

Appendix: Curated LLM: Synergy of LLMs and Data Curation for tabular augmentation in ultra low-data regimes

Table of Contents

A	Extended related work	16
B	Experimental details	18
B.1	Datasets	18
B.2	Data generation.	18
B.3	Data curation	19
B.4	Downstream task	19
B.5	Prompt example	20
C	Additional results	21
C.1	Detailed results for Section 3.1	21
C.2	Full results for performance evaluation	23
C.3	Decoupling prior knowledge and data model	23
C.4	Ablation for contextual information on Compas	24
C.5	Comparison to random noise baseline	25

A EXTENDED RELATED WORK

This paper primarily engages with the work on data augmentation when we have limited data, where our primary goal is synthetic data generation to augment the dataset. Generating synthetic datasets not only helps improve downstream performance, but it is also a flexible solution as it doesn't tie the data consumer to any particular downstream model. That said, beyond the major difference of synthetic data generation, for completeness we contrast our setting of learning from limited data with other seemingly similar settings and highlight their differences.

Contrasting learning w/ limited data vs other settings. The challenge of learning from limited data, while seemingly related to several other learning paradigms, presents distinct differences and unique intricacies that warrant dedicated study.

Transfer learning (Pan & Yang, 2009), *domain adaptation* (Farahani et al., 2021), and *few-shot learning* (Wang et al., 2020) employ additional data resources or rely on specific task-related assumptions to improve learning performance. These methods exploit large labeled data from a source domain, unlabeled data in a target domain, or leverage knowledge from related tasks respectively. For example, Levin et al. (2022) and Jin & Ucar (2023) use models trained on labeled data from a source domain, while Ruiz et al. (2023) and Margeloiu et al. (2022) leverage knowledge-graphs. This is in contrast to our setting, considered of learning with limited data, which must function with whatever scarce labeled data it has, without making any assumptions about the availability of additional data or tasks.

Active learning (Settles, 2009) and *semi-supervised learning* (van Engelen & Hoos, 2019; Chapelle et al., 2006) also operate under the premise of having access to plentiful unlabeled data and the capacity to interactively query labels. However, in our setting, considered learning with limited data does not inherently assume such capabilities, focusing instead on limited labeled data only.

Furthermore, active learning primarily focuses on the iterative process of selecting data samples that, when labeled, are expected to most significantly improve the model's performance. This selection is typically based on criteria such as uncertainty sampling which focuses on **epistemic uncertainty** (Mussmann & Liang, 2018; Houlisby et al., 2011; Kirsch et al., 2019; Nguyen et al., 2022). The primary objective is to minimize labeling effort while maximizing the model's learning efficiency. Additionally, active learning would aims to label instances based on epistemic uncertainty where the model struggles to make accurate predictions, yet the samples themselves are correct. In contrast, CLLM leverage training dynamics based on **aleatoric uncertainty** and confidence and is designed to discard samples that might jeopardize the downstream accuracy. These samples can be considered to have inherent issues or are erroneous, such as being "mis-labeled". To summarize, in active learning, epistemic uncertainty is used to identify data points that, if labeled, would yield the most significant insights for model training. In our approach, they serve to identify and exclude/filter data points that could potentially deteriorate the model's performance.

Self-supervised learning (Liu et al., 2021) leverages large amounts of unlabeled data to learn useful representations for downstream tasks. However, in our setting, considered learning with limited data does not inherently assume such access to vast amounts of unlabeled data.

Data-centric AI. Ensuring high data quality is a critical but often overlooked problem in ML, where the focus is optimizing models (Sambasivan et al., 2021). Even when it is considered, the process of assessing datasets is adhoc or artisanal (Seedat et al., 2022b). However, the recent push of data-centric AI (Liang et al., 2022; Polyzotis & Zaharia, 2021; Zha et al., 2023) aims to develop systematic tools to curate existing datasets. Our work contributes to this nascent body of work (Seedat et al., 2023) – presenting CLLM, which, to the best of our knowledge, is the first systematic data-centric framework looking at how we can tailor synthetic datasets) rather than real datasets) to downstream task use with data curation.

