

SDAGS: SMOTE+FOREST DIFFUSION-BASED DATA AUGMENTATION AND GBT-BASED STACKING ENSEMBLE LEARNING FOR HOLISTIC AI-POWERED DIABETES MELLITUS PREDICTION

AN K. NGUYEN¹, BI Y. GONG¹, BINH T. TRINH², THU H. LE³, MINH T. NGUYEN⁴,
HANH P. DU^{4,*}

¹*The Dewey School, Ha Noi, Viet Nam*

²*VNU Hanoi University of Science, Ha Noi, Viet Nam*

³*Hanoi National University of Education, Ha Noi, Viet Nam*

⁴*VNU University of Engineering and Technology, Ha Noi, Viet Nam*



Abstract. In emerging nations like Vietnam, it's critical to monitor and anticipate diabetes, particularly for those with type 1 diabetes. This article proposes SDAGS, an AI-powered diabetes prediction technique. Our method is based on two ideas: (i) using the SFDM method to make the training data better by combining the oversampling of the Forest Diffusion Model with the SMOTE data balancing method; and (ii) making the GSDP model by stacking different boosting machine learning models together. We also propose a preliminary AI-powered blood glucose monitoring and recommendation system based on SDAGS to provide diabetic patients with all-encompassing assistance with blood glucose monitoring, dietary counseling, physical activity, and the proper use of medications. Our thorough experiments using the Pima Indians diabetes dataset and the 5-fold cross-validation method demonstrate that SDAGS outperforms the state-of-the-art methods. Its prediction performance significantly achieved a sensitivity of 98.3%, a specificity of 99.49%, an F1-score of 98.74%, an accuracy of 98.75%, and a precision of 98.00%.

Keywords. Gradient boosted tree; Stacking ensemble learning; Forest diffusion; Dataset augmentation; Diabetes prediction.

1. INTRODUCTION

Diabetes is a chronic disease that raises blood glucose [23]. Diabetes is a chronic condition caused by the pancreas not producing enough insulin or the body not using it. Hyperglycemia, or high blood glucose, is a typical complication of untreated diabetes and damages nerves and blood vessels [22]. This disorder can cause heart, renal, eyesight, and nerve problems [5]. Diabetes comes in three forms: type 1 diabetes (T1D), type 2 diabetes (T2D), and gestational diabetes [23]. The healthcare industry, particularly in healthcare, has shown great interest in

*Corresponding author.

E-mail addresses: helix6688@outlook.com (A.K. Nguyen); yunvi0808@gmail.com (B.Y. Gong); trinthanhbinh_t67@hus.edu.vn (B.T. Trinh); stu735105099@hnue.edu.vn (T.H. Le); 23020628@vnu.edu.vn (M.T. Nguyen); hanhdp@vnu.edu.vn (H.P. Du);

IoT technology and has widely adopted it. Chronic diseases, such as diabetes (both T1D and T2D), need health support systems. New findings in the field of IoT and sensors enable improved health monitoring efficiency as well as the collection of critical data to be able to process, analyse, and assist doctors in diagnosis, treatment, monitoring, and warning [26]. For diabetes, monitoring blood glucose and other biological and medical indicators is important in determining the severity of the disease as well as enabling early diagnosis and detection of T2D. Advances in biology and IoT have enabled the creation of non-invasive blood glucose monitoring devices. However, these non-invasive devices often give results that are not as accurate as blood tests. T1D is a chronic disorder in which the pancreas produces little or no insulin to manage blood glucose. In an autoimmune condition, the immune system wrongly assaults and destroys pancreatic beta cells that make insulin [23]. For T1D, patients need daily insulin, blood sugar monitoring, and a healthy lifestyle. Hence, blood glucose monitoring and management for diabetes patients, both T1D and T2D, are critical. However, it still faces the following challenges:

- Due to poor infrastructure, limited facilities, and a shortage of skilled healthcare professionals, developing nations' diabetes patients face challenges in diagnosis, treatment, and care. Expensive blood glucose monitoring devices, insulin, and other diabetes supplies are limited, resulting in poor treatment and health consequences [13]. Poverty and malnutrition can make it hard for T1D patients to eat well. Thus, developing nations need blood glucose monitoring, dietary counseling, physical activity, and medication management to fully advise and support them.
- Second, diabetes blood glucose tests require regular blood sampling, including pricking of the patient's finger. Diabetics, particularly T1D patients, may experience discomfort and pain during daily blood sugar monitoring [16, 19]. Frequent blood draws may increase infection risk, particularly if hygiene is neglected [21]. Thus, saliva, tear, and perspiration analysis can now measure blood glucose levels (BGL) non-invasively. These methods require greater accuracy compared to blood glucose analysis, via AI-powered diabetes prediction [33]. However, current AI-powered diabetes prediction models only achieve specificity and sensitivity of 96.56% and 96.01%, respectively [26]. Thus, our second challenge in this study is improving these AI models.

Therefore, the objective of our work in this paper is to enhance specificity and sensitivity in AI-powered diabetes prediction by applying new methods to improve the quality of the training dataset as well as selecting and proposing an optimal AI model. The following is the arrangement of the remaining sections: Section 2. outlines the basic concepts that support our research. Our proposed AI-powered strategy to enhance diabetes prediction performance is presented in Section 3.. This section also specifies a preliminary diabetes counseling and support framework. All the experimental results we conducted to address the novelty of our suggested method are provided in Section 4.. Lastly, Section 5. summarizes contribution highlights and perspectives for further research.

2. BACKGROUNDS AND RELATED WORKS

2.1. Fundamental concepts

In the human body, glucose provides energy for cellular metabolism [23]. Apart from the blood, glucose is also present in intracellular fluid, interstitial fluid (ISF), tears, saliva, and

urine [34]. Currently, blood glucose concentration is the primary basis for diagnosing diabetes. The World Health Organization (WHO) set the fasting blood glucose (FBG) standard for a normal individual at 3.9–6.1 mM and postprandial blood glucose at 7.8 mM or lower in 2009. Patients exhibiting typical symptoms of diabetes (polyuria, excessive thirst, unexplained weight loss) with random BGL ≥ 11.1 mM, $FBG \geq 7.0$ mM, or 2-hour postprandial ≥ 11.1 mM may be diagnosed with diabetes. Besides elevated blood sugar, low blood sugar levels can also be harmful to the human body. Clinically, hypoglycemia is defined as a condition where blood glucose concentration is lower than 3.9 mM and persists for more than 5 minutes [27]. Especially in elderly patients, the risk coefficient for hypoglycemia is even higher. The rate of nocturnal hypoglycemia is relatively high and challenging to promptly monitor using conventional blood glucose detection methods. Tight blood glucose control also has the potential to increase the risk of hypoglycemia [33]. Therefore, continuous glucose monitoring (CGM) in diabetic patients may hold more clinical value and be more suitable for real-world treatment requirements. Blood glucose monitoring can be simplistically categorized as invasive or non-invasive based on whether or not it induces injury to the human epidermis.

