

# MAKING PRE-TRAINED LANGUAGE MODELS GREAT ON TABULAR PREDICTION

Jiahuan Yan<sup>1,2,\*</sup>, Bo Zheng<sup>2</sup>, Hongxia Xu<sup>2</sup>, Yiheng Zhu<sup>2</sup>, Danny Z. Chen<sup>3</sup>, Jimeng Sun<sup>4</sup>, Jian Wu<sup>1,2,†</sup>, Jintai Chen<sup>4,†</sup>

<sup>1</sup>The Second Affiliated Hospital Zhejiang University School of Medicine <sup>2</sup>Zhejiang University

<sup>3</sup>University of Notre Dame <sup>4</sup>University of Illinois at Urbana-Champaign

{jyansir, zjuzhengbo, einstein, zhuyiheng2020, wujian2000}@zju.edu.cn, dchen@end.edu, jimeng@illinois.edu, jtchen721@gmail.com

## ABSTRACT

The transferability of deep neural networks (DNNs) has made significant progress in image and language processing. However, due to the heterogeneity among tables, such DNN bonus is still far from being well exploited on tabular data prediction (e.g., regression or classification tasks). Condensing knowledge from diverse domains, language models (LMs) possess the capability to comprehend feature names from various tables, potentially serving as versatile learners in transferring knowledge across distinct tables and diverse prediction tasks, but their discrete text representation space is inherently incompatible with numerical feature values in tables. In this paper, we present *TP-BERTa*, a specifically pre-trained LM for tabular data prediction. Concretely, a novel *relative magnitude tokenization* converts scalar numerical feature values to finely discrete, high-dimensional tokens, and an *intra-feature attention* approach integrates feature values with the corresponding feature names. Comprehensive experiments demonstrate that our pre-trained TP-BERTa leads the performance among tabular DNNs and is competitive with Gradient Boosted Decision Tree models in typical tabular data regime.

## 1 INTRODUCTION

Tabular data, a common data form, is pivotal in various fields such as medical trial predictions (Hasan et al., 2020) and financial risk detection (Aziz et al., 2022). The remarkable successes of deep neural networks (DNNs) in computer vision (CV) and natural language processing (NLP) have spurred interest in applying DNNs to tabular data for tasks like classification or regression (Popov et al., 2020; Song et al., 2019; Wang et al., 2021; Chen et al., 2023b), which also pave the road for cross-modality processing. However, most current research on tabular data relies on fully supervised paradigms (Arik & Pfister, 2021; Gorishniy et al., 2021; Somepalli et al., 2022; Li et al., 2023; Hollmann et al., 2023), and with typically limited data available for DNNs in this regime, Gradient Boosted Decision Trees (GBDTs) (Chen & Guestrin, 2016; Ke et al., 2017; Prokhorenkova et al., 2018) continue to outperform these paradigms (Grinsztajn et al., 2022).

As widely evidenced in the CV and NLP fields, the transferability of DNNs consistently brought about substantial performance boosts and decreased data demands in downstream tasks (Devlin et al., 2018; Xie et al., 2020; He et al., 2020). However, how to utilize the transferability of DNNs on tabular data is still much under-explored. One major obstacle is the feature heterogeneity among tables (Borisov et al., 2022; Yan et al., 2023; Chen et al., 2022). Unlike images, which often exhibit similar feature distributions (e.g., consistent pixel intensity ranges and color distributions) (Chen et al., 2023a), structured tables inherently contain diverse columns and feature spaces, leading to considerable heterogeneity and feature space shifts between pre-training and downstream datasets.

**Related Work.** Recent studies highlight the importance of tabular transfer learning, with initial efforts like TransTab (Wang & Sun, 2022) and XTab (Zhu et al., 2023) utilizing shared Transformer

\*Work under partial support from the Second Affiliated Hospital Zhejiang University School of Medicine.

†Corresponding authors. Codes are available at <https://github.com/jyansir/tp-berta>.

blocks in the FT-Transformer architecture (Gorishniy et al., 2021) for cross-table learning. TransTab focused on clinical trial tables with common feature names, facilitating partially overlapped feature embeddings, whereas XTab explored a broader domain with dataset-specific encoders. However, neither achieved comprehensive knowledge transfer, resulting in moderate pre-training performance. The advancements in language models (LMs) have demonstrated their capability to act as common-sense knowledge bases (Petroni et al., 2019; Jiang et al., 2020; Gao et al., 2021; Zha et al., 2023). Through self-supervised pre-training on extensive domain-agnostic corpora, LMs can implicitly capture associations among different words or phrases, showcasing potential as tabular transfer agents with their inherent support for feature name processing within a unified language space. Despite this potential, early attempts of applying LMs to tabular prediction were limited to synthetic table generation (e.g., missing value imputation) and faced challenges. GReaT (Borisov et al., 2023) and TapTap (Zhang et al., 2023) fine-tuned GPT-2 (Radford et al., 2019) on simply templated table texts, treating numerical values as strings, which led to insensitivity to such values (Qian et al., 2023). A contemporary work (Ye et al., 2024) developed a BERT-based model (CT-BERT) using a large tabular database and similar techniques to TransTab. However, these studies overlooked the customization of LMs for understanding continuous numerical values, which is a critical aspect of tables and presents challenges to LMs due to their inherent complexity and rarity (Qian et al., 2023).

To unlock LMs’ power and take a pioneering step on LM-based tabular transfer learning, in this paper, we propose a tailored pre-trained LM for tabular prediction based on RoBERTa (Liu et al., 2019), called the *Tabular Prediction adapted BERT approach (TP-BERTa)*. TP-BERTa maintains the strengths of LMs as well as possessing the sensitivity to numeric features. Specifically, TP-BERTa discretizes numerical feature values as relative magnitude tokens (RMT) in order to treat them as some meaningful words in the LM’s vocabulary. The design of relative magnitude tokens enables the LM to perceive relative value magnitudes in the language space. In this way, we decouple the representations of feature names and numerical values (compared to FT-Transformer, TransTab, and CT-BERT), preserving the semantic signal of feature names. Further, we develop a shared intra-feature attention (IFA) module to attentively fuse the embeddings of a feature’s name and value into a single vector. IFA retains the text order in a feature name, and outputs a vector for each feature name-value pair to the subsequent LM process to achieve feature order-agnostic prediction.

We pre-train TP-BERTa on numerous large tabular datasets (101 binary classification and 101 regression datasets), and provide three versions (i.e., pre-trained on only classification tasks, or only regression tasks, or both). We conduct evaluations on extensive downstream datasets: (1) performance comparison with classical GBDTs, advanced deep tabular models, and cross-table models shows that our TP-BERTa (the pre-trained versions on a single task type, with default hyperparameters) outperforms the other tabular DNNs and is competitive with GBDTs in the overall rank on 145 downstream datasets; (2) comparison with two existing numerical encoding strategies (Borisov et al., 2023; Ye et al., 2024) shows that our RMT adaption achieves average AUC improvements of 12.45% and 3.44% on significantly changed (i.e., with AUC variation over 0.5%) downstream binary classification datasets, respectively; (3) ablation on table-specific designs for LM adaption.

**Contributions.** In a nutshell, our work offers: (1) **A pre-trained LM for tabular data:** dealing with fundamental difficulties in LM adaption to tabular data (i.e., numeric feature handling and tabular feature organization), we develop LM-based tabular DNNs and pre-train a tabular-data-tailored LM called TP-BERTa; (2) **superior performances:** comparisons with various existing methods on 145 downstream datasets demonstrate that pre-trained LMs can outperform common tabular DNNs and are competitive with GBDTs in typical tabular regime; (3) **in-depth analysis:** multi-facet comparison implies that TP-BERTa has a data appetite of informative discrete features, and key ablation experiments show that our RMT and IFA adaptations are successful.

## 2 TP-BERTA: TABULAR PREDICTION ADAPTED BERT APPROACH

Our proposed TP-BERTa is built on the basis of RoBERTa (Liu et al., 2019) as default. Its model architecture and key components (the *relative magnitude tokenization* approach and *intra-feature attention* module) are shown in Fig. 1. Below we introduce our novel (i) relative magnitude tokenization (RMT) for numerical value representation, (ii) intra-feature attention (IFA) module for feature name-value matching before the LM processing, and (iii) the overall pre-training paradigm.

## 2.1 RELATIVE MAGNITUDE TOKENIZATION

Tabular features can be roughly categorized into continuous type (i.e., numerical features) and discrete type (categorical, binary, or string features). Although discrete feature values with clear semantics (e.g., “male” and “female” are values of a discrete feature “gender”) can be naturally understood by LMs, it is still difficult to make numerical features fully understandable to LMs due to their wide range of values and counter-intuitive meanings of exact numerical values. In this section, we present a novel Relative Magnitude Tokenization (RMT) approach to boost numerical value understanding.

**Numerical Discretization.** Our RMT process is inspired by classical works on feature binning (Dougherty et al., 1995; Gorishniy et al., 2022) that utilized discretization techniques for numerical features. To deal with diverse labeled tabular datasets, we adopt a target-aware binning method similar to (Gorishniy et al., 2022). Specifically, the “C4.5 Discretization” algorithm (Kohavi & Sahami, 1996) is applied to each numerical feature by recursively splitting its value range guided by its label. This process is equivalent to building a decision tree, and continuous values are grouped into corresponding tree leaves. The boundary values of all the leaves are used to split the value range into multiple bins. Each numerical value is converted to its bin index after discretization, as:

$$\mathbf{e}^{(i)} = \text{C4.5}(\mathbf{x}^{(i),\text{train}}, \mathbf{y}^{(i),\text{train}}), \quad (1)$$

$$\text{BinIndex}(x_j^{(i)}) \equiv k, \quad (2)$$

where  $\mathbf{x}^{(i),\text{train}}$  is the vector of the  $i$ -th numerical feature values in the training set,  $\mathbf{y}^{(i),\text{train}}$  is the corresponding labels,  $\mathbf{e}^{(i)}$  denotes the vector of leaf node boundary values (in ascending order),  $x_j^{(i)}$  is the  $i$ -th feature value of sample  $j$ , and  $k$  is its bin index if  $e_k^{(i)} \leq x_j^{(i)} < e_{k+1}^{(i)}$ . In TP-BERTa, we set the maximum numerical bin (magnitude token) number (denoted as  $n_{\text{bin}}$ ) to 256 (i.e.,  $0 \leq k < 256$ ), unless otherwise specified. A bin index represents a relative magnitude in the value range.

