

# Training Data Soft Selection via Joint Density Ratio Estimation

**Ryuta Matsuno**  
*NEC Corporation.*

RYUTA-MATSUNO@NEC.COM

**Editors:** Hung-yi Lee and Tongliang Liu

## Abstract

This paper studies the training data selection problem, focusing on the selection of effective samples to improve model training using data affected by distributional shifts (i.e., data drifts). Existing drift-detection-based methods struggle with local drifts, while recent drift-localization-based methods lack theoretical support for the problem and are often ineffective. To tackle these issues, this paper proposes TSJD, a training data soft selection method based on joint density ratio estimation. TSJD assigns training weights (i.e., soft selects) to samples based on the estimated joint density ratio to align the selected data with the recent data distribution. By evaluating each sample independently of time, TSJD effectively addresses local data drifts. We also provide theoretical guarantees by deriving an upper bound on the generalization error for models trained with data selected by TSJD. In numerical experiments with four real-world datasets, TSJD shows great versatility, achieving the best or comparable results over baseline methods in all of the experiments.

**Keywords:** Training data selection; data drift; joint density ratio estimation;

## 1. Introduction

Supervised learning aims at training a prediction model to minimize test error, assuming that the data distribution is consistent between the training and test data. However, real-world applications often violate this assumption and the data distribution changes over time, known as data drift. These drifts make it ineffective to directly use the given training data (Awasthi et al., 2024; Shimodaira, 2000; Quionero-Candela et al., 2009). As a result, a problem of training data selection arises, i.e., selecting effective samples from drifting data to improve the prediction performance of models trained with (Hinder et al., 2022; Liu et al., 2017).

A naive approach to the problem just uses recent samples, assuming that the data distribution is approximately consistent for a short time span (Wang et al., 2003; Woźniak, 2013; Brzezinski and Stefanowski, 2014). However, this approach discards all the older samples, which is potentially effective for model training. Drift detection methods (Bifet and Gavalda, 2007; Page, 1954) improve this approach by determining when to separate the recent and old samples, performing concept drift detection from the present to the past, and utilizing the samples up to the time when a drift is detected. However, they still select samples only based on time, failing to adapt to local concept drifts that occurred in a small part of the input space, as well as recurring concept drifts (Hinder et al., 2022).

Recent studies propose drift localization methods (Hinder et al., 2022; Liu et al., 2017), which detect local data drift in the sample space based on recent data, and we can select

samples based on the input and output of each sample, aligned with the recent data distribution for model training. Although these approaches are capable of flexibly selecting samples from old samples, the theoretical properties of models trained with the selected samples remain unknown, limiting their validity to the problem. Indeed, the empirical performance of these methods is inferior to naive baselines, as shown in our experiments.

To address this, we propose TSJD, a Training data soft Selection method based on Joint Density ratio estimation. TSJD first trains a joint density ratio estimator between recent and old data distributions. It then assigns training weights to each sample (i.e., soft-selects) based on the estimated ratio, effectively addressing local data drifts by utilizing both inputs and outputs. In addition, we provide a theoretical upper bound on the generalization error of models trained with our method, ensuring the validity of our method for the training data selection problem. Experiments on four real-world datasets show TSJD consistently achieves the best or comparable results across all 30 settings, demonstrating its effectiveness and versatility.

Our contributions are summarized as follows;

- We propose TSJD, a training data soft selection method using a joint density ratio estimator to effectively handle local data drifts (Section 3).
- We offer theoretical analysis and establish a generalization error upper bound to support the validity of TSJD (Section 4).
- We conduct extensive numerical experiments and provide empirical evidences which highlight the superiority of TSJD over various baseline methods (Section 5).

Due to the space limitation, all proofs of theorems and lemmas are presented in the supplementary material. We also report comprehensive experiments on seven real-world datasets across 126 settings in our supplementary material.

## 2. Preliminary

In this section, we explain the problem formulation as well as related works briefly.

### 2.1. Problem Formulation

We consider a supervised classification problem. The input space is  $\mathcal{X} \subseteq \mathbb{R}^d$  and the output space is  $\mathcal{Y} = [K]$ , where  $[K]$  denotes the set of integers from 1 to  $K$ , i.e.,  $[K] := \{1, \dots, K\}$  and the integer  $K \in \mathbb{Z}_{\geq 2}$  is the number of classes. Let  $p_t(\mathbf{x}, y)$  be a joint distribution over  $\mathcal{X} \times \mathcal{Y}$  at time  $t \in \mathbb{Z}_{\geq 1}$ . A sample  $(\mathbf{x}_t, y_t)$  is sampled from  $p_t(\mathbf{x}, y)$  at every time step  $t \in [T]$ , where  $T \in \mathbb{Z}_{\geq 1}$  is the current time. All samples available at the training phase form a datasets  $D := \{(\mathbf{x}_t, y_t)\}_{t=1}^T$ . A standard approach for classification tasks is to train a model by minimizing the cross-entropy loss, i.e.,

$$\ell_{\text{CE}}(h(\mathbf{x}), y) := -\log(h(y|\mathbf{x})), \quad (1)$$

where  $h : \mathcal{X} \rightarrow \Delta^{K-1}$  is a probabilistic classification model,  $\Delta^{K-1} := \{\mathbf{p} \in \mathbb{R}_{\geq 0}^K \mid \|\mathbf{p}\|_1 = 1\} \subset \mathbb{R}^K$  is the  $(K-1)$  dimensional probability simplex, and  $h(y|\mathbf{x}) = (h(\mathbf{x}))_y$  computes the probability that an input  $\mathbf{x} \in \mathcal{X}$  belongs to a class  $y \in \mathcal{Y}$ . The model  $h$  predicts a class of  $\mathbf{x} \in \mathcal{X}$  by  $\arg \max_{y \in \mathcal{Y}} h(y|\mathbf{x})$ .

Let  $\mathcal{H}$  be the hypothesis space of a classification model  $h$ . We aim to find  $h^* \in \mathcal{H}$  that maximizes accuracy over the next  $M$  steps, i.e.,  $t = T + 1, \dots, T + M$ . To achieve this, we seek to minimize the expected zero-one loss over the distributions  $p_{T+1}(\mathbf{x}, y), \dots, p_{T+M}(\mathbf{x}, y)$  (a.k.a. the zero-one risk), denoted by  $R_{01}$ ;

$$R_{01}(h) := \frac{1}{M} \sum_{t=T+1}^{T+M} \mathbb{E}_{p_t(\mathbf{x}, y)} [\ell_{01}(h(\mathbf{x}), y)], \quad (2)$$

where  $\ell_{01}(h(\mathbf{x}), y) = \mathbb{I}[\arg \max_k h(k|\mathbf{x}) \neq y]$  and  $\mathbb{I}[P]$  is the Iverson bracket, which is 1 if the proposition  $P$  is true and 0 otherwise. The notation  $\mathbb{E}_{p(\mathbf{x}, y)}[f(\mathbf{x}, y)] := \int_{\mathcal{X} \times \mathcal{Y}} f(\mathbf{x}, y) p(\mathbf{x}, y) d\mathbf{x} dy$  is the expectation of a function  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  over the joint distribution  $p(\mathbf{x}, y)$ .

Obtaining  $h^*$  is challenging because the future distribution  $p_t$  for  $t > T$  is unknown. Additionally, if there are concept drifts (changes in the conditional distribution of  $y$  given  $\mathbf{x}$ ,  $p(y|x)$ , a.k.a. conditional shift) with in  $D$ , the naive use of all of  $D$  (a.k.a. *ERM: empirical risk minimization*) is unsuitable for finding  $h^*$ . A common practical approach is to use the most recent  $N$  samples from  $D$ , i.e.,  $\{(\mathbf{x}_t, y_t)\}_{t=T-N+1}^T$ , as training data (Wang et al., 2003; Woźniak, 2013; Brzezinski and Stefanowski, 2014). Here,  $N \in [T]$  is determined based on domain knowledge or set as a hyperparameter. This approach is based on an implicit assumption as follows.

**Assumption 1 (Temporal consistency of the joint distribution)** *There exists an integer  $N \in [T]$  and small constants  $\tau_X, \tau_{Y|X} \geq 0$  such that for any  $t \in \{T - N + 1, \dots, T\}$  and  $t' \in \{T + 1, \dots, T + M\}$  each of the followings holds;*

$$d_X(p_t, p_{t'}) \leq \tau_X \quad (3)$$

$$\forall \mathbf{x} \in \mathcal{X}, \quad d_{Y|X}(p_t(\cdot|\mathbf{x}), p_{t'}(\cdot|\mathbf{x})) \leq \tau_{Y|X}, \quad (4)$$

where  $d_X$  and  $d_{Y|X}$  compute the distances of two marginal ( $p_t(\mathbf{x})$  and  $p_{t'}(\mathbf{x})$ ) and conditional ( $p_t(y|\mathbf{x})$  and  $p_{t'}(y|\mathbf{x})$ ) distributions, respectively.

We define  $p_t(\cdot|\mathbf{x})$  as  $p_t(\cdot|\mathbf{x}) := [p_t(y = 1|\mathbf{x}), \dots, p_t(y = K|\mathbf{x})]^\top \in \Delta^{K-1}$ . Specifically, in our analysis in Section 4, we use the Wasserstein 1-distance (Edwards, 2011) (a.k.a. earth mover's distance)  $W_1$  for  $d_X$  and the  $L^2$ -norm for  $d_{Y|X}$ . Various choices of  $\tau_X$  and  $\tau_{Y|X}$  have been explored in examples such as the following;

**Example 1** *The traditional ERM (Hastie et al., 2001) assumes  $\tau_X = \tau_{Y|X} = 0$ .*

**Example 2** *Covariate shift (Shimodaira, 2000) assumes  $\tau_X > 0$  and  $\tau_{Y|X} = 0$ .*

The naive approach under Assumption 1 selects the recent samples  $\{(\mathbf{x}_t, y_t)\}_{t=T-N+1}^T$  and discards older samples from  $t = 1$  to  $t = T - N$ . Although the older samples might worsen the training of  $h$  due to data drifts, selecting effective ones can enhance its performance. To tackle this, we employ a soft selection method by assigning positive weights to each sample in  $D$ . In summary, this paper formulates the problem as follows.

**Definition 1 (Training Data Soft Selection Problem)** Given a dataset  $D = \{(\mathbf{x}_t, y_t)\}_{t=1}^T$ , where  $(\mathbf{x}_t, y_t)$  is independently sampled from  $p_t$  and assume Assumption 1 holds, the task is to find sample weights  $W^* = \{w_t\}_{t=1}^T \in \mathbb{R}_{\geq 0}^T$  which minimizes the zero-one risk of the trained model with, i.e.,

$$W^* = \arg \min_{W \in \mathbb{R}_{\geq 0}^T} R_{01} \left( \arg \min_{h \in \mathcal{H}} \sum_{t=1}^T w_t \ell_{\text{CE}}(h(\mathbf{x}_t), y_t) \right), \quad (5)$$

where  $\mathcal{H}$  is an arbitrary hypothesis space.

**Remark 2** In our problem formulation,  $\mathcal{H}$  is given after choosing  $W$ . If  $\mathcal{H}$  were given before deciding  $W$ , we could optimize both  $W$  and  $h$  simultaneously, likely improving  $R_{01}$  risk (Zhang et al., 2020; Bassily et al., 2024; Mohri and Muñoz Medina, 2012). In practice, however, AutoML tools, such as `auto-sklearn`<sup>1</sup> and `PyCaret`,<sup>2</sup> often handle the training of  $h$ , limiting customization of the training. In addition,  $\mathcal{H}$  is often composed of different models with different behaviors, including decision trees, linear models, gradient boosting models, and neural networks, making the optimization of  $W$  along with  $h$  unstable. Furthermore, in MLOps frameworks (Kreuzberger et al., 2023; Ruf et al., 2021), data preparation and model training are separate steps. These conditions make joint optimization impractical, while our formulation remains usable.

## 2.2. Related Works

We review drift detection methods, drift localization methods, and density ratio estimation methods as related works as follows;

**Drift Detection Methods.** Drift detection methods identify change points in data distribution and have been studied for over a half century (Page, 1954; Bifet and Gavaldà, 2007; Gama et al., 2004; Mayaki and Riveill, 2022). Concept drift can be detected by applying these drift detection methods to the stream of the prediction losses (Mehmood et al., 2021; Gonçalves et al., 2014) and this can be applied to our problem by detecting concept drift from the present  $t = T$  backward to the past  $t = 1$  and selecting samples until a drift is detected. However, as noted by Hinder et al. (2022), “... if a drift only occurs in a small region of the entire feature space, the other non-drifted regions may also be suspended, thereby reducing the learning efficiency of models.”, these time-based methods often fail to flexibly select samples, which can decrease the efficiency of learning models.

**Drift Localization Methods.** Unlike traditional drift detection methods that determine *when* drift occurs, recent drift localization techniques identify *where* drift happens. Liu et al. (2017) introduce LDD-DIS, which detects local drift by comparing the number of recent and old samples in the  $k$ -nearest neighbors among the data. Building on this, LDD-DSDA is developed to select samples for the problem. Hinder et al. (2022) propose a theoretical framework called LCD, which reframes drift localization as a supervised classification problem, offering improved detection performance over LDD-DIS. However, both methods lack theoretical analysis for model training and often struggle to select samples effectively.

1. <https://automl.github.io/auto-sklearn/master/>

2. <https://pycaret.org/>

**Density Ratio Estimation.** The density ratio, which compares two probability distributions, has been a research focus for over two decades (Shimodaira, 2000). Kernel based methods, such as KLEIP (Sugiyama et al., 2007a), uLSIF (Kanamori et al., 2009), RuLSIF (Yamada et al., 2013), KMM (Schölkopf et al., 2007), and other methods (Sugiyama et al., 2012; Kato and Teshima, 2021; Zhang et al., 2020) have been proposed and utilized not only for covariate shift adaptation (Shimodaira, 2000; Sugiyama et al., 2007a), but also for generative models (Goodfellow et al., 2014), mutual information approximation (Suzuki et al., 2009), and change point detection (Liu et al., 2013). Various extensions exist, such as joint-to-marginal (Matsushita et al., 2022), conditional distribution given input (Sugiyama et al., 2010) and output (Sugiyama, 2010), and continuous covariate shift (Zhang et al., 2023). However, joint density ratio estimation, crucial for addressing our problem, remains insufficiently explored.

### 3. Proposed Method

This section introduces our method, TSJD, a training data soft selection method based on joint density ratio estimation. Section 3.1 explains the notation and assumptions while Section 3.2 present the algorithm of TSJD. Section 3.3 describes training of the joint density estimator, and Section 3.4 offers details on modeling and hyperparameter tuning.

#### 3.1. Notation and Assumption

We define the marginalization of the  $N$ -recent data distribution as

$$\bar{p}_T(\mathbf{x}, y) := \frac{1}{N} \sum_{t=T-N+1}^T p_t(\mathbf{x}, y), \quad (6)$$

where we use the subscript  $T$  to denote *Target*. We consider the data  $D_T := \{(\mathbf{x}_t, y_t)\}_{t=T-N+1}^T$  to be approximately i.i.d. samples from the distribution  $\bar{p}_T(\mathbf{x}, y)$ . Similarly, we define

$$\bar{p}_S(\mathbf{x}, y) := \frac{1}{T-N} \sum_{t=1}^{T-N} p_t(\mathbf{x}, y), \quad (7)$$

as the old data distribution, and  $D_S := \{(\mathbf{x}_t, y_t)\}_{t=1}^{T-N}$  is considered as samples of size  $(T-N)$  from the distribution  $\bar{p}_S(\mathbf{x}, y)$ . Here, the subscript  $S$  is short for *Source*. Moreover, we define the *test* distribution  $p_{te}$  as

$$p_{te}(\mathbf{x}, y) := \frac{1}{M} \sum_{t=T+1}^{T+M} p_t(\mathbf{x}, y). \quad (8)$$

and assume that  $R_{01}(h) = \mathbb{E}_{p_{te}(\mathbf{x}, y)}[\ell_{01}(h(\mathbf{x}), y)]$ .

These formulation and assumption allow us to view the problem of Definition 1 as one that to relate the three distributions  $\bar{p}_S$ ,  $\bar{p}_T$ , and  $p_{te}$ . With this understanding, we present our method in the next section.

### 3.2. Algorithm of TSJD

To derive our proposed method, we establish a key theorem that links the zero-one risk  $R_{01}$  with the squared  $L^2$ -norm of the difference between  $h(\mathbf{x})$  and  $\bar{p}_T(\cdot|\mathbf{x})$  over  $\bar{p}_T$ , as follows.

**Theorem 3** *For any  $h \in \mathcal{H}$ , the following holds.*

$$R_{01}(h) - B_{01} = \mathcal{O}\left(\mathbb{E}_{\bar{p}_T(\mathbf{x})}\left[\|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2\right] + Z(h)\tau_X + \tau_{Y|X}^2\right), \quad (9)$$

where we define  $Z(h) := \sup_{\mathbf{x} \in \mathcal{X}} \left\| \nabla \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right\|_2$  and  $B_{01} := \min_f R_{01}(f)$  as the Bayes error, i.e.,  $B_{01}$  is the lowest value of  $R_{01}$  among any possible classification model  $f : \mathcal{X} \rightarrow \Delta^{K-1}$ .

Theorem 3 indicates that fitting  $h(\mathbf{x})$  to  $\bar{p}_T(\cdot|\mathbf{x})$  is sufficient for the problem, i.e., selecting samples to make  $h$  learn  $\bar{p}_T$  solves the problem of Definition 1.

**Remark 4** *Although we have another  $h$ -related term  $Z(h)\tau_X$  beside the first term  $\mathbb{E}_{\bar{p}_T(\mathbf{x})}\left[\|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2\right]$  in Eq. (9), it can be considered negligible due to the following reasons;*

- *Small  $\tau_X$  Assumption:* The value of  $\tau_X$  is assumed to be a small constant, e.g.,  $\tau_X \ll 1$ , inherently reducing the impact of the term  $Z(h)\tau_X$ .
- *Convergence of  $Z(h)$ :* Even if  $\tau_X$  is not particularly small,  $Z(h) \rightarrow 0$  holds with the first term in Eq. (9) converges to zero, i.e.,  $\mathbb{E}_{\bar{p}_T(\mathbf{x})}\left[\|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2\right] \rightarrow 0$ , further diminishing the significance of the term  $Z(h)\tau_X$ .

This understanding makes us to use the following weighting strategies;

**Weights for the recent samples.** We set the weight  $w_t \propto 1$  for all  $t \in \{T - N + 1, \dots, T\}$ , as the naive approach does. Since  $D_T$  is assumed to be sampled from  $\bar{p}_T$ , this weights enable us to empirically approximate the the expected cross-entropy loss of  $h$  over  $\bar{p}_T$ , denoted by  $L_{\text{CE}}(h)$  as

$$L_{\text{CE}}(h) := \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)}[\ell_{\text{CE}}(h(\mathbf{x}), y)] \approx \frac{1}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} \ell_{\text{CE}}(h(\mathbf{x}), y). \quad (10)$$

By the strict-properness of the the cross-entropy loss (Gneiting and Raftery, 2007) and Theorem 3, the minimization of Eq. (10) leads to minimize  $R_{01}(h)$ .

