## LLM EMBEDDINGS IMPROVE TEST-TIME ADAPTATION TO TABULAR Y|X-Shifts

Anonymous authors

Paper under double-blind review

#### ABSTRACT

For tabular datasets, the change in the relationship between the label and covariates (Y|X-shifts) is common due to missing variables. Since it is impossible to generalize to a completely new and unknown domain, we study models that are easy to adapt to the target domain even with few labeled examples. We focus on building more informative representations of tabular data that can mitigate Y|X-shifts, and propose to leverage the prior world knowledge in LLMs by serializing the tabular data to encode it. We find LLM embeddings alone provide inconsistent improvements in robustness, but models trained on them can be well adapted to the target domain even using 32 labeled observations. Our finding is based on a systematic study consisting of 7650 source-target pairs and benchmark against **261,000** model configurations trained by 20 algorithms. Our observation holds when ablating the size of accessible target data and different adaptation strategies.

#### 1 INTRODUCTION

024 025 026

004

006

007 008 009

010 011

012

013

014

015

016

017

018

019

021

023

Predictive performance degrades when the distribution of target domain shifts from that of source 027 (training) (Bandi et al., 2018; Wong et al., 2021; Hand, 2006; Ding et al., 2021a; Amorim et al., 2018). 028 Distribution shifts can be categorized into shifts in the marginal distribution of covariates (X-shifts) 029 or changes in the relationship between the label and covariates (Y|X-shifts). In computer vision, X-shifts are prevalent since high-quality human labels are consistent across different images (Recht 031 et al., 2019; Miller et al., 2021; Shankar et al., 2019); in contrast, Y|X-shifts are prevalent in tabular 032 data due to missing variables and hidden confounders. There is a large body of work addressing 033 X-shifts due to its dominance in vision and language (Li et al., 2017; Zhuang et al., 2020; Zhou et al., 2022), yet the work on Y|X-shifts remain relatively limited (Liu et al., 2023). 034

The main challenge with addressing Y|X-shifts in tabular tasks is that the source data may provide little insight on the target distribution. Since it is impossible to generalize to a completely new and unknown domain (Arjovsky et al., 2019; Rosenfeld et al., 2021), we focus on leveraging few labeled target examples (on the order of 10 to 100) to address small Y|X-shifts that negatively impact model performance. Our goal is to build a feature representation  $\phi(X)$  such that the difference between  $\mathbb{E}_{source}[Y|\phi(X)]$  and  $\mathbb{E}_{target}[Y|\phi(X)]$  are learnable even based on a few target data.

Using the wealth of world knowledge learned during pre-training, LLMs have the potential to build representations that mitigate the impact of confounders whose distribution changes across source and target. Specifically, we use a LLM encoder (e5-Mistral-7B-Instruct) to featurize tabular data—which we referred to as LLM embeddings—and fit a shallow neural network (NN) on these embeddings for tabular prediction (Figure 1). In contrast to classical numerical encoding of tabular data, our approach automatically incorporates the semantics of each covariate using off-the-shelf LLMs, and can include additional contextual domain-level information that can help account for missing variables whose distribution shifts from source to target.

In this work, we restrict attention to *lightweight* probing approaches instead of expensive end-to-end updates (Hegselmann et al., 2023; Yang et al., 2024) that require full weight access and weight updates to the LLM (we find this infeasible across the thousands of source-target pairs we consider). Throughout our investigation, we use the same LLM encoder to extract the LLM embeddings, and we only "finetune" the shallow NNs we use as the main prediction model across target domains. Investigating how different LLM encoders affect tabular Y|X shifts is left as a future work.

056

058

059

060

061

062

063

064

065

067

068

069

071

072

088

090

095



Figure 1: Overview of methods incorporating LLM embeddings.

073 For rigorous empirical evaluation, we consider 7,650 natural spatial shifts (source $\rightarrow$ target) based on 074 three real-world tabular datasets (ACS Income, Mobility, Pub.Cov (Ding et al., 2021a)). 075 Our testbed serves as a *large-scale* benchmark for Y|X-shifts on tabular data, offering a standardized 076 protocol for training, validation, testing, and finetuning, as well as a consistent hyperparameter 077 selection process. Compared to previous benchmarks on tabular distribution shifts (Liu et al., 2023), this paper not only explores a significantly greater variety of shift settings but also introduces a series of novel approaches to incorporate LLM embeddings as features. Such a comprehensive evaluation 079 ensures the robustness and adaptability of our findings across diverse scenarios, setting a new standard for future work in this domain. We compare our proposed approach with typical methods on Tabular 081 features, including basic models (LR, SVM, NN), gradient-boosting trees (GBDT; XGB, LGBM, 082 GBM), and distributionally robust methods (DRO; KL-DRO,  $\chi^2$ -DRO, Wasserstein DRO, CVaR-DRO, and Unified-DRO). In total, we consider 20 algorithms and **261,000** model configurations. 084

Since it is unrealistic to expect any single method to uniformly dominate over large number of source  $\rightarrow$  target settings, we complement traditional average-case metrics using the *fraction of times each method performs best*. For each method  $\mathcal{M}$ ,

$$FractionBest(\mathcal{M}; \Delta) := \frac{|\mathcal{S}(\Delta) \cap \mathcal{S}_{\mathcal{M}}|}{|\mathcal{S}(\Delta)|},$$
(1)

where  $S(\Delta)$  contains all source  $\rightarrow$  target settings where the performance between the best and second best model is larger than  $\Delta$  (we set  $\Delta = 1\%$  in this paper), and  $S_{\mathcal{M}}$  contains all source  $\rightarrow$  target settings where model  $\mathcal{M}$  performs the best. FractionBest calculates the proportion of source  $\rightarrow$  target settings where (i)  $\mathcal{M}$  outperforms all other methods and (ii) the improvement over the second-best model is meaningful (and significant).

First, we consider LLM embeddings without any adaptation to labeled target data<sup>1</sup>. Shallow networks based on LLM embeddings (LLM|NN) outperform all other methods on tabular features in 85% settings in the ACS Income dataset, and in 78% in the ACS Mobility dataset. However, for the ACS Pub.Cov dataset, the FractionBest drops to 44%, which indicates that LLM embeddings do not always offer a perfect solution (see Figure 2 (a)-(c)). We conclude *LLM embeddings sometimes improve robustness*, but *do not consistently surpass* state-of-the-art tree-ensemble methods.

However, we find that *finetuning* the prediction model (shallow NN) *with few target samples* can *make a big difference* even when using identical LLM embeddings. When finetuning with just 32 target samples, the FractionBest ratio (Equation (1)) remains at 85% on ACS Income, improves from 78% to 86% on ACS Mobility, and from 44% to 56% on ACS Pub.Cov (see Figure 2 (d)-(f)). We find this improvement *surprising*: although the shallow NN has numerous parameters,

<sup>&</sup>lt;sup>1</sup>We do not conduct any target adaptation; however, we use 32 labeled target samples for validation (selection of hyperparameters, etc.).



Figure 2: The FractionBest Ratio in Equation (1) (with  $\Delta = 1\%$ ). We compare our proposed methods—(a)-(c): LLM|NN and (d)-(f): LLM|NN (finetuning)—with methods on Tabular features.

finetuning with only 32 target samples surprisingly improve target performance by a relatively large margin. More importantly, such improvement is observed under Y|X shifts and holds across a wide range of distributional shift settings. This implies the potential of our novel and lightweight approach, and opens up the door for further investigation of using LLM embeddings in tabular classification tasks. Theoretical insights are discussed in Section 4.

Our method also implies multiple additional benefits (See Section 3.2). With the same amount of target samples, the finetuned NNs significantly outperform in-context learning with GPT4-mini, the SOTA decoder model (see Figure 2 (d)-(f)). Moreover, *the performance gain brought by target samples is larger under stronger* Y|X-shifts, where the level of Y|X-shifts is measured by DISDE (Cai et al., 2023). Finetuning with 32 target samples yields an average performance gain of 5.4 percentage points across the worst 500 settings on ACS Pub.Cov, compared to no finetuning. This is notably higher than the 1.2% average gain observed across all 2550 settings (**4.5 times**).

