

ADDRESSING LABEL SHIFT IN DISTRIBUTED LEARNING VIA ENTROPY REGULARIZATION

Zhiyuan Wu*

University of Oslo
zhiyuanw@ifi.uio.no

Changkyu Choi*

UiT The Arctic University of Norway
changkyu.choi@uit.no

Xiangcheng Cao

EPFL
xiangcheng.cao.epfl@gmail.com

Volkan Cevher

LIONS, EPFL
volkan.cevher@epfl.ch

Ali Ramezani-Kebrya

Department of Informatics, University of Oslo
Norwegian Centre for Knowledge-driven Machine Learning (Integreat)
ali@uio.no

ABSTRACT

We address the challenge of minimizing *true risk* in multi-node distributed learning.¹ These systems are frequently exposed to both inter-node and intra-node *label shifts*, which present a critical obstacle to effectively optimizing model performance while ensuring that data remains confined to each node. To tackle this, we propose the Versatile Robust Label Shift (VRLS) method, which enhances the maximum likelihood estimation of the test-to-train label importance ratio. VRLS incorporates Shannon entropy-based regularization and adjusts the importance ratio during training to better handle label shifts at the test time. In multi-node learning environments, VRLS further extends its capabilities by learning and adapting importance ratios across nodes, effectively mitigating label shifts and improving overall model performance. Experiments conducted on MNIST, Fashion MNIST, and CIFAR-10 demonstrate the effectiveness of VRLS, outperforming baselines by up to 20% in imbalanced settings. These results highlight the significant improvements VRLS offers in addressing label shifts. Our theoretical analysis further supports this by establishing high-probability bounds on estimation errors. The code is available at https://github.com/zhiyuan-11/VRLS_main/tree/main.

1 INTRODUCTION

The classical learning theory relies on the assumption that data samples, during training and testing, are *independently and identically distributed (i.i.d.)* drawn from an unknown distribution. However, this *i.i.d.* assumption is often overly idealistic in real-world settings, where the distributions of training and testing samples can differ significantly and change dynamically as the operational environment evolves. In distributed learning (Kim et al., 2022; Wen et al., 2023; Ye et al., 2023; Luo et al., 2023), where nodes retain their own data without sharing, these discrepancies across nodes become more pronounced, further intensifying the learning challenge (Rahman et al., 2023; Wang et al., 2023).

Label shifts (Lipton et al., 2018; Garg et al., 2022; Mani et al., 2022; Zhou et al., 2023) represent a form of distributional discrepancy that arises when the marginal distribution of labels in the training set differs from that in the test set, i.e., $p^{\text{te}}(\mathbf{y}) \neq p^{\text{tr}}(\mathbf{y})$, while the conditional distribution of features given labels, $p(\mathbf{x}|\mathbf{y})$, remains largely stable across both datasets. Label shifts commonly manifest both *inter-node* and *intra-node*, complicating the learning process in real-world distributed learning scenarios. However, a commonly used learning principle in this distributed setting, empirical risk minimization

*These authors contributed equally to this work.

¹We use the term node to refer to a client, FPGA, APU, CPU, GPU, or worker.

(ERM) (Kur et al., 2024), operates under the assumption that the training and test distributions are identical on each node and across nodes. This overlooks these shifts, failing to account for the statistical heterogeneity across decentralized data sources. While the current literature (Yin et al., 2024) addresses statistical heterogeneity across nodes, it often neglects distribution shifts at test or operation time, which has been a significant challenge in the entire data science over decades.

The primary technical challenge in addressing label shifts lies in the efficient and accurate estimation of the test-to-train importance ratios, expressed as $p^{\text{te}}(\mathbf{y})/p^{\text{tr}}(\mathbf{y})$ for all labels. A widely popular solution is Maximum Likelihood Label Shift Estimation (MLLS) (Garg et al., 2020), which frames this estimation as a convex optimization problem, akin to the Expectation-Maximization (EM) algorithm (Saerens et al., 2002). Model calibration refers to the process of ensuring that predicted probabilities reflect the true likelihood of correctness, which is crucial for improving the accuracy of importance ratio estimation (Guo et al., 2017a; Garg et al., 2020). Bias-Corrected Calibration (BCT) (Alexandari et al., 2020) serves as an efficient calibration method that enhances the EM algorithm within MLLS.

While BCT and other post-hoc calibration techniques (Guo et al., 2017c; Kull et al., 2019; Wang et al., 2021; Sun et al., 2024) contribute to improved calibration and may potentially improve model performance, their primary focus remains on refining classification outcomes rather than on accurately approximating the true conditional distribution $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$. The “predictor” in these literature captures the relationship between the input features \mathbf{x} and the corresponding output probabilities across the labels in the discrete label space \mathcal{Y} , with $|\mathcal{Y}| = m$, which should approximate the true distribution of $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$. Despite this goal, training with conventional cross-entropy loss often leads to models that produce predictions that are either highly over-confident or under-confident, resulting in poorly calibrated outputs (Guo et al., 2017a). Consequently, the predictor fails to capture the underlying uncertainty inherent in $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$, which limits its effectiveness in estimating importance ratios (Alexandari et al., 2020; Garg et al., 2020; Guo et al., 2020; 2017b; Pereyra et al., 2017; McMahan et al., 2017).

To address this limitation, we propose a novel Versatile Robust Label Shift (VRLS) method, specifically designed to improve importance ratio estimation for tackling the label shift problem. A key idea of our VRLS method is to approximate $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$ in a way that accounts for the inherent uncertainty over the label space \mathcal{Y} for each input \mathbf{x} . Accordingly, we propose a new objective function incorporating regularization to penalize predictions that lack proper uncertainty calibration. We show that training the predictor in this manner significantly reduces estimation error under various label shift conditions.

