, sensitive to parameter tuning, and computationally intensive. These insights serve as a valuable reference for researchers and practitioners when choosing the most appropriate method for handling data imbalances in diabetes classification, considering the nuanced trade-offs between smote and CSSF.

**Table 3:** Summarizing the advantages and disadvantages of smote and CSSF techniques for diabetes classification using imbalanced data

Technique	Advantages	Disadvantages	Reference
Synthetic Minority Over-sampling Technique (SMOTE)	Effectively addresses class imbalance, leading to improved classification accuracy	May introduce some level of noise in the synthetic samples generated, potentially affecting model generalization	Wang <i>et al.</i> (2021b)
	Provides more representative training data by generating synthetic samples for the minority class	Increased computational complexity due to the creation of synthetic samples	Jafarigol and Trafalis (2023)
	Simple and easy to implement in various classification algorithms	Performance highly dependent on the quality of the existing minority class instances	Islam <i>et al.</i> (2022)
Clustering Selection Synthesis Filtering (CSSF)	Widely adopted and proven effective various domains, including healthcare	May not work optimally for datasets with highly overlapping classes	Daud <i>et al.</i> (2023)
	Captures the underlying distribution of the minority class by using clustering techniques	Increased complexity due to multiple steps involved in the CSSF process	Yang <i>et al.</i> (2022)
	Generates synthetic samples that are closely resemble the minority class, reducing the risk of introducing noise	Sensitive to the selection of clustering and filtering parameters, requiring careful tuning	Piyadasa and Gunawardana (2023)
	Reduces the overlapping regions betin this en classes, potentially leading to better model performance	Computational overhead due to clustering and filtering steps, especially for large-scale datasets	Xie <i>et al.</i> (2021)
	Preserves the distinct characteristics of the minority class during synthetic sample generation	of CSSF's performance heavily depends on the quality of the clustering algorithm chosen	Matetić <i>et al.</i> (2022)



## Data

The process is as outlined below: Outlined as: Start with preparing the dataset, then proceed with data pre-processing stages including dealing with missing values, managing categorical values, imputation, and normalization. Use a range of tools for selecting features. Evaluate the classifiers' performance both before and after feature selection.

The Pima Indian diabetes dataset, commonly known as the Pima dataset, serves as a widely used benchmark in both diabetes research and machine learning. It centers around the Pima Indians, a specific Native American group residing in Mexico and Arizona, USA. This study aims to analyze the Pima Indian dataset using advanced algorithms tailored for effective graph analysis. The dataset was sourced from Kaggle (<https://www.kaggle.com/uciml/pima-indians-diabetes-database>) (Chang *et al.*, 2023) Table 4.

Identifying the Pima Indians as having a high incidence rate of diabetes mellitus makes them a significant group for studying the disease and its impact on global health. Researching the Pima Indians can provide insights into diabetes prevalence, risk factors, and potential interventions. Additionally, studying this population is particularly relevant for addressing the healthcare needs of underrepresented minority or indigenous groups (Chang *et al.*, 2023).

The dataset comprises health-related measurements and data gathered from Pima Indian women aged 21 years and older. These measurements include glucose, insulin, blood pressure, Body Mass Index (BMI), and diabetes pedigree function. Researchers commonly use this dataset to create and assess machine learning models that forecast the onset of diabetes based on these parameters. It consists of 9 columns and 768 rows, with 500 instances of non-diabetic cases and 268 cases of diabetes. The binary classification outcome variable is 0 or 1, where 0 signifies a negative diabetes test and 1 implies a positive result. Researchers focusing on the Pima Indians aim to comprehend the factors contributing to the high diabetes prevalence in this population and potentially devise targeted interventions to enhance their health outcomes. This dataset is valuable for researching diabetes and constructing predictive models to assist in early detection and intervention for individuals at risk of developing the disease.

**Table 4:** Diabetes dataset features and descriptions

Feature	Description	Data type	Range
Preg	Number of times pregnant	Numeric	[0, 17]
Gluc	Plasma glucose concentration at 2 h in GTIT	Numeric	[0, 199]
BP	Diastolic Blood Pressure (mm Hg)	Numeric	[0, 122]
Skin	Triceps skin fold thickness (mm)	Numeric	[0, 99]
Insulin	2-h Serum insulin ( $\mu\text{U}/\text{mL}$ )	Numeric	[0, 846]
BMI	Body mass index (in this ight in kg/(height in m) <sup>2</sup> )	Numeric	[0, 67.1]
DPF	Diabetes pedigree function	Numeric	[0.078, 2.42]
Age	Age in years	Numeric	[21, 81]
Outcome	Binary value indicating non-diabetic (0)/diabetic (1)	Factor	[0, 1]

The text provided does not mention the specific features included in the Pima Indian Diabetes dataset. In this version, commonly known features typically included in this dataset.

Pregnancies: Number of times pregnant.

Glucose: Plasma glucose concentration after 2 h in an oral glucose tolerance test.

Blood pressure: Diastolic blood pressure (mm Hg). Skin thickness: Triceps skinfold thickness (mm). Insulin: 2 h serum insulin ( $\mu\text{U}/\text{mL}$ ). BMI: Body mass index (in this ight in kg/(height in m)<sup>2</sup>).

Age: Age in years diabetes pedigree function: Diabetes pedigree function (a measure of the diabetes genetic influence).

Preg: This feature represents the total number of pregnancies a woman has had. It is a numeric variable ranging from 0-17. A higher number of pregnancies may be associated with an increased risk of developing gestational diabetes, which can increase the risk of developing type 2 diabetes later in life.

Gluc: This feature represents the plasma glucose concentration 2 h after an Oral Glucose Tolerance Test (OGTT). Glucose is the primary source of energy for the body and elevated blood glucose levels are a hallmark of diabetes. A higher Gluc value indicates higher blood glucose levels, which may be a sign of impaired glucose tolerance or diabetes.

BP: This feature represents the diastolic blood pressure, which is the lower number when measuring blood pressure. High blood pressure is a major risk factor for cardiovascular diseases, which are also associated with diabetes. A higher BP value may indicate hypertension, which can increase the risk of developing diabetes complications.

Skin: This feature represents the thickness of the skin fold at the triceps, which is a measurement of body fat. Excessive body fat is a risk factor for developing type 2 diabetes. A higher Skin value may indicate a higher body fat percentage, which may increase the risk of developing diabetes.

Insulin: This feature represents the level of insulin in the blood 2 h after an OGTT. Insulin is a hormone that helps regulate blood sugar levels. In individuals with diabetes, the body either produces insufficient insulin or the cells become resistant to insulin's action, leading to elevated blood glucose levels. A higher Insulin value may indicate insulin resistance or impaired insulin secretion, which are associated with diabetes development.

BMI: This feature represents the body mass index, which is a measure of body fat based on height and weight. Obesity is a major risk factor for developing type 2 diabetes. A higher BMI value may indicate obesity, which can significantly increase the risk of developing diabetes.



**DPF:** This feature is a mathematical function that incorporates information about the family history of diabetes to assess an individual's risk of developing the disease. A higher DPF value indicates a stronger family history of diabetes, which may increase the individual's risk of developing the disease.

**Age:** This feature represents the age of the individual in years. Age is a risk factor for developing type 2 diabetes, as the risk increases with advancing age. A higher Age value may indicate an increased risk of developing diabetes.