**Weights for the old samples.** For the old samples with  $t \in [T - N]$ , we apply the *importance weighting* technique (Shimodaira, 2000; Sugiyama et al., 2007a), initially proposed for covariate shift adaptation; Let  $r(\mathbf{x}, y) := \frac{\bar{p}_T(\mathbf{x}, y)}{\bar{p}_S(\mathbf{x}, y)}$  be the joint density ratio of  $\bar{p}_T$  over  $\bar{p}_S$ . Then, we have

$$\mathbb{E}_{\bar{p}_T(\mathbf{x}, y)}[f(\mathbf{x}, y)] = \int f(\mathbf{x}, y) \frac{\bar{p}_T(\mathbf{x}, y)}{\bar{p}_S(\mathbf{x}, y)} \bar{p}_S(\mathbf{x}, y) d\mathbf{x} dy = \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[r(\mathbf{x}, y) f(\mathbf{x}, y)]. \quad (11)$$

Hence, setting a sample weight  $w_t$  to be  $w_t \propto r(\mathbf{x}_t, y_t)$  converts the expectation over  $\bar{p}_S$  into that over  $\bar{p}_T$ . With the same logic for the recent samples, the minimization of  $\ell_{\text{CE}}(h(\mathbf{x}_t), y_t)$

---

**Algorithm 1** Main algorithm of TSJD
 

---

**Input:** Data  $D = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)) \in (\mathcal{X} \times \mathcal{Y})^T$ , Number of recent samples  $N \in [T-1]$   
 w.r.t. Assumption 1  
 // Step 1  
 1:  $\forall t \in \{T - N + 1, \dots, T\}, w_t \leftarrow \frac{1}{2N}$   
 // Step 2  
 2: Train a joint density ratio estimator  $\hat{g} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$  based on  $D_T$  and  $D_S$   
 3:  $\forall t \in [T - N], w_t \leftarrow \frac{1}{2(T-N)} \hat{g}(\mathbf{x}_t, y_t)$   
**Output:** Sample weights  $W = [w_1, \dots, w_T]^T \in \mathbb{R}_{\geq 0}^T$

---

with the weight  $w_t = r(\mathbf{x}_t, y_t)$  over  $\bar{p}_S$  leads to minimize  $R_{01}(h)$ . Since  $r(\mathbf{x}, y)$  is not explicitly available, we train a joint density ratio estimator  $\hat{g} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$  and use  $w_t \propto \hat{g}(\mathbf{x}_t, y_t)$  for the weight for all  $t \in [T - N]$ .

Algorithm 1 summarizes our method for the training data selection problem, where we normalize the weights using  $N$  and  $T$ . By our Algorithm 1, the model  $h$  will be trained to minimize  $\hat{L}_{\text{CE}}(h; D_S, D_T)$  defined as

$$\hat{L}_{\text{CE}}(h; D_S, D_T) := \frac{1}{2} \left( \frac{1}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} \ell_{\text{CE}}(h(\mathbf{x}), y) + \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} \hat{g}(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y) \right), \quad (12)$$

with our joint density ratio estimator  $\hat{g}$ . Next, we specify how to train  $\hat{g}$  based on  $D_S$  and  $D_T$ .

### 3.3. Training the Density Ratio Estimator $\hat{g}$

We train a density ratio estimator  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$  which approximates the true density ratio  $r(\mathbf{x}, y) = \frac{\bar{p}_T(\mathbf{x}, y)}{\bar{p}_S(\mathbf{x}, y)}$  by minimizing the expected squared error  $J(g)$  over  $\bar{p}_S$  (Kanamori et al., 2009; Zhang et al., 2020);

$$J(g) := \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} \left[ (g(\mathbf{x}, y) - r(\mathbf{x}, y))^2 \right], \quad (13)$$

whose empirical version  $\hat{J}(g; D_S, D_T)$  is defined as

$$\hat{J}(g; D_S, D_T) := \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} g(\mathbf{x}, y)^2 - \frac{2}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} g(\mathbf{x}, y) + C_r, \quad (14)$$

where  $C_r := \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [r(\mathbf{x}, y)^2] = \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [r(\mathbf{x}, y)]$  is an independent constant and can be ignored to train  $g$ . In addition,  $g$  needs to satisfy

$$1 = \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [g(\mathbf{x}, y)] \approx \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} g(\mathbf{x}, y) \quad (15)$$

to be a proper density ratio due to the fact  $\mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[r(\mathbf{x}, y)] = \int \bar{p}_T(\mathbf{x}, y) d\mathbf{x} dy = 1$ . Hence, we add an empirical constraint with a hyperparameter  $\beta > 0$ . The final loss function to train  $g$  is defined as

$$L(g; D_S, D_T) := \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} g(\mathbf{x}, y)^2 - \frac{2}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} g(\mathbf{x}, y) + \beta \left( \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} g(\mathbf{x}, y) - 1 \right)^2, \quad (16)$$

and we denote the minimizer of  $L(g; D_S, D_T)$  by  $\hat{g}$ .

**Remark 5** *The constraint Eq. (15) is often overlooked in existing methods (Kanamori et al., 2009; Yamada et al., 2013; Zhang et al., 2020), as claimed by Sugiyama et al. (2007b) “... the normalization constraint (Eq. (15)) is not generally satisfied exactly ... this may not be critical in practice since the scale of the importance is often irrelevant in subsequent learning algorithms.”. However, since we use both  $D_S$  and  $D_T$ , the correct scale is vital to control and balance the effects of the sample weights. Additionally, in our analysis, we assume  $\bar{p}_S(\mathbf{x}, y)g(\mathbf{x}, y)$  is a probability density. Therefore, contrary to the claim, the constraint term is crucial in our problem.*

### 3.4. Implementation and Hyperparameter Tuning of $\hat{g}$

We employ a linear-in-parameter model (Zhang et al., 2020) (a.k.a. linear basis expansion) with the softplus activation;  $\text{softplus}(x) := \log(1 + \exp(x))$  for  $g$  as

$$g(\mathbf{x}, y) = \text{softplus} \left( \sum_{i=1}^{N_M} a_i \phi_i(\mathbf{x}, y) \right), \quad (17)$$

where  $a_i \in \mathbb{R}$  is the learning parameter,  $\phi_i : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is the  $i$ -th feature mapping (a.k.a. basis function), and  $N_M = 200$  is the number of the feature mappings. The feature mapping  $\phi_i$  is modeled using the Gaussian RBF as  $\phi_i(\mathbf{x}, y) := \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|_2^2}{\sigma_x}\right) \max(\sigma_y, \mathbb{I}[y = y_i])$ , where  $(\mathbf{x}_i, y_i)$  is the kernel center, sampled from  $D_T$  uniformly at random,  $\sigma_x > 0$  and  $\sigma_y \geq 0$  are the hyperparameters. The parameters  $\{a_i\}_{i=1}^{N_M}$  are optimized using gradient descent.

The hyperparameters of TSJD, i.e.,  $\sigma_x, \sigma_y$ , and  $\beta$ , are tuned by a grid search and ones that minimize  $\hat{J}(\hat{g}; D_S, D_T)$ , satisfying the following constraint is selected;

$$\left| \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} \hat{g}(\mathbf{x}, y) - 1 \right| \leq G \sqrt{\frac{\log \frac{2}{\delta}}{2|D_S|}}, \quad (18)$$

where  $\hat{g}$  is obtained by minimizing Eq. (16) with each set of the hyperparameters, and we set  $G = 10$  and  $\delta = 0.05$ . Note this constraint is different from the term inside Eq. (16); This is based on the following lemma, that states that even the training is perfect, i.e.,  $\hat{g} = r$ , the constraint Eq. (15) can only be satisfied with a margin  $G \sqrt{\frac{\log \frac{2}{\delta}}{2|D_S|}}$  in probability. The proof is omitted since it is trivial by Hoeffding’s inequality.



**Lemma 6** Assume  $r(\mathbf{x}, y) \leq G$  for any  $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ . Then, for any  $\delta \in (0, 1)$ , the following holds with probability at least  $1 - \delta$ .

$$\left| \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} r(\mathbf{x}, y) - 1 \right| \leq G \sqrt{\frac{\log \frac{2}{\delta}}{2|D_S|}} \quad (19)$$

## 4. Theoretical Analysis

In this section, we provide our theoretical analysis, bounding the generalization error of models trained with our method. Before presenting our analysis, we clarify the notation and assumption used in our analysis as follows;

- Let  $\mathcal{G}_+$  be the hypothesis space for the joint density ratio predictor  $g : \mathcal{X} \times \mathcal{Y} \rightarrow [0, G]$  with a constant  $G \geq 1$  and assume that  $\forall g \in \mathcal{G}_+, \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[g(\mathbf{x}, y)] = 1$  holds.
- Let  $\mathcal{G}$  be defined as  $\mathcal{G} := \mathcal{G}_+ \cup \{g' : (\mathbf{x}, y) \mapsto -g(\mathbf{x}, y) | g \in \mathcal{G}_+\}$ .
- Assume that  $\forall (\mathbf{x}, y, h) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{H} : \ell_{\text{CE}}(h(\mathbf{x}), y) \leq U$  holds with a constant  $U \geq 0$ .

### 4.1. Main Result

Our analysis yields an upper bound on the zero-one risk of  $\hat{h}$ , trained with weights computed by our method. The main theorem detailing the generalization error bound and its order is presented in Theorem 7 and Corollary 8, respectively. Notably,  $C_4(\delta) = \mathcal{O}(\mathfrak{R}_N(\mathcal{H}) + \mathfrak{R}_{T-N}(\mathcal{H}))$  and  $C_3(\delta) = \mathcal{O}(\mathfrak{R}_N(\mathcal{G}) + \mathfrak{R}_{T-N}(\mathcal{G}))$  are defined using the Rademacher complexity  $\mathfrak{R}$  (Koltchinskii, 2001; Cortes et al., 2016; Maurer, 2016; Ledoux and Talagrand, 2013; Mohri et al., 2018). The exact definitions and notation will be provided in the subsequent sections.

**Theorem 7 (Generalization error bound)** For any  $\delta \in (0, 1)$ , the following inequality holds with probability at least  $1 - \delta$ ;

$$\begin{aligned} R_{01}(\hat{h}) - B_{01} \\ \leq 4K\eta_{\min}^{-2} \left( T_{KL}(h^*) + U\sqrt{J(g^*)} + C_4(\delta/5) + U\sqrt{C_3(\delta/5)} + \frac{1}{K}Z(\hat{h})\tau_X + \frac{1}{K}\tau_{Y|X}^2 \right), \end{aligned} \quad (20)$$

where  $T_{KL}(h)$  is the expected Kullback-Leibler divergence between  $\bar{p}_T(\cdot|\mathbf{x})$  and  $h(\mathbf{x})$  over  $\bar{p}_T(\mathbf{x})$ , i.e.,

$$T_{KL}(h) := \mathbb{E}_{\bar{p}_T(\mathbf{x})} [D_{KL}(\bar{p}_T(\cdot|\mathbf{x}) || h(\mathbf{x}))]. \quad (21)$$

**Corollary 8** Assume that  $\mathfrak{R}_n(\mathcal{G}) = \mathcal{O}(n^{-1/2})$  and  $\mathfrak{R}_n(\mathcal{H}) = \mathcal{O}(n^{-1/2})$ , then following order holds;

$$R_{01}(\hat{h}) - B_{01} = \mathcal{O}\left(T_{KL}(h^*) + \sqrt{J(g^*)} + Z(\hat{h})\tau_X + \tau_{Y|X}^2\right) + \mathcal{O}_p\left(N^{-\frac{1}{4}} + (T - N)^{-\frac{1}{4}}\right), \quad (22)$$

where  $\mathcal{O}_p$  denotes the order in probability.

**Remark 9** Theorem 7 and Corollary 8 show that if  $p(y|\mathbf{x}) = h^*(y|\mathbf{x})$  and  $r(\mathbf{x}, y) = g^*(\mathbf{x}, y)$  hold for any  $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ , then as  $N$  and  $T$  increase, the difference between  $R_{01}(\hat{h})$  and the Bayes error  $B_{01}$  approaches  $\mathcal{O}\left(Z(\hat{h})\tau_X + \tau_{Y|X}^2\right)$ , which is inevitable due to data drift. Therefore, the generalization error of  $\hat{h}$  can be considered optimal, theoretically validating our method of Algorithm 1.

In the following sections, we introduce the key lemmas and theorems for deriving Theorem 7.

## 4.2. Generalization Error Bound of $\hat{g}$

We first establish the generalization error bound of our joint density ratio estimator  $\hat{g} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$  in terms of the expected squared error  $J$ .

**Theorem 10** Let  $\hat{g}$  and  $g^*$  be the minimizers of  $\hat{J}(g; D_S, D_T)$  and  $J(g)$  among  $\mathcal{G}_+$ , respectively. Then, for any  $\delta \in (0, 1)$ , the following inequality holds with probability at least  $1 - \delta$ ;

$$J(\hat{g}) \leq J(g^*) + C_3(\delta), \quad (23)$$

where  $C_3(\delta) := 4G\mathfrak{R}_{|D_S|}(\mathcal{G}) + 4\mathfrak{R}_{|D_T|}(\mathcal{G}) + 4G^2\sqrt{\frac{\log \frac{3}{\delta}}{2}}\left(\frac{1}{\sqrt{|D_S|}} + \frac{1}{\sqrt{|D_T|}}\right)$  and  $\mathfrak{R}_n(\mathcal{G})$  is the Rademacher complexity (Koltchinskii, 2001) of  $\mathcal{G}$  with sampling size  $n$ .

The following corollary is obvious from Theorem 10.

**Corollary 11** Assume that  $\mathfrak{R}_n(\mathcal{G}) = \mathcal{O}(n^{-1/2})$ , then following order holds;

$$J(\hat{g}) = J(g^*) + \mathcal{O}_p\left(N^{-1/2} + (T - N)^{-1/2}\right), \quad (24)$$

**Remark 12** Corollary 11 indicates that if  $\mathcal{G}_+$  is properly chosen and  $r = g^* \in \mathcal{G}_+$ , the right hand of Eq. (24) decreases to 0 at the rate of  $(T - N)^{-1/2} + N^{-1/2}$  in probability. This ensures that  $\hat{g}$  converges to  $r$  as  $N$  and  $T$  approach infinity, confirming the theoretical soundness of our method for training  $g$ .

## 4.3. Generalization Error Bound of $h$ Trained with Our Method

Next, we analyze the generalization error of  $\hat{h}$ , a classification model trained with our selected training sample using  $\hat{g} \in \mathcal{G}_+$ . The following Lemma 13 and Lemma 14 provide the relation between  $L_{\text{CE}}(h)$  and empirical error of a classification model  $h$ .

**Lemma 13** For any  $\delta \in (0, 1)$  and  $h \in \mathcal{H}$ , over the draw of i.i.d. samples  $S_T$  from  $\bar{p}_T$ , the following inequality holds with probability at least  $1 - \delta$ ;

$$L_{\text{CE}}(h) \leq \frac{1}{|S_T|} \sum_{(\mathbf{x}, y) \in S_T} \ell_{\text{CE}}(h(\mathbf{x}), y) + C_1(\delta) \quad (25)$$

where  $C_1(\delta) := 2\sqrt{2}\exp(U)\mathfrak{R}_{|S_T|}(\mathcal{H}) + U\sqrt{\frac{\log \frac{1}{\delta}}{2|S_T|}}$  and  $\mathfrak{R}_n(\mathcal{H})$  is the vector-valued Rademacher complexity (Maurer, 2016; Cortes et al., 2016) of  $\mathcal{H}$  with sampling size  $n$ .

Table 1: Dataset statistics.

Data set	Samples	Features	Classes
Weather	18159	8	2
Smartmeter	22950	96	10
Powersupply	29928	2	24
Forest	581012	54	2

**Lemma 14** For any  $\delta \in (0, 1)$  and any  $h \in \mathcal{H}$ , over the draw of i.i.d. samples  $S$  from  $\bar{p}_S$ , the following inequality holds with probability at least  $1 - \delta$ :

$$L_{CE}(h) \leq \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} g(\mathbf{x}, y) \ell_{CE}(h(\mathbf{x}), y) + C_2(\delta) + U \sqrt{\frac{\mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[(r(\mathbf{x}, y) - g(\mathbf{x}, y))^2]}{2|S|}} \quad (26)$$

where  $C_2(\delta) := 2(2U + G) \exp(U) \mathfrak{R}_{|S|}(\mathcal{H}) + 2(U + 2G) \mathfrak{R}_{|S|}(\mathcal{G}) + MG \sqrt{\frac{\log \frac{1}{\delta}}{2|S|}}$ .

Based on Lemma 13 and Lemma 14, we obtain the generalization error bound w.r.t.  $L_{CE}$ .

**Theorem 15** Let  $\hat{h}$  and  $h^*$  be the minimizers of  $\hat{L}_{CE}(h)$  and  $L_{CE}(h)$  among  $\mathcal{H}$ , respectively. Then, for any  $\delta \in (0, 1)$ , the following inequality holds with probability at least  $1 - \delta$ ;

$$L_{CE}(\hat{h}) - L_{CE}(h^*) \leq U \sqrt{J(g^*)} + C_4(\delta/5) + U \sqrt{C_3(\delta/5)} \quad (27)$$

where we define

$$\begin{aligned} C_4(\delta) := & \sqrt{2} \exp(U) \mathfrak{R}_{|D_T|}(\mathcal{H}) + U \sqrt{\frac{\log \frac{1}{\delta}}{2|D_T|}} + (2U + G) \exp(U) \mathfrak{R}_{|D_S|}(\mathcal{H}) \\ & + (U + 2G) \mathfrak{R}_{|D_S|}(\mathcal{G}) + GU \sqrt{\frac{\log \frac{1}{\delta}}{2|D_S|}}. \end{aligned} \quad (28)$$

Based on Theorem 15 and Theorem 3, we derive the generalization upper bound for the risk  $R_{01}(\hat{h})$ , which is our final target to minimize. The bound is presented in Section 4.1, and we have already discussed its implication, showing theoretical validity of our method.

## 5. Numerical Experiments

We conducted experiments to test the empirical effectiveness of our method on real-world datasets.

**Dataset.** We use four real-world multi-class classification datasets obtained from USP DS Repository (Souza et al., 2020)<sup>3</sup>. We select two severely drifting (Powersupply and Forest) and two relatively stationary datasets (Weather and Smartmeter), as shown in Table 1.

3. <https://sites.google.com/view/uspdsrepository>, Accessed: 2025-06-24

Table 2: Average zero-one loss ( $\downarrow$ ) over 30 random trials. **Boldfaces with star\*** highlight the lowest errors and basic **boldfaces** show comparable results based on the Wilcoxon signed-rank test (Wilcoxon, 1945) with the significance level of 1%.