140 Beyond our primary findings, we also conduct ablation studies to better understand our approach in 141 Section 3.3. Given the large number of model parameters and limited labeled target samples, one 142 might expect parameter-efficient methods like Low-Rank Adaptation (LoRA) (Hu et al.) and Prefix 143 Tuning (Li & Liang, 2021) to offer a clear advantage for target adaptation. However, we find the 144 specific finetuning approach has small impact on target adaptability under tabular Y|X-shifts. In Fig-145 ure 7, all target adaptation methods significantly outperform the non-finetuned version when using LLM embeddings. On the other hand, incorporating the "right" domain information has an outsize 146 *impact* on adaptability to Y|X-shifts. For ACS Pub.Cov, adding additional domain information 147 from Wikipedia shows minimal improvement alone, but significantly enhances performance under 148 target adaptation. 149

Another practical question is how to allocate a fixed number of labeled target samples between target adaptation and validation (selection of finetuning method, hyperparameters, etc). In Figure 8 to come, we compare two allocation schemes of 64 labeled target samples: (i) using all 64 samples for validation (solid bar), and (ii) dividing the samples into 32 for validation and 32 for finetuning (shaded bar). For ACS Mobility and ACS Pub.Cov, *target adaptation provides significant gains over the validation-only approach*, highlighting the need for further investigation into sample allocation.

Related work Tabular data is a common modality in electronic health records, finance, and social and natural sciences. Unlike other modalities like images and text, gradient-boosted trees (GBDT; GBM (Friedman, 2001; 2002), XGBoost (Chen & Guestrin, 2016), LGBM (Ke et al., 2017)) remain the state-of-the-art (Gorishniy et al., 2021; Shwartz-Ziv & Armon, 2022; Gardner et al., 2022) even when compared to neural networks specifically designed for tabular data (Arik & Pfister, 2021; Huang et al., 2020; Kadra et al., 2021; Katzir et al., 2020). GBDTs have recently been observed to provide strong performance under distribution shifts, which forms the basis of its use as a main

baseline (Gardner et al., 2022; Liu et al., 2023). In addition, some recent works employ in-context
learning to address the few-shot classification problem (Hegselmann et al., 2023; Yang et al., 2024).
However, they are significantly more resource-intensive to finetune the LLM itself. Given the huge
number of settings, we only compare with GPT4-mini for this line of research in this work. Besides,
while we use standard LLMs to generate embeddings, LLMs specialized for tabular data (Yan et al.)
can also be used as embedding extractors.

168 A wide range of methods have been proposed to address distribution shifts, notably robust learning 169 methods, balancing methods, and invariant learning methods. Distributionally robust optimization 170 (DRO) construct an uncertainty set around the training distribution and optimizing for the worst-case 171 distribution within this set, thereby mitigating the impact of potential distribution shifts. Variants of 172 DRO methods have been developed using different distance metrics, such as  $\chi^2$ -divergence (Duchi & Namkoong, 2021; Duchi et al., 2021), KL-divergence (Hu & Hong, 2013), and Wasserstein 173 distance (Blanchet et al., 2017; 2018; 2019a; 2023a). However, these approaches have recently been 174 observed to be ineffective in addressing real-world tabular distribution shifts Liu et al. (2023). On 175 the other hand, invariant learning (Peters et al., 2016; Arjovsky et al., 2019; Koyama & Yamaguchi, 176 2020) seeks to learn causally invariant relationships across multiple pre-defined environments. In 177 this work, we include 5 typical DRO methods but do not consider invariant learning methods as 178 they require multiple training environments, which is not the focus of this work. Some statistical 179 works (Li et al., 2022; Tian & Feng, 2023) provide theoretical guarantees for simple linear models 180 in transfer learning, but these guarantees often do not extend to more complex models like decision 181 trees or neural networks used in real-world applications. This version improves clarity and flow while 182 retaining the original meaning. Additionally, many studies have explored domain adaptation (Iwasawa 183 & Matsuo, 2021; Liang et al., 2023; Chen et al., 2023). However, most of these focus on X-shifts in image data, whereas our work addresses Y|X-shifts in tabular data. 184

185 186

187

## 2 Methods

In this section, we introduce a series of methods utilizing LLM embeddings for tabular prediction, as well as different choices to incorporate additional domain information, different model architectures, and target adaptation techniques using a small amount of labeled target samples. To the best of our knowledge, this work is the first to comprehensively explore the impact of LLM embeddings on tabular Y|X-shifts.

193 194

209

210

#### 2.1 LLM EMBEDDINGS FOR TABULAR PREDICTION

We first introduce how we transform tabular data into LLM embeddings, where the key idea is to serialize each sample into a natural language format that the LLM can process. There is a substantial body of research on serialization, including using another LLM to rewrite tabular data into natural language (Hegselmann et al., 2023), adding descriptions of the classification task, training and test examples (Hegselmann et al., 2023; Slack & Singh, 2023), etc. Among these methods, Hegselmann et al. (2023) demonstrate that using a straightforward text template with a task description consistently achieves the best empirical performance.

Using an income prediction problem to illustrate, consider a simple task description such as "Classify whether US working adults' yearly income is above \$50000 in 2018." along with a simple serialization template that enumerates all features in the format "The [feature name] is [value]". Adopting this serialization approach, we employ the encoder model e5-Mistral-7B-Instruct to generate the LLM embedding. Formally, the encoder takes the serialization Serialize(X) of sample X as input and outputs its corresponding embedding  $\Phi(X)$  as

## $X \xrightarrow{\text{serialization}} \text{Serialize}(X) \xrightarrow{\text{eS-Mistral-7B-Instruct}} \Phi(X).$

211 Since e5-mistral-7b-instruct requires input data to be formatted in the following template:

212	Instruct:	description	of	the	classification	task	∖n
213	Query:	description	of	the	data,		

we provide task description in the "Instruct" part, and use the serilization template to format the tabular data in the "Query" part. An illustrative example is provided in Part A-D of Figure 1, with

additional details available in Appendix A.1. Analyzing the impact of different LLM encoders, task
 descriptions, and serialization methods is left for future work.

219 2.2 Additional Domain Information220

Another advantage of using LLM embeddings is their ability to incorporate additional domain information or prior knowledge, denoted by C. As demonstrated in Section 1, incorporating domainspecific information can help address Y|X-shifts and improve generalization performance in the target domain.

225 In this work, we propose a simple yet effective approach for integrating domain knowledge into 226 tabular predictions. Rather than combining the domain information with serialized tabular data 227 and generating a single LLM embedding, we generate separate LLM embeddings for the domain knowledge and the serialized tabular data, and then concatenate them together. The benefits of this 228 approach are twofold: (a) although the domain information may contain significantly more words than 229 the serialized tabular features, our concatenation method ensures a balanced 1:1 ratio between the 230 two, preventing a single embedding that disproportionately focuses on the longer domain information; 231 (b) by separating the tabular features from the domain information, we can efficiently update the 232 domain information without having to regenerate all the embeddings for the entire dataset. 233

234 We explore three sources of domain information: Wikipedia, GPT-4, and labeled target samples. Given that our experiments (see Table 1 and Section 3) focus primarily on socioeconomic factors, 235 we collected "Economy" data for each U.S. state from Wikipedia as C. For GPT-4, we prompt it 236 to provide background knowledge relevant to each prediction task in each state as C. For labeled 237 target samples, we serialize 32 labeled samples from the concerned domain as the prior knowledge 238 C. Further details can be found at Appendix A.2. After obtaining domain information C, we use 239 e5-mistral-7b-instruct to generate an LLM embedding for C. As illustrated in Parts E.2 240 and F.1 of Figure 1, this embedding is then concatenated with the LLM embeddings of the tabular data, 241 which serve as input to the backend neural network models (NN). This approach allows us to generate 242 the LLM embedding for the dataset *just once*, and subsequently concatenate it with embeddings 243 from different prompts as needed. In Section 3, we study whether and how this additional domain 244 information can enhance generalization under Y|X-shifts.

