

# Probabilistic Versus Deep Generative Models: A Fairness Centred Evaluation of Synthetic Healthcare Tabular Data

Received: 8 September 2025

Accepted: 15 January 2026

Published online: 26 February 2026

Cite this article as: Alattal D., Draghi B., Myles P. *et al.* Probabilistic Versus Deep Generative Models: A Fairness Centred Evaluation of Synthetic Healthcare Tabular Data. *Int J Comput Intell Syst* (2026). <https://doi.org/10.1007/s44196-026-01173-7>

Dima Alattal, Barbara Draghi, Puja Myles, Richard Branson & Allan Tucker

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

ARTICLE IN PRESS

## Probabilistic Versus Deep Generative Models: A Fairness Centred Evaluation of Synthetic Healthcare Tabular Data

This Accepted Manuscript (AM) is a PDF file of the manuscript accepted for publication after peer review, when applicable, but does not reflect post-acceptance improvements, or any corrections. Use of this AM is subject to the publisher's embargo period and AM terms of use. Under no circumstances may this AM be shared or distributed under a Creative Commons or other form of open access license, nor may it be reformatted or enhanced, whether by the Author or third parties. By using this AM (for example, by accessing or downloading) you agree to abide by Springer Nature's terms of use for AM versions of subscription articles: <https://www.springernature.com/gp/open-research/policies/accepted-manuscript-terms>

The Version of Record (VOR) of this article, as published and maintained by the publisher, is available online at: <https://doi.org/10.1007/s44196-026-01173-7>. The VOR is the version of the article after copy-editing and typesetting, and connected to open research data, open protocols, and open code where available. Any supplementary information can be found on the journal website, connected to the VOR.

For research integrity purposes it is best practice to cite the published Version of Record (VOR), where available (for example, see ICMJE's guidelines on overlapping publications). Where users do not have access to the VOR, any citation must clearly indicate that the reference is to an Accepted Manuscript (AM) version.

Probabilistic vs Deep Generative Models: A  
 Fairness Centred Evaluation of Synthetic  
 Healthcare Tabular Data

Dima Alattal<sup>1</sup>, Barbara Draghi<sup>2</sup>, Puja Myles,  
 Richard Branson<sup>3,4†</sup>, Allan Tucker<sup>5\*</sup>

<sup>1,2,5\*</sup>Computer Science Department, Brunel University London,  
 London, UK.

<sup>3,4</sup>Medicine and Healthcare products Regulatory Agency, London, UK.

\*Corresponding author(s). E-mail(s): [allan.tucker@brunel.ac.uk](mailto:allan.tucker@brunel.ac.uk);  
 Contributing authors: [dima.alattal2@brunel.ac.uk](mailto:dima.alattal2@brunel.ac.uk);  
[barbara.draghi@brunel.ac.uk](mailto:barbara.draghi@brunel.ac.uk); [puja.myles@mhra.gov.uk](mailto:puja.myles@mhra.gov.uk),  
[richard.branson@mhra.gov.uk](mailto:richard.branson@mhra.gov.uk);

†These authors contributed equally to this work.

#### Abstract

Synthetic data offers a promising avenue for addressing privacy, scarcity, and fairness challenges in healthcare datasets. However, there is limited evaluation of how different generation methods balance fidelity, utility, and fairness, particularly for underrepresented subgroups. This study addresses this gap by comparing representative generative modelling techniques, both probabilistic and deep approaches, that are popular in the research literature. We empirically evaluate BayesBoost, CTGAN, TVAE, CopulaGAN, and DECAF on two healthcare datasets containing numerical, binary, and categorical features. Each model's performance is assessed along three axes: data fidelity, machine learning utility, and fairness, using Accuracy Parity, Equalised Odds, and Predictive Rate Parity. Results show that BayesBoost consistently achieved superior fidelity, utility, and fairness preservation, particularly when paired with Random Forest classifiers, achieving around *60–63% higher* downstream utility than GAN-based deep generative baselines (e.g., Random Forest accuracy up to *0.88* with BayesBoost versus *0.54–0.55* for GAN-based methods). Deep generative models, while

001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046

047 effective in capturing complex structures, often degraded fairness, especially for  
048 underrepresented groups, with equalised odds deviating by over *100%* from the  
049 ideal parity value of *1.0* in some settings. The Variational Autoencoder out-  
050 performed other deep generative models in fairness preservation, especially for  
051 equalised odds, although with some reduction in fidelity and utility. Overall,  
052 these findings suggest that synthetic data generation for healthcare must move  
053 beyond fidelity evaluations to explicitly assess fairness and subgroup impacts,  
054 with probabilistic models such as BayesBoost showing strong potential for eth-  
055 ical deployment, while deep generative models require further adaptation for  
056 fairness-sensitive applications.

057 **Keywords:** Synthetic Data Generation, Tabular Data, Fairness in Machine Learning,  
058 Healthcare Data, Generative Models, Data Fidelity, Bias Mitigation, BayesBoost,  
059 GAN, VAE

060

061

062

063

064

065

## 1 Introduction

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090

091

092

Synthetic tabular data generation has emerged as a promising solution to critical challenges in healthcare, including data privacy, scarcity, and bias mitigation [1]. The growing digitalisation of healthcare and the adoption of electronic health records have created new opportunities for data-driven research but have also intensified concerns around privacy, security, and data-sharing.

Traditional privacy-preserving approaches, such as de-identification and anonymisation, aim to protect sensitive information by removing direct identifiers and introducing random noise [2, 3]. However, these techniques can degrade data fidelity, compromise analytical validity, and risk introducing new biases into downstream analyses.

An alternative that can potentially overcome these limitations involves the generation of fully synthetic data as a substitute for real datasets. Synthetic data generation produces artificial records that reproduce the statistical structure of the original data while avoiding the disclosure of identifiable information. In healthcare, where most information is stored in structured tabular form describing patients through mixed numerical, categorical, and binary variables, generating realistic synthetic

data requires capturing complex conditional dependencies and maintaining clinical  
plausibility.

Previous studies have explored the generation of synthetic patient-level datasets  
from aggregated population statistics [4, 5], balancing confidentiality and statisti-  
cal representativeness. Probabilistic models, particularly Bayesian Networks (BNs),  
have demonstrated success in replicating complex clinical associations and joint  
distributions [6, 7]. Building on this foundation, BayesBoost integrates boosting tech-  
niques to better address representation bias and enhance synthetic data quality for  
underrepresented groups [8].

Recent advances in machine learning have expanded synthetic data generation  
capabilities further. Deep generative models such as Generative Adversarial Networks  
(GANs) and Variational Autoencoders (VAEs) [9, 10] offer powerful tools for capturing  
non-linear feature interactions without requiring explicit distributional assumptions.  
Models like CTGAN, TVAE and CopulaGAN have been adapted to tabular data,  
demonstrating flexibility in complex data settings [11]. Nevertheless, challenges persist.  
Tabular healthcare datasets often contain a mix of continuous, binary, and categorical  
features with high sparsity, imbalance, and complex dependencies, which generative  
models struggle to replicate faithfully.

Fairness in synthetic data generation is a significant concern, particularly within  
healthcare, where biased data can result in unequal outcomes for various patient  
groups. Ensuring fairness involves not only achieving equitable representation of sub-  
groups within the synthetic data but also guaranteeing fair and comparable outcomes  
for underrepresented groups in machine learning models trained on this data. How-  
ever, this is a challenging task, as generative models may unintentionally reproduce,  
amplify biases present in the original datasets or even introduce new biases. To address  
these issues, robust evaluation frameworks are essential, emphasising both fairness and  
ethical considerations.

139 This study addresses these gaps through an empirical evaluation of five syn-  
140 thetic data generation methods—BayesBoost, CTGAN, TVAE, CopulaGAN, and  
141 DECAF—on healthcare datasets comprising numerical, binary, and categorical fea-  
142 tures for binary classification tasks. We systematically assess how these methods  
143 preserve statistical fidelity, predictive utility, and fairness, with a particular focus on  
144 underrepresented groups.

149 While previous surveys have evaluated synthetic data generation techniques with  
150 respect to fidelity, utility, and privacy [9, 12–16], few have explicitly addressed fairness  
151 preservation, especially in healthcare. Moreover, no previous studies have compared  
152 deep generative models against BayesBoost, a probabilistic approach specifically  
153 designed to mitigate bias in tabular data generation. While DECAF and BayesBoost  
154 incorporate bias mitigation mechanisms, CTGAN, TVAE and CopulaGAN are primar-  
155 ily designed to capture complex tabular structures without explicit fairness objectives.  
156 By evaluating fidelity, machine learning utility, and fairness in tandem, our study pro-  
157 vides a comprehensive, comparative assessment of the capabilities and limitations of  
158 these approaches. The findings emphasise that fairness-conscious design, structured  
159 evaluation frameworks, and transparent reporting must be prioritised to ensure the  
160 ethical deployment of synthetic data systems in healthcare applications.

169 In summary, this paper makes three main contributions:

170 **1. Comprehensive evaluation framework:** We introduce an integrated assess-  
171 ment framework to jointly evaluate *fidelity*, *predictive utility*, and *fairness* in synthetic  
172 tabular data, enabling systematic comparison across modelling paradigms.

173 **2. Empirical comparison of probabilistic and deep generative models:**  
174 We conduct a rigorous experimental analysis of five representative approaches (Bayes-  
175 Boost, CTGAN, TVAE, CopulaGAN, and DECAF) using healthcare datasets with  
176 mixed numerical, categorical, and binary features.

177  
178  
179  
180  
181  
182  
183  
184

**3. Insights on fairness preservation and ethical deployment:** We provide new evidence on how probabilistic and deep models differ in their ability to preserve subgroup fairness, highlighting the strengths and limitations of the evaluated approaches for fairness-sensitive healthcare applications.

## 1.1 Related Work

Synthetic data holds considerable value in healthcare, offering a range of significant benefits that address both practical and ethical challenges. By augmenting existing datasets, synthetic data enhances the performance of machine learning models through increased data volume and diversity, which is particularly beneficial when dealing with rare diseases or underrepresented patient populations. It also improves data accessibility, enabling researchers to overcome limitations related to data sharing restrictions. Beyond these advantages, synthetic data provides a powerful privacy-preserving solution by mimicking the statistical patterns of real data while excluding any personally identifiable information [17].

In addition, synthetic health records play a pivotal role in in-silico clinical trials, which use patient-specific models to generate virtual cohorts for evaluating the safety and effectiveness of new drugs and medical devices [18]. By using synthetic data in such trials, researchers can also refine trial designs and make predictions at both population and individual levels, thus improving the chances of success [19].

### 1.1.1 Probabilistic Generative Models

BayesBoost, a type of probabilistic data generation method, addresses biases in the original dataset by targeting difficult-to-predict or under-represented samples. This is particularly important as biases in real data often propagate into synthetic data, potentially compromising its utility and fairness. BayesBoost identifies these under-represented instances and employs a boosting approach to mitigate the issue. Boosting, a machine learning technique that combines multiple weak models to produce a

231 stronger one, is used here to over-sample the under-represented cases. This process  
232 ensures a more balanced and representative distribution of subgroups in the synthetic  
233 dataset. The synthetic data generation step in BayesBoost relies on Bayesian networks  
234 to model the relationships and distributions in the original data. Bayesian networks are  
235 probabilistic graphical models that use directed acyclic graphs to represent variables  
236 and their conditional dependencies [8].

241 Bayesian networks are widely recognised for their capability to model probabilis-  
242 tic relationships and perform inference tasks, allowing the representation of causal  
243 dependencies within data. These characteristics make them particularly valuable for  
244 understanding variable interactions and incorporating uncertainty into predictive mod-  
245 els. Additionally, their flexibility in handling non-linear data patterns enhances their  
246 applicability in domains where decision-making and risk assessment are critical. How-  
247 ever, the application of Bayesian networks to large-scale, high-dimensional datasets  
248 and complex distributions presents significant challenges. Learning intricate depen-  
249 dencies in such contexts is computationally demanding and often requires complex  
250 structuring methods [13].