Data	Model	$N$	$T$	Naive Baseline		Time-based		Cov.shift	Drift Localization		(Ours)
				$D_T$	$D$	PHT	ADWIN	uLSIF	LDD-DSDA	LCD	TSJD
Weather	LGBM	200	2000	29.57	<b>21.97*</b>	<b>22.00</b>	<b>22.30</b>	28.33	24.60	<b>22.53</b>	<b>22.90</b>
		500	5000	19.90	<b>17.73</b>	<b>17.73</b>	<b>17.60*</b>	19.63	<b>18.70</b>	<b>18.33</b>	<b>18.43</b>
		1000	10000	23.70	<b>20.57</b>	<b>21.13</b>	<b>21.30</b>	22.93	<b>21.03</b>	<b>20.43*</b>	<b>21.37</b>
	NN	200	2000	27.50	<b>20.63</b>	<b>21.13</b>	<b>20.47*</b>	26.43	<b>22.73</b>	<b>20.63</b>	<b>22.03</b>
		500	5000	19.90	<b>17.57</b>	<b>17.63</b>	<b>17.53*</b>	<b>17.83</b>	<b>18.47</b>	<b>17.80</b>	<b>18.00</b>
		1000	10000	<b>20.70</b>	<b>18.93*</b>	<b>19.10</b>	<b>19.37</b>	<b>21.03</b>	<b>19.93</b>	<b>19.60</b>	<b>19.23</b>
Smartmeter	LGBM	200	2000	26.67	19.90	20.07	<b>20.10</b>	27.40	21.77	20.13	<b>17.67*</b>
		500	5000	22.97	<b>13.50</b>	17.33	19.90	24.80	16.73	<b>14.33</b>	<b>12.83*</b>
		1000	10000	21.13	<b>13.13</b>	17.40	19.57	21.90	15.40	<b>12.70</b>	<b>12.07*</b>
		2000	20000	22.23	15.20	17.43	17.80	23.23	16.53	15.67	<b>13.37*</b>
	NN	200	2000	36.30	36.80	36.73	<b>33.23</b>	40.00	37.77	36.43	<b>31.00*</b>
		500	5000	35.27	33.33	35.27	36.43	37.43	34.27	33.27	<b>29.20*</b>
		1000	10000	36.10	<b>31.53</b>	39.13	37.83	37.13	34.83	<b>31.80</b>	<b>30.13*</b>
		2000	20000	37.23	<b>32.17</b>	33.80	36.10	37.97	34.33	<b>31.60</b>	<b>30.07*</b>
Powersupply	LGBM	200	2000	<b>80.43*</b>	85.20	85.33	85.60	84.70	85.10	85.53	<b>83.43</b>
		500	5000	<b>79.60*</b>	83.60	83.70	84.10	83.77	83.43	84.73	<b>81.80</b>
		1000	10000	<b>82.57*</b>	<b>84.20</b>	86.17	<b>85.57</b>	85.00	<b>84.10</b>	<b>85.37</b>	<b>83.30</b>
		2000	20000	<b>81.63</b>	<b>80.53*</b>	82.93	<b>82.00</b>	<b>81.87</b>	<b>80.57</b>	<b>81.53</b>	<b>82.07</b>
	NN	200	2000	<b>85.30</b>	<b>83.27</b>	<b>83.27</b>	<b>83.60</b>	<b>84.70</b>	<b>86.20</b>	<b>83.90</b>	<b>82.23*</b>
		500	5000	<b>81.33</b>	82.57	82.77	82.93	<b>80.63</b>	<b>81.87</b>	83.90	<b>78.83*</b>
		1000	10000	<b>80.67</b>	83.77	85.07	83.80	<b>81.17</b>	<b>82.47</b>	<b>84.03</b>	<b>79.77*</b>
		2000	20000	<b>78.40</b>	<b>78.77</b>	<b>79.13</b>	<b>78.37</b>	<b>77.60*</b>	<b>78.60</b>	<b>79.67</b>	<b>78.23</b>
Forest	LGBM	200	2000	34.00	<b>3.13</b>	<b>4.37</b>	12.40	31.83	9.93	<b>2.77*</b>	<b>2.93</b>
		500	5000	<b>14.77</b>	<b>4.40</b>	<b>4.00*</b>	<b>4.33</b>	15.47	<b>5.40</b>	<b>4.20</b>	<b>4.80</b>
		1000	10000	<b>3.07</b>	<b>3.43</b>	<b>3.63</b>	<b>3.00</b>	<b>3.43</b>	<b>3.53</b>	3.60	<b>2.77*</b>
		2000	20000	<b>5.07</b>	6.67	6.50	<b>5.00</b>	<b>6.20</b>	6.40	7.03	<b>4.43*</b>
	NN	200	2000	35.90	<b>3.80</b>	<b>6.17</b>	16.10	32.70	9.40	<b>3.43*</b>	<b>4.07</b>
		500	5000	14.53	<b>4.73</b>	<b>5.03</b>	<b>3.93*</b>	15.40	5.20	<b>4.43</b>	<b>3.97</b>
		1000	10000	<b>4.00</b>	5.07	5.20	<b>4.20</b>	<b>4.27</b>	<b>4.87</b>	<b>4.70</b>	<b>3.60*</b>
		2000	20000	<b>5.83</b>	8.60	8.57	<b>6.27</b>	8.17	8.77	8.87	<b>5.63*</b>
Average Rank				5.27	3.40	4.93	4.43	6.00	5.03	4.37	<b>2.40*</b>
#Best				3	3	1	4	1	0	3	<b>15*</b>
#Best or Comparable				14	19	13	18	10	14	20	<b>30*</b>

**Setting.** We vary  $N$  among 200, 500, 1000, and 2000, setting  $T = 10N$ . The number of test data,  $M$ , is consistently set to 100 across all settings. In each dataset, we select continuous  $T$  samples starting from a randomly chosen index and use them for  $D$ . The subsequent  $M$  samples form the test data,  $D^{\text{te}}$ . Using each baseline and our method, we select the training data from  $D$  and then train a classifier  $\hat{h}$ . The classifier is either modeled by LightGBM (Ke et al., 2017) or a three-layer neural network with 100 hidden units as two representative classification models. The evaluation is based on the average zero-one loss on  $D^{\text{te}}$ , i.e.,  $\hat{R}_{01}(\hat{h}; D^{\text{te}}) = \frac{1}{|D^{\text{te}}|} \sum_{(\mathbf{x}, y) \in D^{\text{te}}} \ell_{01}(\hat{h}(\mathbf{x}), y)$ . We repeat each setting for 30 times with different random seeds, and report the average.

**Comparison methods.** We compare our method with seven various baselines, including naive baselines, drift detection, covariate shift adaptation, and drift localization methods as follows;

- $D_T$ ,  $D$ : Naive baselines. Naively use each of the recent data  $D_T$  and whole data  $D$ .
- PHT (page-hinkley test) (Page, 1954), ADWIN (Bifet and Gavalda, 2007): Representative time-based drift detection methods. First train a LightGBM classification model on  $D_S$  and then apply drift detection to the prediction loss from the present to the past. We select samples until a drift is detected.
- uLSIF (Kanamori et al., 2009): Covariate shift adaptation method; a variant of our approach, not with *joint* density ratio, but with covariate density ratio. Efficient hyperparameter tuning proposed by the authors is conducted for each experiment.
- LDD-DSDA (Liu et al., 2017): An existing method for training data selection, which selects samples based on drift localization method, LDD-DIS. Default parameters provided by the authors are used.
- LCD (Hinder et al., 2022): A drift localization method; we use all  $D_T$  and *no-drifting* samples ( $p$ -value of drift  $\geq 0.05$ ) in  $D_S$ . Parameters provided by the authors are used.
- TSJD: Our method detailed in Section 3. Hyperparameters are pre-tuned for each dataset and  $N$  pair using the entire dataset based on Section 3.4.

**Results.** The results are presented in Table 2. Our method achieves the best average rank of 2.40 and consistently shows the best or comparable results across all datasets and settings. Among the baselines, LCD achieves the best or comparable results 20 times. However, its average rank is 4.37, which is worse than the naive baseline using  $D$ . This highlights the weakness of LCD and underscores the superiority of TSJD. Overall, these findings empirically demonstrate the effectiveness and versatility of our method for the problem of training data selection.

## 6. Conclusion

This paper studied the training data selection problem, focusing on the selection of effective samples to improve model training from drifting data. We proposed TSJD, which assigns training weights for each sample based on joint density ratio estimation. We provide a theoretical analysis that bounds the generalization error of our method. Extensive experiments on real-world datasets demonstrate the superiority of TSJD over baseline methods.

## References

- Pranjal Awasthi, Corinna Cortes, and Mehryar Mohri. Best-effort adaptation. *Annals of Mathematics and Artificial Intelligence*, 92(2), 2024. doi: 10.1007/s10472-023-09917-3.
- Raef Bassily, Corinna Cortes, Anqi Mao, and Mehryar Mohri. Differentially private domain adaptation with theoretical guarantees. In *International Conference on Machine Learning*, 2024.
- Albert Bifet and Ricard Gavalda. Learning from time-changing data with adaptive windowing. In *SDM*, 2007.

- Dariusz Brzezinski and Jerzy Stefanowski. Reacting to different types of concept drift: The accuracy updated ensemble algorithm. *IEEE Transactions on Neural Networks and Learning Systems*, 25(1), 2014.
- Corinna Cortes, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. Structured prediction theory based on factor graph complexity. In *International Conference on Neural Information Processing Systems*, 2016.
- David Albert Edwards. On the kantorovich–rubinstein theorem. *Expositiones Mathematicae*, 29(4), 2011. doi: 10.1016/j.exmath.2011.06.005.
- João Gama, Pedro Medas, Gladys Castillo, and Pedro Rodrigues. Learning with drift detection. In *Advances in Artificial Intelligence*, 2004. doi: 10.1007/978-3-540-28645-5\_29.
- Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477), 2007. doi: 10.1198/016214506000001437.
- Paulo M. Gonçalves, Silas G.T. de Carvalho Santos, Roberto S.M. Barros, and Davi C.L. Vieira. A comparative study on concept drift detectors. *Expert Systems with Applications*, 41(18), 2014. doi: 10.1016/j.eswa.2014.07.019.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *International Conference on Neural Information Processing Systems*, 2014.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*, volume 1. Springer Series in Statistics, 2001.
- Fabian Hinder, Valerie Vaquet, Johannes Brinkrolf, André Artelt, and Barbara Hammer. Localization of concept drift: Identifying the drifting datapoints. In *International Joint Conference on Neural Networks*, 2022. doi: 10.1109/IJCNN55064.2022.9892374.
- Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. A least-squares approach to direct importance estimation. *J. Mach. Learn. Res.*, 10, 2009.
- Masahiro Kato and Takeshi Teshima. Non-negative bregman divergence minimization for deep direct density ratio estimation. In *International Conference on Machine Learning*, volume 139, 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. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- V. Koltchinskii. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory*, 47(5), 2001.
- Dominik Kreuzberger, Niklas Kühl, and Sebastian Hirschl. Machine learning operations (mlops): Overview, definition, and architecture. *IEEE Access*, 11, 2023.

- Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Classics in Mathematics. Springer Berlin Heidelberg, 2013.
- Anjin Liu, Yiliao Song, Guangquan Zhang, and Jie Lu. Regional concept drift detection and density synchronized drift adaptation. In *International Joint Conference on Artificial Intelligence*, 2017. doi: 10.24963/ijcai.2017/317.
- Song Liu, Makoto Yamada, Nigel Collier, and Masashi Sugiyama. Change-point detection in time-series data by relative density-ratio estimation. *Neural Networks*, 43, 2013. doi: 10.1016/j.neunet.2013.01.012.
- Yukitoshi Matsushita, Taisuke Otsu, and Keisuke Takahata. Estimating density ratio of marginals to joint: Applications to causal inference. *Journal of Business and Economic Statistics*, 2(41), 2022.
- Andreas Maurer. A vector-contraction inequality for rademacher complexities. In *Algorithmic Learning Theory*, 2016.
- Mansour Zoubeirou A. Mayaki and Michel Riveill. Autoregressive based drift detection method. *International Joint Conference on Neural Networks*, 2022.
- Hassan Mehmood, Panos Kostakos, Marta Cortes, Theodoros Anagnostopoulos, Susanna Pirttikangas, and Ekaterina Gilman. Concept drift adaptation techniques in distributed environment for real-world data streams. *Smart Cities*, 4(1), 2021. doi: 10.3390/smartcities4010021.
- Mehryar Mohri and Andres Muñoz Medina. New analysis and algorithm for learning with drifting distributions. In *Proceedings of the 23rd International Conference on Algorithmic Learning Theory*, 2012. doi: 10.1007/978-3-642-34106-9\_13.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. MIT Press, second edition, 2018.
- E. S. Page. Continuous inspection schemes. *Biometrika*, 41, 1954.
- Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D. Lawrence. *Dataset Shift in Machine Learning*. The MIT Press, 2009.
- Philipp Ruf, Manav Madan, Christoph Reich, and Djaffar Ould-Abdeslam. Demystifying mlops and presenting a recipe for the selection of open-source tools. *Applied Sciences*, 2021.
- Bernhard Schölkopf, John Platt, and Thomas Hofmann. Correcting sample selection bias by unlabeled data. In *Advances in Neural Information Processing Systems 19: Proceedings of the 2006 Conference*, 2007.
- Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90(2), 2000.

- Vinicius M. A. Souza, Denis M. dos Reis, André Gustavo Maletzke, and Gustavo E. A. P. A. Batista. Challenges in benchmarking stream learning algorithms with real-world data. *Data Mining and Knowledge Discovery*, 34, 2020.
- Masashi Sugiyama. Superfast-trainable multi-class probabilistic classifier by least-squares posterior fitting. *IEICE TRANSACTIONS on Information*, E93-D(10), 2010. doi: 10.1587/transinf.E93.D.2690.
- Masashi Sugiyama, Shinichi Nakajima, Hisashi Kashima, Paul von Büna, and Motoaki Kawanabe. Direct importance estimation with model selection and its application to covariate shift adaptation. In *Neural Information Processing Systems*, 2007a.
- Masashi Sugiyama, Shinichi Nakajima, Hisashi Kashima, Paul Vonbunau, and Motoaki Kawanabe. Kullback-leibler importance estimation procedure for covariate shift adaptation. *JSAT Technical Report, Type 2 SIG*, (DMSM-A702), 2007b. doi: 10.11517/jsaisigtwo.2007.DMSM-A702\_03.
- Masashi Sugiyama, Ichiro Takeuchi, Taiji Suzuki, Takafumi Kanamori, Hirotaka Hachiya, and Daisuke Okanohara. Conditional density estimation via least-squares density ratio estimation. In *International Conference on Artificial Intelligence and Statistics*, volume 9, 2010.
- Masashi Sugiyama, Teruyuki Suzuki, and Takafumi Kanamori. Density-ratio matching under the bregman divergence: a unified framework of density-ratio estimation. *Annals of the Institute of Statistical Mathematics*, 64, 2012.
- Taiji Suzuki, Masashi Sugiyama, and Toshiyuki Tanaka. Mutual information approximation via maximum likelihood estimation of density ratio. In *IEEE International Symposium on Information Theory*, 2009. doi: 10.1109/ISIT.2009.5205712.
- Haixun Wang, Wei Fan, Philip S. Yu, and Jiawei Han. Mining concept-drifting data streams using ensemble classifiers. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2003.
- Frank Wilcoxon. Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6), 1945.
- Michał Woźniak. Application of combined classifiers to data stream classification. In *Computer Information Systems and Industrial Management*, 2013.
- Makoto Yamada, Taiji Suzuki, Takafumi Kanamori, Hirotaka Hachiya, and Masashi Sugiyama. Relative density-ratio estimation for robust distribution comparison. *Neural Computation*, 25(5), 2013. doi: 10.1162/NECO\_a.00442.
- Tianyi Zhang, Ikko Yamane, Nan Lu, and Masashi Sugiyama. A one-step approach to covariate shift adaptation. In *Asian Conference on Machine Learning*, volume 129, 2020.
- Yu-Jie Zhang, Zhen-Yu Zhang, Peng Zhao, and Masashi Sugiyama. Adapting to continuous covariate shift via online density ratio estimation. In *International Conference on Neural Information Processing Systems*, 2023.



## Supplementary Material

This is the supplementary material of the paper “Training Data Soft Selection via Joint Density Ratio Estimation”. We provide useful lemmas in Supplementary A and missing proofs of theorems and lemmas presented in the main body of the paper Supplementary B. Supplementary C presents detailed experimental results, including running time and analysis of hyperparameter sensitivity.

### Supplementary A. Useful Lemmas

**Lemma A.1** *Let  $p, q$  be distributions over  $\mathcal{X}$  and  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a Lipschitz continuous differentiable function. Then, following inequality holds.*

$$\left| \mathbb{E}_{p(\mathbf{x})} [f(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x})} [f(\mathbf{x})] \right| \leq \left( \sup_{\mathbf{x} \in \mathcal{X}} \|\nabla f(\mathbf{x})\|_2 \right) W_1(p, q) \quad (\text{A.1})$$

**Proof.** Obvious by Kantorovich-Rubenstein duality (Edwards, 2011).  $\square$

**Proposition A.2** *For any  $p, q \in [0, 1]$ , the following holds.*

$$(p - q)^2 \leq D_{\text{BKL}}(p||q), \quad (\text{A.2})$$

where  $D_{\text{BKL}}$  is defined as the binary Kullback-Leibler (KL) divergence, defined with the binary cross entropy (CE)  $H_B$  as

$$D_{\text{BKL}}(p||q) := H_B(p, q) - H_B(p, p), \quad (\text{A.3})$$

$$H_B(p, q) := -p \log q - (1 - p) \log(1 - q). \quad (\text{A.4})$$

**Proof.** Let  $p$  be fixed to any number in  $[0, 1]$ . We aim at finding the minimum of the difference  $d(q)$ , defined as

$$d(q) := D_{\text{BKL}}(p||q) - (p - q)^2 = p \log \frac{p}{q} + (1 - p) \log \frac{1 - p}{1 - q} - (p - q)^2. \quad (\text{A.5})$$

Regarding the differential, we have the following.

$$\frac{d}{dq} d(q) = (q - p) \times \frac{1 + (-1 + 2q)^2}{2q(1 - q)} \quad (\text{A.6})$$

By letting  $\frac{d}{dq} d(q) = 0$ , we see that  $d(q)$  is minimized when  $p = q$  and the minimum is 0. Hence we have  $d(q) \geq 0$ , which concludes the proof.  $\square$

**Proposition A.3** *For any  $\mathbf{p}, \mathbf{q} \in \Delta^{K-1}$ ,  $i \in [K]$ , the following inequality holds.*

$$(p_i - q_i)^2 \leq D_{\text{BKL}}(p_i, q_i) \leq D_{\text{KL}}(\mathbf{p}, \mathbf{q}), \quad (\text{A.7})$$

where  $D_{\text{KL}}$  is the KL divergence, defined with the cross entropy (CE)  $H$  as

$$D_{\text{KL}}(\mathbf{p}, \mathbf{q}) := H(\mathbf{p}, \mathbf{q}) - H(\mathbf{p}, \mathbf{p}), \quad (\text{A.8})$$

$$H(\mathbf{p}, \mathbf{q}) := - \sum_{k=1}^K p_k \log q_k. \quad (\text{A.9})$$

**Proof.** By Proposition A.2,  $(p_i - q_i)^2 \leq D_{\text{BKL}}(p_i, q_i)$  holds.

$$D_{\text{KL}}(\mathbf{p}, \mathbf{q}) = \sum_{k=1}^K p_k \log \frac{p_k}{q_k} = p_i \log \frac{p_i}{q_i} + \sum_{k \neq i} p_k \log \frac{p_k}{q_k} \quad (\text{A.10})$$

By the log sum inequality, we have

$$\sum_{k \neq i} p_k \log \frac{p_k}{q_k} \geq (1 - p_i) \log \frac{1 - p_i}{1 - q_i}. \quad (\text{A.11})$$

Combining these we have

$$D_{\text{KL}}(\mathbf{p}, \mathbf{q}) \geq p_i \log \frac{p_i}{q_i} + (1 - p_i) \log \frac{1 - p_i}{1 - q_i} = D_{\text{BKL}}(p_i, q_i) \quad (\text{A.12})$$

as desired to conclude the proof.  $\square$

**Lemma A.4** *For any  $\mathbf{p}, \mathbf{q} \in \Delta^{K-1}$ ,  $i \in [K]$ , the following inequality holds.*

$$\frac{1}{K} \|\mathbf{p} - \mathbf{q}\|_2^2 \leq D_{\text{KL}}(\mathbf{p}, \mathbf{q}) \quad (\text{A.13})$$

**Proof.** By Proposition A.3, we have

$$\frac{1}{K} \|\mathbf{p} - \mathbf{q}\|_2^2 = \frac{1}{K} \sum_{i=1}^K (p_i - q_i)^2 \leq \frac{1}{K} \sum_{i=1}^K D_{\text{KL}}(\mathbf{p}, \mathbf{q}) = D_{\text{KL}}(\mathbf{p}, \mathbf{q}), \quad (\text{A.14})$$

as desired.  $\square$

## Supplementary B. Missing Proofs

### B.1. Proof of Theorem 3

We provide two lemmas, Lemma B.5 and Lemma B.6 below, where the proof of Theorem 3 is based.