Why Data Augmentation? Data augmentation is a flexible approach to address the ultra low-data regime. An alternative might be to resort to a pretrained black-box model for classification, which could be for example via in-context learning for classification (Dong et al., 2022). However, such a solution is inadequate for several reasons, many of which would prevent real-world utility (e.g. in LMICs):

► *Not economical over the long term:* While using an LLM like GPT for classification may seem attractive due to its few-shot capabilities, it is likely not economically viable in real-world settings, especially in LMICs. The reason is classifying each sample will incur a cost to call the LLM, hence scales linearly with the number of test samples. Over time, the cumulative cost of these calls will surpass the once-off fixed cost associated with generating data. With data augmentation, once the dataset is augmented, there are no additional deployment time costs associated with the LLM. Indeed, the downstream models e.g. a random forest or XGBoost have negligible inference costs.

► *Control, interpretability and auditability:* Relying on a large, pre-trained LLM as a black-box classifier raises several concerns. (1) we have no control over our downstream classifier and its architecture, (2) lack of interpretability and auditability of the LLM when issuing predictions. In contrast, training a downstream model on augmented data maintains the ability to understand and explain how the model is making decisions (e.g. feature importance). This is especially crucial in contexts where accountability, transparency, and validation of machine learning processes are paramount.

► *Independence and self-sufficiency:* Relying on third-party services for continuous classification means being dependent on their availability, pricing models, and potential changes in the LLM version. By augmenting data and training a downstream classifier on the augmented dataset, we ensure that there is no external dependencies such as increasing costs or reduced performance with LLM version updates.

► *Hardware and financial constraints:* Even if we opt for an open-source LLM (e.g. Falcon (Penedo et al., 2023) or LLaMA-2 (Touvron et al., 2023)), deploying and running it locally demands significant computational resources. Typically, these models require GPUs with high amounts of VRAM for optimal performance (e.g. needing around 40 GB hence requiring an A100 GPU for Falcon-40b and LLaMA-2 65B). Such high-end GPUs are expensive, and are likely to be inaccessible in a LMIC setting. Furthermore, renting hardware by the hour can quickly become prohibitively expensive. Data augmentation, on the other hand, can often be performed on modest hardware, and once the augmented dataset is created, many classifiers can be trained without the need for high-end GPUs, making the entire process more financially accessible.

In conclusion, while large language models offer vast knowledge, for low-data settings in low-income countries, data augmentation provides a more cost-effective, controllable, and interpretable solution for building robust classifiers.

B EXPERIMENTAL DETAILS

We provide details on our datasets used, as well as, other experimental specifics including: generation, curation, downstream model, prompt template.

B.1 DATASETS

We summarize the different datasets we use in this paper in Table 5. The datasets vary in number of samples, number of features and domain.

Table 5: Summary of the datasets used. * Denotes private/proprietary datasets.

Name	n samples	n features	Domain
Adult Income (Asuncion & Newman, 2007)	30k	12	Finance
Compas (Angwin et al., 2016)	5k	13	Criminal justice
*Covid-19 (Baqui et al., 2020)	7k	29	Healthcare/Medicine
*CUTRACT Prostate (PCUK, 2019)	2k	12	Healthcare/Medicine
Drug (Fehrman et al., 2017)	2k	27	Healthcare/Medicine
*MAGGIC (Pocock et al., 2013)	41k	29	Healthcare/Medicine
*SEER Prostate (Duggan et al., 2016)	20k	12	Healthcare/Medicine

We detail the dataset splits used in Sec. 3.1. For each dataset and number of samples $n \in \{20, 40, 100, 200\}$, we sample a training set D_{train} such that $|D_{\text{train}}| = n$, and each target class has the same number of samples. We then split the remaining samples into two non-overlapping datasets, D_{oracle} and D_{test} , which have the same cardinality. This procedure is repeated $n_{\text{seed}} = 10$ times, thus leading to different training and test sets. Note that the different generative models use the same D_{train} and D_{test} for a given seed.