Building upon our VRLS method, we extend its application to multi-node settings by proposing an Importance Weighted-ERM (IW-ERM) framework (Fang et al., 2023). Within the multi-node distributed environment, our IW-ERM aims to find an unbiased estimate of the overall true risk minimizer across multiple nodes with varying label distributions. By effectively addressing both intra-node and inter-node label shifts with generalization guarantees, our framework handles the statistical heterogeneity inherent in decentralized data sources. Our extensive experiments demonstrate that the IW-ERM framework, which trains predictors exclusively on local node data, significantly improves overall test error. Moreover, it maintains convergence rates and privacy levels comparable to standard ERM methods while achieving minimal communication and computational overhead compared to existing baselines. Our main contributions are as follows:

- We propose VRLS, which enhances the approximation of the probability distribution $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$ by incorporating a novel regularization term based on Shannon entropy (Neo et al., 2024). This regularization leads to more accurate estimation of the test-to-train label importance ratio, resulting in improved predictive performance under various label shift conditions.
- By integrating our VRLS ratio estimation into multi-node distributed learning environment, we achieve performance close to an upper bound that uses true ratios on Fashion MNIST and CIFAR-10 datasets with 5, 100, and 200 nodes. Our IW-ERM framework effectively manages both inter-node and intra-node label shifts while remaining data confined within each node, resulting in up to 20% improvements in average test error over current baselines.
- We establish high-probability estimation error bounds for VRLS, as well as high-probability convergence bounds for IW-ERM with VRLS in nonconvex optimization settings (Section 5, Appendices E,

Table 1: Details of the label shift scenarios. Their IW-ERM formulas are presented in Appendix C.

Scenario	#Nodes	Assumptions on Distributions	Ratio Node i Needs
No-LS in Equation (17)	2	$p_1^{\text{tr}}(\mathbf{y}) = p_1^{\text{te}}(\mathbf{y}), p_1^{\text{tr}}(\mathbf{y}) \neq p_2^{\text{tr}}(\mathbf{y})$	$p_1^{\text{tr}}(\mathbf{y})/p_2^{\text{tr}}(\mathbf{y})$
LS on single in Equation (18)	2	$p_1^{\text{tr}}(\mathbf{y}) \neq p_1^{\text{te}}(\mathbf{y}), p_2^{\text{tr}}(\mathbf{y}) = p_2^{\text{te}}(\mathbf{y})$	$p_1^{\text{te}}(\mathbf{y})/p_1^{\text{tr}}(\mathbf{y}), p_1^{\text{te}}(\mathbf{y})/p_2^{\text{tr}}(\mathbf{y})$
LS on both in Equation (18)	2	$p_1^{\text{tr}}(\mathbf{y}) \neq p_1^{\text{te}}(\mathbf{y}), p_2^{\text{tr}}(\mathbf{y}) \neq p_2^{\text{te}}(\mathbf{y})$	$p_1^{\text{te}}(\mathbf{y})/p_1^{\text{tr}}(\mathbf{y}), p_1^{\text{te}}(\mathbf{y})/p_2^{\text{tr}}(\mathbf{y})$
LS on multi in Equation (19)	K	$p_k^{\text{tr}}(\mathbf{y}) \neq p_1^{\text{te}}(\mathbf{y})$ for all k	$p_1^{\text{te}}(\mathbf{y})/p_k^{\text{tr}}(\mathbf{y})$ for all k

F). Additionally, we demonstrate that incorporating importance weighting does not negatively impact convergence rates or communication guarantees across various optimization settings.

2 IMPORTANCE RATIO ESTIMATION AND IMPORTANCE WEIGHTED-ERM

Importance ratio estimation Importance ratio estimation for label shifts has been addressed by methods such as solving linear systems (Lipton et al., 2018; Azizzadenesheli et al., 2019) and minimizing distribution divergences (Garg et al., 2020), primarily in the context of a single node. Lipton et al. (2018); Azizzadenesheli et al. (2019); Garg et al. (2020) assumed the conditional distribution $p(\mathbf{x}|\mathbf{y})$ remains fixed between the training and test datasets, while the label distribution $p(\mathbf{y})$ changes. Black Box Shift Estimation (BBSE) (Lipton et al., 2018; Rabanser et al., 2019) and Regularized Learning under Label Shift (RLLS) (Azizzadenesheli et al., 2019) are confusion matrix-based methods for estimating importance ratios in label shift problems. While BBSE has been shown consistent even when the predictor is not calibrated, its subpar performance is attributed to information loss inherent in using confusion matrices (Garg et al., 2020). To overcome this, Garg et al. (2020) has introduced the MLLS, resulting in significant improvements in estimation performance, especially when combined with post-hoc calibration methods like BCT (Shrikumar et al., 2019). This EM algorithm based MLLS method (Saerens et al., 2002; Garg et al., 2020) is concave and can be solved efficiently.

Importance Weighted-ERM Classical ERM seeks to minimize the expected loss over the training distribution using finite samples. However, when there is a distribution shift between the training and test data, the objective of ERM is not to minimize the expected loss over the test distribution, regardless of the number of training samples. To address this, IW-ERM is developed (Shimodaira, 2000; Sugiyama et al., 2006; Byrd & C. Lipton, 2019; Fang et al., 2020), which adjusts the training loss by weighting samples according to the importance ratio, i.e., the ratio of test-to-train density. Shimodaira (2000) has shown that the IW-ERM estimator is asymptotically unbiased under certain conditions. Building on this, Ramezani-Kebrya et al. (2023b) have recently introduced Federated IW-ERM, which incorporates importance ratio estimation to handle covariate shifts in distributed learning. However, this approach has limitations, as it does not address label shifts and the importance ratio estimation method poses potential privacy risks.

