

Enhancing Fairness and Accuracy in Diagnosing Type 2 Diabetes in Young Adult Population

Tanmoy Sarkar Pias, Yiqi Su, Xuxin Tang, Haohui Wang, Shahriar Faghani, and Danfeng (Daphne) Yao

Fellow, IEEE

Abstract—While type 2 diabetes is predominantly found in the elderly population, recent publications indicate an increasing prevalence in the young adult population. Failing to diagnose it in the minority younger age group could have significant adverse effects on their health. Several previous works acknowledge the bias of machine learning models towards different gender and race groups and propose various approaches to mitigate it. However, those works failed to propose any effective methodologies to diagnose diabetes in the young population, which is the minority group in the diabetic population. This is the first paper where we mention digital ageism towards the young adult population diagnosing diabetes. In this paper, we identify this deficiency in traditional machine learning models and propose an algorithm to mitigate the bias towards the young population when predicting diabetes. Deviating from the traditional concept of one-model-fits-all, we train customized machine-learning models for each age group. Our pipeline trains a separate machine learning model for every 5-year age band (i.e., age groups 30-34, 35-39, and 40-44). The proposed solution consistently improves recall of diabetes class by 26% to 40% in the young age group (30-44). Moreover, our technique outperforms 7 commonly used whole-group resampling techniques (i.e., random oversampling, random undersampling, SMOTE, ADASYN, Tomek-links, ENN, and Near Miss) by at least 36% in terms of diabetes recall in the young age group. Feature important analysis shows that the age attribute has a significant contribution to the decision of the original model, which was marginalized in the age-personalized model. Our method shows improved performance (e.g., balanced accuracy improved 7-12%) over multiple machine learning models and multiple sampling algorithms.

Index Terms—Ageism, Fairness, Diabetes, Machine Learning, Healthcare, BRFS

I. INTRODUCTION

DIAGNOSING chronic diseases like diabetes is crucial, given the substantial global burden of diabetes-related

Tanmoy Sarkar Pias is with the Department of Computer Science, Virginia Tech, Blacksburg, VA 24060, USA (e-mail: tanmoysarkar@vt.edu).

Yiqi Su is with the Department of Computer Science, Virginia Tech, Blacksburg, VA 24060, USA (e-mail: yiqisu@vt.edu).

Xuxin Tang is with the Department of Computer Science, Virginia Tech, Blacksburg, VA 24060, USA (e-mail: xuxintang@vt.edu).

Haohui Wang is with the Department of Computer Science, Virginia Tech, Blacksburg, VA 24060, USA (e-mail: haohuiw@vt.edu).

Shahriar Faghani is with the Radiology Informatics Lab, Mayo Clinic, Rochester, MN, USA (e-mail: faghani.shahriar@mayo.edu).

Danfeng (Daphne) Yao is with the Department of Computer Science, Virginia Tech, Blacksburg, VA 24060, USA (e-mail: danfeng@vt.edu).

complications and deaths, particularly in low- and middle-income countries [1], [2]. At the present rate of expansion, the International Diabetes Federation predicts that by 2045, a staggering 693 million individuals globally will be affected by diabetes [3]. The prevalence of diabetes, specifically type 2 diabetes, has been steadily increasing over the past few decades [4], [5]. While diabetes has traditionally been associated with the elderly population, recent studies indicate a rising prevalence of diabetes among the younger population as well [6]. According to the Centers for Disease Control and Prevention (CDC), by 2060, type 2 diabetes cases might increase by about 70-700% in the young population [7]. Young adults with diabetes face a significantly elevated risk of early health complications and even premature death compared to their counterparts without diabetes [8]. Younger people diagnosed with diabetes face an increased risk of developing early and severe complications, encompassing microvascular (retinopathy, neuropathy, ulceration, nephropathy) and macrovascular (cardiovascular, cerebrovascular, peripheral vascular) diseases [9]–[11]. We define the population group aged 30-44 years as the “young adult group” relative to the whole age group 30-85 as well as based on multiple studies that highlight this age range as a rising-risk demographic for early-onset type 2 diabetes. [10], [12]–[14] Early detection and awareness can help the young population at risk take steps to prevent or delay type 2 diabetes and early intervention can even reverse prediabetes [15], [16].

Machine learning (ML) has been increasingly integrated into the healthcare systems [17], [18] because of its potential to assist clinicians and medical doctors in taking better care of patients. Many machine learning models have been applied to diagnose diabetes [5]. However, health data can be imbalanced, which could potentially lead machine learning models to learn patterns with existing bias from the provided data [19]–[22] causing low performance [23]. Data bias, if not addressed, can exacerbate and perpetuate inequalities in the performance of algorithms in different subgroups [24], [25], particularly in historically underserved populations like female patients [26], black patients, or those with low socioeconomic status [27]. Artificial intelligence (AI) based models can be susceptible to digital ageism, potentially leading to biased diagnoses that could harm patients [28], [29]. Moreover, AI models existing in the healthcare domain can show faithfulness issues [30].

This paper identifies that the traditional machine learning models such as Logistic regression (LR), Multi-Layer Per-

TABLE I
DATA DISTRIBUTION AMONG DIFFERENT AGE GROUPS

Dataset	Age group	Number of Samples	Diabetes Negative	Diabetes Positive	Ratio (Pos/Neg)
BRFSS 2021	30-34	11,188 (5.59%)	10,864 (97.10%)	324 (2.90%)	0.03
	35-39	13,878 (6.93%)	13,290 (95.76%)	588 (4.24%)	0.04
	40-44	15,724 (7.86%)	14,661 (93.24%)	1,063 (6.76%)	0.07
	45-49	16,077 (8.03%)	14,466 (89.98%)	1,611 (10.02%)	0.11
	50-54	19,357 (9.67%)	16,836 (86.98%)	2,521 (13.02%)	0.15
	55-59	21,757 (10.87%)	18,155 (83.44%)	3,602 (16.56%)	0.20
	60-64	24,672 (12.33%)	20,140 (81.63%)	4,532 (18.37%)	0.23
	65-69	25,526 (12.75%)	20,422 (80.00%)	5,104 (20.00%)	0.25
	70-74	23,121 (11.55%)	17,891 (77.38%)	5,230 (22.62%)	0.29
	75-79	14,740 (7.36%)	11,257 (76.37%)	3,483 (23.63%)	0.31
80-99	14,096 (7.04%)	11,314 (80.26%)	2,782 (19.74%)	0.25	
BRFSS 2019	30-34	12,071 (5.72%)	11,751 (97.35%)	320 (2.65%)	0.03
	35-39	14,024 (6.65%)	13,410 (95.62%)	614 (4.38%)	0.05
	40-44	14,444 (6.85%)	13,528 (93.66%)	916 (6.34%)	0.07
	45-49	16,166 (7.67%)	14,539 (89.94%)	1,627 (10.06%)	0.11
	50-54	19,550 (9.27%)	16,999 (86.95%)	2,551 (13.05%)	0.15
	55-59	23,926 (11.35%)	20,197 (84.41%)	3,729 (15.59%)	0.18
	60-64	27,223 (12.91%)	22,316 (81.97%)	4,907 (18.03%)	0.22
	65-69	27,255 (12.92%)	21,540 (79.03%)	5,715 (20.97%)	0.27
	70-74	24,115 (11.44%)	18,609 (77.17%)	5,506 (22.83%)	0.30
	75-79	15,958 (7.57%)	12,254 (76.79%)	3,704 (23.21%)	0.30
80-99	16,142 (7.65%)	12,960 (80.29%)	3,182 (19.71%)	0.25	

ception (MLP), Naive Bayes (NB), AdaBoost (AB), Random Forest (RF), and K Nearest Neighbor (KNN), trained on imbalance Behavioral Risk Factor Surveillance System (BRFSS) dataset - with only 15% representing the diabetes population [31], tend to misdiagnose diabetes more frequently in the younger population (30-44 years) compared to other subgroups.