Motivation for the choice of datasets.

1. **Open-source:** Adult, Drug and Compas are widely used open-source datasets used in the tabular data literature. Adult and Drug are both UCI datasets that have been used in many papers, while Compas is part of OpenML Vanschoren et al. (2013). Our reason for selecting them is that, despite them being open-source, they are highly reflective of domains in which we might be unable to collect many samples — hence in reality would often be in an ultra-low data regime.
2. **Private datasets:** We wanted to disentangle the possible role of memorization in the strong performance of the LLM. To ensure the datasets are not in the LLMs training corpus, we selected 4 private medical datasets that need an authorization process to access. Hence, these datasets would not be part of the LLMs training corpus given their proprietary nature and hence would be unseen to the LLM. While the private and unseen aspect was the main motivation, we also wish to highlight that these are real-world medical datasets. Consequently, this allows us to test a highly realistic problem setting.

B.2 DATA GENERATION.

GPT-4 and GPT-3.5 We access GPT-4 (OpenAI, 2023) and GPT-3.5-Turbo (Brown et al., 2020) through the API. We use a temperature of 0.9.

GReaT. GReaT Borisov et al. (2023) is a generative model which fine-tunes an LLM based on a training set. We use the implementation provided by authors.

Generative model based approaches. For the other baselines used in 3.1, we use the library SynthCity (Qian et al., 2023), using the defaults. We detail each next.

- TVAE: this is a conditional Variational Auto Encoder (VAE) for tabular data and is based on Xu et al. (2019)

- CTGAN: A conditional generative adversarial network which can handle tabular data and is based on [Xu et al. \(2019\)](#)
- NFLOW: Normalizing Flows are generative models which produce tractable distributions where both sampling and density evaluation can be efficient and exact.
- TabDDPM: A diffusion model that can be universally applied to any tabular dataset, handles any type of feature and is based on [Kotelnikov et al. \(2022\)](#)

Traditional Data Augmentation. We use SMOTE ([Chawla et al., 2002](#)) which augments data by considering nearest neighbors and performing linear interpolations. We use the implementation provided by [Lemaître et al. \(2017\)](#), and set the number of neighbors k to 5.

B.3 DATA CURATION

Learning dynamics computation We train an XGBoost with 100 estimators on D_{train} . We then compute predictive confidence and aleatoric uncertainty for the samples in D_{syn} . **The motivation for the choice of an XGBoost backbone is that we cannot expect good performance by choosing “any” curation model, but rather we require a curation model with enough capacity and generalization properties — where boosting methods like XGBoost used in our work have shown to achieve best performance on tabular data. This leads to our guideline for the curation step: the model used for curation should be at least as flexible as the model that the practitioner intends to use for the downstream task.**

Learning dynamics thresholds Recall that CLLM has two thresholds τ_{conf} and τ_{al} on the predictive confidence and aleatoric uncertainty respectively, as defined in 2.2. We set $\tau_{\text{conf}} = 0.2$, in order to select high confidence samples. We adopt an adaptive threshold for τ_{al} based on the dataset, such that $\tau_{\text{al}} = 0.75 \cdot (\max(v_{\text{al}}(D_{\text{syn}})) - \min(v_{\text{al}}(D_{\text{syn}})))$. Note that by definition $v_{\text{al}}(D_{\text{syn}})$ is bounded between 0 and 0.25.

Example of learning dynamics We include examples of learning dynamics computed for 20 samples in Fig. 7.

