

ADVANCING TABLE UNDERSTANDING OF LARGE LANGUAGE MODELS VIA FEATURE RE-ORDERING

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) exhibit exceptional proficiency in comprehending human language. Despite their significant success across a wide array of tasks, including text generation, translation, question answering, and even code generation, understanding tabular data remains a challenging task. Especially, tabular data lacks an intrinsic order of the different features (table fields), whereas LLMs take only sequential inputs. Consequently, an artificial order is imposed, the impact of which on the performance of LLMs has not yet been thoroughly investigated. Surprisingly, as discovered in this work, this artificially induced order bias dramatically influences the performance of LLMs on tasks related to tabular data. Mitigating the order bias presents a significant challenge. To address this, we propose a simple and cost-effective method, Re-Ordering Tabular feATures fOR LLM (ROTATOR-LLM), to conduct test-time compute without fine-tuning the base LLM. Aiming at optimizing the feature order of tabular data and boosting LLMs’ capability to better understand the data semantics, ROTATOR-LLM re-frames the ordering problem as a feature trajectory generation task. A dynamic programming based meta-controller is trained to auto-regressively generate an individualized feature trajectory for each data instance via accumulative value estimation of the serialized feature input through the LLM’s final performance metrics. Model performance is maximized by iteratively selecting features across different steps. Experimental results on multiple datasets and LLMs show close to or over 20% performance boosts via features reordered by ROTATOR-LLM against the un-ordered counterpart. Also, it outperforms State-Of-The-Art tabular LLM methods with significant margin. Moreover, meta-controller demonstrates strong transferability: the tested LLMs gain performance enhancements when utilizing a meta-controller trained on one of them.

1 INTRODUCTION

Tabular data is prevalent in real-world scientific, medical, biological, sociological, financial, and retail databases, necessitating significant time and effort for humans to process and analyze Dong & Wang (2024); Fang et al. (2024). Fortunately, advancements in large language models (LLMs) have enabled rigorous exploration of their application in various tasks related to tabular data modeling Yuan et al. (2024); Hu et al. (2024). Recent breakthroughs have involved LLMs to handle a wide range of tabular data tasks, such as TabLLM Heggelmann et al. (2023), TableGPT Zha et al. (2023b), and TableLlama Zhang et al. (2023).

Although tabular data can be easily converted into text format, LLMs struggle to effectively analyze the converted data. Since LLMs are primarily pre-trained on natural language, they face challenges in extracting meaningful insights from structured tabular data. To overcome this challenge, existing work primarily focuses on fine-tuning LLMs on tabular dataset to inject the data prior knowledge to the models. For example, TableLlama employs LongLoRA to fine-tune the Llama-2-7B LLM on the extensive TableInstruct datasets. Similarly, TableGPT introduces a table encoder and chain-of-command mechanism, utilizing a Phoenix-7B LLM for inference. Despite these advancements, much of the current research on tabular data analysis overlooks the critical role of feature order in the prompt: due to the sequential nature of transformer decoder based models, an artificial order is inevitably created when feeding the features into the LLM one by one regardless of the detailed prompting schemes. Our extensive studies reveal that this induced ordering of features significantly

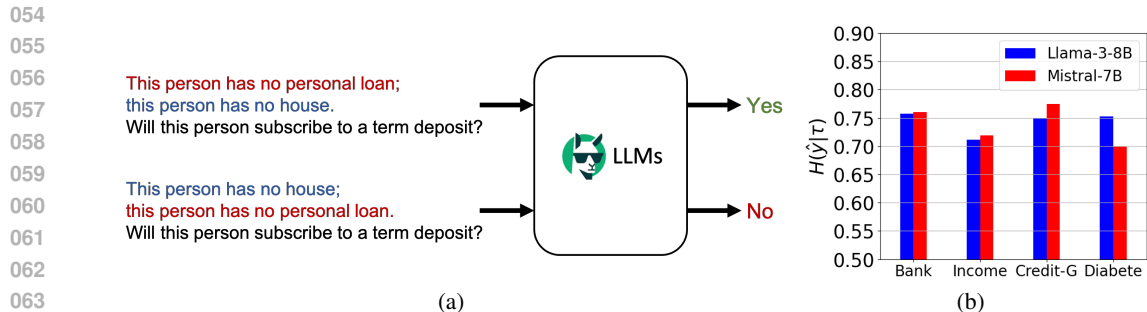


Figure 1: (a) An example of LLM order bias. (b) Order bias generally exist in different LLMs.

impacts LLM’s behavior Chen et al. (2024); Xu et al. (2024). For instance, the LLM prediction on the same data instance can vary just by changing the order of input features, as in Figure 1 (a). Further details are discussed in Section 3.

This problem is mainly rooted in the order bias in the pre-training data, where the collected data follows certain sequences preferred by humans. Such order preference is captured by the LLMs during the pre-training stage, which enables LLMs to better learn the data semantics whose feature importance ranking aligns with the order bias Sagawa et al. (2019); Koh et al. (2021). To tackle this, an intuitive solution is to remove the order bias by fine-tuning the LLMs on unbiased data. However, fine-tuning LLMs is not only time- and resource-consuming due to the billions of updated parameters, but also labor-intensive, requiring collecting high-quality data Yang et al. (2024); Zha et al. (2023a). A more practical approach is to preprocess the data to align with the LLMs’ inherent order bias, enabling them to better grasp the data’s semantics. This alignment offers greater potential for real-world applications due to its feasibility, scalability, and extensibility across diverse datasets.

In this work, we introduce Re-Ordering Tabular feATures FOR LLM (ROTATOR-LLM), a simple and cost-effective method to help LLMs better comprehend data semantics via test-time compute in the input level Snell et al. (2024). Specifically, ROTATOR-LLM converts the feature ordering problem into a task of generating feature trajectories, where each trajectory represents a sequence of features in a specific order. To avoid the high resource consumption of fine-tuning the LLM and the corresponding expensive human labeling, ROTATOR-LLM trains a light-weight neural network as a meta-controller to auto-regressively generates the optimized feature trajectory for each data instance, guided by a value function designed to supervise its training process. It is challenging to define the value function for a specific feature order such that this value aligns with the corresponding LLMs’ performance. We are motivated by dynamic programming to overcome this challenge. Specifically, the value of a feature trajectory is defined as its potential maximal value in the next state within the whole generation path. At the last state, the value of an integral trajectory is determined by the LLMs’ performance. This approach allows us to estimate the value of any feature trajectory, which, in turn, supervises the training of the meta-controller. To evaluate ROTATOR-LLM, we conduct experiments with three LLMs across four tabular datasets. The results demonstrate that LLMs perform significantly better on data reordered by ROTATOR-LLM compared to random or default orders, underscoring the effectiveness of the reordering process. Moreover, ROTATOR-LLM outperforms existing foundational tabular LLMs, further highlighting its potential in real-world applications. In summary, our contributions in this work are as follows:

- **Order Bias of LLMs.** We demonstrate that the order of instance features in a prompt significantly influences LLM predictions, identifying the presence of order bias.
- **Alignment to Order Bias.** We propose ROTATOR-LLM, a cost-effective solution that requires no tuning of LLM parameters. ROTATOR-LLM aligns a data instance with the inherent order bias of LLMs by re-ordering its features.
- **Experimental Evaluation.** Experimental results on four datasets with three popular LLMs demonstrate the superior performance lift brought by ROTATOR-LLM, which improves LLMs’ classification accuracy by 20% in average.