We propose an effective solution which successfully mitigates bias from young groups and increases type 2 diabetes (T2D) diagnosing sensitivity. None of the existing papers developed any effective machine-learning-based approach for effectively diagnosing T2D in the young population (30-44 years). To the best of our knowledge, our work is the first precision T2D diagnosis paper using machine learning and improving the diagnostic performance of machine learning models to diagnose T2D in young populations.

In this study, we train diverse machine learning models, including a transparent model like logistic regression, to uncover the roots of bias and missed detections. This is achieved through a thorough analysis of association coefficients responsible for shaping model decisions. The major contribution is listed here.

- **Contribution:** We identified digital agism in detecting diabetes in the young adult population and mitigated the issue to achieve equitable performance.
- **Method:** We effectively utilized the Enhanced Demographic-Parity (DP) resampling that adjusts class balance only inside the target subgroup, followed by age-specific model selection.
- **Improvement:** The enhanced model shows +26–40% recall and +7–13% balanced accuracy for young adults over the baseline model.

II. METHODS

A. Dataset

The Behavioral Risk Factor Surveillance System (BRFSS) is an annual health survey conducted in the United States to monitor various health behaviors such as cardiovascular diseases, chronic diseases, diabetes, obesity, and other risk factors that contribute to the leading causes of death and disability [31]. The BRFSS datasets have been collected from all 50 states in the U.S. and the District of Columbia and included responses from over 400,000 participants which makes it one of the largest publicly available datasets related to public health. Each record contains an individual's BRFSS survey responses on various health behaviors and risk factors such as tobacco use, physical activity, alcohol consumption, existing chronic diseases, and mental health. The survey also gathered demographic information such as age, gender, and race/ethnicity, which are helpful to explore important correlations and even causation.

This survey is conducted every year, and the CDC makes it publicly available for research. In this study, we selected the BRFSS dataset from 2021, which contains more than 400,000 subject information and 330+ attributes from each subject. The BRFSS dataset is a valuable resource for identifying health disparities and evaluating the effectiveness of public health programs and policies. However, the dataset is not free from challenges because a large portion of attribute values are missing (25%) and this dataset can be a highly imbalanced data imbalance (diabetes class 15%).

B. Data Preprocessing

According to the literature on the risk factors for diabetes, [5], [32], we selected 30 attributes including diabetes labels. Subjects aged over 30 are selected for this study [5], [33], [34]. We identified individuals with type 2 diabetes based on their self-reported response to the question: Has a doctor, nurse, or other health professional ever told you that you have diabetes? Individuals who answered "Yes" were labeled as

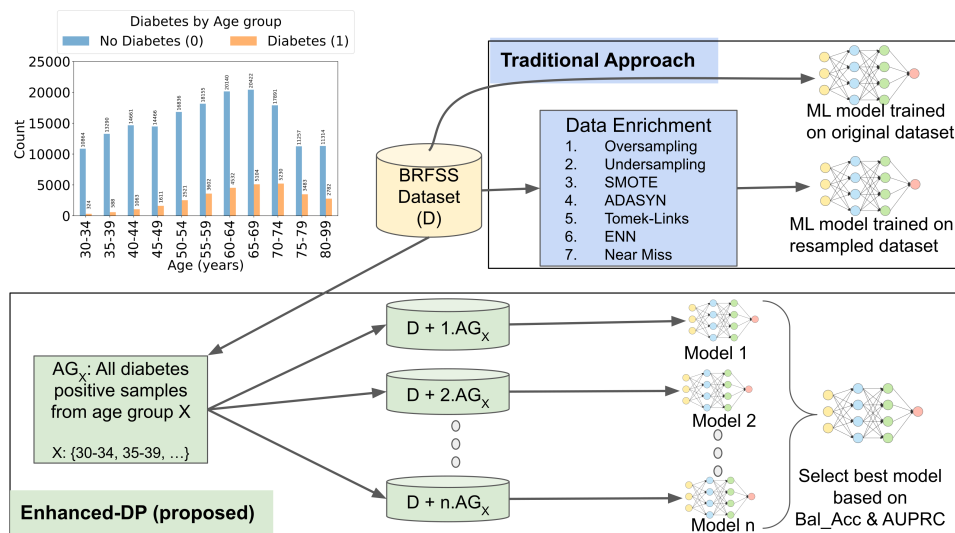


Fig. 1. The Enhanced-DP approach in contrast to traditional approaches, enriches the minority age groups and creates new training sets by replicating diabetic samples (1 to n times) from a minority age group. n machine learning models are trained on each of the n versions of the training sets. The best model is selected based on performance metric balance accuracy (Bal.Acc) and area under precision and recall curve (AUPRC). The top left bar chart represents age distribution (histogram) in the original dataset.

having diabetes. [5] We create age groups spanning 5 years starting from 30 up to 80+ years. The selected cohort contains 200,136 subjects information where 169,296 subjects (84.6%) are diabetes negative and 30,840 subjects (15.4%) are positive.

As we mentioned before, this dataset is not balanced for the age, gender, and racial subgroups. In the age group 30-34, the number of negative cases is 10,864 but the number of positive cases is only 324. The age groups 35-39 and 40-44 are also highly imbalanced as the positive-negative ratio is only 4.4% and 7.2%. All subgroup data distributions are shown in Table I.

The selected variables fall into three distinct types: nominal, ordinal, and binary. Nominal variables lack any inherent order, ordinal variables possess a meaningful order, and binary variables exclusively hold two distinct values. For example, the presence of high blood pressure, cholesterol, or heart disease can be represented by a binary variable. BMI category, education level, and income level are considered ordinal variables. On the other hand, marital status and race are nominal variables. Binary variables are represented using a binary encoding, where 1 signifies a positive outcome and 0 represents a negative outcome. Ordinal variables are encoded using integers, preserving their meaningful order. Nominal variables are transformed into one-hot encoding for appropriate representation. The selected variables are listed in Fig. 9(c).

C. Performance Metrics

We evaluate model performance using standard performance metrics. Starting with Recall of Class 1 (Rec.C1), also referred to as sensitivity, represents the true positive rate for the diabetes-positive class—measuring the proportion of actual diabetes cases correctly identified. Recall of Class 0 (Rec.C0), also referred to as specificity, indicates the true negative rate for the diabetes-negative class—capturing the proportion of

non-diabetes cases correctly identified. We also report the area under the precision-recall curve (AU_PRC) for both classes (PRC.C0 and PRC.C1 for class 0 and 1 respectively) to assess the trade-off between precision and recall. Additional metrics include overall accuracy (ACC) which represents the total proportion of correct predictions). The balanced accuracy (Bal.Acc) represents the average of sensitivity and specificity, which accounts for class imbalance, and area under the receiver operating characteristic curve (AU_ROC), summarizes the model's ability to distinguish between classes. Additionally, we report the F1 score, which is the harmonic mean of precision and recall and is especially useful when class distributions are imbalanced. Finally, we include the Matthews Correlation Coefficient (MCC), a balanced metric that takes into account true and false positives and negatives, offering a more informative measure in scenarios with skewed class distributions. Throughout the manuscript, we use C1 to denote the diabetes-positive class and C0 for the diabetes-negative class.