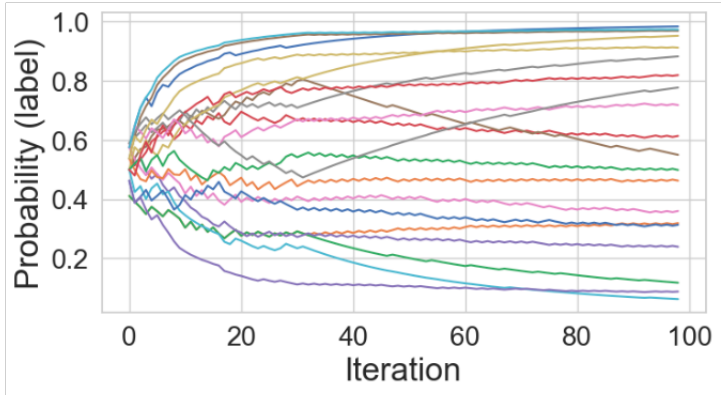


Figure 7: Learning dynamics computed for 20 samples

B.4 DOWNSTREAM TASK

We compute downstream performance in Sec. 3.1 using four different downstream models: XGBoost, Random Forest, Decision tree, and Logistic Regression.

B.5 PROMPT EXAMPLE

We include the template of the prompts used throughout the paper. We show how we include (1) in-context examples (demonstrations), (2) contextual information including dataset background and feature information and (3) the instruction.

```

1   System role: 'You are a tabular synthetic data generation model.'
2
3   You are a synthetic data generator.
4   Your goal is to produce data which mirrors \
5   the given examples in causal structure and feature and label
6   distributions \
7   but also produce as diverse samples as possible.
8
9   I will give you real examples first.
10
11  Context: Leverage your medical knowledge about covid and Brazil to
12  generate 1000 realistic but diverse samples.
13
14  example data: {data}
15
16  The output should be a markdown code snippet formatted in the
17  following schema:
18
19  "Sex_male": string // feature column
20  "Age": string // feature column
21  "Age_40": string // feature column
22  "Age_40_50": string // feature column
23  "Age_50_60": string // feature column
24  "Age_60_70": string // feature column
25  "Age_70": string // feature column
26  "Fever": string // feature column
27  "Cough": string // feature column
28  "Sore_throat": string // feature column
29  "Shortness_of_breath": string // feature column
30  "Respiratory_discomfort": string // feature column
31  "SPO2": string // feature column
32  "Dihareea": string // feature column
33  "Vomitting": string // feature column
34  "Cardiovascular": string // feature column
35  "Asthma": string // feature column
36  "Diabetis": string // feature column
37  "Pulmonary": string // feature column
38  "Immunosuppresion": string // feature column
39  "Obesity": string // feature column
40  "Liver": string // feature column
41  "Neurologic": string // feature column
42  "Renal": string // feature column
43  "Branca": string // feature column
44  "Preta": string // feature column
45  "Amarela": string // feature column
46  "Parda": string // feature column
47  "Indigena": string // feature column
48  "is_dead": string // label if patient dead or not, is_dead
49
50  DO NOT COPY THE EXAMPLES but generate realistic but new and diverse
51  samples which have the correct label conditioned on the features.

```

Listing 1: Template of the prompt

C.2 FULL RESULTS FOR PERFORMANCE EVALUATION

We report full results with standard deviation for the results from the main paper. The performance is AUC averaged over XGBoost, Random forest, Logistic regression, Decision tree.

Table 10: AUC averaged over 4 downstream models on D_{test} where curation improves performance for all methods across all sample sizes n , as indicated by \uparrow . CLLM w/ GPT-4 (Curated) dataset provides the strongest performance for both private/proprietary datasets and public datasets

Table with columns for Dataset, Real data (D_oracle, D_train), CLLM (OURS) (GPT-4, GPT-3.5), and Baselines (CTGAN, TabDDPM, GReT, NFLOW, SMOTE, TVAE). Rows list various datasets like covid, curact, maggic, seer, compas, adult, drug, etc. with performance metrics.