In addition, recent works on prompt engineering have focused on incorporating additional domain information to enhance prediction tasks, often through detailed instructions (Schick & Schütze, 2020; Shin et al., 2020). Our proposed framework introduces a novel approach to leveraging such information and remains fully compatible with these existing methods.

249 250 251

2.3 MODEL TRAINING AND TARGET ADAPTATION

Model architecture For the backend model, we use a vanilla neural network (NN) classifier on both
 tabular features and LLM embeddings for tabular data classification. The NN is a simple feedforward
 neural network with several hidden layers, dropout layer, and ReLU activation functions.

255 When adding additional domain information via an embedding layer, the same embedding is applied to all samples from the same domain. Since the output of e5-mistral-7b-instruct is a 256 4096-dimensional vector, we simply concatenate the LLM embeddings with the embeddings of the 257 domain information. This concatenated vector is then passed through the hidden layers, dropout 258 layer, and ReLU activation functions. For all NNs, the final linear layer an output dimension of 259 2, followed by a softmax layer for binary classification. During training, we use cross-entropy as 260 the loss function, batch size as 128, and use the Adam optimizer. Detailed model architecture and 261 hyper-parameters are provided in the Appendix A.3 and A.4, with a discussion on hyperparameter 262 selection provided in Section 3.1. 263

264Target AdaptationEven with the incorporation of LLM embeddings and domain information, our265model may still experience Y|X-shifts. In practice, it is common to have a small set of samples from<br/>the target domain, which can be leveraged to better adapt the model to the target domain.

For each (source domain, target domain) pair, we begin by selecting the best training hyperparameter
 based on a validation criterion, which will be discussed in the Section 3.1. Using this model trained
 on the source domain, we explore four primary methods for target adaptation: in-context domain info,
 full-parameter fine-tuning, low-rank adaptation (LoRA), and prefix tuning for domain information.

	1			·	1		
#ID	Dataset	#Samples	#Features	Outcome	#Source Domains	#Target Domains	#Source→Target Pair
1	ACS Income	1.60M	9	Income	51 (US States)	50 (US States)	2550
2	ACS Mobility	621K	21	Residential Address	51 (US States)	50 (US States)	2550
3	ACS Pub.Cov	1.12M	18	Public Ins. Coverage	51 (US States)	50 (US States)	2550
	Berformance Drop 10%- Multiple Drop 10%- 0	Mart . J. Julia 100	- Alm M. I.M. D	Mmh	11 Martine and the later of the	by by by do	Prall Drop Y X-Shifts X-Shifts

Table 1: Details of datasets used in this work. "# Source -> Target Pair" denotes the number of 270 distribution shift pair for each dataset, and we consider the natural *spatial* shift between US states.

Figure 3: Shift pattern analysis. For the 2550 source  $\rightarrow$  target distribution shift pairs in ACS Income dataset, we attribute the performance drop for each source  $\rightarrow$  target pair into Y X-shifts (red curve) and X-shifts (blue curve), and sort all pairs according to the drop introduced by Y|X-shifts. We draw the worst-500 settings in each dataset, and the decomposition method used here is DISDE (Cai et al., 2023) with XGBoost as the reference model. Results on other datasets are in Figure 9.

289 For in-context domain info (F.1 of Figure 1), we keep the trained model frozen and only update the domain information, switching it from natural language description of labeled sample from the source 290 domain during training to that of target domain during inference phase. For the other three methods, 291 we conduct further training of the model. In full-parameter fine-tuning (F.2 of Figure 1), the entire 292 neural network is fine-tuned using the target samples. For LoRA, we introduce a low-rank adaptation 293 layer to each linear layer by incorporating two smaller matrices, A and B, both with a rank of 16. Specifically, matrix A has dimensions corresponding to the input size and the rank, while matrix B 295 has dimensions corresponding to the rank and the output size. Matrix A is initialized with a mean 296 of 0 and a standard deviation of 0.02, whereas matrix B is initialized with zeros. These matrices 297 are then multiplied together and added to the original weight matrix. We then fine-tune only these 298 LoRA parameters, while keeping the rest of the model unchanged. In prefix tuning (F.3 of Figure 1), 299 the initial domain information embedding serves as a starting point for further refinement. During 300 training, both the the NN and the domain information embedding of the source domain are trained. For target adaptation, we switch the domain information embedding from the source to the target 301 domain. The NN is kept frozen, and only the domain information embedding of the target domain is 302 updated using these samples from the target domain. We refer to this process as prefix tuning. 303

304 As shown in Table 1, we use different hyperparameters for target adaptation. Detailed hyperparameters 305 are provided in Appendix A.4, and the hyperparameter selection process is discussed in Section 3.1.

306 307 308

321

#### 3 NUMERICAL EXPERIMENTS

In this section, we conduct a thorough investigation of **7650** natural shift settings (source  $\rightarrow$  target 310 domain) in 3 tabular datasets over **261,000** model configurations and summarize the observations. 311 Our findings highlight the potential of incorporating LLM embeddings to enhance the generalization 312 ability in tabular data prediction tasks. 313

314 3.1 TESTBED SETUP 315

**Dataset** In this work, we use the ACS dataset (Ding et al., 2021b) derived from the US-wide ACS 316 PUMS data, where the goal is to predict various socioeconomic factors for individuals. 317

- ACS Income: The goal is to predict whether an individual's income is above \$50K based 318 on individual features. We filter the dataset to only include individuals above 16 years old 319 with usual working hours of at least 1 hour per week in the past year, and an income of at 320 least \$100.
- ACS Mobility: The goal is to predict whether an individual has the same residential 322 address as one year ago. We filter the dataset to only include individuals between the ages of 18 and 35, which increases the difficulty of the prediction task.

274 275 276

278 279

281

283

284

285

286

287

• ACS Public Coverage (abbr. as ACS Pub.Cov): The goal is to predict whether an individual has public health insurance. We focus on low-income individuals who are not eligible for medicare by filtering the dataset to only include individuals under the age of 65 and with an income of less than \$30,000.

The details of datasets are summarized in Table 1.

**Shift Pattern Analysis** Before benchmarking, we first analyze the shift patterns among the 2550 330 source→target pairs in each dataset. Specifically, we utilize DISDE (Cai et al., 2023) (reference 331 model as XGBoost) to decompose the performance degradation from the source domain to the target 332 into two parts: (a) Y|X (concept)-shifts and (b) X (covariate)-shifts. By utilizing tailored shift 333 patterns, we can conduct an in-depth analysis of where the strength of LLM embeddings lies. As 334 shown in Figure 3, we sort all pairs according to the strength of Y|X-shifts, where we find that the 335 natural spatial shifts are mainly comprised of Y|X-shifts. These findings broaden the scope of the 336 analysis in WhyShift (Liu et al., 2023) by examining 7,650 shift pairs, a significant increase from 337 the 169 pairs studied in the original work.

Algorithms As introduced in Section 2, we compare various methods that incorporate LLM embeddings into tabular data prediction, including different finetuning methods (no finetuning, finetuing on full parameters, and low rank adaptation (LoRA)) and different embeddings (w/ or w/o extra information). Besides, in order to fully compare the performances, we also include a wide range of learning strategies that perform on Tabular features, including basic models (LR, SVM, NN), tree ensembles (XGB, LGBM, GBM), robust methods (KL-DRO,  $\chi^2$ -DRO, Wasserstein DRO, CVaR-DRO, and Unified-DRO). All methods are summarized in Table 2.

**Experiment Setup** We conduct experiments with more than **261,000** model configurations on 2550 source  $\rightarrow$  target shift pairs in ACS Income, ACS Mobility, and ACS Pub.Cov datasets respectively (**7650** settings in total). For each source  $\rightarrow$  target shift pairs, we randomly sample 20,000 labeled data from the source and target domain respectively, as the training and test dataset. We evaluate the model trained on the source domain, with or without target adaptation, and report the *Macro F1 score* on the testing dataset.