### 258 **1.1.2 Deep Learning Generative Models**

261 Deep learning-based data generation methods, particularly GANs, have undergone  
262 significant evolution to address the unique challenges posed by various data types,  
263 including tabular data. Vanilla GANs, the foundational GAN architecture, consist of  
264 a generator and discriminator that engage in an adversarial framework to create syn-  
265 thetic data. While effective for image data, Vanilla GANs face limitations when applied  
266 to tabular data due to the complexity of modelling mixed data types (categorical and  
267 numerical) and the intricate dependencies inherent in such datasets [20].

272 To overcome these challenges, specialized GAN variants for tabular data have been  
273 introduced. One notable example is Conditional GANs (CTGANs), which incorporate  
274 specific techniques like mode-specific normalization to better handle the characteristics  
275 and their conditional dependencies [8].

of tabular datasets. CTGANs enable conditional data generation, allowing synthetic data to be generated based on specific labels or features. This feature is particularly advantageous in structured data settings, where controlling the generation process based on attributes such as disease type or patient demographics enhances the realism and practical utility of the synthetic data [21].

Copula GANs represent another advancement in tabular data synthesis, using copula theory to model interdependencies among variables. This approach allows for more accurate capture of relationships in high-dimensional data, particularly when dealing with mixed data types (continuous and categorical), thereby offering improvements over Vanilla GANs [21].

DECAF takes a different approach by prioritizing fairness and bias mitigation in data generation. DECAF explicitly incorporates the data-generating process into its framework as a structural causal model embedded within the generator's input layers. By reconstructing each variable based on its causal parents, DECAF ensures that the generated data faithfully reflects the underlying causal relationships in the original dataset. This feature enables DECAF to produce synthetic data that minimizes biases while preserving critical causal structures [12].

VAEs, another variation of deep generative models, are based on the autoencoder architecture, comprising an encoder and a decoder. The encoder transforms input data into a compressed, continuous latent space represented by probabilistic latent variables. The decoder then reconstructs the data from this latent space, enabling VAEs to generate samples that align with the original data distribution. Unlike traditional autoencoders, VAEs utilise probabilistic encoding, which maps input data into distributions rather than fixed points in the latent space. This probabilistic framework facilitates smooth transitions between generated samples, making VAEs particularly effective in producing diverse synthetic datasets that match the distribution of the original data.

323 However, VAEs, like other deep generative models such as GANs, face several lim-  
324 itations. They can be computationally intensive, particularly when applied to large  
325 and complex datasets. One significant challenge is model collapse, where the generator  
326 fails to capture the full variability of the data, resulting in limited sample diversity and  
327 repetitive outputs. Such issues undermine the model's ability to represent the com-  
328 plexity of real-world data effectively. Moreover, the quality and fidelity of synthetic  
329 data generated by VAEs depend heavily on the adequacy and balance of the training  
330 data. If the training dataset is insufficient, biased, or non-representative, the syn-  
331 thetic data may deviate from the true underlying distribution of the original dataset.  
332 This discrepancy compromises both the fidelity and the utility of the synthetic data,  
333 particularly in sensitive applications.  
334  
335  
336  
337  
338  
339  
340

### 341 342 **1.1.3 Fairness and Bias Mitigation in Synthetic Data** 343

344 While generative models hold significant promise for advancing the healthcare system,  
345 they also raise ethical concerns regarding their generalisability across diverse popula-  
346 tions. These concerns are magnified when data is skewed or underrepresented, leading  
347 to harmful biases related to age, race, or gender. In such cases, synthetic data genera-  
348 tion can serve as a potential mitigation strategy by addressing these gaps and ensuring  
349 better representation within datasets [17].  
350  
351  
352

353 Fairness in AI systems is a major concern, as an AI system may exhibit unfair  
354 behaviour when it extends or withholds opportunities, resources, or information  
355 unequally across different groups. Furthermore, disparities in quality of service—where  
356 the system performs well for one group but inadequately for another—can exacerbate  
357 existing inequities. In healthcare, such biases can significantly impact treatment out-  
358 comes, leading to unequal access and care for specific racial, gender, or demographic  
359 groups. These disparities are not just inequitable but can result in life-threatening  
360 consequences.  
361  
362  
363  
364  
365  
366  
367  
368

Biases in AI systems often stem from multiple sources, including the use of imbalanced or underrepresented data to train generative models, inherent algorithmic biases, and the ways in which generative and predictive AI models are deployed in practice [22, 23]. Biases in the original data are frequently carried over to synthetic datasets generated by AI models, perpetuating the inequities in subsequent analyses. This can lead to predictive machine learning models trained on synthetic data that exhibit a lack of fairness, ultimately causing harm to individuals. For instance, biased models may result in misdiagnoses, inappropriate treatment plans, or unequal access to vital medical resources. Addressing these challenges requires targeted efforts to mitigate biases at every stage of the AI pipeline from data collection and generation to model training and deployment to ensure equitable and ethical healthcare practices [23].

Addressing biases in synthetic data generation can involve multiple strategies, including eliminating biases in the original datasets, balancing data during the generation process, or adapting generative models to account for existing biases. One such approach, BayesBoost, is designed to tackle biases from a representation perspective by applying boosting techniques to underrepresented sub-populations identified through model performance metrics [8]. By focusing on these sub-populations, BayesBoost aims to create a more balanced synthetic dataset. However, its effectiveness and consistency in producing unbiased and informative results require further empirical validation. Decaf, approaches the issue differently by incorporating causal frameworks. Instead of balancing the initial dataset, Decaf employs causal graphs to disconnect the direct relationship between sensitive attributes (e.g., gender, race) and the target variable while preserving the sensitive attributes in the data. This ensures that the generated synthetic data does not encode direct causal relationships between sensitive attributes and the target variable, thereby preventing sensitive attributes from influencing classification outcomes. This approach allows for maintaining fairness without

415 the need to exclude sensitive variables entirely, offering a nuanced way of addressing  
416 biases in synthetic data generation [12].  
417

418 Fairness in machine learning has been a key focus of numerous studies, with par-  
419 ticular attention given to developing measures that identify biases within datasets  
420 and model outcomes. These measures assess fairness either across demographic groups  
421 or at an individual level. For group-based fairness, prominent metrics include Demo-  
422 graphic Parity (DP) [24], Predictive Rate Parity (PRP) [25], Equalized Odds (EO),  
423 Equal Opportunity [26], and Accuracy Equality [27]. These measures evaluate whether  
424 different demographic groups, such as those defined by race, gender, or age, have  
425 an equal opportunities of receiving favorable outcomes—for instance, being correctly  
426 diagnosed in healthcare scenarios.  
427

428 On the individual level, fairness measures assess whether a model provides con-  
429 sistent outputs for individuals differing only in sensitive attributes, such as gender or  
430 ethnicity [24]. This perspective considers an algorithm to be fair if its predictions are  
431 independent of such attributes when all other factors are identical. More extensive  
432 fairness definitions and measures have been formalized and reviewed extensively in  
433 [27].  
434

## 445 **2 Experiments & Methods**

446 This section outlines the experimental framework established to evaluate five state-  
447 of-the-art synthetic tabular data generation methods: the deep generative models  
448 CopulaGAN, CTGAN, DECAF, and VAE, as well as the probabilistic model Bayes-  
449 Boost. The primary focus is to assess these models' performance in terms of fidelity,  
450 utility, and fairness, particularly in healthcare datasets characterised by mixed data  
451 types (numerical, binary, and categorical) and underrepresented sensitive attributes.  
452  
453  
454  
455  
456  
457  
458  
459  
460

We detail the train-test splitting strategy, hardware and software configurations, data preprocessing steps, and model-specific configurations adopted in the experiments.

## 2.1 Datasets

Our experiments employed two tabular datasets, each designed for binary classification tasks within the healthcare domain. The first dataset, derived from synthetic CPRD primary care data [28], is a high-fidelity synthetic dataset focusing on cardiovascular disease risk factors. It includes variables such as smoking behaviour, age, and chronic conditions associated with cardiovascular health. The version used comprises 10,000 synthetic individuals, randomly sampled from the synthetic CPRD data, with a binary target variable indicating the occurrence or absence of a heart attack.

The second dataset is a publicly available resource from Kaggle, focusing on diabetes diagnosis as a binary target variable. It comprises detailed health records for 1,879 patients, including demographic information, lifestyle factors, and medical history [29]. A more detailed description of both datasets is provided in Table 1.

Both datasets were selected specifically for their relevance to subgroup fairness analysis, given their inherent representation biases in sensitive attributes and the variation in subgroup-specific fairness metrics. Notably, they range from datasets exhibiting negligible fairness concerns to those presenting clearer disparities. While maintaining reasonable accuracy parity across demographic subgroups, both datasets show moderate or minimal fairness issues regarding Predictive Rate Parity and Equalised Odds, as detailed in Tables 4 and 5. These characteristics make them particularly suitable for evaluating the ability of synthetic data generation models to preserve, degrade, or enhance fairness relative to the real data.

**Table 1:** Dataset descriptions used in the experiments

Data	CVD	Diabetes
Number of instances	10,000	1,879
Number of features	22	44
Numerical features	5	21
Binary features	13	19
Categorical features	4	4
Sensitive attribute	Ethnicity (6 groups) <sup>1</sup>	Ethnicity (4 groups) <sup>2</sup>
Protected group	Group 0	Group 2
Target variable	Heart attack	Diabetes diagnosis

<sup>1</sup>CVD dataset sensitive attribute groups: 0: Asian; 1: Black; 2: Mixed; 3: Other; 4: Unknown; 5: White.

<sup>2</sup>Diabetes dataset sensitive attribute groups: 0: Caucasian; 1: African American; 2: Asian; 3: Other.

## 2.2 Experimental Settings

Among the most widely used deep learning generative models for tabular data, we evaluated several state-of-the-art approaches extensively explored in recent studies. These models include CopulaGAN, CTGAN, TVAE, and DECAF. In addition, we compared these deep generative models against BayesBoost, a probabilistic approach proposed as a potential solution to address bias issues in synthetic tabular data generation.

To mitigate the risk of information leakage, the ground truth datasets were partitioned into an 80/20 train-test split. The synthetic data generation process was performed exclusively on the training set, ensuring complete independence of the test set for subsequent evaluation. Test data were randomly sampled to preserve the distributions of both the sensitive attribute and the target variable, ensuring adequate subgroup representation for fairness assessment.

Each generative model produced synthetic datasets across 150 independent iterations, each utilising a distinct random seed. Performance metrics were averaged across these iterations to account for variability and to ensure robust and reliable evaluation results.

Various libraries and model-specific implementations were employed to train and evaluate the models. BayesBoost was developed in RStudio using the `bnlearn` library,

while the Synthetic Data Vault (SDV) library was used to implement CTGAN, CopulaGAN, and TVAE models. The original implementation was used for the DECAF model. To ensure consistency across experiments, default settings were adopted for batch size and number of training epochs across all deep learning models. Given the computational demands and complexity associated with hyperparameter tuning in deep generative models, all hyperparameters were kept at their default values, recognising that this choice may not exhaustively explore the parameter search space.

All datasets underwent comprehensive preprocessing to ensure consistency and comparability prior to training the generative models. For BayesBoost, which requires categorical data, all features were discretised before model training. To maintain uniformity in model comparisons, a standardised preprocessing approach was applied across all deep generative models. Categorical features were initially converted into numerical values using the LabelEncoder from the scikit-learn library. Following this, both categorical and numerical features were scaled to the range  $[0, 1]$  using the MinMaxScaler, also from scikit-learn. The MinMaxScaler was selected for data transformation after preliminary trials demonstrated the best performance compared to alternative scaling methods.

To prevent data leakage, the scaler was fitted exclusively on the training data and then applied to both training and test sets. This scaling procedure was consistently repeated for each model across all 150 iterations, with a new scaler implemented in each iteration. The generated data was subsequently reverse-transformed to its original format to enable accurate data fidelity comparisons.