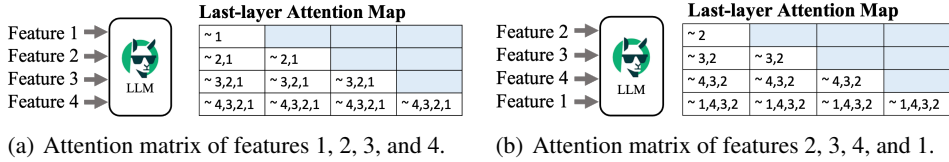


Figure 2: Comparison of the last-layer attention map under different orders of input features. Since each feature is represented by a sentence, i.e. multiple tokens, each cell corresponds to a matrix of attention values between tokens. The notation ‘ $\sim i, j, k$ ’ indicates the attention matrix is computed based on a mixture of information from the token embeddings associated with features i, j and k .

2 PRELIMINARIES

We introduce the notations and data format transition in this section.

2.1 NOTATIONS

We consider aligning the dataset $\mathcal{D} = (\mathbf{x}, y) \mid \mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}$ to the order bias of LLMs $f(\bullet)$. Each instance $\mathbf{x} \in \mathcal{X}$ has M features, $\mathbf{x} = [x_1, x_2, \dots, x_j, \dots, x_M]$, where $j \in \mathcal{J} = \{1, 2, \dots, M\}$ is the feature index in the default order of a particular tabular dataset. Let $\tau = [\tau_1, \tau_2, \dots, \tau_M]$ denote a specific ordering of the features of instance \mathbf{x} , representing a feature trajectory with M positions. For $1 \leq t \leq M$, each $\tau_t \in \{x_1, x_2, \dots, x_M\}$ indicates a feature ranked at position t ; and $\tau_{[0:t]}$ denotes a slice of the trajectory comprising the first t positions $[\tau_1, \dots, \tau_t]$, each containing a feature best suited for the corresponding position. The case $t = 0$ represents the initial state $\tau_{[0:0]} = []$ where no features have been ranked, while $t = M$ denotes the final state $\tau_{[0:M]}$ that all M positions are filled by properly ranked features. For example, if there are in total 3 features, the full trajectory $\tau = [x_2, x_3, x_1]$ represents the features are ordered as 2, 3, and 1 at positions 0, 1, and 2, respectively. In Section 3, we demonstrate the order bias of LLMs by showing that the prediction results $\hat{y} = f(\tau)$ are significantly affected by the order of input features τ . To address this issue, we introduce ROTATOR-LLM in Section 4, which aligns the dataset \mathcal{D} with the order bias of LLMs. ROTATOR-LLM aims to generate the optimal trajectory τ^* for each instance \mathbf{x} , thereby maximizing the accuracy of the LLMs’ predictions.

2.2 TEXT-BASED SERIALIZATION

Text-based Serialization refers to converting tabular data into text data to fit the input modality of LLMs. Existing work explores several methods of text-based serialization. For example, Markdown table Liu et al. (2023); Jaitly et al. (2023), JSON-file format Singha et al. (2023); Sui et al. (2024), and sentence serialization Yu et al. (2023); Jaitly et al. (2023). To maximally leverage the sequence-to-sequence capacity of LLMs, we consider the sentence serialization to convert the data features into text data. The advantage of sentence serialization is its alignment with the natural language data where LLMs are pre-trained. In this work, we use a template given in Appendix A to convert tabular data into text data. For instance, we adopt the sentence “the age of this person is 30; this person has no house” to represent the tabular data $\{\text{Age} : 30, \text{House} : \text{No}\}$. Our method can be easily extended to fit Markdown table and JSON-file formats of serialized data, but their performance is out of the scope of this work.

3 ORDER BIAS OF LLMs ON TABULAR DATA

In this section, we empirically analyze the order bias of LLMs and present the experimental evidence of LLM’s behavior change under the influence of order bias.

3.1 WHY LLMs HAVE ORDER BIAS?

Order bias refers to the impact that the sequence of tabular data features has on the predictions made by LLMs. While from the perspective of how human beings understand the tabular data, the order of features/fields is not meaningful and should not affect the model output, each particular serialization of these features/fields indeed results in a different input sequence for an auto-regressive model and accordingly a difference in the outcome. For LLMs, this difference affects their attention maps. We show an example in Figure 2 to demonstrate the influence of different feature orders on the last-layer attention maps. As each feature is represented by a sentence, i.e. multiple tokens, each cell in Figure 2 corresponds to a matrix of attention values between tokens. The notation ‘ $\sim i, j, k$ ’ indicates the attention matrix is computed based on a mixture of information from the token embeddings associated with features i, j and k . In this example, the sequence of features 1, 2, 3, and 4 in Figure 2 (a) mixes a different set of tokens compared to the feature sequence of 2, 3, 4, and 1 for the computation of last-layer attention map. The variations in last-layer attention maps lead to obvious differences in the prediction results.

3.2 DEMONSTRATIONS OF ORDER BIAS

We demonstrate the presence of order bias in LLMs using real-world tabular datasets. Specifically, we examine the variance in LLMs’ predictions caused by different permutations of data features. The probability of LLMs’ predictions is estimated by $\mathbb{P}(\hat{y} = 1) = \frac{\# \text{ of } 1}{\# \text{ of Permutations}} = \frac{\# \text{ of } 1}{M!}$, and $\mathbb{P}(\hat{y} = 0) = 1 - \mathbb{P}(\hat{y} = 1)$. The variance in predictions is quantified by the entropy $\mathcal{H}(\hat{y}) = -\mathbb{P}(\hat{y} = 0) \log_2 \mathbb{P}(\hat{y} = 0) - \mathbb{P}(\hat{y} = 1) \log_2 \mathbb{P}(\hat{y} = 1)$. For instance, for data instance having two features: age and house, if an LLM outputs $\hat{y} = 1$ for $\{\text{Age}: 30, \text{House}: \text{No}\}$ and $\hat{y} = 0$ for $\{\text{House}: \text{No}, \text{Age}: 30\}$, then $\mathbb{P}(\hat{y} = 1) = \mathbb{P}(\hat{y} = 0) = 0.5$, resulting in an entropy of 1. If the LLM’s predictions show no variance, then either $\mathbb{P}(\hat{y} = 0) = 1$ or $\mathbb{P}(\hat{y} = 1) = 1$, yielding a minimal entropy of 0. Conversely, if the predictions are randomly distributed, $\mathbb{P}(\hat{y} = 0) = 0.5$ and $\mathbb{P}(\hat{y} = 1) = 0.5$, leading to a maximum entropy of 1. Higher entropy indicates greater variance in prediction results, signifying a stronger presence of order bias in the LLMs.