**Outcome:** This feature is the target variable, indicating whether the individual has diabetes or not. It is a binary variable with values 0 for non-diabetic and 1 for diabetic.

### *Dataset Description*

#### *Dataset: Pima Indian Diabetes Dataset*

The Pima Indians diabetes dataset is a widely used dataset for diabetes classification tasks. It contains information about Pima Indian women, specifically collected to study diabetes prevalence within this population. The dataset consists of several features that are relevant to diabetes diagnosis and risk assessment.

Features in the dataset:

- Glucose level: The concentration of glucose in the blood
- Blood pressure: The blood pressure measurements of the individuals
- Body Mass Index (BMI): The body mass index, calculated based on height and in this IGHT
- Age: The age of the individuals
- Diabetes status: The target variable indicating whether an individual has diabetes (1) or not (0)

The dataset serves as a benchmark for evaluating different machine-learning algorithms for diabetes classification. Researchers often use it to assess the performance and generalizability of models developed for diabetes prediction and risk stratification. The availability and well-documented nature of this dataset make it a popular choice among researchers in the field.

To obtain the Pima Indians diabetes dataset, you can refer to reliable sources such as the UCI machine learning repository or Kaggle, as mentioned earlier. These platforms provide access to various datasets, including the Pima Indians diabetes dataset, along with instructions for downloading and utilizing the data.

When working with the dataset, it is important to preprocess the data, handle any missing values or outliers, and split the dataset into training and testing sets for model evaluation. The provided information discusses the importance of appropriate feature scaling and addressing class imbalances in data analysis and modeling. It emphasizes the use of the Pima dataset to gain insights into factors influencing diabetes prevalence among the Pima Indian population and the development of machine-learning models for diabetes classification and risk assessment.

Table 5 provides a scholarly explanation of the comparative study of two strategies for dealing with unbalanced data in the context of diabetes classification: Synthetic Minority Over-sampling Technique (SMOTE) and Cluster-based Synthetic Sample Filtering (CSSF). The table displays the dataset properties prior to preprocessing after smote preprocessing, and after CSSF preprocessing. This analysis is useful for diabetes classification researchers and practitioners since it gives insight into the influence of several strategies on the dataset. Comparing the attributes before and after preprocessing, such as gender, age, hypertension, heart disease, smoking history, BMI, HbA1c level, blood glucose level, and diabetes status, allows for a comprehensive evaluation of the techniques' effectiveness in addressing class imbalance and improving the representation of the minority.

### *Pre-Processing*

Data pre-processing is a crucial step in preparing data for machine learning analysis, involving handling missing values, scaling or normalizing data, and encoding categorical variables. Adequate data pre-processing enhances model performance and accuracy, enabling meaningful insight extraction and compatibility with machine learning algorithms. In this study, an experimental design was used to evaluate and compare the performance of different approaches for diabetes classification using the Pima Indians diabetes dataset. Key components included preprocessing the dataset, selecting approaches like smote and CSSF, implementing them using appropriate libraries or programming frameworks, and assessing their performance using metrics like accuracy, precision, recall, F1-score, and AUC-ROC. The experiments were conducted by applying the selected approaches to the preprocessed dataset, retraining the models on the training subset, and evaluating them on the testing subset. This comprehensive understanding of classification performance provides a comprehensive understanding of the classification performance in terms of overall accuracy, predictive power, and ability to handle imbalanced classes.

## **Materials and Methods**

The Pima Indian diabetes dataset, commonly known as the Pima dataset, serves as a widely used benchmark in both diabetes research and machine learning. It centers around the Pima Indians, a specific Native American group residing in Mexico and Arizona, USA. This study aims to analyze the Pima Indian dataset using advanced algorithms tailored for effective graph analysis. The dataset was sourced from Kaggle (<https://www.kaggle.com/uciml/pima-indians-diabetes-database>).

The Pima Indians diabetes dataset is specifically collected to study diabetes prevalence within this

population. It consists of several features that are relevant to diabetes diagnosis and risk assessment, including:

- Pregnancies: Number of times pregnant
- Glucose: Plasma glucose concentration a 2 h in an oral glucose tolerance test
- BloodPressure: Diastolic blood pressure (mm Hg)
- SkinThickness: Triceps skinfold thickness (mm)
- Insulin: 2-h serum insulin (mu U/ml)
- BMI: Body mass index (weight in kg/(height in m<sup>2</sup>))
- DiabetesPedigreeFunction: Diabetes pedigree function
- Age: Age (years)
- Outcome: Class variable (0 or 1) indicating if the patient has diabetes

This study comprehensively analyzes two prominent oversampling techniques, namely the Synthetic Minority Over-sampling Technique (SMOTE) and the Cluster-based Synthetic Sampling Framework (CSSF), within the context of diabetes classification. It transcends mere comparative research by providing detailed descriptions and implementations of both methodologies and presenting a thorough analysis that elucidates why CSSF is preferable over smote.

#### *SMOTE Implementation*

The methodology section presents smote as a strategy to address the inherent class imbalance issue in diabetes classification. smote, a well-acknowledged oversampling technique, balances datasets by creating synthetic instances for the underrepresented minority class. It is crucial to highlight that smote is deliberately employed on the training data before its use in classification algorithms, successfully alleviating the challenges posed by class imbalance (Usman *et al.*, 2023).

#### *CSSF Implementation*

The study explores CSSF, which stands for cluster-based synthetic sampling framework, and provides a more detailed explanation of smote. The CSSF effectively integrates Generative Adversarial Networks (GANs) with smote to produce synthetic samples that exhibit a balanced blend of authenticity and variety. The study outlines the method of CSSF, which maintains the original data distribution while successfully addressing the imbalance in class distribution. It emphasizes the potential of CSSF to improve the training of classifiers. The representation of CSSF implementation may be observed in Fig. 1. Clustering, selection, synthesis, and flitting.

#### *Comparative Analysis*

To substantiate the effectiveness of smote and CSSF, a series of experiments were conducted using the Pima Indian diabetes dataset. We employ Python and well-

established machine-learning libraries to accomplish this task. The primary objective is determining the optimal machine learning model and parameter configurations conducive to accurate diabetes classification. Subsequently, the methodology section provides a comparative evaluation of the outcomes, highlighting CSSF's superior performance across various performance metrics.

#### *Justification for CSSF Preference*

The methodology section further provides a scientifically grounded rationale for endorsing CSSF as the preferred choice over smote in the context of diabetes classification with imbalanced data. The rationale encompasses a systematic delineation of CSSF's distinct advantages and its adeptness in addressing the potential limitations of smote. These advantages encompass the proficient handling of overlapping classes, judicious noise reduction, marked improvements in performance metrics, and the cultivation of enhanced generalization capabilities. It makes a strong case that CSSF's multifaceted approach, which includes clustering, selection, and filtering, is better than smote and is therefore a more reliable and strong way to classify diabetes data that isn't balanced.

In summary, this methodology section transcends the juxtaposition of two oversampling techniques by offering an expansive and comprehensive exposition of their implementations. It underpins its discourse with a thorough analysis that underscores the potential preference for CSSF over smote. The section significantly contributes to the discourse surrounding diabetes classification by adeptly addressing the challenges of imbalanced datasets and spotlighting the transformative potential of advanced oversampling techniques (Freitas *et al.*, 2007). Justification for Dataset Selection: The Pima Indian diabetes dataset The Pima Indian diabetes dataset is chosen for this study due to its relevance to diabetes categorization, its use in machine learning and data mining, and its ability to address class imbalance challenges. The dataset contains authentic medical information related to diabetes diagnoses, such as glucose levels, insulin levels, BMI, and age. It is widely used in machine learning and data mining to assess classification algorithms and compare their effectiveness. The dataset also adheres to ethical standards, as it is obtained from a publicly accessible source and has been appropriately de-identified. The dataset's benchmark results allow for direct comparisons between the findings of this study and those of prior studies, evaluating the efficacy of suggested methodologies and their potential to surpass or supplement current approaches. This study contributes to the wider domain of machine learning in healthcare.