351 Given the numerous training hyperparameters—learning rate, number of training epochs, hidden 352 layer dimension, dropout ratio—we use a validation set of 32 randomly sampled labeled target 353 domain samples to choose the optimal training hyperparameters, based on the highest F1 score in 354 the validation dataset. Since our metric is Macro F1 score, the validation set is set as balanced 355 between positive and negative classes. Note that the hyperparameter selection is *near*-oracle, as it 356 leverages target samples, albeit in a limited quantity. When doing finetuning, we sample another 357 32 labeled target samples to finetune the model. And we use the same 32-sample validation dataset 358 (for training hyperparameter selection) to select the target adaptation hyperparameters that yield the best Macro F1 score. Note that our testbed allows flexible sample sizes for training, validation, 359 testing, and finetuning. Additionally, we perform an ablation study on different allocations of the 360 overall target samples in validation and finetuning (see Figure 8). See details on hyperparameters in 361 the Appendix A.4. 362

262

324

325

326

327

#### 364 3.2 PRIMARY FINDINGS

We begin by presenting the key observations from our results. In addition to the metric equation 1 introduced in Section 1, we report performance metrics averaged over the source-target pairs.

368 LLM embeddings improve performance, but when applied alone do not consistently outperform 369 tree-ensembles. To better assess the generalization ability when incorporating LLM embeddings, 370 we select the worst 500 settings (out of 2,550 total settings per dataset) based on the severity of 371 Y|X-shifts and report the average Macro F1 Score in Figure 4. Each bar represents the average 372 result across these worst 500 settings, characterized by the most severe Y|X-shifts. Thus, even a 1pp 373 improvement is significant, as it implies consistent gains of about 1pp across each of the worst 500 374 settings.

Comparing "NN on LLM embeddings" to "NN on tabular features" (with the backbone model fixed as
 NN), we observe LLM embeddings significantly enhance generalization under distribution shifts on
 the ACS Income and ACS Mobility datasets, with average improvements of 2.4pp and 9.9pp.
 Notably, "NN on LLM embeddings" even outperforms XGBoost under strong distribution shifts on



Figure 4: Average Macro F1 Score over the worst-500 settings. For each dataset, we sort the 2550 settings according to the magnitude of Y|X-shifts and select the **worst-500** settings. We calculate the average Macro F1 Score for each method. For all methods, we select the best hyper-parameters of the basic model according to 32 samples from the target domain. We use CVaR-DRO based on NN here to represent DRO methods. For finetuning methods, we use 32 target samples.



Figure 5: Average Macro F1 Score over Worst-500 settings with different #target samples. Dotted lines represent methods that do not require finetuning, while solid lines show the performance of finetuning methods relative to the number of target samples. Three figures share the same legend.



Figure 6: Average performance over the worst-*K* pairs. For each dataset, we sort the 2550 pairs according to the magnitude of Y|X-shifts and select the worst-*K* settings ( $K \in \{100, 500, 1000, 1500, 2000, 2550\}$ ). For methods requiring finetuning, we use 32 target samples here. "LLM& Target Samples" represents the in-context domain info in F.1 of Figure 1, where we embed 32 target samples as domain information. Three figures share the same legend.

these datasets. This demonstrates the potential of LLM embeddings in tabular data prediction, where
 they can contribute to more generalizable models.

A different trend is observed on the ACS Pub. Cov dataset, where the inclusion of LLM embeddings results in a performance drop for NN models. This suggests that simply incorporating LLM embeddings does not always resolve distribution shift issues; their effectiveness may vary across datasets, particularly depending on whether the LLM embeddings provide additional relevant information for the specific prediction task.

A small number of target samples can make a big difference. While incorporating LLM embeddings doesn't always yield improvements, we find that even a small number of target samples can
have a significant impact. As shown in 4, finetuning the "NN on LLM embeddings" model with just
32 target samples significantly improves the average performance across the worst 500 settings for
both the ACS Mobility and ACS Pub.Cov datasets. Notably, for the ACS Pub.Cov dataset,
where LLM embeddings alone provided no improvement, finetuning with only 32 target samples
leads to a 5.4pp gain, even surpassing XGBoost by 2.2pp. This highlights the adaptability of LLM



Figure 7: Comparison between full-parameter finetuning, LoRA, and prefix tuning. We show the average Macro F1 Score over the worst-500 settings.



Figure 8: Comparison between different allocation of target samples into validation and finetuning. We report the average Macro F1 Score over the worst-500 settings.

embeddings, making them a promising tool for harnessing their power across various downstream real-world tasks.

458 Furthermore, in Figure 5, we illustrate how the performance of finetuning methods varies with 459 different numbers of target samples. As shown, our conclusions remain consistent regardless of the number of target samples. 460

The performance gain brought by target samples is larger under stronger Y|X-shifts. As 462 shown in Figure 6 (LLM | NN v.s. LLM | NN (finetuning)), for NN using LLM embeddings, 463 finetuning with 32 target samples yields an average performance gain of 5.8pp across the worst-100 settings on ACS Pub. Cov, compared to no finetuning. This is notably higher than the 1.2pp average 465 gain observed across all 2550 settings (4.8 times). Similarly, for ACS Income, it is about 9 times. 466

- AUXILIARY FINDINGS 3.3
- 468 469

467

441

442 443

444

445

446

447 448 449

450

451 452

453

454 455

456

457

461

464

In addition to the primary findings, we have several other noteworthy observations.

470 "**Right**" domain information matters. From Figure 6, we observe that additional domain infor-471 mation from either Wikipedia or 32 target samples (via in-context domain info, F.1 of Figure 1) 472 does not lead to significant improvements (comparing the shallow blue and yellow curves with 473 the shallow red curve). However, when combined with finetuning, this extra domain information 474 performs significantly better on ACS Pub.Cov. Since finetuning can be considered a method for 475 incorporating domain-specific information, this suggests that identifying the *right* domain information 476 is crucial. 477

**Specific finetuning approaches are less crucial than expected.** Given the large number of model 478 parameters and limited labeled target samples, one might expect parameter-efficient methods like 479 Low-Rank Adaptation (LoRA) and Prefix Tuning to offer a clear advantage. However, as shown in 7, 480 all methods significantly outperform the non-finetuned version when using LLM embeddings, and 481 the choice of finetuning method appears less significant in our setting. 482

483 An exception is prefix tuning on the ACS Pub.Cov task, where performance was several percentage points lower. While this requires further investigation, our key takeaway is that under Y|X shifts, 1) 484 finetuning models using LLM embeddings can greatly enhance classification performance; 2) popular 485 finetuning methods yield comparable results.

486 Allocation of labeled target samples matters. In Figure 8, we compare two allocation schemes of 487 64 labeled target samples. For the solid bars, all 64 samples are utilized as the validation dataset for 488 hyperparameter selection. For the shaded bars, we allocate 32 samples for validation and the remaining 489 32 samples for finetuning. Based on this, we initially explore and understand the impact of sample 490 allocation. For ACS Mobility and ACS Pub.Cov, target adaptation using LLM embeddings (the shaded bar) significantly outperforms the validation approach (the solid bar). However, for ACS 491 Income, the improvement from target adaptation is only marginal. This shows that although target 492 adaptation is effective, its improvement is highly dependent on the specific distribution and Y|X shift 493 level. This raises further questions about how to optimally allocate resources. Although our findings 494 endorse considering target adaptation, identifying the best allocation strategy is left as future work. 495

496 497

498

#### 4 DISCUSSION BASED ON THEORETICAL INSIGHTS

We take a brief examination of the theoretical insights that may lie behind our empirical findings. While standard generalization bounds are vacuous for neural networks, we nevertheless find that theoretical results from domain adaptation provide a useful starting point for understanding why finetuning LLM embeddings with small target samples can result in superior generalization performance under substantial Y|X shifts, particularly in comparison to tabular features.

Let  $\phi(X)$  denote a feature map and Y a binary label. We let P and Q be the source and target distributions. We consider a model class  $\mathcal{H}$  of VC dimension d. For any model  $h \in \mathcal{H}$ , we use  $\epsilon_P(h) := \mathbb{E}_P[\mathbb{I}(h(\phi(X)) \neq Y)]$  to denote the expected 0-1 loss on the source domain and  $\hat{\epsilon}_P^m(h)$  its empirical counterpart based on m i.i.d. samples from P. We define  $\epsilon_Q(h)$  and  $\hat{\epsilon}_Q^m(h)$  on the target.