D. Bias Mitigation Approach

We utilized a modified and enhanced version of the prioritized (DP) bias correction method (Fig. 1) which is inspired by [35] prioritizes a specific subgroup, the young age group in this case, that suffers from data imbalance. We incrementally replicate data points of the minority class (diabetes positive indicated by class 1) and choose an optimal unit of replication based on the model performance. As a result, the enriched training set contains the original samples as well as the replicated samples. However, the vanilla DP technique by [35] has several shortcomings, including the selection of replication units.

Our proposed Enhanced-DP technique replicates all samples of the diabetes class from the young population up to n times. Each time the duplicated sample units are merged with the

original dataset, including the original set. In our experiment, we set the maximum DP unit “ n ” based on the lowest subgroup positive-negative sample ratio. In order to achieve a balanced set, the lowest subgroup ratio is for the age group 30-34. The number of negative cases is 10,864 but the number of positive cases is 324. So, the ratio is 34, which is the limit of DP unit n . We also employed early stopping to select the best n quickly as the performance curve contains a global maximum point. This makes our Enhanced-DP algorithm faster compared to the original DP algorithm [35].

Each training set containing the original set and the duplicated n (1 to 34) units is used to train a single model. In this way, we train 35 models (34 models trained on the enriched training sets and 1 model trained on the original training set) and select the best-performing model based on balanced accuracy and the area under the precision-recall curve (AU_PRC) of class 1. We identify the top three models with the highest balanced accuracy values and select the model that gives the highest class 1 AU_PRC for that particular age group. This selected model is used only for diagnosing diabetes for that particular group, e.g., the Enhanced-DP model trained with duplicated age-group 30-34 samples is used to diagnose new patients from age-group 30-34 years.

E. Sampling Algorithms

A frequently employed technique for addressing the challenges of data imbalance is the utilization of sampling methods. To evaluate the effectiveness of the sampling algorithms detecting diabetes in young groups, we first compared multiple sampling algorithms, including (1) random oversampling, randomly duplicating instances from the minority class [36]; (2) random undersampling, randomly removing instances from the majority class [36]; (3) Synthetic Minority Over-sampling Technique (SMOTE), creating synthetic samples for the minority class by interpolating between existing instances [37]; (4) Adaptive Synthetic Sampling (ADASYN), similar to SMOTE, but generating more synthetic samples for difficult-to-learn instances [38]; (5) Tomek Links, removing instances from the majority class that are close to instances in the minority class [39]; (6) Edited Nearest Neighbors (ENN), removing instances from the majority class that are misclassified based on the nearest neighbors from both classes [40]; and (7) NearMiss, selecting instances from the majority class based on their distance to instances in the minority class [41]. These techniques are widely adopted and remain foundational in addressing class imbalance and have been proven to be effective in the literature [42]–[45]; however, these techniques are not tailored to tackle subgroup biases. As a result, diabetes is misdiagnosed in the young population. Therefore, we utilize a new concept of using one model for a single group which deviates from the traditional one model-fits-all ideology.

F. Machine Learning Models

We selected six commonly used ML models including Logistic Regression (LR), Random Forest (RF), Adaboost (AB), Multilayer Perceptron (MLP), and Naive Bayes (NB) to evaluate the effectiveness of the sampling strategy. We

purposely selected simple models such as logistic regression because of their interpretability. It is very important for the science community, especially for healthcare to create interpretable models to find out the root cause of a prediction, however, [46] evaluated 511 scientific papers across different ML domains and identified a notable deficiency in reproducibility metrics, including dataset and code accessibility in clinical ML domain papers. Moreover, complex and high-end models which have recently been applied to the healthcare domain might pose difficulties in reproducibility [47]. Using simple models like logistic regression will pave the way both for explaining the results and making it easy for other researchers to reproduce the model and the same results. We used grid search to tune hyperparameters both for best performance and reproducibility. The hyperparameter search set and selected parameters are listed in Supplementary Table 2.

The dataset undergoes a random division into three disjoint sets namely training (60% or 120,081 samples), validation (20% or 40,027 samples), and testing set (20% or 40,028 samples) through the utilization of the Python sklearn library. The training dataset serves as the foundation for model training, while the testing dataset remains consistent across all the experiments (for each trial), ensuring the robustness and comparability of our results. The validation set is used for hyperparameter tuning and to determine the number of underlying models n in EnhancedDP method, ensuring no data leakage. The performance of the models in each experiment is conducted for 5 independent trials with different splits of train, validation, and test set for consistency. Finally, we calculated the mean and 95% confidence intervals (and standard deviations) from the trials for consistent results.

To examine the potential unfairness in the datasets, we calculated class-based accuracy, recall, AUROC, and PRC. Nonetheless, our primary focus lies in prioritizing the recall of the positive class, also known as sensitivity, as the detection of diabetes holds paramount importance in our specific case.

G. Model Calibration and Threshold Tuning

Before converting probabilities into class labels, we calibrate each model's output on the validation set using Isotonic Regression [48] to improve probability accuracy. We then tune the decision threshold to optimize balanced accuracy and F1 score for Class 1, allowing better control over the trade-off between false positives and false negatives.

We first identify the top 3 thresholds with the highest F1_C1 scores and then select the one with the highest balanced accuracy on the full validation set to avoid overfitting, especially in subgroups with limited samples (<100). The final threshold (0.195) is averaged over 10 trials. Both F1 score and balanced accuracy are used to ensure balanced performance.

III. RESULTS

The recall of the positive class (Rec_C1) in the 30-34 age group is only 30%, whereas the Rec_C1 is 68-72% in the gender group and 66-84% in the ethnic group. Moreover, we bring attention to the fact that solely relying on Area Under the Receiver Operating Characteristic curve (AUROC) can be