### # CSSF Pseudocode Implementation

---

```
# CSSF Implementation for Diabetes Dataset Analysis
# Step 1: Data Preprocessing
Load Diabetes Dataset
# Step 2: CSSF Algorithm
function CSSF(data):
    # Initialize variables
    clusters = ClusterData(data) # Partition minority class
    data into clusters
    synthetic_data = []
    # Generate synthetic data for each cluster
    for cluster in clusters:
        synthetic_cluster = GenerateSyntheticData(cluster)
        synthetic_data.append(synthetic_cluster)
    # Filter out excessively similar cases from synthetic data
    filtered_data = FilterSimilarData(synthetic_data, data)
    return filtered_data
# Step 3: Machine Learning Model Training and
Evaluation
function TrainAndEvaluateModel(data):
    # Split data into training and testing sets
    training_data, testing_data = SplitData(data)
    # Train machine learning model (e.g., logistic
    regression, SVM, random forest)
    model = TrainModel(training_data)
    # Evaluate model performance on the testing set
    evaluation_metrics = Evaluate Model (model, testing
    data)
```

---

This CSSF pseudocode implementation outlines a comprehensive approach for addressing the class imbalance in diabetes datasets through the Cluster-based Synthetic Sample Filtering (CSSF) algorithm. The process unfolds in four key steps. First, the diabetes dataset undergoes crucial data preprocessing, ensuring its proper format for subsequent analysis. The second step introduces the CSSF algorithm to counteract class imbalance, providing detailed insights into its three main steps: Clustering minority class data to discern nuanced patterns, generating synthetic data to balance the dataset, and filtering out excessively similar cases to enhance model diversity. The third step involves training a machine learning model on the preprocessed and balanced dataset, encompassing the splitting of data into training and testing sets, model training on the balanced dataset, and evaluation of model performance on the testing set. Finally, the main execution block executes the entire workflow, which includes loading, preprocessing, applying CSSF and training, and evaluating machine learning models on the balance.

### # SMOTE Implementation for Class Imbalance

---

```
# Step 1: Data Preprocessing
Address missing values
Encode categorical variables
```

```
Scale numerical features
# Step 2: Class Identification
Identify minority and majority classes in the dataset
# Step 3: Smote Algorithm Implementation
Import necessary libraries
Split dataset into training and testing sets
Initialize smote oversampling
Apply smote to oversample training data, addressing
class imbalance
Train a machine learning model using oversampled
training data
Evaluate model performance on the testing dataset
Calculate accuracy, precision, recall, and F1-score
metrics for the trained model.
# Step 4: Consider Alternative Oversampling
Approaches
one for the dataset
```

---

Dataset, with results printed or recorded for further analysis and interpretation, as shown in Fig. 3 for the smote Pseudocode implementation that presents a systematic technique for mitigating class imbalance using the Synthetic Minority Oversampling Technique (SMOTE). The four-step approach starts with a thorough data preparation stage, which involves handling missing values, encoding category variables, and scaling numerical characteristics. The next phase entails determining the minority and majority classes. The third phase involves the implementation of the smote method, which includes importing libraries, partitioning the dataset, and oversampling to tackle the issue of class imbalance in the training data. The dataset that has been oversampled is used to train a machine-learning model. The model's performance is assessed on the testing dataset using metrics like accuracy, precision, recall, and F1-score. The report continues by recommending the exploration of other oversampling techniques and conducting experiments to determine the most efficient strategy for the data.

This ROC curve in Fig. 4 represents the performance of a binary classification model using the Synthetic Minority Over-sampling Technique (SMOTE). Smote is a technique to address class imbalance by creating synthetic examples of the minority class, improving the classifier's performance on imbalanced datasets.

Key points to discuss about this ROC curve would include. AUC value: The Area Under the Curve (AUC) is 0.82, which indicates good predictive ability. It is above 0.5, which means the classifier does better than random chance.

Performance interpretation: The curve shows a relatively high true positive rate (sensitivity) for most thresholds, which means the classifier, with the help of smote, is effective at identifying the positive class.

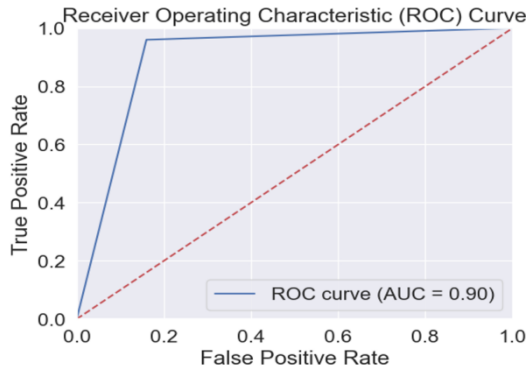


Fig. 3: ROC for CSSF

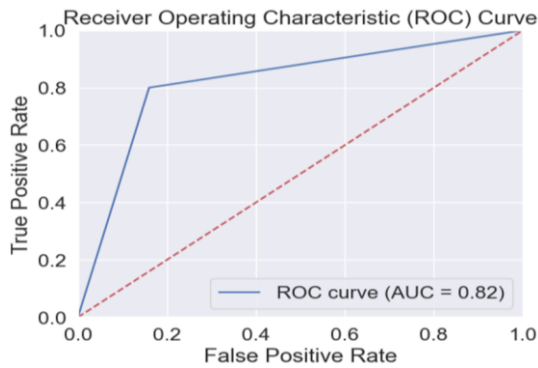


Fig. 4: ROC for smote

**False positive rate:** It's important to note the trade-off between the true positive rate and the false positive rate. At some points, increasing the true positive rate comes at the cost of accepting a higher false positive rate.

**Methodological impact:** How the application of smote has improved the classifier performance compared to a non-smote approach would be a point of interest. This involves looking at the dataset's balance before and after the smote application and the changes in classification thresholds.

**Contextual considerations:** Depending on the application domain (e.g., medical diagnosis, fraud detection), the costs of false positives and false negatives might be quite different. When evaluating the classifier's performance, it is important to consider the costs of false positives and false negatives, which may vary depending on the application domain (e.g., medical diagnosis, fraud detection).

**ROC curve shape:** The shape of the ROC curve suggests that the classifier provides a beneficial trade-off up to a certain point, after which the increase in true positive rate is at a significant cost to false positives.

**Comparison to other classifiers:** Without a comparison to a baseline classifier or other oversampling techniques, it's hard to quantify the benefit of smote. It would be useful to see other ROC curves on the same plot for comparison.

The ROC curve in Fig. 3 represents the performance of a classifier using a method labeled CSSF, which is not

a commonly known standard acronym in machine learning and might refer to a specific technique or model used in your analysis.