**Lemma B.5** *Assume that  $\mathbf{x} \mapsto \bar{p}_T(y|\mathbf{x})$  is differentiable w.r.t.  $\mathbf{x}$  for any  $\mathbf{x} \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . Then, the following inequality holds for any  $h \in \mathcal{H}$ .*

$$\frac{1}{2} \mathbb{E}_{p_{\text{te}}(\mathbf{x})} \left[ \|p_{\text{te}}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] \leq \mathbb{E}_{\bar{p}_T(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] + Z(h) \tau_X + \tau_{Y|X}^2. \quad (\text{B.1})$$

**Proof.** By the Cauchy–Schwarz inequality, we have

$$\|p_{\text{te}}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 = 2 \left( \|p_{\text{te}}(\cdot|\mathbf{x}) - \bar{p}_T(\cdot|\mathbf{x})\|_2^2 + \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right). \quad (\text{B.2})$$

Based on this, we have

$$\frac{1}{2} \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \|p_{te}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] \leq \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \|p_{te}(\cdot|\mathbf{x}) - \bar{p}_T(\cdot|\mathbf{x})\|_2^2 \right] + \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] \quad (\text{B.3})$$

$$\leq \tau_{Y|X}^2 + \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right], \quad (\text{B.4})$$

where  $\|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2 \leq \tau_{Y|X}$  by Assumption 1. By applying Lemma A.1 to the second term in the right hand of Eq. (B.4), we have

$$\begin{aligned} & \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] \\ & \leq \mathbb{E}_{\bar{p}_T(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] + \left( \sup_{\mathbf{x} \in \mathcal{X}} \left\| \nabla \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right\|_2 \right) W_1(p_{te}, \bar{p}_T) \end{aligned} \quad (\text{B.5})$$

$$\leq \mathbb{E}_{\bar{p}_T(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] + Z(h) \tau_X. \quad (\text{B.6})$$

Combining Eq. (B.4) and Eq. (B.6), we obtain

$$\frac{1}{2} \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \|p_{te}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] \leq \mathbb{E}_{\bar{p}_T(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] + Z(h) \tau_X + \tau_{Y|X}^2 \quad (\text{B.7})$$

as desired.  $\square$

**Lemma B.6 (Relation between zero-one loss and  $L^2$ -norm.)** *Let us define  $\max_{(2)}$  as the operator to take the second best, i.e., for some  $f : \mathcal{Y} \rightarrow \mathbb{R}$ ,*

$$\max_{k \in \mathcal{Y}}^{(2)} f(k) := \max_{k \in \{y \in \mathcal{Y} \mid f(y) \neq \max_{k'} f(k')\}} f(k), \quad (\text{B.8})$$

*and use  $\arg \max_{(2)}$  to take the point of the second best. Assume that  $\eta_{min} > 0$  where*

$$\eta_{min} := \min_{\mathbf{x} \in \mathcal{X}} \left( \max_k p_{te}(k|\mathbf{x}) - \max_{(2)} p_{te}(k|\mathbf{x}) \right). \quad (\text{B.9})$$

*Then, for any  $h \in \mathcal{H}$ , the following holds.*

$$R_{01}(h) \leq B_{01} + 2\eta_{min}^{-2} \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \|p_{te}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] \quad (\text{B.10})$$

*where we define  $B_{01}$  as the Bayes error, i.e.,  $B_{01} := \min_f R_{01}(f)$  is the lowest  $R_{01}$  among any possible  $f : \mathcal{X} \rightarrow \Delta^{K-1}$ .*

**Proof.**

$$R_{01}(h) = \mathbb{E}_{p_{te}(\mathbf{x}, y)} \left[ \mathbb{I} \left[ y \neq \arg \max_{k \in \mathcal{Y}} h(k|\mathbf{x}) \right] \right] = \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \mathbb{E}_{p_{te}(y|\mathbf{x})} \left[ \mathbb{I} \left[ y \neq \arg \max_{k \in \mathcal{Y}} h(k|\mathbf{x}) \right] \right] \right] \quad (\text{B.11})$$

It is obvious that

$$\begin{aligned} & (y = \arg \max_k p_{\text{te}}(k|\mathbf{x})) \wedge (\arg \max_k p_{\text{te}}(k|\mathbf{x}) = \arg \max_k h(k|\mathbf{x})) \\ & \Rightarrow (y = \arg \max_k h(k|\mathbf{x})). \end{aligned} \quad (\text{B.12})$$

By taking its contraposition, we have

$$\begin{aligned} & (y \neq \arg \max_k h(k|\mathbf{x})) \\ & \Rightarrow (y \neq \arg \max_k p_{\text{te}}(k|\mathbf{x})) \vee (\arg \max_k p_{\text{te}}(k|\mathbf{x}) \neq \arg \max_k h(k|\mathbf{x})). \end{aligned} \quad (\text{B.13})$$

Hence, we have

$$\begin{aligned} & \mathbb{E}_{p_{\text{te}}(y|\mathbf{x})} \left[ \mathbb{I} \left[ y \neq \arg \max_{k \in \mathcal{Y}} h(k|\mathbf{x}) \right] \right] \\ & \leq \mathbb{E}_{p_{\text{te}}(y|\mathbf{x})} \left[ \mathbb{I} \left[ y \neq \arg \max_k p_{\text{te}}(k|\mathbf{x}) \right] \right] + \mathbb{I} \left[ \arg \max_k p_{\text{te}}(k|\mathbf{x}) \neq \arg \max_k h(k|\mathbf{x}) \right] \end{aligned} \quad (\text{B.14})$$

The first term is the Bayes error. Regarding the second term, the best way to achieve  $\arg \max_k p_{\text{te}}(k|\mathbf{x}) \neq \arg \max_k h(k|\mathbf{x})$  with minimum  $\|p_{\text{te}}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2$  is to make  $h(\mathbf{x})$  to be

$$h(\arg \max_k p_{\text{te}}(k|\mathbf{x})|\mathbf{x}) = \max_k p_{\text{te}}(k|\mathbf{x}) - \frac{\eta(\mathbf{x})}{2} \quad (\text{B.15})$$

$$h(\arg \max_{(2)} p_{\text{te}}(k|\mathbf{x})|\mathbf{x}) = \max_{(2)} p_{\text{te}}(k|\mathbf{x}) + \frac{\eta(\mathbf{x})}{2} \quad (\text{B.16})$$

where

$$\eta(\mathbf{x}) := \max_k p_{\text{te}}(k|\mathbf{x}) - \max_{(2)} p_{\text{te}}(k|\mathbf{x}). \quad (\text{B.17})$$

Hence, we have

$$\mathbb{I} \left[ \arg \max_k p_{\text{te}}(k|\mathbf{x}) \neq \arg \max_k h(k|\mathbf{x}) \right] \leq \frac{2}{\eta(\mathbf{x})^2} \|p_{\text{te}}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2. \quad (\text{B.18})$$

Combining these, we have

$$R_{01}(h) \leq B_{01} + 2\eta_{\min}^{-2} \mathbb{E}_{p_{\text{te}}(\mathbf{x})} \left[ \|p_{\text{te}}(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right], \quad (\text{B.19})$$

as desired.  $\square$

Finally, we prove Theorem 3 as follows.

**Proof** (Proof of Theorem 3). We obtain Theorem 3 by direct combination of Lemma B.5 and Lemma B.6.  $\square$

## B.2. Proof of Theorem 10

Before proving Theorem 10, we establish Lemma B.7. The proof is inspired by Lemma 3 in (Zhang et al., 2020).

**Lemma B.7** *For any  $\delta \in (0, 1)$  and any  $g \in \mathcal{G}_+$ , over the draws of  $S_S$  from  $\bar{p}_S$  and  $S_T$  from  $\bar{p}_T$ , the following inequality holds with probability at least  $1 - \delta$ ;*

$$J(g) - \hat{J}(g; S_S, S_T) \leq 4G\mathfrak{R}_{|S_S|}(\mathcal{G}) + 4\mathfrak{R}_{|S_T|}(\mathcal{G}) + 2G^2 \sqrt{\frac{\log \frac{1}{\delta}}{2}} \left( \frac{1}{\sqrt{|S_S|}} + \frac{1}{\sqrt{|S_T|}} \right). \quad (\text{B.20})$$

**Proof.** Let  $Z := (S_S, S_T)$ ,  $\Phi(Z) := \sup_{g \in \mathcal{G}} (J(g) - \hat{J}(g; S_S, S_T))$  and  $Z' = (S'_p, S'_q)$  differ exactly one sample from  $Z$ , i.e.,  $Z \setminus Z' := (S_S \setminus S'_p) \cup (S_T \setminus S'_q) = \{z\}$  and  $Z' \setminus Z = \{z'\}$ . Since the difference of suprema does not exceed the supremum of the difference, we have

$$\Phi(Z') - \Phi(Z) = \sup_{g \in \mathcal{G}_+} (J(g) - \hat{J}(g; Z')) - \sup_{g \in \mathcal{G}} (J(g) - \hat{J}(g; Z)) \quad (\text{B.21})$$

$$\leq \sup_{g \in \mathcal{G}} \left( (J(g) - \hat{J}(g; Z')) - (J(g) - \hat{J}(g; Z)) \right) \quad (\text{B.22})$$

$$= \sup_{g \in \mathcal{G}} (\hat{J}(g; Z) - \hat{J}(g; Z')). \quad (\text{B.23})$$

If  $z \in S_S$ ,

$$\Phi(Z') - \Phi(Z) = \frac{\sup_{g \in \mathcal{G}_+} g(z')^2 - g(z)^2}{|S_S|} \leq \frac{1}{|S_S|} G^2 \leq \frac{2G^2}{|S_S|} \quad (\text{B.24})$$

otherwise  $z \in S_T$ ,

$$\Phi(Z') - \Phi(Z) = \frac{\sup_{g \in \mathcal{G}_+} 2g(z) - 2g(z')}{|S_T|} \leq \frac{2}{|S_T|} G \leq \frac{2G^2}{|S_T|} \quad (\text{B.25})$$

Hence, by McDiarmid's inequality, the following holds with probability at least  $1 - \delta$ ;

$$\Phi(Z) \leq \mathbb{E}_Z[\Phi(Z)] + \sqrt{\frac{\log \frac{1}{\delta}}{2} \left( |S_S| \frac{4G^4}{|S_S|^2} + |S_T| \frac{4G^4}{|S_T|^2} \right)} \quad (\text{B.26})$$

$$= \mathbb{E}_Z[\Phi(Z)] + 2G^2 \sqrt{\frac{\log \frac{1}{\delta}}{2} \left( \frac{1}{\sqrt{|S_S|}} + \frac{1}{\sqrt{|S_T|}} \right)} \quad (\text{B.27})$$

Next, we bound  $\mathbb{E}_Z[\Phi(Z)]$  from above. By the fact that the supremum of sum is equal to or less than sum of suprema, we have

$$\Phi(Z) = \sup_{g \in \mathcal{G}_+} \left( J(g) - \hat{J}(g; S_S, S_T) \right) \quad (\text{B.28})$$

$$\begin{aligned} &= \sup_{g \in \mathcal{G}_+} \left( \left( \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} \left[ (g(\mathbf{x}, y))^2 \right] - \frac{1}{|S_S|} \sum_{(\mathbf{x}, y) \in S_S} g(\mathbf{x}, y)^2 \right) \right. \\ &\quad \left. + 2 \left( \frac{1}{|S_T|} \sum_{(\mathbf{x}, y) \in S_T} g(\mathbf{x}, y) - \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [g(\mathbf{x}, y)] \right) \right) \end{aligned} \quad (\text{B.29})$$

$$\begin{aligned} &\leq \sup_{g \in \mathcal{G}_+} \left( \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} \left[ (g(\mathbf{x}, y))^2 \right] - \frac{1}{|S_S|} \sum_{(\mathbf{x}, y) \in S_S} g(\mathbf{x}, y)^2 \right) \\ &\quad + 2 \sup_{g \in \mathcal{G}_+} \left( \frac{1}{|S_T|} \sum_{(\mathbf{x}, y) \in S_T} g(\mathbf{x}, y) - \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [g(\mathbf{x}, y)] \right). \end{aligned} \quad (\text{B.30})$$

Hence, We have

$$\begin{aligned} \mathbb{E}_Z[\Phi(Z)] &\leq \mathbb{E}_{S_S} \left[ \sup_{g \in \mathcal{G}_+} \left( \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} \left[ (g(\mathbf{x}, y))^2 \right] - \frac{1}{|S_S|} \sum_{(\mathbf{x}, y) \in S_S} g(\mathbf{x}, y)^2 \right) \right] \\ &\quad + 2 \mathbb{E}_{S_T} \left[ \sup_{g \in \mathcal{G}_+} \left( \frac{1}{|S_T|} \sum_{(\mathbf{x}, y) \in S_T} g(\mathbf{x}, y) - \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [g(\mathbf{x}, y)] \right) \right] \end{aligned} \quad (\text{B.31})$$

$$= A + 2B, \quad (\text{B.32})$$

where we define  $A$  and  $B$  as

$$A := \mathbb{E}_{S_S} \left[ \sup_{g \in \mathcal{G}_+} \left( \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} \left[ (g(\mathbf{x}, y))^2 \right] - \frac{1}{|S_S|} \sum_{(\mathbf{x}, y) \in S_S} g(\mathbf{x}, y)^2 \right) \right] \quad (\text{B.33})$$

$$B := \mathbb{E}_{S_T} \left[ \sup_{g \in \mathcal{G}_+} \left( \frac{1}{|S_T|} \sum_{(\mathbf{x}, y) \in S_T} g(\mathbf{x}, y) - \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [g(\mathbf{x}, y)] \right) \right]. \quad (\text{B.34})$$

By the proof of lemma 3 (II) and (III) in (Zhang et al., 2020), we have upper bounds of  $A$  and  $B$  as

$$A \leq 4G\mathfrak{R}_{|S_S|}(\mathcal{G}) \quad (\text{B.35})$$

$$B \leq 2\mathfrak{R}_{|S_T|}(\mathcal{G}). \quad (\text{B.36})$$

Combined Eq. (B.27), Eq. (B.32), Eq. (B.35), and Eq. (B.36), we have

$$\sup_{g \in \mathcal{G}_+} \left( J(g) - \hat{J}(g; S_S, S_T) \right) \leq 4G\mathfrak{R}_{|S_S|}(\mathcal{G}) + 4\mathfrak{R}_{|S_T|}(\mathcal{G}) + 2G^2 \sqrt{\frac{\log \frac{1}{\delta}}{2}} \left( \frac{1}{\sqrt{|S_S|}} + \frac{1}{\sqrt{|S_T|}} \right) \quad (\text{B.37})$$

as desired.  $\square$

The proof of Theorem 10 is as follows.

**Proof** (Proof of Theorem 10). We decompose  $J(\hat{g}) - J(g^*)$  as

$$\begin{aligned} J(\hat{g}) - J(g^*) &= \left( J(\hat{g}) - \hat{J}(\hat{g}; D_S, D_T) \right) + \left( \hat{J}(\hat{g}; D_S, D_T) - \hat{J}(g^*; D_S, D_T) \right) \\ &\quad + \left( \hat{J}(g^*; D_S, D_T) - J(g^*) \right) \end{aligned} \quad (\text{B.38})$$

The first term is upper bounded by Lemma B.7 with probability at least  $1 - \frac{\delta}{3}$  as

$$\begin{aligned} J(\hat{g}) - \hat{J}(\hat{g}; D_S, D_T) &\leq 4G\mathfrak{R}_{|D_S|}(\mathcal{G}) + 4\mathfrak{R}_{|D_T|}(\mathcal{G}) + 2G^2 \sqrt{\frac{\log \frac{3}{\delta}}{2}} \left( \frac{1}{\sqrt{|D_S|}} + \frac{1}{\sqrt{|D_T|}} \right). \end{aligned} \quad (\text{B.39})$$

The second term is at most 0 by the definition of  $\hat{g}$ . The third term is upper bounded by Hoeffding's inequality and the following holds with probability at least  $1 - \frac{2\delta}{3}$

$$\begin{aligned} \hat{J}(g^*; D_S, D_T) - J(g^*) &= \left( \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} g(\mathbf{x}, y)^2 - \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[g(\mathbf{x}, y)^2] \right) \\ &\quad + 2 \left( \frac{1}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} g(\mathbf{x}, y) - \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[g(\mathbf{x}, y)] \right) \end{aligned} \quad (\text{B.40})$$

$$\leq G^2 \sqrt{\frac{\log \frac{3}{\delta}}{2|D_S|}} + 2G \sqrt{\frac{\log \frac{3}{\delta}}{2|D_T|}} \quad (\text{B.41})$$

$$\leq 2G^2 \sqrt{\frac{\log \frac{3}{\delta}}{2}} \left( \frac{1}{\sqrt{|D_S|}} + \frac{1}{\sqrt{|D_T|}} \right). \quad (\text{B.42})$$

Combining these with the union bound, we obtain

$$J(\hat{g}) - J(g^*) \leq 4G\mathfrak{R}_{|D_S|}(\mathcal{G}) + 4\mathfrak{R}_{|D_T|}(\mathcal{G}) + 4G^2 \sqrt{\frac{\log \frac{3}{\delta}}{2}} \left( \frac{1}{\sqrt{|D_S|}} + \frac{1}{\sqrt{|D_T|}} \right), \quad (\text{B.43})$$

which concludes the proof.  $\square$

### B.3. Proof of Lemma 13

**Proof** (Proof of Lemma 13). Based on Theorem 3.3 in (Mohri et al., 2018), we have

$$\mathbb{E}_{\bar{p}_T(\mathbf{x}, y)}[\ell_{\text{CE}}(h(\mathbf{x}), y)] \leq \frac{1}{|S_T|} \sum_{(\mathbf{x}, y) \in D_T} \ell_{\text{CE}}(h(\mathbf{x}), y) + 2\mathfrak{R}_{|S_T|}(\mathcal{L}) + U \sqrt{\frac{\log \frac{1}{\delta}}{2|S_T|}}, \quad (\text{B.44})$$

where  $\mathcal{L} := \{(\mathbf{x}, y) \mapsto \ell_{\text{CE}}(h(\mathbf{x}), y) : h \in \mathcal{H}\}$ . Then based on the vector-contraction inequality for vector-valued Rademacher complexities (Maurer, 2016; Cortes et al., 2016), we have

$$\mathfrak{R}_{|S_T|}(\mathcal{L}) \leq \sqrt{2}L\mathfrak{R}_{|S_T|}(\mathcal{H}), \quad (\text{B.45})$$

where  $L = \sup_{h, \mathbf{x}, y} \|\nabla \ell_{\text{CE}}(h(\mathbf{x}), y)\|_2 = \sup_{h, \mathbf{x}, y} \frac{1}{h(y|\mathbf{x})} \leq \exp(U)$  by the assumption. Combining these, we conclude the proof.  $\square$

#### B.4. Proof of Lemma 14

We first establish Lemma B.8 for the proof of Lemma 14. As well as Lemma B.7, Lemma B.8 is related to Lemma 3 in (Zhang et al., 2020).