In this work, we focus on label shifts and propose an IW-ERM framework enhanced by our VRLS method. We show that our IW-ERM with VRLS performs comparably to an upper bound that utilizes true importance ratios, all while preserving data privacy across distributed data sources. This approach effectively addresses both intra-node and inter-node label shifts while ensuring convergence in probability to the overall true risk minimizer.

3 VERSATILE ROBUST LABEL SHIFT: REGULARIZED RATIO ESTIMATION

In this section, we introduce the Versatile Robust Label Shift (VRLS) method for importance ratio estimation in a single-node setting, which forms the basis of the IW-ERM framework. To solve the optimization problem of IW-ERM, each node k requires an accurate estimate of the ratio:

$$r_k(\mathbf{y}) = \frac{\sum_{j=1}^K p_j^{\text{te}}(\mathbf{y})}{p_k^{\text{tr}}(\mathbf{y})}, \quad (1)$$

where $p_j^{\text{te}}(\mathbf{y})$ and $p_k^{\text{tr}}(\mathbf{y})$ represent the test and training label densities, respectively. To improve clarity and avoid over-complicating notations, we first consider the scenario where we have only one

Algorithm 1 VRLS Importance Ratio Estimation Algorithm**Require:** Labeled training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{n^{\text{tr}}}$.**Require:** Unlabeled test data $\{\mathbf{x}_j\}_{j=1}^{n^{\text{te}}}$.**Require:** Initial predictor f_{θ} .**Ensure:** Optimized predictor f_{θ^*} and estimated importance ratio \mathbf{r}_{f^*} .

- 1: **Training:**
- 2: Optimize f_{θ} using Equation (4) via SGD.
- 3: Continue until the training loss drops below a threshold or the maximum epochs are reached.
- 4: Obtain the optimized predictor f_{θ^*} .
- 5: **Importance Ratio Estimation:**
- 6: With the optimized predictor f_{θ^*} , estimate the importance ratio \mathbf{r}_{f^*} using equation Equation (3).

node under label shifts and then extend to multiple nodes. Table 1 presents various scenarios. In a single-node label shift scenario, the goal is to estimate the ratio $r(\mathbf{y}) = p^{\text{te}}(\mathbf{y})/p^{\text{tr}}(\mathbf{y})$. Following the seminal work of Garg et al. (2020), we formulate importance ratio estimation as a Maximum Likelihood Estimation (MLE) problem by constructing an optimization problem based on Kullback-Leibler (KL) divergence to directly estimate $r(\mathbf{y})$. We train a predictor f_{θ} to approximate $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$, where θ denotes the parameters of a neural network. After training, we apply the predictor f_{θ^*} to a finite set of unlabeled samples drawn from the test distribution to obtain predicted label probabilities. These predictions are then used to estimate the ratio \mathbf{r}_{f^*} . Further details are provided in Algorithm 1.

One of the novelties of VRLS is its ability to better calibrate the predictor, enabling it to better approximate the true conditional distribution $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$. This approximation faces two main challenges, as highlighted in Theorem 3 of (Garg et al., 2020): finite-sample error and miscalibration error. Entropy-based regularization can directly tackle miscalibration, which occurs when predicted probabilities systematically deviate from true likelihoods. Building on these insights, we introduce an *explicit* entropy regularizer into the training objective, which is based on Shannon’s entropy (Pereyra et al., 2017; Neo et al., 2024). The regularization term $\Omega(f_{\theta})$ is defined as:

$$\Omega(f_{\theta}) = \sum_{c=1}^m \phi(f_{\theta}(\mathbf{x}))_c \log \left(\phi(f_{\theta}(\mathbf{x}))_c \right), \quad (2)$$

where ϕ denotes the softmax function, and c represents the c^{th} element of the softmax output vector. With this regularization to the softmax outputs, VRLS encourages smoother and more reliable predictions that account for inherent uncertainty in the data, leading to more accurate importance ratio estimates and improving the SotA in practice. These improvements are empirically demonstrated in Section 6. Our proposed VRLS objective is formulated as follows:

$$\mathbf{r}_{f^*} = \arg \max_{\mathbf{r} \in \mathbb{R}_+^m} \mathbb{E}_{\text{te}} \left[\log(f_{\theta^*}(\mathbf{x})^{\top} \mathbf{r}) \right], \quad (3)$$

where

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{\text{tr}} \left[\ell_{CE}(f_{\theta}(\mathbf{x}), \mathbf{y}) + \zeta \Omega(f_{\theta}) \right]. \quad (4)$$

The vector \mathbf{r} in Equation (3), representing the importance ratios for all m classes, belongs to the non-negative real space \mathbb{R}_+^m . This constraint set is defined similarly to MLLS (Garg et al., 2022), and we use the expected value \mathbb{E}_{te} for estimation, denoting the optimal importance ratio as \mathbf{r}_{f^*} . To train the predictor θ , we minimize the cross-entropy loss ℓ_{CE} together with a scaled regularization term $\zeta \Omega(f_{\theta})$, where $\zeta > 0$ is a coefficient controlling the regularization strength. Incorporating the regularizer $\Omega(f_{\theta})$ improves the model calibration under the influence of ℓ_{CE} loss.