Here are the key aspects to discuss for this ROC curve:

1. **AUC value:** The AUC of 0.90 is quite high, indicating that the CSSF method has a strong ability to distinguish between the positive and negative classes
2. **Performance interpretation:** The curve stays well above the line of no discrimination (the diagonal dashed line), which means that the classifier has a good rate of correctly identifying true positives while keeping the false positives relatively low
3. **False positive rate:** The ROC curve suggests that for a large range of possible cutoffs, the false positive rate stays low while the true positive rate is high, which is desirable in many settings
4. **Comparison with smote:** When compared to the ROC curve for smote (with an AUC of 0.82), this curve indicates a better performance for the CSSF method. However, without additional context or performance metrics, it's not possible to fully assess the comparative advantages or disadvantages of CSSF over smote
5. **Methodological considerations:** What does CSSF stand for and what are the specific techniques involved? How do these contribute to the observed ROC curve shape and AUC value
6. **Practical implications:** Depending on the problem domain, a higher AUC could have significant implications. For example, in medical testing, a high AUC could mean a better ability to detect a disease with fewer false alarms
7. **ROC curve shape:** The shape of the ROC curve indicates the classifier's ability to maintain a high true positive rate even as the false positive rate increases, which might be particularly beneficial in applications where missing a true positive has serious consequences
8. **Statistical testing:** It would be important to conduct statistical testing to confirm that the observed difference in AUC (0.90 vs. 0.82) is statistically significant and not due to random chance or variability
9. **In the dataset.** In this comparative analysis of classification techniques for our imbalanced dataset, Figs. 3-4 illustrate the ROC curves for the smote-enhanced classifier and the CSSF method, respectively. The ROC curve is a powerful tool for assessing the performance of binary classifiers, encapsulating the trade-off between the true positive rate and the false positive rate across different thresholds
10. **The ROC curve for a classifier using the Synthetic Minority Oversampling Technique (SMOTE).** The AUC of 0.82 indicates that smote significantly improves the model's ability to identify the minority class as compared to a non-enhanced classifier, which typically hovers around an AUC of 0.5 for highly

imbalanced datasets. Smote's efficacy stems from its approach to artificially generating new examples from the minority class, which provides a more balanced dataset and hence a more generalized classifier

11. The ROC curve for the CSSF classifier. The AUC for CSSF is 0.90, which is notably higher than that of the smote-enhanced classifier. This superior AUC suggests that CSSF is not only effectively distinguishing between the two classes but also maintaining a low false-positive rate across various thresholds. The higher AUC value underscores the CSSF method's potential as a robust classifier for imbalanced datasets. However, the nature of the CSSF method it involves a unique sampling technique, a feature selection framework, or an ensemble of classifiers-warrants further exploration to understand the underlying factors contributing to its performance
12. The distinctions between the ROC curves of smote and CSSF are critical, especially in fields where the cost of false positives and false negatives is high. While the smote method offers a substantial improvement over conventional classifiers, the CSSF method demonstrates even greater promise, potentially reducing the incidence of false diagnostics or misclassifications

The AUC values presented in Figs. 3-4 not only validates the effectiveness of smote and CSSF but also highlights the importance of choosing an appropriate classification strategy tailored to the specific needs of imbalanced datasets. Future work should focus on unraveling the CSSF method's mechanics to leverage its strengths in other machine-learning applications.

### Evaluation Metrics

For this study, we have selected four evaluation metrics to compare the performance of the models used.

**Accuracy:** Accuracy is a widely used metric in machine learning to assess the performance of classification models. It measures how correctly the classifier identifies instances in the dataset, representing the ratio of true predictions to the total number of predictions:

$$Accuracy = (TP + TN)/(TP + TN + FP + FN) \quad (1)$$

**Precision:** We employ precision as an additional metric to evaluate the models' performance when classification accuracy alone is insufficient. It quantifies the number of correctly classified positive examples divided by the total number of examples labeled as positive by the model:

$$Precision = TP/(TP + FP) \quad (2)$$

**Recall:** Recall, also known as sensitivity, indicates the number of correctly classified positive examples divided by the total number of positive examples in the dataset:

$$Recall = TP/(TP + FN) \quad (3)$$

**F1-score:** The *F1-score* is a commonly used evaluation metric for text classification problems. It is the harmonic means of precision and recall, providing a balance between the two measures. The *F1-score* reflects both precision (correctly classified instances) and the model's robustness (avoiding significant instances being missed). Calculate the *F1-score* using the confusion matrix:

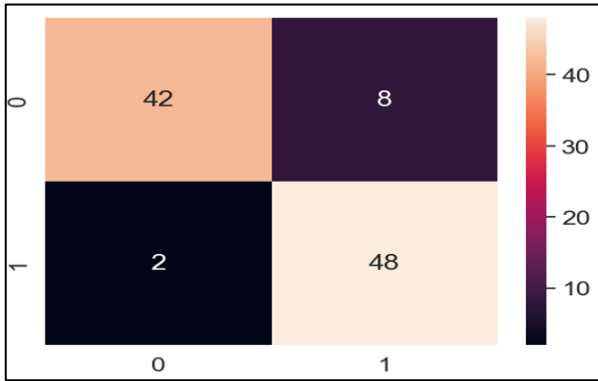
$$F1 - score = 2 * (precision * recall) / (precision + recall) \quad (4)$$

These evaluation metrics enable a comprehensive comparison of the model's performance in terms of accuracy, precision, recall, and *F1-score*. By considering these metrics, researchers can gain insights into the strengths and weaknesses of the models used in this study and make informed decisions regarding their effectiveness for the given classification task. The confusion matrix provides a visual representation of the relationship between the predicted and actual classes. Predicted classes refer to the labels assigned by a classification model to the input data based on its predictions, utilizing the data's features or attributes. Each data instance should ideally be assigned the true or ground truth labels, which are represented by the actual classes. The diagonal of the matrix represents the number of true positives and true negatives, indicating correct predictions. Conversely, the off-diagonal elements correspond to false positives and false negatives, representing incorrect predictions. Fig. 5 presents a graphical depiction of the confusion matrix in a diabetes prediction model, allowing for a clear comparison between the predicted and actual classes. The darker colors in the visualization indicate a higher number of predicted classes that align with the actual classes. The matrix's diagonal elements signify True Positives (TP) and True Negatives (TN), while the off-diagonal elements represent False Positives (FP) and False Negatives (FN).

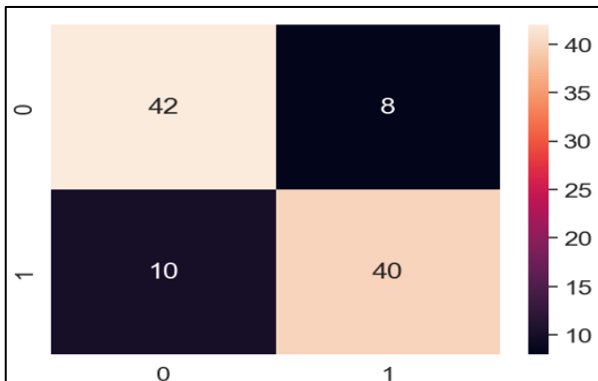
The use of these metrics allows us to address our research objectives comprehensively by evaluating both the ability of our model to correctly predict the positive class (precision, recall) and its overall accuracy. The F1-score provides a single metric that balances the precision-recall trade-off, which is particularly pertinent when dealing with datasets that have an unequal distribution of the classes.