For the DECAF model, we generated the direct causal graph using the PC algorithm in GENie Modeler and applied the no-debias (Decaf ND) approach [30]. This method infers the causal relationships within GANs without removing the causal edges between the target variable and sensitive attributes. We specifically chose not to remove any causal relationships between ethnicity and health outcomes, hence not

553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598

599 employing other DECAF variations such as DECAF FTU, DECAF DP, or DECAF  
600 CU. This choice was based on the critical role that ethnicity plays in healthcare diagno-  
601 sis; by preserving the influence of ethnicity on medical outcomes, we aimed to support  
602 fair and accurate predictions across all ethnic groups, rather than eliminating these  
603 relationships.  
604  
605  
606

607 For each synthetic dataset generated, metrics for fidelity, utility, and fairness were  
608 calculated to comprehensively evaluate the performance of each generative model.  
609 After computing these metrics for each individual dataset, the average scores across  
610 the 150 iterations were used to assess each model's overall performance in terms of the  
611 selected evaluation metrics, yielding a balanced and robust assessment of each model's  
612 capabilities.  
613  
614  
615  
616

617 All experiments were performed on a system equipped with an Intel Evo CPU @  
618 3.80 GHz, 32 GB RAM, and an NVIDIA GeForce RTX 4060 GPU with 16 GB of  
619 memory. Python (version 3.12) was used for the majority of models, while R Studio  
620 was utilised for the BayesBoost experiments.  
621  
622  
623  
624

## 625 **2.3 Evaluation Metrics**

### 626 **2.3.1 Data Fidelity**

627 Data fidelity reflects the extent to which synthetic data accurately replicates the  
628 statistical characteristics and patterns observed in the original dataset. To evaluate  
629 fidelity, we conducted both bivariate and multivariate analyses, assessing categorical  
630 and numerical features using a range of complementary metrics.  
631  
632  
633  
634  
635

636 For bivariate analysis, separate metrics were applied to categorical and numeri-  
637 cal variables. In the case of categorical variables, *Category Coverage* assessed whether  
638 the synthetic data captured all unique categories present in the real data. High cat-  
639 egory coverage indicates that the synthetic data retains the diversity of categorical  
640  
641  
642  
643  
644

values, ensuring representativeness across all subgroups. Additionally, the *TV Complement*, based on the total variation distance, quantified the similarity between the distributions of categorical values in the real and synthetic datasets. A higher TV complement score suggests closer alignment between the two distributions, indicating superior fidelity.

For numerical variables, three metrics were employed. *Boundary Adherence* evaluated whether the synthetic data respected the minimum and maximum values observed in each real data column, ensuring the absence of unrealistic or out-of-bound values. *Range Coverage* assessed whether the synthetic data spanned the full range of real data values, thus preserving the variability and representativeness of numerical features. Finally, the *KS Complement*, derived from the Kolmogorov-Smirnov statistic, measured the similarity between the marginal distributions of numerical variables in the real and synthetic datasets. Higher KS complement values indicate better alignment and therefore higher fidelity.

To ensure robustness, all fidelity scores were averaged across the 150 synthetic datasets generated for each model. This averaging mitigated variability introduced by random seed selection and provided a more reliable measure of model performance.

For multivariate fidelity assessment, we evaluated the ability of synthetic datasets to preserve inter-variable relationships by computing the correlation matrices of numerical features for both real and synthetic datasets. Each synthetic dataset, generated across 150 iterations per model, was compared with the real data correlation matrix using the Frobenius norm distance. This metric quantifies the dissimilarity between two correlation matrices, where lower Frobenius distances indicate greater similarity to the real data structure, and higher values suggest greater divergence. To enhance reliability, Frobenius distances were averaged across all iterations for each model, providing a comprehensive measure of each model's ability to preserve multivariate dependencies.

645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690

### 691 **2.3.2 ML Utility**

692

693 To evaluate the utility of synthetic data, we assessed its predictive utility—namely,  
694 the ability of machine learning models trained on synthetic data to replicate the pre-  
695 dictive performance of models trained on real data. Predictive utility was evaluated  
696 by training three machine learning models—Random Forest (RF), Logistic Regression  
697 (LR), and Naive Bayes (NB)—on synthetic datasets and testing them on the real test  
698 set. The test set, comprising 20% of the original data, was held out from all synthetic  
699 data generation processes to ensure independent and unbiased evaluation.  
700  
701  
702  
703

704 The classification task involved binary prediction of health conditions. Models were  
705 trained using 5-fold cross-validation to mitigate sampling variance and ensure robust  
706 performance estimation. Model performance was evaluated using recall, precision, F1-  
707 score, and Area Under the Receiver Operating Characteristic Curve (AUC-ROC).  
708  
709  
710

711 These metrics were selected due to their clinical relevance: recall reflects the abil-  
712 ity to correctly identify true positive cases, which is crucial in healthcare diagnostics  
713 to avoid missed cases; precision measures the accuracy of positive predictions, helping  
714 to minimise false positives and prevent unnecessary interventions; F1-score balances  
715 precision and recall, providing an overall assessment where both types of errors are con-  
716 sequential; and AUC-ROC evaluates the model's overall discriminative ability, offering  
717 insight into predictive reliability across various decision thresholds.  
718  
719  
720  
721  
722

723

### 724 **2.3.3 Fairness**

725

726 To assess fairness, we employed three widely used group fairness metrics: Accuracy  
727 Parity, Equalised Odds, and Predictive Rate Parity. These metrics were applied to  
728 evaluate and compare the fairness of machine learning models across subgroups within  
729 the sensitive attribute. The analysis specifically focused on the most underrepresented  
730 subgroup—referred to as the "protected subgroup"—in each dataset. In the CVD  
731  
732  
733  
734

735

736

dataset, the protected subgroup was Group 0, while in the Diabetes dataset, it was Group 2, both corresponding to individuals of Asian ethnicity.

Accuracy Parity is achieved when the prediction accuracy for the protected subgroup is equivalent to that of the unprotected subgroups. This metric measures the probability of making correct predictions (whether positive or negative) across different groups. In this context, it ensures that individuals with or without a health condition are classified with comparable overall accuracy, regardless of subgroup membership. Formally, this can be expressed as:  $P(d = Y, SA = PS) = P(d = Y, SA = SG)$  where  $d$  represents the model's prediction,  $Y$  is the true label,  $SA$  is the sensitive attribute,  $PS$  is the protected subgroup, and  $SG$  is any other subgroup. In our evaluation, Accuracy Parity was considered satisfied when the classification accuracy was similar between the protected and unprotected subgroups.

Predictive Rate Parity assesses whether precision (or positive predictive value, PPV) is consistent across subgroups. PPV indicates the likelihood that a positive prediction corresponds to an actual positive case. Fairness under this metric requires that the proportion of true positives among all predicted positives is equivalent for the protected and unprotected subgroups  $P(Y = 1 | d = 1, SA = PS) = P(Y = 1 | d = 1, SA = SG)$ . This metric also indirectly ensures equality in false discovery rates (FDR), which measure the proportion of false positives among all predicted positives:  $P(Y = 0 | d = 1, SA = PS) = P(Y = 0 | d = 1, SA = SG)$ . In our experiments, this implies that for both protected and unprotected subgroups, the probability of a patient being accurately classified as having a health condition should be the same. This metric ensures that individuals receiving positive predictions from different subgroups have equal confidence in their prediction accuracy.

Equalised Odds ensures that a classifier's true positive rate (TPR) and false positive rate (FPR) are approximately equal across subgroups. This metric evaluates the consistency of correct and incorrect predictions between the protected and unprotected

737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782

783 subgroups. Specifically, patients with and without a health condition should have equal  
784 probabilities of being correctly or incorrectly classified, irrespective of their sensitive  
785 attribute values. Mathematically, for a binary outcome  $P(d = 1 | Y = i, SA = PS) =$   
786  $P(d = 1 | Y = i, SA = SG)$ . Where,  $i$  denotes the binary health condition (present or  
787 absent), and  $d=1$  indicates a positive prediction .  
788  
789

791 In our experiments, a model was deemed to maintain fairness if the values of the  
792 fairness metrics (Accuracy Parity, Equalised Odds, and Predictive Rate Parity) for  
793 the protected subgroup were approximately equal to those of other subgroups, ideally  
794 close to a ratio of 1. Significant deviations from this benchmark highlighted fairness  
795 concerns. A ratio greater than 1 indicated that the unprotected subgroup benefited  
796 more from model performance relative to the protected group, while a ratio below 1  
797 suggested that the protected subgroup had relatively better outcomes compared to  
798 the other subgroup.  
799  
800

801 A complete summary of the experimental framework, including datasets, preprocess-  
802 ing, generative model configurations, evaluation metrics, and computational environ-  
803 ment, is provided in Table 2. This table consolidates all methodological components  
804 and serves as a reference for the subsequent results and discussion.  
805  
806  
807

## 812 3 Results

### 815 3.1 Data Fidelity

816 Figure 1 presents the statistical fidelity results for the CVD and Diabetes datasets  
817 across five synthetic data generation models: BayesBoost, CTGAN, VAE, DECAF,  
818 and CopulaGAN. The evaluation includes five complementary metrics: Category Cov-  
819 erage, TV Complement, Range Coverage, Boundary Adherence, and KS Complement.  
820  
821 These metrics assess how closely the synthetic data replicate key statistical proper-  
822 ties of the real datasets. Category Coverage and Range Coverage measure the extent  
823 to which categorical and numerical feature values in the synthetic data overlap with  
824  
825  
826  
827  
828

**Table 2:** Summary of experimental design, models, and evaluation settings

Component	Description	Details / Parameters	Purpose / Output
Datasets	Two binary classification healthcare datasets	CVD (10,000 samples, 22 features); Diabetes (1,879 samples, 44 features)	Assess model generalisability across domains
Sensitive attribute	Ethnicity (categorical, multiple subgroups)	CVD: 6 groups (protected = Group 0, Asian); Diabetes: 4 groups (protected = Group 2, Asian)	Enable subgroup fairness evaluation
Train-test split	Stratified 80/20 split	Preserves distributions of sensitive and target variables	Prevents data leakage and ensures independent evaluation
Preprocessing	Encoding and scaling	LabelEncoder + MinMaxScaler [0,1]; scaler fitted on train set only	Standardised input for all models
Models	Five state-of-the-art methods	CopulaGAN, CTGAN, TVAE, DECAF (Python/SDV); BayesBoost (R/bnlearn)	Compare deep vs probabilistic approaches
Configurations	Default hyperparameters	Default batch size/epochs; DECAF ND (no-debias) causal variant	Ensure comparability and reproducibility
Iterations	150 per model per dataset (distinct random seeds)	Synthetic data regenerated each iteration	Average results to reduce random variability
<b>Evaluation</b>			
metrics	Fidelity	TV/KS Complement, Category/Range Coverage, Boundary Adherence, Correlation (Frobenius norm)	Quantify statistical similarity between real and synthetic data
	Utility	RF, LR, NB classifiers; evaluated via Recall, Precision, F1, AUC-ROC	Assess predictive consistency between real and synthetic data
	Fairness	Accuracy Parity, Predictive Rate Parity, Equalised Odds	Evaluate subgroup equity across sensitive attributes
<b>Hardware/</b>			
Software	Experimental environment	Intel Evo CPU (3.8 GHz), 32 GB RAM, RTX 4060 (16 GB); Python 3.12, R 4.3	Execution environment and library versions

those in the real data. TV Complement and KS Complement evaluate similarity in the marginal distributions. Boundary Adherence measures whether synthetic values remain within the observed real-data limits. Higher values across all metrics indicate closer alignment between synthetic and real data.

