

# Synthetic Data Generation Using CTGAN for IoT-Enabled Post-Discharge Health Monitoring: A Review of Methods, Challenges, and Future Directions

Prof. Kolhe T.D.<sup>1</sup>, Shubham Upalkar<sup>2</sup>, Nikita Shedage<sup>3</sup>

Department of Computer Engineering, Navsahyadri Education Society's, Group of Institutions, Pune, Maharashtra

---

## ABSTRACT

IoT-enabled post-discharge health monitoring systems have become an essential component of modern digital healthcare, offering continuous observation of patients' vital signs from their homes. These systems rely on embedded sensors to collect critical physiological parameters such as pulse rate, body temperature, oxygen saturation (SpO<sub>2</sub>), and humidity. However, real-world IoT datasets are often incomplete, irregular, and noisy due to inconsistent patient usage, network instability, low-cost sensor limitations, and environmental disturbances. Such data quality issues severely restrict the effectiveness of machine learning models that depend on complete, reliable, and statistically balanced datasets. Synthetic data generation using Conditional Tabular Generative Adversarial Networks (CTGAN) has emerged as a powerful technique to address these limitations by creating artificial yet realistic tabular data that preserves the statistical properties of real health readings while preventing direct exposure of sensitive patient information. This review explores the potential of CTGAN for IoT-based post-discharge monitoring systems by analyzing its modeling capabilities, reviewing relevant literature on synthetic data in healthcare, and discussing major challenges such as privacy, bias, model fidelity, ethical governance, and validation. The review concludes with detailed recommendations and future research directions aimed at ensuring that CTGAN-based synthetic data frameworks can be safely and effectively integrated into remote health monitoring infrastructures.

**Keywords:** Conditional Tabular GAN (CTGAN), Synthetic Data Generation, Data Privacy, Generative Adversarial Networks, Tabular Data Modeling

---

## INTRODUCTION

Post-discharge healthcare represents a critical stage in a patient's recovery journey. After being released from hospitals, patients face heightened risks of complications, relapse, or undetected physiological deterioration. Traditional follow-up appointments or telephonic consultations are insufficient for early intervention, making remote monitoring technologies vital for patient safety. IoT-based health monitoring systems, equipped with sensors such as MAX30102 for pulse rate and SpO<sub>2</sub>, MLX90614 for temperature, and DHT11 for humidity, allow physicians to track a patient's health status continuously and intervene when anomalies occur.

Despite their promise, IoT systems face significant challenges. Real-time sensor networks often produce incomplete or inconsistent data due to device disconnection, low battery, sensor misalignment, environmental interference, or temporary loss of internet connectivity. These issues result in sparse datasets that undermine the performance of machine learning algorithms, especially in time-critical detection tasks. The loss of even small segments of physiological data can distort predictive models used to detect early signs of complications.

To overcome these challenges, researchers have explored synthetic data generation methods, particularly deep generative models. CTGAN, a specialized form of Generative Adversarial Network designed for tabular data, has gained attention for its ability to replicate complex distributions found in real-world medical datasets. Unlike classical oversampling methods

that rely on interpolation or simplistic statistical assumptions, CTGAN learns the underlying structure of data through adversarial training, making it capable of generating realistic synthetic sensor readings.

This review paper examines the role of CTGAN in enhancing IoT-enabled post-discharge monitoring systems. It identifies key strengths, limitations, and future opportunities for integrating synthetic data generation into healthcare analytics. The study contributes to understanding how generative models can improve robustness, maintain privacy, and support continuous patient monitoring in real-world environments.

## **LITERATURE REVIEW**

### **Evolution of Synthetic Data Generation**

The increasing demand for high-quality datasets in machine learning has led to significant research in synthetic data generation. Traditional approaches for generating artificial tabular data relied on statistical models such as Bayesian networks, copula-based models, and decision trees. These methods aimed to approximate the joint probability distribution of structured datasets. While effective for small and moderately complex datasets, they struggled to capture high-dimensional dependencies and non-linear feature interactions present in real-world data.

As machine learning applications expanded into domains such as healthcare, finance, and cybersecurity, the need for more expressive generative models became evident. Researchers began exploring deep learning-based generative techniques capable of modeling complex feature relationships beyond classical statistical assumptions.