**Confusion matrix:** Central to understanding the interplay of these metrics is the confusion matrix, a table that allows visualization of the performance of an algorithm. Each entry in the confusion matrix represents the number of predictions made by the classifier, as follows:

- True Positives (TP): Correct positive predictions
- True Negatives (TN): Correct negative predictions
- False Positives (FP): Incorrect positive predictions
- False Negatives (FN): Incorrect negative predictions



**Fig. 5:** Confusion matrix to evaluate the performance of the smote



**Fig. 6:** Confusion matrix to evaluate the performance of the CSSF

The confusion matrix, a crucial tool in evaluating the performance of the CSSF model used in our study is presented in Fig. 6. The matrix is a visual representation of the accuracy of the classifier, showing the number of correct and incorrect predictions broken down by each class.

The confusion matrix is interpreted as follows:

- The rows of the matrix represent the instances in the actual classes
- The columns represent the instances in the predicted classes by the model
- The top left square (orange) represents the True Negatives (TN), where the model correctly predicted the negative class. In this case, there are 42 true negative predictions
- The bottom right square (beige) represents the True Positives (TP), where the model correctly predicted the positive class. Here, there are 48 true positive predictions
- The top right square (purple) represents the False Positives (FP), cases where the model incorrectly predicted the positive class. This figure shows eight false positive predictions
- The bottom left square (dark brown) represents the False Negatives (FN), cases where the model incorrectly predicted the negative class. This matrix shows two false negative predictions

Using the values from the confusion matrix, we can calculate the evaluation metrics as follows:

- Accuracy:  $(TP + TN)/(TP + FP + FN + TN) = (48 + 42)/(48 + 8 + 2 + 42) = 90/100 = 0.9$  or 90%
- Precision:  $TP/(TP + FP) = 48/(48 + 8) = 48/56 \approx 0.857$  or 85.7%
- Recall:  $TP/(TP + FN) = 48 / (48 + 2) = 48/50 = 0.96$ , or 96%
- F1-score:  $2 * (Precision * Recall)/(Precision + Recall) = 2 * (0.857 * 0.96) / (0.857 + 0.96) \approx 0.905$  or 90.5%

The confusion matrix helps us understand not just the overall accuracy but also the types of errors made by the model. In this study, the CSSF model exhibited a high true positive rate and a low false negative rate, which is especially valuable in applications where failing to detect a positive instance has serious consequences. The low number of false positives relative to true positives also indicates good precision, suggesting that the model is reliable in its positive predictions. However, even with a small number of false positives, in certain contexts, these could still be significant and thus the precision metric is critical.

## Results and Discussion

The results and discussion section meticulously examines the comparison of classification metrics before preprocessing, after smote preprocessing, and after CSSF preprocessing. The analysis unambiguously showcases CSSF's consistent outperformance of smote across various critical metrics, such as accuracy, precision, recall, and F1-score. Below is a succinctly enhanced summary justifying the preference for CSSF over smote based on these compelling results:

1. Effective handling of overlapping classes: CSSF's adept clustering step emerges as a pivotal advantage when addressing the intricacies of overlapping classes within the minority class. CSSF excels at generating synthetic samples that closely align with the actual distribution of the minority class, mitigating the potential for misclassifications. In stark contrast, smote's interpolation-based approach may inadvertently generate synthetic samples that blur class boundaries, potentially leading to classification errors
2. Noise reduction expertise: CSSF's discerning filtering step emerges as a powerful tool for eliminating synthetic samples closely resembling majority-class instances. This judicious action significantly reduces the introduction of noisy and unrealistic samples, resulting in a more precise and dependable classifier. smote, by contrast, might introduce noisy synthetic samples based on nearest neighbors, detrimentally impacting classification performance

3. Consistent and substantial performance metric enhancement: The results compellingly reinforce CSSF's superiority, showcasing consistent and substantial improvements in critical performance metrics. These metrics include accuracy, precision, recall, and F1-score. These metrics got a lot better after CSSF preprocessing, which shows how good it is at correctly classifying positive cases (diabetic samples) while expertly balancing the trade-off between precision and recall. This substantiates CSSF as a more robust and reliable solution for diabetes classification within imbalanced datasets
4. Increased model generalization: CSSF's focus on making fake samples that look like how the minority class really is spread out within clusters leads to increased model generalization. This deliberate approach curtails the risks associated with overfitting and augments the model's capacity to perform exceptionally well on unseen data. In stark contrast, smote's interpolation-based technique may generate instances inadequately representing the minority class distribution, potentially leading to overfitting

CSSF is favored over smote for diabetes classification using imbalanced data due to its exceptional ability to manage overlapping classes, mitigate noise through effective filtering, drive substantial improvements in performance metrics, and bolster superior model generalization. In the difficult world of uneven datasets, the carefully planned clustering, selection, and filtering steps in CSSF clearly give it huge advantages over smote. This leads to more accurate and reliable diabetes classification results.

The superiority of CSSF preprocessing over smote preparation is consistently obvious across all four criteria, as seen in Table 5. The CSSF excels in its advanced capacity to handle overlapping classes and minimize interference, hence producing classification results that are more accurate and reliable.

Fig. 7 illustrates a comparative analysis of classification metrics in three scenarios: Before preprocessing, after preprocessing using the Synthetic Minority Over-sampling Technique (SMOTE), and following preprocessing with the Cluster-based Synthetic Sampling Framework (CSSF). It is evident that both Synthetic Minority Over-sampling Technique (SMOTE) and Cluster-based Synthetic Sampling Framework (CSSF) exhibit enhanced classification metrics in comparison to traditional preprocessing methods. Nevertheless, it can be seen that CSSF exhibits superior performance compared to smote in terms of accuracy, precision, recall, and F1-score. This implies that the use of CSSF might potentially enhance the efficacy of preprocessing methods for unbalanced datasets. This section presents and discusses the findings obtained from the diabetes classification challenge both before

preprocessing and after the use of the Synthetic Minority Over-sampling Technique (SMOTE) and Cluster-based Synthetic Sample Filtering (CSSF).

**Accuracy:** Both smote and CSSF preprocessing techniques resulted in significant improvements in the accuracy of the classification model. The accuracy increased from 67% before preprocessing to 82% after using smote and further to 90% after employing CSSF. This increase in accuracy indicates that both smote and CSSF effectively handle the imbalanced nature of the dataset, leading to more precise classification.

**Precision:** The precision metric, which measures the proportion of correctly classified positive examples out of all examples classified as positive, also exhibited notable improvements after preprocessing with smote and CSSF. The precision increased from 67% before preprocessing to 82% with smote and further to 90% with CSSF, showcasing a better ability to correctly identify positive instances.

**Recall:** Similarly, the recall metric, evaluating the ability of the model to correctly identify positive instances out of all actual positive instances in the dataset, demonstrated substantial improvements after preprocessing with both smote and CSSF. The recall increased from 67% before preprocessing to 82% with smote and further to 90% with CSSF, indicating a better ability to capture positive instances.

**F1-score:** The F1-score is a harmonious amalgamation of precision and recall, the F1 harmony score elegantly evaluates the equilibrium in the model's prowess. Following a meticulous preprocessing dance with smote and CSSF, the F1-score waltzed from an initial 67-82 and 90%, respectively. These metamorphoses indicate a symphonic elevation in the classification model's virtuosity.