875 For categorical variables, BayesBoost, CTGAN, and CopulaGAN successfully cap-  
876 tured all subcategories present in the real datasets for both CVD and Diabetes. In  
877 contrast, DECAF lagged slightly behind, covering approximately 87% of the subcate-  
878 gories in the CVD dataset. TV Complement scores remained consistently high across  
879 all models for the CVD dataset and exceeded 80% for the Diabetes dataset, suggesting  
880 that the marginal distributions of categorical variables in the synthetic data closely  
881 resembled those of the real data.  
882

886 Regarding numerical variables, Range Coverage showed a significant decline in  
887 the CVD dataset for the deep generative models CopulaGAN, CTGAN, VAE, and  
888 DECAF, indicating challenges in capturing the complete range of feature distributions.  
889 DECAF also demonstrated poor range coverage in the Diabetes dataset, whereas the  
890 other models successfully covered the full numerical range. KS Complement scores  
891 exhibited moderate variation, with CTGAN and VAE performing slightly worse than  
892 BayesBoost, which achieved the best preservation of the numerical distributions,  
893 maintaining KS Complement values above 97%. DECAF recorded the lowest KS Com-  
894 plement scores, with approximately 80% fidelity in the CVD dataset and even lower  
895 in the Diabetes dataset, further underscoring its limitations in accurately replicating  
896 real data distributions.  
897

905 Boundary Adherence remained consistently high across all models for both  
906 datasets, demonstrating that the synthetic samples adhered well to the expected  
907 numerical feature boundaries and no new values were introduced in the generated data.  
908

910 Figure 2 illustrates examples of generated data distributions for selected features,  
911 based on randomly chosen synthetic datasets of identical size to the real datasets. For  
912 numerical features, such as systolic blood pressure (CVD dataset) and age (Diabetes  
913 dataset), BayesBoost and VAE most closely approximated the distributions observed  
914 in the real data. In contrast, DECAF exhibited the greatest divergence, generating  
915 distributions that noticeably deviated from the original datasets.  
916  
917  
918  
919  
920



**Fig. 1:** Statistical similarity between synthetic data and real data for Diabetes and CVD datasets. Higher values indicate closer alignment between synthetic and real data.

For categorical features, particularly Ethnicity in both datasets, both DECAF and VAE struggled to accurately replicate the distributions of underrepresented subgroups. DECAF failed to generate any instances for Groups 3 and 4 in the CVD dataset, and for Group 3 in the Diabetes dataset, all of which were already minority subgroups in the original data. Similarly, VAE failed to generate instances for the protected subgroup (Group 2) in both datasets. Furthermore, VAE distorted subgroup proportions, notably over-representing Group 4 in the Diabetes dataset and failing to generate instances for Group 5, the most populous group in the real data.

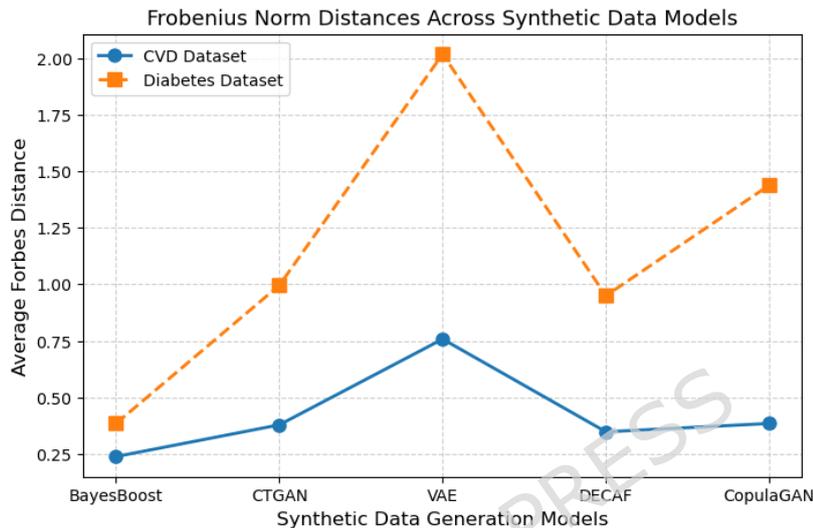
By contrast, CopulaGAN and CTGAN maintained subgroup distributions more faithfully, with only minor deviations from the real datasets. BayesBoost, however, over-sampled certain minority subgroups, nearly doubling their representation relative to the real data. This behaviour can be attributed to BayesBoost's internal mechanism, which deliberately amplifies minority subgroup instances to enhance class balance during generation.

For the multivariate analysis, Figure 3 presents the average Frobenius norm distances between the correlation matrices of the real and synthetic datasets across

921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966



For the Diabetes dataset, discrepancies were more pronounced. BayesBoost again achieved the lowest distance (0.3845), reinforcing its ability to maintain the correlation structure of the real data. In contrast, CTGAN (0.9953) and DECAF (0.9495) displayed substantial deviations, while CopulaGAN (1.4379) and VAE (2.0172) exhibited the highest distances, indicating significant loss of inter-variable relationships within the synthetic data.



**Fig. 3:** Average Frobenius Norm Distances across synthetic data generation models

### 3.1.1 ML Utility

This section presents a comparative analysis of the performance of machine learning classifiers trained on real and synthetic datasets for two use cases: CVD prediction and diabetes prediction. The effectiveness of five synthetic data generation methods—CopulaGAN, CTGAN, VAE, DECAF, and BayesBoost—is evaluated by comparing performance metrics, including accuracy, precision, recall, F1-score, and AUC-ROC, against those achieved on real datasets. All performance metrics are summarised in Table 3.

1059 **Table 3:** Machine Learning Performance Metrics on Real and Synthetic Data for CVD  
 1060 and Diabetes Datasets  
 1061

		CVD						
M	m <sup>1</sup>	a	b	c	d	e	f	
1064	RF	i	0.914 ± 0.000	0.887 ± 0.001	0.889 ± 0.001	0.904 ± 0.001	0.799 ± 0.006	0.919 ± 0.000
1065		ii	0.588 ± 0.002	0.474 ± 0.007	0.473 ± 0.007	0.597 ± 0.002	0.436 ± 0.006	0.619 ± 0.001
1066		iii	0.459 ± 0.002	0.408 ± 0.012	0.396 ± 0.011	0.535 ± 0.004	0.558 ± 0.008	0.493 ± 0.001
1067		iv	0.463 ± 0.002	0.666 ± 0.010	0.687 ± 0.010	0.682 ± 0.004	0.384 ± 0.010	0.829 ± 0.001
1068		v	0.884 ± 0.001	0.853 ± 0.001	0.854 ± 0.001	0.869 ± 0.001	0.798 ± 0.003	0.896 ± 0.000
1069	LR	i	0.912 ± 0.000	0.889 ± 0.001	0.891 ± 0.001	0.901 ± 0.001	0.873 ± 0.002	0.922 ± 0.000
1070		ii	0.585 ± 0.002	0.481 ± 0.007	0.479 ± 0.007	0.596 ± 0.002	0.533 ± 0.004	0.642 ± 0.001
1071		iii	0.463 ± 0.002	0.408 ± 0.012	0.395 ± 0.011	0.552 ± 0.004	0.536 ± 0.005	0.525 ± 0.001
1072		iv	0.795 ± 0.002	0.680 ± 0.009	0.700 ± 0.009	0.657 ± 0.005	0.548 ± 0.008	0.825 ± 0.001
1073		v	0.878 ± 0.001	0.864 ± 0.001	0.866 ± 0.001	0.866 ± 0.001	0.849 ± 0.002	0.891 ± 0.000
1074	NB	i	0.763 ± 0.002	0.854 ± 0.002	0.851 ± 0.005	0.845 ± 0.005	0.851 ± 0.005	0.850 ± 0.009
1075		ii	0.476 ± 0.002	0.542 ± 0.003	0.545 ± 0.003	0.543 ± 0.004	0.491 ± 0.006	0.543 ± 0.006
1076		iii	0.805 ± 0.002	0.648 ± 0.005	0.653 ± 0.005	0.673 ± 0.004	0.524 ± 0.008	0.590 ± 0.011
1077		iv	0.340 ± 0.003	0.473 ± 0.005	0.477 ± 0.005	0.463 ± 0.006	0.485 ± 0.009	0.577 ± 0.014
1078		v	0.868 ± 0.001	0.842 ± 0.002	0.846 ± 0.002	0.830 ± 0.002	0.828 ± 0.003	0.835 ± 0.003
		Diabetes						
M	m	a	b	c	d	e	f	
1079	RF	i	0.906 ± 0.001	0.542 ± 0.006	0.545 ± 0.006	0.840 ± 0.002	0.601 ± 0.000	0.882 ± 0.001
1080		ii	0.873 ± 0.001	0.221 ± 0.018	0.221 ± 0.018	0.786 ± 0.003	0.130 ± 0.000	0.839 ± 0.002
1081		iii	0.812 ± 0.002	0.281 ± 0.028	0.275 ± 0.027	0.738 ± 0.005	0.151 ± 0.001	0.775 ± 0.003
1082		iv	0.945 ± 0.001	0.352 ± 0.019	0.382 ± 0.020	0.848 ± 0.004	0.191 ± 0.000	0.916 ± 0.001
1083		v	0.951 ± 0.001	0.503 ± 0.004	0.500 ± 0.004	0.900 ± 0.001	0.512 ± 0.002	0.947 ± 0.000
1084	LR	i	0.831 ± 0.001	0.539 ± 0.006	0.544 ± 0.006	0.767 ± 0.002	0.489 ± 0.008	0.763 ± 0.001
1085		ii	0.781 ± 0.002	0.268 ± 0.017	0.255 ± 0.017	0.730 ± 0.002	0.512 ± 0.014	0.653 ± 0.002
1086		iii	0.757 ± 0.002	0.316 ± 0.026	0.297 ± 0.026	0.784 ± 0.005	0.798 ± 0.027	0.561 ± 0.003
1087		iv	0.808 ± 0.002	0.406 ± 0.014	0.411 ± 0.012	0.688 ± 0.004	0.421 ± 0.011	0.785 ± 0.002
1088		v	0.905 ± 0.001	0.505 ± 0.005	0.504 ± 0.005	0.854 ± 0.002	0.640 ± 0.007	0.854 ± 0.001
1089	NB	i	0.790 ± 0.002	0.525 ± 0.005	0.534 ± 0.005	0.517 ± 0.003	0.601 ± 0.000	0.671 ± 0.002
1090		ii	0.717 ± 0.003	0.379 ± 0.010	0.370 ± 0.008	0.534 ± 0.003	0.140 ± 0.012	0.549 ± 0.002
1091		iii	0.667 ± 0.004	0.407 ± 0.017	0.377 ± 0.015	0.692 ± 0.007	0.180 ± 0.001	0.502 ± 0.003
1092		iv	0.777 ± 0.003	0.408 ± 0.006	0.420 ± 0.004	0.439 ± 0.003	0.201 ± 0.003	0.608 ± 0.003
1093		v	0.875 ± 0.001	0.504 ± 0.005	0.504 ± 0.004	0.564 ± 0.002	0.500 ± 0.000	0.758 ± 0.002

1090 <sup>1</sup>Performance metrics; i: accuracy; ii: f1-score; iii: recall; iv: precision; v: auc-roc;  
 1091 Training Data; a: real data; b: CopulaGAN; c: CTGAN; d: VAE; e: DECAF; f: BayesBoost

1092  
 1093 Across both datasets, the RF classifier consistently achieved the highest perfor-  
 1094 mance when trained on real data, demonstrating its robustness. For CVD prediction,  
 1095 RF achieved an accuracy of  $0.914 \pm 0.000$  and an AUC-ROC of  $0.884 \pm 0.001$ , while  
 1096 for diabetes, the corresponding metrics were  $0.906 \pm 0.001$  and  $0.951 \pm 0.001$ . Logistic  
 1097 Regression (LR) also performed strongly on the CVD data, recording an accuracy of  
 1098  $0.912 \pm 0.000$  and an AUC-ROC of  $0.878 \pm 0.001$ . However, its performance declined  
 1099 for diabetes prediction (accuracy of  $0.831 \pm 0.001$  and AUC-ROC of  $0.905 \pm 0.001$ )  
 1100  
 1101  
 1102  
 1103  
 1104

compared to RF. Naive Bayes (NB) demonstrated the lowest performance across both datasets, achieving an accuracy of  $0.763 \pm 0.002$  for CVD and  $0.790 \pm 0.002$  for diabetes.