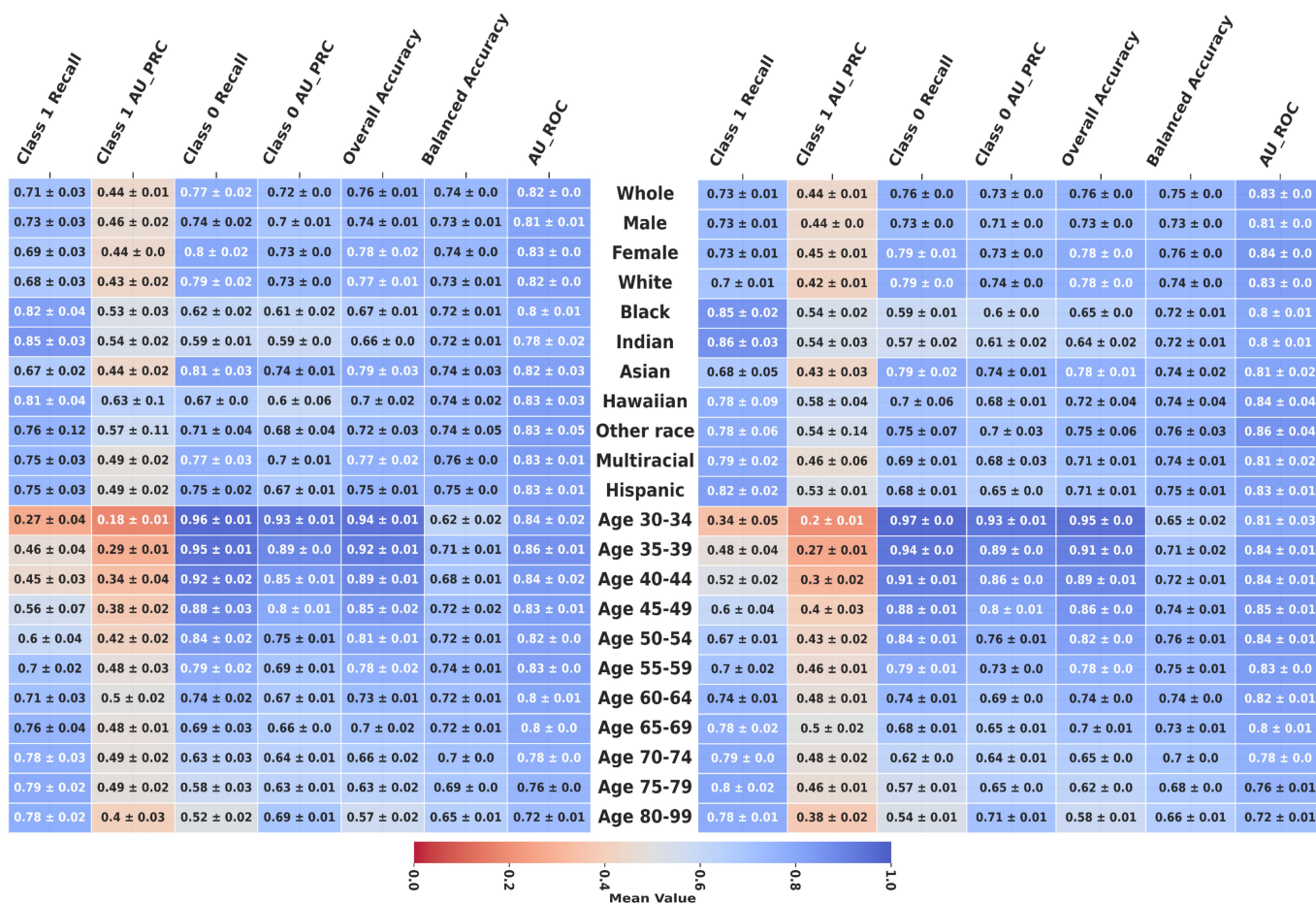


Fig. 2. Performance of logistic regression model trained and tested on the original BRFSS 2021 (left) and BRFSS 2015 (right). The x-axis represents performance metrics and the y-axis represents groups (whole population and subpopulations). C1 and C0 stand for class 1 (diabetes positive) and class 0 (diabetes negative) respectively. Performance scores are reported using mean +/- 95% confidence interval from multiple independent trials.

misleading – while the overall group’s AUROC is 82%, the age group 30-34 exhibits a seemingly high 84% AUROC that masks the poor diabetes detection rate (Rec.C1) within that specific age group. So, we focus mostly on recall of class 1, balanced accuracy (average recall of class 1 and class 0), and AUROC for a fair comparison. Recall class 1 (Rec.C1) is also known as Sensitivity or True Positive Rate (TPR) and recall class 0 (Rec.C0) is known as Specificity or True Negative Rate (TNR).

A. Performance of Original Model

The machine learning models trained on the original training set show different performances for the minority diabetes-positive group (C1) and the majority-negative group (C0). Fig. 2 (a) and (b) show the performance of the original logistic regression model on 2021 and 2015 datasets. The whole group Rec.C1 is 0.70-0.72 whereas Rec.C0 is 0.77-0.78. For males, Rec.C1 is 0.72 and for females Rec.C1 is 0.68. The ethnic groups such as White, Asian, and multicultural show Rec.C1 of 0.66, 0.66, and 0.73, respectively. On the other hand, Black, Indian, and Hawaiian show Rec.C1 of 0.84, 0.82, and 0.84, respectively.

However, the young age group suffers the most from missed

detection of diabetes. For age groups [30-34], [35-39], and [40-44] Rec.C1 is only 0.30-0.31, 0.42-0.45, and 0.45-0.50, respectively. This means out of 100 diabetic patients aged 30-34, 70 patients are misdiagnosed. The Adaboost, Naive bias, random forest, and KNN models also show similar performance, represented in Fig. 2 (c). The standard deviation of each metric is less than 0.04 (from multiple experimental trials) unless mentioned otherwise. The model behaves similarly in the BRFSS 2015 dataset.

In imbalanced datasets, commonly employed metrics like AUROC and accuracy can be misleading and do not accurately represent the performance of the minority class. Despite potentially poor performance in the minority class, these metrics might indicate falsely elevated values. The AUROC of the whole population, age groups [30-34], [35-39], and [40-44] are 0.82, 0.84, 0.87, and 0.84. However, these age groups show very poor Rec.C1 which is not reflected with AUROC. In contrast to the AUROC and accuracy metrics which can be overly optimistic, we use Rec.C1 to measure the true detection rate reflecting the performance of the model.

We also statistically analyze the performance difference (Supplementary Table IV) [49]. For each metric, we tested residual normality (Shapiro–Wilk) and homogeneity of vari-

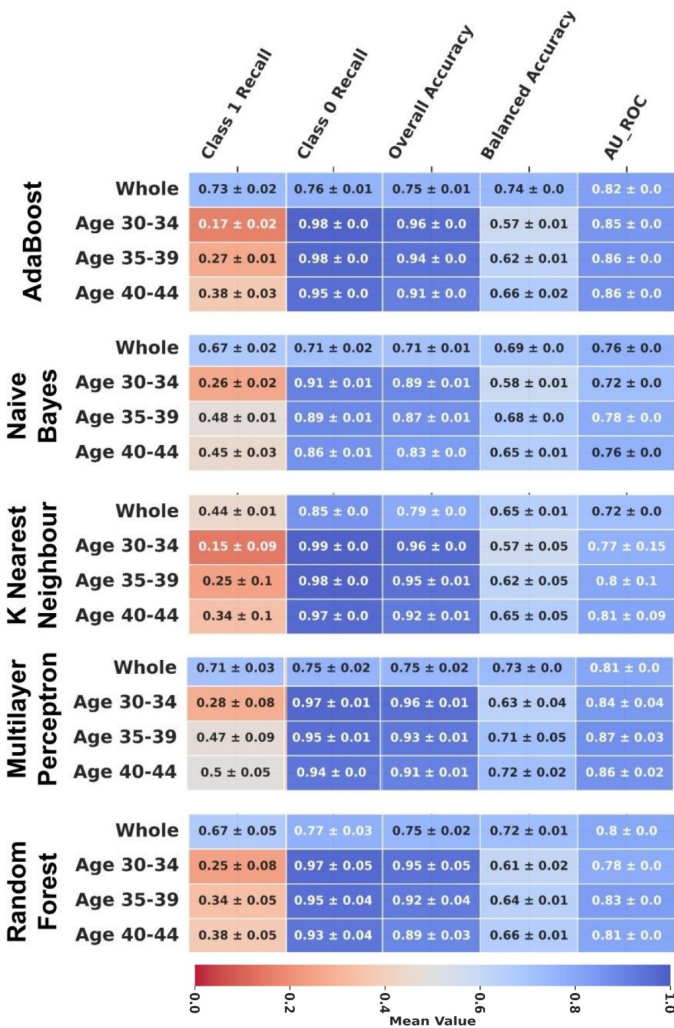
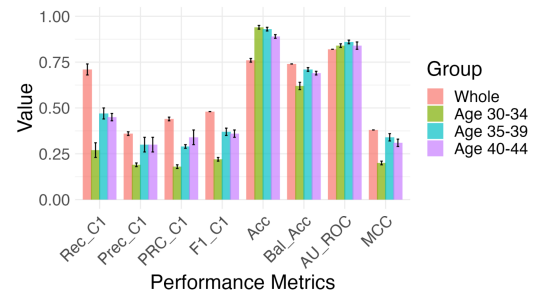


Fig. 3. Comparing whole group sampling method performance with Enhanced-DP (age group 35-39) for the young adult age groups 30-34, 35-39, and 40-44 along with the whole group. Performance scores are reported using mean +/- 95% confidence interval from multiple independent trials.