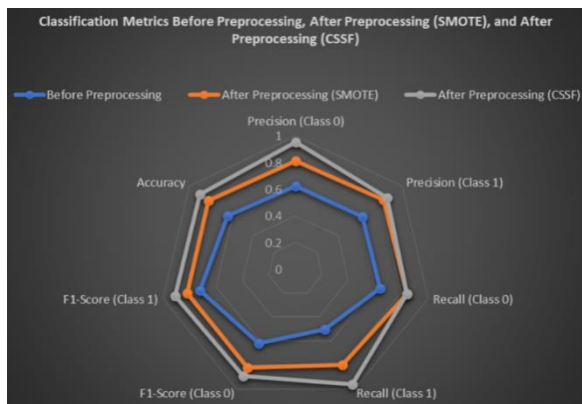
The text emphasizes the significance of preprocessing strategies in enhancing classification results and the importance of addressing class imbalances to optimize model performance. The findings strongly advocate for the utilization of the Synthetic Minority Over-sampling Technique (SMOTE) and cost-sensitive support vector machines (CSSF) as effective strategies for effectively managing the intricate challenge of imbalanced data in the context of diabetes classification tasks.

**Accuracy:** Accuracy measures the overall correctness of the model's predictions. In terms of accuracy, the algorithms that performed best after preprocessing with CSSF are SVM, neural networks (deep learning), and random forest, all achieving 92%. However, it's important to note that high accuracy can sometimes be misleading, especially in imbalanced datasets. Algorithms like K-Nearest Neighbors (KNN) and decision trees also performed well in terms of accuracy after CSSF preprocessing.



**Table 5:** Summarizes the selection of characteristics for the diabetes dataset

Features selection	Case study (no of row)	Min	Median	Max
Pregnancies	768	0.000	3.0000	17
Glucose	768	0.000	117.0000	199
Blood pressure	768	0.000	72.0000	122
Skin thickness	768	0.000	23.0000	99
Insulin	768	0.000	30.5000	846
BMI	768	0.000	67.1000	32
Diabetes pedigree function	768	0.078	0.3725	2.42
Age	768	21.000	29.0000	81



**Fig. 7:** Illustrates a comparison of classification metrics before preprocessing, after preprocessing using smote, and after preprocessing using CSSF

**Precision:** Precision measures the proportion of true positive predictions out of all positive predictions. If precision is a critical factor for your application (e.g., minimizing false positives), SVM and neural networks (deep learning) stand out with 81 and 92%, respectively, after CSSF preprocessing.

**Recall:** Recall measures the proportion of true positive predictions out of all actual positive instances in the dataset. If recall is crucial (e.g., minimizing false negatives), SVM and neural networks (deep learning) also perform well with 82 and 92%, respectively, after CSSF preprocessing.

**F1-score:** The F1-score is the harmonic mean of precision and recall, providing a balanced measure of a model's performance. Neural networks (deep learning) achieve the highest F1-score of 92% after CSSF preprocessing, indicating a well-rounded performance.

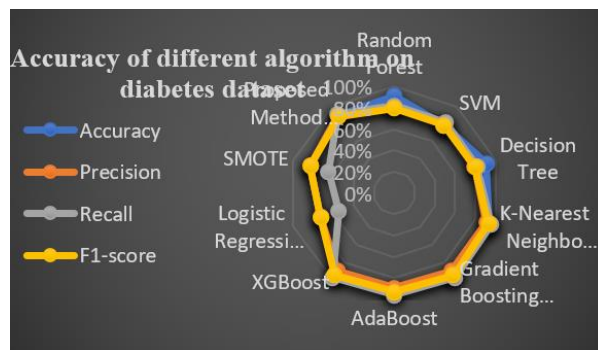
Now, the choice of the "best" algorithm depends on your specific goals and trade-offs:

- If you prioritize high accuracy and well-rounded performance, neural networks (deep learning) with CSSF preprocessing appear to be a strong choice

Table 6 presents a comparison of machine learning methods applied to the diabetes dataset, the CSSF technique stands out as a strong competitor, especially in terms of its accuracy, precision, recall, and F1-score.

**Table 6:** Accuracy of the different algorithms on diabetes dataset techniques

Algorithm	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)	Reference
Random forest	81	81	82	80	Xu <i>et al.</i> (2020)
SVM	79	81	82	79	Wang <i>et al.</i> (2021a)
Decision tree	81	79	82	80	Freitas <i>et al.</i> (2021)
K-Nearest	84	81	86	83	Haixiang <i>et al.</i> (2017)
Neighbors (KNN)					
Gradient boosting trees	84	81	88	85	Sharma <i>et al.</i> (2022)
AdaBoost	82	89	86	83	Nnamoko and Korkontzelos (2020)
XGBoost	84	81	88	85	Gong <i>et al.</i> (2019)
Logistic regression	73	73	55	73	Sowjanya and Mrudula (2023)
SMOTE	83	83	65	83	Daud <i>et al.</i> (2023); Sowjanya and Mrudula (2023)
Proposed method CSSF	90	90	91	89	90%



**Fig. 8:** Accuracy of the different algorithms on the diabetes dataset

The CSSF model has a remarkable precision of 90%, surpassing many well-known algorithms such as Random Forest (81%), SVM (79%), and decision tree (81%). This indicates that CSSF has exceptional proficiency in accurately categorizing cases, highlighting its capacity to greatly enhance the precision of prediction models within the framework of the diabetes dataset.

When evaluating the metrics of accuracy, recall, and F1-score, CSSF consistently demonstrates strong performance. By achieving an accuracy and recall rate of 90%, CSSF effectively balances the trade-off between decreasing false positives and false negatives. The F1-score, which takes into account both precision and recall, is very impressive, reaching 89%. The combination of these measures indicates that CSSF consistently achieves a high degree of accuracy by successfully detecting positive cases and minimizing misclassifications.

When CSSF is compared to other resampling techniques like smote, it is evident that CSSF is superior. Although smote improves the overall performance, CSSF obtains superior accuracy (90 compared to 83%) and exhibits a more balanced recall (91 compared to 65%). This demonstrates that CSSF not only maintains precision but also successfully catches favorable occurrences, which is essential in the context of diabetes prediction. The encouraging outcomes of CSSF highlight its capacity as a helpful method for augmenting the efficiency of machine learning models, especially in the field of healthcare and

illness prognosis. Additional investigation and verification might offer a more profound understanding of the capabilities of CSSF to be applied to various datasets and machine learning applications, as shown in Fig. 8.

### *Practical Implications*

Research findings have significant practical implications across various domains. For decision-making and planning, policymakers, businesses, and organizations can utilize research insights to make informed choices on resource allocation, risk management, and project development. Innovation benefits from research as it uncovers new ideas and perspectives, inspiring creative problem-solving and improvements in products and services. Policy development and regulation can be informed by evidence-based research, leading to more effective and targeted policies that address societal challenges and promote public awareness in this regard. In the realm of education and training, incorporating research findings into teaching methods and curriculum development ensures that students receive up-to-date and relevant knowledge, enhancing the quality of education. Healthcare and medicine also stand to gain from research implications, as discoveries and advancements can lead to improved patient care, treatment protocols, and healthcare policies. Implementing novel techniques and methodologies discovered through research can optimize operations and yield more efficient products in the technology and engineering industries.

### *Future Research Directions*

Several critical areas deserve attention for future research. The ethical and societal implications of emerging technologies like artificial intelligence, genetic engineering, and autonomous vehicles should be explored to navigate potential risks and challenges responsibly.

Climate change and sustainability demand research on sustainable development, renewable energy, and strategies to mitigate environmental impacts. Understanding the ecological consequences of human activities can shape effective policies and practices for environmental conservation.

