DISTPFN: TEST-TIME POSTERIOR ADJUSTMENT FOR TABULAR FOUNDATION MODELS UNDER LABEL SHIFT

Anonymous authors

Paper under double-blind review

ABSTRACT

TabPFN has recently gained attention as a foundation model for tabular datasets, achieving strong performance by leveraging in-context learning on synthetic data. However, we find that TabPFN is vulnerable to *label shift*, often overfitting to the majority class in the training distribution. To address this limitation, we propose *DistPFN*, the first test-time posterior adjustment method designed for in-context tabular foundation models. DistPFN rescales predicted class probabilities by downweighting the influence of the training prior (i.e., the class distribution of the context) and emphasizing the contribution of the model's predicted posterior, without modifying the architecture or requiring additional training. We further introduce *DistPFN-T*, which incorporates temperature scaling to adaptively control the adjustment strength based on the discrepancy between prior and posterior. We evaluate our methods on over 250 OpenML datasets, demonstrating substantial improvements for various TabPFN-based models in classification tasks under label shift, while maintaining strong performance in standard settings without label shift.

1 Introduction

Tabular data is among the most prevalent data formats across various domains, such as healthcare (Johnson et al., 2016) and finance (Arun et al., 2016). Tree-based models (Chen & Guestrin, 2016; Ke et al., 2017) have consistently demonstrated strong performance on tabular tasks, owing to their ability to handle heterogeneous feature types with minimal hyperparameter tuning. Recently, deep learning (DL) methods, especially transformer-based models (Huang et al., 2020; Gorishniy et al., 2021), have emerged as strong alternatives by capturing complex feature interactions.

Among these methods, TabPFN (Hollmann et al., 2023) introduces in-context learning (ICL) to tabular classification by pretraining on synthetic datasets and producing predictions for test samples in a single forward pass. While TabPFN achieves strong performance on small-scale datasets, it suffers from scalability issues due to the quadratic complexity of self-attention (Vaswani et al., 2017). To address this limitation, several extensions have been proposed to improve inference efficiency on larger datasets (Thomas et al., 2024; Xu et al., 2025; Zeng et al., 2025).

In this paper, we highlight an overlooked limitation of TabPFN, namely its vulnerability to *label shift*, which is a critical scenario in tabular learning and frequently arises in real-world tasks (Kim et al., 2024). We observe that TabPFN tends to overfit to the majority class in the training dataset (i.e., *majority-class bias*), resulting in poor performance when the class distribution in the test dataset differs. As shown in Figure 1 and Table 1, TabPFN-v2 (Hollmann et al., 2025) exhibits a strong majority-class bias, making incorrect predictions even when trained and tested on the *same* dataset.

To this end, we propose *DistPFN*, a simple yet effective test-time adaptation method that improves the robustness of TabPFN-based models to label shift, without modifying the architecture or updating any parameters. Specifically, it adjusts the model's output distribution (i.e., **posterior**) by reweighting class probabilities based on the ratio between the posterior and the class distribution of the training dataset (i.e., **prior**). Intuitively, this adjustment downweights the influence of the training distribution and amplifies the impact of the observed test samples. Furthermore, to make the adjustment more adaptive to distributional mismatch between the labels of the training and test datasets, we propose *DistPFN-T*, which adjusts the reweighting intensity via temperature scaling, where the temperature is determined based on the discrepancy between the posterior and the prior. Unlike classical correction

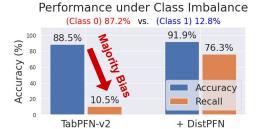


Figure 1: Majority-class bias of TabPFN
TabPFN suffers from majority-class bias, re-
sulting in poor recall for the minority class.

	abPFN-v ature 202		+	Ours)	N
(%)	$\hat{y} = 0$	$\hat{y} = 1$	(%)	$\hat{y} = 0$	$\hat{y} = 1$
y = 0 $y = 1$	87.2 11.1	0.0 1.7	y = 0 $y = 1$	81.1 3.0	6.1 9.8
Total	98.3	1.7	Total	84.1	15.9

Table 1: **Confusion matrices.** TabPFN exhibits severe *majority-class bias*, predicting **98.3%** of samples as the **majority class**, whereas DistPFN alleviates this issue through a simple *test-time adjustment*.

methods that require explicit estimation of test priors, our approach is novel in that it leverages the in-context inference structure to enable an adjustment without additional training or prior estimation.

We conduct extensive experiments on over 250 classification datasets to evaluate the effectiveness of DistPFN in both standard and label-shifted classification scenarios. As shown in Figure 2, DistPFN significantly improves accuracy of TabPFN-v2 (Hollmann et al., 2025) under label shift, with the x-axis indicating the degree of shift (see Sec. 4.4) and the y-axis showing average accuracy over 253 datasets. Our main contributions are summarized as follows:

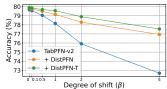


Figure 2: Robustness to shift.

- We identify an overlooked limitation of TabPFN—its vulnerability under *label shift*, as it tends to overfit to the majority class in training dataset. The performance of various TabPFN-based models degrades drastically as the degree of shift increases.
- We propose **DistPFN**, a novel and simple test-time adaptation method that adjusts TabPFN's output distribution based on the ratio between the predicted posterior and the prior representing the class distribution of the training dataset. We further introduce **DistPFN-T**, which extends this approach by applying temperature scaling to adaptively control the strength of adjustment according to the degree of distributional mismatch between the prior and the posterior.
- We present extensive evaluations on over 250 classification datasets, demonstrating that our
 methods significantly improve the performance of various TabPFN-based models under label shift,
 surpassing baseline methods and achieving state-of-the-art (SoTA) performance.
- We provide a theoretical interpretation of our method, showing that it can be viewed from both 1) classical label shift correction and 2) Bayesian inference, with details provided in Appendix D.

2 RELATED WORKS

Gradient Boosting Decision Trees (GBDTs). GBDTs (Chen & Guestrin, 2016; Ke et al., 2017; Prokhorenkova et al., 2018) are widely used for tabular data due to their strong inductive biases, minimal preprocessing, and robust performance across diverse datasets (Grinsztajn et al., 2022; McElfresh et al., 2023). Despite advances in deep learning, GBDTs remain dominant in tabular tasks, as gradient-based models struggle to encode suitable inductive biases and often underperform in several benchmarks (Grinsztajn et al., 2022; Shwartz-Ziv & Armon, 2022; McElfresh et al., 2023).

Tabular Deep Learning (DL). Transformer-based models have been proposed to capture complex feature interactions in tabular data, where TabTransformer (Huang et al., 2020) applies contextual embeddings to categorical features using self-attention. Several other architectures have also been introduced (Arik & Pfister, 2021; Somepalli et al., 2021), but these methods typically require extensive hyperparameter tuning and often fail to generalize across datasets (Kadra et al., 2021; Grinsztajn et al., 2022). Recently, TabM (Gorishniy et al., 2024a) proposes an MLP-based ensemble model that generates multiple predictions per instance with shared parameters, and ModernNCA (Ye et al., 2024) introduces a differentiable kNN approach based on neighborhood components analysis. RealMLP (Holzmüller et al., 2024) simplifies MLPs with improved design and meta-tuned default parameters, and TabR (Gorishniy et al., 2024b) augments inputs by retrieving similar training examples.

Tabular Foundation Models. TabPFN (Hollmann et al., 2023) is a foundation model for tabular classification that uses in-context learning (ICL) via large-scale synthetic pretraining. It predicts test instances by conditioning on training inputs and labels, without gradient updates. TabPFN-v2 (Hollmann et al., 2025) enhances scalability and generalization through dual-axis attention over samples and features, and structurally diverse synthetic pretraining. Limited by its computational

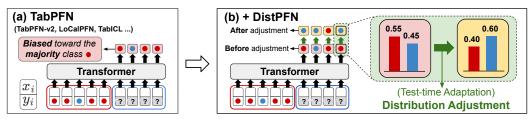


Figure 3: **Overall framework of DistPFN.** (a) **TabPFN** exhibits a *majority-class bias* under label shift, predicting test instances toward the dominant class in the training dataset. (b) **DistPFN** mitigates this bias via a simple *test-time adaptation* method that rescales the predicted class probabilities for each test instance. This scaling factor is computed from the 1) the training class distribution (i.e., prior, $p_{\text{train}}(y)$) and 2) the predicted class distribution of the test instances (i.e., posterior, $p_{\text{test}}(y)$).

complexity, TabPFN has led to several extensions, where LoCalPFN (Thomas et al., 2024) improves TabPFN by retrieving neighbors of test samples and fine-tuning on this local context, and MixturePFN (Xu et al., 2025) scales TabPFN to larger datasets by combining nearest-neighbor sampling with bootstrapped fine-tuning at inference time. TuneTable (Feuer et al., 2024) scales TabPFN to larger datasets by learning dataset-specific contexts through fine-tuning, and TabFlex (Zeng et al., 2025) replaces softmax attention of TabPFN with linear attention to improve efficiency. TabICL (Qu et al., 2025) uses a two-stage architecture with column-then-row attention to embed rows.

Label shift. Several studies have addressed label shift by rescaling classifier outputs, typically requiring estimation of the test distribution (Elkan, 2001; Lipton et al., 2018; Azizzadenesheli et al., 2019). In contrast, our method avoids test prior estimation and instead leverages only the training prior (e.g., the in-context dataset) and the predicted distribution, yielding a simple and efficient plug-in adjustment that requires no architectural modifications. Moreover, unlike Drift-Resilient TabPFN (Helli et al., 2024), which addresses temporal shift through causal pretraining, our method specifically addresses label shift and enables test-time adaptation of pretrained models without retraining.

3 PRELIMINARIES

Tabular classification. A tabular dataset consists of instances $x_i \in \mathbb{R}^d$, where each x_i is a d-dimensional feature vector composed of numerical, ordinal, or (one-hot encoded) categorical attributes. Each instance is associated with a label $y_i \in \{1, \dots, C\}$ indicating one of C predefined classes. The training dataset is denoted as $\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^{N_{\text{train}}}$, and the test dataset as $\mathcal{D}_{\text{test}} = \{x_j\}_{j=1}^{N_{\text{test}}}$. In a tabular classification task, a model f predicts the class labels y_j given the x_j from the test set.

Tabular ICL. TabPFN-based methods (Hollmann et al., 2023; 2025) predict test labels by conditioning on the labeled $\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^{N_{\text{train}}}$ and the unlabeled $\mathcal{D}_{\text{test}} = \{x_j\}_{j=1}^{N_{\text{test}}}$. Note that these models infer the corresponding test labels y_j in a single forward pass without any gradient updates.

Label shift. In this paper, we aim to address label shift, which is a distribution shift scenario where the marginal label distribution differs between training and test datasets, i.e., $p_{\text{train}}(y) \neq p_{\text{test}}(y)$. Specifically, for each class $c_k \in \{1, \dots, C\}$, the label distributions are estimated as

$$p_{\text{train}}(y=c_k) = \frac{|\{i \in [N_{\text{train}}]: y_i = c_k\}|}{N_{\text{train}}}, \quad p_{\text{test}}(y=c_k) = \frac{|\{j \in [N_{\text{test}}]: y_j = c_k\}|}{N_{\text{test}}},$$

where $p_{\text{test}}(y = c_k)$ cannot be directly computed, as test labels are not observed.

4 METHODOLOGY

In this section, we build upon **TabPFN** ¹ (Sec. 4.1), which performs ICL by conditioning on the training dataset but suffers from label shift. To address this, we propose **DistPFN** (Sec. 4.2), which reweights predicted class probabilities using the ratio between the posterior and the prior. We further introduce **DistPFN-T** (Sec. 4.3), which adaptively controls the adjustment strength using temperature scaling. Lastly, we propose a **inverse-frequency-based oversampling** method (Sec. 4.4) that oversamples rare classes in the training dataset based on inverse frequency to enable controlled evaluation under label shift. The overall framework of our method is described in Figure 3 and 4.

¹We use **TabPFN** to refer broadly to the *family of ICL-based tabular foundation models*, including variants such as TabPFN-v2 (Hollmann et al., 2025), TabICL (Qu et al., 2025), and LoCalPFN (Thomas et al., 2024)

•	b) Average accuracy by sillic								
	253 OpenML	Label Shift							
	Datasets	Low	Mid	High					
	TabPFN	81.8	78.2	72.7					
	+ DistPFN	81.8	79.1	76.9					
	+ DistPFN-T	81.8	79.6	77.5					

(h) Average accuracy by shift

Figure 4: **TabPFN vs. DistPFN vs. DistPFN-T.** (a) As TabPFN suffers from majority-class bias under label shift, DistPFN mitigates this by adjusting its posterior distribution based on its own predictions. DistPFN-T further refines the adjustment by dynamically scaling its strength based on the discrepancy between the prior and posterior distributions. (b) Under label shift, both DistPFN and DistPFN-T outperform TabPFN, with DistPFN-T showing the most consistent improvements.