**Lemma B.8** *For any  $\delta \in (0, 1)$  with probability at least  $1 - \delta$  over the draw of i.i.d. samples  $S$  from  $\bar{p}_S$ , the following inequality holds for any  $h \in \mathcal{H}$  and  $g \in \mathcal{G}_+$ :*

$$\mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y)] \leq \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y) + C_2(\delta). \quad (\text{B.46})$$

**Proof** (Proof of Lemma B.8). Let us define  $V$  and  $\hat{V}$  as

$$V(h, g) := \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y)] \quad (\text{B.47})$$

$$\hat{V}(h, g; S) := \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y). \quad (\text{B.48})$$

Let  $\Phi(S) := \sup_{g \in \mathcal{G}} (V(h, g) - \hat{V}(h, g; S))$  and  $S'$  differ exactly one sample from  $S$ , i.e.,  $S \setminus S' = \{s\}$  and  $S' \setminus S = \{s'\}$ . Since the difference of suprema does not exceed the supremum of the difference, we have

$$\Phi(S) - \Phi(S') = \sup_{h \in \mathcal{H}, g \in \mathcal{G}_+} \hat{V}(g; S) - \sup_{h, g} \hat{V}(g; S') \quad (\text{B.49})$$

$$\leq \sup_{h \in \mathcal{H}, g \in \mathcal{G}_+} (\hat{V}(h, g; S) - \hat{V}(h, g; S')) \quad (\text{B.50})$$

$$= \frac{\sup_{h \in \mathcal{H}, g \in \mathcal{G}_+} g(s) \ell_{\text{CE}}(s) - g(s') \ell_{\text{CE}}(s')}{|S|} \leq \frac{MG}{|S|} \quad (\text{B.51})$$

Hence, by McDiarmid's inequality, the following holds with probability at least  $1 - \delta$ ;

$$\Phi(S) \leq \mathbb{E}_S[\Phi(S)] + MG \sqrt{\frac{\log \frac{1}{\delta}}{2|S|}} \quad (\text{B.52})$$



Next, we upper bound  $\mathbb{E}_S[\Phi(S)]$ . Let  $\bar{S}$  be another i.i.d. samples from  $\bar{p}_S$  whose size is  $|S|$  and  $\boldsymbol{\sigma} = \{\sigma_i\}_{i=1}^{|S|} \in \{-1, 1\}^{|S|}$  is a set of Rademacher random variables. We have

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[ \sup_{g \in \mathcal{G}_+} \left( \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y)] - \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y) \right) \right] \quad (\text{B.53})$$

$$= \mathbb{E}_S \left[ \sup_{g \in \mathcal{G}_+} \left( \mathbb{E}_{\bar{S}} \left[ \frac{1}{|\bar{S}|} \sum_{(\mathbf{x}, y) \in \bar{S}} g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y) \right] - \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y) \right) \right] \quad (\text{B.54})$$

$$\leq \mathbb{E}_{S, \bar{S}} \left[ \sup_{g \in \mathcal{G}_+} \left( \frac{1}{|\bar{S}|} \sum_{(\mathbf{x}, y) \in \bar{S}} g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y) - \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} g(\mathbf{x}, y) \ell_{\text{CE}}(h(\mathbf{x}), y) \right) \right] \quad (\text{B.55})$$

(The expectation of suprema exceeds the supremum of expectation.)

$$= \mathbb{E}_{S, \bar{S}} \left[ \sup_{g \in \mathcal{G}_+} \left( \frac{1}{|S|} \sum_{i=1}^{|S|} (g(\bar{\mathbf{x}}_i, \bar{y}_i) \ell_{\text{CE}}(h(\bar{\mathbf{x}}_i), \bar{y}_i) - g(\mathbf{x}_i, y_i) \ell_{\text{CE}}(h(\mathbf{x}_i), y_i)) \right) \right] \quad (\text{B.56})$$

$$\leq \mathbb{E}_{S, \bar{S}} \left[ \sup_{g \in \mathcal{G}_+} \left| \frac{1}{|S|} \sum_{i=1}^{|S|} (g(\bar{\mathbf{x}}_i, \bar{y}_i) \ell_{\text{CE}}(h(\bar{\mathbf{x}}_i), \bar{y}_i) - g(\mathbf{x}_i, y_i) \ell_{\text{CE}}(h(\mathbf{x}_i), y_i)) \right| \right] \quad (\text{B.57})$$

$$= \mathbb{E}_{S, \bar{S}, \boldsymbol{\sigma}} \left[ \sup_{g \in \mathcal{G}_+} \left| \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i (g(\bar{\mathbf{x}}_i, \bar{y}_i) \ell_{\text{CE}}(h(\bar{\mathbf{x}}_i), \bar{y}_i) - g(\mathbf{x}_i, y_i) \ell_{\text{CE}}(h(\mathbf{x}_i), y_i)) \right| \right] \quad (\text{B.58})$$

(By the property of Rademacher variables.)

$$\leq 2 \mathbb{E}_{S, \boldsymbol{\sigma}} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i g(\mathbf{x}_i, y_i) \ell_{\text{CE}}(h(\mathbf{x}_i), y_i) \right] \quad (\text{B.59})$$

$$= \mathbb{E}_{S, \boldsymbol{\sigma}} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i \left( (\ell_{\text{CE}}(h(\mathbf{x}_i), y_i) + g(\mathbf{x}_i, y_i))^2 - (\ell_{\text{CE}}(h(\mathbf{x}_i), y_i))^2 - (g(\mathbf{x}_i, y_i))^2 \right) \right] \quad (\text{B.60})$$

(By the identity:  $2ab = (a + b)^2 - a^2 - b^2$ .)

$$\begin{aligned} &\leq \mathbb{E}_{S, \boldsymbol{\sigma}} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i (\ell_{\text{CE}}(h(\mathbf{x}_i), y_i) + g(\mathbf{x}_i, y_i))^2 \right] \\ &+ \mathbb{E}_{S, \boldsymbol{\sigma}} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i (\ell_{\text{CE}}(h(\mathbf{x}_i), y_i))^2 \right] + \mathbb{E}_{S, \boldsymbol{\sigma}} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i (g(\mathbf{x}_i, y_i))^2 \right] \end{aligned} \quad (\text{B.61})$$

(To the next page.)

(From the previous page.)

$$\begin{aligned} &\leq 2(U + G) \left( \mathbb{E}_{S, \sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i \ell_{\text{CE}}(h(\mathbf{x}_i), y_i) + \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i g(\mathbf{x}_i, y_i) \right] \right) \\ &\quad + 2U \mathbb{E}_{S, \sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i \ell_{\text{CE}}(h(\mathbf{x}_i), y_i) \right] + 2G \mathbb{E}_{S, \sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{|S|} \sum_{i=1}^{|S|} \sigma_i g(\mathbf{x}_i, y_i) \right] \end{aligned} \quad (\text{B.62})$$

(By the Ledoux-Talagrand contraction lemma ([Ledoux and Talagrand, 2013](#)).)

$$\leq 2(U + G)(\exp(U)\mathfrak{R}_{|S|}(\mathcal{H}) + \mathfrak{R}_{|S|}(\mathcal{G})) + 2U \exp(U)\mathfrak{R}_{|S|}(\mathcal{H}) + 2G\mathfrak{R}_{|S|}(\mathcal{G}) \quad (\text{B.63})$$

(By the Ledoux-Talagrand contraction lemma.)

$$= 2(2U + G) \exp(U)\mathfrak{R}_{|S|}(\mathcal{H}) + 2(U + 2G)\mathfrak{R}_{|S|}(\mathcal{G}). \quad (\text{B.64})$$

Combining these, we have

$$\begin{aligned} &\sup_{g \in \mathcal{G}} (V(h, g) - \hat{V}(h, g; S)) \\ &\leq 2(2U + G) \exp(U)\mathfrak{R}_{|S|}(\mathcal{H}) + 2(U + 2G)\mathfrak{R}_{|S|}(\mathcal{G}) + MG \sqrt{\frac{\log \frac{1}{\delta}}{2|S|}}, \end{aligned} \quad (\text{B.65})$$

as desired.  $\square$

**Lemma B.9** *Assume that we have  $g \in \mathcal{G}_+$  such that  $\mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[g(\mathbf{x}, y)] = 1$  holds and  $\bar{p}_S(\mathbf{x}, y)g(\mathbf{x}, y)$  is a distribution over  $\mathcal{X} \times \mathcal{Y}$ . The following inequality holds for any  $h \in \mathcal{H}$ ;*

$$\mathbb{E}_{\bar{p}_T(\mathbf{x}, y)}[\ell_{\text{CE}}(h(\mathbf{x}), y)] - \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[\ell_{\text{CE}}(h(\mathbf{x}), y)g(\mathbf{x}, y)] \leq U \sqrt{\mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[(r(\mathbf{x}, y) - g(\mathbf{x}, y))^2]} \quad (\text{B.66})$$

**Proof.** By direct calculation, we have

$$\mathbb{E}_{\bar{p}_T(\mathbf{x}, y)}[\ell_{\text{CE}}(h(\mathbf{x}), y)] - \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[\ell_{\text{CE}}(h(\mathbf{x}), y)g(\mathbf{x}, y)] \quad (\text{B.67})$$

$$= \int (\bar{p}_T(\mathbf{x}, y) - \bar{p}_S(\mathbf{x}, y)g(\mathbf{x}, y)) \ell_{\text{CE}}(h(\mathbf{x}), y) d\mathbf{x}dy \quad (\text{B.68})$$

$$\leq \int |\bar{p}_T(\mathbf{x}, y) - \bar{p}_S(\mathbf{x}, y)g(\mathbf{x}, y)| \ell_{\text{CE}}(h(\mathbf{x}), y) d\mathbf{x}dy \quad (\text{B.69})$$

$$\leq U \int |\bar{p}_S(\mathbf{x}, y)r(\mathbf{x}, y) - \bar{p}_S(\mathbf{x}, y)g(\mathbf{x}, y)| d\mathbf{x}dy \quad (\text{B.70})$$

$$= U \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[|r(\mathbf{x}, y) - g(\mathbf{x}, y)|] \quad (\text{B.71})$$

$$\leq U \sqrt{\mathbb{E}_{\bar{p}_S(\mathbf{x}, y)}[(r(\mathbf{x}, y) - g(\mathbf{x}, y))^2]} \quad (\text{B.72})$$

as desired.  $\square$

Finally, we prove Lemma 14 as follows.

**Proof** (Proof of Lemma 14). The proof is obvious by the combination of Lemma B.8 and Lemma B.9.  $\square$

### B.5. Proof of Theorem 15

**Proof** (Proof of Theorem 15).

We decompose  $L_{CE}(\hat{h}) - L_{CE}(h^*)$  as

$$L_{CE}(\hat{h}) - L_{CE}(h^*) = \left( L_{CE}(\hat{h}) - \widehat{L}_{CE}(\hat{h}) \right) + \left( \widehat{L}_{CE}(\hat{h}) - \widehat{L}_{CE}(h^*) \right) + \left( \widehat{L}_{CE}(h^*) - L_{CE}(h^*) \right) \quad (\text{B.73})$$

$$\leq A + B, \quad (\text{B.74})$$

where the second term in the right hand of Eq. (B.73) does not exceed 0 by definition of  $\hat{h}$  and we denote  $A$  and  $B$  as

$$A := L_{CE}(\hat{h}) - \widehat{L}_{CE}(\hat{h}), \quad (\text{B.75})$$

$$B := \widehat{L}_{CE}(h^*) - L_{CE}(h^*). \quad (\text{B.76})$$

We upper bound each of  $A$  and  $B$  as follows.

The term  $A$  is further decomposed as

$$\begin{aligned} A &= \frac{1}{2} \left( L_{CE}(\hat{h}) - \frac{1}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} \ell_{CE}(\hat{h}(\mathbf{x}), y) \right) \\ &\quad + \frac{1}{2} \left( L_{CE}(\hat{h}) - \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} \widehat{g}(\mathbf{x}, y) \ell_{CE}(\hat{h}(\mathbf{x}), y) \right). \end{aligned} \quad (\text{B.77})$$

By Lemma 13, the first term is upper bounded with probability at least  $1 - \delta$  by

$$L_{CE}(\hat{h}) - \frac{1}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} \ell_{CE}(\hat{h}(\mathbf{x}), y) \leq C_1(\delta). \quad (\text{B.78})$$

By Lemma 14, the second term is upper bounded with probability at least  $1 - \delta$  by

$$L_{CE}(\hat{h}) - \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} \widehat{g}(\mathbf{x}, y) \ell_{CE}(\hat{h}(\mathbf{x}), y) \leq C_2(\delta) + U \sqrt{\frac{\mathbb{E}}{\bar{p}_S(\mathbf{x}, y)} \left[ (r(\mathbf{x}, y) - \widehat{g}(\mathbf{x}, y))^2 \right]}. \quad (\text{B.79})$$

By Theorem 10, with probability at least  $1 - \delta$ , we have

$$\frac{\mathbb{E}}{\bar{p}_S(\mathbf{x}, y)} \left[ (r(\mathbf{x}, y) - \widehat{g}(\mathbf{x}, y))^2 \right] \leq J(g^*) + C_3(\delta). \quad (\text{B.80})$$

Combining Eq. (B.78), Eq. (B.79), and Eq. (B.80), the term  $A$  is upper bounded with probability at least  $1 - 3\delta$  as

$$A \leq \frac{1}{2} \left( C_1(\delta) + C_2(\delta) + U \sqrt{J(g^*)} + U \sqrt{C_3(\delta)} \right). \quad (\text{B.81})$$

Next, we upper bound the term  $B$ .

$$\begin{aligned}
B &= \frac{1}{2} \left( \frac{1}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} \ell_{\text{CE}}(h^*(\mathbf{x}), y) - L_{\text{CE}}(h^*) \right) \\
&\quad + \frac{1}{2} \left( \frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} \widehat{g}(\mathbf{x}, y) \ell_{\text{CE}}(h^*(\mathbf{x}), y) - \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [\widehat{g}(\mathbf{x}, y) \ell_{\text{CE}}(h^*(\mathbf{x}, y))] \right) \\
&\quad + \frac{1}{2} \left( \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [\widehat{g}(\mathbf{x}, y) \ell_{\text{CE}}(h^*(\mathbf{x}, y))] - L_{\text{CE}}(h^*) \right) \tag{B.82}
\end{aligned}$$

The first term is bounded by Hoeffding's inequality. With probability at least  $1 - \delta$ , we have

$$\frac{1}{|D_T|} \sum_{(\mathbf{x}, y) \in D_T} \ell_{\text{CE}}(h^*(\mathbf{x}), y) - L_{\text{CE}}(h^*) \leq U \sqrt{\frac{\log \frac{1}{\delta}}{2|D_T|}}. \tag{B.83}$$

The second term is similarly bounded with probability at least  $1 - \delta$  as

$$\frac{1}{|D_S|} \sum_{(\mathbf{x}, y) \in D_S} \widehat{g}(\mathbf{x}, y) \ell_{\text{CE}}(h^*(\mathbf{x}), y) - \mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [\widehat{g}(\mathbf{x}, y) \ell_{\text{CE}}(h^*(\mathbf{x}, y))] \leq GU \sqrt{\frac{\log \frac{1}{\delta}}{2|D_S|}}. \tag{B.84}$$

The third term is bounded similarly with Eq. (B.80), i.e., by combining Lemma B.9 and Theorem 10, we have

$$\mathbb{E}_{\bar{p}_S(\mathbf{x}, y)} [\widehat{g}(\mathbf{x}, y) \ell_{\text{CE}}(h^*(\mathbf{x}, y))] - \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [\ell_{\text{CE}}(h^*(\mathbf{x}, y))] \leq U \sqrt{J(g^*)} + U \sqrt{C_3(\delta)}. \tag{B.85}$$

Combining Eq. (B.83), Eq. (B.84), and Eq. (B.85), the term  $B$  is bounded with probability at least  $1 - 2\delta$  as

$$B \leq \frac{1}{2} \left( U \sqrt{\frac{\log \frac{1}{\delta}}{2|D_T|}} + GU \sqrt{\frac{\log \frac{1}{\delta}}{2|D_S|}} + U \sqrt{J(g^*)} + U \sqrt{C_3(\delta)} \right). \tag{B.86}$$

Combining Eq. (B.81), Eq. (B.86), and the definition of  $C_1$  and  $C_2$ , we have

$$\begin{aligned}
&2(A + B) \\
&\leq C_1(\delta) + C_2(\delta) + 2U \sqrt{J(g^*)} + 2U \sqrt{C_3(\delta)} + U \sqrt{\frac{\log \frac{1}{\delta}}{2|D_T|}} + GU \sqrt{\frac{\log \frac{1}{\delta}}{2|D_S|}} \tag{B.87}
\end{aligned}$$

$$\begin{aligned}
&= 2U \sqrt{J(g^*)} + 2\sqrt{2} \exp(U) \Re_{|D_T|}(\mathcal{H}) + 2U \sqrt{\frac{\log \frac{1}{\delta}}{2|D_T|}} \\
&\quad + 2(2U + G) \exp(U) \Re_{|D_S|}(\mathcal{H}) + 2(U + 2G) \Re_{|D_S|}(\mathcal{G}) + 2GU \sqrt{\frac{\log \frac{1}{\delta}}{2|D_S|}} + 2U \sqrt{C_3(\delta)} \tag{B.88}
\end{aligned}$$

$$= 2U \sqrt{J(g^*)} + 2C_4(\delta) + 2U \sqrt{C_3(\delta)}, \tag{B.89}$$

which holds with probability at least  $1 - 5\delta$  by union bound. By replacing  $\delta$  with  $\delta/5$ , we obtain

$$L_{CE}(\hat{h}) - L_{CE}(h^*) \leq U\sqrt{J(g^*)} + C_4(\delta/5) + U\sqrt{C_3(\delta/5)}, \quad (\text{B.90})$$

which holds with probability at least  $1 - \delta$ , as desired.  $\square$

### B.6. Proof of Theorem 7

The following lemma relates the cross-entropy loss with  $L^2$  norm. Note the definition of the KL divergence  $D_{KL}$  and the cross entropy  $H$  is provided in Proposition A.3.

**Lemma B.10** *For any  $h \in \mathcal{H}$ , the following inequality holds.*

$$\frac{1}{K} \mathbb{E}_{\bar{p}_T(\mathbf{x})} \left[ \|\bar{p}_T(\cdot|\mathbf{x}) - h(\mathbf{x})\|_2^2 \right] \leq \mathbb{E}_{\bar{p}_T(\mathbf{x})} [D_{KL}(\bar{p}_T(\cdot|\mathbf{x})||h(\mathbf{x}))] \quad (\text{B.91})$$

$$= \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [\ell_{CE}(h(\mathbf{x}), y)] - \mathbb{E}_{\bar{p}_T(\mathbf{x})} [H(\bar{p}_T(\cdot|\mathbf{x}), \bar{p}_T(\cdot|\mathbf{x}))]. \quad (\text{B.92})$$

**Proof.** By direct calculation, we have

$$\mathbb{E}_{\bar{p}_T(\mathbf{x})} [D_{KL}(\bar{p}_T(\cdot|\mathbf{x})||h(\mathbf{x}))] = \mathbb{E}_{\bar{p}_T(\mathbf{x})} [H(\bar{p}_T(\cdot|\mathbf{x}), h(\mathbf{x})) - H(\bar{p}_T(\cdot|\mathbf{x}), \bar{p}_T(\cdot|\mathbf{x}))] \quad (\text{B.93})$$

$$= \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [\ell_{CE}(h(\mathbf{x}), y)] - \mathbb{E}_{\bar{p}_T(\mathbf{x})} [H(\bar{p}_T(\cdot|\mathbf{x}), \bar{p}_T(\cdot|\mathbf{x}))] \quad (\text{B.94})$$

where we use the fact that

$$H(\bar{p}_T(\cdot|\mathbf{x}), h(\mathbf{x})) = - \sum_{k=1}^K \bar{p}_T(\cdot|\mathbf{x})_k \log h(\mathbf{x})_k = \mathbb{E}_{\bar{p}_T(y|\mathbf{x})} [\ell_{CE}(h(\mathbf{x}), y)]. \quad (\text{B.95})$$

By combining Lemma A.4 and Eq. (B.94), we conclude the proof.  $\square$

We prove Theorem 7 as follows.