4 VRLS FOR MULTI-NODE ENVIRONMENT

We now extend VRLS to the multi-node environment, taking into account the privacy and communication requirements. This extension naturally aligns with the concept of IW-ERM, effectively integrating these considerations into the multi-node learning paradigm. We consider multiple nodes where each node has distinct training and test distributions. The goal here is to train a global model

Algorithm 2 IW-ERM with VRLS in Distributed Learning

Require: Labeled training data $\{(\mathbf{x}_{k,i}^{\text{tr}}, \mathbf{y}_{k,i}^{\text{tr}})\}_{i=1}^{n_k^{\text{tr}}}$ at each node k , for $k = [K]$.
Require: Unlabeled test data $\{\mathbf{x}_{k,j}^{\text{te}}\}_{j=1}^{n_k^{\text{te}}}$ at each node k , for $k = [K]$.
Require: Initial global model h_w .
Ensure: Trained global model h_w optimized with IW-ERM.

- 1: **Phase 1: Importance Ratio Estimation with VRLS**
- 2: **for each node** $k = 1$ to K **in parallel do**
- 3: Train local predictor $f_{k,\theta}$ on local training data $\{(\mathbf{x}_{k,i}^{\text{tr}}, \mathbf{y}_{k,i}^{\text{tr}})\}$.
- 4: Use f_{k,θ^*} to estimate the importance ratio r_{k,f^*} on unlabeled test data $\{\mathbf{x}_k^{\text{te}}\}$ at node k .
- 5: **end for**
- 6: **Phase 2: Importance Ratio Aggregation**
- 7: **for each node** $k = 1$ to K **do**
- 8: Aggregate importance ratio using Equation (1).
- 9: **end for**
- 10: **Phase 3: Global Model Training with IW-ERM**
- 11: Train global model h_w using Equation (IW-ERM) with the aggregated importance ratios.

that utilizes local data and addresses overall test error. In this setup, each node uses its local data to estimate the required importance ratios, as outlined in Section 3, and shares only low-dimensional ratio information, without the need to share any local data.

The process begins with each node training a global model on its local data, independently estimating its importance ratios. These locally computed ratios are then shared amongst the nodes, allowing for the aggregated ratio required for IW-ERM to be computed centrally. This aggregated ratio is then used to further refine the global model in a second round of global training. This approach ensures minimal communication overhead and preserves node data privacy, as detailed in Section 5. Our experimental results in Section 6 demonstrate that the IW-ERM framework significantly improves test error performance while minimizing communication and computation overhead compared to baseline ERM. The importance ratio estimation and IW-ERM are described in Algorithm 2.

To provide a more comprehensive understanding of the multi-node environment, the following discussion delves into its details. Let $\mathcal{X} \subseteq \mathbb{R}^{d_0}$ be a compact metric space for input features, \mathcal{Y} be a discrete label space with $|\mathcal{Y}| = m$, and K be the number of nodes in a multi-node setting.² Let $\mathcal{S}_k = \{(\mathbf{x}_{k,i}^{\text{tr}}, \mathbf{y}_{k,i}^{\text{tr}})\}_{i=1}^{n_k^{\text{tr}}}$ denote the training set of node k with n_k^{tr} samples drawn i.i.d. from a probability distribution p_k^{tr} on $\mathcal{X} \times \mathcal{Y}$. The test data of node k is drawn from another probability distribution p_k^{te} on $\mathcal{X} \times \mathcal{Y}$. We assume that the class-conditional distribution $p_k^{\text{te}}(\mathbf{x}|\mathbf{y}) = p_k^{\text{tr}}(\mathbf{x}|\mathbf{y}) := p(\mathbf{x}|\mathbf{y})$ remains the same for all nodes k . This is a common assumption and holds when label shifts primarily affect labels' prior distribution of the labels $p(\mathbf{y})$ rather than the underlying feature distribution given the labels, e.g., when features that are generated given a label remains constant (Zadrozny, 2004; Huang et al., 2006; Sugiyama et al., 2007). Note that $p_k^{\text{tr}}(\mathbf{y})$ and $p_k^{\text{te}}(\mathbf{y})$ can be arbitrarily different, which gives rise to intra- and inter-node *label shifts* (Zadrozny, 2004; Huang et al., 2006; Sugiyama et al., 2007; Garg et al., 2023).

In this multi-node environment, the aim is to find an unbiased estimate of the overall *true risk* minimizer across multiple nodes under both intra-node and inter-node *label shifts*. Specifically, we aim to find a hypothesis $h_w \in \mathcal{H} : \mathcal{X} \rightarrow \mathcal{Y}$, represented by a neural network parameterized by w , such that $h_w(\mathbf{x})$ provides a good approximation of the label $\mathbf{y} \in \mathcal{Y}$ corresponding to a new sample $\mathbf{x} \in \mathcal{X}$ drawn from the aggregated *test* data. Let $\ell : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ denote a loss function. Node k aims to learn a hypothesis h_w that minimizes its true (expected) risk:

$$R_k(h_w) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_k^{\text{te}}(\mathbf{x}, \mathbf{y})} [\ell(h_w(\mathbf{x}), \mathbf{y})]. \quad (\text{Local Risk})$$

²Sets and scalars are represented by calligraphic and standard fonts, respectively. We use $[m]$ to denote $\{1, \dots, m\}$ for an integer m . We use \lesssim to ignore terms up to constants and logarithmic factors. We use $\mathbb{E}[\cdot]$ to denote the expectation and $\|\cdot\|$ to represent the Euclidean norm of a vector. We use lower-case bold font to denote vectors.