The experiments are conducted on the Bank, Income, German Credit, and Diabetes datasets Asuncion et al. (2007), using the Llama-2-8B-instruct Touvron et al. (2023) and Mistral-7B-Instruct Jiang et al. (2024) LLMs as predictors. The entropy of predictions resulting from feature reordering is shown in Figure 1 (b). Notably, all LLMs applied to the tabular datasets exhibit an entropy exceeding 0.7, approaching the maximum value of 1. This clearly indicates the presence of order bias.

4 RE-ORDERING TABULAR FEATURES FOR LLM (ROTATOR-LLM)

In this section, we introduce Re-Ordering Tabular feATures FOR LLM (ROTATOR-LLM) in details. Specifically, ROTATOR-LLM adopts a meta-controller to generate the reordered feature trajectory; then converts the features to text data following the template in Appendix A; finally inputs the data features in text format to LLMs for inference. The overall objective is to maximize the accuracy of the LLM predictions for tabular data classification tasks. We discuss the details as follows.

4.1 FEATURE TRAJECTORY GENERATION

ROTATOR-LLM maintains a meta-controller $g(\bullet | \theta) : \mathcal{T} \rightarrow \mathbb{R}$ to estimate the ranking value of each feature at each location. Specifically, for $0 \leq t \leq M$, with a slice of trajectory $\tau_{[0:t]}$ as input, the value of $g([\tau_{[0:t]}, x_j] | \theta) \in \mathbb{R}$ represents the value of trajectory $[\tau_{[0:t]}, x_j]$, which also indicates the ranking value of feature j at position t , given the feature ordering of first t positions $\tau_{[0:t]}$. We consider a higher value $g(\tau | \theta)$ as indicative of better ranking results for feature orders that align more closely with the preferences of the LLMs. Therefore, ROTATOR-LLM can recursively generate a trajectory of M data features by

$$\tau_t = \arg \max_{j \in \mathcal{J}} g([\tau_{[0:t-1]}, x_j] | \theta). \quad (1)$$

We define a value function $v(\tau)$ to compute the classification loss of LLMs’ prediction over input data crafted with the feature trajectory τ . We believe a feature ordering that is more aligned with

Algorithm 1 Re-Ordering Tabular feATures fOR LLM (ROTATOR-LLM)**Input:** Training dataset \mathcal{D} and LLM $f(\bullet)$.**Output:** Meta-controller $g(\bullet | \theta)$.

- 1: **for** $(x, y) \sim \mathcal{D}$ **do**
- 2: Generate trajectory τ by Equation (1) based on initial value $\tau_{[0:0]} = []$.
- 3: Estimate the loss value of LLMs' prediction $L_f(f(\tau), y)$.
- 4: Estimate the value function $v(\tau_{[0:t]})$ for $1 \leq t \leq M$ based on Equation (6).
- 5: Update the parameters of $g(\bullet | \theta)$ to minimize Equation (5).
- 6: **end for**

LLMs' pre-training can lead to better prediction result. Therefore, $v(\tau)$ is defined as follows:

$$v(\tau) = -L_f(f(\tau), y) \quad (2)$$

where L_f denotes the cross-entropy; $f(\tau)$ is the prediction output of the base LLM; trajectory value function $v(\tau)$ is opposite to the cross-entropy loss such that the optimal trajectory τ^* can minimize the classification error while maximizing the corresponding value function.

Note that Equation (2) only defines the value of a complete trajectory $v(\tau)$, it is important to extend its definition to a slice of trajectory $v(\tau_{[0:t]})$, for the purpose of training the controller $g(\bullet | \theta)$. However, the value function is strictly defined on the full trajectory τ (not on its slices) and the final LLM output after feeding τ into it, so that $v(\tau_{[0:t]})$ cannot be directly obtained via Equation (2). To overcome this challenge, we employ dynamic programming to define $v(\tau_{[0:t]})$, where $0 \leq t < M$. Specifically, for a slice of trajectory $\tau_{[0:t]}$, its value function $v(\tau_{[0:t]})$ is defined as the maximal value of $v(\tilde{\tau})$ such that $\tilde{\tau}_{[0:t]} = \tau_{[0:t]}$, which is given by

$$v(\tau_{[0:t]}) = \max_{\tilde{\tau}_{[t-1:M]}} \gamma^{M-t} v([\tau_{[0:t-1]}, \tilde{\tau}_{[t-1:M]}]), \quad (3)$$

$$= \max_{j \in \mathcal{J}} \gamma v([\tau_{[0:t-1]}, x_j]), \quad (4)$$

where $0 < \gamma < 1$ denotes a discounting factor. The discounting factor regulates how features ranked at different positions cumulatively contribute to the final cross entropy and full trajectory value. This is inspired by the observation in previous studies that tokens closer to the end contribute relatively more to the output of LLMs Jin et al. (2024).

According to Equation (4), we have an iterative property of the value function given by $v(\tau_{[0:t]}) = \gamma v(\tau_{[0:t+1]})$ running backwards from positions $t = M$ to $t = 0$ with the last state value given by $v(\tau) = -L_f(f(\tau), y)$ at $t = M$. Therefore, the parameters of $g(\tau_{[0:t]} | \theta)$ is updated to minimize the mean-square error aligned with the value function $v(\tau_{[0:t]})$ as follows:

$$L_\theta = \frac{1}{M} \sum_{t=0}^M [g(\tau_{[0:t]} | \theta) - v(\tau_{[0:t]})]^2, \quad (5)$$

where $v(\tau_{[0:t]})$ can be estimated based on its iterative property as follows:

$$v(\tau_{[0:t]}) = \begin{cases} \gamma \max_j g([\tau_{[0:t]}, x_j] | \theta) & \text{if } t < M, \\ -L_f(f(\tau), y) & \text{if } t = M. \end{cases} \quad (6)$$

4.2 ALGORITHM OF ROTATOR-LLM

Algorithm 1 shows one epoch of ROTATOR-LLM. Specifically, for each mini-batch of instances, ROTATOR-LLM first generate an order of features following Equation (1) (line 2); then estimate the loss function of LLMs' prediction, where the input data of LLMs follows the generated feature order (line 3); then estimate the value functions based on Equation (6) (line 4); finally updates the parameters of meta-controller to minimize the loss function given in Equation (5) (line 5).

5 EXPERIMENTS

In this section, we conduct experiments to evaluate ROTATOR-LLM, aiming to answer the following research questions: **RQ1:** Does ROTATOR-LLM effectively align the data with the LLMs for better

Table 1: Balance accuracy of ROTATOR-LLM on the Bank, Income, German Credit, and Diabetes datasets.