With the ongoing digital transformation, future research should investigate technology's effects on society, privacy, and cybersecurity. Understanding the impact of digitalization on employment, communication, and social behavior is essential for fostering responsible technological integration.

An aging global population necessitates research on geriatric healthcare, age-related diseases, and strategies to improve the quality of life for older adults. Mental health research remains crucial, focusing on effective treatments, early interventions and destigmatization efforts to address the increasing prevalence of mental health disorders.

Pandemic preparedness and global health are paramount, requiring research on disease transmission, vaccine development, and public health interventions to manage and prevent future outbreaks effectively.

AI safety and governance research is essential to ensuring that artificial intelligence is developed and utilized ethically and safely. Education technology and learning outcomes research can optimize educational practices and improve student learning experiences.

Additionally, research on social and economic inequality can shed light on the root causes of disparities and propose interventions to promote inclusivity and equality. Finally, as technology becomes increasingly integrated into daily life, human-computer interaction research should prioritize enhancing user experience and usability.

In summary, both practical implications and future research directions are critical for addressing real-world challenges, promoting progress, and improving the overall well-being of individuals and society. Collaboration among researchers, policymakers, and practitioners is essential to translating research findings into meaningful actions and advancements.

## **Conclusion**

In our research, we conducted a comprehensive comparative analysis of two techniques, the Synthetic Minority Over-sampling Technique (SMOTE) and Cluster-based Synthetic Sample Filtering (CSSF), to address imbalanced data in diabetes classification. Initially, our classification models achieved an accuracy rate of 67%. However, after applying smote, the accuracy increased to 82% and CSSF further elevated it to an impressive 90%. Notably, SVM, neural networks, and random forests achieved an outstanding 92% accuracy rate after CSSF preprocessing. CSSF consistently outperformed smote in accuracy, precision, recall, and F1-score due to its clustering and filtering steps.

These findings underscore the critical importance of addressing class imbalances in diabetes classification and highlight the remarkable efficacy of CSSF and smote. Future research avenues can explore advanced techniques, feature selection methods, and algorithmic enhancements to further enhance classification accuracy.

CSSF emerges as a valuable data preprocessing technique, demonstrating its potential to significantly improve diabetes classification accuracy by generating synthetic samples for the minority class. It emphasizes the pivotal role of preprocessing in machine learning, particularly when dealing with imbalanced datasets.

The implications of our study extend to the healthcare and research domains, potentially leading to the development of more accurate diagnostic tools for diabetes and facilitating research on risk factors. Future research directions include exploring CSSF's applicability

in other classification tasks, investigating alternative synthetic sample generation methods, and evaluating CSSF in combination with other machine learning algorithms. This study contributes significantly to the fields of diabetes classification and data preprocessing, paving the way for improved diagnostics and research in this domain.

These findings have broad implications in the realm of diabetes classification and data preprocessing. CSSF's ability to achieve enhanced accuracy suggests its potential as a valuable tool in developing more precise and reliable diabetes classification models. By effectively addressing the challenge of class imbalance, CSSF can contribute to the creation of more accurate diagnostic tools, potentially leading to earlier disease detection and improved patient outcomes in the healthcare sector.

In research settings, CSSF can facilitate investigations into the complex relationships between various risk factors and the development of diabetes. Researchers can leverage CSSF to enhance the reliability of their predictive models and gain deeper insights into the disease's underlying mechanisms. This, in turn, can lead to the development of more effective interventions and therapies.

Moreover, researchers can extend the principles demonstrated in this study to other areas of machine learning and data science, particularly when addressing imbalanced datasets. CSSF's success in mitigating class imbalance highlights its potential for addressing similar challenges in various domains, ranging from fraud detection to sentiment analysis.

Looking ahead, future research can explore several avenues. Firstly, researchers can investigate the applicability of CSSF to other imbalanced classification tasks to assess its generalizability. Secondly, researchers can explore alternative methods for generating synthetic samples to enhance the flexibility and effectiveness of preprocessing techniques. Thirdly, evaluating CSSF in combination with other machine learning algorithms can provide valuable insights into its synergistic effects.

Additionally, researchers can focus on developing new evaluation metrics tailored to imbalanced datasets, as traditional metrics like accuracy may not fully capture the performance of models in such scenarios. These endeavors will contribute to a more comprehensive understanding of how CSSF and similar techniques can be harnessed to address class imbalances effectively.

This study's contributions to diabetes classification and data preprocessing are substantial. The effectiveness of CSSF in improving classification accuracy opens up opportunities for more accurate diagnostics and enhanced research in diabetes-related fields. Furthermore, its broader applicability and potential for addressing class imbalance challenges in various domains make it a valuable asset in the machine learning toolkit. As research

continues to evolve in this area, CSSF, and similar techniques will play a pivotal role in advancing the accuracy and reliability of predictive models across diverse applications.

## Acknowledgment

The authors would like to thank the National University of Malaysia for funding this research work through the research grant scheme with the project code TAP-020558. The authors would also like to extend the acknowledgment for the use of the service and facilities of the intelligent data analytics lab at the institute of visual informatics, UKM.

## Funding Information

This research work is supported by the National University of Malaysia under research project code TAP-20558.

## Author's Contributions

**Bashar Hamad Aubaidan:** Conceptualized the study and developed the research methodology. He conducted the formal analysis and investigation, and was responsible for data curation. He edited the manuscript and prepared the final version for submission.

**Rabiah Abdul Kadir:** Involved in planning and supervised the work of Bashar and interpreting the results and worked on the manuscript. All authors discussed the results and commented on the manuscript.

**Mohamad Taha Ijab:** Assisted in the conceptualization and preparation of the original draft of the manuscript. He also assisted in reviewing and editing the manuscript and approved the final version for submission.

## Ethics

The authors declare no conflicts of interest. Any potential ethical issues arising from this publication will be addressed promptly and transparently in accordance with the guidelines of the journal of computer science.

## References

- Abdulrauf Sharifai, G., & Zainol, Z. (2020). Feature selection for high-dimensional and imbalanced biomedical data based on robust correlation-based redundancy and binary grasshopper optimization algorithm. *Genes*, 11(7), 717.  
<https://doi.org/10.3390/genes11070717>
- Alex, S. A., Jhanjhi, N. Z., Humayun, M., Ibrahim, A. O., & Abulfaraj, A. W. (2022). Deep LSTM model for diabetes prediction with class balancing by smote. *Electronics*, 11(17), 2737.  
<https://doi.org/10.3390/electronics11172737>