**Proof** (Proof of Theorem 7). By Lemma B.5 and Lemma B.10, we have

$$\begin{aligned} & \frac{1}{2K} \mathbb{E}_{p_{te}(\mathbf{x})} \left[ \left\| p_{te}(\cdot|\mathbf{x}) - \hat{h}(\mathbf{x}) \right\|_2^2 \right] - T_{KL}(h^*) \\ & \leq \frac{1}{K} \left( \mathbb{E}_{\bar{p}_T(\mathbf{x})} \left[ \left\| \bar{p}_T(\cdot|\mathbf{x}) - \hat{h}(\mathbf{x}) \right\|_2^2 \right] + Z(\hat{h})\tau_X + \tau_{Y|X}^2 \right) - T_{KL}(h^*) \end{aligned} \quad (\text{B.96})$$

$$\begin{aligned} & \leq \left( \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [\ell_{CE}(\hat{h}(\mathbf{x}), y)] - \mathbb{E}_{\bar{p}_T(\mathbf{x})} [H(\bar{p}_T(\cdot|\mathbf{x}), \bar{p}_T(\cdot|\mathbf{x}))] \right) + \frac{1}{K} \left( Z(\hat{h})\tau_X + \tau_{Y|X}^2 \right) \\ & \quad - \left( \mathbb{E}_{\bar{p}_T(\mathbf{x}, y)} [\ell_{CE}(h^*(\mathbf{x}), y)] - \mathbb{E}_{\bar{p}_T(\mathbf{x})} [H(\bar{p}_T(\cdot|\mathbf{x}), \bar{p}_T(\cdot|\mathbf{x}))] \right) \end{aligned} \quad (\text{B.97})$$

$$= L_{CE}(\hat{h}) - L_{CE}(h^*) + \frac{1}{K} Z(\hat{h})\tau_X + \frac{1}{K} \tau_{Y|X}^2. \quad (\text{B.98})$$

By combining Eq. (B.98) with Lemma B.6, we obtain

$$R_{01}(\hat{h}) - B_{01} \leq 4K\eta_{min}^{-2} \left( T_{KL}(h^*) + U\sqrt{J(g^*)} + C_4(\delta/5) + U\sqrt{C_3(\delta/5)} \right. \\ \left. + \frac{1}{K}Z(\hat{h})W_1(p_{te}, \bar{p}_T) + \frac{1}{K}\tau_{Y|X}^2 \right), \quad (\text{B.99})$$

which concludes the proof.  $\square$

## Supplementary C. Experimental Details

This section reports the full experimental results over seven real-world datasets. In addition, details of the experiments are provided.

### C.1. Environment

The implementation of our method is based on PyTorch<sup>4</sup>, NumPy<sup>5</sup>, and scikit-learn<sup>6</sup>. All experiments are carried out on a computational server equipping four Intel Xeon Platinum 8260 CPUs with 192 logical cores in total and 1TB RAM.

### C.2. Full Experimental Results

We conduct extensive experiments to verify the effectiveness of our method over seven real-world datasets. While we have reported a part of the results in the main part of this paper, we present the full results for completeness.

**Dataset.** We use seven real-world datasets in our extensive experiment. The statistics are presented in Table C.1. All datasets are obtained from USP DS Repository (Souza et al., 2020)<sup>7</sup>. We exclude datasets whose number of samples are less than 15000 or whose number of features are more than 100 for feasibility of the experiments. We also exclude datasets with highly imbalanced data in terms of overall class-balance and class-balance w.r.t. the sample orders; these nature of imbalance often make training data only contain a single labels and/or test data contain novel labels which do not appear in training data. For Forest dataset, we convert the original 7 class classification into binary classification by letting labels be whether the original label is 1 (majority class) or not for better label balances. The inputs and outputs are normalized to have a mean of zero and a variance of one as a preprocess.

---

4. <https://pytorch.org/>

5. <https://numpy.org/>

6. <https://scikit-learn.org/stable/>

7. <https://sites.google.com/view/uspdsrepository>, Accessed: 2025-06-24

Table C.1: Dataset statistics.

Dataset	Samples	Features	Classes
Weather	18159	8	2
Smartmeter	22950	96	10
Powersupply	29928	2	24
Electricity	45312	8	2
Rialto	82250	27	10
Airlines	539383	7	2
Forest	581012	54	2

**Settings.** The experimental settings are generally the same with Section 5. On the other hand, we have tested  $T = 5N$  for each value of  $N \in \{200, 500, 1000, 2000, 5000\}$ . It should be noted that when the ratio  $T/N < 5$  where there are less flexibility for selecting old samples, the sample selection computed by comparison methods becomes almost the same and the differences among the methods cannot be observed.

**Results.** Before discussing on the results, we would like to remark following notes w.r.t. the results and evaluation.

- Variances of the results are generally very high by nature. Real-world datasets are not always drifting and there exist relatively stationary periods. Then, the performance over drifting periods and stationary periods differs significantly. This indicates that no matter how many experiments are conducted (we run 30 independent experiments for stable results), high variances certainly appear. Thus, we evaluate the results based on the statistical tests to examine the significances.
- Results are very competitive. The performance of the methods has saturated and there might less spaces for improvement. Hence, we mainly interpret the results based on average ranks over the experiments to take the versatility of the methods into account.

We present all results in Table C.2 to Table C.5. We conduct experiments under 126 settings in total. Our TSJD achieves 37 best results and 119 best or comparable results, with an average rank of 3.06, which is the best among all baselines. These results confirm the effectiveness and superiority of TSJD. All other baseline methods perform worse than the naive baseline of using  $D$ , highlighting the challenge of selecting effective training samples from drifting data.

Table C.2: Average zero-one loss ( $\downarrow$ ) over 30 random trials. Decoration follows Table 2.

Data	Model	$N$	$T$	Naive $D_T$	Baseline $D$	Time-based		Cov.shift uLSIF	Drift Localization		(Ours)
						PHT	ADWIN		LDD-DSDA	LCD	TSJD
Weather	LGBM	200	1000	30.23 (9.25)	<b>23.07</b> (6.50)	<b>23.13</b> (6.57)	<b>22.47*</b> (5.73)	29.13 (10.07)	24.93 (7.63)	<b>23.10</b> (7.32)	24.47 (7.73)
			2000	29.57 (9.18)	<b>21.97*</b> (7.60)	<b>22.00</b> (7.48)	<b>22.30</b> (7.38)	28.33 (9.50)	24.60 (9.20)	<b>22.53</b> (8.17)	<b>22.90</b> (8.06)
		500	2500	25.80 (6.93)	<b>22.70</b> (6.51)	<b>22.87</b> (6.55)	<b>22.67*</b> (6.40)	25.67 (5.93)	<b>23.63</b> (6.22)	<b>23.20</b> (6.88)	<b>24.33</b> (5.77)
			5000	19.90 (6.36)	<b>17.73</b> (5.02)	<b>17.73</b> (5.02)	<b>17.60*</b> (4.85)	19.63 (6.09)	<b>18.70</b> (5.84)	<b>18.33</b> (5.73)	<b>18.43</b> (5.72)
		1000	5000	<b>19.33</b> (6.16)	<b>18.27</b> (5.94)	<b>18.60</b> (6.12)	<b>18.33</b> (5.55)	<b>18.97</b> (6.09)	<b>18.37</b> (5.19)	<b>18.10*</b> (5.84)	<b>18.27</b> (5.95)
			10000	23.70 (6.75)	<b>20.57</b> (4.58)	<b>21.13</b> (5.53)	<b>21.30</b> (5.74)	22.93 (6.29)	<b>21.03</b> (6.00)	<b>20.43*</b> (5.08)	<b>21.37</b> (5.70)
		2000	10000	<b>21.67</b> (5.19)	<b>20.43</b> (4.38)	<b>20.70</b> (4.40)	<b>20.93</b> (5.25)	<b>21.17</b> (4.77)	<b>21.03</b> (4.78)	<b>20.13*</b> (5.66)	<b>21.10</b> (5.44)
	NN	200	1000	27.97 (9.50)	<b>21.87</b> (6.32)	<b>21.70*</b> (6.35)	<b>21.70*</b> (6.42)	29.17 (11.73)	<b>23.07</b> (6.57)	<b>21.73</b> (6.26)	<b>23.17</b> (7.45)
			2000	27.50 (8.70)	<b>20.63</b> (7.15)	<b>21.13</b> (7.97)	<b>20.47*</b> (7.90)	26.43 (8.52)	<b>22.73</b> (9.19)	<b>20.63</b> (7.68)	<b>22.03</b> (7.86)
		500	2500	23.63 (6.83)	<b>21.23*</b> (5.93)	<b>21.40</b> (5.94)	<b>21.27</b> (5.83)	22.93 (6.16)	<b>21.87</b> (6.14)	<b>21.47</b> (5.67)	<b>21.27</b> (5.19)
			5000	19.90 (6.23)	<b>17.57</b> (5.44)	<b>17.63</b> (5.55)	<b>17.53*</b> (5.51)	<b>17.83</b> (4.75)	<b>18.47</b> (5.35)	<b>17.80</b> (5.09)	<b>18.00</b> (5.75)
		1000	5000	<b>18.27</b> (4.70)	<b>17.27</b> (4.34)	<b>17.17*</b> (4.19)	<b>17.67</b> (4.46)	<b>18.33</b> (5.54)	<b>17.93</b> (5.35)	<b>17.60</b> (4.42)	<b>17.40</b> (5.47)
			10000	<b>20.70</b> (5.44)	<b>18.93*</b> (4.55)	<b>19.10</b> (4.85)	<b>19.37</b> (5.12)	<b>21.03</b> (5.48)	<b>19.93</b> (5.16)	<b>19.60</b> (5.28)	<b>19.23</b> (5.41)
		2000	10000	20.20 (5.50)	<b>19.27</b> (5.32)	<b>19.07</b> (5.23)	<b>18.80*</b> (4.91)	20.23 (5.11)	<b>20.33</b> (5.99)	<b>19.57</b> (5.25)	20.20 (5.29)
Smartmeter	LGBM	200	1000	28.23 (9.78)	<b>19.27</b> (8.19)	<b>19.27</b> (8.19)	<b>19.53</b> (8.44)	28.53 (14.52)	23.83 (10.81)	<b>19.63</b> (8.81)	<b>18.50*</b> (8.60)
			2000	26.67 (8.10)	19.90 (9.40)	20.07 (9.32)	<b>20.10</b> (7.81)	27.40 (10.49)	21.77 (10.14)	20.13 (9.72)	<b>17.67*</b> (8.36)
		500	2500	23.83 (8.91)	<b>19.90</b> (9.83)	20.33 (10.07)	<b>21.07</b> (9.96)	25.23 (11.14)	21.33 (11.02)	<b>20.13</b> (9.27)	<b>18.80*</b> (9.63)
			5000	22.97 (8.40)	<b>13.50</b> (7.23)	17.33 (9.09)	19.90 (10.60)	24.80 (13.80)	16.73 (8.79)	<b>14.33</b> (8.09)	<b>12.83*</b> (6.40)
		1000	5000	21.73 (8.45)	<b>13.87</b> (7.54)	17.17 (9.73)	17.73 (7.86)	21.20 (10.97)	15.73 (7.82)	<b>13.83</b> (7.53)	<b>13.30*</b> (6.89)
			10000	21.13 (9.02)	<b>13.13</b> (7.23)	17.40 (9.46)	19.57 (10.26)	21.90 (11.97)	15.40 (9.14)	<b>12.70</b> (7.68)	<b>12.07*</b> (7.61)
		2000	10000	19.83 (11.38)	<b>12.97*</b> (7.37)	14.93 (7.32)	17.17 (9.87)	19.67 (9.61)	14.87 (7.87)	<b>13.20</b> (7.68)	<b>13.20</b> (7.38)
			20000	22.23 (9.32)	15.20 (5.69)	17.43 (7.09)	17.80 (6.22)	23.23 (9.85)	16.53 (5.42)	15.67 (5.94)	<b>13.37*</b> (5.40)
	NN	200	1000	37.30 (10.22)	34.63 (9.98)	34.63 (9.98)	<b>34.50</b> (9.85)	40.10 (12.46)	36.80 (10.77)	<b>34.30</b> (10.12)	<b>32.10*</b> (10.44)
			2000	36.30 (10.49)	36.80 (10.66)	36.73 (10.65)	<b>33.23</b> (10.27)	40.00 (9.41)	37.77 (9.12)	36.43 (9.96)	<b>31.00*</b> (8.96)
		500	2500	36.97 (12.23)	<b>35.47</b> (12.26)	35.97 (12.48)	<b>35.20</b> (12.17)	38.03 (12.48)	37.33 (12.56)	35.73 (11.91)	<b>32.20*</b> (10.85)
			5000	35.27 (9.98)	33.33 (9.60)	35.27 (9.80)	36.43 (10.66)	37.43 (8.82)	34.27 (8.49)	33.27 (8.55)	<b>29.20*</b> (6.88)
		1000	5000	36.87 (10.44)	<b>33.27</b> (9.09)	34.73 (10.75)	35.37 (8.92)	37.20 (8.97)	35.90 (10.29)	<b>32.90</b> (8.64)	<b>31.17*</b> (8.09)
			10000	36.10 (10.80)	<b>31.53</b> (10.13)	39.13 (8.69)	37.83 (9.29)	37.13 (10.85)	34.83 (11.07)	<b>31.80</b> (9.19)	<b>30.13*</b> (8.61)
		2000	10000	38.53 (9.58)	<b>32.00</b> (10.25)	35.20 (7.92)	38.43 (7.68)	39.37 (10.14)	34.37 (9.20)	<b>32.03</b> (10.60)	<b>31.20*</b> (9.79)
			20000	37.23 (8.39)	<b>32.17</b> (5.45)	33.80 (5.97)	36.10 (7.50)	37.97 (8.72)	34.33 (4.92)	<b>31.60</b> (5.12)	<b>30.07*</b> (4.45)



Table C.3: Average zero-one loss ( $\downarrow$ ). Continued from Table C.2.

Data	Model	$N$	$T$	Naive Baseline		Time-based		Cov.shift	Drift Localization		(Ours)		
				$D_T$	$D$	PHT	ADWIN	uLSIF	LDD-DSDA	LCD	TSJD		
Powersupply	LGBM	200	1000	<b>81.50*</b> (7.20)	<b>81.97</b> (7.62)	<b>81.97</b> (7.62)	<b>82.13</b> (7.63)	<b>83.83</b> (7.33)	<b>83.03</b> (6.58)	<b>82.47</b> (7.28)	<b>81.87</b> (6.52)		
			2000	<b>80.43*</b> (9.52)	85.20 (7.24)	85.33 (7.22)	85.60 (7.08)	84.70 (7.17)	85.10 (7.06)	85.53 (7.55)	<b>83.43</b> (7.27)		
		500	2500	<b>83.37*</b> (8.72)	87.33 (5.80)	87.47 (5.66)	86.57 (6.17)	<b>85.43</b> (7.11)	87.60 (5.15)	88.10 (5.58)	<b>85.27</b> (6.93)		
			5000	<b>79.60*</b> (8.44)	83.60 (6.92)	83.70 (6.34)	84.10 (5.91)	83.77 (7.69)	83.43 (6.49)	84.73 (6.82)	<b>81.80</b> (6.72)		
		1000	5000	<b>81.57*</b> (7.66)	<b>83.63</b> (7.44)	<b>83.90</b> (6.95)	<b>83.33</b> (7.54)	<b>84.10</b> (7.07)	<b>83.37</b> (8.28)	<b>84.33</b> (7.77)	<b>82.37</b> (7.60)		
			10000	<b>82.57*</b> (7.02)	<b>84.20</b> (7.19)	86.17 (7.33)	<b>85.57</b> (6.80)	85.00 (6.41)	<b>84.10</b> (7.37)	<b>85.37</b> (5.52)	<b>83.30</b> (7.41)		
		2000	10000	<b>84.93</b> (7.55)	<b>83.87*</b> (7.14)	<b>85.03</b> (6.67)	<b>85.53</b> (6.76)	<b>84.73</b> (7.27)	<b>84.37</b> (6.37)	<b>84.10</b> (7.13)	<b>84.13</b> (7.74)		
			20000	<b>81.63</b> (5.57)	<b>80.53*</b> (4.22)	82.93 (4.66)	<b>82.00</b> (5.12)	<b>81.87</b> (4.77)	<b>80.57</b> (4.28)	<b>81.53</b> (3.86)	<b>82.07</b> (5.16)		
		5000	25000	<b>82.63</b> (4.87)	<b>82.40</b> (4.26)	<b>81.03*</b> (4.03)	<b>82.00</b> (4.97)	<b>81.97</b> (4.63)	<b>82.47</b> (3.97)	<b>82.23</b> (4.70)	<b>81.50</b> (5.02)		
		NN	200	1000	83.70 (6.11)	<b>81.93</b> (7.65)	<b>81.93</b> (7.65)	<b>81.97</b> (7.63)	<b>82.33</b> (7.22)	84.83 (6.13)	<b>81.93</b> (7.01)	<b>81.10*</b> (6.94)	
				2000	<b>85.30</b> (4.90)	<b>83.27</b> (9.21)	<b>83.27</b> (9.21)	<b>83.60</b> (8.94)	<b>84.70</b> (7.71)	<b>86.20</b> (7.10)	<b>83.90</b> (9.04)	<b>82.23*</b> (8.74)	
			500	2500	<b>83.77*</b> (7.39)	87.37 (6.29)	87.57 (5.74)	87.53 (5.90)	<b>84.17</b> (6.29)	87.00 (5.30)	<b>87.00</b> (5.85)	<b>84.07</b> (6.89)	
	5000			<b>81.33</b> (8.77)	82.57 (8.32)	82.77 (7.18)	82.93 (7.85)	<b>80.63</b> (7.77)	<b>81.87</b> (8.90)	83.90 (8.82)	<b>78.83*</b> (8.36)		
	1000		5000	<b>81.33</b> (7.43)	<b>82.50</b> (9.38)	<b>83.07</b> (8.70)	<b>81.57</b> (9.06)	<b>81.77</b> (8.68)	<b>82.20</b> (8.39)	<b>82.63</b> (9.26)	<b>80.60*</b> (9.00)		
			10000	<b>80.67</b> (9.01)	83.77 (7.47)	85.07 (7.70)	83.80 (8.16)	<b>81.17</b> (8.69)	<b>82.47</b> (7.52)	<b>84.03</b> (6.83)	<b>79.77*</b> (9.95)		
	2000		10000	<b>83.67</b> (7.75)	<b>83.57</b> (7.18)	<b>84.37</b> (7.36)	<b>84.03</b> (7.66)	<b>82.53</b> (8.48)	<b>83.67</b> (7.28)	<b>83.53</b> (7.51)	<b>82.17*</b> (8.87)		
			20000	<b>78.40</b> (7.13)	<b>78.77</b> (4.50)	<b>79.13</b> (6.26)	<b>78.37</b> (6.13)	<b>77.60*</b> (6.64)	<b>78.60</b> (5.07)	<b>79.67</b> (4.58)	<b>78.23</b> (7.00)		
	5000		25000	<b>80.87</b> (5.37)	<b>80.17</b> (5.35)	<b>79.07*</b> (5.95)	80.50 (6.91)	80.80 (4.95)	<b>80.73</b> (5.53)	80.33 (5.79)	<b>80.93</b> (5.45)		
	Electricity		LGBM	200	1000	<b>23.77</b> (15.13)	<b>22.80</b> (12.49)	<b>22.53</b> (12.53)	<b>21.57*</b> (13.20)	<b>21.60</b> (13.63)	<b>24.10</b> (14.11)	<b>22.43</b> (13.07)	<b>23.60</b> (16.37)
					2000	<b>20.27</b> (11.58)	<b>18.57</b> (11.11)	<b>18.53</b> (11.52)	<b>21.23</b> (14.74)	<b>21.03</b> (11.09)	<b>20.00</b> (11.52)	<b>18.07*</b> (11.24)	<b>18.43</b> (10.75)
				500	2500	<b>21.90</b> (13.05)	<b>20.80</b> (12.37)	<b>19.77*</b> (13.11)	<b>20.20</b> (11.94)	<b>22.23</b> (12.79)	<b>20.10</b> (12.79)	<b>21.07</b> (12.75)	<b>21.30</b> (12.43)
		5000			<b>21.07</b> (14.52)	<b>19.80</b> (13.51)	<b>19.97</b> (12.77)	<b>20.80</b> (14.19)	<b>22.17</b> (14.43)	<b>19.97</b> (13.50)	<b>19.53*</b> (13.67)	<b>19.77</b> (14.63)	
		1000		5000	20.47 (13.26)	<b>19.13</b> (13.43)	<b>20.43</b> (13.32)	<b>20.20</b> (13.31)	<b>19.67</b> (12.91)	<b>19.03</b> (13.17)	<b>19.23</b> (13.58)	<b>18.53*</b> (13.69)	
				10000	<b>21.67</b> (12.60)	<b>19.87*</b> (11.30)	<b>21.10</b> (11.78)	<b>21.70</b> (12.39)	<b>21.30</b> (12.29)	<b>21.30</b> (10.59)	<b>21.07</b> (11.11)	<b>20.67</b> (11.51)	
2000		10000		<b>20.50</b> (11.50)	<b>20.63</b> (10.56)	<b>21.00</b> (11.16)	<b>21.23</b> (10.43)	<b>21.23</b> (11.99)	<b>20.20</b> (9.00)	<b>20.03*</b> (11.12)	<b>21.37</b> (10.66)		
		20000		<b>20.50</b> (15.45)	<b>20.83</b> (14.57)	<b>21.00</b> (14.68)	<b>20.47</b> (14.73)	<b>20.13*</b> (14.30)	<b>22.30</b> (14.61)	<b>20.60</b> (14.29)	<b>20.53</b> (14.57)		
5000		25000		<b>19.33*</b> (10.75)	<b>23.33</b> (14.96)	<b>21.80</b> (11.20)	<b>21.93</b> (11.49)	<b>23.03</b> (12.92)	<b>21.20</b> (11.29)	<b>23.07</b> (13.59)	<b>20.93</b> (12.31)		
NN		200		1000	<b>24.67</b> (15.31)	<b>26.77</b> (17.18)	<b>26.33</b> (16.98)	<b>24.40*</b> (14.20)	<b>24.87</b> (13.25)	<b>27.03</b> (15.07)	<b>26.97</b> (17.60)	<b>26.73</b> (15.85)	
				2000	<b>27.33</b> (13.38)	<b>23.60</b> (12.82)	<b>22.43</b> (12.12)	<b>26.73</b> (16.34)	<b>22.67</b> (10.47)	<b>21.73*</b> (14.48)	<b>24.77</b> (13.70)	<b>24.00</b> (11.69)	
		500		2500	<b>24.53</b> (12.89)	<b>24.80</b> (14.80)	<b>24.40</b> (15.46)	<b>22.10*</b> (14.80)	<b>26.67</b> (15.95)	<b>24.47</b> (16.15)	<b>24.53</b> (16.10)	<b>25.53</b> (15.63)	
			5000	<b>24.03</b> (14.41)	<b>23.50</b> (12.67)	<b>22.00*</b> (12.98)	<b>24.10</b> (14.90)	<b>24.03</b> (11.29)	<b>24.03</b> (13.32)	<b>23.00</b> (13.08)	<b>22.83</b> (13.80)		
		1000	5000	<b>23.77</b> (14.47)	24.00 (13.03)	<b>21.20</b> (11.46)	<b>22.23</b> (12.52)	<b>23.50</b> (12.40)	<b>20.27*</b> (11.62)	23.17 (12.16)	<b>23.43</b> (14.16)		
			10000	<b>24.47</b> (12.79)	<b>23.23</b> (11.39)	<b>22.20</b> (11.83)	<b>22.97</b> (13.09)	<b>21.80</b> (12.39)	<b>22.80</b> (12.46)	<b>23.23</b> (11.50)	<b>20.60*</b> (10.03)		
		2000	10000	<b>20.07*</b> (11.41)	<b>22.90</b> (11.87)	<b>21.23</b> (12.33)	<b>21.43</b> (11.13)	<b>22.93</b> (11.46)	<b>22.30</b> (11.67)	<b>22.97</b> (11.65)	<b>22.83</b> (10.76)		
			20000	<b>24.50</b> (16.19)	<b>21.97</b> (12.70)	<b>19.67*</b> (11.58)	<b>21.33</b> (13.41)	<b>23.37</b> (14.36)	<b>20.90</b> (12.28)	<b>22.23</b> (13.10)	<b>20.23</b> (12.04)		
		5000	25000	24.93 (9.67)	<b>23.43</b> (11.98)	<b>25.17</b> (10.57)	<b>21.87*</b> (9.41)	30.40 (15.01)	<b>23.33</b> (10.98)	<b>23.83</b> (11.91)	<b>23.03</b> (12.17)		

