

Replication of LLM-Based Tabular Data Augmentation for Cognitive Biomarker Prediction

Aaron Dumont

Department of Computer Science, Tulane University

November 11, 2025

Abstract

Clinical research is often constrained by limited sample sizes, particularly when studying specific populations like "SuperAgers" or early-stage Alzheimer's patients. This study replicates the methodology of Lee et al. [1] to assess the efficacy of Large Language Model (LLM) based data augmentation in data-scarce scenarios. Using a proxy Alzheimer's disease dataset restricted to a total available sample of $N = 40$, we generated synthetic patient profiles using the GReaT (Generation of Realistic Tabular Data) framework. By fine-tuning a DistilGPT-2 Transformer and employing a conditional sampling strategy, we generated a class-balanced synthetic dataset. Our results demonstrate that augmenting the real data with synthetic samples improved the Random Forest classifier's Accuracy by 26% (0.62 to 0.88) and AUC by 4% (0.83 to 0.87). These findings validate the hypothesis that generative augmentation can effectively recover signal and regularize decision boundaries in small-sample clinical studies.

1 Introduction

The transition from reactive to predictive precision medicine relies heavily on biomarkers to forecast health trajectories. However, a pervasive challenge in this domain is the scarcity of high-quality, labeled clinical data. As highlighted by Lee et al. [1], clinical cohorts for specialized groups often consist of fewer than 100 subjects ($n = 81$ in their study). Such small sample sizes are insufficient for modern machine learning algorithms to generalize without significant risk of overfitting.

This project investigates a novel solution to the "small n " problem: using Large Language Models (LLMs) to generate synthetic tabular data. Unlike traditional oversampling techniques (e.g., SMOTE [5]) which interpolate in geometric space, LLM-based methods like GReaT treat patient data as textual sequences, capturing complex, non-linear feature dependencies [2]. The objective of this study is to replicate the data augmentation pipeline of Lee et al. [1] on an independent Alzheimer's dataset to verify if synthetic data generation can recover predictive performance in an extreme data-scarce environment ($N = 40$).

2 Related Work

The primary basis for this replication is Lee et al. [1], who utilized the GReaT framework to enhance a dataset of $n = 81$ subjects for SuperAger prediction using blood biomarkers. They demonstrated that synthetic augmentation combined with feature selection (BORUTA) significantly improved model accuracy.

In the broader context of neuroimaging, Siegel et al. [3] demonstrated that deep learning models (CNNs and Transformers) typically require massive datasets ($N > 40,000$) to achieve stability in Brain Age prediction. This contrast highlights the critical need for augmentation in biochemical studies where N is orders of magnitude smaller. Additionally, Zhou et al. [4] emphasized the importance of integrating heterogeneous data types (multi-omics). The flexible nature of Transformer-based generative models makes them uniquely suited for such complex data fusion tasks.

This approach represents a paradigm shift from traditional methods. While geometric oversampling techniques like SMOTE [6] have long been the standard for imbalance, and gradient boosting systems like XGBoost [6] are the state-of-the-art for tabular data, neither effectively solves the issue of extreme data scarcity ($N = 40$) as well as generative semantic augmentation.

3 Problem Setting

We define the problem as a binary classification task $y = f(x)$, where $x \in \mathbb{R}^d$ represents a vector of clinical and demographic biomarkers, and $y \in \{0, 1\}$ represents the diagnosis (Healthy vs. Alzheimer's).

We operate under the constraint of **Extreme Data Scarcity**, defined as a training set D_{real} where $|D_{real}| \approx d$. We selected a sample of 40 patients ($|D_{real}|$) and there were 33 variables (d). In this regime, discriminative models suffer from high variance. The goal is to learn a generative distribution $P_\theta(x, y)$ from D_{real} , sample a synthetic dataset D_{syn} from this distribution, and train a classifier on $D_{real} \cup D_{syn}$ such that generalization error is minimized.

4 Methodology

To ensure reproducibility, the experiment was conducted using the publicly available *Alzheimer's Disease*

Dataset (Kaggle). The code and resources are hosted in a public GitHub repository (<https://github.com/adumont2/CMPS-7010-Final-Presentation>).

4.1 Data Description and Preprocessing

The dataset consists of health records for 2,149 patients. After removing administrative identifiers (**PatientID**, **DoctorInCharge**), the feature space consisted of 33 variables spanning five domains:

- **Demographics:** Age, Gender, Ethnicity, Education.
- **Lifestyle/Vitals:** BMI, Smoking, Alcohol Consumption, DietQuality, SleepQuality, Systolic/Diastolic BP.
- **Clinical History:** CardiovascularDisease, Diabetes, Depression, HeadInjury, Hypertension.
- **Biochemistry:** Total Cholesterol, LDL, HDL, Triglycerides.
- **Assessments:** MMSE, FunctionalAssessment, ADL, MemoryComplaints, BehavioralProblems, Confusion.