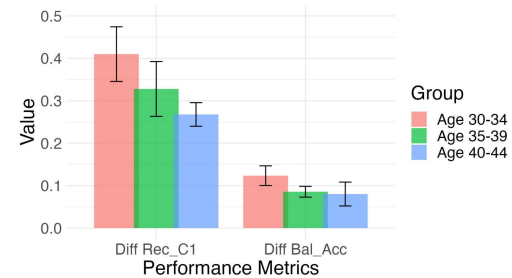
ances (Levene's test). Normality was consistently rejected, and variances were sometimes unequal. Given these violations, we applied the Kruskal-Wallis test (non-parametric ANOVA). This revealed highly significant subgroup differences across all four primary metrics ($p < 0.001$). We conducted Dunn's test with Holm correction to identify specific subgroup pairs that differed, confirming that younger adults (ages 30-44) underperformed in Class 1 Recall and AU-PRC compared to the older subgroups ($p < 0.05$).

B. Enhanced-DP model improves diagnostic accuracy

The Enhanced-DP model for three age groups improves the diabetes detection accuracy significantly (Fig. 4 (a)). Fig. 4 (b) shows the performance difference between the Enhanced-DP model and the original model in terms of positive recall and balanced accuracy. Diff Rec_C1 means subtracting the recall of class 1 of the original model from the Enhanced-DP model. The Enhanced-DP model for age groups [30-34], [35-39], and [40-44] improves the positive recall by 41% (SD 6.4%), 32%



(a) Original model performance



(b) Performance difference between DP and original model

Fig. 4. Performance of the original logistic regression model (a) when tested on the whole population and minority age group. (b) Performance difference between Enhanced-DP and original model Logistic Regression model. "Diff Rec_C1" means subtracting the recall of class 1 of the original model from the Enhanced-DP model and "Diff Bal_Acc" means subtracting the balanced accuracy of the original model from the Enhanced-DP model. Positive values indicate performance improvement from the original model. The error bars represent the standard deviation of the experiment results.

TABLE II
PERFORMANCE COMPARISON OF BASELINE XGBOOST MODEL AND ENHANCED DP XGBOOST MODEL. PERFORMANCE IS REPORTED USING AVERAGE ± 95% CI.

Model	Group	Class 1 Recall	Balanced Accuracy
Baseline XGB	Age 30-34	0.186 ± 0.078	0.582 ± 0.038
	Age 35-39	0.36 ± 0.033	0.66 ± 0.013
	Age 40-44	0.426 ± 0.038	0.677 ± 0.013
Enhanced DP XGB	Age 30-34	0.63 ± 0.072	0.73 ± 0.017
	Age 35-39	0.702 ± 0.06	0.723 ± 0.037
	Age 40-44	0.774 ± 0.06	0.751 ± 0.032

(SD 6.4%), and 24% (SD 2.8%), respectively. This means the Enhanced-DP models are better at detecting diabetes in the young population.

On the other hand, Diff Bal_Acc means subtracting the balanced accuracy of the original model from the Enhanced-DP model. The balanced accuracy (average recall of both C0 and C1 classes) is also improved by 13% (SD 2.3%), 10.5% (SD 1%), and 7.7% (SD 2.8%) for the Enhanced-DP model.

Similarly, the enhanced DP method substantially improves the detection accuracy for the XGboost model as well (Figure II).

C. Whole-population sampling

The whole group-based sampling approach doesn't improve the detection rate in the young group Fig. 5. Moreover, AdaSyn (Rec_C1 0.19), SMOTE (Rec_C1 0.19), Tomek-Link (Rec_C1

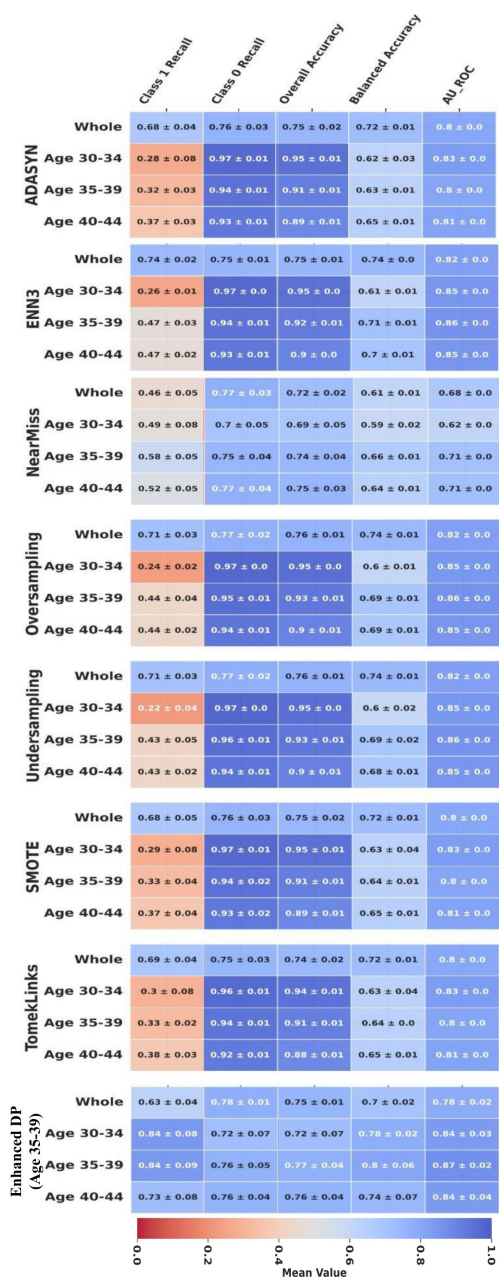


Fig. 6. Models trained on age-segmented training sets composed of age groups spanning 5 years or 15 years. The models are trained and tested on the same or overlapping age group. We utilized all 7 resampling techniques (Fig. 5) and reported the best performance for each model. The baseline model represents logistic regression model trained on the whole training set whereas DP is our proposed model.

Fig. 6. Models trained on age-segmented training sets composed of age groups spanning 5 years or 15 years. The models are trained and tested on the same or overlapping age group. We utilized all 7 resampling techniques (Fig. 5) and reported the best performance for each model. The baseline model represents logistic regression model trained on the whole training set whereas DP is our proposed model.

age group [35-39] all of the methods decrease the balanced accuracy from the original value. On the other hand, the Enhanced-DP model increases balanced accuracy by 3-9%.

D. Segmented training

We compared the performance of models trained on age-segmented datasets, where the original dataset was divided into groups with 5-year and 15-year spans. Each model was trained and tested on data from the same or overlapping age groups. For example, a model trained on individuals aged 30-34 was evaluated using a test set from the same age range. Similarly, a model trained on the 30-44 age group was tested across the 30-34, 35-39, and 40-44 age groups. Each model was calibrated using the Isotonic Regression and model-specific decision threshold was calculated. We applied all seven resampling techniques outlined in Fig. 5 and reported the best performance for each model. As shown in Fig. 6, the DP method consistently performed the best. Interestingly, models trained on 15-year spans outperformed the baseline model trained on the entire dataset. These findings were also consistent with results obtained from the BRFSS'15 data.

E. Fairness aware ML performance

To address fairness across age groups, we evaluated two standard fairness-aware ML methods: reweighting (using balanced class weights) and subgroup-specific thresholding. These experiments were conducted using logistic regression for consistency with our resampling and Enhanced-DP experiments. Reweighting did not improve recall for the 30-44 age group, which remained low, similar to baseline and resampling methods. Subgroup-specific thresholding, implemented using Microsoft's Fairlearn toolkit [50]. Performance is shown in Fig. 7 demonstrates the improvement was still very limited compared to our proposed Enhanced-DP method.