Datasets	Order	Bank	Income	German Credit	Diabetes	Average
Llama-3-8B	Default	0.522	0.516	0.521	0.312	0.468
	Random	0.510	0.520	0.535	0.385	0.488
	ROTATOR-LLM	0.791	0.752	0.665	0.738	0.737
Mistral-7B	Default	0.599	0.540	0.493	0.699	0.585
	Random	0.574	0.577	0.546	0.676	0.593
	ROTATOR-LLM	0.782	0.801	0.701	0.722	0.752
Phi-3-mini	Default	0.504	0.510	0.405	0.634	0.513
	Random	0.481	0.521	0.440	0.655	0.524
	ROTATOR-LLM	0.712	0.771	0.665	0.743	0.723

Table 2: F1 score of ROTATOR-LLM on the Bank, Income, German Credit, and Diabetes datasets.

Datasets	Order	Bank	Income	German Credit	Diabetes	Average
Llama-3-8B	Default	0.466	0.674	0.600	0.191	0.483
	Random	0.555	0.676	0.605	0.353	0.547
	ROTATOR-LLM	0.811	0.796	0.732	0.774	0.778
Mistral-7B	Default	0.428	0.678	0.145	0.691	0.486
	Random	0.456	0.692	0.365	0.695	0.552
	ROTATOR-LLM	0.774	0.808	0.734	0.765	0.770
Phi-3-mini	Default	0.245	0.664	0.182	0.505	0.399
	Random	0.439	0.660	0.512	0.632	0.561
	ROTATOR-LLM	0.658	0.776	0.622	0.763	0.705

performance? **RQ2:** Can the controller be transferred between different LLMs? **RQ3:** How does the reordering intrinsically impact the LLMs?

5.1 EXPERIMENT SETUP

We specify the datasets, LLMs, baseline methods, evaluation metrics, and implementation details.

Datasets. The evaluation of ROTATOR-LLM is based on the Bank, Income, German Credit, and Diabetes datasets from the areas of social media, finance and healthcare. The datasets source from the UC Irvine machine learning repository Asuncion et al. (2007). On each dataset, the data features are first reordered; then converted into text data following the template in Appendix A; and finally being input to LLMs for classification.

LLMs. We evaluate ROTATOR-LLM using three popular model families: Llama-3-8B Touvron et al. (2023), Mistral-7B Jiang et al. (2024), and Phi-3-mini-4k Li et al. (2023). These LLMs are employed due to their leadership among open-sourced LLMs according to existing leaderboards Chiang et al. (2024). We download their instruct-tuned version from the Huggingface Wolf et al. (2019).

Baseline Methods. We consider four baseline methods compared with ROTATOR-LLM. **Default order.** The features of each data instance follow the default order provided by the datasets. **Random order.** The features of each data instance are randomly ordered. **TableLlama.** A Llama-based foundational tabular LLM fine-tuned on large-scale tabular datasets Zhang et al. (2023). **TableLLM.** A GPT-2-based foundational tabular LLM fine-tuned on large-scale tabular datasets Zha et al. (2023b).

Evaluation Metrics. Due to the imbalance of positive and negative examples in the datasets, the regular accuracy metric is not sufficient to truly reflect the classification performance. Therefore, we evaluate the balance accuracy (\uparrow) and F1 score (\uparrow) of LLMs' classification on the datasets. To estimate the balance accuracy, the instances of the minority class are first duplicated to align with the size of the majority class. Then the accuracy is calculated.

Table 3: Transfer-ability of ROTATOR-LLM, where the meta controller is trained with a source LLM and tested on a different target LLM.

Metric	Configuration	Bank	Income	Germen Credit	Diabetes	Average
Balance accuracy	Default-Llama	0.522	0.516	0.521	0.312	0.468
	Random-Llama	0.510	0.520	0.535	0.385	0.488
	Mistral→Llama	0.544	0.622	0.627	0.670	0.616
	Default-Mistral	0.599	0.540	0.500	0.699	0.585
	Random-Mistral	0.574	0.577	0.546	0.676	0.593
	Llama→Mistral	0.581	0.756	0.581	0.756	0.669
F1 score	Default-Llama	0.466	0.674	0.600	0.191	0.483
	Random-Llama	0.555	0.676	0.605	0.353	0.547
	Mistral→Llama	0.598	0.714	0.675	0.722	0.677
	Default-Mistral	0.428	0.678	0.145	0.691	0.486
	Random-Mistral	0.456	0.692	0.365	0.695	0.552
	Llama→Mistral	0.504	0.743	0.414	0.690	0.588

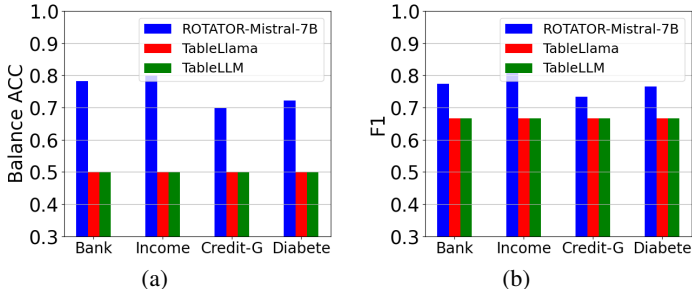


Figure 3: Comparison of ROTATOR-LLM with state-of-the-art foundational Table LLMs.

Implementation Details. The meta-controller takes a three-layer MLP that is trained using Adam optimizer with learning rate 10^{-3} for 200 epochs. An early stop is implemented on the validation datasets. The training and evaluation processes follow the same template of text serialization given in Appendix A. The detailed hyper-parameter setting of ROTATOR-LLM is given in Appendix B.

5.2 ALIGNMENT PERFORMANCE (RQ1)

We evaluate the performance of ROTATOR-LLM by examining the classification of LLMs after the alignment. For fair comparison, ROTATOR-LLM and baseline methods adopt the same prompt given in Appendix A for text serialization. The balanced accuracy and F1 score are shown in Tables 1 and 2, respectively. The comparison with baseline foundational tabular LLMs is illustrated in Figure 3. According to the experimental results, we have the following observations:

- **Effectiveness of Alignment.** LLMs show much better performance based on ROTATOR-LLM than the data with default and random feature orders. This indicates that ROTATOR-LLM effectively align the data feature to LLMs, and thereafter enhances LLMs’ understanding on the tabular data by optimally reordering the features.
- **Competitive Performance.** ROTATOR-LLM outperforms foundational tabular LLMs, e.g., TableLLM and TableLlama. Compare to these costly fine-tuning methods, ROTATOR-LLM not only saves resources effectively but also shows performance superiority.
- **Consistent Performance.** ROTATOR-LLM is consistently competitive over baseline methods across various LLMs and tabular datasets, indicating its stability and generalizability for real-world applications.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

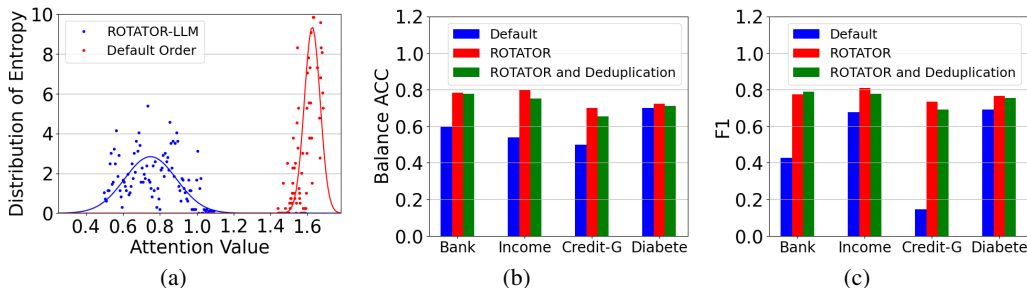


Figure 4: (a) Entropy of last layer attention. The lower the entropy, the more focus of attention. (b) Balanced accuracy and (c) F1 score of shrinking the duplicated features in the prompts.