**Magnitude Tokenization.** To transform numerical values into the language space, we treat the numerical bins as new words. Specifically,  $n_{\text{bin}}$  additional tokens are added to the RoBERTa vocabulary with randomly initialized token embeddings. Each numerical value is discretized with a feature-specific C4.5 process and mapped to these shared magnitude tokens. Since there may be a large number of values in a single numerical bin, the final token embedding of a numerical value is its corresponding bin token embedding multiplied with the value itself, i.e.,  $\text{RMT}(x_j^{(i)}) \equiv \mathbf{E}_{:,k}^{\text{extra}} \times x_j^{(i)}$ , where  $\mathbf{E}_{:,k}^{\text{extra}}$  denotes the  $k$ -th embedding of the RoBERTa additional vocabulary for the numerical magnitude. These embeddings are shared across any numerical features or datasets that purely represent relative magnitudes with word vectors. Just as LMs show general capability of language modeling based on reasonable pair-wise word similarity, we seek to make the designed “magnitude embeddings” follow a similar relationship. Hence, we devise a magnitude-aware triplet loss to regularize the learning process of the magnitude embeddings. We formulate the regularization as:

$$L_{\text{reg}} = \max(d(f(k_1), f(k_2)) - d(f(k_1), f(k_3)) + m(k_1, k_2, k_3), 0), \quad (3)$$

$$s.t. \ |k_1 - k_2| < |k_1 - k_3|, \quad (4)$$

$$m(k_1, k_2, k_3) = \frac{|k_1 - k_3| - |k_1 - k_2|}{n_{\text{bin}}}, \quad (5)$$

where  $k_1, k_2$ , and  $k_3$  are three bin indices, and  $d(\mathbf{x}, \mathbf{y})$  is the  $L_2$  distance between vectors  $\mathbf{x}$  and  $\mathbf{y}$ . In a nutshell, this regularization process assists to pull the embedding of a bin close to the embedding of a nearby one, while pushing away from the embedding of a bin far away from it, serving as an auxiliary loss to help embedding learning for magnitude tokens.

**Tabular Feature Pre-processing.** A tabular sample may contain features of different types. We process each feature  $i$  by simply concatenating the embeddings of its feature name ( $E_i^{\text{name}} \in \mathbb{R}^{l_1 \times d}$ ) and value ( $E_i^{\text{value}} \in \mathbb{R}^{l_2 \times d}$ ), i.e.,  $E_i = E_i^{\text{name}} \otimes E_i^{\text{value}}$ , where  $d$  is the hidden dimension of the RoBERTa embeddings,  $l_1$  is the token length of the feature name, and  $l_2$  is the length of the feature value. Notably,  $l_2 \equiv 1$  for numerical features. As for categorical features, we convert their values into structured texts (e.g., value “0” of the feature “gender” is mapped to “male”). Note that we do not distinguish binary and categorical ones in this paper since they are both converted to meaningful

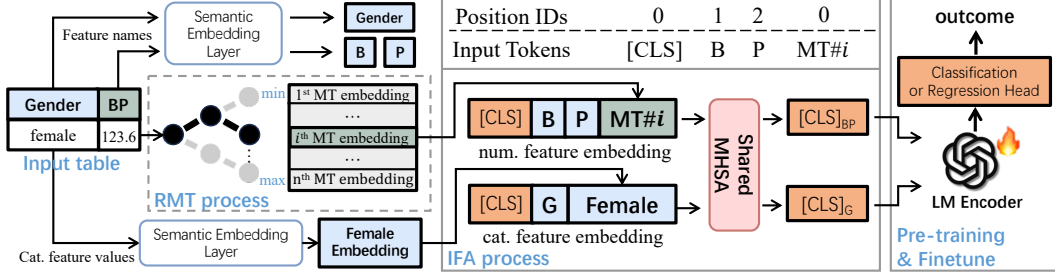


Figure 1: Illustrating the TP-BERTa workflow. “BP” in the input table denotes the feature name text “blood pressure”. The rectangles with “B”, “P”, and “Gender” (“G”) represent word embedding of “blood”, “pressure”, and “gender”, respectively. In the RMT process, numerical values are discretized by the feature-specific C4.5 decision tree. In the IFA process, “MT# $i$ ” indicates the  $i$ -th magnitude token. All numerical features share these MT embeddings for magnitude representation. “MHSa” is a shared multi-head self-attention across all features for feature refinement.

texts. Some datasets contain string features, such as a feature “movie comment” with unstructured texts. We process the values of these feature types in the same way as for feature names.

## 2.2 INTRA-FEATURE ATTENTION MODULE

Previous attempts of using LMs to process tables still face three lingering issues. (1) Targets in tabular predictions are independent of feature permutations, while LMs inherently process texts with positional encoding since positions of linguistic units matter. (2) When we simply feed all feature values with names into a vanilla LM (e.g., “[Gender] is female, [Blood Pressure] is 123.8”), it likely increases the training difficulty of LMs since they have to understand the correctly matched name-value pairs of features and learn to alleviate interference from other features. However, fully connected attention mechanism (commonly adopted in auto-encoder LMs) makes it inevitable to generate mismatched name-value signal. (3) Feeding the whole templated text can incur computation burden caused by excessively long sequences when the feature amount is large. Recently, a solution was given for issue (1) by augmenting a sample with copies of different feature permutations (Borisov et al., 2023), and position encoding was directly dropped and text order of feature names was ignored (Ye et al., 2024). But, they all neglected issues (2) and (3). Hence, we develop the intra-feature attention (IFA) module for feature refinement before feeding features to the LM.

IFA is essentially a single multi-head self-attention (MHSa) module shared across all features and datasets. It accepts embeddings of a feature name-value pair and fuses them into a single vector. We formulate the process of IFA fusion on a single feature  $i$  as:

$$\mathbf{H}^{(i)} = e^{\text{CLS}} \otimes \mathbf{E}^{(i)}, \quad (6)$$

$$\mathbf{Q}^{(i)} = \mathbf{W}_q^T(\mathbf{H}^{(i)} + \mathbf{P}^{(i)}), \quad \mathbf{K}^{(i)} = \mathbf{W}_k^T(\mathbf{H}^{(i)} + \mathbf{P}^{(i)}), \quad \mathbf{V}^{(i)} = \mathbf{W}_v^T \mathbf{H}^{(i)}, \quad (7)$$

$$\hat{\mathbf{H}}^{(i)} = \text{MHSa}(\mathbf{Q}^{(i)}, \mathbf{K}^{(i)}, \mathbf{V}^{(i)}), \quad \hat{\mathbf{h}}^{(i)} \equiv \hat{\mathbf{H}}^{(i)}_{:, \text{Index}(\text{CLS})}, \quad (8)$$

where  $\mathbf{E}^{(i)} \in \mathbb{R}^{(l_1+l_2) \times d}$  is concatenation of name-value embeddings,  $e^{\text{CLS}} \in \mathbb{R}^{1 \times d}$  is the [CLS] embedding,  $\mathbf{W}_q$ ,  $\mathbf{W}_k$ , and  $\mathbf{W}_v \in \mathbb{R}^{d \times d}$  are transformations for query, key, and value vectors, and  $\mathbf{P}^{(i)} \in \mathbb{R}^{(1+l_1+l_2) \times d}$  is position embeddings. IFA uses the output vector at the [CLS] position  $\hat{\mathbf{h}}^{(i)}$  as refined feature information and feeds it to the subsequent RoBERTa. It can be clearly seen that information from both the name and value is included in  $\hat{\mathbf{h}}^{(i)}$ , and information from other feature names or values cannot corrupt feature  $i$ ’s representation. As shown in Fig. 1(IFA process), the positions of the [CLS] token and magnitude token are assigned to id 0, and those of feature names are from 1 to  $l_1$ . This design aims to make the [CLS] token pay more attention to values (which are probably more important for prediction) as well as keeping the text order of feature names. Notably, we remove position encoding on value vectors (see Eq. (7)); the key reason for this is to protect magnitude token embeddings from the impact of embeddings at a constant id position (e.g., position id 0). Since magnitude embeddings are randomly initialized and intentionally regularized to represent the meaning of the relative magnitude carefully, a constant signal may distort the representations and thus make the embedding learning process more difficult.

### 2.3 OVERALL TRAINING PARADIGM

After features are processed by the IFA module, an  $n$ -feature sample is organized as the concatenation of feature vectors and a [CLS] embedding to be the RoBERTa input, i.e.,  $\mathbf{X} \equiv e^{\text{CLS}} \otimes \hat{\mathbf{h}}_1 \otimes \hat{\mathbf{h}}_2 \otimes \dots \otimes \hat{\mathbf{h}}_n \in \mathbb{R}^{(1+n) \times d}$ , which is computation-friendly. Since the text order of feature names has been considered in  $\hat{h}_i$ , we can avoid position encoding in this step, and achieve feature order-agnostic prediction. The prediction is based on the [CLS] output of the *RoBERTa-Encoder*, as:

$$\hat{\mathbf{y}}_m = \text{PredictionHead}^{(m)}(\text{RoBERTa-Encoder}(\mathbf{X}^{(m)})_{\cdot, \text{Index}(\text{CLS})}), \quad (9)$$

$$\text{PredictionHead}(\mathbf{x}) = \text{Dropout}(\text{Linear}_1(\text{Tanh}(\text{Linear}_2(\mathbf{x})))), \quad (10)$$

where  $\mathbf{X}^{(m)}$  represents the input from the  $m$ -th task (dataset), and we use task-specific prediction heads  $\text{PredictionHead}^{(m)}$ , the shared RoBERTa, and the IFA module (constituting TP-BERTa) to perform supervised pre-training on extensively large tabular datasets. The final pre-training loss consists of supervised loss and regularization loss (see Eq. (3)), as:

$$L = L_{\text{sup}} + \lambda L_{\text{reg}}, \quad (11)$$

where for the supervised loss  $L_{\text{sup}}$ , we use binary cross entropy loss for binary classification tasks and mean squared error loss for regression tasks. We keep a constant weight  $\lambda \equiv 0.1$  in pre-training. For downstream tasks, ordinary finetune is adopted only with  $L_{\text{sup}}$ . We exclude multi-class datasets in this work as in (Grinsztajn et al., 2022), for the reasons: (1) they can be decomposed into multiple binary classification tasks, (2) the trends on binary classification can essentially reflect the classification ability, and (3) multi-class datasets are not very common in tabular dataset collections.

## 3 EXPERIMENTS

We first compare our TP-BERTa with classical and advanced tabular prediction models, including (1) the dominating GBDTs, (2) advanced deep tabular models, and (3) recent open-source cross-table models or pre-trained tabular models. We utilize extensive downstream datasets, and analyze the huge potential of our pre-trained LM, TP-BERTa, as a powerful tabular prediction learner from the data perspective (Sec. 3.2). Based on that, we further demonstrate how the encoding strategy of numerical values impacts the LMs’ performances, and discuss why they were neglected in previous tabular prediction research (Sec. 3.3). Transferability evaluations (Sec. 3.4) and design ablations (Sec. 3.5) are conducted to reflect the generalization capability and rational adaption of TP-BERTa.

### 3.1 EXPERIMENTAL DETAILS

**Datasets.** We leverage the high-quality large semantic tabular database TabPertNet (Ye et al., 2024). Datasets with at least 10,000 samples and no more than 32 features are taken for pre-training, and datasets with fewer than 10,000 samples are collected as downstream tasks (following the same “typical tabular data” settings of “medium-sized dataset regime” and “not high dimensional” in (Grinsztajn et al., 2022)). We strictly remove the same datasets in the database and make sure that no subset of pre-training datasets (e.g., a small version of a large dataset) appears in the downstream ones. Since LMs are fueled by meaningful texts, we manually exclude datasets with uninformative feature names (e.g., feature names like “v1, v2, x1, x2”) or unmappable categorical features (e.g., a feature “job” with values “0, 1, 2”). Note that those excluded datasets can still benefit from our model with simple feature preprocessing with their corresponding data dictionaries. In total, our pre-training datasets consist of 101 binary classification datasets and 101 regression datasets with about 10 million samples, and our downstream datasets consist of 80 binary classification datasets and 65 regression datasets. Detailed dataset statistics are provided in Appendix B.