Fig. 5. Comparing whole group sampling method performance with Enhanced-DP (age group 35-39) for the young adult age groups 30-34, 35-39, and 40-44 along with the whole group. Performance scores are reported using mean +/- 95% confidence interval from multiple independent trials.

0.28), Random oversampling (Rec_C1 0.25) and undersampling (Rec_C1 0.26) methods decrease the original detection accuracy (Rec_C1 0.30) for age group [30-34]. A similar performance decrease is observed in the other two age groups [35-39] and [40-44]. Near Miss and ENN merely improve the Rec_C1 by 9% and 2%. On the other hand, the respective Enhanced-DP models (i.e., trained on age group [30-34] and applied to age group [30-34]) improve the positive recall by at least 24%.

The whole group sampling methods also show poor performance in terms of balanced accuracy. For example, for the

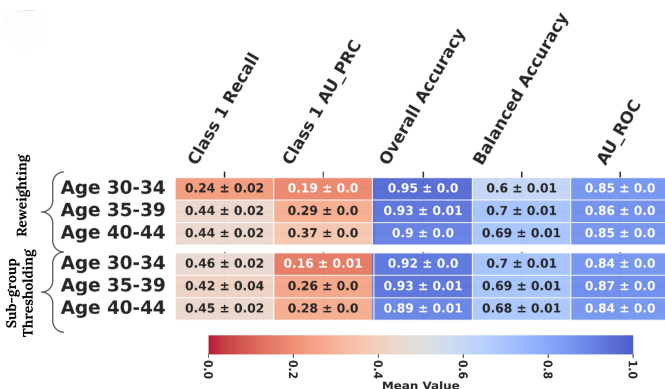


Fig. 7. Reweighting (balanced) and subgroup-specific thresholding model performance. Performance scores are reported using mean +/- 95% confidence interval from multiple independent trials.

F. Cross-group performance

Enhanced-DP model trained for a specific age group is applied to all other age groups Fig. 8 and Supplementary Fig. 10 shows the performance. For example, a DP model trained for age group [30-34], is tested on age groups [30-34], [35-39], and [40-44]. We also compare the whole group's performance with the DP models. The original model shows very poor recall C1 for all three age groups. Interestingly, the DP model trained for the age group [30-34] can also be applied to age groups [35-39] and [40-44]. DP model for age group [35-39] shows the highest recall C1 of 73% - 84% for all age groups. The balanced accuracy is also high when one DP model is applied to another age group.

On the other hand, the DP model trained for the young age group shows poor performance when applied to gender or ethnic groups for which the DP model is not trained. The recall of class 1 goes down by 8% and 19% when the DP model for the age group [40-44] is applied to the female and Asian groups, respectively. The balanced accuracy also declines if age group DP models are applied to female or Asian groups.

G. Feature analysis

One of the major motivations for selecting a logistic regression model is its interpretability. [51] showed that deep learning models are mostly black-box and cannot always be correctly interpretable.

Fig. 9 shows the logistic regression model coefficients of the original and the Enhanced-DP models. It shows that the original model had a strong correlation with the subject's race (-0.40), employment status (-0.40), and marital status (-0.54). The Enhanced-DP model reduces the strength of these attributes by at least 19%. Moreover, age was positively correlated in the original model while it was negatively correlated in the Enhanced-DP models. Other attribute coefficients show minor changes.

IV. DISCUSSION

The machine learning trained on the original dataset shows poor diabetes diagnosis performance in the young group. Because the diabetes data is limited in the original dataset the

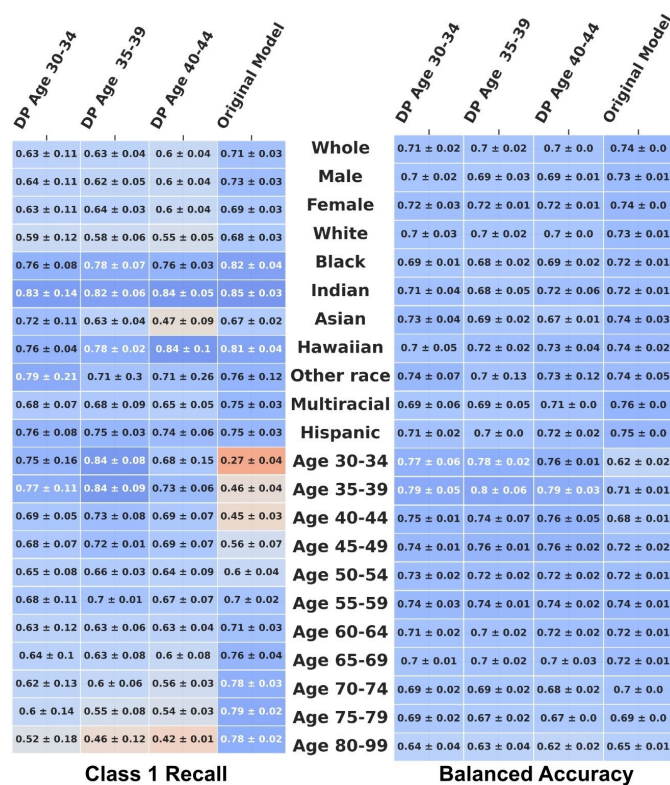


Fig. 8. Cross-group performance analysis using class 1 recall (Rec.C1) and balanced accuracy (Bal.Acc) on BRFSS 2021. Each column corresponds to an Enhanced-DP model trained for a specific subgroup. Each row represents a subgroup that a model is evaluated on. Performance scores are reported using mean +/- 95% confidence interval from multiple independent trials.

traditional machine learning model tends to pick the general statistics to build a diagnostic model. As a result, the models show poor performance in the minor young group. Whole group-based data enrichment such as sampling methods cannot overcome the problem of poor performance as it doesn't enhance the young group. Some of them such as SMOTE, a popular sampling method which is well known for removing bias from imbalanced datasets, decrease the performance. This is because these sampling methods are not equipped to reduce disparate ratios in the minority group. Moreover, we tested multiple and diverse machine learning models. However, all the models show consistently poor performance in this scenario.

We proposed an enhanced version of the double prioritized bias correction technique (DP) to make the model effective and useful for the young age group. By replicating the minor group population dynamically, the technique improves the model's performance. However, one DP model can be targeted for one particular group. The results show that the DP model trained for the young age group is not applicable for both genders or ethnic groups. This limitation can be easily overcome by using multiple DP models for each minor group. However, one of the limitations of this approach is the use of multiple models compared to a single model. We need three separate models for each of the underperforming young adult (age 30-34, 35-39, 40-44 years) subgroups and another model for the rest

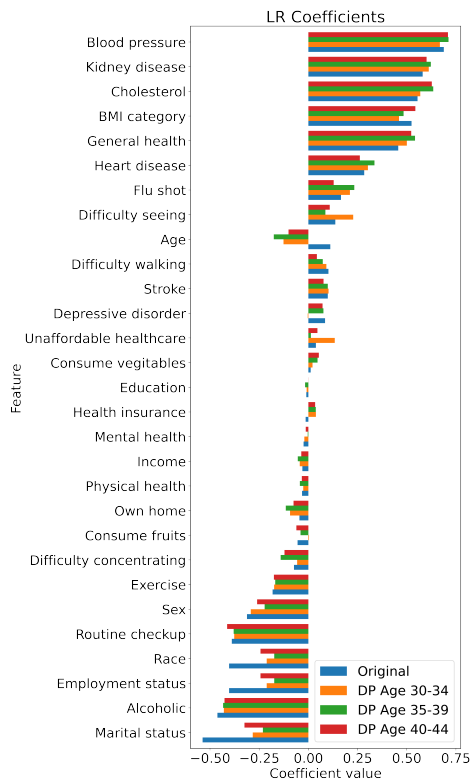


Fig. 9. Logistic regression model coefficient values which are associated with each feature from the original model and Enhanced-DP models. The one hot encoded feature (marital status, employment status, and race) coefficients are averaged

of the population. It requires more computation time to train these separate models; however, it doesn't affect the inference time. Training multiple DP models will take more time than a single model approach but the DP model will take the same time during diagnosis as only one model is going to be used for predicting a particular sample. Interestingly, this technique has broader applications beyond the specific dataset evaluated and can be utilized in any domain to overcome bias or data limitation challenges. Our proposed method also aligns with recent CDC data showing a rising incidence of type 2 diabetes among individuals under 45, highlighting the need for age-aware modeling strategies. [52].