4.1 TABPFN

In tabular in-context learning, TabPFN predicts the label of a test instance x_j by conditioning on the entire training dataset $\mathcal{D}_{\text{train}}$. Let $f(x_j, \mathcal{D}_{\text{train}}) \in \mathbb{R}^C$ denote the model output logits for the C classes. The posterior distribution of TabPFN is computed via softmax:

$$\widehat{p}_{\text{TabPFN}}(y \mid x_j, \mathcal{D}_{\text{train}}) = \frac{\exp\left(f(x_j, \mathcal{D}_{\text{train}})[y]\right)}{\sum_{c=1}^{C} \exp\left(f(x_j, \mathcal{D}_{\text{train}})[c]\right)},\tag{1}$$

where $[\cdot]$ denotes indexing over class logits. For simplicity, we denote the posterior distribution of any model $\widehat{p}(y \mid x_j, \mathcal{D}_{train})$ as $\widehat{p}(y)$, omitting the notations of input and training dataset.

4.2 DISTPFN: TEST-TIME POSTERIOR ADJUSTMENT

To improve the robustness under label shift, DistPFN extends TabPFN by adjusting the model's output, or the predicted distribution of x_j . Specifically, it introduces an adjustment factor based on the ratio between the predicted distribution $\widehat{p}_{\text{TabPFN}}(y)$ and the training prior $p_{\text{train}}(y)$ as:

$$\widetilde{p}_{\text{DistPFN}}(y) = \text{Norm}(\widehat{p}_{\text{TabPFN}}(y) \cdot \underbrace{\frac{\widehat{p}_{\text{TabPFN}}(y)}{p_{\text{train}}(y)}}_{\text{Adjustment factor } (\alpha)}) = \text{Norm}(\frac{\widehat{p}_{\text{TabPFN}}(y)^{2}}{p_{\text{train}}(y)}), \tag{2}$$

where $Norm(\cdot)$ denotes normalization over classes to ensure the probabilities sum to one. This adjustment down-weights the influence of the training prior and instead emphasizes the model's own prediction at test time, without requiring any modification to the model architecture or parameters.

Note that unlike conventional methods where the training distribution is only *indirectly* encoded in model parameters, ICL-based TabPFN *directly* conditions on the training dataset (i.e., context) at inference time, thereby enabling access to training data during prediction. This property makes our method specifically tailored to mitigate such bias in tabular foundation models.

4.3 DISTPFN-T: TEMPERATURE-SCALED ADJUSTMENT

While DistPFN corrects for label shift using the training prior, the optimal strength of adjustment may vary depending on the deviation of test-time predictions from the training prior. To make posterior adjustment more adaptive to the discrepancy between the training and the test dataset, we propose DistPFN-T, which introduces temperature scaling to control the sharpness of adjustment based on the discrepancy between the training prior $p_{\text{train}}(y)$ and the predicted distribution $\widehat{p}_{\text{TabPFN}}(y)$. Specifically, a temperature value τ is computed using the cross-entropy (CE) between the two distributions, where the predicted distribution $\widehat{p}_{\text{TabPFN}}(y)$ is then passed through a temperature-scaled softmax as:

$$\widehat{p}_{\text{TabPFN-T}}(y=c) = \frac{\exp\left(\widehat{p}_{\text{TabPFN}}(y=c)/\tau\right)}{\sum_{c'=1}^{C} \exp\left(\widehat{p}_{\text{TabPFN}}(y=c')/\tau\right)}, \quad \text{where} \quad \tau = \text{CE}(\widehat{p}_{\text{TabPFN}}(y), p_{\text{train}}(y)).$$
(3)

When the predicted distribution strongly deviates from the training prior (i.e., high cross-entropy), a high temperature is applied to smooth the predictions, thereby preventing over-adjustment. This temperature-scaled distribution $\widehat{p}_{\text{TabPFN-T}}(y)$ is used as the numerator of the adjustment factor in DistPFN-T, replacing $\widehat{p}_{\text{TabPFN}}(y)$ used in DistPFN as:

$$\widetilde{p}_{\text{DistPFN-T}}(y) = \text{Norm}(\widehat{p}_{\text{TabPFN}}(y) \cdot \underbrace{\frac{\widehat{p}_{\text{TabPFN-T}}(y)}{p_{\text{train}}(y)}}_{\text{Adjustment factor }(\alpha)}). \tag{4}$$

Algorithm 1 Pseudocode for DistPFN

216

217

218

219

220

221

222

224

225

226

227

228

229

230

231232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250 251

252 253

254

255

256

257

258

259260

261 262

264

265 266

267

268

269

```
# x_test: test instance(s)
# D_train: training dataset
# p_train: training class prior
# f: TabPFN-based model
# alpha: adjustment factor
# method: ["tabpfn", "distpfn", "distpfn-t"]
logits = f(x_test, D_train)
p_hat = softmax(logits, dim=0)
if method == "tabpfn":
    alpha = 1
elif method == "distpfn":
    alpha = p_hat/p_train
elif method == "distpfn-t":
    tau = cross_entropy(p_hat, p_train)
    p_hat_scaled = softmax(p_hat/tau, dim=0)
    alpha = p_hat.scaled/p_train
p_hat = alpha * p_hat
return p_hat/p_hat.sum(dim=0)
```

: Stronger label shift (Label 0:1) 10%:90% (Label 0:1) 30%:70% 1.0 0.8 Train Proportion 0.6 0.2 Proportion 0.6 Label 0 Label 0 Label 1 Label 1 0.0 0.0 1.0 1.0 0.8 0.8 est Proportion 0.6 Proportion 9.0 Label 0 Label 1 Label 1 Shift Strength (β) Shift Strength (β)

Figure 5: Inverse-frequency-based oversampling. As β increases, the training distribution becomes increasingly biased toward rare classes, while the test distribution remains fixed, enabling fair comparison across varying degrees of shift.

Table 2 shows the example where DistPFN-T smooths the adjusted probabilities of DistPFN based on the prediction of TabPFN and the training prior. When the prediction of TabPFN aligns with the majority class (i.e., leans toward the majority class in the training prior), as in **Case 1**) **Majority**, DistPFN-T amplifies the *minority* class more than DistPFN, potentially inducing overprediction. Conversely, when the prediction of

Prediction	Case 1) Majority	Case 2) Minority		
Class	A	В	A	В	
Prior	0.80	0.20	0.80	0.20	
TabPFN + DistPFN + DistPFN-T	0.60 0.36 0.33	0.40 (-) 0.64 (↑) 0.67 (↑↑)	0.40 0.10 0.12	0.60 (-) 0.90 (↑↑) 0.88 (↑)	

Table 2: Prediction w/ DistPFN and DistPFN-T.

TabPFN aligns with the minority class, as in **Case 2**) **Minority**, DistPFN-T amplifies the *majority* prediction to mitigate overprediction. This design is intuitive, as it counterbalances the bias introduced by the label distribution and the model's own prediction, leading to more calibrated and stable outputs. The effectiveness of DistPFN-T over DistPFN is demonstrated in Table 3 and Figure 7, using three different foundation models under varying degrees of distribution shift.

As TabPFN takes the entire test dataset (i.e., multiple instances) as input at once rather than processing each instance individually (i.e., single instance), the proposed methods apply their adjustment based on either the predicted distribution of a single instance or the average distribution of the test dataset, with the latter used by default in our experiments. A comparison of these two strategies is presented in Table 9, demonstrating the robustness of this design choice.

4.4 BENCHMARK FOR LABEL SHIFT: INVERSE-FREQUENCY-BASED OVERSAMPLING

To enable controlled evaluation under label shift, we propose *inverse-frequency-based oversampling*, which modifies the label distribution of $\mathcal{D}_{\text{train}}$ by oversampling each class according to its inverse frequency, while keeping $\mathcal{D}_{\text{test}}$ unchanged for direct comparison with the non-shifted setting. Note that we adopt oversampling rather than undersampling to avoid potential performance degradation caused by the removal of training instances. We define the class frequency in a dataset \mathcal{D} consisting of N samples as $p(y=c_k)=\frac{|i\in[N]:y_i=c_k|}{N}$, where the class-wise sampling weights are computed using their inverse frequency, assigning higher weights to rarer classes in the original distribution as:

$$w_k = \left(\frac{1}{p(y=c_k)}\right)^{\beta}, \qquad \tilde{w}_k = \frac{w_k}{\sum_{j=1}^C w_j},$$

with $\beta \geq 0$ controlling the strength of the shift.

Figure 5 illustrates the class distributions of the training and test datasets after oversampling, with respect to 1) the shift strength (β) and 2) the class distribution of the entire dataset. Higher values of β assign higher sampling probabilities to rare classes, inducing a stronger shift, whereas $\beta=0$ corresponds to uniform sampling (i.e., equal proportions across all classes). Note that $\beta=0$ does not imply the absence of label shift, as the dataset itself may be class-imbalanced.