**Pre-training Details.** Since our work does not focus on the curriculum learning issue, we warp all the datasets into a large data-loader, which provides a data batch from a randomly selected dataset per training step. Each dataset is learned with a dataset-specific prediction head and the shared TP-BERTa (see Sec. 2.3). Because the massive LM is likely to overfit a single dataset, we use 5% of the training data as the validation set. For binary classification, we keep the same label distributions for the training set and validation set. Pre-training is conducted on four NVIDIA A100 Tensor Core GPUs, with a total batch size of 512 per step. We reuse the weights of the RoBERTa-base as the

Table 1: The average values (standard deviations) of all method ranks on the dataset collections of two task types. “(d)” in the “Baselines” means using default hyperparameters, and “(t)” for using tuned ones. “Ours<sub>j</sub>” is TP-BERTa pre-trained on both binary classification and regression tasks, and “Ours<sub>s</sub>” contains two models pre-trained on the corresponding single-type tasks separately. “All” denotes rank information calculated on all the datasets,  $\alpha$  is the amount ratio of categorical features and numerical ones in a dataset, and  $\beta$  is the ratio of the Shapley value sums between the two feature types.  $\alpha$  or  $\beta$  provides a reference on the dominating feature type in tabular data: “ $\alpha \geq 1$ ” represents that only the datasets with their  $\alpha \geq 1$  are considered (similar denotations are for the others). The top performances are marked in **bold**, and the second best ones are underlined. We present feature type distribution statistics,  $\alpha$  and  $\beta$  formulation, and the original performances in Appendix B.

Baselines	80 downstream binary classification tasks						65 downstream regression tasks					
	All	$\alpha > 0$	$\alpha \geq 1$	$\alpha = 0$	$\beta > 0$	$\beta > 0.5$	All	$\alpha > 0$	$\alpha \geq 1$	$\alpha = 0$	$\beta > 0$	$\beta > 0.5$
XGBoost(d)	7.7(4.0)	7.8(4.1)	9.2(4.0)	6.8(3.5)	8.2(4.1)	8.3(3.9)	7.7(4.4)	7.7(4.6)	7.3(4.1)	7.8(4.0)	8.0(4.7)	9.2(4.3)
CatBoost(d)	6.7(4.1)	6.8(4.0)	7.4(4.0)	6.0(4.6)	7.0(4.1)	6.8(4.2)	5.5(2.7)	5.5(2.6)	5.5(2.7)	5.6(3.0)	5.5(2.7)	5.8(3.2)
FTT(d)	7.1(3.5)	7.0(3.5)	6.6(3.5)	6.9(3.6)	6.9(3.6)	7.2(3.6)	7.8(2.7)	7.8(2.5)	8.2(3.0)	7.6(3.2)	8.0(2.6)	8.3(1.3)
TransTab(d)	11.0(4.5)	11.2(4.5)	11.2(4.1)	10.2(4.6)	11.6(4.3)	11.7(4.2)	12.1(4.0)	12.1(3.8)	13.3(2.2)	12.4(4.5)	12.0(4.0)	13.6(1.2)
XGBoost(t)	6.2(4.1)	6.3(4.1)	6.5(4.3)	5.9(4.2)	6.5(4.2)	6.7(4.5)	4.5(3.7)	4.3(3.8)	<b>3.3(3.3)</b>	5.0(3.5)	4.7(3.9)	4.1(3.2)
CatBoost(t)	5.9(3.8)	6.3(3.9)	7.1(4.1)	<b>4.9(3.1)</b>	6.4(3.9)	6.4(4.1)	5.5(3.6)	5.7(3.6)	5.8(3.5)	4.9(3.7)	5.7(3.7)	6.1(3.8)
MLP(t)	8.6(4.0)	8.9(3.9)	8.7(4.1)	8.5(4.1)	8.5(3.9)	8.3(4.1)	8.5(3.6)	8.8(3.4)	9.3(3.2)	<u>7.6(4.1)</u>	9.0(3.4)	7.5(3.8)
AutoInt(t)	8.0(3.5)	7.8(3.3)	7.4(3.4)	8.6(4.0)	7.7(3.4)	7.7(3.2)	8.3(3.0)	8.6(3.0)	8.5(2.7)	7.4(3.1)	8.3(3.0)	8.2(3.2)
DCNv2(t)	7.9(3.9)	8.0(3.9)	8.4(3.8)	7.9(4.0)	7.7(3.9)	8.8(3.3)	8.4(3.4)	8.4(3.5)	8.5(3.1)	8.5(3.2)	8.4(3.5)	7.2(3.5)
TabNet(t)	12.1(3.5)	12.4(3.3)	12.7(2.7)	11.5(4.2)	12.3(3.4)	12.3(3.8)	12.6(3.6)	13.2(2.6)	13.1(2.4)	10.5(5.1)	13.5(1.9)	14.1(1.4)
SAINT(t)	8.2(3.8)	8.0(3.7)	8.1(4.1)	8.7(4.2)	7.9(3.8)	7.5(3.9)	7.6(3.8)	7.3(3.9)	7.7(3.3)	8.4(3.7)	6.6(3.6)	7.2(3.0)
FTT(t)	6.8(3.5)	6.8(3.6)	6.5(3.4)	6.2(3.3)	6.9(3.6)	6.9(3.9)	7.9(3.4)	7.6(3.3)	7.7(3.1)	9.0(3.4)	7.2(3.0)	6.8(3.2)
XTab(t)	9.8(4.0)	9.7(4.0)	8.9(3.8)	10.5(4.1)	9.4(4.0)	9.9(3.7)	12.4(2.8)	12.5(2.8)	13.3(1.6)	12.0(3.0)	12.4(2.9)	13.1(1.8)
Ours <sub>j</sub> (d)	8.4(4.5)	7.7(4.5)	7.0(5.0)	9.9(4.1)	7.9(4.6)	7.0(4.7)	6.9(4.6)	6.3(4.4)	4.8(3.9)	8.5(5.0)	6.5(4.5)	5.2(3.9)
Ours <sub>s</sub> (d)	<b>5.8(4.0)</b>	<b>5.1(3.9)</b>	<b>4.4(3.3)</b>	7.5(3.7)	<b>5.2(4.1)</b>	<b>4.5(3.4)</b>	<b>4.3(2.8)</b>	<b>4.1(2.6)</b>	<u>3.9(2.4)</u>	<b>4.8(3.4)</b>	<b>4.3(2.7)</b>	<b>3.6(2.8)</b>

starting point, and follow similar pre-training settings of RoBERTa (Liu et al., 2019): We use a total of 30 training epochs, with a linear warm-up for the first 6% of steps, followed by a linear decay to 0. The best checkpoint is saved by the average validation loss over all the datasets. We provide three TP-BERTa versions: pre-trained on only binary classification tasks, or only regression tasks, or both types. More detailed pre-training information and analysis are given in Appendix D.

**Compared Methods.** We compare our TP-BERTa with (1) the representative non-deep learning models XGBoost (Chen & Guestrin, 2016) and CatBoost (Prokhorenkova et al., 2018); (2) known DNNs including MLP, TabNet (Arik & Pfister, 2021), AutoInt (Song et al., 2019), DCNv2 (Wang et al., 2021), FT-Transformer (FTT) (Gorishniy et al., 2021), and SAINT (Somepalli et al., 2022); (3) the recent open-source cross-table model TransTab (Wang & Sun, 2022) and pre-trained model XTab (Zhu et al., 2023). We split each finetune dataset ((64%, 16%, 20%) for training, validation, and testing separately), and keep the same label distribution in each split on binary classification.

**Hyperparameter Tuning & Finetune.** We implement our TP-BERTa with PyTorch and the HuggingFace Transformers package on Python 3.8. All the models are finetuned on NVIDIA RTX 3090. In training, we uniformly use a training batch size of 64 for all the DNNs. Since the LM takes an increased training time, we directly set fixed hyperparameters on the pre-trained TP-BERTa across all the downstream datasets without tuning. For the other DNNs, the optimizer is AdamW (Loshchilov & Hutter, 2019) with the default configuration except for the learning rate and weight decay rate. We follow the hyperparameter spaces from the original work for SAINT. For TransTab, we use its default hyperparameters without cross-table pre-training because it originally required partially overlapped medical tables. For XTab, we follow its settings that report the score of the best pre-trained checkpoint on the validation set, with other hyperparameters kept fixed. For XGBoost, CatBoost, and the other DNNs, we follow the default (for GBDTs and FT-Transformer) and tuning settings provided in (Gorishniy et al., 2021). Hyperparameter search is performed with the Optuna library (Akiba et al., 2019). More detailed information of hyperparameters is provided in Appendix E.

### 3.2 ARE PRE-TRAINED TP-BERTA GREAT TABULAR PREDICTION LEARNERS?

**Overall Comparison.** Table 1 reports the means and standard deviations of model ranks on two dataset collections. As expected, a similar trend as shown in (Grinsztajn et al., 2022) is attained: GBDTs (i.e., XGBoost and CatBoost) still outperform classical and advanced DNNs in typical tabu-

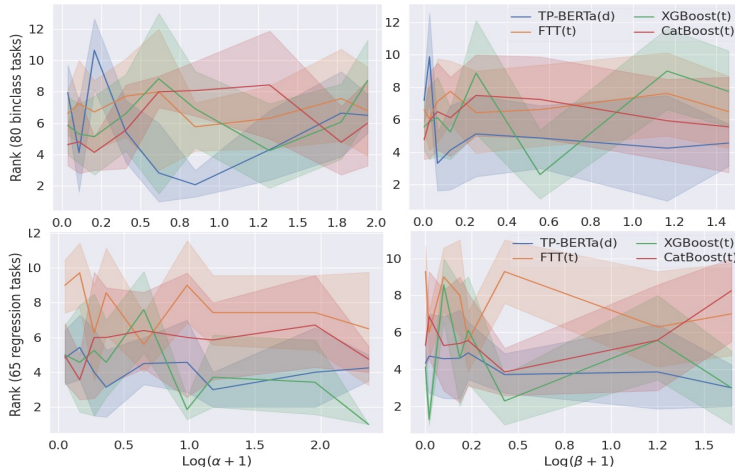


Figure 2: Rank variation curve plots of several representative models with respect to variations of some feature type characteristics. Each point represents a set of datasets in a range of  $\alpha$  or  $\beta$ .

lar regime (specified in “**Datasets**” of Sec. 3.1). Yet, it is worth noting that the pre-trained TP-BERTa exhibits a significantly different progress and competitive performances. This notable improvement may be attributed to the generalization ability of the pre-trained LMs (e.g., GPT-3 (Brown et al., 2020)). A medium-sized dataset (with  $< 10K$  points and low-dimensional features) may not have sufficient information for non-pre-trained DNNs, while LMs are able to leverage semantic information from feature names and structured values. Besides, our RMT approach further enables the LMs to handle numerical values in the language space (Sec. 3.3 discusses the necessity of RMT). Few previous deep tabular models were evaluated in such data settings, and this is the first time an extensive comparison on typical tabular data is brought to the forefront. As for cross-table models, TransTab was inspired by overlapped columns between pre-training datasets and downstream ones, which can benefit on domain datasets (e.g., medical tables), but general tables can contain many features from various domains, thus constraining its application. XTab adopted dataset-specific featurizers, though a Transformer backbone is shared; it misses the inherent relationship between features of different datasets and learns the feature embeddings from scratch, which may be trapped in insufficiently generalized data patterns. TP-BERTa is able to exploit feature semantics, e.g., the patterns learned on feature values “male & female” in pre-training can be inherently transferred to “boy & girl” by LMs without compulsory need for overlapped features or dataset-specific encoders.