We now modify the classical ERM and formulate IW-ERM to find a predictor that minimizes the overall true risk over all nodes under label shifts:

$$\min_{h_w \in \mathcal{H}} \sum_{k=1}^K \frac{1}{n_k^{\text{tr}}} \sum_{i=1}^{n_k^{\text{tr}}} \frac{\sum_{j=1}^K p_j^{\text{te}}(\mathbf{y}_{k,i}^{\text{tr}})}{p_k^{\text{tr}}(\mathbf{y}_{k,i}^{\text{tr}})} \ell(h_w(\mathbf{x}_{k,i}^{\text{tr}}), \mathbf{y}_{k,i}^{\text{tr}}), \quad (\text{IW-ERM})$$

where n_k^{tr} is the number of training samples at node k .

To incorporate our VRLS importance ratio estimation method into the IW-ERM framework, we replace the ratio term $\frac{\sum_{j=1}^K p_j^{\text{te}}(\mathbf{y}_{k,i}^{\text{tr}})}{p_k^{\text{tr}}(\mathbf{y}_{k,i}^{\text{tr}})}$ with our estimated importance ratios. This modification aims to align the empirical risk minimization with the true risk minimization over all nodes. We formalize the convergence of this approach in Proposition 4.1.

Proposition 4.1. *Under the label shift setting described in Section 1, equation IW-ERM is consistent and the learned function h_w converges in probability towards the optimal function that minimizes the overall true risk across nodes, $\sum_{k=1}^K R_k$.*

Proof. Due to space limitations, the proof is provided in Appendix C. Convergence in probability is established by applying the law of large numbers following (Shimodaira, 2000)[Section 3] and (Sugiyama et al., 2007)[Section 2.2]. \square

5 RATIO ESTIMATION BOUNDS AND CONVERGENCE RATES

In this section, we present bounds on ratio estimation and convergence rates for the finite sample errors incurred during the estimation, as further discussed in Appendices E, F. In practice, we only have access to a finite number of labeled training samples, $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{n^{\text{tr}}}$, and a finite number of unlabeled test samples, $\{\mathbf{x}_j\}_{j=1}^{n^{\text{te}}}$. These samples serve to compute the following estimates:

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{n^{\text{tr}}} &= \arg \min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n^{\text{tr}}} \sum_{i=1}^{n^{\text{tr}}} \left(\ell_{CE}(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i) + \zeta \Omega(f_{\boldsymbol{\theta}}) \right), \\ \text{and } \hat{\mathbf{r}}_{n^{\text{te}}} &= \arg \max_{\mathbf{r} \in \mathbb{R}_+^m} \frac{1}{n^{\text{te}}} \sum_{j=1}^{n^{\text{te}}} \log(f_{\hat{\boldsymbol{\theta}}_{n^{\text{tr}}}}(\mathbf{x}_j)^{\top} \mathbf{r}). \end{aligned}$$

We will show that the errors of these estimates can be controlled. The following assumptions are necessary to establish our results.

Assumption 5.1 (Boundedness). *The data and the parameter space Θ are bounded, i.e., there exists $b_{\mathcal{X}}, b_{\Theta} > 0$ such that*

$$\forall \mathbf{x} \in \mathcal{X}, \|\mathbf{x}\|_2 \leq b_{\mathcal{X}} \quad \text{and} \quad \forall \boldsymbol{\theta} \in \Theta, \|\boldsymbol{\theta}\|_2 \leq b_{\Theta}.$$

Assumption 5.2 (Calibration). *Let $\boldsymbol{\theta}^*$ be as defined in Equation (4). There exists $\mu > 0$ such that*

$$\mathbb{E} [f_{\boldsymbol{\theta}^*}(\mathbf{x}) f_{\boldsymbol{\theta}^*}(\mathbf{x})^{\top}] \succeq \mu \mathbf{I}_m.$$

The calibration Assumption 5.2 first appears in (Garg et al., 2020). It is necessary for the ratio estimation procedure to be consistent and we refer the reader to Section 4.3 of Garg et al. (2020) for more details. We further need Assumption 5.1 because, unlike (Garg et al., 2020), the empirical estimator $\hat{\mathbf{r}}_{n^{\text{te}}}$ is estimated using another estimator $\hat{\boldsymbol{\theta}}_{n^{\text{tr}}}$. Uniform bounds are therefore needed to control finite sample error as we cannot directly apply concentration inequalities, as is done in the proof of (Garg et al., 2020, Lemma 3), since we do not have independence of the terms appearing in the empirical sums. We nonetheless prove a similar result in the following theorem.

Theorem 5.1 (Ratio Estimation Error Bound). *Let $\delta \in (0, 1)$ and $\mathcal{F} := \{\mathbf{x} \mapsto \mathbf{r}^{\top} f_{\boldsymbol{\theta}}(\mathbf{x}), (\mathbf{r}, \boldsymbol{\theta}) \in \mathcal{R} \times \Theta\}$. Under Assumptions 5.1-5.2, there exist constants $L > 0, B > 0$ such that with probability at least $1 - \delta$:*

$$\|\hat{\mathbf{r}}_{n^{\text{te}}} - \mathbf{r}_{f^*}\|_2 \leq \frac{2}{\mu p_{\min}} \left(\frac{4}{\sqrt{n^{\text{te}}}} \text{Rad}(\mathcal{F}) + 4B \sqrt{\frac{\log(4/\delta)}{n^{\text{te}}}} \right) + \frac{4L}{\mu p_{\min}} \mathbb{E} [\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2]. \quad (5)$$

Here, $p_{\min} = \min_y p(y)$ and

$$\text{Rad}(\mathcal{F}) = \frac{1}{\sqrt{n^{tr}}} \mathbb{E}_{\sigma_1, \dots, \sigma} \left[\sup_{(\mathbf{r}, \boldsymbol{\theta}) \in \mathcal{R} \times \Theta} \left| \sum_{i=1}^{n^{tr}} \sigma_i \mathbf{r}^\top f_{\boldsymbol{\theta}}(\mathbf{x}_i) \right| \right], \quad (6)$$

where σ_1, \dots, σ are Rademacher variables uniformly chosen from $\{-1, 1\}$.