-		Methods	vyla ahift			Shif	\hat{t} strength (β)			
		Methods	w/o shift	0.0	0.1	0.5	1.0	2.0	5.0	Avg.
		LogReg. + HPO	$\begin{array}{c} 0.765_{\pm 0.002} \\ 0.771_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.719_{\pm 0.002} \\ 0.697_{\pm 0.005} \end{array}$	$\begin{array}{c} 0.709_{\pm 0.002} \\ 0.687_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.674_{\pm 0.004} \\ 0.653_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.634_{\pm 0.004} \\ 0.616_{\pm 0.006} \end{array}$	$\begin{array}{c} 0.597_{\pm 0.002} \\ 0.586_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.566_{\pm 0.004} \\ 0.550_{\pm 0.003} \end{array}$	0.650 0.631
		SVM + HPO	$\begin{array}{c} 0.780_{\pm 0.003} \\ 0.784_{\pm 0.003} \end{array}$	$0.684_{\pm 0.002} \\ 0.731_{\pm 0.001}$	$\begin{array}{c} 0.646_{\pm 0.004} \\ 0.689_{\pm 0.008} \end{array}$	$\begin{array}{c} 0.560_{\pm 0.005} \\ 0.626_{\pm 0.003} \end{array}$	$0.531_{\pm 0.003} \\ 0.597_{\pm 0.005}$	$\begin{array}{c} 0.486_{\pm 0.004} \\ 0.570_{\pm 0.005} \end{array}$	$0.448_{\pm 0.004} \\ 0.541_{\pm 0.007}$	$0.559 \\ 0.626$
in:	9	MLP + HPO	$\begin{array}{c} 0.778_{\pm 0.004} \\ 0.795_{\pm 0.005} \end{array}$	$0.658_{\pm 0.005} \ 0.706_{\pm 0.006}$	$0.647_{\pm 0.004} \\ 0.688_{\pm 0.006}$	$\substack{0.613_{\pm 0.005}\\0.654_{\pm 0.008}}$	$0.574_{\pm 0.006} \\ 0.615_{\pm 0.006}$	$\begin{array}{c} 0.527_{\pm 0.004} \\ 0.580_{\pm 0.005} \end{array}$	$0.493_{\pm 0.001} \\ 0.545_{\pm 0.005}$	$0.585 \\ 0.631$
i du	37	kNN + HPO	$\begin{array}{c} 0.765_{\pm 0.004} \\ 0.783_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.663_{\pm 0.004} \\ 0.693_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.657_{\pm 0.004} \\ 0.684_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.629_{\pm 0.003} \\ 0.644_{\pm 0.004} \end{array}$	$\begin{array}{c} 0.589_{\pm 0.003} \\ 0.588_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.538_{\pm 0.003} \\ 0.540_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.501_{\pm 0.004} \\ 0.498_{\pm 0.002} \end{array}$	$0.596 \\ 0.608$
Machine Learning		Random Forest + HPO	$0.796_{\pm 0.003} \\ 0.803_{\pm 0.002}$	$\begin{array}{c} 0.768_{\pm 0.003} \\ 0.771_{\pm 0.002} \end{array}$	$0.765_{\pm 0.003} \\ 0.767_{\pm 0.001}$	$0.748_{\pm 0.005} \\ 0.743_{\pm 0.004}$	$0.718_{\pm 0.004} \\ 0.701_{\pm 0.004}$	$\begin{array}{c} 0.665_{\pm 0.005} \\ 0.627_{\pm 0.008} \end{array}$	$0.618_{\pm 0.005} \\ 0.578_{\pm 0.006}$	0.714 0.698
		LightGBM + HPO	$0.789_{\pm 0.003}$ $0.790_{\pm 0.006}$	$\begin{array}{c} 0.758_{\pm 0.004} \\ 0.726_{\pm 0.008} \end{array}$	$\begin{array}{c} 0.753_{\pm 0.002} \\ 0.661_{\pm 0.005} \end{array}$	$\begin{array}{c} 0.734_{\pm 0.003} \\ 0.655_{\pm 0.008} \end{array}$	$\begin{array}{c} 0.705_{\pm 0.004} \\ 0.608_{\pm 0.008} \end{array}$	$\begin{array}{c} 0.657_{\pm 0.005} \\ 0.577_{\pm 0.015} \end{array}$	$\begin{array}{c} 0.618_{\pm 0.005} \\ 0.551_{\pm 0.004} \end{array}$	0.704 0.630
		CatBoost + HPO	$\begin{array}{c} 0.803_{\pm 0.001} \\ 0.802_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.774_{\pm 0.002} \\ 0.774_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.771_{\pm 0.002} \\ 0.771_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.751_{\pm 0.004} \\ 0.752_{\pm 0.004} \end{array}$	$\begin{array}{c} 0.718_{\pm 0.004} \\ 0.719_{\pm 0.004} \end{array}$	$\begin{array}{c} 0.665_{\pm 0.005} \\ 0.665_{\pm 0.006} \end{array}$	$\begin{array}{c} 0.621_{\pm 0.005} \\ 0.621_{\pm 0.005} \end{array}$	$0.717 \\ 0.717$
۵۵	Non-found.	FT-Transformer TabM TabulaRNN MambaTab RealMLP	$\begin{array}{c} 0.784_{\pm 0.002} \\ 0.794_{\pm 0.002} \\ 0.749_{\pm 0.003} \\ 0.719_{\pm 0.004} \\ 0.794_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.748_{\pm 0.004} \\ 0.762_{\pm 0.004} \\ 0.699_{\pm 0.003} \\ 0.629_{\pm 0.006} \\ 0.760_{\pm 0.004} \end{array}$	$\begin{array}{c} 0.746_{\pm 0.004} \\ 0.757_{\pm 0.004} \\ 0.684_{\pm 0.003} \\ 0.603_{\pm 0.004} \\ 0.758_{\pm 0.005} \end{array}$	$\begin{array}{c} 0.718_{\pm 0.005} \\ 0.735_{\pm 0.003} \\ 0.641_{\pm 0.004} \\ 0.525_{\pm 0.002} \\ 0.745_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.674_{\pm 0.005} \\ 0.694_{\pm 0.005} \\ 0.585_{\pm 0.009} \\ 0.466_{\pm 0.010} \\ 0.720_{\pm 0.005} \end{array}$	$\begin{array}{c} 0.610_{\pm 0.003} \\ 0.624_{\pm 0.006} \\ 0.522_{\pm 0.011} \\ 0.430_{\pm 0.005} \\ 0.677_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.551_{\pm 0.007} \\ 0.565_{\pm 0.006} \\ 0.465_{\pm 0.008} \\ 0.394_{\pm 0.002} \\ 0.643_{\pm 0.004} \end{array}$	0.675 0.690 0.599 0.508 0.717
Deep Learning	u	LoCalPFN + DistPFN + DistPFN-T	$\begin{array}{c} \textbf{0.816}_{\pm 0.002} \\ \textbf{0.816}_{\pm 0.002} \\ \textbf{0.816}_{\pm 0.002} \end{array}$	$0.794_{\pm 0.003} \\ \underline{0.797}_{\pm 0.001} \\ 0.798_{\pm 0.002}$	$0.793_{\pm 0.004} \\ \underline{0.796}_{\pm 0.002} \\ 0.797_{\pm 0.002}$	$0.788_{\pm 0.003} \\ \underline{0.794}_{\pm 0.002} \\ 0.796_{\pm 0.002}$	$0.778_{\pm 0.002} \\ \underline{0.790}_{\pm 0.002} \\ 0.794_{\pm 0.002}$	$0.753_{\pm 0.004} \\ \underline{0.782}_{\pm 0.001} \\ 0.787_{\pm 0.001}$	$0.719_{\pm 0.000} \\ \underline{0.770}_{\pm 0.003} \\ 0.776_{\pm 0.003}$	0.771 0.788 0.791
Dee	Foundation	TabICL + DistPFN + DistPFN-T	$\begin{array}{c} \textbf{0.806}_{\pm 0.002} \\ \textbf{0.806}_{\pm 0.002} \\ \textbf{0.806}_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.783_{\pm 0.003} \\ \underline{0.786}_{\pm 0.002} \\ 0.786_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.781_{\pm 0.003} \\ \underline{0.786}_{\pm 0.002} \\ 0.786_{\pm 0.003} \end{array}$	$0.770_{\pm 0.003} \\ \underline{0.781}_{\pm 0.002} \\ \underline{0.783}_{\pm 0.002}$	$0.747_{\pm 0.003} \\ \underline{0.776}_{\pm 0.002} \\ \underline{0.780}_{\pm 0.002}$	$0.704_{\pm 0.006} \\ \underline{0.763}_{\pm 0.002} \\ \underline{0.771}_{\pm 0.001}$	$\begin{array}{c} 0.664_{\pm 0.006} \\ \underline{0.746}_{\pm 0.004} \\ 0.755_{\pm 0.004} \end{array}$	0.742 0.773 0.777
	F	TabPFN-v2 + DistPFN + DistPFN-T	$\begin{array}{c} \textbf{0.818}_{\pm 0.004} \\ \textbf{0.818}_{\pm 0.002} \\ \textbf{0.818}_{\pm 0.002} \end{array}$	$\begin{array}{c} \underline{0.797}_{\pm 0.003} \\ \textbf{0.799}_{\pm 0.001} \\ \textbf{0.799}_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.796_{\pm 0.004} \\ \underline{0.797}_{\pm 0.002} \\ 0.798_{\pm 0.002} \end{array}$	$0.790_{\pm 0.002} \\ \underline{0.795}_{\pm 0.002} \\ 0.797_{\pm 0.002}$	$\begin{array}{c} 0.782_{\pm 0.002} \\ \underline{0.791}_{\pm 0.003} \\ 0.796_{\pm 0.003} \end{array}$	$0.759_{\pm 0.003}$ $0.783_{\pm 0.003}$ $0.789_{\pm 0.003}$	$0.727_{\pm 0.003} \\ \underline{0.769}_{\pm 0.003} \\ 0.775_{\pm 0.003}$	0.775 0.789 0.792

Table 3: **Tabular classification results.** While most baselines suffer substantial performance degradation under label shift, our methods significantly improve the accuracy of ICL-based tabular foundation models (e.g., TabPFN-v2) across varying degrees of shift (β), averaged over 253 datasets.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Task and metrics. We evaluate our methods on tabular classification tasks both with and without label shift to assess their robustness to label shift. For the evaluation metrics, we employ accuracy (Acc.), ROC-AUC, and average rank (Rank), following the previous works (Hollmann et al., 2023). Further details regarding the experimental setups are provided in Appendix A.

Datasets. We evaluate our methods on 253 tabular classification datasets from OpenML (Bischl et al., 2017), which span a wide range of feature dimensions, class cardinalities, sample sizes, and domains. Unless otherwise specified, performance is reported as the mean accuracy across all datasets, averaged over five different random seeds. Additionally, to specifically assess performance under label shift, we construct synthetic variants by modifying the test set while keeping the training set fixed, following the standard setup for fair comparison, as described in Section 4.4. Following the previous works (Hollmann et al., 2023; 2025), all methods are evaluated using the fixed train/test splits, where each dataset is randomly split into 50% training and 50% test data.

Baseline models. We categorize a total of 15 baseline tabular models into the following three groups:

- ML models (7): Logistic Regression (LR), Support Vector Machines (SVM), Random Forest (Liaw & Wiener, 2002), k-nearest neighbors (kNN), Multi-layer Perceptrons (MLP), LightGBM (Ke et al., 2017), CatBoost (Prokhorenkova et al., 2018)
- **DL** (**non-foundation**) **models** (**5**): FT-Transformer (Gorishniy et al., 2021), TabM (Gorishniy et al., 2024a), TabulaRNN (Thielmann & Samiee, 2024), MambaTab (Ahamed & Cheng, 2024), RealMLP (Holzmüller et al., 2024)
- DL (foundation) models based on ICL (3): TabPFN-v2 (Hollmann et al., 2025), LoCalPFN (Thomas et al., 2024), TabICL (Qu et al., 2025)

Additionally, we perform hyperparameter optimization (HPO) for ML models² for stronger baselines, using the search space provided in a public implementation, with details provided in Appendix F.

²While MLP can be regarded as a DL model, we categorize it as an ML model for the purpose of HPO.

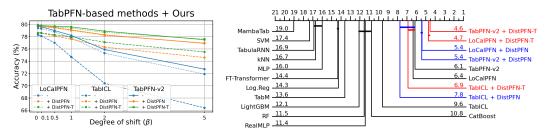


Figure 7: Performance by shift.

Figure 8: Rank (CD Diagram) under shift ($\beta = 2$).

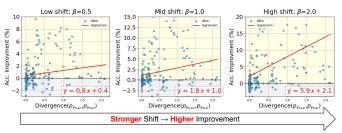


Figure 9: **Per-dataset improvement.** The figure shows the accuracy improvement for each dataset under varying β values with DistPFN-T applied, shown against the KL-divergence between the train and test label distributions of the original dataset.

Figure 10: vs. Oracle. The figure compares the performance of our method to DistPFN-Oracle, which uses the ground-truth test label distribution as the adjustment factor.

5.2 CLASS-IMBALANCED BENCHMARK DATASETS

We examine the degree of class imbalance across 253 OpenML datasets by defining the *balance ratio* as the number of samples in the minority class ($N_{\rm minority}$) divided by the number of samples in the majority class ($N_{\rm majority}$). A balance ratio of 100% corresponds to a perfectly balanced dataset, while lower values indicate increasing imbalance. As shown in Figure 6, approximately 85% of the datasets exhibit class imbalance, highlighting the importance of addressing the majority-class bias.

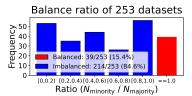


Figure 6: Balance ratio of datasets.

5.3 TABULAR CLASSIFICATION

Table 3 reports the average accuracy over 253 datasets across six levels of label shift (β) , comparing our method against 16 baselines, including three tabular foundation models to which our method is applied. For LoCalPFN, we use the k=10 nearest neighbors for each test sample, with robustness to k demonstrated in Table 5. While maintaining the original performance under the standard setting without label shift, our method substantially improves all three foundation models under shift, without any parameter update or additional computational cost, as shown in Table 11.

Figure 7 shows that as label shift strength β increases, our methods improve all three foundation models, yielding larger gains at higher β , with DistPFN-T providing additional improvements over DistPFN. For instance, DistPFN and DistPFN-T improve the accuracy of TabICL by 12.3% and 13.7% under $\beta=5$, respectively. Additionally, Figure 8 shows the average rank across datasets under $\beta=2$ using a critical difference (CD) diagram, with DistPFN-T applied to TabPFN-v2 achieving SoTA performance. A comparison based on ROC-AUC is provided in Appendix L.

6 ANALYSIS

In this section, we conduct various analyses on the effectiveness of our methods, DistPFN and DistPFN-T, which are applied to TabPFN-v2 (Hollmann et al., 2025) unless otherwise stated.

Performance gain of each dataset. Figure 9 illustrates the accuracy improvement across 253 datasets under three different β values when DistPFN-T is applied. Each point represents a dataset, with the x-axis showing the KL-divergence between the train and test label distributions of the original dataset, and the y-axis indicating the corresponding accuracy improvement. As β increases, datasets with larger original discrepancies become more imbalanced, with those exhibiting stronger divergence (i.e., larger induced shifts) benefiting more from our method.