**Comparison from the Feature Perspectives.** Since LMs typically operate on discrete texts, we further investigate from the perspective of feature type distributions. We report ranks among the datasets with various feature type distributions in Table 1. One can see that TP-BERTa achieves stably better performances (both “Ours<sub>j</sub>” and “Ours<sub>s</sub>”) when the categorical feature type gradually becomes dominating (a larger  $\alpha$  or  $\beta$ ). This can be intuitively explained by LMs’ ability to understand meaningful structured values (as discrete strings). Even among the datasets with at least one categorical feature ( $\alpha > 0$ ), TP-BERTa still leads the performances on both task types (also shown in Fig. 2). However, if all features are numerical ( $\alpha = 0$ ), TP-BERTa performs inferiorly. This may be due to its LM nature that precision loss in numerical representation is inevitable. For a more detailed illustration, we show in Fig. 2 rank variation curve plots across datasets grouped in different ranges of feature type distributions. Overall, TP-BERTa is stably promising when discrete features begin to dominate in the datasets, while for purely numerical datasets, GBDTs or FTT are still better choices (especially for classification tasks). Since there can exist useless tabular features, we introduce the ratio of Shapley value sums between categorical and numerical features (i.e.,  $\beta$ , the right column of Fig. 2); an expected smoother trend that TP-BERTa performs better on datasets with larger  $\beta$  is observed. Additionally, we empirically observe: (1) XGBoost highly relies on hyperparameter tuning and thus performs unstably (shown in its standard deviations and Fig. 2); (2) in contrast, CatBoost, just using default hyperparameters, is often a good choice, especially on datasets in which categorical features dominate (the same was suggested in (Prokhorenkova et al., 2018)).

**Comparison from the Data Volume Perspective.** We show rank variations on data volumes in the Appendices (Fig. 5). Similar to studying the feature type distribution effects, we examine the trend in two dataset groups ( $\beta < 0.1$  and  $\beta \geq 0.1$ ). In the typical tabular regime, the choices are mostly influenced by the distributions, while from the data scale dimension, no special trend is observed.

Table 2: Performance changes on encoding strategy substitution and IFA ablation using 80 binary classification datasets. The column “ $|\Delta| \leq 0.5\%$ ” denotes the number of datasets with AUC variation less than 0.5% (these datasets are called “insignificantly changed datasets” due to different random seeds); the other “ $\Delta$ ” columns use similar denotations. “Avg. diff.” means the average performance difference on significantly changed datasets. “Avg. training time ratio” is the average ratio of training time compared to using the IFA module. Appendix 11 gives more detailed performances.

Comparison (numerical encoding strategies)				
Substitution	$ \Delta  \leq 0.5\%$	$\Delta < -0.5\%$	$\Delta > 0.5\%$	Avg. diff.
Value2Str (Borisov et al., 2023)	16	54	10	-12.45%
VMFE (Ye et al., 2024)	34	36	10	-3.44%
Ablation (w/o IFA module)				
Avg. training time ratio	$ \Delta  \leq 0.5\%$	$\Delta < -0.5\%$	$\Delta > 0.5\%$	Avg. diff.
1.32	14	52	14	-4.17%

**Joint Task Type Pre-training.** A more expensive TP-BERTa version is conducted on “Ours<sub>j</sub>” by pre-training on binary classification and regression tasks jointly. A similar trend as “Ours<sub>s</sub>” is observed, with a stably inferior performance. This may be due to: (1) incompatible natures of classification and regression tasks, (2) the dataset-specific head pre-training strategy is unsuitable for the combined patterns of classification and regression, or (3) a more powerful base LM is needed for such complicated data configuration. This will become a part of our future study.

**Takeaway.** We empirically show a strong potential for well-adapted LMs to fill the void of previous DNN-based tabular learners under typical tabular data settings, and demonstrate the capability of TP-BERTa to handle tabular prediction tasks, especially those with informative categorical features, which can help architecture selection based on feature type characteristics in future studies.

### 3.3 WHY WERE LMS NEGLECTED ON TABULAR PREDICTION?

Only a few previous studies directly employed LM-based learners for tabular prediction tasks. Their numerical encoding strategies can be categorized into: (1) value string methods (directly treating numerical values as strings, e.g., GReaT (Borisov et al., 2023) and TapTap (Zhang et al., 2023)); (2) value-multiplied feature name embeddings (e.g., CT-BERT (Ye et al., 2024), FTT (Gorishniy et al., 2021), and TransTab (Wang & Sun, 2022)). Besides, they all fed a longer templated table text to LMs, which may further incur a heavy training burden. In this section, we compare our RMT approach with two other strategies, and conduct ablation study on the intra-feature attention (IFA) module to demonstrate that refining single-feature information before the LM processing is a better and computation-friendly adaption. Since each pre-training round is costly and is equivalent to evaluation under the same condition, we use the non-pre-trained TP-BERTa (initialized with the RoBERTa weights) for the subsequent comparison and ablation studies (marked in the Appendix tables). To show a more clear and quantifiable comparison, we conduct our analysis on the binary classification datasets (Table 2). To alleviate the impact of performance fluctuations caused by random seeds, we exclude the datasets on which the performance changes are insignificant (the column “ $|\Delta| \leq 0.5\%$ ”) in average difference calculation (the column “Avg. diff.”). The following sections use a similar analysis method. We present the detailed results in the Appendices (Table 11).

**Numerical Encoding Strategy Comparison.** After directly substituting RMT with “value string” (Value2Str) or “value-multiplied feature name embeddings” (VMFE), changes in performance are observed (the upper half of Table 2). Both the previous strategies hurt AUC scores on most significantly changed datasets with average declines of 12.45% and 3.44%, respectively. There are still 10 datasets with better performances on both substitutions. This is probably due to insufficient embedding learning: Since in non-pre-trained TP-BERTa, magnitude embeddings are randomly initialized, direct finetune on downstream tasks faces a risk of data inadequacy to learn precise representations.

**IFA Module Ablation.** Performances and training time changes are reported in the lower half of Table 2, by removing IFA and directly feeding all feature names and values to the LM as done in previous works. A noticeable performance degradation occurs on 52 datasets ( $\Delta < -0.5\%$ ) with an average AUC decline of 4.17% (on  $52 + 14 = 66$  datasets). This indicates that LMs are likely to be confused when they process a pile of unmatched feature name-value texts, giving



Table 3: Performance changes by comparing the pre-trained TP-BERTa with (1) TP-BERTa randomly initialized and (2) TP-BERTa initialized with the RoBERTa weights. ‘‘Avg. diff.’’ is calculated by excluding the datasets with  $|\Delta| \leq 0.5\%$ .

Comparison (w/ no pre-training) using 80 binary classification datasets						
Initialization	$ \Delta  \leq 0.5\%$	$\Delta < -0.5\%$	$\Delta > 0.5\%$	$\Delta < -3\%$	$\Delta > 3\%$	Avg. diff.
Random	29	41	10	26	5	-3.16%
RoBERTa	26	35	19	21	6	-2.79%

them additional burden to learn correct matchings while fully connected attention in Transformer-based LMs interacts a name with values from other features. The IFA module explicitly fuses the name-value pair of a single feature into a vector before passing it to the LM, guarding against noisy interactions from other features. Besides, a shorter input sequence length (equal to the feature amount) accelerates learning (e.g., see the 1.32 average training time ratio without IFA in Table 2).

Since TP-BERTa adopts both magnitude-aware numerical encoding and intra-feature pre-processing before the LM process, it acquires a significantly better ability as well as friendly computation cost, becoming competitive with GBDTs and other deep models. Our comparison shows that simply treating tables as normal texts can pose difficulties for LMs to understand structured tabular features, thus decreasing their potential on high-precision demanding tasks such as tabular predictions.

### 3.4 TP-BERTa TRANSFERABILITY ON TABULAR DATA

Table 3 reports performance changes by comparing the non-pre-trained TP-BERTa (initialized by random weights or RoBERTa weights) with the pre-trained one (‘‘Ours<sub>s</sub>’’ in Table 1). Overall, over 3% AUC increase is attained on 26 (comparing to random weights) and 21 (comparing to RoBERTa weights) datasets using pre-training, and the average improvement on significantly changed datasets is 3.16% and 2.79%, respectively. It seems that using the RoBERTa weights is better than random weights, as LM weights have inherently entailed meaningful semantic knowledge. A more significant leap can be achieved by further pre-training on extensive tabular data. This indicates that LMs are also effective in transferring tabular data knowledge and suitable for cross-table pre-training.

### 3.5 THE NECESSITY OF OTHER DESIGN DETAILS

Since there are several key differences in our TP-BERTa design compared to the common Transformer-based LMs, we further examine their necessity by ablation studies. By evaluating different **magnitude token numbers** ( $n_{\text{bin}} = 256$  as default), we find that using 128 and 32 tokens yields an average AUC decline of 2.06% and 3.59%, respectively. This indicates that a larger  $n_{\text{bin}}$  for numerical value representation benefits performances on most datasets, while a few tables favor a smaller  $n_{\text{bin}}$ , which may be due to over-representation of excessive magnitude tokens. The detailed results with analysis are presented in Appendix C, which further discusses the regularization efforts on magnitude embeddings of the **magnitude-aware triplet loss function** (Eq. (3)) and the function of **removing position encoding for value vectors** (Eq. (7)). We find that the magnitude-aware triplet loss function potentially facilitates fast convergence and over-fitting reduction.

## 4 CONCLUSIONS AND FUTURE WORK

This paper undertook the first study of the substantial difficulties of continuous value representation and tabular feature organization in building LM-based tabular DNNs. We designed and deployed two bespoke adaptations, relative magnitude tokenization and intra-feature attention, to explore the possibilities of using pre-trained LMs on tabular prediction tasks. Our proposed TP-BERTa exhibits unprecedented progress over various non-LM DNNs, and is competitive with GBDTs under the typical tabular prediction regime, contributing a powerful DNN alternative for typical tabular data.

While our approach has significantly improved the performance of language models in handling numerical features in tables, TP-BERTa currently excels more on tables dominated by categorical features. Besides, it was witnessed that some tables prefer a small magnitude token number. We will conduct more endeavors on better numerical representation in future tabular prediction studies.

## ACKNOWLEDGMENTS

The research of Jiahuan Yan, Hongxia Xu and Jian Wu was supported in part by the National Natural Science Foundation of China under grants No. 62176231 and No. 82202984, and the Zhejiang Key R&D Program of China under grant No. 2023C03053. The research of Jintai Chen and Jimeng Sun was supported by NSF award SCH-2205289, SCH-2014438, and IIS-2034479.