C.3 DECOUPLING PRIOR KNOWLEDGE AND DATA MODEL

Two components can be attributed to the good performances of CLLM: the background knowledge of the LLM, and its capacity to build a strong data model. In this subsection, we provide insights to understand the effect of the LLM’s background knowledge (e.g. prior). We considered the Covid dataset (private medical dataset, to avoid memorization issues) and generated data with GPT-4 (same as Section 2.1). We ablate the prompt used in our work (detailed in Appendix B.5), and solely provide one in-context example in the prompt, in order to give the LLM the minimal amount of information about the desired structure of the dataset. This lack of examples forces the LLM to rely on its own prior (background knowledge), and removes the effect of in-context examples which could be used to build a data model. We report the results in Table 11, along with the results for CLLM. From these results, we conclude the following:

- 1. The LLM prior permits to obtain good downstream performance, but it is outperformed by D_oracle by a margin of 4.4%. Hence, we cannot solely rely on the prior.
2. Downstream performance increases as the number of in-context samples increases. This shows it is indeed important to include the in-context examples if we wish to obtain downstream performance close to D_oracle, as the LLM can build a good data model.

This implies that while the LLM does use background knowledge of similar datasets, it still requires in-context samples to refine its prior by creating a good data model.

Table 11: Downstream accuracy when varying the number of in-context samples in the prompt to generate the augmented datasets.

Table with 2 columns: In-context samples and Downstream accuracy. Rows show n=1 (Prior), n=20, n=40, n=100, and D_oracle with their respective accuracy values.

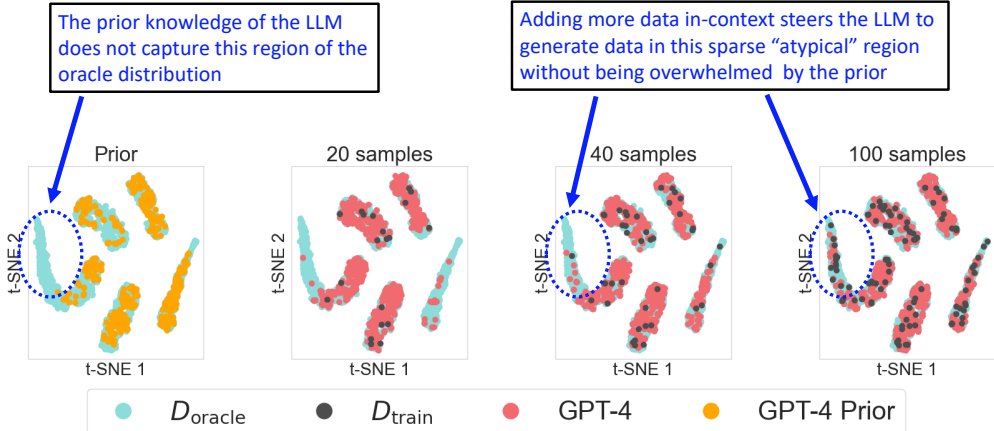


Figure 8: The data generated by the LLM captures the distinct features in "atypical" regions of the Oracle manifold, as in-context samples are added to the prompt. This shows that it is flexible enough to adapt its prior knowledge to the nuances of the data. The group encircled in blue represents patients who are > 88 years old, representing around 3.5% of the Oracle. This illustrates the added in-context samples can successfully guide the LLM to generate these rare samples.

Next, we quantify and visualize the strength of the prior, by studying how much the LLMs output distribution adapts to the in-context samples provided. We evaluate data generated by the prior of the LLM ($n = 1$), and for $n = 20, 40, 100$ on the Covid dataset.

In particular, we observe in Figure 8 that there is a region in the oracle data which is not captured by the LLM’s prior output (the left part of the leftmost blob, circled in blue in Figure 8). However, as the number of in-context real examples increases in the prompt of the LLM, we observe that this steers the LLM to generate data which covers this region. This region is associated to the subgroup of people older than 87 years old, and having many severe comorbidities (e.g. Diabetes, Cardiovascular diseases) and many respiratory symptoms. This subgroup, in the Oracle dataset, represents less than 3.5% of the data, and is completely ignored by the GPT-4 prior. In particular, the prior defaults to more typical patients in the range 70-80 years old. On the contrary, as n increases, the LLM is guided by the in-context samples and generates samples from this subgroup, which are "rarer" or different from the general population.