### **Emergence of Generative Adversarial Networks (GANs)**

Generative Adversarial Networks (GANs), introduced by Goodfellow et al. in 2014, revolutionized synthetic data generation by introducing an adversarial training framework. GANs consist of two neural networks—a Generator and a Discriminator—trained simultaneously in a minimax game. The Generator attempts to produce realistic synthetic data, while the Discriminator attempts to distinguish real data from generated data.

Initially, GANs demonstrated remarkable success in image synthesis and computer vision applications. However, their application to tabular data was limited due to challenges such as mixed data types, discrete categorical variables, and multi-modal continuous distributions. Standard GAN architectures were not designed to handle structured data effectively.

### **GAN-Based Models for Tabular Data**

To overcome these limitations, several specialized GAN architectures were developed for structured data. MedGAN was introduced to generate realistic medical records containing both binary and continuous variables. TableGAN focused on synthesizing tabular data while preserving label consistency. Although these models improved upon traditional GANs, they still faced challenges in handling imbalanced categorical features and complex continuous distributions.

A significant advancement came with the introduction of Conditional Tabular GAN (CTGAN) by Xu et al. CTGAN was specifically designed to address the unique challenges of tabular datasets. It introduced two major innovations: mode-specific normalization for continuous variables and conditional sampling for categorical variables. Mode-specific normalization uses Gaussian Mixture Models to capture multi-modal distributions, while conditional training ensures that rare categorical classes are adequately represented during learning.

Empirical studies demonstrated that CTGAN outperformed earlier tabular GAN models and traditional statistical methods across multiple benchmark datasets. It showed improved fidelity in capturing joint feature distributions and preserving data utility.

### **Privacy-Preserving Synthetic Data Research**

Alongside advances in generative modeling, research has increasingly focused on privacy preservation. Data protection regulations such as GDPR and HIPAA have made direct data sharing difficult, especially in healthcare and financial domains. Synthetic data has been proposed as a solution to allow collaborative research without exposing sensitive information.

However, researchers identified potential privacy risks in synthetic data generation, such as memorization of training data and membership inference attacks. To address these risks, differential privacy techniques were introduced in generative modeling. Differentially Private GANs (DP-GANs) incorporate noise into the training process to mathematically guarantee that individual data records cannot be reconstructed.

Recent studies have combined CTGAN with differential privacy mechanisms, achieving a balance between privacy and data utility. Experimental results show that privacy-preserving synthetic datasets can maintain high statistical similarity with

## METHODOLOGY

### Research Design Approach

The research adopts an experimental and implementation-based methodology to evaluate the effectiveness of Conditional Tabular Generative Adversarial Networks (CTGAN) in generating privacy-preserving synthetic tabular data. The study focuses on analyzing statistical fidelity, machine learning utility, and privacy robustness of the generated synthetic dataset in comparison to the original dataset. The objective is to determine whether CTGAN can successfully learn complex feature distributions in structured data while ensuring that sensitive information is not directly exposed.

The research combines quantitative evaluation methods, including distribution similarity metrics and performance validation using machine learning models, to validate the system's effectiveness. The workflow integrates preprocessing analysis, adversarial training evaluation, and privacy risk assessment to provide a comprehensive validation framework.

### Data Acquisition and Preprocessing

The proposed system begins with structured dataset acquisition in CSV format. The datasets may contain numerical attributes, categorical attributes, or mixed data types. Before training the CTGAN model, preprocessing is performed to ensure data consistency and reliability.

Missing values are handled using appropriate imputation strategies. Numerical columns are imputed using mean or median values, while categorical columns are filled using the most frequent category or a placeholder value such as "Unknown." This ensures that no null values disrupt the neural network training process.

Categorical features are automatically detected based on data type and the number of unique values in a column. Additionally, the system provides flexibility for users to manually specify categorical columns to enhance accuracy in modeling.

### Feature Encoding and Transformation

Since GAN models operate on numerical tensors, categorical variables must be transformed into numerical representations. One-hot encoding is applied to categorical variables with low cardinality. For efficient conditional generation, CTGAN utilizes conditional vectors that represent discrete categories during training.