To replicate the constraint of "small n " studies, we isolated a subset of exactly $N = 40$ subjects. This subset was further stratified into a **training set** ($n_{train} = 32$) and a **held-out test set** ($n_{test} = 8$). The remaining 2,109 subjects were sequestered.

4.2 Generative Modeling (GReaT)

We employed the GReaT framework [2], which fine-tunes a pre-trained DistilGPT-2 Transformer on the textual representation of the tabular data (e.g., "Age is 80, BMI is 22, Diagnosis is 1..."). We trained for 150 epochs with a batch size of 8.

To address class imbalance, we utilized a **Conditional Sampling** strategy. We prompted the fine-tuned LLM to generate 300 candidate profiles, then applied a filter to select exactly 50 synthetic subjects with **Diagnosis=1** and 50 with **Diagnosis=0**. This resulted in a perfectly balanced augmented dataset D_{syn} of 100 synthetic patients.

4.3 Evaluation

We trained two Random Forest Classifiers ($n_{estimators} = 100$) to evaluate the impact of augmentation: a **Baseline Model** trained only on the **32 real training samples**, and an **Augmented Model** trained on the union of real and synthetic data ($N = 132$). Both models were evaluated on the same isolated 8 test samples.

5 Results

The experiment yielded a significant improvement in predictive performance, confirming the hypothesis that synthetic data can regularize models in data-scarce regimes.

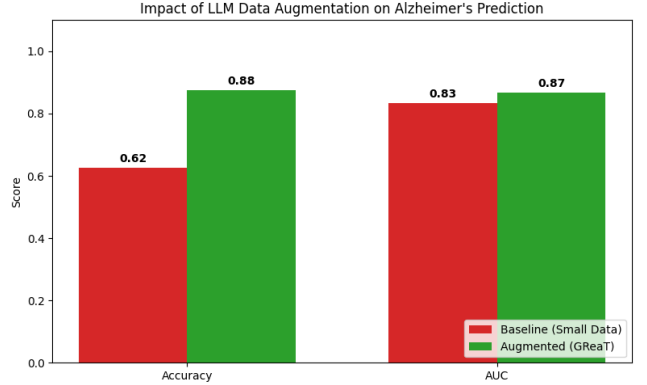


Figure 1: Impact of LLM Augmentation on Model Performance. The Augmented model (green) outperforms the Baseline (red) significantly in both Accuracy (+26%) and AUC (+4%).

Metric	Baseline (N=32)	Augmented (N=132)
Accuracy	0.6250	0.8750
AUC	0.8333	0.8667

Table 1: Performance comparison on the held-out test set ($n = 8$).

As shown in Figure 1 and Table 1, the Baseline model struggled with an accuracy of 62.5%, reflecting the difficulty of generalizing from only 32 samples. The Augmented model achieved an accuracy of 87.5%, a substantial improvement of 25 percentage points. The AUC also improved from 0.83 to 0.87.

6 Discussion and Conclusion

This project successfully replicated the findings of Lee et al. [2], demonstrating that LLM-based tabular augmentation is a powerful tool for small-sample biomedical research. By generating realistic, class-balanced synthetic patients, we transformed a weak classifier into a much stronger predictor.

Moreover, our experiments revealed a boundary condition: augmentation was most effective when the baseline model was "starved" for data ($N = 40$). In preliminary experiments with larger sample sizes ($N = 80$), the baseline was already highly proficient at classification due to strong clinical features (e.g., ADL scores), and augmentation yielded diminishing returns. This suggests that generative AI is best deployed in the early phases of clinical research or for subtle tasks (like SuperAger prediction) where data scarcity is the primary bottleneck.

Several limitations of the present work include a small test set (although we used a standard 80/20 train/test split), "seed bias" and lack of generalization (due to training on a small "seed" of only 32 patients), the possibility of hallucinations with the synthetic data generation, and limitation to tabular data only (and not multi-modal data including genetics and imaging).

Further work in examining data augmentation as outlined herein for other clinical scenarios where patient specimens may be scarce is warranted. Also, examining the efficacy of these methods in studies examining data fusion would also add significantly to this work.

References

- [1] Lee, H. B., et al. (2025). *Machine learning based prediction of cognitive metrics using major biomarkers in SuperAgers*. Scientific Reports, 15:18735.
- [2] Borisov, V., et al. (2022). *Language models are realistic tabular data generators*. ICLR.
- [3] Siegel, N. T., et al. (2025). *Do Transformers and CNNs Learn Different Concepts of Brain Age?*. Human Brain Mapping, 46:e70243.
- [4] Zhou, S., et al. (2025). *A novel sequence-based transformer model architecture for integrating multi-omics data*. npj Digital Medicine, 8:536.
- [5] Chawla, N. V., et al. (2002). *SMOTE: Synthetic Minority Over-sampling Technique*. Journal of Artificial Intelligence Research, 16:321-357.
- [6] Chen, T., & Guestrin, C. (2016). *XGBoost: A Scalable Tree Boosting System*. KDD '16.