New studies allow non-invasive blood glucose monitoring without tissue damage. Non-invasive blood glucose detection methods include optical, microwave, and electrochemical. Numerous optical methods exist, such as NIRS, MIDI, polarimetry, Raman, fluorescence, OCT, and others [34]. They are divided into optical, microwave, and electrochemical methods. Each method for non-invasive blood glucose monitoring shows promise, but technical barriers remain in improving accuracy, sensitivity, and correlation with BGL across diverse population groups. Interdisciplinary efforts in materials, nanotechnology, and data analysis will improve non-invasive glucose monitoring technology for clinical and consumer needs. Diabetic tattoos measure ISF glucose levels with biosensors and special ink, changing colors or patterns. Continuous monitoring with temporary tattoos allows T1D patients to act quickly and avoid complications. Google’s wireless chips and sensors in smart lenses track tear glucose levels in real-time and send data to mobile devices. When not weeping, these lenses are inactive. Biosensors coated with glucose oxidase generate electrical currents when in contact with saliva glucose to measure glucose levels. Saliva analysis yields immediate results but requires repetition to detect trends. These non-invasive methods have promise, but they need more research and validation to determine accuracy and operational procedures. Further research is needed to improve reliability and establish clear usage guidelines [18].

2.2. Data augmentation

Currently, modern machine learning methods allow for data augmentation to enhance the quality of training datasets, especially for tabular data such as the Pima diabetes dataset (or PIMA for short) [26]. For imbalanced datasets, where the number of samples varies greatly across classes and labels, the Synthetic minority oversampling technique (SMOTE) [7] is one of the traditional methods that can be used to balance the number of samples across classes. Recently, a variant of GAN, the Wasserstein generative adversarial network (WGAN), was designed to enhance learning stability, ameliorate issues such as mode collapse, and generate meaningful learning curves that are practical for debugging and searching hyperparameters [36]. A WGAN minimizes the Wasserstein distance, which is an approximation of the Earth-Mover’s distance, as opposed to the Jensen-Shannon divergence. By changing the objective function in this way, there is less evidence of mode collapse and training is more stable than with

the original GANs [6]. Thus, WGAN demonstrates significant improvements in data quality through experiments with specialized tabular datasets [35,38].

Recently, some data augmentation techniques based on diffusion models have also been emphasized in research [30]. In [15], the authors describe a new way to create and fill in mixed-type tabular data that includes both continuous and categorical variables. They do this by using score-based diffusion and conditional flow matching, which is also known as FDM. Instead of using neural networks to train the scoring function or the vector field, they use XGBoost, a popular Gradient boosted tree (GBT) method. Empirical assessment across many benchmarks demonstrates that the FDM technique beats deep-learning generation methods in data creation tasks and stays competitive in data imputation.

2.3. Ensemble learning

New machine learning methods, exemplified by algorithms based on GBT, such as XGBoost, CatBoost, LightGBM, ExtraTrees, and Deep learning (DL) techniques, such as Convolutional neural networks (CNN), Graph neural networks (GNN), Long short-term memory (LSTM), and Transformer, have demonstrated high performance for classification tasks in general, as well as for diabetes mellitus datasets specifically [26]. However, despite the use of data augmentation methods such as SMOTE, WGAN, or FDM, the performance of a single GBT/DL model for diabetes prediction still requires improvement, as these models often excel only in specific problem domains. Consequently, employing ensemble learning to combine multiple individual classifiers allows us to achieve stable results across various data domains [38]. Currently, there are three main ensemble methods: Stacking, Hard Voting, and Soft Voting [29,37].

2.4. Related works

Almed et al. [3] propose buying non-invasive wearables to estimate BGL. The method uses AI models to link glycemic measurements to non-invasive WD traits. Their research with SVR, RF, MLP, and ANFIS AI models proves this. Alain et al. [14] propose an automated IoT-edge-AI-blockchain diabetes prediction system. This system uses medical sensors and devices to assess risk factors and predict diabetes. They used PIMA Indian, Sylhet, and MIMIC III diabetes datasets to test the system. According to Zhouyu et al. [12], an AI system can comprehend biomedical data about a user. Medtronic's Guardian Connect System alerts users an hour before hypoglycemia attacks. An XGBoost model that uses many predictive elements to estimate BGL can predict hyperglycemia in T2D patients, but not hypoglycemia. A study by Rajalakshmi et al. [26] aims to enhance diabetes prediction accuracy in real-time IoT datasets based on the PIMA dataset. They developed KST-BLSTM, which uses synthetic minority oversampling to fix the dataset's imbalance, the k-means clustering-based sailfish optimization algorithm to fill in missing values and choose the most important features, and the bidirectional long short-term memory (BLSTM) architecture to train the AI model. The PIMA dataset test yielded 98.26% accuracy.

T1D self-management requires hundreds of daily choices [32]. Zhang et al. [40] developed a food recognition system called snap-n-eat, which can identify food items from a photo of a user's full meal. This system could identify 15 food categories with 85% accuracy, but it must be scaled up to identify hundreds of other food items to be realistic and practical. A study by Robert et al. [1] examines a system that uses AI to suggest a personalized

diet Classifying meals using food image recognition, recommending meals using K-Nearest Neighbor, and answering diabetes questions with a chatbot. An intelligent system to help with diet, medication, and lifestyle for personalized diabetes management is the goal. G. Prabhakar et al. [24] propose a user-cloud-based ensemble framework for T2D prediction and diet planning. Although it predicts 87.41% accuracy, its diet recommendations' effects on users are not specified.

The Sweetch app by Sweetch Health, Ltd. effectively improved physical activity and weight reduction in prediabetic adults after 3 months of use [11]. Results include 1.6kg weight loss, 0.6 kg/m² BMI reduction, and 0.1% HbA1c reduction. T1D and T2D patients' physical activity recommendations haven't been studied. Chen et al. [9] recommend the right exercise mode for users in one-month, two-month, and three-month periods with 95.80%, 100.00%, and 95.00% accuracy. This study recommended 7,000–12,000 METs per minute (Metabolic Equivalent of Task) per month, while high-risk diabetics recommend 9,000–16,000 METs.

In their study, Aileen et al. [39] utilized sequential pattern mining to predict the next medication for diabetic patients. In three attempts, they predicted drug classes with 90.0% accuracy and generic drugs with 64.1% accuracy, demonstrating that sequential pattern mining may help prescribers choose future drug classes for diabetic patients. Shinji et al. [2] suggest using AI, including predictive analytics, to enhance pharmacotherapy recommendations for chronic conditions like T2D. New Treatment Pathway Graph-based Estimation outperformed previous ML algorithms for predicting treatment outcomes in this diabetes type.