508 Although in practice we finetune on the target, we use a mixture problem as a rough approximation. 509 Suppose that we have  $(1 - \beta)m$  i.i.d. samples from source domain P and  $(\beta m)$  i.i.d. samples from 510 target domain Q. Let  $\hat{h}_{\alpha,\beta}$  be the minimizer of the  $\alpha$ -weighted empirical error

$$\hat{h}_{\alpha,\beta} := \operatorname*{arg\,min}_{h \in \mathcal{H}} \left\{ \alpha \hat{\epsilon}_Q^{\beta m}(h) + (1-\alpha) \hat{\epsilon}_P^{(1-\beta)m}(h) \right\}.$$

The following classical result from Ben-David et al. (2010, Theorem 3) bounds the generalization error on the target. Brenesition 1. For any  $\delta \in (0, 1)$  with probability at least  $1 - \delta$ 

**Proposition 1.** For any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\epsilon_{Q}(\hat{h}_{\alpha,\beta}) - \inf_{h \in \mathcal{H}} \epsilon_{Q}(h) \leq 4\sqrt{\frac{\alpha^{2}}{\beta} + \frac{(1-\alpha)^{2}}{1-\beta}} \sqrt{\frac{2d\log(2(m+1)) + 2\log\frac{8}{\delta}}{m}} + 2(1-\alpha)\underbrace{\frac{d_{\mathcal{H}\Delta\mathcal{H}}(P_{X},Q_{X})}{X-shifts}}_{Y|X-shifts} + 2(1-\alpha)\underbrace{\inf_{h \in \mathcal{H}} \{\epsilon_{P}(h) + \epsilon_{Q}(h)\}}_{Y|X-shifts},$$

$$(2)$$

520 521 522

523

519

516 517 518

511 512

where  $d_{\mathcal{H}\Delta\mathcal{H}}(\cdot, \cdot)$  denotes the  $\mathcal{H}\Delta\mathcal{H}$ -distance between two (marginal) distributions.

Since we use limited (32) labeled target samples for finetuning,  $\beta$  is small. Also, we use a smaller 524 learning rate for finetuning than for training on the source, implying that  $\alpha$  is even smaller than 525  $\beta$ . Thus, we expect the constants  $\sqrt{\frac{\alpha^2}{\beta} + \frac{(1-\alpha)^2}{1-\beta}}$  and  $(1-\alpha)$  to both be close to 1. Although 526 527 the VC-dimension d can be vacuously large for neural networks (a well-known defect in statistical 528 learning theory), having a large enough sample size m can generally make the first term comparable 529 to the next two terms. For the raw features  $\phi_{tabular}(X) = X$ , we have observed empirically that Y|X-shifts are salient, implying that the third term (related to Y|X-shifts) dominates the second term 530 (related to X-shifts), and it can also dominate the first term when Y|X shifts are particular significant. 531 In comparison, we conjecture that LLM embeddings  $\phi_{\text{LLM}}$  can reduce the gap between  $\mathbb{E}_P[\mathbb{I}(Y \neq$ 532  $h(z) \mid \phi_{\text{LLM}}(X) = z$  and  $\mathbb{E}_Q[\mathbb{I}(Y \neq h(z)) \mid \phi_{\text{LLM}}(X) = z]$  by incorporating prior knowledge 533 encoded during LLM pre-training. Since  $\epsilon_P(h) = \int \mathbb{E}_P[\mathbb{I}(Y \neq h(z)) \mid \phi(X) = z] dP_{\phi(X)}(z)$ 534 (and similarly for Q), we thus expect the third term to be much smaller than using raw features 535  $\phi_{\text{tabular}}(X)$ . This suggests that when using LLM embeddings  $\phi_{\text{LLM}}$ , the generalization bound can 536 be smaller than that of raw features  $\phi_{\text{tabular}}(X)$ , especially when Y |X shift is more significant under 537 raw features  $\phi_{\text{tabular}}(X) = X$ ; recall Figure 6. 538

539 We hope our empirical findings spur future theoretical investigations into the foundations of LLMbased target adaptation.

# 540 REFERENCES 541

542 543	Evelin Amorim, Marcia Cançado, and Adriano Veloso. Automated essay scoring in the presence of biased ratings. In <i>Association for Computational Linguistics (ACL)</i> , pp. 229–237, 2018. 1
544 545 546	Sercan Ö Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 35, pp. 6679–6687, 2021. <b>3</b>
547 548	Martin Arjovsky, Leon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. arXiv:1907.02893 [stat.ML], 2019. 1, 4
549 550 551 552 553	Peter Bandi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong, et al. From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge. <i>IEEE transactions on medical imaging</i> , 38(2):550–560, 2018. 1
554 555	Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. <i>Machine learning</i> , 79:151–175, 2010. 10
556 557 558	Jose Blanchet, Yang Kang, Fan Zhang, and Karthyek Murthy. Data-driven optimal transport cost selection for distributionally robust optimizatio. <i>arXiv:1705.07152 [stat.ML]</i> , 2017. 4
559 560	Jose Blanchet, Karthyek Murthy, and Fan Zhang. Optimal transport based distributionally robust optimization: Structural properties and iterative schemes. <i>arXiv:1810.02403 [stat.ML]</i> , 2018. 4
561 562 563	Jose Blanchet, Yang Kang, and Karthyek Murthy. Robust Wasserstein profile inference and applica- tions to machine learning. <i>Journal of Applied Probability</i> , 56(3):830–857, 2019a. 4
564 565 566	Jose Blanchet, Yang Kang, Karthyek Murthy, and Fan Zhang. Data-driven optimal transport cost selection for distributionally robust optimization. In <i>2019 winter simulation conference (WSC)</i> , pp. 3740–3751. IEEE, 2019b. 17
568 569	Jose Blanchet, Daniel Kuhn, Jiajin Li, and Bahar Taskesen. Unifying distributionally robust optimiza- tion via optimal transport theory. <i>arXiv:2308.05414 [math.OC]</i> , 2023a. 4
570 571 572	Jose Blanchet, Daniel Kuhn, Jiajin Li, and Bahar Taskesen. Unifying distributionally robust optimiza- tion via optimal transport theory. <i>arXiv preprint arXiv:2308.05414</i> , 2023b. 17
573 574	Tiffany Tianhui Cai, Hongseok Namkoong, and Steve Yadlowsky. Diagnosing model performance under distribution shift. <i>arXiv preprint arXiv:2303.02011</i> , 2023. 3, 6, 7, 16
575 576 577 578	Liang Chen, Yong Zhang, Yibing Song, Ying Shan, and Lingqiao Liu. Improved test-time adaptation for domain generalization. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 24172–24182, 2023. 4
579 580	Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System. In ACM SIGKDD International Conference on Knowledge Discovery, pp. 785–794. ACM, 2016. 3, 17
581 582 583	Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. <i>Advances in Neural Information Processing Systems</i> 34, 34, 2021a. 1, 2, 15
584 585	Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. <i>Advances in neural information processing systems</i> , 34:6478–6490, 2021b. 6
587 588	John C. Duchi and Hongseok Namkoong. Variance-based regularization with convex objectives. Journal of Machine Learning Research, 20(68):1–55, 2019. 17
589 590 591	John C. Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization. <i>Annals of Statistics</i> , 49(3):1378–1406, 2021. 4
592 593	John C. Duchi, Peter W. Glynn, and Hongseok Namkoong. Statistics of robust optimization: A generalized empirical likelihood approach. <i>Mathematics of Operations Research</i> , 46:946–969, 2021. 4