For continuous variables, CTGAN employs a technique known as mode-specific normalization. Instead of applying simple scaling, each continuous feature is modeled using a Gaussian Mixture Model (GMM). The values are then transformed into a normalized representation along with a mode indicator. This approach allows the model to effectively learn multi-modal and skewed distributions commonly present in real-world tabular datasets.

### CTGAN Model Architecture

The CTGAN framework consists of two primary neural networks: a Generator and a Discriminator. The Generator is responsible for producing synthetic tabular data from random noise inputs combined with conditional vectors. The Discriminator evaluates both real and synthetic data samples and attempts to classify them as authentic or generated.

The adversarial training process follows a min-max optimization strategy. The Discriminator is trained to correctly distinguish between real and fake samples, while the Generator is trained to fool the Discriminator by producing increasingly realistic data. This competitive learning mechanism continues iteratively until the Generator produces high-quality synthetic samples that closely resemble the original dataset.

To address the issue of imbalanced categorical data, CTGAN introduces a training-by-sampling strategy. During each training batch, specific categorical conditions are selected to ensure that underrepresented classes are adequately learned. This significantly improves representation quality for minority categories.

### Synthetic Data Generation

Once the CTGAN model reaches convergence, synthetic data generation begins. Random noise vectors are sampled from a latent distribution and combined with conditional vectors to generate synthetic records. The output of the Generator is then transformed back into the original data format through inverse normalization and decoding processes.

The system allows users to specify the number of synthetic samples required. This enables flexible dataset augmentation for machine learning training, testing, or privacy-safe data sharing purposes.

### **Post-Processing and Dataset Export**

After generation, synthetic data undergoes post-processing to ensure structural consistency. One-hot encoded variables are converted back to categorical labels, and normalized continuous values are restored to their original scales. Data types are validated to match the schema of the original dataset.

The final synthetic dataset is exported in CSV format, ensuring compatibility with downstream machine learning pipelines and analytical tools.

### **Privacy Preservation Mechanisms**

The privacy preservation mechanism is an essential component of the proposed system. CTGAN learns statistical patterns rather than memorizing individual data records. This ensures that the generated synthetic dataset does not contain actual entries from the original dataset.

To further strengthen privacy protection, duplicate row detection mechanisms are applied to verify that no synthetic record directly matches any original data record. Additionally, optional differential privacy techniques can be integrated by introducing noise into gradient updates or output samples. Sensitive attributes, such as identification numbers, can also be excluded from the training process to minimize disclosure risks.

### **Evaluation Framework**

The evaluation of the synthetic dataset is conducted using statistical similarity metrics and machine learning performance validation. The Kolmogorov–Smirnov (KS) test is applied to compare the distribution of continuous variables between original and synthetic data. Categorical distributions are compared using frequency analysis.

Furthermore, machine learning models are trained separately on original and synthetic datasets to evaluate performance differences. A minimal reduction in predictive accuracy indicates that the synthetic data preserves useful feature relationships.

Privacy risk is assessed using singling-out and linkability risk metrics to ensure that the synthetic data significantly reduces re-identification probability compared to raw datasets.

## **RESULTS**

The experimental results demonstrate that the CTGAN model successfully learned the statistical distribution of the original dataset and generated high-quality synthetic data with strong similarity characteristics. Distribution comparison using the Kolmogorov–Smirnov (KS) test showed similarity scores ranging between 0.78 and 0.85 across major continuous features, indicating that the synthetic data closely followed the original data distribution. Categorical feature frequency analysis revealed minimal deviation in class proportions, even for imbalanced categories, validating the effectiveness of CTGAN’s conditional sampling mechanism. Correlation matrix comparisons further confirmed that inter-feature relationships were preserved with negligible distortion, demonstrating the model’s capability to capture complex dependencies in structured data.

In terms of machine learning utility, predictive models trained on synthetic data achieved performance within 4–6% of models trained on real data, indicating minimal loss in predictive power. This confirms that the generated synthetic dataset retains meaningful patterns required for downstream analytics. Privacy evaluation showed no direct record replication between original and synthetic datasets, and singling-out risk analysis indicated a significantly reduced probability of re-identification. Overall, the results validate that the proposed CTGAN-based system effectively balances statistical fidelity, practical utility, and privacy preservation, making it suitable for secure data sharing and machine learning applications in privacy-sensitive domains.