Table C.4: Average zero-one loss ( $\downarrow$ ). Continued from Table C.3.

Data	Model	$N$	$T$	Naive $D_T$	Baseline $D$	Time-based PHT	ADWIN	Cov.shift uLSIF	Drift Localization LDD-DSDA	LCD	(Ours) TSJD
Rialto	LGBM	200	1000	<b>29.80*</b> (20.34)	<b>32.13</b> (20.89)	<b>32.33</b> (21.02)	<b>30.83</b> (20.30)	39.53 (19.46)	<b>32.67</b> (20.69)	<b>32.37</b> (21.58)	<b>33.20</b> (21.86)
			2000	38.80 (24.47)	<b>32.20</b> (20.66)	<b>32.13</b> (20.45)	<b>35.00</b> (23.40)	41.83 (19.92)	<b>35.20</b> (23.04)	<b>32.43</b> (20.70)	<b>31.70*</b> (20.20)
		500	2500	<b>36.53</b> (22.25)	<b>33.10</b> (19.88)	<b>33.63</b> (19.40)	<b>36.43</b> (21.23)	<b>40.43</b> (21.92)	37.13 (21.84)	<b>32.73*</b> (20.15)	<b>33.13</b> (19.73)
			5000	40.47 (23.26)	<b>29.03*</b> (22.07)	34.93 (23.04)	38.13 (24.16)	42.83 (23.15)	<b>32.03</b> (24.06)	<b>29.43</b> (22.07)	<b>30.10</b> (22.50)
		1000	5000	40.03 (24.49)	<b>28.70</b> (22.01)	<b>31.53</b> (19.50)	38.60 (24.29)	39.47 (21.64)	<b>30.10</b> (23.76)	<b>27.50*</b> (20.71)	<b>28.83</b> (21.89)
			10000	34.13 (16.84)	<b>16.10*</b> (16.79)	24.87 (17.99)	32.37 (16.49)	30.67 (20.95)	20.67 (21.12)	<b>16.50</b> (16.78)	<b>16.37</b> (16.73)
		2000	10000	31.40 (17.60)	<b>16.20</b> (17.29)	21.80 (15.88)	31.07 (18.60)	30.43 (20.79)	20.17 (21.16)	<b>15.30*</b> (16.36)	<b>16.03</b> (16.72)
			20000	38.13 (22.03)	<b>17.33*</b> (16.12)	28.50 (22.67)	36.53 (24.36)	39.60 (23.10)	23.33 (22.44)	<b>17.70</b> (16.37)	<b>17.50</b> (16.29)
		5000	25000	22.23 (18.53)	<b>13.60</b> (14.59)	20.17 (18.10)	20.83 (18.15)	19.90 (19.56)	<b>14.37</b> (15.79)	<b>13.63</b> (15.00)	<b>13.00*</b> (14.10)
			50000	20.93 (21.88)	<b>12.03</b> (19.42)	19.83 (21.07)	20.83 (21.67)	21.50 (22.35)	<b>12.50</b> (19.63)	<b>12.47</b> (19.32)	<b>11.57*</b> (18.11)
	NN	200	1000	<b>30.37*</b> (25.13)	<b>32.43</b> (24.60)	<b>32.50</b> (24.77)	<b>31.47</b> (24.33)	<b>37.40</b> (27.50)	<b>32.40</b> (23.25)	<b>32.90</b> (26.01)	<b>33.10</b> (26.02)
			2000	<b>36.47</b> (25.59)	<b>36.40</b> (22.84)	<b>36.73</b> (23.52)	<b>36.40</b> (26.17)	45.40 (25.94)	<b>35.33*</b> (24.90)	<b>39.20</b> (23.94)	<b>35.50</b> (22.88)
		500	2500	<b>38.10</b> (26.29)	<b>42.90</b> (24.53)	<b>41.03</b> (24.65)	<b>37.50*</b> (24.98)	46.50 (24.62)	<b>40.80</b> (26.96)	<b>43.40</b> (25.49)	<b>39.93</b> (26.37)
			5000	<b>44.67</b> (27.26)	<b>43.03</b> (28.01)	51.63 (26.77)	<b>44.73</b> (27.25)	51.23 (27.70)	<b>41.77*</b> (29.56)	<b>44.37</b> (27.91)	<b>42.83</b> (27.93)
		1000	5000	48.67 (26.73)	<b>42.50</b> (27.07)	<b>49.43</b> (25.73)	49.90 (27.44)	55.47 (28.86)	<b>42.63</b> (29.44)	<b>42.67</b> (27.11)	<b>42.40*</b> (27.09)
			10000	38.20 (22.22)	<b>30.33</b> (25.49)	39.70 (21.97)	41.43 (23.18)	45.90 (24.69)	<b>31.37</b> (25.25)	<b>32.80</b> (24.55)	<b>30.10*</b> (25.77)
		2000	10000	41.87 (24.20)	<b>34.27</b> (26.06)	41.97 (24.85)	45.90 (23.83)	45.50 (28.43)	<b>34.10</b> (28.05)	<b>33.10</b> (26.78)	<b>32.60*</b> (25.18)
			20000	46.20 (27.70)	<b>35.60*</b> (29.34)	<b>41.73</b> (29.86)	45.67 (30.12)	52.83 (28.50)	<b>39.43</b> (30.71)	<b>37.17</b> (29.09)	<b>36.63</b> (29.27)
		5000	25000	<b>34.20</b> (25.97)	<b>30.17</b> (26.53)	35.17 (27.20)	<b>34.53</b> (26.15)	40.03 (27.13)	<b>29.20*</b> (24.74)	<b>32.10</b> (28.95)	<b>31.03</b> (27.85)
			50000	32.53 (25.90)	<b>27.47</b> (27.52)	<b>32.07</b> (24.27)	<b>31.77</b> (26.30)	37.50 (27.16)	<b>25.60*</b> (25.55)	<b>27.07</b> (26.65)	<b>25.70</b> (25.15)
Airlines	LGBM	200	1000	36.80 (8.12)	<b>33.73</b> (6.41)	<b>33.73</b> (6.41)	<b>33.60</b> (6.40)	<b>36.60</b> (8.43)	<b>35.03</b> (7.55)	<b>33.33*</b> (6.49)	<b>35.13</b> (7.59)
			2000	37.27 (8.27)	<b>33.73</b> (6.43)	<b>33.33*</b> (6.20)	<b>33.63</b> (6.28)	37.33 (6.80)	<b>35.23</b> (6.50)	<b>33.77</b> (6.38)	36.57 (5.84)
		500	2500	<b>35.30</b> (8.68)	<b>32.43</b> (5.44)	<b>32.67</b> (5.40)	<b>31.90*</b> (5.61)	37.43 (10.33)	<b>32.50</b> (6.98)	<b>32.57</b> (5.90)	<b>33.57</b> (5.71)
			5000	35.73 (8.64)	<b>32.57*</b> (9.64)	<b>33.33</b> (9.76)	<b>32.83</b> (8.38)	36.17 (8.76)	<b>32.97</b> (7.19)	<b>33.07</b> (9.73)	36.53 (9.86)
		1000	5000	<b>34.17</b> (7.23)	<b>34.43</b> (9.94)	<b>34.47</b> (10.05)	<b>33.70</b> (7.87)	<b>34.37</b> (7.41)	<b>34.97</b> (8.81)	<b>33.47*</b> (9.39)	<b>34.77</b> (9.89)
			10000	<b>33.50</b> (7.03)	<b>30.53</b> (6.19)	<b>31.20</b> (5.67)	<b>30.77</b> (5.05)	35.23 (7.89)	<b>31.70</b> (7.00)	<b>30.43*</b> (6.10)	32.90 (6.96)
		2000	10000	<b>31.47</b> (5.86)	<b>31.20</b> (6.13)	<b>30.93</b> (5.34)	<b>30.77</b> (5.30)	<b>32.13</b> (5.42)	<b>31.93</b> (7.24)	<b>30.33*</b> (6.14)	<b>31.37</b> (6.22)
			20000	35.23 (7.10)	<b>32.33</b> (5.74)	<b>33.17</b> (5.64)	<b>33.10</b> (6.02)	34.90 (7.42)	<b>31.70*</b> (5.47)	<b>32.30</b> (5.46)	33.90 (6.21)
		5000	25000	<b>32.63</b> (8.19)	<b>30.97</b> (8.09)	<b>31.07</b> (7.88)	<b>31.67</b> (7.89)	<b>32.33</b> (7.58)	<b>31.93</b> (8.27)	<b>30.73</b> (8.27)	<b>30.67*</b> (7.85)
			50000	<b>35.37</b> (7.57)	<b>33.47*</b> (6.40)	<b>34.23</b> (6.70)	<b>35.10</b> (6.94)	35.80 (7.06)	<b>34.10</b> (6.83)	<b>33.67</b> (6.32)	<b>33.87</b> (7.01)

Table C.5: Average zero-one loss ( $\downarrow$ ). Continued from Table C.4.

Data	Model	$N$	$T$	Naive $D_T$	Baseline $D$	Time-based PHT	ADWIN	Cov.shift uLSIF	Drift Localization LDD-DSDA	LCD	(Ours) TSJD		
Airlines	NN	200	1000	<b>39.10</b> (9.21)	<b>37.93*</b> (8.63)	<b>37.93*</b> (8.63)	<b>37.93*</b> (8.61)	42.53 (9.90)	<b>38.00</b> (10.37)	<b>38.23</b> (8.53)	<b>38.97</b> (9.34)		
			2000	<b>38.90</b> (8.61)	<b>38.87</b> (7.08)	<b>38.83</b> (7.12)	<b>38.60*</b> (7.66)	42.90 (8.08)	<b>39.17</b> (7.86)	<b>38.70</b> (7.21)	<b>41.40</b> (6.95)		
		500	2500	<b>38.63</b> (10.32)	<b>37.17*</b> (8.39)	<b>37.23</b> (8.40)	<b>37.30</b> (8.26)	41.20 (8.97)	<b>38.37</b> (10.60)	<b>37.23</b> (8.07)	<b>37.50</b> (8.92)		
			5000	<b>39.30</b> (9.74)	<b>36.03*</b> (7.99)	<b>36.20</b> (8.36)	<b>36.23</b> (8.23)	41.57 (10.14)	<b>38.37</b> (9.17)	<b>36.30</b> (8.65)	<b>37.13</b> (9.03)		
		1000	5000	<b>36.27</b> (8.99)	36.10 (8.45)	<b>36.17</b> (8.34)	<b>36.30</b> (8.51)	<b>37.87</b> (9.35)	<b>37.93</b> (10.02)	<b>35.03*</b> (8.20)	<b>36.30</b> (8.29)		
			10000	<b>37.40</b> (9.00)	<b>33.93</b> (6.93)	<b>34.57</b> (7.43)	<b>34.50</b> (7.73)	38.97 (9.89)	<b>35.47</b> (7.98)	<b>33.70*</b> (7.30)	37.10 (10.14)		
		2000	10000	37.30 (8.66)	<b>33.80</b> (6.95)	<b>34.63</b> (7.39)	<b>34.27</b> (7.38)	36.13 (8.28)	<b>35.87</b> (8.42)	<b>33.93</b> (6.77)	<b>33.67*</b> (6.82)		
			20000	39.30 (7.46)	<b>36.40</b> (7.29)	<b>35.57*</b> (7.06)	<b>36.37</b> (6.69)	39.37 (7.52)	<b>36.43</b> (8.05)	<b>35.63</b> (6.95)	<b>36.10</b> (7.10)		
		5000	25000	<b>36.27</b> (9.07)	<b>34.43*</b> (8.27)	<b>34.43*</b> (7.71)	<b>35.80</b> (9.09)	<b>36.17</b> (9.33)	<b>35.53</b> (8.32)	<b>34.77</b> (8.40)	<b>34.67</b> (8.15)		
			50000	39.13 (7.21)	<b>36.43</b> (8.37)	<b>36.27*</b> (7.83)	<b>37.37</b> (8.19)	40.63 (7.98)	<b>38.53</b> (7.75)	<b>36.80</b> (8.61)	<b>36.50</b> (6.27)		
		Forest	LGBM	200	1000	35.07 (27.43)	<b>4.73</b> (4.38)	<b>3.97*</b> (3.40)	14.83 (20.42)	34.27 (24.90)	11.63 (11.48)	<b>4.53</b> (4.26)	<b>4.67</b> (4.05)
					2000	34.00 (27.24)	<b>3.13</b> (3.95)	<b>4.37</b> (7.99)	12.40 (16.17)	31.83 (25.22)	9.93 (15.15)	<b>2.77*</b> (3.09)	<b>2.93</b> (4.24)
500	2500			9.37 (9.52)	<b>3.70</b> (4.48)	<b>3.80</b> (4.54)	<b>4.00</b> (4.88)	13.27 (17.71)	<b>3.63*</b> (4.81)	<b>3.97</b> (5.13)	<b>3.63*</b> (4.80)		
	5000			<b>14.77</b> (20.18)	<b>4.40</b> (4.39)	<b>4.00*</b> (4.20)	<b>4.33</b> (5.48)	15.47 (19.89)	<b>5.40</b> (5.24)	<b>4.20</b> (4.45)	<b>4.80</b> (6.64)		
1000	5000			<b>4.33</b> (6.56)	<b>4.27</b> (4.26)	<b>4.27</b> (4.26)	<b>4.07*</b> (4.98)	<b>5.13</b> (7.38)	<b>4.97</b> (4.92)	<b>4.30</b> (4.42)	<b>4.83</b> (7.33)		
	10000			<b>3.07</b> (3.92)	<b>3.43</b> (4.56)	<b>3.63</b> (4.67)	<b>3.00</b> (3.56)	<b>3.43</b> (4.58)	<b>3.53</b> (5.01)	3.60 (4.68)	<b>2.77*</b> (3.45)		
2000	10000			<b>2.93*</b> (3.57)	<b>3.33</b> (4.57)	<b>3.40</b> (4.64)	<b>3.07</b> (3.66)	<b>3.50</b> (3.99)	<b>3.43</b> (4.27)	<b>3.47</b> (4.39)	<b>2.97</b> (3.88)		
	20000			<b>5.07</b> (4.58)	6.67 (6.07)	6.50 (5.77)	<b>5.00</b> (4.34)	<b>6.20</b> (7.61)	6.40 (5.20)	7.03 (5.94)	<b>4.43*</b> (3.87)		
5000	25000			<b>5.57</b> (7.19)	7.90 (7.90)	7.40 (7.44)	<b>4.93*</b> (5.56)	<b>6.23</b> (6.64)	7.80 (8.19)	8.07 (7.65)	<b>6.10</b> (6.12)		
	50000			<b>3.93*</b> (4.01)	9.03 (9.18)	7.43 (7.80)	<b>4.27</b> (4.71)	6.03 (6.79)	10.23 (10.13)	9.07 (9.02)	<b>4.70</b> (5.01)		
200	1000			26.40 (20.64)	<b>5.43</b> (5.16)	<b>4.47</b> (4.26)	<b>16.47</b> (22.93)	33.53 (24.08)	10.63 (9.33)	<b>4.23*</b> (3.58)	<b>5.90</b> (5.38)		
	2000			35.90 (24.98)	<b>3.80</b> (4.75)	<b>6.17</b> (9.32)	16.10 (19.68)	32.70 (21.27)	9.40 (12.52)	<b>3.43*</b> (4.35)	<b>4.07</b> (4.49)		
500	2500			13.87 (16.94)	<b>4.33</b> (5.77)	<b>4.17</b> (5.77)	<b>4.50</b> (5.86)	13.20 (15.25)	<b>4.30</b> (4.70)	<b>4.23</b> (5.61)	<b>3.90*</b> (4.17)		
	5000			14.53 (18.94)	<b>4.73</b> (3.72)	<b>5.03</b> (5.03)	<b>3.93*</b> (3.35)	15.40 (18.62)	5.20 (3.91)	<b>4.43</b> (3.86)	<b>3.97</b> (3.12)		
1000	5000			<b>4.07*</b> (3.31)	<b>5.03</b> (4.44)	<b>5.03</b> (4.44)	<b>4.63</b> (4.98)	<b>5.23</b> (5.06)	<b>5.00</b> (4.60)	<b>5.07</b> (4.95)	<b>4.43</b> (4.32)		
	10000			<b>4.00</b> (4.47)	5.07 (5.79)	5.20 (6.21)	<b>4.20</b> (5.40)	<b>4.27</b> (5.75)	<b>4.87</b> (5.66)	<b>4.70</b> (6.60)	<b>3.60*</b> (4.64)		
2000	10000			<b>4.17</b> (5.48)	<b>4.70</b> (5.73)	<b>4.77</b> (5.76)	<b>3.87*</b> (5.10)	<b>4.30</b> (5.71)	<b>4.70</b> (5.82)	<b>4.83</b> (6.10)	<b>4.00</b> (5.20)		
	20000			<b>5.83</b> (5.99)	8.60 (7.76)	8.57 (6.90)	<b>6.27</b> (6.19)	8.17 (8.95)	8.77 (7.93)	8.87 (8.40)	<b>5.63*</b> (5.80)		
5000	25000	<b>4.80*</b> (5.00)	7.90 (7.46)	7.60 (7.55)	<b>5.37</b> (6.06)	6.90 (7.84)	8.23 (8.66)	7.50 (7.28)	<b>6.10</b> (7.11)				
	50000	<b>4.60*</b> (5.69)	8.67 (8.53)	7.50 (7.62)	<b>5.03</b> (5.76)	<b>7.27</b> (8.49)	10.07 (10.27)	7.60 (7.70)	<b>4.87</b> (5.18)				
Average Rank				5.59	3.37	4.27	4.21	6.48	4.84	3.88	<b>3.06*</b>		
#Best				16	17	14	19	2	8	18	<b>37*</b>		
#Best or Comparable				70	106	86	95	53	88	108	<b>119*</b>		

### C.3. Running Time

We evaluate the running time of each method on Rialto dataset with  $T/N = 10$  as a representative setting. Actual running times are shown in Figure C.1.