(T)		w/o shift			Shif	t strengtl	ı (β)		
$\alpha = \frac{1}{p_{\text{train}}(y)}$	1	w/o silit	0.0	0.1	0.5	1.0	2.0	5.0	Avg.
TabPFN-v2	-	0.818	0.797	0.796	0.790	0.782	0.759	0.727	0.775
+ DistPFN + DistPFN-T + DistPFN-Oracle	$\begin{array}{c} \widehat{p}_{\text{TabPFN}}(y) \\ \widehat{p}_{\text{DistPFN-T}}(y) \\ p_{\text{test}}(y) \end{array}$	0.818 0.818 0.818	0.799 0.799 0.803	0.797 0.798 0.802	0.795 0.797 0.800	0.791 0.796 0.797	0.783 0.789 0.792	0.769 0.775 0.784	0.789 0.792 0.796

Table 4: **Comparison with oracle.** While our method computes the adjustment factor α using the *predicted* test label distribution, we also compare it to *DistPFN-Oracle*, which instead uses the *true* test label ratio. The results demonstrate that our methods achieve performance close to this oracle.

k	Methods	Avg.
3	LoCalPFN + DistPFN + DistPFN-T	0.750 0.782 0.785
10	LoCalPFN + DistPFN + DistPFN-T	0.770 0.787 0.789
20	LoCalPFN + DistPFN + DistPFN-T	0.771 0.788 0.791

$\alpha = \frac{\widehat{p}_{\text{DistPFN}(-T)}(y)}{\textcircled{2}}$	2	w/o shift			Shif	t strengtl	$\mathbf{n}\left(eta ight)$		
	2	w/o silit	0.0	0.1	0.5	1.0	2.0	5.0	Avg. 0.775 0.789 0.789 0.792
TabPFN-v2	-	0.818	0.797	0.796	0.790	0.782	0.759	0.727	0.775
+ DistPFN	$p_{ ext{train}}(y) \ \widehat{p}_{ ext{train}}(y)$	0.818 0.818	0.799 0.799	0.797 0.797	0.795 0.795	0.791 0.792	0.783 0.783	0.769 0.768	
+ DistPFN-T	$p_{ ext{train}}(y) \ \widehat{p}_{ ext{train}}(y)$	0.818 0.818	0.799 0.800	0.798 0.800	0.797 0.798	0.796 0.797	0.789 0.791	0.775 0.777	0.792 0.793

Table 5: Application to LoCalPFN.

Table 6: **Training prior vs. Training prediction.** Replacing the training prior $p_{\text{train}}(y)$ with the predicted distribution $\widehat{p}_{\text{train}}(y)$ in α shows negligible performance difference, validating the use of $p_{\text{train}}(y)$ as a simple and reliable choice, as $\widehat{p}_{\text{train}}(y)$ requires additional computation.

①: p_{train} , ②: $\hat{p}_{\text{TabPFN}}(y)$		w/o shift			Shif	t strengtl	n (β)		
		W/O SHIIL	0.0	0.1	0.5	1.0	2.0	5.0	Avg.
TabPFN-v2		0.818	0.797	0.796	0.790	0.782	0.759	0.727	0.775
+ DistPFN-T	$\tau = CE(1, 2)$ $\tau = CE(2, 1)$	0.818 0.818	0.799 0.799	0.799 0.798	0.797 0.797	0.795 0.796	0.788 0.789	0.769 0.775	0.791 0.792

Table 7: **Asymmetric cross-entropy.** DistPFN-T employs asymmetric cross-entropy to compute the temperature τ for temperature scaling, where both directions outperform TabPFN-v2.

Comparison with oracle. The adjustment factor α in our method is computed using the *predicted* label distribution. We compare this to a variant that uses the *true* label ratio of the test set, referred to as *DistPFN-Oracle*, which is unavailable in practice. This replaces the *predicted* distribution with the *true* distribution as the numerator of α . Table 4 and Figure 10 show the results, indicating that our methods achieve performance close to the oracle even without access to the (ground truth) test label.

Robustness to k **for LoCalPFN.** We apply our methods to LoCalPFN (Thomas et al., 2024), which improves the efficiency of TabPFN by retrieving k nearest neighbors for each test sample to construct the training dataset. For a fair comparison, we do not fine-tune the model and instead use the pretrained weights of TabPFN-v2 (Hollmann et al., 2025), providing a stronger baseline than the original TabPFN (Hollmann et al., 2023) used in LoCalPFN. Table 5 reports results across different values of number of neighbors (k), averaged over six β values. The results indicate that our methods consistently improve performance, with full results provided in Appendix I.

Training prior vs. Training prediction. The adjustment factor α of our method uses the ground-truth label distribution of the training dataset (i.e., $training\ prior$ or $p_{train}(y)$) as the denominator, while the numerator is based on the predicted distribution of the test set, introducing a mismatch between true and predicted quantities. To assess the impact of this discrepancy, we analyze an alternative that replaces the training prior with the model's average predicted distribution on the training dataset (i.e., $training\ prediction$ or $\widehat{p}_{train}(y)$), which requires additional inference. As shown in Table 6, the performance difference is negligible, validating the choice of the training prior.

Asymmetric cross-entropy (CE). DistPFN-T employs asymmetric cross-entropy to compute the temperature τ for temperature scaling, where the value differs depending on whether it is computed as $\text{CE}(\widehat{p}_{\text{TabPFN}}(y), p_{\text{train}})$ or $\text{CE}(p_{\text{train}}, \widehat{p}_{\text{TabPFN}}(y))$. Table 7 shows that both directions outperform TabPFN-v2, demonstrating robustness to the choice of direction.

447

448

449 450

451

452

453

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

Methods	LoCal	PFN	TabICL		TabPFN-v2		
Methods	w/o shift	w/ shift	w/o shift	w/ shift	w/o shift	w/ shift	
-	0.816	0.771	0.806	0.742	0.818	0.775	
+ EME	0.801	0.783	0.798	0.766	0.801	0.786	
+ BBE	0.805	0.787	0.802	0.768	0.805	0.789	
+ DistPFN	0.816	0.788	0.806	0.773	0.818	0.789	
+ DistPFN-T	0.816	0.791	0.806	0.777	0.818	0.792	

	Pred. distn.	w/o shift	w/ shift
TabPFN-v2	-	0.818	0.775
+ DistPFN	Single	0.818	0.789
	Multiple	0.818	0.789
+ DistPFN-T	Single	0.818	0.791
	Multiple	0.818	0.792

Table 8: Comparison with methods for label shift.

Table 9: Pred distn: single vs. multiple.

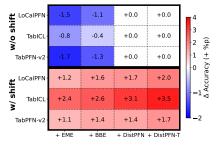
P	Methods	Shift strength (β)						
Г	Wichiods	0.0	0.1	0.5	1.0	2.0	5.0	Avg.
0.05	TabPFN-v2	0.644	0.639	0.591	0.547	0.505	0.460	0.554
	+ DistPFN	0.659	0.662	0.627	0.586	0.548	<u>0.493</u>	0.589
	+ DistPFN-T	0.662	0.667	0.640	0.601	0.562	0.504	0.605
0.10	TabPFN-v2	0.663	0.664	0.617	0.582	0.534	0.481	0.620
	+ DistPFN	0.679	0.685	0.650	0.618	0.577	0.515	0.642
	+ DistPFN-T	0.685	0.691	0.656	0.629	0.588	0.527	0.652
0.20	TabPFN-v2	0.697	0.689	0.651	0.619	0.561	0.510	0.638
	+ DistPFN	0.713	0.711	0.682	0.663	0.610	0.553	0.676
	+ DistPFN-T	0.716	0.715	0.689	0.670	0.620	0.563	0.688

	Pred. time	Avg. Acc.
LoCalPFN	0.618	0.771
+ DistPFN	0.619	0.788
+ DistPFN-T	0.619	0.791
TabICL	0.620	0.742
+ DistPFN	0.622	0.773
+ DistPFN-T	0.622	0.777
TabPFN-v2	1.002	0.775
+ DistPFN	1.003	0.789
+ DistPFN-T	1.003	0.792

Table 10: Dataset selection with K-means clustering. Our method remains effective when training subsets are formed by sampling a proportion (P) from each of the K=10 clusters, demonstrating robustness to the choice of training subsets.

Table 11: Efficiency analysis. Average prediction time (in seconds) average across 253 datasets with three different backbones.

Comparison label shift correction methods. To demonstrate the effectiveness of our methods, we compare it with other techniques handling label shift: EM-based Estimation (EME) (Saerens et al., 2002) and Black-box Estimation (**BBE**) (Lipton et al., 2018). Table 8 demonstrates that our method outperforms these approaches without requiring estimation of the test prior. In particular, as shown in Figure 11, while other methods suffer from performance degradation when no shift is present, our method maintains stable performance across both settings. Details of each method and full results are provided in Appendix E and J. Figure 11: vs. Other label shift methods.



Predicted distribution of single vs. multiple instance(s). As TabPFN allows test instances to be evaluated either individually or in batches, with both modes yielding identical predictions, DistPFN and DistPFN-T can apply their adjustment based on either the 1) prediction of a *single* instance or 2) the average prediction across *multiple* instances in the test dataset. As shown in Table 9, both choices consistently improve TabPFN-v2, averaged across six β s for w/shift. The results demonstrate that our method is robust to the choice of distribution source, with full results shown in Appendix K.

Training set selection. TabPFN suffers from quadratic complexity, making it inefficient for large training datasets (Thomas et al., 2024; Zeng et al., 2025; Qu et al., 2025), and several works mitigate this by training on *selected* subsets. One common approach is to use only the local neighbors of each test sample, as validated in Table 3 with LoCalPFN (Thomas et al., 2024). Another approach clusters the training dataset and selects the centroid and a few nearby samples per cluster. As shown in Table 10, our method remains effective across different percentages of samples per cluster (P) under K=10, where K is the number of clusters. Results for various Ks are provided in Appendix H.

Efficiency analysis. To evaluate the efficiency of our method, we compare the average prediction time (seconds) across 253 datasets using TabPFN-v2. Note that our method does not require any additional parameters, and only applies a simple multiplication of an adjustment factor to the predicted results. Table 11 summarizes the results, including the average performance across six β s, highlighting that our method achieves superior performance gains with negligible computational burden.

7 Conclusion

In this work, we introduce DistPFN, a test-time adjustment method to mitigate label shift in tabular foundation models using ICL. We further propose DistPFN-T to stabilize the adjustment via temperature scaling based on distributional divergence between training prior and predicted distribution. While our method effectively handles label shift without retraining, it does not address feature shift, which can also occur in practice. A potential direction for future work is to design foundation models that are inherently robust to both label and feature shift, beyond post-hoc adjustment. We hope that this work encourages further exploration of robustness to distribution shifts in tabular ICL.

REFERENCES

- Md Atik Ahamed and Qiang Cheng. Mambatab: A plug-and-play model for learning tabular data. In 2024 IEEE 7th International Conference on Multimedia Information Processing and Retrieval (MIPR), pp. 369–375. IEEE, 2024.
 - Naomi S Altman. An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician*, 46(3):175–185, 1992.
 - Sercan Ö Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In AAAI, 2021.
 - Kumar Arun, Garg Ishan, and Kaur Sanmeet. Loan approval prediction based on machine learning approach. *IOSR J. Comput. Eng*, 18(3):18–21, 2016.
 - Kamyar Azizzadenesheli, Anqi Liu, Fanny Yang, and Animashree Anandkumar. Regularized learning for domain adaptation under label shifts. *arXiv preprint arXiv:1903.09734*, 2019.
 - Bernd Bischl, Giuseppe Casalicchio, Matthias Feurer, Pieter Gijsbers, Frank Hutter, Michel Lang, Rafael G Mantovani, Jan N van Rijn, and Joaquin Vanschoren. Openml benchmarking suites. *arXiv preprint arXiv:1708.03731*, 2017.
 - Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794, 2016.
 - Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
 - David R Cox. The regression analysis of binary sequences. *Journal of the Royal Statistical Society: Series B (Methodological)*, 20(2):215–232, 1958.
 - Charles Elkan. The foundations of cost-sensitive learning. In *IJCAI*, volume 17, pp. 973–978. Lawrence Erlbaum Associates Ltd, 2001.
 - Benjamin Feuer, Robin T Schirrmeister, Valeriia Cherepanova, Chinmay Hegde, Frank Hutter, Micah Goldblum, Niv Cohen, and Colin White. Tunetables: Context optimization for scalable prior-data fitted networks. *NeurIPS*, 37:83430–83464, 2024.
 - Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting deep learning models for tabular data. *NeurIPS*, 34:18932–18943, 2021.
 - Yury Gorishniy, Akim Kotelnikov, and Artem Babenko. Tabm: Advancing tabular deep learning with parameter-efficient ensembling. *ICLR*, 2024a.
 - Yury Gorishniy, Ivan Rubachev, Nikolay Kartashev, Daniil Shlenskii, Akim Kotelnikov, and Artem Babenko. Tabr: Tabular deep learning meets nearest neighbors in 2023. *ICLR*, 2024b.
 - Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? *NeurIPS*, 35:507–520, 2022.
 - Simon Haykin. Neural Networks: A Comprehensive Foundation. Prentice Hall PTR, 1994.
 - Kai Helli, David Schnurr, Noah Hollmann, Samuel Müller, and Frank Hutter. Drift-resilient tabpfn: In-context learning temporal distribution shifts on tabular data. *NeurIPS*, 37:98742–98781, 2024.
 - Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. Tabpfn: A transformer that solves small tabular classification problems in a second. *ICLR*, 2023.
- Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor Schirrmeister, and Frank Hutter. Accurate predictions on small data with a tabular foundation model. *Nature*, 637(8045):319–326, 2025.
 - David Holzmüller, Léo Grinsztajn, and Ingo Steinwart. Better by default: Strong pre-tuned mlps and boosted trees on tabular data. *NeurIPS*, 37:26577–26658, 2024.

- Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data modeling using contextual embeddings. *arXiv preprint arXiv:2012.06678*, 2020.
- Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
 - Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. Well-tuned simple nets excel on tabular datasets. *NeurIPS*, 34:23928–23941, 2021.
 - 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. *NeurIPS*, 30, 2017.
 - Changhun Kim, Taewon Kim, Seungyeon Woo, June Yong Yang, and Eunho Yang. Adaptable: Test-time adaptation for tabular data via shift-aware uncertainty calibrator and label distribution handler. *arXiv preprint arXiv:2407.10784*, 2024.
 - Andy Liaw and Matthew Wiener. Classification and regression by randomforest. *R News*, 2(3):18–22, 2002. URL https://CRAN.R-project.org/doc/Rnews/.
 - Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift with black box predictors. In *ICML*, pp. 3122–3130. PMLR, 2018.
 - Duncan McElfresh, Sujay Khandagale, Jonathan Valverde, Vishak Prasad C, Ganesh Ramakrishnan, Micah Goldblum, and Colin White. When do neural nets outperform boosted trees on tabular data? *NeurIPS*, 36:76336–76369, 2023.
 - Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. Catboost: unbiased boosting with categorical features. *NeurIPS*, 31, 2018.
 - Jingang Qu, David HolzmAvZller, GaAl Varoquaux, and Marine Le Morvan. Tabicl: A tabular foundation model for in-context learning on large data. *ICML*, 2025.
 - Marco Saerens, Patrice Latinne, and Christine Decaestecker. Adjusting the outputs of a classifier to new a priori probabilities: a simple procedure. *Neural computation*, 14(1):21–41, 2002.
 - Ravid Shwartz-Ziv and Amitai Armon. Tabular data: Deep learning is not all you need. *Information Fusion*, 81:84–90, 2022.
 - Gowthami Somepalli, Micah Goldblum, Avi Schwarzschild, C Bayan Bruss, and Tom Goldstein. Saint: Improved neural networks for tabular data via row attention and contrastive pre-training. *arXiv preprint arXiv:2106.01342*, 2021.
 - Anton Frederik Thielmann and Soheila Samiee. On the efficiency of nlp-inspired methods for tabular deep learning. *arXiv preprint arXiv:2411.17207*, 2024.
 - Valentin Thomas, Junwei Ma, Rasa Hosseinzadeh, Keyvan Golestan, Guangwei Yu, Maks Volkovs, and Anthony L Caterini. Retrieval & fine-tuning for in-context tabular models. *NeurIPS*, 37: 108439–108467, 2024.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017.
 - Derek Xu, Olcay Cirit, Reza Asadi, Yizhou Sun, and Wei Wang. Mixture of in-context prompters for tabular pfns. *ICLR*, 2025.
 - Han-Jia Ye, Huai-Hong Yin, De-Chuan Zhan, and Wei-Lun Chao. Revisiting nearest neighbor for tabular data: A deep tabular baseline two decades later. *ICLR*, 2024.
 - Yuchen Zeng, Tuan Dinh, Wonjun Kang, and Andreas C Mueller. Tabflex: Scaling tabular learning to millions with linear attention. *ICML*, 2025.

L Other Metrics

APPENDIX **A Experimental Setups B** Baseline Methods **C** Baseline Implementations **D** Theoretical Justification **E** Classical Methods for Label Shift Correction E.3 **Hyperparameter Tuning for Stronger Baselines G** Dataset Statistics **H** K-means Clustering for Dataset Selection **Application to LoCalPFN** Ι Comparison with Methods for Label Shift Correction

K Predicted Distribution of Single vs. Multiple Instances

A EXPERIMENTAL SETUPS

Experimental setups. We use the official implementation of TabPFN³ and adopt all default settings without modification. This includes architectural choices such as the number of layers and hidden dimensions, where we use 12 transformer layers, each with a hidden size of 192 and 6 attention heads. The feedforward layer dimension is implicitly set to 768 via a hidden factor of 4. For inference, we load the pretrained weights from TabPFN-v2⁴ available on Hugging Face.

Dataset. We evaluate on 250+ tabular datasets from OpenML (Bischl et al., 2017). The dataset list is retrieved from the benchmark configuration provided in this repository⁵, which is built on top of the official TabPFN evaluation setup. Dataset statistics are summarized in Appendix G.

B BASELINE METHODS

We categorize 15 baseline tabular models into three groups:

• ML models (7): Logistic Regression (LR), Support Vector Machines (SVM), Random Forest (Liaw & Wiener, 2002), k-nearest neighbors (kNN), Multi-layer Perceptrons (MLP), LightGBM (Ke et al., 2017), CatBoost (Prokhorenkova et al., 2018)

• **DL** (**non-foundation**) **models** (**5**): FT-Transformer (Gorishniy et al., 2021), TabM (Gorishniy et al., 2024a), TabulaRNN (Thielmann & Samiee, 2024), MambaTab (Ahamed & Cheng, 2024), RealMLP (Holzmüller et al., 2024)

• DL (foundation) models based on ICL (3): TabPFN-v2 (Hollmann et al., 2025), LoCalPFN (Thomas et al., 2024), TabICL (Qu et al., 2025)

Details of each method are provided below.

B.1 MACHINE LEARNING (ML) MODELS

• Logistic Regression (LR) (Cox, 1958): A simple linear model commonly used for binary and multiclass classification tasks in tabular data.

• Support Vector Machine (SVM) (Cortes & Vapnik, 1995): A kernel-based classifier that aims to find the optimal decision boundary with maximum margin between classes.

 • Multilayer Perceptron (MLP) (Haykin, 1994): A feedforward neural network consisting of multiple fully connected layers with non-linear activations, trained via backpropagation.

 k-Nearest Neighbors (kNN) (Altman, 1992): A non-parametric method that classifies a sample based on the majority class among its k nearest neighbors in the feature space.
 Random Forest (Liaw & Wiener, 2002): An ensemble learning method based on bagging over

 decision trees, which improves robustness and generalization.
LightGBM (Ke et al., 2017): A fast and efficient GBDT model using histogram-based algorithms and leaf-wise tree growth.

 • CatBoost (Prokhorenkova et al., 2018): A GBDT model that handles categorical features efficiently and mitigates prediction shift via ordered boosting.

³https://github.com/PriorLabs/TabPFN

⁴https://huggingface.co/Prior-Labs/TabPFN-v2-clf

⁵https://github.com/carteakey/tabpfn-eval/blob/main/src/data/openml_ list.csv

704 705

706

708 709 710

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

744 745 746

747 748 749

750 751

752 754

755

https://auto.gluon.ai/

B.2 DEEP LEARNING (NON-FOUNDATION) MODELS

- FT-Transformer (Gorishniy et al., 2021): A transformer-based architecture tailored for tabular data, providing a simple yet powerful baseline that outperforms many prior DL models on classification and regression tasks.
- TabM (Gorishniy et al., 2024a): An MLP-based model that leverages an efficient ensemble mechanism to approximate deep ensembles, enabling multiple predictions per instance while maintaining computational efficiency.
- TabulaRNN (Thielmann & Samiee, 2024): An RNN-inspired architecture for tabular data that emphasizes efficiency, addressing limitations of NLP-style models in terms of scalability and training cost.
- MambaTab (Ahamed & Cheng, 2024): A scalable and efficient model built on structured state-space models (SSMs), capturing long-range dependencies with fewer parameters while maintaining strong predictive performance.
- RealMLP (Holzmüller et al., 2024): An enhanced MLP variant with meta-tuned hyperparameters, achieving competitive accuracy-efficiency trade-offs compared to gradient boosting methods in tabular benchmarks.

B.3 DEEP LEARNING (FOUNDATION) MODELS BASED ON ICL

- LoCalPFN (Thomas et al., 2024): A lightweight PFN variant that reduces computational cost by leveraging local task priors and architectural simplifications.
- TabICL (Ou et al., 2025): A two-stage model that first applies column attention to capture feature dependencies and then row attention to encode sample interactions.
- TabPFN-v2 (Hollmann et al., 2025): A state-of-the-art foundation model for tabular classification that leverages a pretrained transformer for zero-shot prediction on small datasets.

BASELINE IMPLEMENTATIONS

The baseline results are obtained from the following publicly available repositories:

- [1] **TabPFN Evaluation** framework⁶ was used to evaluate all ML models, as well as the foundation models *TabPFN* and *LoCalPFN*. Since *LoCalPFN* does not have an official implementation, we reimplemented it based on the TabPFN codebase.
- [2] AutoGluon v1.4.07 was used to benchmark *TabICL* and several non-foundation models such as FT-Transformer, TabM, and RealMLP.
- [3] Mambular⁸ provided implementations for additional non-foundation models including MambaTab, TabulaRNN, FT-Transformer, and TabM.

For models implemented in both [2] and [3] (e.g., FT-Transformer and TabM), we use the [2] versions as they yield stronger performance for a stronger baseline.

6https://github.com/carteakey/tabpfn-eval

⁸https://github.com/OpenTabular/DeepTabular

D THEORETICAL JUSTIFICATION

 We provide theoretical grounding for our posterior adjustment to clarify that DistPFN is not merely a heuristic trick, but a principled approximation derived from existing theory. We present two complementary perspectives: 1) connection to classical label shift correction as a plug-in reweighting (Section D.1) and 2) Bayesian view that replaces the mismatched prior with a self-consistent estimate from model predictions (Section D.2).

D.1 RELATION TO LABEL SHIFT CORRECTION

The label shift setting assumes that the conditional distribution p(x|y) remains invariant while the marginal priors differ:

$$p_{\text{train}}(y) \neq p_{\text{test}}(y), \quad p(x|y) \text{ is fixed.}$$

Under this assumption, the Bayes-optimal posterior is given by

$$p_{ ext{test}}(y|x) \propto rac{p_{ ext{train}}(y|x)}{p_{ ext{train}}(y)} \, p_{ ext{test}}(y).$$

Classical approaches such as EM-based reweighting (Saerens et al., 2002; Lipton et al., 2018) estimate $p_{\text{test}}(y)$ explicitly by matching marginal predictions to unlabeled test data. DistPFN instead uses the predictive marginal $\hat{p}(y)$ obtained directly from the model, and constructs the adjustment factor

$$\alpha(y) = \frac{\hat{p}(y)}{p_{\text{train}}(y)}.$$

This yields the corrected posterior

$$\hat{p}(y|x) \propto \frac{p_{\mathrm{train}}(y|x)}{p_{\mathrm{train}}(y)}\,\hat{p}(y),$$

which can be seen as a plug-in realization of the classical correction rule, avoiding iterative estimation while remaining theoretically consistent with label shift correction.

D.2 BAYESIAN INTERPRETATION

From a Bayesian perspective, TabPFN models the posterior under the training distribution:

$$p_{\text{train}}(y|x) \propto p(x|y) p_{\text{train}}(y)$$
.

At test time, however, the desired posterior is

$$p_{\text{test}}(y|x) \propto p(x|y) p_{\text{test}}(y)$$
.

The difference comes solely from the prior. DistPFN addresses this gap by substituting $p_{\text{train}}(y)$ with $\hat{p}(y)$, the average predictive distribution obtained on the test set:

$$\hat{p}_{ ext{DistPFN}}(y|x) \propto rac{p_{ ext{train}}(y|x)}{p_{ ext{train}}(y)}\,\hat{p}(y).$$

This interpretation shows that DistPFN is not an ad-hoc adjustment but a Bayesian posterior correction where the unknown test prior is approximated in a self-consistent manner from model outputs. The method therefore inherits a principled justification while retaining the efficiency of a simple, training-free plug-in procedure.

E CLASSICAL METHODS FOR LABEL SHIFT CORRECTION

In this section, we summarize three representative approaches for handling label shift. All of these methods directly adjust classifier outputs, but they differ in how the test prior π_{test} is obtained.