3. PROPOSED METHOD

3.1. Approach direction

Based on the challenges described in Section 1, we propose a novel method called SDAGS, which stands for SMOTE+Forest diffusion-based data augmentation and GBTs-based stacking ensemble learning for holistic AI-powered diabetes mellitus prediction. Here, we aim to propose a comprehensive research solution for predicting diabetes based on two main ideas:

- Augment the training dataset quality by evaluating and selecting the best sample generating method from WGAN, TGAN, SMOTE, and Forest diffusion.
- Propose an AI-powered diabetes prediction method by evaluating, selecting, and optimizing modern machine learning algorithms.

Thus, in this work, we strongly aim to further enhance both the quality of the training dataset and diabetes detection performance. To improve the quality of the training dataset, we propose the SDFM algorithm, which is designed to combine the SMOTE technique for balancing between two diabetes/normal classes and generate more realistic samples based on the Forest Diffusion Model algorithm. To enhance diabetes detection performance, we propose the GSDP algorithm with a GBT-based stacking ensemble learning approach, employing, for instance, XGB, CBT, GBM, ET, and RF. The following subsections describe our SDAGS in detail.

3.2. Improving the training data quality

As outlined in Section 2, there are various methods for generating additional samples to enhance the quality of training data for AI models in classification and prediction tasks, particularly for tabular data, such as SMOTE, WGAN, FDM, etc. Based on our detailed

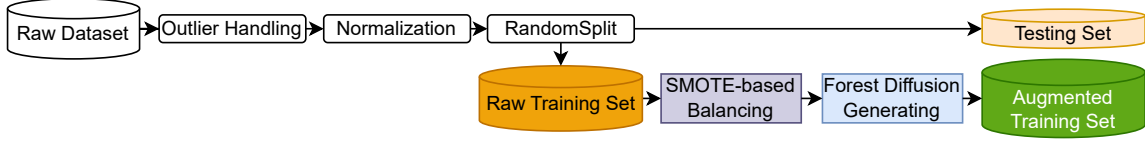


Figure 1: SFDM-based Training Data Augmentation

Algorithm 1 SFDM: Create the training & testing sets by fusing SMOTE and FDM

Input: DF - Raw dataset, represented by a list of feature vectors; r - the ratio between training and testing sets; default is 7:3; τ - maximum samples in a label.

Output: T - Training set; V - Testing set.

- 1: $F \leftarrow \text{OutlierHandling}(F)$ \triangleright Determine the outlier threshold and apply it to handle the detected outliers.
 - 2: $F \leftarrow \text{Normalize}(F)$ \triangleright Normalize all feature vectors by scaling each feature to a given range.
 - 3: $(RT, V) \leftarrow \text{SplitTrainTest}(F, r)$ \triangleright Split F randomly into the raw training set RT and testing set V with the ratio of r .
 - 4: $T_{SMOTE} \leftarrow \text{SMOTE_Sampling}(RT)$ \triangleright Perform over-sampling using SMOTE from RT .
 - 5: $T \leftarrow \text{FDM_Generate}(T_{SMOTE}, \tau)$ \triangleright Generate up to τ realistic samples by employing the forest diffusion model based on the balanced training set T_{SMOTE} .
 - 6: **return** (T, V) .
-

experiments with the PIMA dataset, SMOTE effectively addresses the issue of data imbalance between labels by generating additional samples for minority classes to achieve a balance between classes. However, for datasets with a small number of samples, such as PIMA (a total of 768 samples, with 500 normal and 268 diabetes samples), using SMOTE to generate extra samples for the majority class is considered to potentially lead to overfitting and reduced accuracy [7].

Regarding WGAN and FDM models, both are highly effective methods for generating high-quality samples. WGAN exhibits stability in training but has a slow training speed and requires a large amount of data for effective training. Similarly, FDM also generates high-quality data samples, providing enhanced stability during training and excelling in matching the distribution of real samples with precision. Our experiments with the PIMA dataset show that FDM has a faster data generation speed and higher quality compared to WGAN. These experimental results have allowed us to develop the data quality enhancement method by combining both SMOTE and FDM models, namely **SFDM**. Consequently, the SFDM is depicted in Fig. 1 and described formally in Alg. 1. In this algorithm, to augment the training set, we perform the following steps:

Step 1 - Dataset preparation: Eliminate noise and duplicated raw data on the dataset; replace missing values by using the median of each feature; determine a threshold for outlier detection and replace it with its value; carry out dataset normalizing to transform features by scaling each feature to a given range; and split it into the training and testing sets, preset by a ratio of 7:3, respectively. We use τ as a constant for determining the maximum samples of the label class in the training set. Hence, we use the testing set to evaluate our AI models.

Step 2 - Training set balancing: The raw training set is augmented by generating more samples in the minority class employing SMOTE.

Step 3 - Augmenting more realistic samples: The balanced training set will be augmented

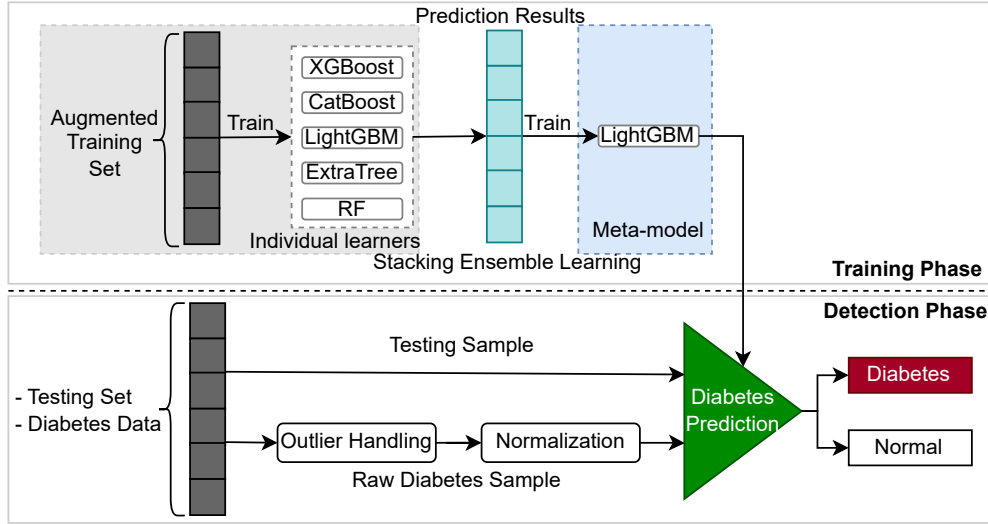


Figure 2: Architecture of GBT-based stacking ensemble learning-powered diabetes prediction.

by employing the forest diffusion model to generate realistic samples.

Step 4 - Results: Finally, we obtain a new training set, which contains the number of every label class that is the same as τ , and use this training set to train and the testing set in Step 1 to test our models. Thus, it helps us to obtain a better AI model.