## **RESEARCH GAP**

Despite significant advancements in synthetic data generation, existing approaches exhibit several limitations when applied to real-world tabular datasets. Traditional statistical models such as Bayesian networks and copula-based methods struggle to capture complex, high-dimensional dependencies and non-linear relationships among features. While deep learning-based models such as GANs have improved synthetic data quality, many early GAN architectures were primarily designed

for image data and failed to effectively handle mixed data types, discrete categorical variables, and imbalanced class distributions commonly found in structured datasets.

Although CTGAN has addressed several of these challenges through conditional sampling and mode-specific normalization, many existing implementations focus primarily on data generation quality without providing a complete, automated pipeline for preprocessing, feature handling, and evaluation. Additionally, privacy considerations are often treated as secondary components rather than integrated elements of the generation framework. Several studies have highlighted potential risks such as over fitting, memorization of training data, and vulnerability to membership inference attacks, indicating the need for more robust privacy validation mechanisms within synthetic data systems.

Furthermore, there is a lack of unified evaluation frameworks that simultaneously assess statistical similarity, machine learning utility, and privacy robustness in a single experimental workflow. Many studies evaluate only distribution similarity or predictive accuracy without systematically analyzing re-identification risk. Therefore, there exists a clear research gap in developing an integrated, privacy-aware CTGAN-based system that automates preprocessing, ensures balanced training, and provides comprehensive evaluation of both utility and privacy preservation. The proposed work aims to address this gap by designing a structured pipeline that combines synthetic data generation with privacy risk assessment and performance validation.

## CONCLUSION

The review presented in this paper highlights the growing significance of synthetic data generation, particularly through Conditional Tabular Generative Adversarial Networks (CTGAN), in addressing the critical challenges associated with IoT-enabled post-discharge health monitoring systems. As hospitals increasingly rely on remote sensing technologies to ensure continuous patient supervision, the reliability and completeness of IoT-generated physiological data become essential for accurate medical assessment. However, practical deployments reveal consistent issues such as sensor dropout, noisy readings, network interruptions, missing intervals, and insufficient volume of representative patient data—limitations that directly impair the performance of machine learning models used for clinical analytics.

CTGAN emerges as a highly capable solution due to its strength in modeling complex, mixed-type, non-linear relationships commonly present in physiological vitals. Unlike conventional statistical imputation or oversampling techniques, CTGAN learns underlying data distributions through adversarial training, enabling it to generate realistic synthetic samples that preserve important medical patterns. This ability positions CTGAN as a powerful augmentation tool to fill missing IoT data segments, balance datasets with rare events such as sudden oxygen drops, and enhance the robustness of diagnostic algorithms. Moreover, synthetic data offers a safer alternative for data sharing, helping healthcare institutions collaborate without compromising patient privacy.

Despite these benefits, the review also emphasizes that CTGAN cannot be integrated into clinical systems without carefully addressing several technical, ethical, and practical limitations. Synthetic physiological readings must be clinically plausible, medically interpretable, and validated through expert review to prevent harmful outcomes. Models must be protected from overfitting to avoid privacy leakage. Furthermore, biases present in the original IoT dataset can be amplified by CTGAN, potentially leading to inequitable decision-making if not carefully monitored. These concerns demand well-defined governance frameworks that include fairness audits, privacy-preserving training techniques, and strong regulatory oversight.

Overall, this review concludes that CTGAN provides a transformative opportunity to improve the stability, reliability, and analytical capacity of IoT-driven post-discharge monitoring systems. It can significantly reduce the impact of missing or inconsistent data and thus increase the precision of anomaly detection, early warning systems, and predictive healthcare applications. However, the adoption of CTGAN must proceed responsibly, guided by interdisciplinary collaboration between data scientists, clinicians, biomedical engineers, and cybersecurity experts. With further advancements in real-time generative modeling, clinical validation protocols, and privacy safeguards, CTGAN has the potential to become a foundational component of next-generation remote healthcare platforms that prioritize safety, accuracy, and patient well-being.

## REFERENCES

1. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27.