Proof. The proof of Theorem 5.1 is provided in Appendix E. The Rademacher complexity appearing in the bound will depend on the function class chosen for f . Moreover as regularization often encourages lower complexity functions, this complexity can be reduced because of the presence of the regularization term in the estimation of $\boldsymbol{\theta}$ in our setting. \square

By estimating the ratios locally and incorporating them into local losses, the properties of the modified loss with respect to neural network parameters \mathbf{w} remain unchanged, with data-dependent parameters like Lipschitz constants scaled linearly by r_{\max} . Our approach trains the predictor using only local data, ensuring IW-ERM with VRLS retains the same privacy guarantees as baseline ERM-solvers. Communication involves only the marginal label distribution, adding negligible overhead, as it is far smaller than model parameters and requires just one round of communication. Overall, importance weighting does not impact communication guarantees during optimization.

Theorem 5.2 (Convergence-communication). *Let $\max_{\mathbf{y} \in \mathcal{Y}} \sup_f r_f(\mathbf{y}) = r_{\max}$. Suppose Algorithm 2, e.g., IW-ERM with VRLS for multi-node environment, is run for T iterations. Then Algorithm 2 achieves a convergence rate of $\mathcal{O}(r_{\max} h(T))$, where $\mathcal{O}(h(T))$ denotes the rate of ERM-solver baseline without importance weighting. Throughout the course of optimization, Algorithm 2 has the same overall communication guarantees as the baseline.*

In the following, we establish tight convergence rates and communication guarantees for IW-ERM with VRLS in a broad range of importance optimization settings including upper and lower bounds for convex optimization (Theorems 5.3- 5.4), second-order differentiability, composite optimization with proximal operator (Theorem F.3), optimization with adaptive step-sizes, and nonconvex optimization (Theorems F.1- F.2), along the lines of e.g., (Woodworth et al., 2020; Haddadpour et al., 2021; Glasgow et al., 2022; Liu et al., 2023; Hu & Huang, 2023; Wu et al., 2023; Liu et al., 2023).

Assumption 5.3 (Convex and Smooth). *1) A minimizer \mathbf{w}^* exists with bounded $\|\mathbf{w}^*\|_2$; 2) The $\ell \circ h_{\mathbf{w}}$ is β -smoothness and convex w.r.t. \mathbf{w} ; 3) The stochastic gradient $\mathbf{g}(\mathbf{w}) = \tilde{\nabla}_{\mathbf{w}} \ell(h_{\mathbf{w}})$ is unbiased, i.e., $\mathbb{E}[\mathbf{g}(\mathbf{w})] = \nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})$ for any $\mathbf{w} \in \mathcal{W}$ with bounded variance $\mathbb{E}[\|\mathbf{g}(\mathbf{w}) - \nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})\|_2^2]$.*

For convex and smooth optimization, we establish convergence rates for IW-ERM with VRLS and local updating along the lines of e.g., (Woodworth et al., 2020, Theorem 2).

Theorem 5.3 (Upper Bound for Convex and Smooth). *Let $D = \|\mathbf{w}_0 - \mathbf{w}^*\|$, τ denote the number of local steps (number of stochastic gradients per round of communication per node), R denote the number of communication rounds, and $\max_{\mathbf{y} \in \mathcal{Y}} \sup_f r_f(\mathbf{y}) = r_{\max}$. Under Assumption 5.3, suppose Algorithm 2 with τ local updates is run for $T = \tau R$ total stochastic gradients per node with an optimally tuned and constant step-size. Then we have the following upper bound:*

$$\mathbb{E}[\ell(h_{\mathbf{w}_T}) - \ell(h_{\mathbf{w}^*})] \lesssim \frac{r_{\max} \beta D^2}{\tau R} + \frac{(r_{\max} \beta D^4)^{1/3}}{(\sqrt{\tau R})^{2/3}} + \frac{D}{\sqrt{K \tau R}}. \quad (7)$$

Assumption 5.4 (Convex and Second-order Differentiable). *1) The $\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})$ is β -smoothness and convex w.r.t. \mathbf{w} for any (\mathbf{x}, \mathbf{y}) ; 2) The stochastic gradient $\mathbf{g}(\mathbf{w}) = \tilde{\nabla}_{\mathbf{w}} \ell(h_{\mathbf{w}})$ is unbiased, i.e., $\mathbb{E}[\mathbf{g}(\mathbf{w})] = \nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})$ for any $\mathbf{w} \in \mathcal{W}$ with bounded variance $\mathbb{E}[\|\mathbf{g}(\mathbf{w}) - \nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})\|_2^2]$.*

Theorem 5.4 (Lower Bound for Convex and Second-order Differentiable). *Let $D = \|\mathbf{w}_0 - \mathbf{w}^*\|$, τ denote the number of local steps, R denote the number of communication rounds, and $\max_{\mathbf{y} \in \mathcal{Y}} \sup_f r_f(\mathbf{y}) = r_{\max}$. Under Assumption 5.4, suppose Algorithm 2 with τ local updates is run for $T = \tau R$ total stochastic gradients per node with a tuned and constant step-size. Then we have the following lower bound:*

$$\mathbb{E}[\ell(h_{\mathbf{w}_T}) - \ell(h_{\mathbf{w}^*})] \gtrsim \frac{r_{\max} \beta D^2}{\tau R} + \frac{(r_{\max} \beta D^4)^{1/3}}{(\sqrt{\tau R})^{2/3}} + \frac{D}{\sqrt{K \tau R}}. \quad (8)$$

We finally establish high-probability convergence bounds for IW-ERM with VRLS along the lines of e.g., (Liu et al., 2023, Theorem 4.1). To show the impact of importance weighting on convergence rate decoupled from the impact of number of nodes and obtain the current SotA *high-probability* bounds for nonconvex optimization, we focus on IW-ERM with $K = 1$.