<p>Prompts: You are a data analyst. Given information of a person, you should predict whether this person will subscribe to a term deposit. <Data Features> Will this person subscribe to a term deposit?\n\n[Your Response Format]: “Yes / No”</p> <p>Label: Yes</p>
<p>Default features: This person’s age is 33.0. The type of this person’s job is technician. This person’s marital status is single. This person’s education is secondary. This person has no credit in default. This person’s average yearly balance in euros is 2979.0. This person has no house. This person has no personal loan. This person’s contact communication type is cellular. This person’s last contact day of the month is 5.0. This person’s last contact month of year is aug. This person’s last contact duration is 326.0 seconds. This person has 2.0 contacts performed during this campaign. 437.0 days have passed since this person was last contacted from a previous campaign. This person has 1.0 contacts performed before this campaign. The outcome of this person’s previous marketing campaign is failure.</p> <p>LLM prediction: No</p>
<p>Reordered features: This person’s last contact month of year is aug. This person’s last contact month of year is aug. This person’s last contact month of year is aug. 437.0 days have passed since this person was last contacted from a previous campaign. This person has 1.0 contacts performed before this campaign. The type of this person’s job is technician. The type of this person’s job is technician. This person has no personal loan. This person’s average yearly balance in euros is 2979.0. This person’s last contact day of the month is 5. This person has no personal loan. This person’s age is 33. This person has no house. This person has no house. The outcome of this person’s previous marketing campaign is failure. This person has no personal loan.</p> <p>LLM prediction: Yes</p>
<p>Reorder and Deduplication: This person’s last contact month of year is aug. 437.0 days have passed since this person was last contacted from a previous campaign. This person has 1.0 contacts performed before this campaign. The type of this person’s job is technician. This person has no personal loan. This person’s average yearly balance in euros is 2979.0. This person’s last contact day of the month is 5.0. This person has no personal loan. This person’s age is 33.0. This person has no house. The outcome of this person’s previous marketing campaign is failure. This person has no personal loan.</p> <p>LLM prediction: Yes</p>

Figure 5: Examples of LLM’s predictions based on default ordered features, reordered features, and reordered and deduplicated features.

5.3 TRANSFER-ABILITY OF CONTROLLER (RQ2)

In this section, we evaluate the transferability of the learned controller. The meta-controller is trained based on a source LLM and tested on a target LLM, marked as “source LLM→target LLM”. We take Llama-2-8B, Mistral-7B for the source LLMs, and Mistral-7B, Llama-2-8B for the target LLMs, respectively. The results of the controller transfer are shown in Table 3. It is observed that transferring the controller from one LLM to another achieves better performance than inputting the data instance following the default or random order. The results validate the transferability of our learned con-

432 troller, which meets our expectations as different LLMs could have similar order bias due to the fact
 433 that they all focus on learning the large human-generated content in pre-training.
 434

435 5.4 ATTENTION CONCENTRATION BY FEATURE RE-ORDERING (RQ3) 436

437 It has been widely shown in existing work Xiao et al. (2023); Zhang et al. (2024b) that the attention
 438 of LLM-generated tokens should focus on some key input tokens. Uniform patterns of attention can
 439 potentially lead to hallucinations. We conducted experiments to evaluate ROTATOR-LLM in terms
 440 of attention concentration. Specifically, this experiment is with Llama-3-8B on the bank dataset us-
 441 ing the prompts in Appendix A. The attention is estimated by $\text{softmax}(\mathbf{Q}[:, -1]\mathbf{K}^T/\sqrt{d})$, where \mathbf{Q} ,
 442 \mathbf{K} take the last-layer activations; d takes the hidden dimension value; and the index -1 of \mathbf{Q} indicates
 443 the attention is estimated for the answer token. To study the concentration of attention, we show
 444 the entropy of last layer attention in Figure 4 (a). The entropy is calculated by $-\sum_{p_j} p_j \log p_j$,
 445 where $p_j \sim \text{softmax}(\mathbf{Q}[:, -1]\mathbf{K}^T/\sqrt{d})$ are the attention weights obtained from the softmax opera-
 446 tion. Lower entropy corresponds to higher concentrations of attention on a small number of input
 447 tokens. It is observed that the last layer attention shows lower entropy after the feature re-ordering
 448 than utilizing the default order, indicating more focused attentions on the particular input tokens,
 449 rather than uniformly sprout to the whole prompt sequence. This contributes to a better aligned
 450 results in Tables 1 and 2.

451 5.5 CASE STUDIES (RQ3) 452

453 In this section, we show the data features reordered by ROTATOR-LLM. The data features in natural
 454 language sentences are shown in Figure 5, where the place holder `<Data Features>` takes the “Data
 455 features”, “Reordered features”, and “Reorder and Deduplication” below. We further investigate
 456 the affect of deduplication to LLMs’ performance in Figure 4, where the deduplication removes the
 457 duplicated features from the reordered data. Overall, we have the following insights:
 458

- 459 • **Significance of Feature Order.** A good feature order benefits LLMs more than a high number of
 460 features. The data instance has 16 features, and only 10 features left after reordering. However,
 461 LLMs show more accurate predictions based on the reordered data features.
- 462 • **Feature Order is Robust to Deduplication.** The features may be duplicated after the reordering
 463 because the features are reordered without replacement. As shown in Figure 4, LLMs maintain
 464 the performance to high-levels after removing the redundant features from the input context. This
 465 indicates the feature order is robust to the deduplication of redundant features.
 466

467 6 RELATED WORK 468

469 We discuss related work on tabular data understanding in this section. Existing work that leverages
 470 LLMs to process tabular data is primarily viewed from three perspectives: feature serialization,
 471 large-scale fine-tuning, and prompt engineering. We give more details as follows.
 472

473 **Feature Serialization.** Feature serialization is a simple way to let LLMs understand tabular data.
 474 Specifically, a straightforward way would be to directly input a programming-language readable data
 475 structure, such as Markdown format Liu et al. (2023); Jaitly et al. (2023), JSON-file format Singha
 476 et al. (2023); Sui et al. (2024), HTML format Singha et al. (2023), and Python dictionary Wang et al.
 477 (2023). Another way is to convert the tables into natural language sentence using templates based on
 478 the column headers and cell values Yu et al. (2023); Jaitly et al. (2023). This method can maximally
 479 leverage the sequence-to-sequence capacity of LLMs to understand tabular data.