608

624

629

- Jerome H Friedman. Greedy function approximation: a gradient boosting machine. Annals of statistics, pp. 1189–1232, 2001. 3
- Jerome H Friedman. Stochastic gradient boosting. *Computational statistics & data analysis*, 38(4):
   367–378, 2002. 3
- Josh Gardner, Zoran Popovic, and Ludwig Schmidt. Subgroup robustness grows on trees: An empirical baseline investigation. *Advances in Neural Information Processing Systems*, 35:9939– 9954, 2022. 3, 4, 17
- Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting deep learning
   models for tabular data. *Advances in Neural Information Processing Systems*, 34:18932–18943,
   2021. 3
- David J Hand. Classifier technology and the illusion of progress. *Statistical Science*, 21(1):1–14, 2006. 1
- 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. 1, 4, 15
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
  et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*. 3
- <sup>615</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>617</sup>
   <sup>619</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>618</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>611</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>613</sup>
   <sup>614</sup>
   <sup>615</sup>
   <sup>615</sup>
   <sup>616</sup>
   <sup>615</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>618</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>611</sup>
   <sup>611</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>613</sup>
   <sup>614</sup>
   <sup>615</sup>
   <sup>615</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>618</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>611</sup>
   <sup>611</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>613</sup>
   <sup>614</sup>
   <sup>615</sup>
   <sup>615</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>618</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>611</sup>
   <sup>611</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>612</sup>
   <sup>613</sup>
   <sup>614</sup>
   <sup>615</sup>
   <sup>615</sup>
   <sup>616</sup>
   <sup>617</sup>
   <sup>617</sup>
   <sup>618</sup>
   <sup>618</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>619</sup>
   <sup>611</sup>
   <sup>611</sup>
   <sup>612</sup>
   <sup>611</sup>
   <sup>612</sup>
- Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data
   modeling using contextual embeddings. *arXiv preprint arXiv:2012.06678*, 2020. 3
- Yusuke Iwasawa and Yutaka Matsuo. Test-time classifier adjustment module for model-agnostic domain generalization. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 2427–2440.
   Curran Associates, Inc., 2021. 4
- Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. Well-tuned simple nets excel on tabular datasets. *Advances in neural information processing systems*, 34:23928–23941, 2021. 3
- Liran Katzir, Gal Elidan, and Ran El-Yaniv. Net-dnf: Effective deep modeling of tabular data. In International conference on learning representations, 2020. 3
- 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. *Advances in neural information processing systems*, 30, 2017. 3, 17
- Masanori Koyama and Shoichiro Yamaguchi. When is invariance useful in an out-of-distribution
   generalization problem? *arXiv preprint arXiv:2008.01883*, 2020. 4
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pp. 5542–5550, 2017. 1
- Sai Li, T Tony Cai, and Hongzhe Li. Transfer learning for high-dimensional linear regression:
   Prediction, estimation and minimax optimality. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(1):149–173, 2022. 4
- Kiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*,
  pp. 4582–4597, 2021. 3
- 647 Jian Liang, Ran He, and Tieniu Tan. A comprehensive survey on test-time adaptation under distribution shifts. *International Journal Of Computer Vision*, 2023. 4

672

680

686

687

688 689

690

- Jiashuo Liu, Tianyu Wang, Peng Cui, and Hongseok Namkoong. On the need for a language describing distribution shifts: Illustrations on tabular datasets. In *Advances in Neural Information Processing Systems 36*, 2023. 1, 2, 4, 7
- John Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar, Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation between out-of-distribution and in-distribution generalization. In *Proceedings of the 38th International Conference on Machine Learning*, 2021. 1
- Alexey Natekin and Alois Knoll. Gradient boosting machines, a tutorial. *Frontiers in neurorobotics*, 7:21, 2013. 17
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier
  Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake VanderPlas,
  Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay.
  Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.*, 12:2825–2830, 2011. 17
- Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):947–1012, 2016. 4
- Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do ImageNet classifiers
   generalize to ImageNet? In *Proceedings of the 36th International Conference on Machine Learning*,
   2019. 1
- R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42, 2000. 17
- Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The risks of invariant risk minimization.
  In Proceedings of the Ninth International Conference on Learning Representations, 2021. 1
- Timo Schick and Hinrich Schütze. Exploiting cloze questions for few shot text classification and natural language inference. *arXiv preprint arXiv:2001.07676*, 2020. 5
- Vaishaal Shankar, Achal Dave, Rebecca Roelofs, Deva Ramanan, Benjamin Recht, and Ludwig
  Schmidt. Do image classifiers generalize across time? *arXiv:1906.02168 [cs.LG]*, 2019.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*, 2020. 5
- Ravid Shwartz-Ziv and Amitai Armon. Tabular data: Deep learning is not all you need. *Information Fusion*, 81:84–90, 2022. 3
  - Dylan Slack and Sameer Singh. Tablet: Learning from instructions for tabular data. *arXiv preprint arXiv:2304.13188*, 2023. 4
  - Ye Tian and Yang Feng. Transfer learning under high-dimensional generalized linear models. *Journal* of the American Statistical Association, 118(544):2684–2697, 2023. 4
- Andrew Wong, Erkin Otles, John P Donnelly, Andrew Krumm, Jeffrey McCullough, Olivia DeTroyer Cooley, Justin Pestrue, Marie Phillips, Judy Konye, Carleen Penoza, et al. External validation of a
   widely implemented proprietary sepsis prediction model in hospitalized patients. JAMA Internal
   Medicine, 181(8):1065–1070, 2021. 1
- Jiahuan Yan, Bo Zheng, Hongxia Xu, Yiheng Zhu, Danny Chen, Jimeng Sun, Jian Wu, and Jintai
   Chen. Making pre-trained language models great on tabular prediction. In *The Twelfth International Conference on Learning Representations*. 4
- Yazheng Yang, Yuqi Wang, Sankalok Sen, Lei Li, and Qi Liu. Unleashing the potential of large language models for predictive tabular tasks in data science. *arXiv preprint arXiv:2403.20208*, 2024. 1, 4

702 703 704	Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 45(4):4396–4415, 2022.
704	I
705	Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and
707	Qing He. A comprehensive survey on transfer learning. Proceedings of the IEEE, 109(1):43–76,
708	2020. 1
700	
710	
711	
712	
713	
714	
715	
716	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
740	
740	
7/18	
740	
750	
751	
752	
753	
754	
755	

#### 756 MODEL TRAINING DETAILS А

In this section, we outline the serialization scheme, the generation of additional domain information, 759 our model architecture, and the hyperparameters used for training and target adaptation.

### A.1 SERIALIZATION SCHEME

758

760 761

762

765

769

770

771

772 773

774

775

786

787

788

789

791

792

793 794

796

797

804

805

809

763 As discussed in Section 2.1, serializing a row of tabular data, such as X, requires two components: a 764 task description and a description of the data.

766 **Task description** For the ACS Income, ACS Mobility, and ACS Public Coverage datasets, we 767 consider the same binary classification task as described in Ding et al. (2021a). We adopt concise task descriptions as suggested by Hegselmann et al. (2023), as follows: 768

- ACS Income: Classify whether US working adults' yearly income is above \$50000 in 2018.
- ACS Mobility: Classify whether a young adult moved addresses in the last year.
- ACS Public Coverage : Classify whether a low-income individual, not eligible for Medicare, has coverage from public health insurance.

776 **Description of the data** For all three ACS datasets, we utilize features as shown in (Ding et al., 777 2021a, Appendix B), along with the domain name (state) to characterize the data. Specifically, for data from the source state, we apply the source domain name, and for data from the target state, we 778 use the target domain name. 779

780 As illustrated in Figure 1 and recommended by Hegselmann et al. (2023), we adopt a straightforward 781 text template: "The feature name is value." However, for certain features with less common 782 or more complex feature names, we opt for a template as "The person (appropriate verb) value" to 783 more clearly convey the data description. For example, for the feature "Gave birth to child within the past 12 months" with the value "No", the serialization would be "The person did not give birth to a 784 785 child within the past 12 months." The features that use this alternative template are as follows:

- ACS Income: class of worker;
- ACS Mobility: class of worker, disability, employment status of parents, citizenship, military service, hearing difficulty, vision difficulty, cognitive difficulty, grandparent living with grandchildren, employment status;
- ACS Public Coverage: disability, employment status of parents, citizenship, mobility, military service, hearing difficulty, vision difficulty, cognitive difficulty, gave birth to child within the past 12 months, employment status.
- For features without associated values, we omit them during serialization.