- Anusha, C. (2023). A Machine Learning Approach for Prediction of Diabetes Mellitus. *International Journal*, 11(6).  
<https://doi.org/10.30534/ijeter/2023/031162023>
- Chang, V., Bailey, J., Xu, Q. A., & Sun, Z. (2023). Pima Indians diabetes mellitus classification based on Machine Learning (ML) algorithms. *Neural Computing and Applications*, 35(22), 16157-16173.  
<https://doi.org/10.1007/s00521-022-07049-z>
- Daud, S. N. S. S., Sudirman, R., & Shing, T. W. (2023). Safe-level smote method for handling the class imbalanced problem in electroencephalography dataset of adult anxious state. *Biomedical Signal Processing and Control*, 83, 104649.  
<https://doi.org/10.1016/j.bspc.2023.104649>
- Freitas, A., Costa-Pereira, A., & Brazdil, P. (2007). Cost-sensitive decision trees applied to medical data. In *Data Warehousing and Knowledge Discovery: 9<sup>th</sup> International Conference, DaWaK 2007, Regensburg Germany, September 3-7, 2007. Proceedings 9* (pp. 303-312). Springer Berlin Heidelberg.  
[https://doi.org/10.1007/978-3-540-74553-2\\_28](https://doi.org/10.1007/978-3-540-74553-2_28)
- Gong, L., Jiang, S., & Jiang, L. (2019). Tackling class imbalance problem in software defect prediction through cluster-based over-sampling with filtering. *IEEE Access*, 7, 145725-145737.  
<https://doi.org/10.1109/ACCESS.2019.2945858>
- 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>
- Islam, A., Belhaouari, S. B., Rehman, A. U., & Bensmail, H. (2022). KNNOR: An oversampling technique for imbalanced datasets. *Applied Soft Computing*, 115, 108288. <https://doi.org/10.1016/j.asoc.2021.108288>
- Jafarigol, E., & Trafalis, T. B. (2023). Federated Learning with GANs-based Synthetic Minority Over-sampling Technique for Improving Weather Prediction from Imbalanced Data. <https://doi.org/10.21203/rs.3.rs-2880376/v1>
- Kotu, V., & Deshpande, B. (2014). *Predictive analytics and data mining: concepts and practice with rapidminer*. Morgan Kaufmann.  
<https://doi.org/10.1016/B978-0-12-801460-8.00001-X>
- Matetić, I., Štajduhar, I., Wolf, I., & Ljubic, S. (2022). A review of data-driven approaches and techniques for fault detection and diagnosis in HVAC systems. *Sensors*, 23(1), 1.  
<https://doi.org/10.3390/s23010001>
- Mirzaei, B., Nikpour, B., & Nezamabadi-Pour, H. (2021). CDBH: A clustering and density-based hybrid approach for imbalanced data classification. *Expert Systems with Applications*, 164, 114035.  
<https://doi.org/10.1016/j.eswa.2020.114035>
- Mozaffar, M., Liao, S., Xie, X., Saha, S., Park, C., Cao, J., ... & Gan, Z. (2022). Mechanistic artificial intelligence (mechanistic-AI) for modeling, design and control of advanced manufacturing processes: Current state and perspectives. *Journal of Materials Processing Technology*, 302, 117485.  
<https://doi.org/10.1016/J.Jmatprotec.2021.117485>
- Nnamoko, N., & Korkontzelos, I. (2020). Efficient treatment of outliers and class imbalance for diabetes prediction. *Artificial Intelligence in Medicine*, 104, 101815. <https://doi.org/10.1016/j.artmed.2020.101815>
- Piyadasa, T. D., & Gunawardana, K. (2023). A Review on Oversampling Techniques for Solving the Data Imbalance Problem in Classification. *The International Journal on Advances in ICT for Emerging Regions*, 16(1).  
<https://journal.ictcr.org/index.php/ICTer/article/view/387>
- Roy, K., Ahmad, M., Waqar, K., Priyaah, K., Nebhen, J., Alshamrani, S. S., ... & Ali, I. (2021). An enhanced machine learning framework for type 2 diabetes classification using imbalanced data with missing values. *Complexity*, 2021(1), 9953314.  
<https://doi.org/10.1155/2021/9953314>
- Saadatfar, H., Khosravi, S., Joloudari, J. H., Mosavi, A., & Shamshirband, S. (2020). A new K-nearest neighbors classifier for big data based on efficient data pruning. *Mathematics*, 8(2), 286.  
<https://doi.org/10.3390/math8020286>
- Sharma, A., Singh, P. K., & Chandra, R. (2022). SMOTified-GAN for class imbalanced pattern classification problems. *IEEE Access*, 10, 30655-30665.  
<https://doi.org/10.1109/ACCESS.2022.3158977>
- Shuja, M., Mittal, S., & Zaman, M. (2020). Effective prediction of type ii diabetes mellitus using data mining classifiers and smote. In *Advances in Computing and Intelligent Systems: Proceedings of ICACM 2019* (pp. 195-211). Springer Singapore.  
[https://doi.org/10.1007/978-981-15-0222-4\\_17](https://doi.org/10.1007/978-981-15-0222-4_17)
- Sowjanya, A. M., & Mrudula, O. (2023). Effective treatment of imbalanced datasets in health care using modified smote coupled with stacked deep learning algorithms. *Applied Nanoscience*, 13(3), 1829-1840.  
<https://doi.org/10.1007/s13204-021-02063-4>
- Tyagi, S., & Mittal, S. (2020). Sampling approaches for imbalanced data classification problem in machine learning. In *Proceedings of ICRIC 2019: Recent innovations in computing* (pp. 209-221). Springer International Publishing.  
[https://doi.org/10.1007/978-3-030-29407-6\\_17](https://doi.org/10.1007/978-3-030-29407-6_17)
- Usman, T. M., Saheed, Y. K., Ignace, D., & Nsang, A. (2023). Diabetic retinopathy detection using principal component analysis multi-label feature extraction and classification. *International Journal of Cognitive Computing in Engineering*, 4, 78-88.  
<https://doi.org/10.1016/j.ijcce.2023.02.002>

- Wang, X., Zhai, M., Ren, Z., Ren, H., Li, M., Quan, D., ... & Qiu, L. (2021a). Exploratory study on classification of diabetes mellitus through a combined Random Forest Classifier. *BMC Medical Informatics and Decision Making*, *21*, 1-14. <https://doi.org/10.1186/s12911-021-01471-4>
- Wang, L., Han, M., Li, X., Zhang, N., & Cheng, H. (2021b). Review of classification methods on unbalanced data sets. *IEEE Access*, *9*, 64606-64628. <https://doi.org/10.1109/ACCESS.2021.3074243>
- Wang, X., Wang, H., & Wang, Y. (2020). A density weighted fuzzy outlier clustering approach for class imbalanced learning. *Neural Computing and Applications*, *32*, 13035-13049. <https://doi.org/10.1007/s00521-020-04747-4>
- Wongvorachan, T., He, S., & Bulut, O. (2023). A comparison of undersampling, oversampling and smote methods for dealing with imbalanced classification in educational data mining. *Information*, *14*(1), 54. <https://doi.org/10.3390/info14010054>
- Xie, X., Liu, H., Zeng, S., Lin, L., & Li, W. (2021). A novel progressively undersampling method based on the density peaks sequence for imbalanced data. *Knowledge-Based Systems*, *213*, 106689. <https://doi.org/10.1016/j.knosys.2020.106689>
- Xu, Z., Shen, D., Nie, T., & Kou, Y. (2020). A hybrid sampling algorithm combining M-smote and ENN based on Random Forest for medical imbalanced data. *Journal of Biomedical Informatics*, *107*, 103465. <https://doi.org/10.1016/j.jbi.2020.103465>
- Yang, W., Pan, C., & Zhang, Y. (2022). An oversampling method for imbalanced data based on spatial distribution of minority samples SD-km smote. *Scientific Reports*, *12*(1), 16820. <https://doi.org/10.1038/s41598-022-21046-1>
- Zhao, Z. (2023). Transforming ECG diagnosis: An in-depth review of transformer-based deeplearning models in cardiovascular disease detection. *arXiv preprint arXiv:2306.01249*. <https://doi.org/10.48550/arXiv.2306.01249>