480 **Large-scale Fine-tuning.** Fine-tuning on tabular datasets is a straightforward way to inject the
 481 data prior knowledge to LLMs. There are several existing work of fine-tuning. TableLlama adopts
 482 LongLoRA to fine-tune the Llama-2-7B LLM on the extensive TableInstruct datasets Zhang et al.
 483 (2023). TableGPT introduces a table encoder and chain-of-command mechanism and performs
 484 instruction tunings for Phoenix-7B LLMs on collections of tabular datasets Li et al. (2024). Different
 485 from existing work, TabLLM considers few-shot examples for prompts during the fine-tuning, and
 updates the Bigscience/T0-3B LLMs on single domain tabular datasets Zhang et al. (2024a).

In-context Learning. Existing work has demonstrated that LLMs are few-shot learners of tabular data Chen (2022); Narayan et al. (2022); Guo et al. (2023). Leveraging few-shot examples in the prompts, LLMs can better understand the data semantics through in-context learning. Other prompt engineering methods include chain-of-thoughts Wei et al. (2022), tree-of-thoughts Yao et al. (2024), self-consistency Wang et al. (2022), and others Sui et al. (2023).

7 CONCLUSION

In this work, we demonstrate novel discoveries and thoroughly explore the order bias of LLMs on tabular data, where the arrangement of data features can mislead LLM predictions. To address this issue, we propose ROTATOR-LLM, an approach designed to align tabular data with this order bias, enabling LLMs to better comprehend the data semantics. Specifically, ROTATOR-LLM employs a meta-controller to learn the optimal feature order. It estimates the value function for each feature order using dynamic programming, which guides the training of the meta-controller. Our experimental results on four datasets across three LLMs show that ROTATOR-LLM achieves superior performance compared to baseline methods and state-of-the-art foundational tabular LLMs when applied to reordered data. Additionally, ROTATOR-LLM exhibits strong transferability across multiple LLMs, indicating its adaptability to diverse tasks. Without requiring fine-tuning of LLMs, ROTATOR-LLM proves to be a more cost-effective solution than traditional debiasing methods, underscoring its potential for real-world applications.

REFERENCES

- Arthur Asuncion, David Newman, et al. Uci machine learning repository, 2007.
- Wenhu Chen. Large language models are few (1)-shot table reasoners. [arXiv preprint arXiv:2210.06710](#), 2022.
- Xinyun Chen, Ryan A Chi, Xuezhi Wang, and Denny Zhou. Premise order matters in reasoning with large language models. [arXiv preprint arXiv:2402.08939](#), 2024.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. [arXiv preprint arXiv:2403.04132](#), 2024.
- Haoyu Dong and Zhiruo Wang. Large language models for tabular data: Progresses and future directions. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2997–3000, 2024.
- Xi Fang, Weijie Xu, Fiona Anting Tan, Jiani Zhang, Ziqing Hu, Yanjun Jane Qi, Scott Nickleach, Diego Socolinsky, Srinivasan Sengamedu, Christos Faloutsos, et al. Large language models (llms) on tabular data: Prediction, generation, and understanding-a survey. 2024.
- Jiayan Guo, Lun Du, Hengyu Liu, Mengyu Zhou, Xinyi He, and Shi Han. Gpt4graph: Can large language models understand graph structured data? an empirical evaluation and benchmarking. [arXiv preprint arXiv:2305.15066](#), 2023.
- Stefan Hegselmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaoyi Jiang, and David Sontag. Tabllm: Few-shot classification of tabular data with large language models. In *International Conference on Artificial Intelligence and Statistics*, pp. 5549–5581. PMLR, 2023.
- Yaojie Hu, Ilias Fountalis, Jin Tian, and Nikolaos Vasiloglou. Annotatedtables: A large tabular dataset with language model annotations. [arXiv preprint arXiv:2406.16349](#), 2024.
- Sukriti Jaitly, Tanay Shah, Ashish Shugani, and Razik Singh Grewal. Towards better serialization of tabular data for few-shot classification. [arXiv preprint arXiv:2312.12464](#), 2023.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. [arXiv preprint arXiv:2401.04088](#), 2024.

- 540 Hongye Jin, Xiaotian Han, Jingfeng Yang, Zhimeng Jiang, Zirui Liu, Chia-Yuan Chang, Huiyuan
541 Chen, and Xia Hu. Llm maybe longlm: Self-extend llm context window without tuning. [arXiv](#)
542 [preprint arXiv:2401.01325](#), 2024.
- 543
544 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
545 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A
546 benchmark of in-the-wild distribution shifts. In [International conference on machine learning](#), pp.
547 5637–5664. PMLR, 2021.
- 548 Peng Li, Yeye He, Dror Yashar, Weiwei Cui, Song Ge, Haidong Zhang, Danielle Rifinski Fainman,
549 Dongmei Zhang, and Surajit Chaudhuri. Table-gpt: Table fine-tuned gpt for diverse table tasks.
550 [Proceedings of the ACM on Management of Data](#), 2(3):1–28, 2024.
- 551 Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee.
552 Textbooks are all you need ii: **phi-1.5** technical report. [arXiv preprint arXiv:2309.05463](#), 2023.
- 553
554 Tianyang Liu, Fei Wang, and Muhao Chen. Rethinking tabular data understanding with large lan-
555 guage models. [arXiv preprint arXiv:2312.16702](#), 2023.
- 556
557 Avanika Narayan, Ines Chami, Laurel Orr, Simran Arora, and Christopher Ré. Can foundation
558 models wrangle your data? [arXiv preprint arXiv:2205.09911](#), 2022.
- 559 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust
560 neural networks for group shifts: On the importance of regularization for worst-case generaliza-
561 tion. [arXiv preprint arXiv:1911.08731](#), 2019.
- 562
563 Ananya Singha, José Cambrero, Sumit Gulwani, Vu Le, and Chris Parnin. Tabular representa-
564 tion, noisy operators, and impacts on table structure understanding tasks in llms. [arXiv preprint](#)
565 [arXiv:2310.10358](#), 2023.
- 566
567 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
568 can be more effective than scaling model parameters. [arXiv preprint arXiv:2408.03314](#), 2024.
- 569
570 Yuan Sui, Jiaru Zou, Mengyu Zhou, Xinyi He, Lun Du, Shi Han, and Dongmei Zhang. Tap4llm:
571 Table provider on sampling, augmenting, and packing semi-structured data for large language
572 model reasoning. [arXiv preprint arXiv:2312.09039](#), 2023.
- 573
574 Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. Table meets llm: Can
575 large language models understand structured table data? a benchmark and empirical study. In
576 [Proceedings of the 17th ACM International Conference on Web Search and Data Mining](#), pp.
577 645–654, 2024.
- 578
579 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
580 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
581 tion and fine-tuned chat models. [arXiv preprint arXiv:2307.09288](#), 2023.
- 582
583 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
584 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
585 [arXiv preprint arXiv:2203.11171](#), 2022.
- 586
587 Zifeng Wang, Chufan Gao, Cao Xiao, and Jimeng Sun. Meditab: Scaling medical tabular data
588 predictors via data consolidation, enrichment, and refinement. [arXiv preprint arXiv:2305.12081](#),
589 2023.
- 590
591 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
592 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. [Advances in](#)
593 [neural information processing systems](#), 35:24824–24837, 2022.
- 594
595 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
596 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers:
597 State-of-the-art natural language processing. [arXiv preprint arXiv:1910.03771](#), 2019.
- 598
599 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
600 language models with attention sinks. [arXiv preprint arXiv:2309.17453](#), 2023.