Recall varied significantly across classifiers and datasets. For the CVD data, NB achieved the highest recall ( $0.805 \pm 0.003$ ), outperforming RF ( $0.459 \pm 0.002$ ) and LR ( $0.463 \pm 0.002$ ). In contrast, for diabetes prediction, RF achieved the highest recall ( $0.812 \pm 0.002$ ), outperforming LR ( $0.757 \pm 0.002$ ) and NB ( $0.667 \pm 0.004$ ).

For CVD prediction, synthetic data generated by CopulaGAN and CTGAN yielded reasonable performance for RF and LR in terms of accuracy, AUC-ROC, and recall, with only minor deviations from the performance observed on real data. However, for precision and F1-score, larger differences were observed, reaching up to 0.2 in some cases. NB's performance on CopulaGAN- and CTGAN-generated data remained suboptimal, with decreased recall but slight improvements in accuracy and precision, reflecting a trade-off between sensitivity and specificity.

Conversely, for diabetes prediction, CopulaGAN and CTGAN exhibited significant performance deterioration across all classifiers. RF accuracy on CopulaGAN-generated data dropped sharply to 0.542, with AUC-ROC falling to 0.503. Similarly, CTGAN yielded an RF accuracy of  $0.545 \pm 0.006$ .

The VAE-generated CVD data demonstrated high utility, closely resembling the predictive capabilities achieved with real data across all classifiers. Specifically, RF trained on VAE-generated data achieved an accuracy of  $0.904 \pm 0.001$  and an AUC-ROC of  $0.869 \pm 0.001$ . LR similarly achieved strong results with an accuracy of  $0.901 \pm 0.001$  and an AUC-ROC of  $0.866 \pm 0.001$ .

For the diabetes dataset, VAE outperformed CopulaGAN, CTGAN, and DECAF, delivering results closest to real data across most metrics and classifiers. RF demonstrated particularly strong performance on VAE data, achieving an accuracy of  $0.840 \pm$

1151 0.002 and an AUC-ROC of  $0.900 \pm 0.001$ . LR also performed well, achieving an accu-  
1152 racy of  $0.767 \pm 0.002$  and an AUC-ROC of  $0.854 \pm 0.002$ . These results demonstrate  
1153 that VAE maintained utility and complexity relatively well for challenging datasets,  
1154 enabling robust ML model performance.  
1155

1157 DECAF, however, demonstrated the weakest performance overall. For CVD, RF  
1158 and LR accuracies dropped to 0.799 and 0.873, respectively, with corresponding AUC-  
1159 ROC values of 0.798 and 0.849. Although DECAF slightly improved recall, this was  
1160 overshadowed by significant reductions in precision and overall model utility. On the  
1161 diabetes dataset, DECAF's performance collapsed almost entirely, with RF and NB  
1162 accuracy falling to  $0.601 \pm 0.000$  and metrics such as precision and recall approaching  
1163 0.151 and 0.191, indicating a major loss of predictive signal.  
1164

1166 Among all synthetic data generation methods, BayesBoost consistently achieved  
1167 the best performance across most evaluation metrics. For CVD data, RF showed a  
1168 slight improvement when trained on BayesBoost-generated data, achieving an accu-  
1169 racy of 0.919 compared to 0.914 on real data, and a higher F1-score (0.619 vs 0.588).  
1170 Similarly, LR trained on BayesBoost data achieved an accuracy of 0.922 and an F1-  
1171 score of 0.642, outperforming the results achieved with real data (accuracy of 0.912  
1172 and F1-score of 0.585). NB also showed noticeable improvements in accuracy (0.850  
1173 vs 0.763) and F1-score (0.543 vs 0.476), although recall slightly declined.  
1174

1176 For the diabetes dataset, BayesBoost maintained relatively strong predictive util-  
1177 ity compared to other synthetic data methods. RF performance remained high, with  
1178 only a marginal decline in accuracy (0.882 vs 0.906) and AUC-ROC (0.947 vs 0.951).  
1179 However, LR and NB experienced more pronounced declines: LR's accuracy fell from  
1180 0.831 to 0.763, and NB's accuracy decreased from 0.790 to 0.671.  
1181

### 1192 3.1.2 Fairness Analysis

1193 This section presents an evaluation of the fairness of training data by comparing fair-  
1194 ness metrics across different classifiers and subgroups defined by the sensitive attribute.  
1195

The primary objective is to assess the baseline fairness of the real datasets and to examine the extent to which synthetic data generated by various models preserves or distorts this fairness. Fairness is quantified using three established metrics: Accuracy Parity, Equalised Odds, and Predictive Rate Parity, as detailed in Tables 4 and 5. Real data serves as the benchmark against which the performance of each synthetic data generation model is evaluated.

- **CVD Data**

**Accuracy parity:** The real data consistently demonstrates fairness across all classifiers (LR, Rf and NB) and ethnic subgroups. The protected subgroup (Group 0 - Asians) exhibits classification accuracy comparable to all other groups (Groups 1–5), with the ratio differences not exceeding 0.04 across classifiers. This finding highlights the capacity of real data to maintain equitable classification performance across diverse subgroups within the sensitive attribute of ethnicity.

Synthetic data generated by all methods largely preserves the accuracy parity observed in the real data for groups 1, 2, 4, and 5, with only minor deviations that can be considered negligible. However, for group 3, the ratio deviation from real data is slightly higher when using the NB classifier, reaching more than 0.1. Despite this, the overall fairness remains well preserved relative to the baseline.

**Equalised odds:** For models trained on real data, NB achieved the best fairness performance, with deviations from the ideal parity value of 1 remaining within 0.05 across all ethnic groups. LR and RF also maintained reasonable parity overall; however, noticeable deviations emerged, particularly for Group 5, where parity ratios reached 1.19 (LR) and 1.23 (RF), indicating a systematic disadvantage to the protected subgroup (Group 0).

When classifiers were trained on synthetic datasets, it became evident that all generative models struggled to maintain equalised odds parity, particularly for Groups

1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242

1243 **Table 4:** Fairness Metrics for Ethnic Groups Measured Relative to Sensitive Group 0  
 1244 - CVD Data

1246		Accuracy Parity			Equalised Odds			PRP			
1247	G <sup>1</sup>	D <sup>2</sup>	LR	NB	RF	LR	NB	RF	LR	NB	RF
1248	1	a	1.01	0.99	1.02	1.06	0.95	1.11	0.81	0.77	0.89
1249		b	1.06	1.08	1.05	1.63	1.32	1.44	1.08	1.33	1.04
1250		c	1.06	1.08	1.05	1.70	1.33	1.50	1.10	1.29	1.03
1251		d	1.07	1.10	1.08	1.84	1.29	1.57	1.08	1.47	1.29
1252		e	1.08	1.09	1.04	1.81	1.74	1.89	1.36	1.45	1.04
1253		f	1.06	1.05	1.05	1.66	1.71	1.50	1.01	1.18	1.00
1254	2	a	0.98	0.99	0.99	1.10	1.00	1.16	0.81	1.07	0.86
1255		b	0.97	0.98	0.97	0.73	0.85	0.79	0.88	0.93	0.88
1256		c	0.97	0.98	0.98	0.73	0.85	0.79	0.89	0.92	0.89
1257		d	0.98	1.00	0.98	0.98	0.88	0.92	0.88	1.03	0.92
1258		e	1.00	1.00	0.98	1.10	1.13	1.20	1.01	1.06	1.02
1259		f	0.98	0.98	0.98	0.94	1.08	1.04	0.83	0.89	0.83
1260	3	a	0.96	1.01	0.97	0.88	1.02	0.98	0.82	1.20	0.88
1261		b	0.92	0.89	0.92	0.35	0.69	0.37	0.53	0.66	0.55
1262		c	0.92	0.89	0.92	0.35	0.65	0.37	0.51	0.63	0.55
1263		d	0.94	0.93	0.95	0.68	0.72	0.73	0.76	0.79	0.87
1264		e	0.93	0.94	0.91	0.55	0.52	0.64	0.66	0.68	0.68
1265		f	0.93	0.87	0.94	0.46	0.55	0.60	0.73	0.44	0.74
1266	4	a	1.00	0.98	1.01	0.97	0.97	1.05	0.70	0.74	0.77
1267		b	0.99	0.98	1.00	0.96	1.10	1.03	0.67	0.73	0.76
1268		c	0.99	0.98	1.00	0.93	1.10	0.96	0.67	0.74	0.73
1269		d	1.01	0.97	1.00	1.41	1.09	1.23	0.79	0.74	0.79
1270		e	1.03	1.04	0.99	1.45	1.45	1.76	0.95	1.02	0.89
1271		f	1.02	1.04	1.04	1.39	1.53	1.59	0.81	1.00	0.93
1272	5	a	0.99	0.96	0.99	1.19	0.95	1.23	0.85	0.97	0.89
1273		b	0.98	0.99	0.99	1.03	1.03	1.05	0.86	0.96	0.88
1274		c	0.99	0.99	0.99	1.02	1.02	1.04	0.87	0.95	0.88
1275		d	0.99	1.00	1.00	1.22	1.03	1.16	0.86	1.04	0.99
1276		e	1.01	1.01	0.99	1.37	1.33	1.41	1.05	1.09	1.05
1277		f	1.00	1.03	1.00	1.26	1.33	1.29	0.83	0.96	0.82

1274 Fairness metrics are measured relative to Group 0, with values closest to 1 being optimal and high-  
 1275 lighted in white, while darker red shades indicate greater deviation from 1.

1276 <sup>1</sup>G: Ethnicity groups; 0: Asian; 1: Black; 2: Mixed; 3: Others; 4: unknown; 5: White

1277 <sup>2</sup>Data; a: real data; b: CopulaGAN; c: CTGAN; d: VAE; e: DECAF; f: BayesBoost

1278

1279 1 and 3. Synthetic data often introduced higher parity values for Group 1 and sig-  
 1280 nificantly lower parity values for Group 3, indicating a pronounced amplification of  
 1281 existing biases.

1284 Synthetic datasets produced by CopulaGAN, CTGAN, and DECAF led to sub-  
 1285 stantial fairness distortions. These models caused dramatic parity shifts, particularly  
 1286 in Groups 1, 2 and 3, where parity ratios fell sharply—reaching approximately 50–70%  
 1287  
 1288

lower than Group 0 for Group 3, and conversely achieving 40–90% better parity for Group 1 across LR and RF classifiers. These extreme deviations highlight that CopulaGAN, CTGAN, and DECAF severely amplified fairness concerns, disproportionately benefiting or disadvantaging certain groups relative to the protected subgroup. Although fairness discrepancies in Groups 4 and 5 were somewhat less severe, their deviations remained notable compared to the real data baselines and 1.

BayesBoost also exhibited considerable fairness disparities across most groups. While it maintained relatively stable parity for Group 2, significant fairness violations were observed for Groups 1, 3, 4, and 5. In Group 1, parity ratios rose as high as 1.66, 1.71, and 1.50 for LR, NB, and RF respectively, suggesting substantial over-prediction compared to the protected subgroup. In contrast, Group 3 exhibited severe under-prediction, with parity ratios falling between 0.46 and 0.60 across classifiers. Furthermore, fairness deviations in Groups 4 and 5 remained consistently elevated, with parity values exceeding 1.25 in several cases.

VAE, on the other hand, achieved equalised odds values closest to those observed in the real data compared to the other generation methods. While it did introduce some fairness concerns, these deviations were substantially less severe, particularly when combined with NB and RF. This was most evident across Groups 2, 3, 4, and 5, where VAE consistently maintained parity values nearer to those in the baseline.

**Predictive rate parity:** In the real dataset, the PRP outcomes do not conform to the expectations of ideal fairness. Across all classifiers, models consistently demonstrated better predictive parity for the protected subgroup (Group 0) compared to other ethnic groups. This consistent advantage in prediction rates for Group 0 reflects a fairness imbalance inherently embedded within the original data.