3.3. Stacking ensemble-based diabetes prediction

As mentioned in Subsection 3.1 regarding our approach to improving the quality of diabetes prediction, we focus on combining multiple BGT-based models to develop our diabetes detection method. Based on the evaluation results of related research in Section 2, recent machine learning methods have achieved relatively high effectiveness in predicting diabetes. However, the rates of false positives and false negatives are still relatively high. Therefore, in this study, we propose combining multiple AI models to further reduce false positives and improve specificity.

We thoroughly tested 18 AI models in Subsection 4.3 using the PIMA dataset and the PyCaret library, version 3.0.4. The models with the best specificity, accuracy, and F1 scores were LightGBM (GBM), CatBoost (CBT), XGBoost (XGB), ExtraTrees (ET), and Random forest (RF), in that order. In addition to the aforementioned AI models, we also experimented with all three ensemble methods: stacking, hard voting, and soft voting. To assess the influence of each individual model when combined, we conducted experiments combining 3 to 5 individual models mentioned above. We evaluated the prediction results of all five models to select the meta-model with the best outcome for the stacking method. As a result, the diabetes prediction approach suggested in our study is known as GSDP, which stands for “GBT-based stacking ensemble-powered diabetes prediction”, with these 5 individual models. Its architecture is specifically illustrated in Fig. 2, with two phases: (i) The training phase handles the training tasks of the five individual models and subsequently trains the meta-model to form the diabetes prediction model GSDP; (ii) The detection phase is responsible for predicting diabetes disease from both the normalized testing set and the

Algorithm 2 GSDP: GBT-based stacking ensemble learning-powered diabetes prediction

Learner models: *GBM*, *XGB*, *CBT*, *ET*, *RF* - individually trained models.

Meta model: *MetaModel* - best of individual models.

Input: f - feature vector collected from a patient.

Output: 0 (normal) or 1 (diabetes).

```

1:  $F_{in} \leftarrow \text{OutlierHandling}(F)$       ▷ Determine the outlier threshold and apply it to handle the
   detected outliers.
2:  $F_{in} \leftarrow \text{Normalize}(F_{in})$       ▷ Normalize the feature vector.
3:  $L_{GBM} \leftarrow \text{GBM.predict}(F_{in}, \text{Conts})$       ▷ Perform the prediction using GBM.
4:  $L_{XGB} \leftarrow \text{XGB.predict}(F_{in}, \text{Conts})$       ▷ Perform the prediction using XGB.
5:  $L_{CBT} \leftarrow \text{CBT.predict}(F_{in}, \text{Conts})$       ▷ Perform the prediction using CBT.
6:  $L_{BME} \leftarrow \text{ET.predict}(F_{in}, \text{Conts})$       ▷ Perform the prediction using BME.
7:  $L_{RF} \leftarrow \text{RF.predict}(F_{in}, \text{Conts})$       ▷ Perform the prediction using BLM.
8:  $F_{meta} \leftarrow [L_{XGB}, L_{GBM}, L_{CBT}, L_{BME}, L_{RF}]$  ▷ Build the feature vector for meta-model-based
   prediction.
9:  $\text{score} \leftarrow \text{MetaModel.predict}(F_{meta})$       ▷ Perform the prediction using MetaModel.
10:  $\text{label} \leftarrow \text{score.argmax}(\text{axis} = 1)$       ▷ Get the predicted label, 0 or 1.
11: return  $\text{label}$ 

```

data collected from the clinical diagnosis process (diabetes data).

With this approach and the GSDP method, the process of analyzing and determining whether a patient may have diabetes (returning a result of 1) or not (returning a result of 0) from a dataset consisting of 8 features (including: Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, and Age) is described specifically in Alg. 2. It is also worth emphasizing here that our proposed method, SDAGS (combining SFDM for data augmentation and GSDP for prediction), uses all 8 features instead of applying additional attribute selection techniques as in other studies, such as [4, 26]. This is derived from analyzing the features using the SHAP results described in Subsection 4.3, where all 8 features have a significant impact on all AI prediction models.

3.4. AI-powered diabetes counseling and support framework

Based on the GSDP model, we propose to build an AI-powered diabetes counseling and support framework. Our framework consists of two components: (i) A social channel on Instagram that provides general information to support diabetic patients; (ii) A mobile app that stores input data collected from the non-invasive blood glucose sensor and suggests medical, dietary, and physical activity recommendations to the patients.

We have created an Instagram channel dedicated to providing information and knowledge about diabetes mellitus at https://www.instagram.com/diabetes_manager/. We plan to compile videos and articles discussing four main issues that diabetic patients must consider: (i) Blood glucose monitoring: We will provide information today about the most common types of controlled blood glucose deficiency. We also talk about different non-invasive devices that monitor blood glucose. (ii) Dietary counseling: We will provide information about how nutrition counseling can aid T2D. We will highlight the importance of a tailored diet and the role of a dietitian in optimizing BGL. (iii) Medical recommendations from doctors: We will provide recommendations and information about clinical practices for diagnosing and treating

diabetic patients at home or in hospital settings. We also talk about prescribing the right medications and amounts of insulin, and (iv) Physical activity advice: We will provide tips and recommendations regarding physical exercise for people with diabetes and its benefits. We will also link to some exercise videos that these patients can do at home.

Based on the analyses and evaluations of new methods in non-invasive blood glucose monitoring outlined in Section 2, we aim to develop a system that enables continuous, non-invasive blood glucose monitoring and management. This equipment system also needs to ensure the triple constraint of convenience, low cost, and high accuracy. Our blood glucose monitoring device targets individuals with T1D in developing countries, particularly in climates like Viet Nam which are characterized by hot, humid weather and high environmental pollution [17]. We anticipate that this device will use an ISF biosensor for two primary reasons: (i) ISF displays a wide range of blood glucose measurement values before and after meals, second only to blood analysis and significantly higher than saliva, tears, and sweat; (ii) Given the hot, humid climate and high environmental pollution, sensors using tears, sweat, and even saliva might not be suitable or easily influenced by the environment.

In addition to using ISF sensors, we also developed a mobile app to analyze the data collected from the ISF sensors. This app utilizes our GSDP to establish a blood glucose monitoring and management system for diabetes prediction. Furthermore, it will provide functionalities for monitoring and managing physical activity and diet for patients, alerting T1D patients of abnormally high blood sugar levels, enabling connectivity with diabetes specialists, and more.

4. EXPERIMENTS AND EVALUATION

To demonstrate the performance of our proposed method, SDAGS, we conduct a comprehensive experiment to answer the following research questions:

RQ1: Which data augmentation method performs best with the diabetes dataset?

RQ2: Is SFDM a superior method for enhancing the quality of the diabetes dataset, resulting in better outcomes compared to other data augmentation methods?