Assumption 5.5 (Sub-Gaussian Noise). 1) A minimizer \mathbf{w}^* exists; 2) The stochastic gradients $\mathbf{g}(\mathbf{w}) = \tilde{\nabla}_{\mathbf{w}}\ell(h_{\mathbf{w}})$ is unbiased, i.e., $\mathbb{E}[\mathbf{g}(\mathbf{w})] = \nabla_{\mathbf{w}}\ell(h_{\mathbf{w}})$ for any $\mathbf{w} \in \mathcal{W}$; 3) The noise $\|\mathbf{g}(\mathbf{w}) - \nabla_{\mathbf{w}}\ell(h_{\mathbf{w}})\|_2$ is σ -sub-Gaussian (Vershynin, 2018).

Theorem 5.5 (High-probability Bound for Nonconvex Optimization). Let $\delta \in (0, 1)$ and $T \in \mathbb{Z}_+$. Let $K = 1$ and $\max_{\mathbf{y} \in \mathcal{Y}} \sup_f r_f(\mathbf{y}) = r_{\max}$. Under Assumption 5.5 and β -smoothness of nonconvex $\ell \circ h_{\mathbf{w}}$, suppose IW-ERM is run for T iterations with a step-size $\min\left\{\frac{1}{r_{\max}\beta}, \sqrt{\frac{1}{\sigma^2 r_{\max}\beta T}}\right\}$. Then with probability $1 - \delta$, gradient norm squareds satisfy:

$$\frac{1}{T} \sum_{t=1}^T \|\nabla_{\mathbf{w}}\ell(h_{\mathbf{w}_t})\|_2^2 = O\left(\sigma \sqrt{\frac{r_{\max}\beta}{T}} + \frac{\sigma^2 \log(1/\delta)}{T}\right). \quad (9)$$

Proof. We note that importance ratios do not depend on the model parameters \mathbf{w} and the Lipschitz and smoothness constants for $\ell \circ h_{\mathbf{w}}$ w.r.t. \mathbf{w} are scaled by r_{\max} . The rest of the proof follows the arguments of (Liu et al., 2023, Theorem 4.1). \square

Theorem 5.5 shows that when the stochastic gradients are too noisy $\sigma = \Omega(\sqrt{r_{\max}\beta}/\log(1/\delta))$ such that the second term in the rate dominates, then importance weighting does not have any negative impact on the convergence rate.

6 EXPERIMENTS

The experiments are divided into two main parts: evaluating VRLS’s performance on a single node focusing on intra-node label shifts, and extending it to multi-node distributed learning scenarios with 5, 100, and 200 nodes. In the multi-node cases, we account for both inter-node and intra-node label shifts. Further experimental details, results, and discussions are provided in Appendix J.

Importance ratio estimation. We begin by evaluating VRLS on the MNIST (LeCun et al., 1998) and CIFAR-10 (Krizhevsky) datasets in a single-node setting. Following the common experimental setup in the literature (Lipton et al., 2018), we simulate the test dataset using a Dirichlet distribution with varying α parameters. In this context, a higher α value indicates smoother transitions in the label distribution, while lower values reflect more abrupt shifts. The training dataset is uniformly distributed across all classes. Initially, using a sample size of 5,000, we investigate 20 α values within the range $[10^{-1}, 10^1]$. Next, we fix α at either 1.0 or 0.1 and explore 50 different sample sizes ranging from 200 to 10,000. For each experiment, we run 100 trials and compute the mean squared error (MSE) between the true ratios and the estimated ratios. A two-layer MLP is used for MNIST, while ResNet-18 (He et al., 2016) is applied for CIFAR-10.

Method	Avg. accuracy
Our IW-ERM	0.7520 \pm 0.0209
Our IW-ERM (small)	0.7376 \pm 0.0099
FedAvg	0.5472 \pm 0.0297
FedBN	0.5359 \pm 0.0306
FedProx	0.5606 \pm 0.0070
SCAFFOLD	0.5774 \pm 0.0036
Upper Bound	0.8273 \pm 0.0041

Table 2: We utilize LeNet on Fashion MNIST to address label shifts across 5 nodes. For the baseline methods—FedAvg, FedBN, FedProx, and SCAFFOLD—we run 15,000 iterations, while both the Upper Bound (IW-ERM with true ratios) and our IW-ERM with VRLS are limited to 5,000 iterations. Notably, we employ a simple MLP with dropout for training the predictor. The model labeled *Our IW-ERM (small)* refers to our approach where the black-box predictor is trained using only 10% of the available training data, balancing computational efficiency with competitive performance.