## REFERENCES

- Takuya Akiba, Shotaro Sano, et al. Optuna: A next-generation hyperparameter optimization framework. In *KDD*, 2019.
- Sercan Ö Arik and Tomas Pfister. TabNet: Attentive interpretable tabular learning. In *AAAI*, 2021.
- Saqib Aziz, Michael Dowling, Helmi Hammami, and Anke Piepenbrink. Machine learning in finance: A topic modeling approach. *European Financial Management*, 2022.
- Vadim Borisov, Tobias Leemann, Kathrin Sessler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. Deep neural networks and tabular data: A survey. *TNNLS*, 2022.
- Vadim Borisov, Kathrin Sessler, Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Language models are realistic tabular data generators. In *ICLR*, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020.
- Jintai Chen, Kuanlun Liao, et al. DANets: Deep abstract networks for tabular data classification and regression. In *AAAI*, 2022.
- Jintai Chen, KuanLun Liao, Yanwen Fang, Danny Z. Chen, and Jian Wu. TabCaps: A capsule neural network for tabular data classification with BoW routing. In *ICLR*, 2023a.
- Jintai Chen, Jiahuan Yan, Danny Ziyi Chen, and Jian Wu. Excelformer: A neural network surpassing GBDTs on tabular data. *arXiv preprint arXiv:2301.02819*, 2023b.
- Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In *KDD*, 2016.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional Transformers for language understanding. In *NAACL*, 2018.
- James Dougherty, Ron Kohavi, and Mehran Sahami. Supervised and unsupervised discretization of continuous features. In *ICML*, 1995.
- Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. In *ACL*, pp. 3816–3830, 2021.
- Yury Gorishniy, Ivan Rubachev, et al. Revisiting deep learning models for tabular data. In *NeurIPS*, 2021.
- Yury Gorishniy, Ivan Rubachev, and Artem Babenko. On embeddings for numerical features in tabular deep learning. In *NeurIPS*, 2022.
- Leo Grinsztajn, Edouard Oyallon, and Gael Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? In *NeurIPS*, 2022.
- Md Rafiul Hassan, Sadiq Al-Insaf, et al. A machine learning approach for prediction of pregnancy outcome following IVF treatment. *Neural Computing and Applications*, 2020.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. TabPFN: A transformer that solves small tabular classification problems in a second. In *ICLR*, 2023.

- Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *TACL*, 8:423–438, 2020.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. In *NeurIPS*, 2017.
- Ron Kohavi and Mehran Sahami. Error-based and entropy-based discretization of continuous features. In *KDD*, 1996.
- Liyao Li, Haobo Wang, Liangyu Zha, Qingyi Huang, Sai Wu, Gang Chen, and Junbo Zhao. Learning a data-driven policy network for pre-training automated feature engineering. In *ICLR*, 2023.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pre-training approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? In *EMNLP*, pp. 2463–2473, 2019.
- Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision ensembles for deep learning on tabular data. In *ICLR*, 2020.
- Liudmila Prokhorenkova, Gleb Gusev, et al. CatBoost: Unbiased boosting with categorical features. In *NeurIPS*, 2018.
- Jing Qian, Hong Wang, Zekun Li, Shiyang Li, and Xifeng Yan. Limitations of language models in arithmetic and symbolic induction. In *ACL*, 2023.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI Blog*, 2019.
- Gowthami Somepalli, Avi Schwarzschild, Micah Goldblum, C Bayan Bruss, and Tom Goldstein. SAINT: Improved neural networks for tabular data via row attention and contrastive pre-training. In *NeurIPS*, 2022.
- Weiping Song et al. AutoInt: Automatic feature interaction learning via self-attentive neural networks. In *CIKM*, 2019.
- Ruoxi Wang, Rakesh Shivanna, et al. DCN V2: Improved deep & cross network and practical lessons for web-scale learning to rank systems. In *WWW*, 2021.
- Zifeng Wang and Jimeng Sun. TransTab: Learning transferable tabular Transformers across tables. In *NeurIPS*, 2022.
- Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves ImageNet classification. In *CVPR*, 2020.
- Jiahuan Yan, Jintai Chen, Yixuan Wu, Danny Z Chen, and Jian Wu. T2G-Former: Organizing tabular features into relation graphs promotes heterogeneous feature interaction. In *AAAI*, 2023.
- Chao Ye, Guoshan Lu, Haobo Wang, Liyao Li, Sai Wu, Gang Chen, and Junbo Zhao. Towards cross-table masked pretraining for web data mining. In *WWW*, 2024.
- Liangyu Zha, Junlin Zhou, Liyao Li, et al. Tablegpt: Towards unifying tables, nature language and commands into one gpt. *arXiv preprint arXiv:2307.08674*, 2023.
- Tianping Zhang, Shaowen Wang, Shuicheng Yan, Jian Li, and Qian Liu. Generative table pre-training empowers models for tabular prediction. *arXiv preprint arXiv:2305.09696*, 2023.
- Bingzhao Zhu, Xingjian Shi, Nick Erickson, Mu Li, George Karypis, and Mahsa Shoaran. XTab: Cross-table pretraining for tabular Transformers. In *ICML*, 2023.

## A LIMITATIONS

Since our TP-BERTa relies on the semantic knowledge of LMs to transfer feature patterns from pre-training feature names to downstream ones, it implicitly requires meaningful and clear feature semantics. However, in the real world, there always exist tables with unclear feature names or values, such as quantum physics experiment tables containing lots of uncommon quantum descriptor feature names and medical domain substituting specific feature values with meaningless encoding to protect patient privacy. This suggests that LM-based tabular DNNs cannot take their inherent advantages of feature semantic understanding in privacy-sensitive (e.g., federated learning) or semantic-incomplete (missing original meanings in data collection) tables. For an uncommon domain, LMs pre-trained on domain corpora may be utilized as base models. Besides, LMs own a larger space of parameters and hyperparameters, making them more time-consuming in hyperparameter tuning with a potentially higher performance ceiling compared to non-LM tabular DNNs. Hence, we directly finetune TP-BERTa with default hyperparameters for fairness of time, in which case it can achieve adequate results with less time than other tuned methods.

## B DATASET INFORMATION AND MAIN EXPERIMENTAL DETAILS

We provide detailed information of the pre-training datasets in Table 6 and detailed information of the downstream datasets in Table 7 and Table 9. Detailed baseline performances are given in Table 8 and Table 10. To represent distributions of feature types in a dataset, we define  $\alpha$  as the feature amount ratio between the categorical type and numerical type. Since there exist datasets with usefulness features, we further define  $\beta$  as the Shapley value (calculated with default XGBoost in Appendix E) sum ratio between categorical features and numerical features. These are used as references of feature type characteristics. Specifically,  $\alpha$  and  $\beta$  of the  $i$ -th dataset are formulated as:

$$\alpha_i = \frac{\#Cat.^{(i)}}{\#Num.^{(i)}}, \quad (12)$$

$$\beta_i = \frac{\sum_{f \in Cat.^{(i)}} \text{Shapley value}(f)}{\sum_{f \in Num.^{(i)}} \text{Shapley value}(f)}. \quad (13)$$

We exclude datasets with pure categorical features to avoid zero division, while as an LM, TP-BERTa can inherently handle discrete features. We provide the dataset frequency in different feature type distribution ranges in Table 4.

Table 4: Dataset frequency of two downstream collections in several feature type distribution ranges.

Collection	$\beta = 0$	$\beta \in (0, 0.5)$	$\beta \in [0.5, 1.0)$	$\beta \geq 1.0$
80 binary classification datasets	24	33	7	16
65 regression datasets	24	28	4	9

## C RESULTS AND ANALYSIS OF OTHER DESIGN COMPARISONS

**The Number of Magnitude Tokens.** In Sec. 2.1, we set the maximum number of magnitude tokens,  $n_{\text{bin}} = 256$ , for TP-BERTa. In fact, the C4.5 decision tree splits the value range in a greedy fashion, and the actual number of leaves can be less than 256 (e.g., a small dataset with less than 256 points). Hence, we choose a conservative method to balance between “not too many new words to learn” and “enough precision for relative magnitude representation”. In Table 5, we present the performance changes by setting  $n_{\text{bin}}$  to 32 and 128 separately. As expected, the overall performance gradually drops when using a smaller token number to represent the numerical value magnitude. We find that some datasets favor a smaller  $n_{\text{bin}}$ , which may be attributed to over-representation of too many magnitude tokens. Thus, a better solution is to set a large  $n_{\text{bin}}$  for pre-training in order to enhance the upper limit of magnitude representation capability, and search for a reasonable dataset-specific  $n_{\text{bin}}$  on downstream tasks.

**Regularization on Magnitude Embeddings.** Since all the magnitude embeddings are randomly initialized, in Sec. 2.1, we propose a magnitude-aware triplet loss (see Eq. (3)) to assist the learning process. Fig. 3 presents the validation AUC curves of several datasets on using our regularization or not using it during finetuning, which shows that the designed regularization provides a potential of fast convergence and overfitting reduction. We use the triplet loss only in pre-training for a smoother and accelerated learning, and exclude it in actual finetune because the loss has converged in pre-training.

**No Position Encoding for Value Vectors.** In Eq. (7), we explicitly remove position encoding in value vectors of self-attention since position embeddings may distort randomly initialized magnitude embedding learning. Table 5 shows a major performance decline when adding position encoding. The reason for this probably lies in semantic corruption of magnitude tokens, since numerical values need precise representations to convey magnitude information.

Table 5: Performance changes with respect to using different magnitude token numbers (the default is  $n_{\text{bin}} = 256$ ; see Sec. 2.1) and position encoding in value vectors (see Eq. (7)).

Comparison (magnitude token numbers) using 80 binary classification datasets				
Substitution	$ \Delta  \leq 0.5\%$	$\Delta < -0.5\%$	$\Delta > 0.5\%$	Avg. diff.
$n_{\text{bin}} = 32$	27	40	13	-3.59%
$n_{\text{bin}} = 128$	42	26	12	-2.06%
Ablation (w/ value vector position encoding)				
80 binary classification datasets	$ \Delta  \leq 0.5\%$	$\Delta < -0.5\%$	$\Delta > 0.5\%$	Avg. diff.
	36	31	13	-2.35%

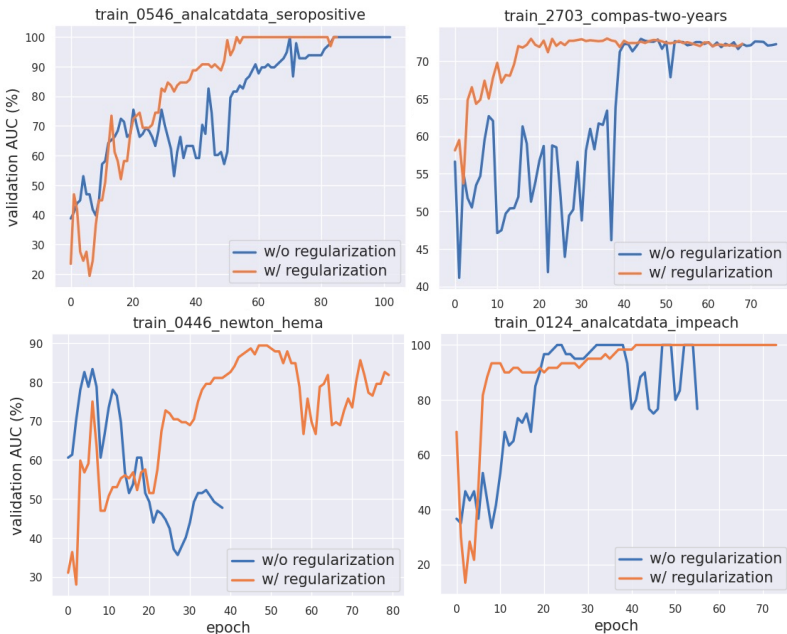


Figure 3: Comparison of using regularization or not using it during finetuning on the non-pre-trained TP-BERTa. The validation AUC curves of several representative binary classification datasets show that the effect of the magnitude-aware triplet loss (see Eq. (3)) is to help quick convergence and avoid potential overfitting of TP-BERTa. In experiments, we use this regularization only in pre-training to smooth and accelerate the learning process.