Upon training classifiers on synthetic datasets, pronounced fairness variations emerged across generative models. In particular, all synthetic models reversed the fairness trends observed in real data for Group 1. Whereas Group 0 initially exhibited

1335 higher PRP values, synthetic data disproportionately favoured Group 1. This shift was  
1336 especially pronounced with the NB classifier, where parity ratios indicated over a 40%  
1337 advantage for Group 1 relative to Group 0. Although fairness concerns persisted across  
1338 classifiers, the extent of fairness deviation was notably reduced with RF, where PRP  
1339 values were shifted closer to 1, thereby rendering the predictions for Group 1 more  
1340 comparable to those for Group 0. In addition, across all generative models, additional  
1341 fairness disparities were introduced in Group 3 compared to the real data baseline.

1342 For Groups 2, 4, and 5, CopulaGAN, CTGAN, and VAE largely preserved the  
1343 fairness structure observed in the real data, without inducing substantial deviations.  
1344

1345 While DECAF - like other generative models- introduced pronounced fairness dis-  
1346 tortions, particularly affecting Groups 1 and 3. Nonetheless, in Groups 2, 4, and 5,  
1347 DECAF marginally improved fairness by narrowing the gap in parity ratios between  
1348 the protected subgroup and the other groups.

1349 Among the evaluated methods, BayesBoost demonstrated the most consistent  
1350 preservation of fairness. It either retained the original PRP values or yielded ratios  
1351 that remained closest to 1 across all groups. Particularly, RF trained on BayesBoost-  
1352 generated data exhibited highly stable parity, achieving a PRP of 1.00 for Group 1  
1353 and 0.74 for Group 3, outperforming all other model-classifier combinations by min-  
1354 imising deviations from 1 or from real data parity. Logistic Regression (LR) combined  
1355 with BayesBoost similarly demonstrated competitive fairness, albeit with a slight  
1356 underestimation for Group 3 (PRP = 0.73) compared to real data.

#### 1357 • **Diabetes Data**

1358 **Accuracy parity:** On the real data, AP results were consistently close to the  
1359 ideal value of 1 across all ethnic groups and classifiers (LR, NB, RF). Deviations were  
1360 minimal, with values ranging between 0.97 and 1.04, indicating that classification  
1361 accuracy was relatively balanced between the protected subgroup (Group 2, Asian  
1362 ethnicity) and the other ethnic groups.

**Table 5:** Fairness Metrics for Ethnic Groups Measured Relative to Sensitive Group 2 - Diabetes Dataset

G <sup>1</sup>	D <sup>2</sup>	Accuracy Parity			Equalised Odds			PRP		
		LR	NB	RF	LR	NB	RF	LR	NB	RF
0	a	1.02	1.02	0.99	1.15	1.12	1.01	1.02	1.08	1.00
	b	1.00	1.01	0.98	2.17	1.48	1.47	9.82	1.96	10.27
	c	1.00	1.03	0.99	1.44	1.91	1.07	10.89	3.53	9.49
	d	1.14	1.05	1.04	1.06	1.10	0.98	1.27	1.15	1.16
	e	0.94	0.94	0.94	0.91	0.12	0.15	1.08	0.11	0.13
	f	1.02	1.03	1.04	1.20	1.36	1.21	1.05	1.18	1.02
1	a	1.02	1.04	0.97	1.21	1.22	0.99	1.02	1.10	0.96
	b	1.01	1.00	0.96	1.86	1.34	1.24	8.42	2.59	7.82
	c	0.99	1.02	0.98	1.13	1.44	1.08	8.72	2.94	8.69
	d	1.17	1.11	1.05	1.05	1.14	1.01	1.35	1.26	1.17
	e	0.92	0.92	0.92	0.32	0.11	0.14	1.90	0.20	0.18
	f	1.04	1.13	1.02	1.09	1.53	1.13	1.21	1.42	1.04
3	a	1.01	1.04	0.99	1.13	1.11	0.97	0.98	1.09	0.99
	b	1.03	1.03	1.00	1.75	1.52	1.31	4.90	2.33	5.54
	c	0.99	1.01	1.00	1.04	1.38	1.09	3.98	3.29	4.37
	d	1.06	0.99	1.00	0.86	1.12	0.84	1.17	1.03	1.11
	e	0.99	0.99	0.99	0.91	0.23	0.21	1.01	0.16	0.19
	f	1.00	1.02	0.95	0.91	1.51	0.90	1.08	1.07	0.95

Fairness metrics are measured relative to Group 2, with values closest to 1 being optimal and highlighted in white, while darker red shades indicate greater deviation from 1.

<sup>1</sup>G: Ethnicity groups; 0: Caucasian; 1: African American; 2:Asian; 3:Others.

<sup>2</sup>Data; a: real data; b: CopulaGAN; c: CTGAN; d: VAE; e: DECAF; f: BayesBoost

When synthetic data was used, CopulaGAN, CTGAN, and BayesBoost largely preserved the accuracy parity observed in real data across all groups and classifiers. AP ratios remained between 0.96 and 1.04 for most combinations, effectively mirroring real data performance. Although BayesBoost combined with NB introduced some disparities in group 1, the deviations were not substantial.

VAE exhibited a slightly different trend: while fairness was mostly preserved, it introduced noticeable parity inflation, particularly in Groups 0 and 1. For instance, under VAE, Group 0 with LR reached an AP of 1.14 (compared to 1.02 in real data), while Group 1 showed AP values of 1.17 (LR), 1.11 (NB), and 1.05 (RF). Although these deviations are also moderate, they indicate that VAE may increase classification accuracy disproportionately for some groups, potentially shifting fairness dynamics.

1427 DECAF displayed a slight reduction in AP values across groups, most notably  
1428  
1429 for Groups 0 and 1, where AP dropped to 0.94 and 0.92, respectively. However, the  
1430 magnitude of these deviations remained below 10%, suggesting that while DECAF  
1431  
1432 introduced a slight disadvantage for these groups relative to Group 2, it did not  
1433  
1434 significantly compromise overall fairness.

1435 **Equalised odds:** On the real data, EO results demonstrated reasonably good  
1436  
1437 fairness, particularly for RF, where parity values remained close to the ideal value  
1438  
1439 of 1 across all groups (ranging between 0.97 and 1.01). However, larger deviations  
1440 were observed for LR and NB. Notably, in Group 1, parity ratios reached 1.21 (LR)  
1441  
1442 and 1.22 (NB), indicating that the protected subgroup (Group 2) was systematically  
1443  
1444 disadvantaged relative to Group 1. Similar patterns were seen for Group 0 and Group  
1445  
1446 3, where parity ratios reached 1.15 (LR) and 1.12 (NB).

1447 When classifiers were trained on synthetic datasets, substantial fairness distortions  
1448  
1449 became apparent, particularly for CopulaGAN and CTGAN. These models intro-  
1450  
1451 duced severe parity inflations, especially in Group 0 (e.g., CopulaGAN\_LR = 2.17,  
1452 CTGAN\_LR = 1.44) and Group 3 (CopulaGAN\_LR = 1.75, CopulaGAN\_NB = 1.52).  
1453  
1454 Such high ratios suggest significant over-prediction relative to the protected subgroup,  
1455  
1456 thus amplifying existing fairness gaps in the real data.

1457 DECAF exhibited the poorest performance regarding equalised odds. In both  
1458  
1459 Group 0, 1 and 2, DECAF produced near-zero values (between 0.11 and 0.23 with NB  
1460  
1461 and RF), indicating that correct and incorrect prediction rates between the protected  
1462  
1463 group and others became dramatically imbalanced. This critical fairness breakdown  
1464  
1465 renders DECAF unsuitable for fairness-sensitive applications in this context.

1466 Among all models, VAE showed the best preservation of equalised odds. Across  
1467  
1468 Groups 0 and 1, VAE maintained EO ratios close to 1, achieving values of 1.06 (LR),  
1469  
1470 1.10 (NB), and 0.98 (RF) for Group 0, and 1.05 (LR), 1.14 (NB), and 1.01 (RF) for  
1471  
1472 Group 1. These results represent only minor fairness deviations relative to the real

data baseline. For Group 3, VAE achieved parity values of 0.86 (LR), 1.12 (NB), and 0.84 (RF), which, although slightly lower, still performed substantially better than CopulaGAN, CTGAN, and DECAF.

BayesBoost, by contrast, demonstrated a more mixed fairness profile. While it generally performed better than CopulaGAN, CTGAN, and DECAF, it consistently introduced larger deviations from 1 than VAE. In Group 0, BayesBoost recorded values of 1.20 (LR), 1.36 (NB), and 1.21 (RF), notably inflating the advantage compared to VAE. Particularly concerning was the combination of BayesBoost with NB, where the parity ratio reached 1.36, higher than the baseline real data value of 1.12, indicating a significant worsening of fairness compared to LR and RF.

When comparing LR and RF across VAE and BayesBoost, VAE consistently pushed equalised odds ratios closer to 1 for Groups 0 and 1. RF and LR trained on VAE-synthetic data achieved EO values notably closer to the fairness ideal compared to those trained on BayesBoost-synthetic data. However, for Group 3, a different pattern emerged: VAE lowered parity values to 0.86 (LR) and 0.84 (RF), slightly below the ideal, whereas BayesBoost combined with LR achieved a parity value closer to 1 for this group. This suggests that for preserving equalised odds parity in Group 3, BayesBoost combined with LR slightly outperformed VAE. Nevertheless, overall, VAE combined with RF achieved the best overall equalised odds parity across all groups, with BayesBoost combined with LR providing a complementary solution for achieving more balanced fairness.

**Predictive Rate parity:** On the real data, PRP values across LR, NB, and RF classifiers were consistently close to the ideal value of 1 for all ethnic groups, suggesting reasonably good fairness. Group 0 exhibited PRP values of 1.02 (LR), 1.08 (NB), and 1.00 (RF), while Group 1 ranged between 0.96 and 1.10, and Group 3 displayed similarly balanced values between 0.98 and 1.09. These results indicate a slight systematic advantage for Group 0 that does not exceed 2% but without severe deviations from ideal fairness.

1519 When classifiers were trained on synthetic datasets, significant fairness dispar-  
1520 ities emerged depending on the generative model. Similar to the Equalised Odds  
1521 findings, CopulaGAN, CTGAN, and DECAF introduced major fairness violations,  
1522 substantially inflating PRP values in Groups 0, 1, and 3, often exceeding a ratio of  
1523 2 across all classifiers. These extreme inflations reflect severe fairness issues, particu-  
1524 larly as the models performed markedly worse on the protected Group 2. Moreover,  
1525 DECAF caused PRP values to collapse entirely to near zero for NB and RF classifiers,  
1526 indicating a complete breakdown in fairness for these groups.

1532 In contrast, VAE demonstrated considerably greater stability. Although over-  
1533 predictions persisted, PRP values remained moderately close to 1 across all groups  
1534 and classifiers. For instance, in Group 0, VAE achieved PRP values of 1.27 (LR),  
1535 1.15 (NB), and 1.16 (RF), and for Group 3, ratios remained within the range of  
1536 1.11–1.17. These results suggest that VAE effectively mitigated the extreme fairness  
1537 distortions observed with other deep generative models while maintaining relatively  
1538 stable predictive parity compared to other models.

1544 BayesBoost achieved the best preservation of predictive rate parity among all gen-  
1545 erative models, particularly when paired with LR and RF classifiers. PRP values  
1546 remained consistently close to 1 across groups, with Group 0 reaching 1.05 (LR) and  
1547 1.02 (RF), and Group 3 achieving 1.08 (LR) and 0.95 (RF). BayesBoost outperformed  
1548 VAE in several scenarios, especially when used with Random Forest, maintaining  
1549 predictive parity closest to both the real data baseline and the fairness ideal.