E.1 PRIOR-RATIO ADJUSTMENT

 Prior-ratio Adjustment (Elkan, 2001) introduces a simple correction under changing class priors in the binary setting. The method assumes that the new prior π_{test} is available from external knowledge or domain statistics. Given a posterior $p_{\text{train}}(1|x)$ trained under π_{train} , the corrected posterior is

$$p_{\text{test}}(1|x) \ = \ \frac{p_{\text{train}}(1|x) \cdot \frac{\pi_{\text{test}}(1)}{\pi_{\text{train}}(1)}}{p_{\text{train}}(1|x) \cdot \frac{\pi_{\text{test}}(1)}{\pi_{\text{train}}(1)} + (1 - p_{\text{train}}(1|x)) \cdot \frac{1 - \pi_{\text{test}}(1)}{1 - \pi_{\text{train}}(1)}}.$$

This approach directly modifies posterior probabilities by scaling them with prior ratios. The same principle naturally extends to the multiclass case by applying the ratio $\pi_{\text{test}}(y)/\pi_{\text{train}}(y)$ to each class posterior.

E.2 EM-BASED ESTIMATION

EM-based Estimation (Saerens et al., 2002) proposes an iterative procedure to estimate unknown test priors when they are not directly given. At iteration t, the posterior is updated by

$$p^{(t+1)}(y|x) \propto p_{\text{train}}(y|x) \cdot \frac{\pi_{\text{test}}^{(t)}(y)}{\pi_{\text{train}}(y)}.$$

The updated posteriors provide a new estimate of π_{test} by averaging across the test set. Repeating this E-step and M-step allows the estimated test prior to gradually converge. The final corrected posterior then follows the standard prior-ratio adjustment, but with π_{test} estimated rather than assumed.

E.3 BLACK-BOX ESTIMATION

Black-box Estimation (Lipton et al., 2018) employs a validation dataset with true labels to construct a confusion matrix $C(s|y) = P(\hat{y} = s \mid y)$ that characterizes prediction errors of the classifier. On an unlabeled test set, it collects predicted labels to obtain the empirical distribution $p_{\text{test}}(s)$. These quantities are related through the equation

$$p_{\text{test}}(s) \; pprox \; \sum_{y} C(s|y) \, \pi_{\text{test}}(y).$$

By solving this linear system, the method estimates the test prior π_{test} . Once the test prior is recovered, the posterior correction is applied using the prior ratio:

$$p_{ ext{test}}(y|x) \, \propto \, p_{ ext{train}}(y|x) \cdot rac{\pi_{ ext{test}}(y)}{\pi_{ ext{train}}(y)}.$$

This approach is considered black-box as it does not require access to classifier internals, only its predicted outputs and a validation set to estimate the confusion matrix.

F HYPERPARAMETER TUNING FOR STRONGER BASELINES

To ensure strong and fair baselines, we perform hyperparameter tuning for each conventional ML method using the search space provided in a public implementation⁹. The search spaces are manually designed to cover commonly used ranges for each model class, including both optimization-related parameters (e.g., learning rate, max iterations) and regularization or structural options (e.g., penalty, tree depth, number of neighbors). We conduct random search over these spaces and tune the models on validation datasets that are kept separate from the final test splits. The details of the hyperparameter search spaces are provided in Table F.1.

Model	Hyperparameter	Type	Log-scale	Range
	max_iter	int	no	{50, 100, 200, 500, 1000}
	solver	categorical	no	{newton-cg, lbfgs, liblinear, sag, saga}
Logistic Regression	fit_intercept	boolean	no	{True, False}
	penalty	categorical	no	{11, 12, elasticnet, none}
	С	float	no	{0.1, 1.0, 10.0, 100.0}
	n_estimators	int	no	{10, 50, 100, 200, 500}
	criterion	categorical	no	{gini, entropy}
D 1 E	probability	boolean	no	{True}
Random Forest	max_depth	int / None int	no	{None, 10, 50, 100, 200}
	min_samples_split min_samples_leaf	int	no	$\{2, 5, 10\}$ $\{1, 2, 4\}$
	max_features	categorical	no no	{1, 2, 4} {auto, sqrt, log2}
	C	float	no	{0.1, 1.0, 10.0, 100.0}
SVM	kernel	categorical boolean	no	{linear, poly, rbf, sigmoid}
S V IVI	probability degree	int	no no	{True} {2, 3, 4, 5}
	gamma	categorical	no	{scale, auto}
	_			
	max_iter	int	no	{50, 100, 200, 500, 1000}
MLP	activation	categorical	no	{identity, logistic, tanh, relu}
	solver	categorical	no	{lbfgs, sgd, adam}
	alpha	float	no	{0.0001, 0.001, 0.01, 0.1}
	learning_rate	categorical float	no no	{constant, invscaling, adaptive}
	learning_rate_init			{0.001, 0.01, 0.1}
	n_neighbors	int	no	$\{3, 5, 11, 19\}$
1-NINI	weights	categorical	no	{uniform, distance}
kNN	algorithm leaf_size	categorical int	no	{auto, ball_tree, kd_tree, brute}
		int	no no	{30, 50, 100} {1, 2}
	p			
	n_estimators	int	no	{50, 100, 200}
	max_depth	int	no	{6, 10, 15, 20}
VCD	learning_rate	float	no	$\{0.001, 0.01, 0.1\}$
XGBoost	subsample	float	no	{0.5, 0.6, 0.7, 0.8, 0.9, 1.0}
	colsample_bytree	float	no	{0.4, 0.5,, 1.0}
	colsample_bylevel	float	no	{0.4, 0.5,, 1.0}
	n_estimators	int	no	{50, 100, 200}
Li-14CDM	max_depth	int	no	{6, 10, 15, 20}
LightGBM	learning_rate	float	no	$\{0.001, 0.01, 0.1\}$
	num_leaves	int	no	{31, 60, 120, 240, 480, 960}
	min_child_samples	int	no	{10, 20, 30, 40, 50}
	iterations	int	no	{50, 100, 200}
CatBoost	depth	int	no	{6, 8, 10}
	learning_rate	float	no	$\{0.001, 0.01, 0.1\}$
	12_leaf_reg	float	no	{1, 3, 5, 7, 9}

Table F.1: **Hyperparameter search spaces for each conventional ML baseline.** All hyperparameter values are tuned via random search over manually defined discrete sets.

⁹https://github.com/carteakey/tabpfn-eval

G DATASET STATISTICS

We evaluate on 253 tabular datasets from OpenML (Bischl et al., 2017). Summary statistics for all datasets are provided in Table G.1, G.2, and G.3. Each dataset is described using the following attributes: the dataset name (Name), the total number of input features (#Features), the number of categorical features among them (#Cat. Feat.), the number of data instances (#Instances), the number of class labels (#Classes), the number of missing values (#NaNs), and the number of samples belonging to the smallest class (Minority Class Size).

Name	#Features	#Cat. Feat.	#Instances	#Classes	#NaNs	Minority Class S
oollen	6	1	3848	2	0	1924
Sick_numeric	30	1	3772	2	0	231
ungle_chess_2pcs_endgame_rat_rat	47	27	3660	2	0	1605
ICI.churn	21	1	3333	2	0	483
ed24 ed7	25 8	25 8	3200 3200	10 10	0	296 270
curvs-kp	37	37	3196	2	0	1527
plice	61	61	3190	3	0	767
pace_ga	7	1	3107	2	0	1541
tackOverflow-polarity-train	2	i	3097	3	0	842
eismic-bumps	19	5	2584	2	0	170
zone-level-8hr	73	1	2534	2	0	160
ingle_chess_2pcs_endgame_lion_lion	47	27	2352	2	0	949
ingle_chess_2pcs_endgame_elephant_elephant	47	27	2351	2	0	1035
egment	20	1	2310	7	0	330
itanic	4	1	2201	2	0	711
uake	4	1	2178	2	0	969
cl cl	22	1	2109	2	0	326
alloon	2	1	2001	2	0	482
nfeat-fourier	77	1	2000	10	0	200
zone-level-8hr_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	73	1	2000	2	0	126
ıfeat-karhunen	65	1	2000	10	0	200
nnis_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	55	1	2000	2	0	1000
overtype_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	55	45	2000	2	0	1000
rst-order-theorem-proving_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	52	1	2000	6	0	159
finiBooNE_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	51	1	2000	2	0	1000
DDCup09_upselling_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	50	16	2000	2	0	1000
da_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	49	1	2000	2	0	496
feat-zernike	48	1	2000	10	0	200
onnect-4_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	43	43	2000	3	0	191
r-vs-kp_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	37	37	2000	2	0	956
oad-safety_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	33	4	2000	2	0	1000
esturePhaseSegmentationProcessed_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	33	1	2000	5	0	202
hishingWebsites_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	31	31	2000	2	0	886
ol_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	27	1	2000	2	0	1000
liggs_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	25	1	2000	2	0	1000
ye_movements_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	24	4	2000	2	0	1000
umerai28.6_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	22	1	2000	2	0	990
c1_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	22	1	2000	2	0	309
dd_ipums_la_97-small_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	21	1	2000	2	0	1000
hurn_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	21	5	2000	2	0	283
ompass_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	18	10	2000	2	0	1000
ouse_16H_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	17	1	2000	2	0	1000
egment_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	17	1	2000	7	0	285
dult_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	15	9	2000	2	242	479
dult_seed_1_nrows_2000_nclasses_10_ncols_100_stratify_True	15	9	2000	2	248	479
dult_seed_2_nrows_2000_nclasses_10_ncols_100_stratify_True	15	9	2000	2	279	479
dult_seed_3_nrows_2000_nclasses_10_ncols_100_stratify_True	15	9	2000	2	254	479
dult_seed_4_nrows_2000_nclasses_10_ncols_100_stratify_True	15	9	2000	2	253	479
_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	13	8	2000	2	0	1000
vine_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	12	1	2000	2	0	1000
lick_prediction_small_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	12	7	2000	2	0	337
.mazon_employee_access_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	10	10	2000	2	0	116
alifornia_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	9	1	2000	2	0	1000
f-police-incidents_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	9	6	2000	2	0	243
lectricity_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	8	1	2000	2	0	1000
irlines_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	8	5	2000	2	0	891
feat-morphological	7	1	2000	10	0	200
ingle_chess_2pcs_raw_endgame_complete_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	7	1	2000 2000	3 2	0	194
honeme_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	6	1 1		2 2		1000
rilt_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	6 34	1	2000 1941		0	108
eel-plates-fault eel-plates-fault_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	28	1	1941	2 7	0	673 55
	28	21		2	0	800
AMETES_Epistasis_2-Way_20atts_0.1H_EDM-1_1	38	1	1600 1563	2	0	160
nc	10	8	1473	3	0	333
nc_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	10	8	1473	3	0	333
me_seed_0_mows_2000_ncrasses_10_ncors_100_stratify_frue m-employee-performance	34	1	1473	2	0	226
24	38	i	1458	2	0	178
24_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	38	i	1458	2	0	178
anknote-authentication	5	1	1372	2	0	610
nalcatdata_halloffame	17	2	1340	2	20	125
ofn-3-7-10	11	11	1324	2	0	292
ocmob	6	5	1156	2	0	256
arity5_plus_5	11	11	1124	2	0	557
ieChart3	38	1	1077	2	0	134
sar-biodeg	42	1	1055	2	0	356
sar-biodeg_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True	42	i	1055	2	0	356
izzaCutter3	38	i	1043	2	0	127
nftsa_sleepdata	3	i	1024	4	0	94
ntsa.sieepdata redit-g	21	14	1024	2	0	300
ummy	7	14	1000	2	0	273
ummy d6	10	10	973	2	0	322
oo kyol	45	3	973	2	0	346
c-tac-toe	10	10	959	2	0	332
	7			2		
our-and-Travels-Customer-Churn-Prediction		5	954		60	224 462
tock	10	1 1	950 846	2 4	0	462 199
					()	199
ehicle ehicle_reproduced	19 19	i	846	4	0	199