We considered improving diabetes detection (or diagnosis) as the primary task, as recall class 1 was poor in the original model for the young adult age group, while other metrics, such as AUROC, are similar across the subgroups. Using our method, the recall of young adults improved substantially, which means the diabetes detection rate improved. We also show the balanced accuracy, which is the average of positive and negative class recall.

We also investigated the root cause of bias in the original model by visualizing the coefficients of the logistic regression model. Being a white box model, we can easily understand the feature importance and the correlation of each feature with the detection probability. The original model shows a strong correlation with non-medical attributes such as marital status, employment status, and race. The DP model decreases the strength of correlation with these factors significantly.

Moreover, the original model is positively correlated with age. It means the diabetes-positive probability increases with age. However, this positive correlation was one of the key factors why the original shows poor performance in the young age group.

We can also use real-world data in the form of the observational medical outcomes partnership (OMOP) common data model (CDM) or any other data sources. The core methodology will be the same as the DP method is model-agnostic as well as data-type agnostic. The OMOP-CDM might need to be preprocessed (e.g., relevant cohort selection, data cleaning, data scaling) before feeding to our proposed DP model. Our approach can be applied to any kind of data including OMOP-CDM.

V. CONCLUSION

We address a critical gap in the literature by developing a machine learning approach specifically designed for precision type 2 diabetes (T2D) diagnosis in a young adult population (30-44 years old). This is a new study to target this specific age group for T2D diagnosis using machine learning. We also demonstrate that multiple machine learning models and a number of resampling techniques fail to achieve fair performance in terms of detecting T2D in young adults. Our proposed bias correction technique is specifically for improving T2D diagnosis in young adults, and demonstrates the effectiveness of subgroup-focused bias correction, promoting fairer and more accurate machine learning models in healthcare settings. This approach has the potential to mitigate bias issues in the diagnosis of diabetes among young adults, offering a pathway toward more equitable and accurate healthcare practices in this demographic.

Data and code availability: The dataset used in the study is publicly available from [31] and our code is publicly available at an anonymous repository - <https://github.com/PiasTanmoy/Diabetes-BRFSS-DP>.

ACKNOWLEDGMENT

This work was partially supported by the Virginia Commonwealth Cyber Initiative (CCI) and the National Science Foundation under Grant No. 2231002.

REFERENCES

- [1] K. De Silva, W. K. Lee, A. Forbes, R. T. Demmer, C. Barton, and J. Enticott, "Use and performance of machine learning models for type 2 diabetes prediction in community settings: A systematic review and meta-analysis," *International journal of medical informatics*, vol. 143, p. 104268, 2020.
- [2] WHO, "Diabetes — who.int." <https://www.who.int/health-topics/diabetes>. [Accessed 04-23-2023].
- [3] N. H. Cho, J. Shaw, S. Karuranga, Y. Huang, J. da Rocha Fernandes, A. Ohlogge, and B. Malanda, "Idf diabetes atlas: Global estimates of diabetes prevalence for 2017 and projections for 2045," *Diabetes research and clinical practice*, vol. 138, pp. 271–281, 2018.
- [4] E. G. Krug, "Trends in diabetes: Sounding the alarm," *The Lancet*, vol. 387, no. 10027, pp. 1485–1486, 2016.
- [5] Z. Xie, O. Nikolayeva, J. Luo, and D. Li, "Building risk prediction models for type 2 diabetes using machine learning techniques," *Preventing chronic disease*, vol. 16, 2019.