594 Zhichao Xu, Daniel Cohen, Bei Wang, and Vivek Srikumar. In-context example ordering guided by
595 label distributions. [arXiv preprint arXiv:2402.11447](#), 2024.

596

597 Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Shaochen
598 Zhong, Bing Yin, and Xia Hu. Harnessing the power of llms in practice: A survey on chatgpt and
599 beyond. *ACM Transactions on Knowledge Discovery from Data*, 18(6):1–32, 2024.

600 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
601 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances
602 in Neural Information Processing Systems*, 36, 2024.

603

604 Bowen Yu, Cheng Fu, Haiyang Yu, Fei Huang, and Yongbin Li. Unified language representation for
605 question answering over text, tables, and images. [arXiv preprint arXiv:2306.16762](#), 2023.

606 Jiayi Yuan, Hongyi Liu, Yu-Neng Chuang, Songchen Li, Guanchu Wang, Duy Le, Hongye Jin,
607 Vipin Chaudhary, Zhaozhuo Xu, Zirui Liu, et al. Kv cache compression, but what must we
608 give in return? a comprehensive benchmark of long context capable approaches. [arXiv preprint
609 arXiv:2407.01527](#), 2024.

610 Daochen Zha, Zaid Pervaiz Bhat, Kwei-Herng Lai, Fan Yang, Zhimeng Jiang, Shaochen Zhong, and
611 Xia Hu. Data-centric artificial intelligence: A survey. [arXiv preprint arXiv:2303.10158](#), 2023a.

612

613 Liangyu Zha, Junlin Zhou, Liyao Li, Rui Wang, Qingyi Huang, Saisai Yang, Jing Yuan, Changbao
614 Su, Xiang Li, Aofeng Su, et al. Tablegpt: Towards unifying tables, nature language and commands
615 into one gpt. [arXiv preprint arXiv:2307.08674](#), 2023b.

616 Tianshu Zhang, Xiang Yue, Yifei Li, and Huan Sun. Tablellama: Towards open large generalist
617 models for tables. [arXiv preprint arXiv:2311.09206](#), 2023.

618

619 Xiaokang Zhang, Jing Zhang, Zeyao Ma, Yang Li, Bohan Zhang, Guanlin Li, Zijun Yao, Kangli Xu,
620 Jinchang Zhou, Daniel Zhang-Li, et al. Tablellm: Enabling tabular data manipulation by llms in
621 real office usage scenarios. [arXiv preprint arXiv:2403.19318](#), 2024a.

622 Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song,
623 Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient generative
624 inference of large language models. *Advances in Neural Information Processing Systems*,
625 36, 2024b.

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

```

648 1 table2text_template = {
649 2   "age": "This_person's_age_is_{}.",
650 3   "job": "The_type_of_this_person's_job_is_{}.",
651 4   "marital": "This_person's_marital_status_is_{}.",
652 5   "education": "This_person's_education_is_{}.",
653 6   "default": {"no": "This_person_has_no_credit_in_default.",
654 7               "yes": "This_person_has_credit_in_default."},
655 8   "balance": "This_person's_average_yearly_balance_in_euros_is_{}.",
656 9   "housing": {"no": "This_person_has_no_house.",
657 10              "yes": "This_person_owns_houses."},
658 11   "loan": {"no": "This_person_has_no_personal_loan.",
659 12            "yes": "This_person_has_personal_loan."},
660 13   "contact": "This_person's_contact_communication_type_is_{}.",
661 14   "day": "This_person's_last_contact_day_of_the_month_is_{}.",
662 15   "month": "This_person's_last_contact_month_of_year_is_{}.",
663 16   "duration": "This_person's_last_contact_duration_is_{}_seconds.",
664 17   "campaign": "This_person_has_{}_contacts_performed_during_this_
665 18               campaign.",
666 19   "pdays": "{_}days_have_passed_since_this_person_was_last_contacted_
667 20               from_a_previous_campaign.",
668 21   "previous": "This_person_has_{}_contacts_performed_before_this_
669 22               campaign.",
670 23   "outcome": "The_outcome_of_this_person's_previous_marketing_campaign_
671 24               is_{}."
672 25 }

```

Figure 6: Table to Text data template on the bank dataset.

```

671 1 table2text_template = {
672 2   "workclass": "The_class_of_this_person's_job_is_{}.",
673 3   "marital_status": "This_person's_marital_status_is_{}.",
674 4   "education": "This_person's_education_is_{}.",
675 5   "occupation": "This_person's_job_is_{}.",
676 6   "relationship": "This_person's_relationship_in_family_is_{}.",
677 7   "sex": "This_person's_gender_is_{}.",
678 8   "race": "This_person's_race_is_{}.",
679 9   "native_country": "The_native_country_of_this_person_is_{}.",
680 10   "age": "This_person's_age_is_{}.",
681 11   "fnlwgt": "The_final_analysis_weight_of_this_person_is_{}.",
682 12   "education_num": "The_education_duration_of_this_person_is_{}.",
683 13   "capital_gain": "The_capital_gain_of_this_person_is_{}.",
684 14   "capital_loss": "The_capital_loss_of_this_person_is_{}.",
685 15   "hours_per_week": "The_person_works_{}_hours_per_week_in_average.",
686 16 }

```

Figure 7: Table to Text data template on the Income dataset.

APPENDIX

A TEMPLATE OF TEXT-BASED SERIALIZATION

We give the template of text-based serialization in this work. The templates for the bank, Income, German Credit, and Diabete datasets are given in Figures 6, 7, 8, and 9, respectively.