Table G.1: Dataset statistics - Part 1

Name	#Features	#Cat. Feat.	#Instances	#Classes	#NaNs	Minority Class Siz
analcatdata_dmft diabetes	5 9	5 1	797 768	6 2	0	123 268
blood-transfusion-service-center	5	1	748	2	0	178
blood-transfusion-service-center_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True doa_bwin_balanced	5 14	1 3	748 708	2 2	0	178 354
PieChart1	38	1	705	2	0	61
breast-w credit-approval	10 16	1 10	699 690	2 2	16 67	241 307
credit-approval_reproduced	16	10	690	2	67	307
Australian	15	9	690	2 2	0	307
Australian_seed_0_nrows_2000_nclasses_10_ncols_100_stratify_True disclosure_x_bias	15 4	1	690 662	2	0	307 317
disclosure_x_tampered	4	1	662	2	0	327
disclosure_x_noise disclosure_z	4 4	1	662 662	2 2	0	329 314
PizzaCutter1	38	1	661	2	0	52
balance-scale monks-problems-2	5 7	1 7	625 601	3 2	0	49 206
synthetic_control	61	1	600	6	0	100
sensory wdbc	12 31	12 1	576 569	2 2	0	239 212
arsenic-female-bladder	5	2	559	2	0	80
monks-problems-1	7	7	556	2	0	278
monks-problems-3 climate-model-simulation-crashes	7 21	7 1	554 540	2 2	0	266 46
doa_bwin	14	3	530	2	0	176
CPMP-2015-runtime-classification	23	1	527	4	0	78
kc2 threeOf9	22 10	1 10	522 512	2 2	0	107 238
rmftsaladata	11	1	508	2	0	222
boston_corrected	21	4 2	506	2 2	0	223
boston collins	14 23	3	506 500	2	0	209 80
pm10	8	1	500	2	0	246
no2 LED-display-domain-7digit	8	1 1	500 500	2 10	0	249 37
irish	6	4	500	2	32	222
PopularKids analcatdata_apnea2	11 4	5 3	478 475	3 2	0	90 64
analcatdata_apnea1	4	3	475	2	0	61
thoracic-surgery	17 4	14 2	470	2 2	0	70 208
analcatdata_vineyard chscase_vine2	3	1	468 468	2	0	212
sa-heart	10	2	462	2	0	160
analcatdata_apnea3 wholesale-customers	4 9	3 2	450 440	2 2	0	55 142
mw1	38	1	403	2	0	31
user-knowledge chscase_census5	6 8	1 1	403 400	5 2	0	24 193
chscase_census4	8	1	400	2	0	194
chscase_census3 chscase_census2	8	1 1	400 400	2 2	0	192 197
chscase_census6	7	1	400	2	0	165
analcatdata_germangss	6	5	400	4	0	100
calendarDOW autoMpg	33 8	21 4	399 398	5 2	0 6	44 189
vinnie	3	1	380	2	0	185
jEdit_4.2_4.3 dermatology	9 35	1 34	369 366	2 6	0 8	165 20
analcatdata_draft	5	3	366	2	1	32
analcatdata_birthday	4	3	365	2	30	53
ionosphere SPECTF	35 45	1 1	351 349	2 2	0	126 95
penguins	7	3	344	3	18	68
CastMetal1 visualizing_galaxy	38 5	1 1	327 323	2 2	0	42 148
plasma_retinol	14	4	315	2	0	133
solar-flare	13	13	315	5 9	0	21
diggle_table_a2 vertebra-column	7	1 1	310 310	3	0	18 60
haberman	4	2	306	2	0	81
heart-c cleveland	14 14	8 8	303 303	2 2	7 6	138 139
cholesterol	14	8	303	2	6	137
cleve cleveland-nominal	14 8	9 8	303 303	2 5	0	138 13
CostaMadre1	38	1	296	2	0	38
Heart_disease_prediction_20	14	1	296	2	0	137
breast-cancer breastTumor	10 10	10 9	286 286	2 2	9	85 120
analcatdata_broadwaymult	8	5	285	7	27	21
mu284 DiabeticMellitus	11 98	1 1	284 281	2 2	0 2	142 99
breast-cancer-dropped-missing-attributes-values	10	10	277	2	0	81
jEdit_4.0_4.2 heart-statlog	9 14	1 1	274 270	2 2	0	134 120
SPECT	23	23	267	2	0	55
Touch2	11	1	265	8	0	27
analcatdata_lawsuit	5	2	264	2	0	19

Table G.2: Dataset statistics - Part 2

Name	#Features	#Cat. Feat.	#Instances	#Classes	#NaNs	Minority Class Size
MegaWatt1	38	1	253	2	0	27
bodyfat	15 7	1	252 250	2 2	0	124 107
qualitative-bankruptcy prnn_synth	3	7 1	250	2	0	125
conference_attendance	7	7	246	2	Ö	31
chatfield_4	13	1	235	2	0	93
chscase_whale lungcancer_GSE31210	9 24	1 3	228 226	2 2	20 0	111 35
chscase_geyser1	3	1	222	2	0	88
thyroid-new	6	1	215	3	0	30
glass	10	1	214	6	0	9
prnn_fglass seeds	10 8	1 1	214 210	2 3	0	76 70
biomed	9	2	209	2	15	75
cpu	8	2	209	2	0	53
machine_cpu sonar	7 61	1 1	209 208	2 2	0	56 97
regime_alimentaire	20	17	202	2	17	41
heart-long-beach	14	1	200	5	0	10
pwLinear	11 8	1 2	200 200	2 2	0	97 100
prnn_crabs parkinsons	23	1	195	2	0	48
pharynx	11	10	195	2	2	74
KnuggetChase3	40	1	194	2	0	36
wisconsin lowbwt	33 10	1 8	194 189	2 2	0	90 90
triazines	61	1	186	2	0	77
chscase_funds	3	1	185	2	0	87
planning-relax Smartphone-Based_Recognition_of_Human_Activition	13 es 68	1 2	182 180	2 6	0	52 30
backache	es 68 32	27	180	2	0	30 25
wine	14	1	178	3	0	48
servo	5	5	167	2	0	38
robot-failures-lp5 analcatdata_wildcat	91 6	1 3	164 163	5 2	0	21 47
mc2	40	1	161	2	0	52
corral	7	7	160	2	0	70
hayes-roth auto_price	5 16	1 2	160 159	3 2	0	31 54
auto-price autoPrice	16	1	159	2	0	54 54
analcatdata_gsssexsurvey	10	6	159	2	6	35
TuningSVMs	81 9	1 7	156 155	2 4	0	54 19
grub-damage teachingAssistant	7	5	155	3	0	49
tae	6	3	151	3	0	49
iris	5	1	150	3 3	0	50 50
iris-example sleuth_case2002	5 7	1 5	150 147	2	0	50 69
kc1-top5	95	1	145	2	0	8
kc1-binary	95	1	145	2	0	60
newton_hema veteran	4 8	2 5	140 137	2 2	0	70 43
analcatdata_boxing2	4	4	137	2	0	61
analcatdata_seropositive	4	2	132	2	0	46
transplant	4 9	1 1	131 130	2 2	0	48 11
datatrieve visualizing_livestock	3	2	130	5	0	26
humandevel	2	1	130	2	0	65
mux6	7	7	128	2	0	64
MindCave2 fruitfly	40 5	1 3	125 125	2 2	0 0	44 49
KungChi3	40	1	123	2	0	16
heart-switzerland	13	1	123	5	0	5
arl	30	1	121	2	0	9
analcatdata_boxing1 rabe_266	4 3	4 1	120 120	2 2	0 0	42 57
robot-failures-lp4	91	1	117	3	0	21
visualizing_environmental	4	1	111	2	0	53
cloud	8	2	108	2	0	32
analcatdata_michiganacc ar4	4 30	3 1	108 107	2 2	0 0	48 20
molecular-biology_promoters	58	58	106	2	0	53
breast-tissue	10	1	106	6	0	14
ar6	30	1	101	2	0	15
zoo fertility	17 10	16 1	101 100	7 2	0 0	4 12
analcatdata_creditscore	7	4	100	2	0	27
blogger	6	6	100	2	0	32
analcatdata_chlamydia	4	4	100	2	0	19
analcatdata_neavote	3	2	100	2	0	7

Table G.3: Dataset statistics - Part 3

H K-MEANS CLUSTERING FOR DATASET SELECTION

Table H.1 reports the extended results of our K-means clustering-based training set selection under different numbers of clusters $K \in \{3,5,10\}$, where a proportion $(P \in \{0.05,0.10,0.20\})$ of samples is drawn from each cluster. Across all settings, our method demonstrates stable performance regardless of K, confirming its robustness when applied with clustering-based selection.

	D	M.d. 1	/ 1°C			Shif	t strengtl	n (β)		
K	P	Methods	w/o shift	0.0	0.1	0.5	1.0	2.0	5.0	Avg.
	0.05	TabPFN-v2 DistPFN DistPFN-T	0.668 0.661 0.657	0.596 0.622 0.625	0.591 0.614 0.616	0.548 0.579 0.588	0.500 0.532 0.540	0.439 0.454 0.459	0.408 0.428 0.433	0.513 0.538 0.543
3	0.10	TabPFN-v2 DistPFN DistPFN-T	0.699 0.692 0.687	0.626 0.641 0.643	0.632 0.653 0.657	0.570 0.620 0.628	0.528 0.564 0.569	0.465 0.498 0.504	0.424 0.450 0.454	0.541 0.573 0.584
		TabPFN-v2 DistPFN DistPFN-T	0.732 0.727 0.722	0.673 0.692 0.691	0.668 0.688 0.692	0.626 0.669 0.674	0.576 0.614 0.620	0.505 0.556 0.568	0.468 0.509 0.516	0.591 0.639 0.661
	0.05	TabPFN-v2 DistPFN DistPFN-T	0.676 0.673 0.672	0.605 0.628 0.629	0.606 0.630 0.634	0.550 0.587 0.594	0.503 0.534 0.539	0.459 0.487 0.493	0.429 0.453 0.460	0.529 0.561 0.565
5		TabPFN-v2 DistPFN DistPFN-T	0.699 0.696 0.693	0.631 0.654 0.655	0.644 0.665 0.670	0.586 0.624 0.630	0.540 0.583 0.591	0.485 0.528 0.538	0.446 0.475 0.483	0.569 0.609 0.620
	0.20	TabPFN-v2 DistPFN DistPFN-T	0.732 0.736 0.731	0.670 0.687 0.690	0.679 0.697 0.698	0.628 0.662 0.670	0.582 0.625 0.631	0.523 0.576 0.584	0.481 0.521 0.531	0.618 0.645 0.667
	0.05	TabPFN-v2 DistPFN DistPFN-T	0.708 0.706 0.701	0.644 0.659 0.662	0.639 0.662 0.667	0.591 0.627 0.640	0.547 0.586 0.601	0.505 0.548 0.562	0.460 0.493 0.504	0.554 0.589 0.605
10	0.10	TabPFN-v2 DistPFN DistPFN-T	0.727 0.723 0.718	0.663 0.679 0.685	0.664 0.685 0.691	0.617 0.650 0.656	0.582 0.618 0.629	0.534 0.577 0.588	0.481 0.515 0.527	0.620 0.642 0.652
	0.20	TabPFN-v2 DistPFN DistPFN-T	0.749 0.749 0.748	0.697 0.713 0.716	0.689 0.711 0.715	0.651 0.682 0.689	0.619 0.663 0.670	0.561 0.610 0.620	0.510 0.553 0.563	0.638 0.676 0.688

Table H.1: **K-means-based training dataset selection.** Our method remains effective when training subsets are selected by clustering the data and sampling a percentage (P) of samples from each of K clusters.

I APPLICATION TO LOCALPFN

Table I.1 provides the full results for LoCalPFN under different values of k across $\sin \beta$ values. The results confirm that our methods yield consistent improvements regardless of the choice of k, demonstrating robustness of the approach.

k	Methods	Shift strength (β)							
κ	Methods	0.0	0.1	0.5	1.0 2.0 5.0 4 0.758 0.711 0.67 2 0.786 0.772 0.75 4 0.790 0.779 0.75 5 0.775 0.744 0.71 3 0.790 0.777 0.76 5 0.778 0.752 0.72 6 0.778 0.752 0.72 3 0.791 0.779 0.76 6 0.794 0.785 0.77 8 0.778 0.753 0.71 4 0.790 0.782 0.772	5.0	Avg.		
3	LoCalPFN + DistPFN + DistPFN-T	0.789 0.794 0.794	0.787 0.794 0.794	0.774 0.792 0.794	<u>0.786</u>	0.772	0.679 0.752 0.759	0.750 0.782 0.785	
5	LoCalPFN + DistPFN + DistPFN-T	0.792 0.794 0.795	0.791 0.795 0.796	0.785 0.793 0.795	0.790	0.777	0.714 0.766 0.770	0.767 0.786 0.789	
10	LoCalPFN + DistPFN + DistPFN-T	0.794 0.796 0.797	0.792 0.795 0.797	0.786 0.793 0.796	0.791	0.779	0.720 0.768 0.774	0.770 0.787 0.789	
20	LoCalPFN + DistPFN + DistPFN-T	0.794 0.797 0.798	0.793 0.796 0.797	0.788 0.794 0.796	0.790	0.782	0.719 0.770 0.776	0.771 0.788 0.791	

Table I.1: **Application to LoCalPFN.** DistPFN and DistPFN-T applied to LoCalPFN show consistent improvements across varying numbers of neighbors (k).