This demonstrates that the LLM captures the distinct features of this particular region, and hence is not overwhelmed by the prior, but instead the data in the form of in-context samples adapts it, hence aligning the augmented dataset with the ground-truth distribution.

C.4 ABLATION FOR CONTEXTUAL INFORMATION ON COMPAS

We conduct a similar experiment as in Table 2, and use the dataset Compas. We report the results in Table 12.

Table 12: Including contextual information in the prompt improves precision (P), recall (R), and utility (U) in low-sample settings (results shown for Compas).

n_{samples} in D_{train}	GPT-4 w/ context			GPT-4 no context			TVAE		
	P	R	U	P	R	U	P	R	U
20	0.69 _(0.02)	0.88 _(0.02)	0.69 _(0.02)	0.27 _(0.03)	0.89 _(0.03)	0.60 _(0.03)	0.43 _(0.02)	0.43 _(0.05)	0.55 _(0.04)
40	0.70 _(0.0)	0.92 _(0.01)	0.65 _(0.03)	0.31 _(0.06)	0.84 _(0.03)	0.57 _(0.01)	0.54 _(0.02)	0.80 _(0.02)	0.50 _(0.04)
100	0.69 _(0.02)	0.89 _(0.02)	0.69 _(0.01)	0.34 _(0.1)	0.85 _(0.05)	0.62 _(0.01)	0.60 _(0.03)	0.86 _(0.02)	0.59 _(0.03)
200	0.70 _(0.01)	0.89 _(0.02)	0.69 _(0.01)	0.31 _(0.05)	0.87 _(0.03)	0.58 _(0.05)	0.65 _(0.02)	0.88 _(0.01)	0.63 _(0.01)

These results highlight the importance of incorporating contextual information in the prompt, as it enables to exploit the prior knowledge of the LLM.

C.5 COMPARISON TO RANDOM NOISE BASELINE

We now compare to a random noise baseline. Specifically, where we augment the dataset with random additive Gaussian noise. In order to capture the correlations between the different features, we fit a Kernel Density Estimator with a Gaussian kernel and bandwidth given by Scott’s rule. We then sample 1000 points to create an augmented dataset D_{syn} . We report the performance gap

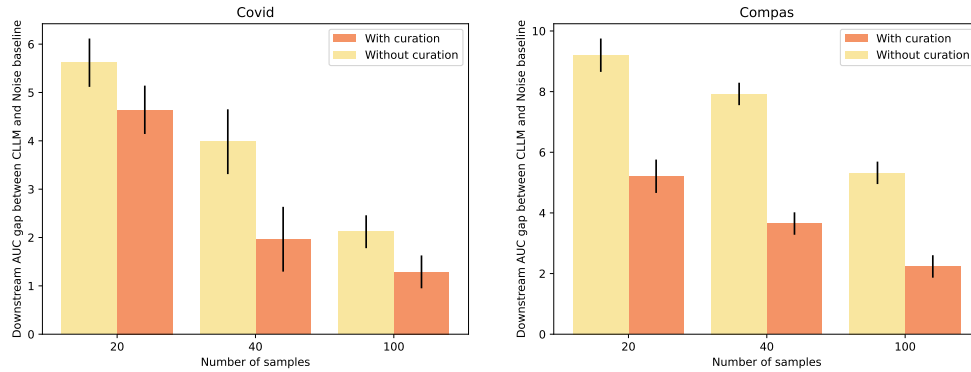


Figure 9: The random noise baseline does not match the performance of CLLM

between CLLM and this baseline (with and without curation) for the Covid and Compas datasets in Figure 9.

We observe that the random noise baseline does not match the performance of CLLM (i.e. has a performance gap), although the baseline naturally improves as the dataset D_{train} grows in size.