1554 Overall, the results indicate that BayesBoost, particularly when paired with RF,  
1555 offers the most robust approach for maintaining fairness in predictive rate parity. In  
1556 contrast, VAE combined with RF and NB emerged as the most effective option for  
1557 preserving fairness in equalised odds. Conversely, CopulaGAN, CTGAN, and DECAF  
1558 require cautious consideration before use in fairness-critical healthcare applications.  
1559 These findings highlight the importance of rigorous fairness evaluation when selecting  
1560  
1561  
1562  
1563  
1564

synthetic data generation techniques for machine learning models involving protected attributes. 1565  
1566  
1567  
1568

## 4 Discussion 1569

This study presents a comprehensive evaluation of five prominent synthetic data generation models applied to healthcare tabular datasets, assessing their fidelity, machine learning utility, and fairness preservation. While synthetic data is widely promoted as a solution to privacy concerns and data scarcity [3, 31], our results highlight that many models still face substantial challenges in producing clinically reliable, fair, and structurally faithful synthetic datasets—especially when sensitive attributes such as ethnicity are involved. 1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583

After assessing data fidelity and structural preservation, BayesBoost emerged as the most robust method overall. It achieved the highest statistical fidelity, with the lowest Frobenius norm distances and superior boundary and range coverage. Its probabilistic boosting mechanism improved the representation of minority subgroups, but we caution that overboosting, while improving minority representation, can also distort the underlying data distribution and inadvertently introduce new forms of bias or alter the structural integrity of the dataset yielding to the butterfly effect [32]. Careful calibration is therefore necessary when applying such techniques, particularly in fairness-sensitive domains. These results support recent work showing the effectiveness of probabilistic approaches in tabular domains with complex feature types [33]. 1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599

In contrast, deep generative models—particularly TVAE and DECAF—struggled to preserve feature interdependencies. Although TVAE produced reasonably realistic marginal distributions, it exhibited the highest multivariate deviations in both datasets, highlighting limitations in capturing complex relational structures. This fidelity degradation is particularly concerning given the relatively modest size of the CVD dataset (10,000 instances), especially considering that in healthcare contexts, 1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610

1611 larger datasets are often unattainable due to strict privacy regulations or limited  
1612 resources hence the need to data generation. One might expect that with larger  
1613 datasets, where overfitting risks are reduced, models would generalise relational struc-  
1614 tures more effectively. However, our findings indicate otherwise: even at these moderate  
1615 sample sizes, deep generative models failed to robustly reconstruct multivariate depen-  
1616 dencies. Prior studies have suggested that generative models require substantially  
1617 larger and more balanced datasets to adequately capture multimodal feature inter-  
1618 actions [14]. Yet, our results show that differences in dataset size between CVD and  
1619 diabetes data (10,000 vs 1,879 instances) did not significantly influence fidelity out-  
1620 comes, suggesting that model architecture and learning dynamics, rather than dataset  
1621 size alone, critically limit fidelity in healthcare tabular data synthesis.

1622  
1623 In predictive utility, BayesBoost again led performance, matching or exceeding  
1624 real data benchmarks in some cases, particularly in CVD prediction. This means that  
1625 simpler models can outperform deep networks when diversity and structure preser-  
1626 vation matter. VAE showed the best performance among deep models, particularly  
1627 for diabetes, but CTGAN and CopulaGAN exhibited severe degradation, especially  
1628 for AUC-ROC and accuracy where classifier performance collapsed to near-random  
1629 levels. This significant deterioration likely stems from the inability of GAN-based  
1630 models to effectively handle the sparsity, heterogeneity, and small sample subgroup  
1631 distributions characteristic of tabular healthcare data without extensive tuning [11].

1632  
1633 Classifier choice also influenced results. RF consistently achieved strong and  
1634 balanced performance across both datasets, real and synthetic, reinforcing its suit-  
1635 ability for healthcare tasks involving complex, heterogeneous data types. In contrast,  
1636 NB underperformed overall but occasionally delivered higher recall; highlighting the  
1637 importance of selecting classifiers based on clinical priorities such as sensitivity ver-  
1638 sus specificity. This observation underscores the need for a context-driven model  
1639 evaluation approach rather than a reliance on a single global performance metric.

1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656

Fairness, however, remains the most pressing concern since algorithmic decisions can directly influence patient outcome. Our findings reveal that fairness imbalances already exist in the real datasets—particularly for equalised odds and PRP—. Rather than mitigating these imbalances, generative models exacerbated them. Moreover, across both datasets, all generative models introduced extreme fairness violations primarily affecting the most underrepresented groups in the real data. These minority groups, suffered the largest distortions in fairness metrics post data generation. This systematic deterioration highlights that models not only failed to address existing disparities but actively amplified them—a failure that could seriously undermine equity in clinical AI systems if unchecked.

CopulaGAN, CTGAN, and DECAF consistently introduced severe fairness violations across both datasets. In the diabetes data, PRP and equalised odds ratios inflated well above 2 or collapsed to near-zero. DECAF performed worst, producing near zero predictions for some subgroups, even among majority populations. Such extreme distortions likely stem from these models' failure to adequately reproduce the full distribution of sensitive subgroups within the generated data.

BayesBoost exhibited strong and consistent fairness preservation, particularly in predictive rate parity. When paired with RF classifiers, BayesBoost maintained PRP values remarkably close to 1 across all ethnic groups, often improving upon the fairness observed in real data. However, its performance on equalised odds was less stable—particularly when combined with Naive Bayes—resulting in inflated parity ratios for certain subgroups. These findings suggest that while BayesBoost is highly effective in correcting representation bias and maintaining overall class balance and preserving data structure, it still struggle with fairness metrics tied to subgroup-specific decision thresholds, which require finer calibration.

VAE, while weaker in utility and structural fidelity, preserved EO best across all generative models. Its architecture, based on minimising instance-level reconstruction

1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702

1703 loss rather than adversarial optimisation, may inherently support more equitable fea-  
1704 ture distributions. This aligns with previous work highlighting the relative fairness  
1705 stability of VAEs in healthcare settings [].

1706  
1707  
1708 These results reinforce the need for synthetic data pipelines to integrate robust,  
1709 multi-dimensional evaluation frameworks. Utility or fidelity alone is insufficient: fair-  
1710 ness must be explicitly assessed, audited, and optimised. Model–classifier interactions  
1711 should be systematically tested, as fairness and utility outcomes can shift dramatically  
1712 depending on pairing. Above all, synthetic data systems must disclose known limita-  
1713 tions—including subgroup risks—and define safe use cases where fairness or fidelity  
1714 thresholds are met.

1715  
1716  
1717  
1718  
1719  
1720

#### 1721 **4.1 Towards Robust and Accountable Synthetic Data**

1722

##### 1723 **Generation in Healthcare**

1724

1725 The findings of this study underscore the urgent need for a structured, multi-  
1726 dimensional framework when developing machine learning-based synthetic data gen-  
1727 eration systems—particularly in fairness-critical domains such as healthcare. Rather  
1728 than relying on a single generative approach or isolated metric, we advocate for an  
1729 iterative, modular pipeline that integrates multiple synthetic data generation strate-  
1730 gies—such as probabilistic models, GAN-based methods, and autoencoders—and  
1731 evaluates them holistically across three core dimensions: statistical fidelity, predictive  
1732 utility, and fairness.

1733  
1734  
1735  
1736  
1737  
1738 Crucially, this framework must be cyclic and adaptive. If a particular generative  
1739 model fails to meet acceptable thresholds for fidelity or fairness, it should trigger a  
1740 revision of the generation process—be it through reparameterisation, retraining, or  
1741 the integration of fairness-aware learning objectives. Furthermore, to capture the full  
1742 range of model behaviour, multiple classification models (e.g. RF, LR, NB) should  
1743 be paired systematically with each synthetic data generator. This is essential to  
1744  
1745  
1746  
1747  
1748

identify performance variations that may arise due to interactions between the inductive biases of classifiers and the structural properties of the synthetic data. Certain model-generator combinations may obscure or amplify disparities, and only by evaluating them together can we ensure robust, generalisable conclusions about the utility and fairness of synthetic datasets.

While current evaluation practices largely focus on statistical similarity and classifier performance, they often lack explicit thresholds for determining when fairness violations become unacceptable or clinically dangerous. Our findings reveal that even models preserving overall fidelity or utility may introduce extreme subgroup disparities in precision, recall, or decision thresholds. To address this, the field needs standardised, context-aware fairness deviation metrics—such as acceptable tolerance bands around a baseline (e.g., parity ratios between 0.9 and 1.1)—tailored to specific applications and aligned with ethical guidelines. This would enable developers and regulators to distinguish between benign and harmful deviations, particularly in high-stakes settings like diagnostics or risk stratification.

Beyond traditional metrics, we also recommend incorporating causal structure validation into synthetic data evaluation. While existing fairness metrics assess statistical parity, they may miss latent structural shifts that alter the interpretability or utility of synthetic data. Approaches such as structural causal models (SCMs) or conditional independence testing can help assess whether synthetic data preserves the underlying causal pathways of real-world variables—essential for ensuring synthetic data supports safe and valid inference [34, 35].

Finally, we stress that synthetic data systems must aim to address both representation bias (e.g. over- or under-generation of sensitive subgroups) and algorithmic performance fairness. Our results suggest that most current generative models struggle to satisfy both simultaneously, with some amplifying subgroup disparities. Consequently, when a model cannot meet both objectives, its limitations must be explicitly

1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794

1795 documented—including subgroup-specific fairness risks and clearly defined safe use  
1796 cases. Transparent reporting of synthetic data limitations should be treated as  
1797 essential, not optional, particularly in regulated domains such as healthcare and  
1798 finance.  
1800

1801

1802

## 1803 5 Conclusion

1804

1805

1806 This study demonstrates that while synthetic data generation holds significant  
1807 promise for healthcare applications, achieving robust fidelity, predictive utility, and  
1808 fairness remains a substantial challenge. By evaluating five representative mod-  
1809 els—BayesBoost, CTGAN, TVAE, CopulaGAN, and DECAF—on two healthcare  
1810 datasets, we provided an integrated comparison across statistical, predictive, and  
1811 fairness dimensions.  
1812

1813

1814

1815 Among the evaluated methods, BayesBoost, particularly when combined with Ran-  
1816 dom Forest, consistently delivered the most reliable outcomes across all dimensions,  
1817 preserving clinical validity while mitigating subgroup disparities. VAE also emerged  
1818 as a valuable alternative among deep generative models, achieving superior fairness  
1819 preservation, particularly for equalised odds, albeit at some cost to overall utility.  
1820

1821

1822

1823 Conversely, GAN-based models frequently amplified existing biases and failed to  
1824 replicate critical multivariate dependencies, raising serious concerns regarding their  
1825 deployment in fairness-sensitive healthcare settings. The pronounced fairness viola-  
1826 tions among underrepresented groups underscore that fidelity and predictive accuracy  
1827 alone are insufficient markers of synthetic data quality.  
1828

1829

1830

1831

1832

1833

1834 Overall, the strength of this work lies in its unified, multi-criteria evaluation frame-  
1835 work that jointly examines fidelity, utility, and fairness—dimensions often studied in  
1836 isolation. The findings highlight the importance of incorporating fairness auditing as  
1837 a core, mandatory component of synthetic data pipelines, evaluated alongside tra-  
1838 ditional utility and fidelity metrics. The proposed multi-criteria framework provides  
1839

1840

a practical foundation for the responsible selection and benchmarking of synthetic data models for healthcare applications, emphasising the importance of structured evaluation frameworks and transparency regarding model limitations as prerequisites for safe and ethical deployment. Until such frameworks are standardised and rigorously applied, the clinical use of synthetic data—particularly for fairness-sensitive applications—should proceed only under strict regulatory oversight and with explicit subgroup-specific validation.

### 5.1 Limitations and Future Work

This study assumes that the datasets used (CVD and Diabetes) are representative of broader healthcare data structures and that model performance on these datasets can generalise to similar tabular healthcare domains. Although the evaluation covers key indicators of fidelity, predictive utility, and fairness, it remains limited to two datasets and a selection of five established models. The analysis does not include emerging generative paradigms such as diffusion or transformer-based models, nor does it account for longitudinal or multi-modal data, as it focuses exclusively on tabular structures.