#### A.2 ADDITIONAL DOMAIN INFORMATION

798 As shown in Section 2.2, we study three sources of domain information: Wikipedia, GPT-4, and la-799 beled target samples. As we use the e5-mistral-7b-instruct to generate an LLM embedding 800 for the domain information, we need to specify both the "Instruct" and "Query" components.

801 For the "Instruct" part, we apply the same task description as outlined in Section A.1. For the "Query" 802 part, we utilize various descriptions of the additional domain information: 803

- Wikipedia: For all three datasets, we use the "Economy" section of each state's Wikipedia page as the additional domain information.
- GPT4: For each state, we pose the following question to GPT4, and use its response as the additional domain:
  - ACS Income: "We aim to develop a classifier to determine whether U.S. individuals earned over \$50000 in 2018, using features such as age, sex, educational attainment,

```
810
                     race, class of worker, marital status, occupation, and hours worked per week over the
811
                     past 12 months. Given the unique economic and demographic profile of state_name,
812
                     how might these factors influence income levels differently compared to other U.S.
813
                     states? Please provide a 2000-word summary detailing these differences."
814
                   - ACS Mobility: "We aim to develop a classifier to determine whether a young adult
815
                     moved addresses in the last year, using features such as age, sex, educational attainment,
816
                     race, class of worker, marital status, occupation, total income, and hours worked per
                     week over the past 12 months. Given the unique economic and demographic profile of
817
                     state_name, how might these factors influence mobility levels differently compared
818
                     to other U.S. states? Please provide a 2000-word summary detailing these differences."
819
                   - ACS Public Coverage: "We aim to develop a classifier to determine whether
820
                     a low-income individual, not eligible for Medicare, has coverage from public health
821
                     insurance, using features such as age, sex, educational attainment, race, disability,
822
                     marital status, occupation, citizenship status, mobility status, military service, nativity,
823
                     total income, and employment status. Given the unique economic and demographic
824
                     profile of state_name, how might these factors influence public coverage levels
825
                     differently compared to other U.S. states? Please provide a 2000-word summary
                     detailing these differences."
827
              • labeled target samples: We use the following template for 32 labeled target samples:
828
                Here are some examples of the data:
                                                                    ١n
829
                description of one target sample \n
830
                Answer: (Yes or No).
                                                \n \n
831
                description of another target sample \n
832
                Answer:
                             (Yes or No). \n \n
833
```

We use the same data descriptions as in Section A.1, with the exception that the state/domain name is omitted, as it is represented by the labeled target samples provided.





Figure 9: Shift pattern analysis. For the 2550 source  $\rightarrow$  target distribution shift pairs in ACS Income, ACS Pub.Cov, ACS Mobility datasets respectively, we attribute the performance drop for each source  $\rightarrow$  target pair into Y|X-shifts (red curve) and X-shifts (blue curve), and sort all pairs according to the drop introduced by Y|X-shifts. We draw the *worst-500* settings in each dataset, and the decomposition method used here is DISDE (Cai et al., 2023) with XGBoost as the reference model.

850 A.3 MODEL ARCHITECTURE

...

834

835

836 837

838 839

840

841 842

843

844

845

846

847

848

849

851

852

We detail the baselines used in our paper.

Fully-connected Neural Networks (NN) Given the varying input dimensions for Tabular features,
 LLM embeddings, and LLM embeddings with additional domain information, we employ three neural networks with similar architectures. We then train these networks and conduct target adaptation using Empirical Risk Minimization (ERM).

For all three datasets using Tabular features, we use a hidden layer with output dimension being hidden layer dim as a hyperparameter, a dropuput layer with dropout ratio being a hyperparameter, ReLu activation. We then have another hidden layer with input and output dimension both being hidden layer dim, and then softmax layer with dimension 2 as the output.

For datasets using Tabular features, the network includes a hidden layer where the output dimension is set by the hyperparameter hidden layer dim. It is followed by a dropout layer (dropout ratio as a hyperparameter), ReLU activation, another hidden layer the input and output dimensions are both equal to hidden layer dim. The network concludes with a softmax layer with an output dimension of 2.

For datasets using only LLM embeddings, the input dimension is 4096, which is the output dimen-867 sion of e5-mistral-7b-instruct. The network consists of three hidden layers with fixed 868 dimensions, where the input and output dimensions are (4096, 1024), (1024, 256), and (256, 128), respectively. Each layer uses ReLU activation. Next, there's a hidden layer with an input dimension 870 of 128 and an output dimension set by the hidden layer dim hyperparameter. This is followed 871 by a dropout layer (with dropout ratio as a hyperparameter), ReLU activation, and another 872 hidden layer where both the input and output dimensions are hidden layer dim. A ReLU 873 activation follows this final hidden layer, and the network concludes with a softmax layer that outputs 874 a dimension of 2.

For datasets using LLM embeddings concatenated with additional domain information, the input dimension is 8192. We slightly update the first three hidden layers, with input and out dimensions being (8192, 2048), (2048, 512), and (512, 128), respectively. All other neural network structure and hyperparameters are the same as NNs with LLM embeddings only.

When applying low-rank adaptation (LoRA) for target adaptation, we introduce a low-rank adaptation layer to each linear layer by incorporating two smaller matrices, A and B, both with a rank of 16.
Specifically, matrix A has dimensions corresponding to the input size and the rank, while matrix B has dimensions corresponding to the rank and the output size. Matrix A is initialized with a mean of 0 and a standard deviation of 0.02, whereas matrix B is initialized with zeros. These matrices are then multiplied together and added to the original weight matrix.

Tree Ensemble Models Gardner et al. (2022) show that several tree-based methods are competitive on tabular datasets. And gradient-boosted trees (e.g., XGB (Chen & Guestrin, 2016), LGBM (Ke et al., 2017), GBM (Natekin & Knoll, 2013)) are widely considered as the state-of-the-art methods on tabular data. Therefore, we compare XGB, LGBM, and GBM in this work. For GBM, we use the standard implementations in scikit-learn (Pedregosa et al., 2011). For XGB and LGBM, we use the standard implementations in the xgboost package<sup>2</sup> and the lightgbm package<sup>3</sup>. All these methods are trained on CPUs.

893

896 897

899 900

901 902

903

904 905

906 907

B94 **DRO Methods** Distributionally robust optimization (DRO) methods optimize the worst-case loss over an ambiguity set  $\mathcal{P}$ :

$$\min_{f \in \mathcal{F}} \sup_{Q \in \mathcal{P}} \mathbb{E}_Q[\ell(f(X), Y)].$$
(3)

The ambiguity set is typically chosen as a "ball" around the training distribution  $P_{\rm tr}$ 

$$\mathcal{P}(d,\epsilon) = \{Q : d(Q, P_{\rm tr}) \le \epsilon\},\tag{4}$$

where  $d(\cdot, \cdot)$  is a distance metric between probability measures and  $\epsilon$  is the radius of set. When d is set as the generalized f-divergence (including CVaR) as:

$$d(P,Q) = E_Q \left[ f\left(\frac{dP}{dQ}\right) \right],\tag{5}$$

for the KL-DRO method (Hu & Hong, 2013), we use  $f(x) = x \log x - (x - 1)$ ; for the  $\chi^2$ -DRO method (Duchi & Namkoong, 2019), we use  $f(x) = (x - 1)^2$ ; for the CVaR-DRO problem (Rockafellar et al., 2000), we use f(x) = 0 if  $x \in [\frac{1}{\alpha}, \alpha]$  and  $\infty$  otherwise, and  $\alpha$  controls the worst-case ratio.

For Wasserstein DRO (Blanchet et al., 2019b), we choose  $d(\cdot, \cdot)$  as the Wasserstein distance. Unified-DRO (Blanchet et al., 2023b) combines Wasserstein distance and KL-divergence as  $d(\cdot, \cdot)$ , and we follow their initial Github codebases and hyperparameter selection when implementing these methods.