RQ3: How do AI models, including DL and GBT, perform in predicting diabetes?

RQ4: Is GSDP a more effective combined method for predicting diabetes compared to other ensemble methods and individual models?

RQ5: Does the SDAGS method, combining SDFM and GSDP, demonstrate superior diabetes prediction performance compared to recent state-of-the-art (SOTA) studies?

Our rigorous experiments were conducted to answer the above research questions. The following sections, in turn, detail the results we obtained while experimenting with and evaluating our SDAGS method.

4.1. Dataset collection

We empirically assessed the SDAGS method using the five specified RQs, employing the widely used Pima Indians diabetes dataset, namely PIMA, extensively used in tutorials on machine learning methods [31]. The dataset is derived from the National Institute of Diabetes and Digestive and Kidney Diseases and includes data on 768 women from a population surrounding Phoenix, Arizona, USA. The outcome under investigation is diabetes, with 258

individuals confirmed as positive and 500 individuals confirmed as negative. The dataset comprises nine features, as follows: (1) Pregnancies: Number of times pregnant, in integers; (2) Glucose: plasma glucose concentration at 2 hours in an oral glucose tolerance test (mg/dl) (integer); (3) BloodPressure: diastolic blood pressure (mm Hg) (integer); (4) SkinThickness: triceps skin fold thickness (mm) (integer); (5) Insulin: 2-h serum insulin (μ U/ml) (float); (6) BMI: body mass index (weight in kg/height in m) (float); (7) DiabetesPedigreeFunction: diabetes pedigree function (integer); (8) Age: age years (integer); (9) Outcome: class label (0: normal or 1: diabetes mellitus) (binary).

4.2. Experiment setup

All of our experiments deploy on an NVIDIA Tesla T4 (16 GB) GPU, 2 x Xeon Platinum 8160, and 256 GB of RAM. Besides that, we implement our method using the PyCaret v3.0.4, Tensorflow v2.11, Scikit-learn v1.2.2, PyTorch v1.13.1+cu117 frameworks, and the following libraries: Matplotlib V3.7.1, Numpy V1.20.2.

The selection of optimal hyperparameters for all GBT models was performed using the PyCaret framework. We determined the *learning_rate* to be 0.01, the *n_estimators* to be 100, and the *max_depth* to be 8. To evaluate the diabetes mellitus prediction method, we use standard metrics computed from the confusion matrix with the values of true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN), such as Specificity (SPEC), Sensitivity (SENS), F1-score (F1), Accuracy (ACC), Precision (PREC) [10, 26].

4.3. Evaluation of SFDM-based dataset augmentation

From the results analyzed in Section 2, we conducted experiments to evaluate various modern data augmentation methods, specifically SMOTE, WGAN, CTGAN, FDM, and our proposed method, SFDM. We assessed the enhanced quality of the PIMA dataset using the following approach: (i) Randomly dividing the PIMA dataset into 5 folds using k-fold cross-validation with $k = 5$; (ii) For each fold, performing data augmentation on the training set and using the test set for evaluation; (iii) Calculating evaluation metrics based on the average values across all folds. Within each fold, the original training set had 614 samples, and the test set had 154 samples.

From these folds, we will proceed to test and evaluate all 5 single models mentioned in Subsection 3.3. Table 1 illustrates the results obtained by synthesizing the average measurement value of all 5 folds. According to this table, the GBM method gives the best diabetes prediction results. Therefore, we will choose this method to evaluate the results of enhancing data quality.

For the SMOTE method, each training set had a different distribution of positive and negative samples, but the resulting count always doubled the largest sample count, ensuring an equal number of positive and negative samples. Because the total samples in PIMA are only 768, we empirically set each fold to generate 1,000 samples for each label for WGAN. Meanwhile, for FDM and SFDM, the total number of generated samples was analytically fixed at 2,000. The evaluation metrics assessing the effectiveness of the data augmentation process were synthesized in Table 2.

We also employ the SHapley Additive exPlanations (SHAP) technique to visualize the influence levels of all eight features in the PIMA dataset on the LightGBM model. SHAP is a game-theoretic approach to explaining the output of any machine learning model [20].

Table 1: Evaluation of AI models on raw training sets (bold is the best, %).

| Method | SPEC | SENS | F1 | ACC | PRE | ROC |
|--------|--------------|--------------|--------------|--------------|--------------|--------------|
| GBM | 92.66 | 80.75 | 83.28 | 88.28 | 86.12 | 87.69 |
| CBT | 90.78 | 82.28 | 82.43 | 87.89 | 82.74 | 86.72 |
| XGB | 90.24 | 82.33 | 81.87 | 87.63 | 81.49 | 86.19 |
| ET | 89.8 | 80.36 | 81.08 | 86.36 | 81.82 | 85.35 |
| RF | 88.93 | 83.24 | 80.83 | 87.11 | 78.73 | 85.2 |
| LSTM | 89.58 | 77.59 | 79.65 | 85.06 | 81.82 | 84.34 |
| CNN | 85.84 | 92.11 | 79.39 | 87.63 | 69.8 | 83.44 |

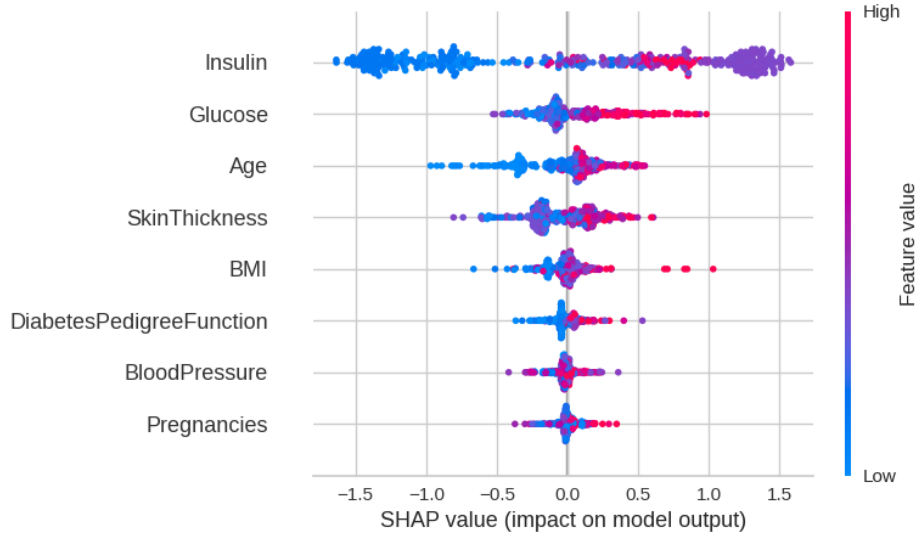


Figure 3: Impact of features on diabetes prediction model.

Specifically, Fig. 3 illustrates the impact of each feature on the LightGBM model’s prediction of diabetes. Consequently, these results also enable us to decide whether to utilize all eight features in training AI models to predict diabetes.