B HYPER-PARAMETER SETTING OF ROTATOR-LLM

The hyper-parameter setting of ROTATOR-LLM in Table 4. The discounting factor for meta-controller training is given in Table 5.

```

702
703
704 1 table2text_template = {
705 2   "checking_status": "The_status_of_this_person's_checking_account_is_
706 3   {},",
707   "credit_history": "The_status_of_this_person's_historical_credits_is_
708 4   {},",
709 5   "purpose": "This_person's_purpose_to_apply_for_credits_is_{}",
710   "savings_status": "The_status_of_this_person's_saving_account_is_{}."
711 6   ,
712   "employment": "The_present_employment_of_this_person_is_{}",
713   "personal_status": "The_marital_status_of_this_person_is_{}",
714   "other_parties": {"none": "This_person_does_not_have_other_debtors.",
715   "co_applicant": "This_person_has_co-applicants.",
716   "guarantor": "This_person_has_guarantors."} ,
717   "property_magnitude": "The_property_magnitude_of_this_person_is_{}",
718   "other_payment_plans": {"none": "This_person_does_not_have_other_
719   installment_plans.",
720   "stores": "This_person_has_installment_plans_for_stores.",
721   "bank": "This_person_has_installment_plans_for_banks."},
722   "housing": {"own": "This_person_owns_houses.",
723   "rent": "This_person_rents_a_house.",
724   "for_free": "This_person_lives_in_a_free_house."},
725   "job": "The_type_of_this_person's_job_is_{}",
726   "own_telephone": {"none": "This_person_does_not_have_a_telephone.",
727   "yes": "This_person_owns_a_telephone."},
728   "foreign_worker": {"yes": "This_person_is_a_foreign_worker.",
729   "no": "This_person_is_not_a_foreign_worker."},
730   "duration": "The_duration_of_this_person_is_{}_months.",
731   "credit_amount": "The_amount_of_this_person's_credit_is_{}",
732   "installment_commitment": "This_person_has_a_installment_rate_of_{}_
733   of_disposable_income.",
734   "residence_since": "This_person_has_been_a_residence_for_{}_years.",
735   "age": "This_person's_age_is_{}",
736   "existing_credits": "This_person_already_has_{}_credits.",
737   "num_dependents": "This_person_supports_{}_dependents.",
738 }

```

Figure 8: Table to Text data template on the German Credit dataset.

Name	Value
Layer Number	3
Hidden Dimension	512
Optimizer	Adam
Learning Rate	0.001
Epoch	200
Mini-batch Size	128

Table 4: Hyper-parameter setting of ROTATOR-LLM.

	Bank	Income	German Credit	Diabete
Llama-3-8B-Instruct	0.75	0.8	0.8	0.8
Mistral-7B-Instruct	0.85	0.9	0.85	0.9
Phi-3-Mini-Instruct	0.9	0.8	0.8	0.8

Table 5: Discounting factor on meta-controller training.

```

756
757 1 table2text_template = {
758 2   "HighBP": {0: "This_person_has_a_normal_blood_pressure.",
759 3     1: "This_person_has_a_high_blood_pressure."},
760 4   "HighChol": {0: "This_person_has_normal_cholesterol.",
761 5     1: "This_person_has_high_cholesterol."},
762 6   "CholCheck": {0: "This_person_has_no_cholesterol_check_in_5_years.",
763 7     1: "This_person_has_cholesterol_checks_in_5_years."},
764 8   "BMI": "This_person's_Body_Mass_Index_is{}",
765 9   "Smoker": {0: "This_person_smoked_less_than_100_cigarettes_in_the_
766 10     entire_life.",
767 11     1: "This_person_smoked_at_least_100_cigarettes_in_the_entire_life
768 12     ."},
769 13   "Stroke": {0: "This_person_does_not_have_a_stroke.",
770 14     1: "This_person_has_a_stroke."},
771 15   "HeartDiseaseorAttack": {0: "This_person_does_not_have_coronary_heart
772 16     disease_(CHD)_or_myocardial_infarction.",
773 17     1: "This_person_has_a_coronary_heart_disease_(CHD)_or_myocardial_
774 18     infarction."},
775 19   "PhysActivity": {0: "This_person_did_not_have_physical_activities_in_
776 20     the_past_30_days.",
777 21     1: "This_person_had_physical_activities_in_the_past_30_days."},
778 22   "Fruits": {0: "This_person_does_not_consume_fruit_every_day.",
779 23     1: "This_person_consumes_fruit_one_or_more_times_every_day."},
780 24   "Veggies": {0: "This_person_does_not_consume_vegetables_every_day.",
781 25     1: "This_person_consumes_vegetables_one_or_more_times_every_day."
782 26   },
783 27   "HvyAlcoholConsump": {0: "This_person_is_not_a_heavy_drinker_(adult_
784 28     men_having_more_than_14_drinks_per_week_and_adult_women_having_
785 29     more_than_7_drinks_per_week).",
786 30     1: "This_person_is_a_heavy_drinker_(adult_men_having_more_than_14
787 31     _drinks_per_week_and_adult_women_having_more_than_7_drinks_
788 32     per_week)."},
789 33   "AnyHealthcare": {0: "This_person_does_not_Have_any_kind_of_health_
790 34     care_coverage,_including_health_insurance,_prepaid_plans_such_as_
791 35     HMO.",
792 36     1: "This_person_has_any_kind_of_health_care_coverage,_including_
793 37     health_insurance,_prepaid_plans_such_as_HMO."},
794 38   "NoDocbcCost": {0: "This_person_never_misses_a_doctor_because_of_cost_
795 39     in_the_past_12_months.",
796 40     1: "This_person_once_needed_to_see_a_doctor_but_could_not_because_
797 41     of_cost_in_the_past_12_months."},
798 42   "GenHlth": "This_person's_general_health_score_is{}__(1_represents_
799 43     the_best,_and_5_represents_the_worst).",
800 44   "MentHlth": "This_person_had_stress,_depression,_or_problems_with_
801 45     emotions_in{}_days_of_the_past_30_days.",
802 46   "PhysHlth": "This_person_had_a_physical_illness_or_injury_in{}_days_
803 47     of_the_past_30_days.",
804 48   "DiffWalk": {0: "This_person_does_not_have_serious_difficulty_
805 49     walking_or_climbing_stairs.",
806 50     1: "This_person_has_serious_difficulty_walking_or_climbing_stairs
807 51     ."},
808 52   "Sex": {0: "This_person_is_a_female.",
809 53     1: "This_person_is_a_male."},
810 54   "Age": "This_person's_age_is{}.",
811 55   "Education": {
812 56     1: "This_person_never_attended_school_or_only_kindergarten.",
813 57     2: "This_person_has_grades_1_through_8_(Elementary).",
814 58     3: "This_person_has_grades_9_through_11_(Some_high_school).",
815 59     4: "This_person_has_grade_12_or_GED_(High_school_graduate).",
816 60     5: "This_person_has_college_1_year_to_3_years_(Some_college_or_
817 61     technical_school).",
818 62     6: "This_person_has_college_4_years_or_more_(College_graduate).",
819 63   },
820 64   ...

```

Figure 9: Table to Text data template on the Diabete dataset (i).

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

```
1  ...  
2  "Income": {  
3    1: "This_person's_income_is_less_than_10000_dollars.",  
4    2: "This_person's_income_is_more_than_10000_dollars_but_less_than  
5      _15000_dollars.",  
6    3: "This_person's_income_is_more_than_15000_dollars_but_less_than  
7      _20000_dollars.",  
8    4: "This_person's_income_is_more_than_20000_dollars_but_less_than  
9      _25000_dollars.",  
10   5: "This_person's_income_is_more_than_25000_dollars_but_less_than  
11     _35000_dollars.",  
12   6: "This_person's_income_is_more_than_35000_dollars_but_less_than  
     _55000_dollars.",  
     7: "This_person's_income_is_more_than_55000_dollars_but_less_than  
     _75000_dollars.",  
     8: "This_person's_income_is_more_than_75000_dollars.",  
  },  
}
```

Figure 10: Table to Text data template on the Diabete dataset (ii).