## D PRE-TRAINING DETAILS

**Starting Point.** We reuse the weights of the RoBERTa-base with the HuggingFace Transformers API. The additional  $n_{\text{bin}}$  magnitude embeddings and the IFA module are randomly initialized.

**Runtime Environment.** Pre-training is conducted with PyTorch version 1.9.0, CUDA version 11.3, and HuggingFace Transformers package version 4.18.0, using 4 NVIDIA A100 PCIe 40GB and 2 Intel 6248R 48C@3.0GHz.

**Pre-training Process.** We pre-train three versions of TP-BERTa: pre-training only on binary classification datasets, or only on regression datasets (these two versions constitute “Ours<sub>s</sub>” in Table 1), or jointly pre-training on both types (“Ours<sub>j</sub>” in Table 1). These three versions share the same maximum epoch number of 30 and the total batch size of 512, using the best average validation loss across all the datasets to save the checkpoint. The same pre-training learning rate and linear decay in (Liu et al., 2019) are used. The pre-training on binary classification tasks and regression tasks took 72 hours and 98 hours respectively, and the one on both types of tasks took 143 hours. We provide several loss curves during pre-training in Fig. 6 and validation metric curves on several pre-training datasets in Fig. 7.

**Analysis.** In most cases, the TP-BERTa version pre-trained jointly on both task types yields slightly lower validation scores on the pre-training datasets compared to the two TP-BERTa versions pre-trained on a single task type (see Fig. 7), and the gap is more noticeable on binary classification datasets. This probably leads to a smaller overall rank difference between Ours<sub>j</sub> and Ours<sub>s</sub> on regression tasks (see Table 1).

## E HYPERPARAMETER TUNING

For the baselines of XGBoost, CatBoost, MLP, AutoInt, DCNv2, TabNet, and FT-Transformer, we reuse the implementations, default settings, and hyperparameter search spaces in (Gorishniy et al., 2021). For SAINT, the hyperparameter space in (Somepalli et al., 2022) is used. For TransTab and XTab, we follow the same settings as in (Zhu et al., 2023), using the default hyperparameters in (Wang & Sun, 2022) for TransTab and the best checkpoint on the validation set for XTab. As for TP-BERTa (including the joint pre-trained version “Ours<sub>j</sub>” and the single pre-trained one “Ours<sub>s</sub>”), we keep the default hyperparameters of 1e-5 learning rate without weight decay. All the baselines use AdamW (Loshchilov & Hutter, 2019) as the optimizer and the Optuna-driven tuning.

## F INTERPRETABILITY OF RMT

In Fig. 4, we visualize the TP-BERTa’s 256 magnitude tokens by directly applying t-SNE algorithm using scikit-learn package (with default function parameters). Interestingly, even if we have randomly placed them in the language space at first, after pre-training a clear inherent distribution is captured, all tokens lie on a highly regular manifold and successfully maintain an intuitive assumption that the embedding of a numerical value should close to the embedding of a nearby one. This empirically demonstrates the TP-BERTa is sensitive to the numerical value magnitudes and benefits from the captured regular relationship among the relative magnitudes, which interprets its significant progress on the existing LM-based tabular models (e.g., directly treating values as raw strings).

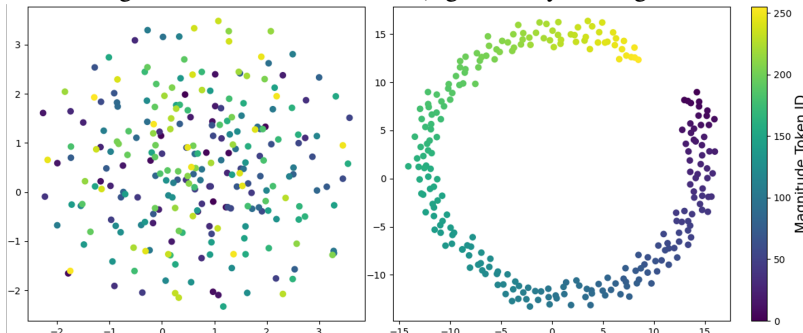


Figure 4: The t-SNE visualization of 256 magnitude token embeddings before and after pre-training.

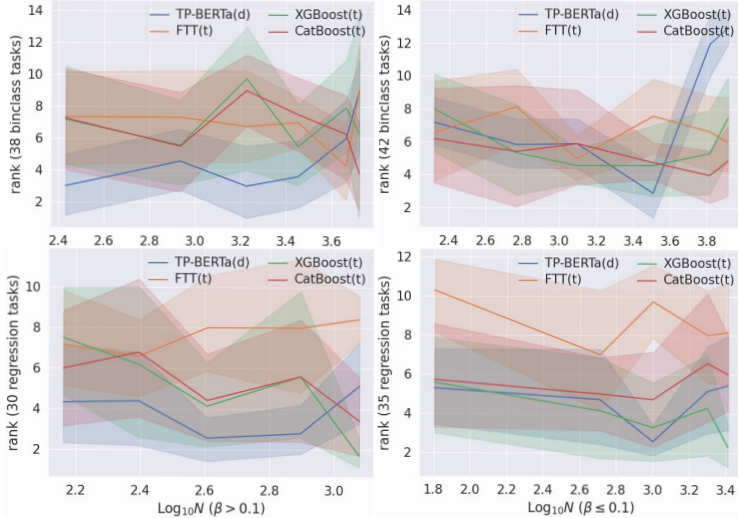


Figure 5: Rank change curve plots of several representative models with variations of data volume ( $N$ ). We divide the datasets into two groups (the first column is for “ $\beta > 0.1$ ” and the second column is for “ $\beta \leq 0.1$ ”) to alleviate the impact from the feature type distributions. The split value 0.1 is chosen by keeping a roughly equal number of datasets in both groups.

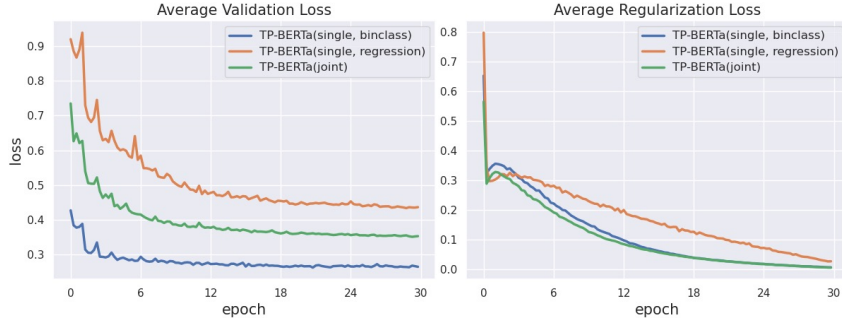


Figure 6: The average validation loss curves and regularization loss (Eq. (3)) curves in pre-training. “TP-BERTa(single, binclass)” and “TP-BERTa(single, regression)” are the two versions separately pre-trained on binary classification datasets and regression ones (constituting “Ours<sub>s</sub>” in Table 1); “TP-BERTa(joint)” is the version pre-trained on both task types (“Ours<sub>j</sub>” in Table 1).

## G EVALUATIONS ON OTHER DATA SCENARIOS

### G.1 IMBALANCED LABEL DISTRIBUTION

To inspect the performances on imbalanced datasets, we create pivot tables by filtering datasets whose minor-class proportions, i.e.,  $p = \min(\#positive, \#negative) / \#sample$ , are less than  $1/3$  (32 datasets),  $1/5$  (18 datasets),  $1/8$  (12 datasets),  $1/20$  (4 datasets) from 80 binary classification datasets. The results are reported in the upper part of Table 15. It is obvious that TP-BERTa outperforms baselines in moderate class-imbalance situations. In the extremely imbalanced situations (4 datasets, with  $p < 0.05$ ), GBDTs showcase dominating performances, but TP-BERTa still outperforms most DNN approaches. Perhaps this is attributed to TP-BERTa leveraging the transferability and semantic understanding capabilities of the language model, thus resulting in consistent performance across various levels of data imbalance.

### G.2 MULTI-CLASS CLASSIFICATION

We additionally experiment on 32 downstream multi-class datasets from the database, the data statistics and results are reported in Table 12 and the middle part of Table 15 respectively, and the TP-BERTa used here is the version pre-trained on 101 binary classification datasets.

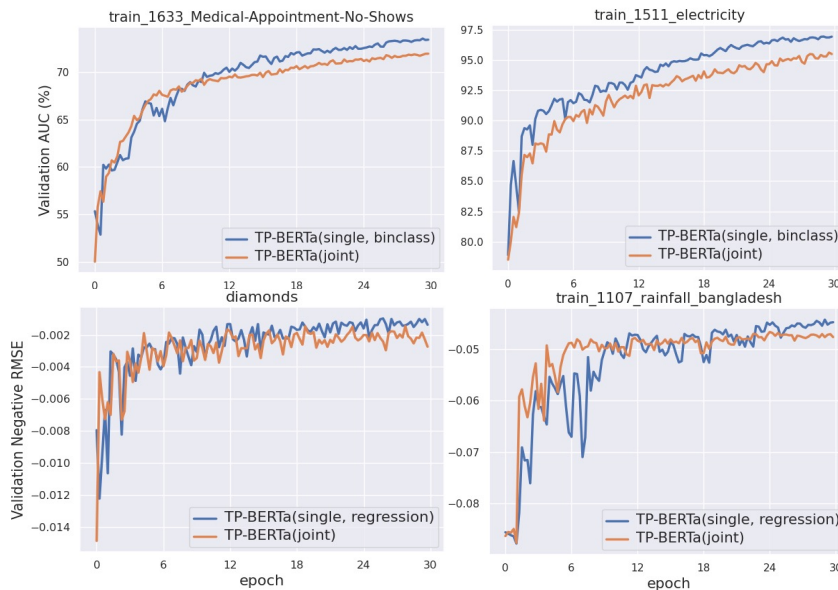


Figure 7: Validation score curves in AUC and RMSE on several pre-training datasets.

### G.3 MEDICAL APPLICATIONS

Medical data and labels hold inherently greater value than those from other domains. Therefore, we further inspect whether our pre-trained model can achieve notable performance on medical domain tasks, such as patient risk prediction and clinical trial outcome prediction, by leveraging knowledge learned and transferred from diverse domains, which typically offer more cost-effective data sources. We identified and filtered out all 25 medical tasks from 145 main experiment downstream datasets (statistics and results are given in Table 13 and Table 14). Refer to the bottom row of Table 14, the results indicate that the pre-trained TP-BERTa model significantly outperforms Gradient Boosting Decision Trees (GBDTs) and other Deep Neural Networks (DNNs) that are trained in the supervision manner and undergo meticulous hyperparameter tuning. However, our TP-BERTa achieved this without the need for hyperparameter tuning.

The noteworthy performance suggests that we can harness more affordable data sources to alleviate reliance on costly healthcare data in clinical practice. This illuminates a promising pathway towards reducing the need for extensive medical data and expediting the development of algorithms for healthcare applications.





Table 7: Statistics of 80 downstream binary classification datasets.