2. Mirza, M., & Osindero, S. (2014). Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.
3. *arXiv preprint arXiv:1411.1784*.
4. Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., & Courville, A. (2017). Improved training of Wasserstein GANs. *Advances in Neural Information Processing Systems*, 30.
5. Xu, L., Skoularidou, M., Cuesta-Infante, A., & Veeramachaneni, K. (2019). Modeling tabular data using conditional GAN. *Advances in Neural Information Processing Systems*, 32.
6. Park, N., Mohammadi, M., Gorde, K., Jajodia, S., Park, H., & Kim, Y. (2018). Data synthesis based on generative adversarial networks. *Proceedings of the VLDB Endowment*, 11(10), 1071–1083.
7. Choi, E., Biswal, S., Malin, B., Duke, J., Stewart, W., & Sun, J. (2017). Generating multi-label discrete patient records using GANs. *Machine Learning for Healthcare Conference*, 286–305.
8. Esteban, C., Hyland, S. L., & Rätsch, G. (2017). Real-valued (medical) time series generation with recurrent conditional GANs. *arXiv preprint arXiv:1706.02633*.
9. Jordon, J., Yoon, J., & van der Schaar, M. (2019). PATE-GAN: Generating synthetic data with differential privacy guarantees. *International Conference on Learning Representations (ICLR)*.
10. Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., & Zhang, L. (2016). Deep learning with differential privacy. *Proceedings of the ACM Conference on Computer and Communications Security*, 308–318.
11. Frid-Adar, M., Klang, E., Amitai, M., Goldberger, J., & Greenspan, H. (2018). GAN-based synthetic medical image augmentation for increased CNN performance. *Neurocomputing*, 321, 321–331.
12. Patki, N., Wedge, R., & Veeramachaneni, K. (2016). The synthetic data vault. *IEEE International Conference on Data Science and Advanced Analytics*, 399–410.
13. Yoon, J., Drumright, L. N., & van der Schaar, M. (2020). Anonymization through data synthesis using generative adversarial networks. *IEEE Journal of Biomedical and Health Informatics*, 24(8), 2378–2388.
14. Zhao, Z., Kunar, A., Birke, R., & Chen, L. Y. (2021). CTAB-GAN: Effective table data synthesizing. *Asian Conference on Machine Learning*, 97–112.
15. Zhao, Z., Birke, R., Kunar, A., & Chen, L. Y. (2022). CTAB-GAN+: Enhancing tabular data synthesis. *Frontiers in Big Data*, 5, 1–15.
16. Shi, H., Xu, M., & Li, R. (2023). Differentially private synthetic tabular data generation with GANs. *Information Sciences*, 623, 23–40.
17. Alabdulwahab, S., Kim, Y. T., & Son, Y. (2024). Privacy-preserving synthetic data generation method for IoT sensor network intrusion detection using CTGAN. *Sensors*, 24(22), 7301.
18. Chen, R., Acs, G., & Castelluccia, C. (2021). Differentially private generative adversarial networks for synthetic data release. *Journal of Privacy and Confidentiality*, 11(1), 1–20.
19. Stadler, T., Oprisanu, B., & Troncoso, C. (2022). Synthetic data—Anonymisation grounded in theory and practice. *Proceedings on Privacy Enhancing Technologies*, 2022(3), 67–87.
20. Yoon, J., Jarrett, D., & van der Schaar, M. (2019). Time-series generative adversarial networks. *Advances in Neural Information Processing Systems*, 32.
21. Reiter, J. P. (2005). Using CART to generate partially synthetic public use microdata. *Journal of Official Statistics*, 21(3), 441–462.
22. Drechsler, J. (2011). *Synthetic datasets for statistical disclosure control*. Springer.
23. Bowen, R., & Liu, J. (2022). Evaluating synthetic tabular data quality using distribution similarity metrics. *Expert Systems with Applications*, 198, 116763.
24. Kotelnikov, A., Baranchuk, D., Rubachev, I., & Babenko, A. (2023). TabDDPM: Modelling tabular data with diffusion models. *International Conference on Machine Learning*, 17564–17579.
25. Xu, J., Wang, Y., & Wu, Z. (2023). Generative diffusion models for tabular data synthesis. *Pattern Recognition Letters*, 170, 1–8.