- [6] N. Lascar, J. Brown, H. Pattison, A. H. Barnett, C. J. Bailey, and S. Bellary, "Type 2 diabetes in adolescents and young adults," *The lancet Diabetes & endocrinology*, vol. 6, no. 1, pp. 69–80, 2018.
- [7] CDC, "Diabetes in young people is on the rise — cdc.gov." <https://www.cdc.gov/diabetes/data-research/research/young-people-diabetes-on-rise.html>. [Accessed 10-09-2024].
- [8] CDC, "How to make the leap from type 1 teen to adult — cdc.gov." <https://www.cdc.gov/diabetes/about/type-1-teen-adult.html>. [Accessed 10-09-2024].
- [9] R. G. McCoy, R. S. Kidney, D. Holzngel, T. Peters, and V. Madzura, "Challenges for younger adults with diabetes," *Minnesota medicine*, vol. 102, no. 2, p. 34, 2019.
- [10] Diabetes, "Type 2 diabetes and young adults — diabetes.org.uk." <https://www.diabetes.org.uk/diabetes-the-basics/types-of-diabetes/type-2/young-adults>. [Accessed 12-09-2024].
- [11] G. NIH, "Serious complications from youth-onset type 2 diabetes arise by young adulthood," *National Institutes of Health*, Jul 2021.
- [12] T. A. Hillier and K. L. Pedula, "Complications in young adults with early-onset type 2 diabetes: losing the relative protection of youth," *Diabetes care*, vol. 26, no. 11, pp. 2999–3005, 2003.
- [13] G. Fremlin, S. Orpin, and M. Kaur, "Clarithromycin, rifampicin and fusidic acid triple combination therapy for chronic folliculocentric pustulosis of the scalp," *Clinical and experimental dermatology*, vol. 42, no. 8, pp. 913–914, 2017.
- [14] D. Morris, "Early-onset type 2 diabetes: Clinical implications, diagnosis and management," *Journal of Diabetes Nursing*, vol. 26, no. 6, 2022.
- [15] CDC, "Preventing type 2 diabetes — cdc.gov." <https://www.cdc.gov/diabetes/prevention-type-2>. [Accessed 10-09-2024].
- [16] Diabetes, "Diabetes prevention — ADA — diabetes.org." <https://diabetes.org/about-diabetes/diabetes-prevention>. [Accessed 10-09-2024].
- [17] W.-C. Lin, J. S. Chen, M. F. Chiang, and M. R. Hribar, "Applications of artificial intelligence to electronic health record data in ophthalmology," *Translational vision science & technology*, vol. 9, no. 2, pp. 13–13, 2020.
- [18] T. H. Davenport, T. Hongsermeier, and K. A. Mc Cord, "Using AI to improve electronic health record," *Harvard Business Review*, vol. 12, pp. 1–6, 2018.
- [19] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A survey on bias and fairness in machine learning," *ACM Computing Surveys (CSUR)*, vol. 54, no. 6, pp. 1–35, 2021.
- [20] M. Ghassemi, T. Naumann, P. Schulam, A. L. Beam, I. Y. Chen, and R. Ranganath, "A review of challenges and opportunities in machine learning for health," *AMIA Summits on Translational Science Proceedings*, vol. 200, no. 01, pp. 191–200, 2020.
- [21] S. S. Gervasi, I. Y. Chen, A. Smith-McLallen, D. Sontag, Z. Obermeyer, M. Vennera, and R. Chawla, "The potential for bias in machine learning and opportunities for health insurers to address it: Article examines the potential for bias in machine learning and opportunities for health insurers to address it.," *Health Affairs*, vol. 41, no. 2, pp. 212–218, 2022.
- [22] I. Chen, F. D. Johansson, and D. Sontag, "Why is my classifier discriminatory?," *Advances in neural information processing systems*, vol. 31, 2018.
- [23] T. S. Pias, S. Afrose, M. D. Tuli, I. H. Trisha, X. Deng, C. B. Nemeroff, and D. D. Yao, "Low responsiveness of machine learning models to critical or deteriorating health conditions," *Communications Medicine*, vol. 5, no. 1, p. 62, 2025.
- [24] I. Y. Chen, *Machine Learning Approaches for Equitable Healthcare*. PhD thesis, Massachusetts Institute of Technology, 2022.
- [25] L. Seyyed-Kalantari, G. Liu, M. McDermott, I. Y. Chen, and M. Ghassemi, "Chexclusion: Fairness gaps in deep chest x-ray classifiers," in *BIOCOMPUTING 2021: proceedings of the Pacific symposium*, pp. 232–243, World Scientific, 2020.
- [26] I. Y. Chen, P. Szolovits, and M. Ghassemi, "Can AI help reduce disparities in general medical and mental health care?," *AMA journal of ethics*, vol. 21, no. 2, pp. 167–179, 2019.
- [27] L. Seyyed-Kalantari, H. Zhang, M. B. McDermott, I. Y. Chen, and M. Ghassemi, "Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations," *Nature medicine*, vol. 27, no. 12, pp. 2176–2182, 2021.
- [28] C. H. Chu, S. Donato-Woodger, S. S. Khan, R. Nyrup, K. Leslie, A. Lyn, T. Shi, A. Bianchi, S. A. Rahimi, and A. Grenier, "Age-related bias and artificial intelligence: A scoping review," *Humanities and Social Sciences Communications*, vol. 10, no. 1, pp. 1–17, 2023.
- [29] W. Tan, Q. Wei, Z. Xing, H. Fu, H. Kong, Y. Lu, B. Yan, and C. Zhao, "Fairer AI in ophthalmology via implicit fairness learning for mitigating sexism and ageism," *Nature Communications*, vol. 15, no. 1, p. 4750, 2024.
- [30] Q. Xie, E. J. Schenck, H. S. Yang, Y. Chen, Y. Peng, and F. Wang, "Faithful AI in healthcare and medicine," *medRxiv*, pp. 2023–04, 2023.
- [31] CDC, "CDC - BRFSS — cdc.gov." <https://www.cdc.gov/brfss/>. [Accessed 23-Apr-2023].
- [32] G. S. Collins, S. Mallett, O. Omar, and L.-M. Yu, "Developing risk prediction models for type 2 diabetes: A systematic review of methodology and reporting," *BMC medicine*, vol. 9, no. 1, pp. 1–14, 2011.
- [33] P. W. Sullivan, E. H. Morrato, V. Ghushchyan, H. R. Wyatt, and J. O. Hill, "Obesity, inactivity, and the prevalence of diabetes and diabetes-related cardiovascular comorbidities in the us, 2000–2002," *Diabetes care*, vol. 28, no. 7, pp. 1599–1603, 2005.
- [34] D. Noble, R. Mathur, T. Dent, C. Meads, and T. Greenhalgh, "Risk models and scores for type 2 diabetes: Systematic review," *Bmj*, vol. 343, 2011.
- [35] S. Afrose, W. Song, C. B. Nemeroff, C. Lu, and D. Yao, "Subpopulation-specific machine learning prognosis for underrepresented patients with double prioritized bias correction," *Communications medicine*, vol. 2, no. 1, p. 111, 2022.
- [36] J. Van Hulse, T. M. Khoshgoftaar, and A. Napolitano, "Experimental perspectives on learning from imbalanced data," in *Proceedings of the 24th international conference on Machine learning*, pp. 935–942, 2007.
- [37] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [38] H. He, Y. Bai, E. A. Garcia, and S. Li, "ADASYN: Adaptive synthetic sampling approach for imbalanced learning," in *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*, pp. 1322–1328, Ieee, 2008.
- [39] I. Tomek, "Two modifications of CNN," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-6, no. 11, pp. 769–772, 1976.
- [40] D. L. Wilson, "Asymptotic properties of nearest neighbor rules using edited data," *IEEE Transactions on Systems, Man, and Cybernetics*, no. 3, pp. 408–421, 1972.
- [41] I. Mani and I. Zhang, "kNN approach to unbalanced data distributions: A case study involving information extraction," in *Proceedings of workshop on learning from imbalanced datasets*, vol. 126, pp. 1–7, ICML, 2003.
- [42] M. Carvalho, A. J. Pinho, and S. Brás, "Resampling approaches to handle class imbalance: a review from a data perspective," *Journal of Big Data*, vol. 12, no. 1, p. 71, 2025.
- [43] S. F. Taskiran, B. Turkoglu, E. Kaya, and T. Asuroglu, "A comprehensive evaluation of oversampling techniques for enhancing text classification performance," *Scientific Reports*, vol. 15, no. 1, p. 21631, 2025.
- [44] W. Chen, K. Yang, Z. Yu, Y. Shi, and C. P. Chen, "A survey on imbalanced learning: latest research, applications and future directions," *Artificial Intelligence Review*, vol. 57, no. 6, p. 137, 2024.
- [45] Z. Zhao, T. Cui, S. Ding, J. Li, and A. G. Bellotti, "Resampling techniques study on class imbalance problem in credit risk prediction," *Mathematics*, vol. 12, no. 5, p. 701, 2024.
- [46] M. B. McDermott, S. Wang, N. Marinsek, R. Ranganath, L. Foschini, and M. Ghassemi, "Reproducibility in machine learning for health research: Still a ways to go," *Science Translational Medicine*, vol. 13, no. 586, p. eabb1655, 2021.
- [47] I. Y. Chen, E. Pierson, S. Rose, S. Joshi, K. Ferryman, and M. Ghassemi, "Ethical machine learning in healthcare," *Annual review of biomedical data science*, vol. 4, pp. 123–144, 2021.
- [48] B. Zadrozny and C. Elkan, "Transforming classifier scores into accurate multiclass probability estimates," in *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 694–699, 2002.
- [49] L. B. Iantovics, M. Dehmer, and F. Emmert-Streib, "Metrintsimil—an accurate and robust metric for comparison of similarity in intelligence of any number of cooperative multiagent systems," *Symmetry*, vol. 10, no. 2, 2018.
- [50] S. Bird, M. Dudík, R. Edgar, B. Horn, R. Lutz, V. Milan, M. Sameki, H. Wallach, and K. Walker, "Fairlearn: A toolkit for assessing and improving fairness in ai," *Microsoft, Tech. Rep. MSR-TR-2020-32*, 2020.
- [51] M. Ghassemi, L. Oakden-Rayner, and A. L. Beam, "The false hope of current approaches to explainable artificial intelligence in health care," *The Lancet Digital Health*, vol. 3, no. 11, pp. e745–e750, 2021.
- [52] CDC, "Centers for disease control and prevention - national diabetes statistics report." <https://www.cdc.gov/diabetes/php/data-research/index.html>. [Accessed 5-7-2025].