ID	Dataset name	# samples	# num.	# cat.	$\alpha$	Num. Shapely sum	Cat. Shapely sum	$\beta$
0	BankNoteAuthentication	1372	4	0	0.00	9.21	0.00	0.00
1	bt_dataset_t3	1644	17	0	0.00	6.05	0.00	0.00
2	0292_cpu_small	8192	12	0	0.00	7.61	0.00	0.00
3	0312_cpu_act	8192	21	0	0.00	7.39	0.00	0.00
4	0345_delta_aileron	7129	5	0	0.00	3.98	0.00	0.00
5	0356_delta_elevators	9517	6	0	0.00	3.83	0.00	0.00
6	0406_visualizing	111	3	0	0.00	0.69	0.00	0.00
7	0419_pm10	500	7	0	0.00	1.06	0.00	0.00
8	0435_strikes	625	6	0	0.00	4.73	0.00	0.00
9	0437_quake	2178	3	0	0.00	0.33	0.00	0.00
10	0445_arsenic-male	559	4	0	0.00	1.35	0.00	0.00
11	0446_arsenic-female	559	4	0	0.00	1.28	0.00	0.00
12	0447_arsenic-female	559	4	0	0.00	1.70	0.00	0.00
13	0509_pollen	3848	5	0	0.00	0.11	0.00	0.00
14	1201_Gender	3168	20	0	0.00	6.34	0.00	0.00
15	1600_VulNoneVul	5692	16	0	0.00	0.91	0.00	0.00
16	1006_Titanic	2201	3	0	0.00	1.57	0.00	0.00
17	1592_Diabetes-Data	768	8	0	0.00	1.69	0.00	0.00
18	0424_autoPrice	159	14	1	0.07	2.99	0.00	0.00
19	audit_data	776	23	3	0.13	5.29	0.00	0.00
20	audit_risk	776	23	3	0.13	5.29	0.00	0.00
21	new_model	400	10	3	0.30	5.83	0.00	0.00
22	trial	776	7	10	0.39	5.29	0.00	0.00
23	1333_ricci_vs	118	3	2	0.67	3.66	0.00	0.00
24	1619_NBA-2k20-player	439	4	10	2.50	1.51	0.01	0.01
25	1458_kdd_ipums_la_97	5188	17	3	0.18	4.73	0.03	0.01
26	1736_combined-wine	6497	10	2	0.20	11.77	0.09	0.01
27	1759_Red-White-wine	6497	10	2	0.20	11.77	0.09	0.01
28	1578_kdd_ipums_la_97	5188	17	3	0.18	5.10	0.05	0.01
29	0526_colleges_aaup	1161	13	2	0.15	8.10	0.09	0.01
30	1408_national	4908	6	10	1.67	7.25	0.10	0.01
31	0284_bank8FM	8192	7	1	0.14	6.78	0.15	0.02
32	Customer_Behaviour	400	2	1	0.50	2.94	0.15	0.05
33	2306_electricity_seed	2000	7	1	0.14	3.94	0.22	0.06
34	0546_analcatdata	132	2	1	0.50	3.81	0.22	0.06
35	2304_electricity_seed	2000	7	1	0.14	3.25	0.23	0.07
36	2308_electricity_seed	2000	7	1	0.14	4.37	0.32	0.07
37	1461_heart-failure	299	7	5	0.71	4.47	0.37	0.08
38	2392_airlines_seed_3	2000	3	4	1.33	0.45	0.04	0.09
39	2390_airlines_seed_1	2000	3	4	1.33	0.32	0.03	0.09
40	0124_analcatdata	100	2	7	3.50	3.11	0.27	0.09
41	2305_electricity_seed	2000	7	1	0.14	4.08	0.37	0.09
42	2389_airlines_seed_0	2000	3	4	1.33	0.34	0.04	0.11
43	2393_airlines_seed_4	2000	3	4	1.33	0.44	0.05	0.12
44	0541_plasma_retinol	315	10	3	0.30	2.62	0.37	0.14
45	NFL	3477	10	7	0.70	1.69	0.25	0.15
46	1142_Sick_numeric	3772	6	23	3.83	4.07	0.68	0.17
47	2703_compas-two	4966	2	9	4.50	0.86	0.16	0.19
48	0885_compas-two	5278	2	11	5.50	1.00	0.19	0.19
49	2391_airlines_seed_2	2000	3	4	1.33	0.20	0.04	0.19
50	0400_analcatdata	4052	1	6	6.00	4.62	1.01	0.22
51	2621_sf-police	2000	4	4	1.00	0.49	0.11	0.23
52	b_depressed	1429	12	9	0.75	0.89	0.25	0.28
53	2619_sf-police	2000	4	4	1.00	0.83	0.23	0.28
54	2620_sf-police	2000	4	4	1.00	0.99	0.34	0.34
55	Breast_Cancer	4024	5	10	2.00	2.15	0.81	0.37
56	1512_eye_movements	7608	17	5	0.29	1.94	0.74	0.38
57	0472_analcatdata	364	23	9	0.39	1.21	0.68	0.56
58	0555_socmob	1156	1	4	4.00	2.69	1.51	0.56
59	2622_sf-police	2000	4	4	1.00	0.55	0.36	0.66
60	loan_train	614	4	7	1.75	0.34	0.22	0.67
61	TravelInsurancePrediction	1987	1	7	7.00	0.78	0.56	0.72
62	1742_Loan	614	4	7	1.75	0.60	0.48	0.79
63	1635_Is-this-a-good	1723	4	9	2.25	0.62	0.56	0.89
64	0948_Ishwar	2205	17	4	0.24	3.34	3.70	1.11
65	piracydataset	1423	1	3	3.00	0.18	0.27	1.47
66	1752_Wisconsin	699	2	8	4.00	1.85	3.06	1.66
67	1413_shill-bidding	6321	7	4	0.57	1.59	3.58	2.25
68	Employee	500	1	11	11.00	0.16	0.44	2.81
69	0446_newton_hema	140	2	1	0.50	0.34	1.00	2.92
70	Bank_Personal_Loan	5000	6	6	1.00	0.23	0.68	2.93
71	UniversalBank	5000	6	6	1.00	0.23	0.68	2.93
72	1011_cleve	303	5	8	1.60	0.84	2.76	3.29
73	1898_Personal-Loan	5000	6	6	1.00	0.19	0.72	3.85
74	1564_Mammographic	830	1	4	4.00	0.48	2.22	4.65
75	1692_Gender	5001	1	6	6.00	0.55	6.31	11.57
76	diabetes_data_upload	520	1	15	15.00	0.24	5.48	22.64
77	1451_early-stage	520	1	15	15.00	0.22	5.75	26.47
78	1774_Early-Stage	520	1	15	15.00	0.22	5.75	26.47
79	0408_pharynx	195	2	8	4.00	0.04	1.32	33.20



Table 9: Statistics of the 65 downstream regression datasets.

ID	Dataset name	# samples	# num.	# cat.	$\alpha$	Num. Shapely sum	Cat. Shapely sum	$\beta$
0	my_csv-3-10-2022-10-35	10	7	1	0.14	0.95	0.00	0.00
1	0237_arsenic-female	559	4	0	0.00	0.20	0.00	0.00
2	0251_arsenic-male	559	4	0	0.00	0.19	0.00	0.00
3	0258_no2	500	7	0	0.00	1.46	0.00	0.00
4	0259_strikes	625	6	0	0.00	0.34	0.00	0.00
5	ph-data	653	3	0	0.00	1.07	0.00	0.00
6	0925_Concrete_Data	1030	8	0	0.00	1.57	0.00	0.00
7	1065_hungarian	522	19	0	0.00	1.21	0.00	0.00
8	1260_optical	640	5	4	0.80	0.23	0.00	0.00
9	1331_dataset_time_9	2178	3	0	0.00	0.17	0.00	0.00
10	1464_dow-jones-index	750	6	9	1.50	0.63	0.00	0.00
11	1564_Game Concrete	1005	8	0	0.00	1.56	0.00	0.00
12	1623_GameStop	4773	5	0	0.00	0.90	0.00	0.00
13	1624_Alcohol	1549	5	1	0.20	1.07	0.00	0.00
14	1769_Facebook	2076	5	0	0.00	2.07	0.00	0.00
15	1845_Predict-Amazon	349	6	0	0.00	0.48	0.00	0.00
16	1869_Bitcoin-Stock	2397	5	0	0.00	0.98	0.00	0.00
17	1874_Goodreads	1234	5	3	0.60	0.33	0.00	0.00
18	1901_Netflix-10-Year	4581	5	0	0.00	3.52	0.00	0.00
19	2644_concrete	1030	8	0	0.00	1.57	0.00	0.00
20	User Knowledge	403	4	0	0.00	0.85	0.00	0.00
21	wines_SPA	6331	4	6	1.50	0.38	0.00	0.00
22	World Population Live	234	11	3	0.27	0.89	0.00	0.00
23	1222_premier_league	2565	19	0	0.00	0.73	0.00	0.00
24	real_estate_listings	4942	7	2	0.29	1.66	0.00	0.00
25	1228_Premier_League	2961	13	3	0.23	0.94	0.01	0.02
26	1594_Spotify—All-Time	1994	9	4	0.44	0.59	0.01	0.02
27	1878_COVID	2580	2	3	1.50	0.42	0.01	0.03
28	1848_Minneapolis-Air	4790	4	12	3.00	0.84	0.02	0.03
29	1712_Running-Log	689	1	15	15.00	0.74	0.02	0.03
30	1900_Another-Dataset	1538	4	3	0.75	1.02	0.04	0.04
31	0988_test_data	200	10	1	0.10	0.99	0.04	0.04
32	0274_kdd_coil_3	316	8	3	0.38	0.83	0.04	0.05
33	0235_plasma_retinol	315	10	3	0.30	1.09	0.08	0.07
34	0272_kdd_coil_1	316	8	3	0.38	0.86	0.08	0.10
35	0279_kdd_coil_5	316	8	3	0.38	0.42	0.05	0.12
36	0364_sleuth_case2002	147	2	4	2.00	0.44	0.06	0.14
37	0117_fruitfly	125	2	2	1.00	0.18	0.03	0.16
38	0273_kdd_coil_2	316	8	3	0.38	0.75	0.13	0.17
39	1755_Detailed	148	5	8	1.60	0.14	0.02	0.17
40	1417_ibm-employee	1470	11	21	1.91	1.02	0.18	0.17
41	1118_jura	359	8	9	1.12	0.88	0.18	0.20
42	1266_CSM	196	10	2	0.20	0.66	0.14	0.21
43	0911_forest_fires	517	8	4	0.50	0.48	0.11	0.24
44	1528	1197	6	8	1.33	1.14	0.30	0.27
45	1449_garments-worker	1197	6	8	1.33	1.14	0.30	0.27
46	0907_UCI-student	395	3	29	9.67	0.90	0.24	0.27
47	1267_autoMpg	392	4	1	0.25	0.77	0.21	0.27
48	1787_Lisbon-House	246	5	10	2.00	0.78	0.22	0.28
49	0149_socmob	1156	1	4	4.00	0.47	0.13	0.28
50	1640_Calculate	1030	7	1	0.14	1.20	0.37	0.31
51	mechanical_analysis	927	7	3	0.43	0.96	0.39	0.40
52	1890_ECDC-daily-data	9370	3	6	2.00	0.45	0.27	0.58
53	thyroidDF	9172	7	23	3.29	0.45	0.33	0.72
54	0261_analcatdata	108	1	2	2.00	0.46	0.36	0.78
55	2168_Intersectional	1000	13	5	0.38	0.68	0.55	0.80
56	0130_breastTumor	286	2	7	3.50	0.28	0.30	1.07
57	1616_myiris	150	2	2	1.00	0.38	0.63	1.67
58	0226_analcatdata	468	1	2	2.00	0.31	0.59	1.95
59	1660_Swiss-banknote	200	5	1	0.20	0.44	0.89	2.00
60	0225_veteran	137	2	5	2.50	0.08	0.20	2.68
61	1591_Superstore	9800	1	15	15.00	0.06	0.25	4.19
62	0125_pharynx	195	2	8	4.00	0.14	0.93	6.76
63	1872_Forest-Surfaces	8111	1	2	2.00	0.03	0.21	6.83
64	0211_analcatdata	365	1	2	2.00	0.07	0.88	12.12





Table 12: Statistics of additional 32 downstream multi-class classification datasets.