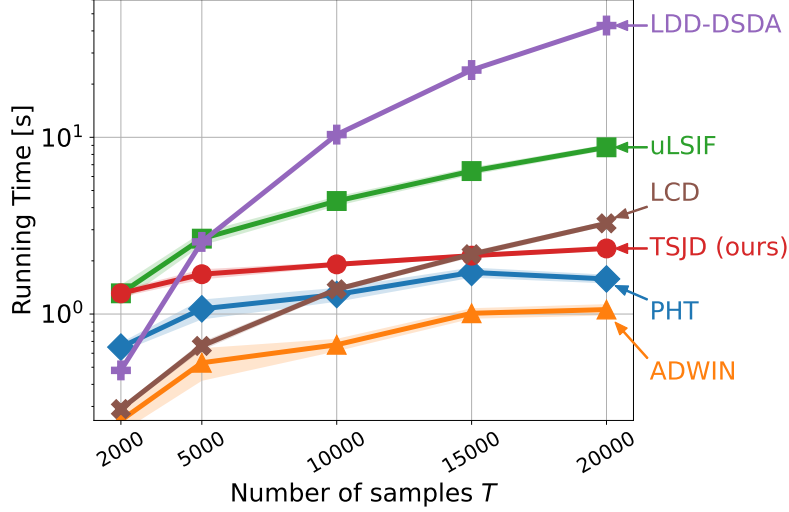


Figure C.1: Average running time [s] ( $\downarrow$ ) over 30 trials. Shades represent the 90% confidence intervals.

PHT, TSJD, and ADWIN have similar running times, less affected by the sample size  $T$ . uLSIF takes longer due to hyperparameter tuning, but speeds up if the hyperparameters are fixed. Theoretically, LDD-DSDA has the highest computational complexity at  $\mathcal{O}(T^3)$ , since it performs a  $\mathcal{O}(T)$ -nearest neighbor search for each sample in  $D$ . Indeed, its running time increases significantly with larger  $T$ . LCD is also influenced by  $T$  and slows down for large  $T$ , but remains faster than LDD-DSDA.

### C.4. Hyperparameters

Table C.6 shows the all hyperparameters of TSJD used in Section 5 and Supplementary C.2. Each of them are tuned in the way described in Section 3.4 using each entire data set.

Table C.6: Hyperparameters

Data	$N$	$T$	$\sigma_x$	$\sigma_y$	$\beta$
Weather	200	1000	0.1	0.2	10.0
	200	2000	0.05	0.05	100.0
	500	2500	0.05	0.5	100.0
	500	5000	0.2	0.0	10.0
	1000	5000	2.0	0.5	10.0
	1000	10000	0.2	0.0	10.0
	2000	10000	2.0	0.5	10.0
Smartmeter	200	1000	0.5	0.01	0.1
	200	2000	1.0	0.01	1.0
	500	2500	2.0	0.0	0.1
	500	5000	1.0	0.01	1.0
	1000	5000	2.0	0.0	0.1
	1000	10000	5.0	0.0	10.0
	2000	10000	10.0	0.0	0.1
	2000	20000	10.0	0.0	1.0
Powersupply	200	1000	0.0002	0.1	100.0
	200	2000	0.0002	0.1	100.0
	500	2500	0.01	0.02	0.1
	500	5000	0.0001	0.5	10.0
	1000	5000	0.1	0.2	0.1
	1000	10000	0.002	0.05	0.1
	2000	10000	0.1	0.2	0.1
	2000	20000	0.01	0.1	1.0
Electricity	5000	25000	0.02	0.2	1.0
	200	1000	0.2	0.0	10.0
	200	2000	0.2	0.0	10.0
	500	2500	0.5	0.1	10.0
	500	5000	0.5	0.0	100.0
	1000	5000	1.0	0.5	10.0
	1000	10000	0.5	0.5	100.0
	2000	10000	2.0	0.5	10.0
Rialto	2000	20000	5.0	0.5	10.0
	5000	25000	5.0	0.5	10.0
	200	1000	0.2	0.5	1.0
	200	2000	1.0	0.1	10.0
	500	2500	1.0	0.5	100.0
	500	5000	0.5	0.02	100.0
	1000	5000	10.0	0.5	10.0
	1000	10000	0.2	0.0	100.0
Airlines	2000	10000	1.0	0.05	10.0
	2000	20000	0.1	0.0	10.0
	5000	25000	0.5	0.01	1.0
	5000	50000	0.2	0.05	1.0
	200	1000	0.05	0.5	100.0
	200	2000	10.0	0.5	1.0
	500	2500	10.0	0.5	10.0
	500	5000	0.2	0.5	1.0
Forest	1000	5000	1.0	0.5	10.0
	1000	10000	5.0	0.0	10.0
	2000	10000	0.5	0.5	10.0
	2000	20000	10.0	0.5	10.0
	5000	25000	1.0	0.5	100.0
	5000	50000	2.0	0.01	100.0
	200	1000	0.05	0.5	0.1
	200	2000	0.05	0.1	10.0
Forest	500	2500	0.1	0.5	0.1
	500	5000	0.2	0.2	10.0
	1000	5000	0.5	0.02	1.0
	1000	10000	0.2	0.5	10.0
	2000	10000	0.2	0.0	0.1
	2000	20000	0.5	0.0	100.0
	5000	25000	1.0	0.01	10.0
	5000	50000	2.0	0.01	100.0

### C.5. Hyperparameter Sensitivity

We analyze the parameter sensitivities of our method on seven real-world datasets. Using the tuned hyperparameters from the previous section, we individually vary  $\sigma_x$ ,  $\sigma_y$ , and  $\beta$  to evaluate their impacts. For all experiments,  $N = 1000$  and  $T = 10000$  are fixed, and both NN and LGBM models are tested. Each hyperparameter setting is repeated 30 times, and we report the mean and 90% confidence interval.

Results are displayed in Figure C.2 to Figure C.8. Sensitivity varies across datasets and models; for example, hyperparameter changes have little effect on the Rialto dataset, whereas the Airlines dataset shows significant fluctuation. Additionally, sensitivities are higher for NN models compared to LGBM models.

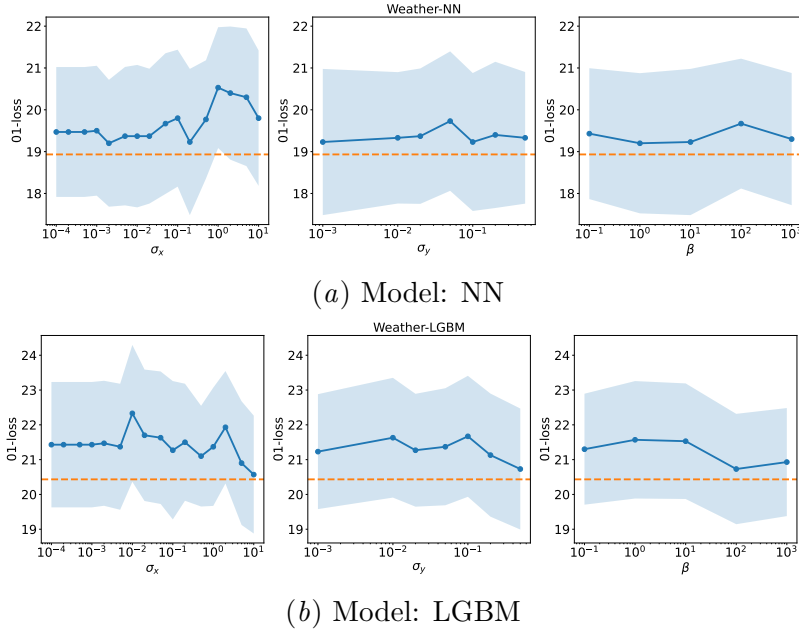
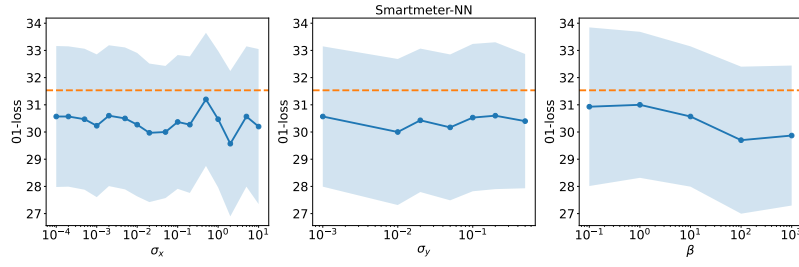
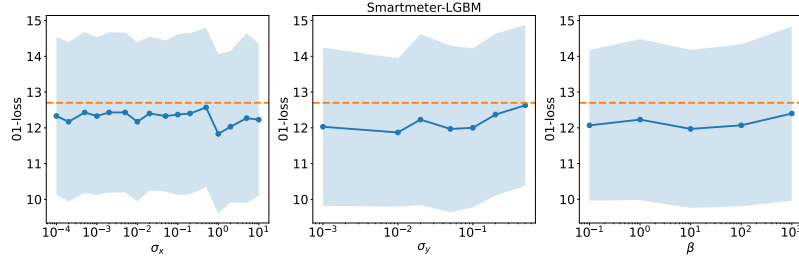


Figure C.2: Parameter sensitivity of our method for Weather dataset. Orange dashed lines indicate the best results among the baselines, reported in Supplementary C.2. Shades represent the 90% confidence intervals over 30 trials.

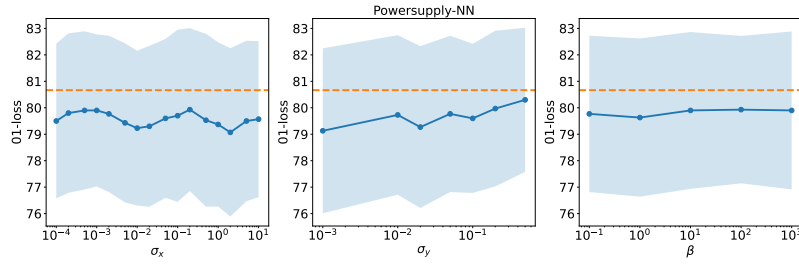


(a) Model: NN

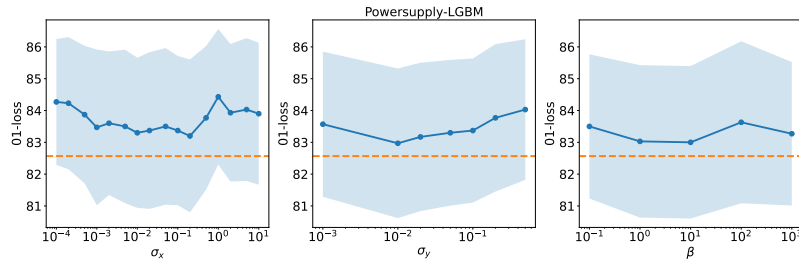


(b) Model: LGBM

Figure C.3: Parameter sensitivity for Smartmeter dataset.



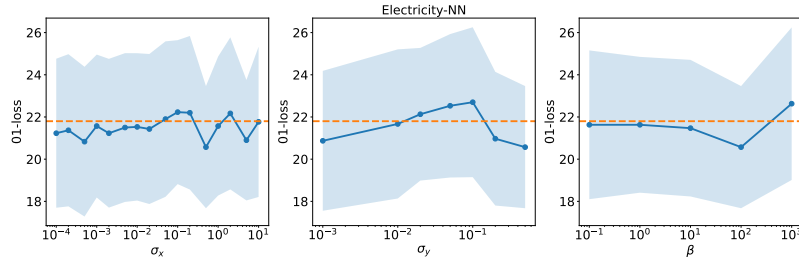
(a) Model: NN



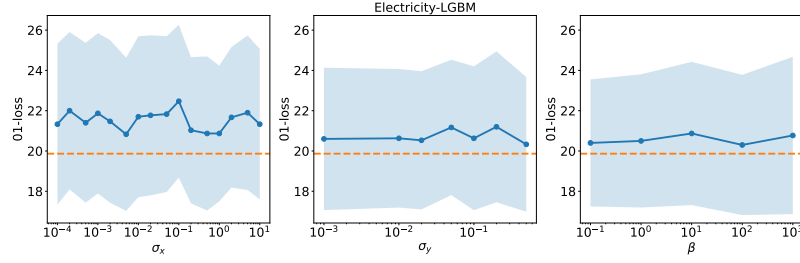
(b) Model: LGBM

Figure C.4: Parameter sensitivity for Powersupply dataset.

# MATSUNO

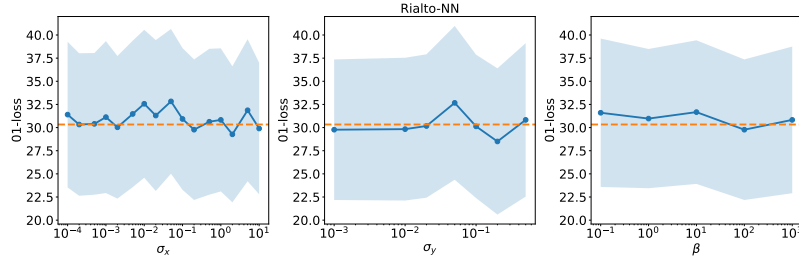


(a) Model: NN

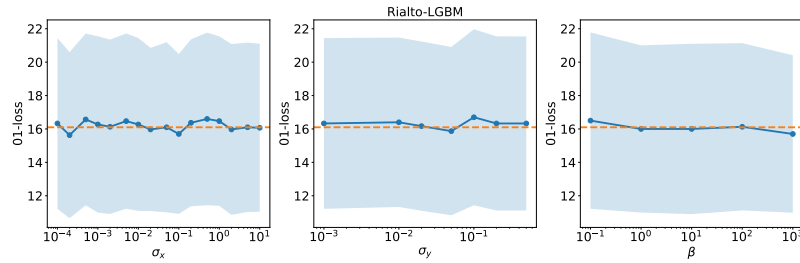


(b) Model: LGBM

Figure C.5: Parameter sensitivity for Electricity dataset.



(a) Model: NN



(b) Model: LGBM

Figure C.6: Parameter sensitivity for Rialto dataset.



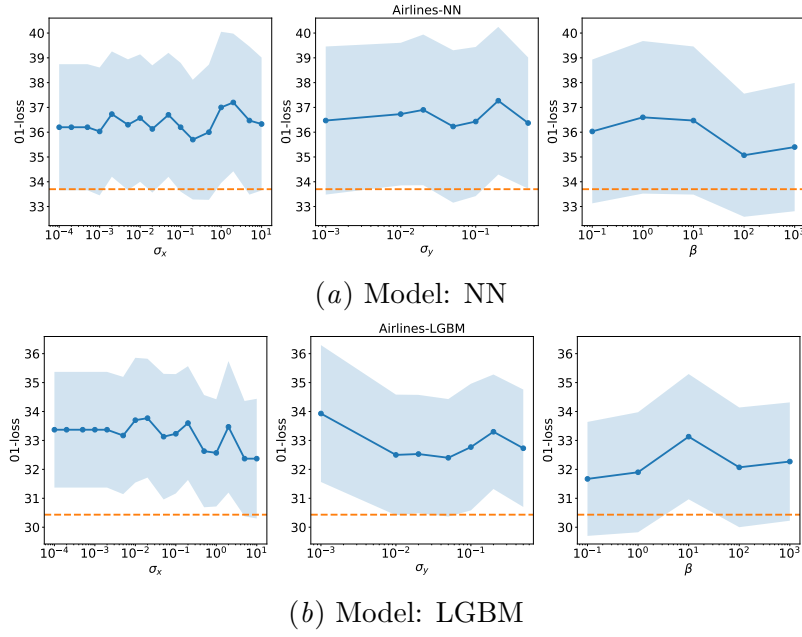


Figure C.7: Parameter sensitivity for Airlines dataset.

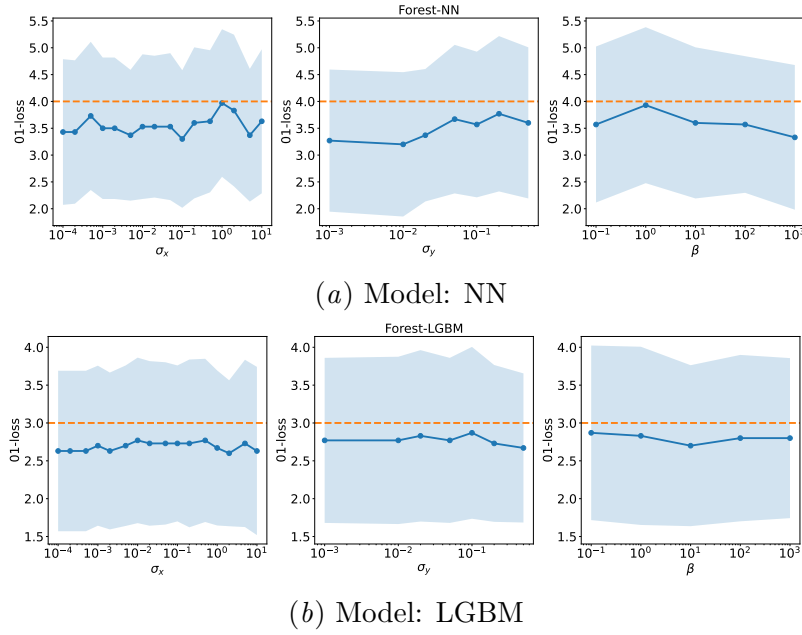


Figure C.8: Parameter sensitivity for Forest dataset.