Future work should extend these evaluations to larger and more heterogeneous datasets and explore recent generative approaches such as diffusion and transformer-based architectures to further validate and generalise the findings. Additionally, translating the proposed framework into practical bias-auditing and model certification tools could support its application in clinical trial simulation, digital health research, and regulatory evaluation of synthetic data systems.

### Acknowledgements.

1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886

## 1887 **Declarations**

1888

1889 • Funding: This work was funded by the Regulators Pioneer Fund, Department for  
1890

1891 Science, Innovation and Technology. This work was also supported by the UK Reg-  
1892  
1893 ulatory Science and Innovation Networks – Implementation Phase: Human Health

1894 CERSIs programme through the project RADIANT: Regulatory Science Empower-  
1895

1896 ing Innovation in Transformative Digital Health and AI (Grant Ref: MCPC24031),  
1897  
1898 funded by the Medical Research Council (MRC) and Innovate UK.

1899 • Conflict of interest/Competing interests : Not applicable  
1900

1901 • Ethics approval and consent to participate For the use of CPRD data to generate  
1902  
1903 synthetic data: this was covered by CPRD’s Database Research Ethics Approval  
1904  
1905 (IRAS number: 242149)

1906 • Consent for publication: Not applicable

1907 • Data availability: CPRD cardiovascular disease synthetic dataset used in this  
1908  
1909 paper can be requested from CPRD ([https://cprd.com/cprd-cardiovascular-](https://cprd.com/cprd-cardiovascular-disease-synthetic-dataset)  
1910  
1911 [disease-synthetic-dataset](https://cprd.com/cprd-cardiovascular-disease-synthetic-dataset)). The diabetes dataset is publicly available on Kaggle  
1912  
1913 ([https://www.kaggle.com/datasets/rabieelkharoua/diabetes-health-dataset-](https://www.kaggle.com/datasets/rabieelkharoua/diabetes-health-dataset-analysis)  
1914  
1915 [analysis](https://www.kaggle.com/datasets/rabieelkharoua/diabetes-health-dataset-analysis))

1916 • Materials availability: Not applicable

1917 • Code availability : Not applicable  
1918

1919 • Author contribution: P.M. and R.B. initiated the project and conceived the overall  
1920  
1921 study design. A.T. supervised the project and experiments, providing critical input  
1922  
1923 on result interpretation and the discussion of their implications. D.A. designed the  
1924  
1925 methodology, carried out the experiments, performed the data analysis, and served  
1926  
1927 as the lead author of the manuscript. B.D contributed expertise on fairness metrics  
1928  
1929 and supported the implementation of the BayesBoost model. All authors contributed  
1930  
1931 to revising the manuscript and approved the final version for submission.  
1932

## References

- [1] Jordon, J., Szpruch, L., Houssiau, F., Bottarelli, M., Cherubin, G., al.: Synthetic Data – what, why and how? Preprint at arXiv:2205.03257 (2022). <https://doi.org/10.48550/arXiv.2205.03257>
- [2] Quintana, D.S.: A synthetic dataset primer for the biobehavioural sciences to promote reproducibility and hypothesis generation. *Elife* **9**, 53275 (2020)
- [3] Hernandez, M., Epelde, G., Alberdi, A., Cilla, R., Rankin, D.: Synthetic data generation for tabular health records: A systematic review. *Neurocomputing* **493**, 28–45 (2022)
- [4] Harper, P.R., Moore, J.W., Woolley, T.E.: Covid-19 transmission modelling of students returning home from university. *Health Systems* **10**(1), 31–40 (2021)
- [5] Caramelo, F., Ferreira, N., Oliveiros, B.: Estimation of risk factors for covid-19 mortality-preliminary results. *MedRxiv*, 2020–02 (2020)
- [6] Wang, Z., Myles, P., Tucker, A.: Generating and evaluating synthetic uk primary care data: preserving data utility & patient privacy. In: 2019 IEEE 32nd International Symposium on Computer-Based Medical Systems (CBMS), pp. 126–131 (2019). IEEE
- [7] Wang, Z., Myles, P., Tucker, A.: Generating and evaluating cross-sectional synthetic electronic healthcare data: Preserving data utility and patient privacy. *Computational Intelligence* **37**(2), 819–851 (2021)
- [8] Draghi, B., Wang, Z., Myles, P., Tucker, A.: Bayesboost: Identifying and handling bias using synthetic data generators. In: Third International Workshop on Learning with Imbalanced Domains: Theory and Applications, pp. 49–62 (2021). PMLR

1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978

- 1979 [9] Yadav, P., Gaur, M., Madhukar, R.K., Verma, G., Kumar, P.: Rigorous exper-  
1980 imental analysis of tabular data generated using tvae and ctgan. *International*  
1981 *Journal of Advanced Computer Science & Applications* **15**(4) (2024)  
1982  
1983  
1984  
1985 [10] Marecha, P., Ye, L.: Generation and evaluation of tabular data in different  
1986 domains using gans. *Asian Journal of Research in Computer Science* **16**, 15–27  
1987 (2023) <https://doi.org/10.9734/ajrcos/2023/v16i1331>  
1988  
1989  
1990 [11] Miletic, M., Sariyar, M.: Challenges of using synthetic data generation methods  
1991 for tabular microdata. *Applied Sciences* **14**(14), 5975 (2024)  
1992  
1993  
1994 [12] Van Breugel, B., Kyono, T., Berrevoets, J., Schaar, M.: Decaf: Generating fair  
1995 synthetic data using causally-aware generative networks. *Advances in Neural*  
1996 *Information Processing Systems* **34**, 22221–22233 (2021)  
1997  
1998  
1999  
2000 [13] Liu, T., Qian, Z., Berrevoets, J., Schaar, M.: Goggle: Generative modelling for  
2001 tabular data by learning relational structure. In: *The Eleventh International*  
2002 *Conference on Learning Representations* (2023)  
2003  
2004  
2005 [14] Wang, A.X., Chukova, S.S., Simpson, C.R., Nguyen, B.P.: Challenges and oppor-  
2006 tunities of generative models on tabular data. *Applied Soft Computing*, 112223  
2007 (2024)  
2008  
2009  
2010 [15] Stoian, M.C., Dymishi, S., Cordy, M., Lukasiewicz, T., Giunchiglia, E.: How  
2011 realistic is your synthetic data? constraining deep generative models for tabular  
2012 data. arXiv preprint arXiv:2402.04823 (2024)  
2013  
2014  
2015 [16] Nik, A.H.Z., Riegler, M.A., Halvorsen, P., Storås, A.M.: Generation of synthetic  
2016 tabular healthcare data using generative adversarial networks. In: *International*  
2017 *Conference on Multimedia Modeling*, pp. 434–446 (2023). Springer  
2018  
2019  
2020  
2021  
2022  
2023  
2024

- [17] Pezoulas, V.C., Zaridis, D.I., Mylona, E., Androutsos, C., Apostolidis, K., Tachos, N.S., Fotiadis, D.I.: Synthetic data generation methods in healthcare: A review on open-source tools and methods. *Computational and Structural Biotechnology Journal* (2024) 2025  
2026  
2027  
2028  
2029  
2030  
2031
- [18] Zand, R., Abedi, V., Hontecillas, R., Lu, P., Noorbakhsh-Sabet, N., Verma, M., Leber, A., Tubau-Juni, N., Bassaganya-Riera, J.: Development of synthetic patient populations and in silico clinical trials. *Accelerated Path to Cures*, 57–77 (2018) 2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039
- [19] Pappalardo, F., Russo, G., Tshinanu, F.M., Viceconti, M.: In silico clinical trials: concepts and early adoptions. *Briefings in bioinformatics* **20**(5), 1699–1708 (2019) 2040  
2041  
2042  
2043
- [20] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. *Communications of the ACM* **63**(11), 139–144 (2020) 2044  
2045  
2046  
2047  
2048  
2049
- [21] Xu, L., Skoularidou, M., Cuesta-Infante, A., Veeramachaneni, K.: Modeling tabular data using conditional gan. *Advances in neural information processing systems* **32** (2019) 2050  
2051  
2052  
2053  
2054  
2055
- [22] Shahul Hameed, M.A., Qureshi, A.M., Kaushik, A.: Bias mitigation via synthetic data generation: A review. *Electronics* (2079-9292) **13**(19) (2024) 2056  
2057  
2058  
2059
- [23] Jain, A., Brooks, J.R., Alford, C.C., Chang, C.S., Mueller, N.M., Umscheid, C.A., Bierman, A.S.: Awareness of racial and ethnic bias and potential solutions to address bias with use of health care algorithms. In: *JAMA Health Forum*, vol. 4, pp. 231197–231197 (2023). American Medical Association 2060  
2061  
2062  
2063  
2064  
2065  
2066
- [24] Dwork, C., Hardt, M., Pitassi, T., Reingold, O., Zemel, R.: Fairness through awareness. In: *Proceedings of the 3rd Innovations in Theoretical Computer* 2067  
2068  
2069  
2070

- 2071 Science Conference, pp. 214–226 (2012)  
2072
- 2073 [25] Chouldechova, A.: Fair prediction with disparate impact: A study of bias in  
2074 recidivism prediction instruments. *Big data* **5**(2), 153–163 (2017)  
2075  
2076
- 2077 [26] Hardt, M., Price, E., Srebro, N.: Equality of opportunity in supervised learning.  
2078 *Advances in neural information processing systems* **29** (2016)  
2079  
2080
- 2081 [27] Verma, S., Rubin, J.: Fairness definitions explained. In: *Proceedings of the*  
2082 *International Workshop on Software Fairness*, pp. 1–7 (2018)  
2083  
2084
- 2085 [28] Datalink, C.P.R.: CPRD cardiovascular disease synthetic dataset (Version  
2086 2020.06.001) [Data set]. *Clinical Practice Research Datalink* (2020). <https://doi.org/10.11581/YK6N-B652> . <https://doi.org/10.11581/YK6N-B652>  
2087  
2088  
2089  
2090
- 2091 [29] kharoua, R.E.: *Diabetes Health Dataset Analysis*. Kaggle (2024). <https://doi.org/10.34740/KAGGLE/DSV/8665939> . <https://www.kaggle.com/dsv/8665939>  
2092  
2093  
2094
- 2095 [30] Fusion, B.: *GeNIe Modeller*. *Bayes Fusion* (2025). <https://www.bayesfusion.com/>  
2096  
2097
- 2098 [31] Alattal, D.R., Wang, Z., Myles, P., Tucker, A.: Creating synthetic geospatial  
2099 patient data to mimic real data whilst preserving privacy: \*2022 35th interna-  
2100 tional symposium on computer-based medical systems (cbms). In: *2023 IEEE*  
2101 *36th International Symposium on Computer-Based Medical Systems (CBMS)*,  
2102 pp. 7–12 (2023). <https://doi.org/10.1109/CBMS58004.2023.00183>  
2103  
2104  
2105  
2106
- 2107 [32] Ferrara, E.: The butterfly effect in artificial intelligence systems: Implications for  
2108 ai bias and fairness. *Machine Learning with Applications* **15**, 100525 (2024)  
2109  
2110
- 2111 [33] Draghi, B., Wang, Z., Myles, P., Tucker, A.: Identifying and handling data bias  
2112 within primary healthcare data using synthetic data generators. *Heliyon* **10**(2)  
2113  
2114  
2115  
2116

- [34] Makhlouf, K., Zhioua, S., Palamidessi, C.: When causality meets fairness: A survey. *Journal of Logical and Algebraic Methods in Programming*, 101000 (2024)
- [35] Binkyte-Sadauskiene, R., Makhlouf, K., Pinzón, C., Zhioua, S., Palamidessi, C.: Causal discovery for fairness. *CoRR abs/2206.06685* (2022)
- 2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159  
2160  
2161  
2162