The experimental results evaluating the performance of the diabetes prediction method using GBM with training data augmented from the above methods are summarized in Table 2. This experimental result also allows us to answer the first two RQs. For RQ1, the data augmentation methods on the PIMA dataset are SFDM, FDM, WGAN, CTGAN, and SMOTE, respectively. Consequently, RQ2 also has an answer: SFDM demonstrates superior performance compared to other methods. The training data set quality has increased from 92.66% and 83.28% to 96.05% and 92.10% for Specificity and F1-score, respectively, thanks to SFDM.

4.4. Evaluation of GSDP-based diabetes prediction

We conducted extensive experiments to evaluate AI models, including DL and GBT, using the well-known PIMA dataset. As shown in Table 1, both the LightBGM and CBT models performed very well compared to others, with F1 scores of 83.28% and 82.43%, respectively. These results also allow us to partially answer RQ3, demonstrating the ability to predict

Table 2: Evaluation of training set augmentation using SMOTE, WGAN, CTGAN, FDM, and SFDM with GBM-based prediction (bold is the best, %).

| Training Set | SPEC | SENS | F1 | ACC | PRE | ROC |
|-------------------------|--------------|--------------|-------------|--------------|--------------|-------------|
| SFDM-based (our) | 96.05 | 91.88 | 92.1 | 94.67 | 92.34 | 94.1 |
| FDM-based | 93.76 | 84.38 | 86.47 | 90.37 | 88.81 | 94.4 |
| WGAN-based | 89.8 | 87.27 | 84.96 | 88.89 | 82.76 | 87.7 |
| CTGAN-based | 91.09 | 86.83 | 84.56 | 89.71 | 82.58 | 88.02 |
| SMOTE-based | 91 | 85.19 | 84.4 | 88.96 | 83.64 | 87.78 |
| Original | 92.66 | 80.75 | 83.28 | 88.28 | 86.12 | 87.69 |

Table 3: Assessment of ensemble learning using various combinations of GBT methods (bold is the best, %).

| Combination | SPEC | SENS | F1 | ACC | PRE | ROC |
|--------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| GBM+CBT+XGB+ET+RF | 97.80 | 90.32 | 93.33 | 94.77 | 96.55 | 98.58 |
| GBM+CBT+XGB+ET | 93.75 | 94.83 | 92.44 | 94.16 | 90.16 | 99.44 |
| GBM+CBT+XGB | 94.79 | 92.98 | 92.17 | 94.12 | 91.38 | 98.71 |

diabetes with all seven AI models. We experimented with merging individual AI models from the enhanced training dataset using the SFDM approach. The outcomes of merging a subset of the seven AI models we trained are presented in Table 3. This lets us choose the optimal diabetes prediction outcomes by combining five top models, e.i., GBM, CBT, XGB, ET, and RF. Furthermore, it provides a framework for applying the GSDP approach to these models. Remarkably, the results of our experiments using the soft voting combination approach to determine the optimal combinations for diabetes prediction outcomes serve as a foundation for conducting experiments with the three ensemble techniques: stacking, soft voting, and hard voting.

From there, we conducted experiments, evaluated, and compared the GSDP method based on stacking with the other two ensemble techniques across all four datasets: (i) The original PIMA dataset, (ii) PIMA balanced with SMOTE, namely SMOTE-DS, (iii) PIMA augmented with FDM, called FDM-DS, and (iv) PIMA enhanced with our SFDM, namely SFDM-DS. Table 4 presents the complete experimental results with all three ensemble learning methods on these four datasets. It is worth noting that these results are averaged from experiments conducted using 5-fold cross-validation, as mentioned in Subsection 4.3. Table 4 also allows us to assert that the SFDM data augmentation method yields superior results compared to the other methods. Additionally, when using our combined GSDP machine learning method, the diabetes prediction quality is very high. In particular, the evaluation metrics for F1, Accuracy, Precision, Sensitivity, and Specificity were 99.49%, 98.74%, 98.75%, and 98.00%, respectively. As can be seen, the SDAGS method combining both data augmentation using SFDM and GSDP has significantly improved performance, increasing Specificity by 9.03%, Sensitivity by 6.9%, F1 by 11.02%, Accuracy by 7.84%, and Precision by 14.67% in comparison to the original PIMA dataset (with corresponding results of 90.00%, 92.59%, 87.72%, 90.91%, and 83.33%). Consequently, these results also allow us to affirm that our proposed GSDP yields excellent prediction results, surpassing all single AI models as well as both soft-voting and hard-voting ensemble methods, thereby addressing RQ4.

Table 4: Assessment of ensemble learning methods (bold is the best).

| Dataset | Method | SPEC | SENS | F1 | ACC | PRE | ROC |
|---------------|-------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Original PIMA | Hardvoting | 90.00 | 92.59 | 87.72 | 90.91 | 83.33 | 89.54 |
| | Softvoting | 90.00 | 92.59 | 87.72 | 90.91 | 83.33 | 89.54 |
| | GSDP(our) | 90.00 | 92.59 | 87.72 | 90.91 | 83.33 | 89.54 |
| SMOTE-DS | Hardvoting | 88.00 | 93.00 | 90.73 | 90.50 | 88.57 | 90.60 |
| | Softvoting | 89.00 | 93.00 | 91.18 | 91.00 | 89.42 | 91.07 |
| | GSDP (our) | 89.00 | 94.00 | 91.71 | 91.50 | 89.52 | 91.60 |
| FDM-DS | Hardvoting | 93.60 | 96.45 | 95.00 | 95.00 | 93.60 | 95.02 |
| | Softvoting | 93.60 | 96.45 | 95.00 | 95.00 | 93.60 | 95.02 |
| | GSDP (our) | 95.57 | 95.94 | 95.70 | 95.75 | 95.45 | 95.75 |
| SFDM-DS | Hardvoting | 98.02 | 95.00 | 95.72 | 96.93 | 96.45 | 96.82 |
| | Softvoting | 98.31 | 94.50 | 95.70 | 96.93 | 96.92 | 96.93 |
| | GSDP (our) | 98.03 | 99.49 | 98.74 | 98.75 | 98.00 | 98.75 |

Table 5: Comparison with SOTA methods (bold is the best, %).