ID	Dataset name	# samples	# num.	# cat.
0	Iris	150	4	0
1	Al_index_db	62	8	3
2	all_data_updated	1275	11	11
3	milknew	1059	2	5
4	fitz_undersampled	4515	0	3
5	0181_bridges	105	2	9
6	Life_expectancy	223	3	1
7	0901_iris-example	150	4	0
8	0238_pbcseq	1945	12	6
9	1420_burst-header	1075	19	2
10	0540_MyIris	150	4	0
11	0659	151	3	2
12	1400_iriisiiiiis	150	4	0
13	0941_TEST10e627dcde	150	4	0
14	0829	150	4	0
15	0968_iris	150	4	0
16	1261_Heart-Disease	303	5	5
17	1748_Sales_DataSet_of	1000	2	9
18	2754	124	0	19
19	1310_penguins	344	4	2
20	1402_iris.test	150	4	0
21	1604_iris_reproduced	150	4	0
22	1607_Indian-Liver	583	9	2
23	1620_myiris	150	3	0
24	2215_Iris	150	4	0
25	0882_JuanFeldmanIris	150	4	0
26	1191_Students	480	4	12
27	1127_mom	140	2	1
28	1394_IRIS-flower	150	4	0
29	1618_aaaa	150	4	0
30	1738_StocksData	600	3	2
31	1811_Pokemon-with	1017	9	3

Table 13: Statistics of 25 medical domain datasets. “bin@10” denotes the dataset corresponds to the one with ID 10 in the downstream binary classification collection, and “reg@1” means the dataset corresponds to the one with ID 1 in the downstream regression collection.

ID	Dataset name	# samples	# num.	# cat.
bin@10	0445_arsenic-male	559	4	0
bin@11	0446_arsenic-female	559	4	0
bin@12	0447_arsenic-female	559	4	0
bin@17	1592_Diabetes-Data	768	8	0
bin@44	0541_plasma_retinol	315	10	3
bin@52	b_depressed	1429	12	9
bin@55	Breast_Cancer	4024	5	10
bin@66	1752_Wisconsin	699	2	8
bin@69	0446_newton_hema	140	2	1
bin@72	1011_cleve	303	5	8
bin@74	1564_Mammographic	830	1	4
bin@76	diabetes_data_upload	520	1	15
bin@77	1451_early-stage	520	1	15
bin@78	1774_Early-Stage	520	1	15
bin@79	0408_pharynx	195	2	8
reg@1	0237_arsenic-female	559	4	0
reg@2	0251_arsenic-male	559	4	0
reg@33	0235_plasma_retinol	315	10	3
reg@41	1118_jura	359	8	9
reg@52	1890_ECDC-daily-data	9370	3	6
reg@53	thyroidDF	9172	7	23
reg@55	2168_Intersectional	1000	13	5
reg@56	0130_breastTumor	286	2	7
reg@60	0225_veteran	137	2	5
reg@62	0125_pharynx	195	2	8

Table 14: Results of the baselines on 25 medical domain datasets. The average rank verifies that our approach performs the best.

ID	XGB(d)	Cat(d)	FTT(d)	Trans(d)	XGB(t)	Cat(t)	MLP(t)	Auto(t)	DCN(t)	Tab(t)	SAI(t)	FTT(t)	XTab(t)	Ours <sub>s</sub> (d)
bin@10	0.8000	0.9009	0.7514	0.8860	0.8636	0.9009	0.6897	0.9028	0.6374	0.5495	0.5617	0.7626	0.8262	0.7607
bin@11	0.8333	0.7747	0.7682	0.8216	0.7435	0.8542	0.7650	0.7878	0.8333	0.7041	0.7266	0.7643	0.5798	0.7806
bin@12	0.8738	0.8912	0.9144	0.8981	0.9398	0.9062	0.8275	0.9815	0.8333	0.9537	0.8241	0.7708	0.8796	0.9248
bin@17	0.7672	0.7556	0.8031	0.3363	0.7778	0.8134	0.7887	0.7830	0.8002	0.7087	0.7894	0.7957	0.7841	0.7913
bin@44	0.4393	0.5278	0.4352	0.3796	0.4928	0.5149	0.5319	0.4630	0.5484	0.3858	0.4115	0.5473	0.3405	0.6235
bin@52	0.4950	0.4454	0.4959	0.4417	0.4472	0.4775	0.4767	0.4699	0.5004	0.5194	0.5371	0.4529	0.4813	0.5722
bin@55	0.8442	0.8492	0.8534	0.5585	0.8539	0.8559	0.8478	0.8524	0.8507	0.8341	0.8359	0.8552	0.8532	0.8500
bin@66	0.9955	0.9946	0.9962	0.5170	0.9903	0.9975	0.9966	0.9952	0.9946	0.9665	0.9966	0.9982	0.9921	0.9993
bin@69	0.6429	0.6480	0.6276	0.6837	0.6378	0.6990	0.5204	0.6429	0.6429	0.5561	0.5204	0.5714	0.4847	0.7908
bin@72	0.8290	0.8474	0.8160	0.8398	0.7668	0.8139	0.8409	0.8561	0.8398	0.6742	0.8539	0.8593	0.8182	0.8939
bin@74	0.8780	0.8833	0.8887	0.8718	0.8792	0.8930	0.7521	0.8833	0.8776	0.8417	0.8784	0.8791	0.8630	0.8813
bin@76	0.9973	0.9984	0.9984	0.2641	0.9902	0.9941	0.9891	0.9973	0.9785	0.9031	0.9988	0.9988	0.9742	0.9973
bin@77	0.9842	0.9924	0.9984	0.7633	0.9906	0.9912	0.9723	0.9668	0.9621	0.9676	0.9902	0.9840	0.9828	0.9906
bin@78	0.9842	0.9924	0.9984	0.9859	0.9891	0.9988	0.9777	0.9902	0.9773	0.8660	0.9965	0.9863	0.9828	0.9906
bin@79	0.7667	0.7306	0.7778	0.7639	0.7819	0.7639	0.7778	0.6972	0.8028	0.5306	0.6972	0.7556	0.5194	0.8306
reg@1	0.0437	0.0210	0.0326	0.0638	0.0469	0.0178	0.0159	0.0121	0.0438	0.0600	0.0402	0.0356	0.0551	0.0087
reg@2	0.0485	0.0230	0.0402	0.0640	0.0293	0.0262	0.0732	0.0292	0.0420	0.0508	0.0540	0.0418	0.0580	0.0124
reg@33	0.1661	0.1376	0.1341	0.1322	0.1374	0.1415	0.1214	0.1428	0.1264	0.5017	0.1280	0.1279	0.1277	0.1219
reg@41	0.0820	0.0860	0.1074	0.1286	0.0779	0.0759	0.0970	0.0872	0.1175	0.3678	0.0901	0.0958	0.1544	0.0839
reg@52	0.0226	0.0203	0.1333	0.1496	0.0006	0.0137	0.1360	0.1362	0.1367	0.1387	0.1349	0.1350	0.1413	0.0033
reg@53	0.2271	0.2252	0.2396	0.2509	0.2226	0.2255	0.2454	0.2406	0.2452	0.2449	0.2399	0.2399	0.2463	0.2310
reg@55	0.2009	0.1874	0.1881	0.2098	0.1863	0.1887	0.1883	0.1917	0.1852	0.2159	0.1894	0.1858	0.1884	0.1825
reg@56	0.2165	0.1997	0.1848	0.2058	0.1931	0.2070	0.1770	0.1876	0.1825	0.8456	0.1792	0.1780	0.1961	0.1722
reg@60	0.1412	0.1102	0.1115	0.1119	0.1108	0.1083	0.1109	0.1073	0.1110	0.2060	0.1116	0.1106	0.1909	0.1117
reg@62	0.2379	0.2117	0.2059	0.2656	0.1957	0.2078	0.1921	0.2054	0.2015	0.3927	0.1880	0.1976	0.2418	0.1815
rank	8.3(3.0)	6.2(3.4)	6.1(2.8)	10.9(3.7)	6.5(3.6)	<u>4.8(3.7)</u>	8.1(3.8)	6.6(3.5)	7.5(3.8)	12.0(3.1)	7.7(3.7)	6.3(3.3)	10.6(3.0)	<b>3.5(2.8)</b>

Table 15: The average values (standard deviations) of all method ranks under several data scenarios. We use the TP-BERTa version pre-trained on 101 binary classification datasets (i.e., “TP-BERTa(single, binclass)”) under imbalanced binary classification and multi-class classification settings.

Models:	XGB(d)	Cat(d)	FTT(d)	Trans(d)	XGB(t)	Cat(t)	MLP(t)	Auto(t)	DCN(t)	Tab(t)	SAI(t)	FTT(t)	Xtab(t)	Ours(d)
Imbalanced binary classification														
$p < 1/3$	7.9(4.0)	8.0(4.3)	6.6(3.6)	10.6(4.7)	6.9(4.6)	6.4(4.2)	8.5(4.3)	7.4(3.5)	6.6(4.0)	11.2(4.2)	8.6(3.5)	6.6(3.7)	9.0(4.0)	<b>6.2(4.0)</b>
$p < 1/5$	8.1(4.2)	8.1(4.4)	6.8(2.9)	10.3(4.4)	7.0(5.1)	6.4(4.8)	10.1(4.0)	6.0(3.3)	7.2(4.3)	10.9(4.2)	10.1(3.6)	7.1(3.6)	7.8(4.2)	<b>5.2(3.2)</b>
$p < 1/8$	8.0(4.1)	7.3(4.6)	6.9(3.2)	10.2(4.2)	6.0(4.9)	6.7(4.7)	10.1(4.6)	<b>5.5(3.5)</b>	8.3(4.4)	10.8(3.9)	11.0(2.3)	6.8(3.5)	7.7(4.2)	<b>5.5(3.1)</b>
$p < 1/20$	10.0(2.5)	5.8(4.5)	8.3(4.0)	8.0(4.8)	<b>3.1(1.4)</b>	7.6(5.3)	9.0(5.4)	4.8(4.8)	8.6(4.5)	11.0(6.0)	11.9(1.7)	8.5(4.1)	6.8(1.7)	6.3(3.2)
32 additional multi-class classification tasks														
rank	5.8(2.3)	7.0(3.1)	6.4(2.9)	6.7(3.7)	<b>4.9(3.6)</b>	7.1(3.5)	7.3(4.2)	8.7(3.2)	7.1(3.3)	13.4(1.6)	7.8(3.3)	5.9(3.4)	12.7(2.4)	<u>5.1(2.2)</u>