J COMPARISON WITH METHODS FOR LABEL SHIFT CORRECTION

To demonstrate the effectiveness of our approach, we compare it with classical methods for handling label shift by rescaling classifier outputs, which typically require estimating the test distribution: EM-based Estimation (EME) (Saerens et al., 2002) and Black-box Estimation (BBE) (Lipton et al., 2018). Table J.1 presents the results, showing that our method is effective without requiring estimation of the test prior.

Methods	w/o shift	Shift strength (β)							
Methods	W/O SHIIT	0.0	0.1	0.5	1.0	2.0	5.0	Avg.	
LoCalPFN	0.816	0.794	0.793	0.788	0.778	0.753	0.719	0.771	
+ EME	0.801	0.792	0.790	0.786	0.785	0.778	0.769	0.783	
+ BBE	0.805	0.798	0.795	0.792	0.789	0.782	0.770	0.787	
+ DistPFN	0.816	0.797	0.796	0.794	0.790	0.782	0.770	0.788	
+ DistPFN-T	0.816	0.798	0.797	0.796	0.794	0.787	0.776	0.791	
TabICL	0.806	0.783	0.781	0.770	0.747	0.704	0.664	0.742	
+ EME	0.798	0.776	0.776	0.770	0.769	0.761	0.747	0.766	
+ BBE	0.802	0.783	0.785	0.780	0.774	0.754	0.734	0.768	
+ DistPFN	0.806	0.786	0.786	0.781	0.776	0.763	0.746	0.773	
+ DistPFN-T	0.806	0.786	0.786	0.783	0.780	0.771	0.755	0.777	
TabPFN-v2	0.818	0.797	0.796	0.790	0.782	0.759	0.727	0.775	
+ EME	0.801	0.793	0.793	0.790	0.787	0.783	0.768	0.786	
+ BBE	0.805	0.799	0.797	0.797	0.791	0.783	0.768	0.789	
+ DistPFN	0.818	0.799	0.797	0.795	0.791	0.783	0.769	0.789	
+ DistPFN-T	0.818	0.799	0.798	0.797	0.796	0.789	0.775	0.792	

Figure J.1: Comparison with other label shift methods.

K PREDICTED DISTRIBUTION OF SINGLE VS. MULTIPLE INSTANCES

As TabPFN produces identical predictions whether test instances are evaluated individually or in batches, DistPFN and DistPFN-T can adjust based on either 1) the prediction of a *single* instance or 2) the average prediction across *multiple* instances. As shown in Table K.1, both choices consistently improve TabPFN-v2 (Hollmann et al., 2025), averaged across six β s for *w/ shift*, demonstrating robustness to the choice of distribution source.

	Pred. distn.	w/o shift		Shift strength (β)						
	Pred. distil.	W/O SIIIIt	0.0	0.1	0.5	1.0	2.0	5.0	Avg.	
TabPFN-v2	-	0.818	0.797	0.796	0.790	0.782	0.759	0.727	0.775	
+ DistPFN	Single Multiple	0.818 0.818	0.797 0.799	0.796 0.797	0.795 0.795	0.793 0.791	0.784 0.783	0.770 0.770	0.789 0.789	
+ DistPFN-T	Single Multiple	0.818 0.818	0.797 0.799	0.797 0.798	0.796 0.797	0.795 0.796	0.788 0.789	0.773 0.775	0.791 0.792	

Table K.1: **Predicted distributions: Single vs. Multiple.** The proposed methods consistently improves TabPFN-v2 regardless of whether the adjustment is based on single or aggregated distribution.

L OTHER METRICS

Table L.1 reports the comparison of our methods and baselines under $\beta=2$ in terms of ROC-AUC, demonstrating the effectiveness of our method. The results demonstrate that our method shows nearly the same values as the backbone, as the adjustment only rescales predicted probabilities without altering their order.

		Methods	w/o shift			Shif	t strength (β)			
		Methods	w/o sniit	0.0	0.1	0.5	1.0	2.0	5.0	Avg.
		LogReg. + HPO	$\begin{array}{c} 0.813_{\pm 0.002} \\ 0.817_{\pm 0.002} \end{array}$	$0.789_{\pm 0.002} \\ 0.806_{\pm 0.001}$	$\begin{array}{c} 0.789_{\pm 0.002} \\ 0.806_{\pm 0.001} \end{array}$	$0.790_{\pm 0.002} \\ 0.805_{\pm 0.002}$	$\begin{array}{c} 0.788_{\pm 0.002} \\ 0.803_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.784_{\pm 0.003} \\ 0.797_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.777_{\pm 0.002} \\ 0.791_{\pm 0.001} \end{array}$	0.786 0.801
		SVM + HPO	$\begin{array}{c} 0.815_{\pm 0.002} \\ 0.840_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.744_{\pm 0.004} \\ 0.804_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.747_{\pm 0.005} \\ 0.804_{\pm 0.003} \end{array}$	$0.750_{\pm 0.004} \\ 0.800_{\pm 0.001}$	$\begin{array}{c} 0.744_{\pm 0.004} \\ 0.799_{\pm 0.004} \end{array}$	$\begin{array}{c} 0.733_{\pm 0.005} \\ 0.794_{\pm 0.002} \end{array}$	$0.725_{\pm 0.006} \\ 0.785_{\pm 0.002}$	0.741 0.798
rning	MLP + HPO	$\begin{array}{c} 0.821_{\pm 0.003} \\ 0.849_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.747_{\pm 0.002} \\ 0.799_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.750_{\pm 0.003} \\ 0.799_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.746_{\pm 0.005} \\ 0.796_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.735_{\pm 0.004} \\ 0.788_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.719_{\pm 0.002} \\ 0.781_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.702_{\pm 0.004} \\ 0.772_{\pm 0.002} \end{array}$	0.733 0.789	
Machine Learning		kNN + HPO	$\begin{array}{c} 0.789_{\pm 0.001} \\ 0.828_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.728_{\pm 0.003} \\ 0.775_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.728_{\pm 0.004} \\ 0.775_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.727_{\pm 0.004} \\ 0.773_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.722_{\pm 0.004} \\ 0.768_{\pm 0.002} \end{array}$	$0.709_{\pm 0.003} \\ 0.756_{\pm 0.001}$	$\begin{array}{c} 0.693_{\pm 0.003} \\ 0.742_{\pm 0.002} \end{array}$	$0.718 \\ 0.765$
Machi		Random Forest + HPO	$0.836_{\pm 0.003} \\ 0.849_{\pm 0.003}$	$\begin{array}{c} 0.824_{\pm 0.003} \\ 0.836_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.823_{\pm 0.003} \\ 0.835_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.821_{\pm 0.003} \\ 0.834_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.818_{\pm 0.002} \\ 0.832_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.812_{\pm 0.001} \\ 0.826_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.802_{\pm 0.001} \\ 0.818_{\pm 0.002} \end{array}$	$0.817 \\ 0.830$
		LightGBM + HPO	$\begin{array}{c} 0.824_{\pm 0.002} \\ 0.845_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.813_{\pm 0.001} \\ 0.776_{\pm 0.009} \end{array}$	$\begin{array}{c} 0.812_{\pm 0.001} \\ 0.767_{\pm 0.006} \end{array}$	$\begin{array}{c} 0.809_{\pm 0.002} \\ 0.781_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.805_{\pm 0.003} \\ 0.774_{\pm 0.010} \end{array}$	$\begin{array}{c} 0.797_{\pm 0.002} \\ 0.775_{\pm 0.009} \end{array}$	$0.785_{\pm 0.003} \\ 0.779_{\pm 0.004}$	0.805 0.775
		CatBoost + HPO	$0.847_{\pm 0.003} \\ 0.843_{\pm 0.003}$	$\begin{array}{c} 0.833_{\pm 0.003} \\ 0.833_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.832_{\pm 0.002} \\ 0.832_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.830_{\pm 0.003} \\ 0.830_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.827_{\pm 0.002} \\ 0.827_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.820_{\pm 0.002} \\ 0.819_{\pm 0.002} \end{array}$	$\substack{0.810_{\pm 0.002}\\0.810_{\pm 0.001}}$	$0.825 \\ 0.825$
60	Non-found.	FT-Transformer TabM TabulaRNN MambaTab RealMLP	$\begin{array}{c} 0.821_{\pm 0.003} \\ 0.824_{\pm 0.003} \\ 0.774_{\pm 0.003} \\ 0.743_{\pm 0.005} \\ 0.821_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.818 {\pm} 0.003 \\ 0.824 {\pm} 0.003 \\ 0.699 {\pm} 0.003 \\ 0.629 {\pm} 0.006 \\ 0.805 {\pm} 0.003 \end{array}$	$\begin{array}{c} 0.819_{\pm 0.002} \\ 0.824_{\pm 0.002} \\ 0.684_{\pm 0.003} \\ 0.603_{\pm 0.004} \\ 0.806_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.816_{\pm 0.003} \\ 0.823_{\pm 0.003} \\ 0.641_{\pm 0.004} \\ 0.525_{\pm 0.002} \\ 0.807_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.812_{\pm 0.002} \\ 0.821_{\pm 0.001} \\ 0.585_{\pm 0.009} \\ 0.466_{\pm 0.010} \\ 0.804_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.795_{\pm 0.003} \\ 0.808_{\pm 0.003} \\ 0.522_{\pm 0.011} \\ 0.430_{\pm 0.005} \\ 0.795_{\pm 0.000} \end{array}$	$\begin{array}{c} 0.771_{\pm 0.002} \\ 0.791_{\pm 0.002} \\ 0.465_{\pm 0.008} \\ 0.394_{\pm 0.002} \\ 0.781_{\pm 0.004} \end{array}$	0.805 0.815 0.599 0.508 0.800
Deep Learning	п	LoCalPFN + DistPFN + DistPFN-T	$\begin{array}{c} 0.858_{\pm 0.002} \\ 0.858_{\pm 0.002} \\ 0.858_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.842_{\pm 0.002} \\ 0.842_{\pm 0.002} \\ 0.842_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.840_{\pm 0.001} \\ 0.840_{\pm 0.002} \\ 0.840_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.839_{\pm 0.000} \\ 0.839_{\pm 0.001} \\ 0.839_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.836_{\pm 0.000} \\ 0.836_{\pm 0.001} \\ 0.837_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.830_{\pm 0.000} \\ 0.830_{\pm 0.001} \\ 0.830_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.826_{\pm 0.001} \\ 0.826_{\pm 0.002} \\ 0.826_{\pm 0.003} \end{array}$	0.836 0.836 0.836
Deer	Foundation	TabICL + DistPFN + DistPFN-T	$\begin{array}{c} 0.845_{\pm 0.003} \\ 0.845_{\pm 0.003} \\ 0.845_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.832_{\pm 0.003} \\ 0.832_{\pm 0.003} \\ 0.832_{\pm 0.003} \end{array}$	$\begin{array}{c} 0.832_{\pm 0.001} \\ 0.832_{\pm 0.001} \\ 0.832_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.830_{\pm 0.002} \\ 0.830_{\pm 0.002} \\ 0.830_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.826_{\pm 0.002} \\ 0.826_{\pm 0.002} \\ 0.826_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.821_{\pm 0.002} \\ 0.821_{\pm 0.002} \\ 0.821_{\pm 0.002} \end{array}$	$\begin{array}{c} 0.813_{\pm 0.003} \\ 0.814_{\pm 0.003} \\ 0.814_{\pm 0.003} \end{array}$	0.826 0.826 0.826
	<u> </u>	TabPFN-v2 + DistPFN + DistPFN-T	$0.859_{\pm 0.002} \\ 0.859_{\pm 0.002} \\ 0.859_{\pm 0.002}$	$0.843_{\pm 0.002} \\ 0.843_{\pm 0.002} \\ 0.843_{\pm 0.002}$	$0.842_{\pm 0.003} \\ 0.843_{\pm 0.003} \\ 0.842_{\pm 0.003}$	$0.841_{\pm 0.002} \\ 0.841_{\pm 0.002} \\ 0.841_{\pm 0.002}$	$0.838_{\pm 0.002} \\ 0.838_{\pm 0.002} \\ 0.838_{\pm 0.002}$	$0.833_{\pm 0.001} \\ 0.833_{\pm 0.001} \\ 0.833_{\pm 0.001}$	$0.826_{\pm 0.002}$ $0.826_{\pm 0.003}$ $0.826_{\pm 0.003}$	0.837 0.837 0.837

Table L.1: Tabular classification results: ROC-AUC comparisons.