<sup>&</sup>lt;sup>2</sup>https://pypi.org/project/xgboost/

<sup>&</sup>lt;sup>3</sup>https://pypi.org/project/lightgbm/

Table 2: Summary of methodologies: As baselines, we use basic models (LR, SVM, and NN), GBDTs (XGB, LGBM, and GBM) and robust learning methods (KL-DRO,  $\chi^2$ -DRO, etc.). For methods utilizing LLM embeddings, we consider different ways to incorporate domain information and different model architectures and finetuning techniques.

920							
921	Name	Feature	Domain Info	Model	Adaptation	# of HPs	Part of Fig. 1
922	LR	Tabular	-	LR	No Finetuning	34	-
923	SVM	Tabular	-	SVM	No Finetuning	34	-
004	GBDT	Tabular	-	XGB, LGBM, GBM	No Finetuning	200	-
924	KL-DRO	Tabular	-	SVM	No Finetuning	138	-
925	$\chi^2$ -DRO	Tabular	-	SVM	No Finetuning	138	-
926	Wasserstein-DRO	Tabular	-	SVM	No Finetuning	17	-
007	Unified-DRO	Tabular	-	SVM	No Finetuning	150	-
927	CVaR-DRO	Tabular	-	NN	No Finetuning	200	-
928					No Finetuning	96	-
929	Tabular   NN	Tabular	-	NN	Finetuning	12	-
930	·				LoRA	15	-
931		IIM			No Finetuning	96	E.1
551	LLM   NN	Embeddings	-	NN	Finetuning	12	F.2
932		Embeddings			LoRA	15	F.2
933					No Finetuning	48	E.2
934	LLM &	LLM	Wikipedia or GPT/	NN	Finetuning	12	F.2
035	Wiki/GPT4   NN	Embeddings	wikipedia of Of 14	1111	LoRA	15	F.2
000					Prefix Tuning	18	F.3
930	LLM & In-Context	LLM	Labeled		N. D	0.6	
937	Domain Info   NN	Embeddings	Target Samples	NN	No Finetuning	96	F.1

#### A.4 HYPERPARAMETERS FOR TRAINING AND TARGET ADAPTATION

For each algorithm, we maintain a grid of candidate hyperparameters as shown in Tables 3 and 4. and perform hyperparameter selection as described in Section 3.1. When the number of hyperparameter configurations exceeds 200, we randomly select 200 configurations to reduce computational cost and maintain fairness in the comparison across all algorithms.

Model	Feature	Domain Info	# of HPs	Hyperparameter	Value Range
SVM	Tabular	-	96	C Kernel Loss Y	$ \begin{array}{c} \{1e^{-2}, 1e^{-1}, 1, 1e^{1}, 1e^{2}, 1e^{3}\} \\ \{\text{linear}, \text{RBF}\} \\ \text{Squared Hinge} \\ \{0.1, 0.3, 0.5, 1, 1.5, 2, \text{scale, aut} \end{array} $
LR	Tabular	-	23	$L_2$ penalty	$ \begin{array}{c} \{1e^{-3}, 3e^{-3}, 5e^{-3}, 7e^{-3}, 1e^{-3}, 3e^{-2}, 5e^{-2}, \dots, 1.3, 1.7, 5, 3e^{-2}, 5e^{-2}, \dots, 1.3, 1.7, 5, 1e^{1}, 5e^{1}, 1e^{2}, 5e^{2}, 1e^{3}, 5e^{3}, 1e^{-3}, 5e^{-3}, 1e^{$
Tabular   NN	Tabular	-	96	Learning Rate Hidden Layer Dim Dropout Ratio Train Epoch	$\{ \begin{array}{c} \{0.001, 0.003, 0.005, 0.01\} \\ \{16, 32, 64, 128\} \\ \{0, 0.1\} \\ \{50, 100, 200\} \end{array}$
GBM	Tabular	-	1680°	Learning Rate Num. Estimators Max Depth Min. Child Samples	$\{ \begin{array}{c} \{1e^{-2}, 1e^{-1}, 5e^{-1}, 1\} \\ \{32, 64, 128, 256\} \\ \{2, 4, 8, 16\} \\ \{1, 2, 4, 8\} \end{array}$
LGBM	Tabular	-	1680°	Learning Rate Num. Estimators $L_2$ -reg. Min. Child Samples Column Subsample Ratio (tree)	$\{ \begin{array}{c} \{1e^{-2}, 1e^{-1}, 5e^{-1}, 1\} \\ \{64, 128, 256, 512\} \\ \{0, 1e^{-3}, 1e^{-2}, 1e^{-1}, 1\} \\ \{1, 2, 4, 8, 16, 32, 64\} \\ \{0.5, 0.8, 1.\} \end{array}$
XGB	Tabular	-	1944°	Learning Rate Min. Split Loss Max. Depth Column Subsample Ratio (tree) Column Subsample Ratio (level) Max. Bins Growth Policy	$ \{ 0.1, 0.3, 1.0, 2.0 \} \\ \{ 0, 0.1, 0.5 \} \\ \{ 4, 6, 8 \} \\ \{ 0.7, 0.9, 1 \} \\ \{ 0.7, 0.9, 1 \} \\ \{ 128, 256, 512 \} \\ \{ \text{Depthwise, Loss Guide} \} $
KL-DRO $\chi^2$ -DRO Wasserstein-DRO	Tabular Tabular Tabular	-	117 117 138	Uncertainty Set Size $\epsilon$ Uncertainty Set Size $\epsilon$ Uncertainty Set Size $\epsilon$	$ \begin{array}{c} 1e^{-4}, \dots, 0.01, \dots, 0.99\\ \{1e^{-4}, \dots, 0.01, \dots, 0.99\\ f1e^{-4}, \dots, 0.01, \dots, 0.99 \end{array} $
CVaR-DRO	Tabular	-	1620°	Worst-case Ratio $\alpha$ Underlying Model Class	{0.01, 0.1, 0.2, 0.3, 0.5, 1.0 NN
Unified-DRO	Tabular	-	180	Distance Type Uncertainty Set Size $\epsilon$ $\theta_1$	$ \begin{array}{c} L_{\inf} \\ \{1e^{-3}, \dots, 9e^{-1}\} \\ {}_{\{1.001,  1.01,  1.1,  1.5,  2,  3,  5,  10,  50, \end{array} $
LLM   NN	LLM	-	48	Learning Rate Hidden Layer Dim Dropout Ratio Train Epoch	$\{ \begin{array}{c} \{0.001, 0.01\} \\ \{32, 64, 128\} \\ \{0, 0.1\} \\ \{100, 200, 300, 500\} \end{array}$
LLM & In-Context Domain Info/Wiki/GPT4   NN	LLM	Labeled Target Samples or Wikipedia or GPT4	48	Learning Rate Hidden Layer Dim Dropout Ratio Train Epoch	$\{ \begin{array}{c} \{0.001, 0.01\} \\ \{32, 64, 128\} \\ \{0, 0.1\} \\ \{100, 200, 300, 500\} \end{array}$

Table 3: Training hyperparameter grids used in all experiments.  $\diamond$  : for methods with the total grid size above 200, we randomly sample 200 configurations for fair comparisons.

## Table 4: Target adaptation hyperparameter grids used in all experiments.

Model	Target Adaptation Method	# of HPs	Hyperparameter	Value Range	
Tabular   NN or LLM   NN or LLM & Wiki/GPT4   NN	Finetuning	-	12	Learning Rate Target adaptation Epoch	${ \{10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}\} \\ \{25, 50, 100\} }$
Tabular   NN or LLM   NN or LLM & Wiki/GPT4   NN	LoRA	-	12	Learning Rate Target adaptation Epoch	${ \{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 0.01\} \\ \{25, 50, 100\} }$
LLM & Wiki/GPT4   NN	Prefix Tuning	-	18	Learning Rate Target adaptation Epoch	$ \begin{array}{c} \{10^{-5}, 10^{-4}, 10^{-3}, 0.01, 0.05, 0.1\} \\ \{25, 50, 100\} \end{array} $