| Method | Venue/Year | SPEC | SENS | F1 | ACC | PREC |
|----------------------------|---|--------------|--------------|--------------|--------------|--------------|
| SDAGS (our) | - | 98.03 | 99.49 | 98.74 | 98.75 | 98.00 |
| KST-BLSTM + KMCSFO [26] | Springer J. of Supercomputing/24 | 96.56 | 96.01 | 96.98 | 97.26 | 97.89 |
| Stacking + Local Data [28] | Heliyon/24 | - | 97.00 | 96.00 | 95.50 | 94.00 |
| Stacking Ensemble [28] | Computer Methods & Programs in Biomedicine/24 | 95.00 | 98.00 | 87.00 | 93.50 | 96.00 |
| Fused ML [4] | IEEE Access/22 | 94.38 | 95.52 | 94.12 | 94.87 | 92.76 |
| CNN-LSTM [8] | Biomedical Signal Processing & Control/24 | - | 94.55 | 93.43 | 91.25 | 92.33 |
| KFPredict [25] | Biomedical Signal Processing & Control/23 | 98.00 | 85.00 | - | 93.50 | - |

4.5. Comparison with SOTAs

To assess our proposed method, we compare the experimental results of SDAGS with those of existing SOTA approaches using the same popular PIMA dataset. The comparison results are shown in Table 5 and were taken straight from the authors’ published papers. Thus, the findings demonstrate that SDAGS outperforms SOTA methods and obtains the top ratings across all evaluation categories. For example, SDAGS surpasses specificity by 1.47%, sensitivity by 3.48%, F1 by 1.76%, and accuracy by 1.49% compared to KST-BLSTM + KMCSFO [26] in 2024. Consequently, we can validate the effectiveness of our SDAGS method, and RQ5 has been answered.

5. CONCLUSIONS

In this paper, we investigate a holistic AI-powered diabetes mellitus prediction to effectively improve diabetes prediction specificity and sensitivity. Our proposed method, SDAGS, is

based on two key concepts: first, we enhance the quality of training data by integrating the SMOTE data balancing technique and FDM oversampling in the SFDM method, and second, we design the GSDP method by stacking five GBT AI models (GBM, CBT, XGB, ET, and RF) to further improve diabetes prediction performance. SDAGS outperforms existing SOTA methods, as demonstrated by our exhaustive experiments on the PIMA dataset utilizing the 5-fold cross-validation method. The model achieved sensitivity, specificity, F1, accuracy, and precision of 98.03%, 99.49%, 98.74%, 98.75%, and 98.00%, respectively. Subsequently, we put forth a preliminary framework for diabetes counseling and support that is driven by SDAGS. This framework aims to offer comprehensive assistance to patients with diabetes by addressing dietary counseling, physical activity, proper medication usage, and blood glucose monitoring.

Moving forward, our research will persist in enhancing the accuracy of diabetes prognostications through the development of novel machine learning techniques. Furthermore, we are committed to enhancing the capabilities of our AI-powered diabetes counseling and support framework to deliver more effective assistance to individuals with diabetes, particularly those with T1D residing in developing nations such as Viet Nam.

REFERENCES

- [1] “Design and development of diabetes management system using machine learning - hindawi.com,” <https://www.hindawi.com/journals/ijta/2020/8870141/>, [Accessed 16-02-2024].
- [2] “Leveraging artificial intelligence to improve chronic disease care: Methods and application to pharmacotherapy decision support for type-2 diabetes mellitus - ncbi.nlm.nih.gov,” <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8294941/>, [Accessed 09-03-2024].
- [3] A. Ahmed, S. Aziz, U. Qidwai, A. Abd-Alrazaq, and J. Sheikh, “Performance of artificial intelligence models in estimating blood glucose level among diabetic patients using non-invasive wearable device data,” *Computer Methods and Programs in Biomedicine Update*, vol. 3, p. 100094, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2666990023000034>
- [4] U. Ahmed, G. F. Issa, M. A. Khan, S. Aftab, M. F. Khan, R. A. T. Said, T. M. Ghazal, and M. Ahmad, “Prediction of diabetes empowered with fused machine learning,” *IEEE Access*, vol. 10, pp. 8529–8538, 2022.
- [5] K. Anthony and L. Caizhi, “Salivary glucose measurement: A holy ground for next generation of non-invasive diabetic monitoring,” *Hybrid Advances*, vol. 3, p. 100052, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2773207X23000350>
- [6] M. Arjovsky, S. Chintala, and L. Bottou, “Wasserstein gan,” 2017.
- [7] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, “Smote: synthetic minority over-sampling technique,” *J. Artif. Int. Res.*, vol. 16, no. 1, p. 321–357, jun 2002.
- [8] L. Z. Chee, S. Sivakumar, K. H. Lim, and A. A. Gopalai, “Gait acceleration-based diabetes detection using hybrid deep learning,” *Biomedical Signal Processing and Control*, vol. 92, p. 105998, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1746809424000569>
- [9] H.-K. Chen, F.-H. Chen, and S.-F. Lin, “An ai-based exercise prescription recommendation system,” *Applied Sciences*, vol. 11, no. 6, 2021. [Online]. Available: <https://www.mdpi.com/2076-3417/11/6/2661>

- [10] H. P. Du, D. H. Pham, and H. N. Nguyen, “An efficient parallel method for performing concurrent operations on social networks,” in *Computational Collective Intelligence*, N. T. Nguyen, G. A. Papadopoulos, P. Jedrzejowicz, B. Trawinski, and G. Vossen, Eds. Cham: Springer International Publishing, 2017, pp. 148–159.
- [11] E. Everett, B. Kane, A. Yoo, A. Dobs, and N. Mathioudakis, “A novel approach for fully automated, personalized health coaching for adults with prediabetes: Pilot clinical trial,” *J Med Internet Res*, vol. 20, no. 2, p. e72, Feb 2018. [Online]. Available: <http://www.ncbi.nlm.nih.gov/pubmed/29487046>
- [12] Z. Guan, H. Li, R. Liu, C. Cai, Y. Liu, J. Li, X. Wang, S. Huang, L. Wu, D. Liu, S. Yu, Z. Wang, J. Shu, X. Hou, X. Yang, W. Jia, and B. Sheng, “Artificial intelligence in diabetes management: Advancements, opportunities, and challenges,” *Cell Reports Medicine*, vol. 4, no. 10, p. 101213, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2666379123003804>
- [13] Guidance, “Type 1 diabetes in adults: diagnosis and management,” 2023, [Online]; Retrieved 19-Nov-2023]. [Online]. Available: <https://www.nice.org.uk/guidance/ng17>
- [14] A. Hennebelle, L. Ismail, H. Materwala, J. Al Kaabi, P. Ranjan, and R. Janardhanan, “Secure and privacy-preserving automated machine learning operations into end-to-end integrated iot-edge-artificial intelligence-blockchain monitoring system for diabetes mellitus prediction,” *Computational and Structural Biotechnology Journal*, vol. 23, pp. 212–233, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2001037023004531>
- [15] A. Jolicoeur-Martineau, K. Fatras, and T. Kachman, “Generating and imputing tabular data via diffusion and flow-based gradient-boosted trees,” 2024.
- [16] A. Ko and C. Liao, “Salivary glucose measurement: A holy ground for next generation of non-invasive diabetic monitoring,” *Hybrid Advances*, vol. 3, p. 100052, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2773207X23000350>
- [17] H. V. Le, O. V. Phung, and H. N. Nguyen, “Information security risk management by a holistic approach: a case study for vietnamese e-government,” *International Journal of Computer Science and Network Security*, vol. 20, no. 6, pp. 72–82, 08 2020.
- [18] H. Lee, Y. J. Hong, S. Baik, T. Hyeon, and D.-H. Kim, “Enzyme-based glucose sensor: From invasive to wearable device,” *Adv Healthc Mater*, vol. 7, no. 8, p. e1701150, Jan. 2018.
- [19] S. Liu, Z. Shen, L. Deng, and G. Liu, “Smartphone assisted portable biochip for non-invasive simultaneous monitoring of glucose and insulin towards precise diagnosis of prediabetes/diabetes,” *Biosensors and Bioelectronics*, vol. 209, p. 114251, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0956566322002913>
- [20] S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS’17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 4768–4777.
- [21] M. Nayak, S. Gupta, J. Sunitha, G. Dawar, N. Sinha, and N. Rallan, “Correlation of salivary glucose level with blood glucose level in diabetes mellitus,” *Journal of Oral and Maxillofacial Pathology*, vol. 21, no. 3, p. 334, 2017. [Online]. Available: http://dx.doi.org/10.4103/jomfp.jomfp_222_15
- [22] W. H. Organization, “Diabetes,” [Online]; Retrieved 19-Jan-2024]. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/diabetes>

- [23] A. Petersmann, D. Müller-Wieland, U. A. Müller, R. Landgraf, M. Nauck, G. Freckmann, L. Heinemann, and E. Schleicher, "Definition, classification and diagnosis of diabetes mellitus," *Experimental and Clinical Endocrinology & Diabetes : Official Journal, German Society of Endocrinology and German Diabetes Association*, vol. 127, no. S 01, pp. S1–S7, December 2019. [Online]. Available: <http://www.thieme-connect.de/products/ejournals/pdf/10.1055/a-1018-9078.pdf>
- [24] G. Prabhakar, V. R. Chintala, T. Reddy, and T. Ruchitha, "User-cloud-based ensemble framework for type-2 diabetes prediction with diet plan suggestion," *e-Prime - Advances in Electrical Engineering, Electronics and Energy*, vol. 7, p. 100423, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2772671124000056>
- [25] H. Qi, X. Song, S. Liu, Y. Zhang, and K. K. Wong, "Kfpredict: An ensemble learning prediction framework for diabetes based on fusion of key features," *Computer Methods and Programs in Biomedicine*, vol. 231, p. 107378, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0169260723000457>
- [26] R. Rajalakshmi, P. Sivakumar, L. K. Kumari, and M. C. Selvi, "A novel deep learning model for diabetes mellitus prediction in iot-based healthcare environment with effective feature selection mechanism," *The Journal of Supercomputing*, vol. 80, no. 1, pp. 271–291, Jan 2024. [Online]. Available: <https://doi.org/10.1007/s11227-023-05496-6>
- [27] RetinaRisk, "Saliva-based glucose tests for diabetes management," [Online]; Retrieved 16-Nov-2023]. [Online]. Available: <https://www.retinarisk.com/blog/saliva-based-glucose-test-for-diabetes-management/>
- [28] M. S. Reza, R. Amin, R. Yasmin, W. Kulsum, and S. Ruhi, "Improving diabetes disease patients classification using stacking ensemble method with pima and local healthcare data," *Heliyon*, vol. 10, no. 2, p. e24536, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S240584402400567X>
- [29] L. Rokach, *Ensemble Learning: Pattern Classification Using Ensemble Methods (Second Edition)*, 2nd ed. Singapore: World Scientific Publishing Co Pte Ltd, 2019.
- [30] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 10 684–10 695.
- [31] V. Sigillito, "Pima indians diabetes database," *Applied Physics laboratory, The Johns Hopkins University, laurel, MD1990*, 1990.
- [32] K. Stawarz, D. Katz, A. Ayobi, P. Marshall, T. Yamagata, R. Santos-Rodriguez, P. Flach, and A. A. O’Kane, "Co-designing opportunities for human-centred machine learning in supporting type 1 diabetes decision-making," *International Journal of Human-Computer Studies*, vol. 173, p. 103003, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1071581923000095>
- [33] M. Sun, X. Pei, T. Xin, J. Liu, C. Ma, M. Cao, and M. Zhou, "A flexible microfluidic chip-based universal fully integrated nanoelectronic system with point-of-care raw sweat, tears, or saliva glucose monitoring for potential noninvasive glucose management," *Analytical Chemistry*, vol. 94, no. 3, pp. 1890–1900, Jan 2022. [Online]. Available: <https://doi.org/10.1021/acs.analchem.1c05174>
- [34] L. Tang, S. J. Chang, C.-J. Chen, and J.-T. Liu, "Non-invasive blood glucose monitoring technology: A review," *Sensors (Basel, Switzerland)*, vol. 20, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:228081145>

- [35] H. V. Vo, H. P. Du, and H. N. Nguyen, "Ai-powered intrusion detection in large-scale traffic networks based on flow sensing strategy and parallel deep analysis," *Journal of Network and Computer Applications*, vol. 220, p. 103735, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1084804523001546>
- [36] H. V. Vo, D. H. Nguyen, T. T. Nguyen, H. N. Nguyen, and D. V. Nguyen, "Leveraging ai-driven realtime intrusion detection by using wgan and xgboost," in *Proceedings of the 11th International Symposium on Information and Communication Technology*. New York, NY, USA: Association for Computing Machinery, 2022, p. 208–215.
- [37] H. V. Vo, P. H. Nguyen, H. T. Nguyen, D. B. Vu, and H. N. Nguyen, "Enhancing AI - powered malware detection by parallel ensemble learning," in *2023 RIVF International Conference on Computing and Communication Technologies (RIVF)*, 2023, pp. 503–508.
- [38] H. V. Vo, H. P. Du, and H. N. Nguyen, "Apelid: Enhancing real-time intrusion detection with augmented wgan and parallel ensemble learning," *Computers & Security*, vol. 136, p. 103567, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167404823004777>
- [39] A. P. Wright, A. T. Wright, A. B. McCoy, and D. F. Sittig, "The use of sequential pattern mining to predict next prescribed medications," *Journal of Biomedical Informatics*, vol. 53, pp. 73–80, 2015. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1532046414002007>
- [40] W. Zhang, Q. Yu, B. Siddiquie, A. Divakaran, and H. Sawhney, "'snap-n-eat': Food recognition and nutrition estimation on a smartphone," *Journal of Diabetes Science and Technology*, vol. 9, no. 3, pp. 525–533, 2015, PMID: 25901024. [Online]. Available: <https://doi.org/10.1177/1932296815582222>

Received on March 17, 2024

Accepted on May 14, 2024