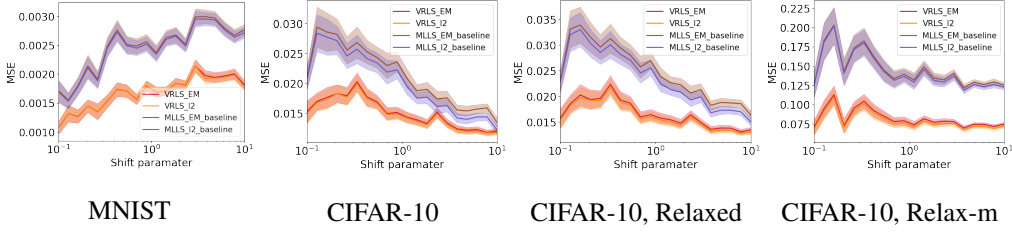


Figure 1: MSE analysis across different datasets and settings for VRLS (ours) compared to baselines, focusing on **shift parameter** (α) experiments. These subfigures include results from MNIST, CIFAR-10, and relaxed label shift, illustrating the consistent superiority of VRLS. In the ‘relaxed’ setting, Gaussian blur (kernel size: 3; σ : 0.1–0.5) and brightness adjustment (factor: ± 0.1) are applied with a 30% probability to introduce real-world variability. In the ‘relax-m’ scenario, augmentations are applied with a 50% probability, with Gaussian blur (σ : 0.1–0.7) and brightness (factor: ± 0.2).

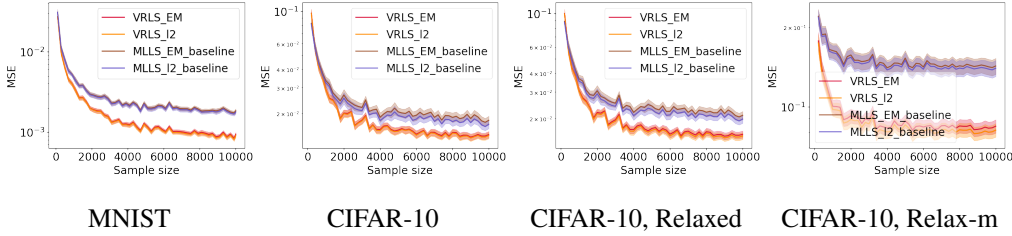


Figure 2: MSE analysis across different datasets and settings for VRLS (ours) compared to baselines, focusing on **sample size** experiments. These subfigures include results from MNIST, CIFAR-10, and relaxed label shift conditions, highlighting VRLS’s superior performance across varying test set sizes.

Table 3: We deploy ResNet-18 on CIFAR-10 to address label shifts across 5 nodes. The predictor is also a ResNet-18, ensuring consistency with the single-node scenario. For a fair comparison, we limit IW-ERM with VRLS and the true ratios to 5,000 iterations, while FedAvg and FedBN are run for 10,000 iterations. Detailed results are provided in Table 7.

CIFAR-10	Our IW-ERM	FedAvg	FedBN	Upper Bound
Avg. accuracy	0.5640 ± 0.0241	0.4515 ± 0.0148	0.4263 ± 0.0975	0.5790 ± 0.0103

Table 4: We present the average node accuracies from the CIFAR-10 target shift experiment conducted with 100 and 200 nodes, where 5 nodes are randomly sampled to participate in each training round. Our IW-ERM with VRLS is run for 5,000 and 10,000 iterations, respectively, while both FedAvg and FedBN are run for 10,000 iterations each.

CIFAR-10	Our IW-ERM	FedAvg	FedBN
Avg. accuracy (100 nodes)	0.5354	0.3915	0.1537
Avg. accuracy (200 nodes)	0.6216	0.5942	0.1753

Figure 1 and Figure 2 compares our proposed VRLS with baselines (Garg et al., 2020; Saerens et al., 2002) under label shifts. MLLS.L2 refers to the MLLS method using convex optimization via SGD (Garg et al., 2020), while MLLS.EM employs the same objective function but is optimized using the EM algorithm (Saerens et al., 2002). Our proposed VRLS is optimized in a similar manner, resulting in VRLS.L2 and VRLS.EM, as shown in the figure. Our method consistently achieves lower MSE across different label shift intensities (α) and test sample sizes on both datasets. Notably, our importance ratio estimation experiments align with the error bound in Theorem 5.1, demonstrating that increasing the number of test samples improves estimation error at a rate proportional to the square root of the sample size. Additionally, the regularization term constrains the parameter space and reduces Rademacher complexity, leading to smoother predictions and improved model calibration, as supported by Section S2 in (Guo et al., 2017a). Both of them contribute to reduced estimation error.

We also tested importance ratio estimation under relaxed label shift conditions and found VRLS to exhibit greater robustness (see Appendix J.2 for detailed settings). Although this assumption holds broader potential for real-world applications, its precise alignment with real-world datasets requires further investigation—an important direction for future research that extends beyond the scope of this work.

Distributed learning settings. We apply VRLS in a distributed learning context, addressing both intra- and inter-node label shifts. The initial experiments involve 5 nodes, using predefined label distributions on Fashion MNIST (Xiao et al., 2017) and CIFAR-10, as shown in Tables 8- 9 in Appendix J.

We employ a simple MLP with dropout as the predictor for Fashion MNIST. For global training with IW-ERM, LeNet (LeCun et al., 1998) is used on Fashion MNIST, and ResNet-18 (Ramezani-Kebrya et al., 2023b) on CIFAR-10. All experiments are run with three random seeds, reporting the average accuracy across nodes. We compare IW-ERM with VRLS against baseline methods, including FedAvg (McMahan et al., 2017), FedBN (Li et al., 2021b), FedProx (Li et al., 2020), and SCAFFOLD (Karimireddy et al., 2020a), as well as IW-ERM with true ratios serving as an upper bound. Hyperparameters are kept consistent with those in (McMahan et al., 2017; Li et al., 2021b; Ramezani-Kebrya et al., 2023b).

Each node’s stochastic gradients are computed with a batch size of 64 and aggregated using the Adam optimizer. All experiments are run on a single GPU within an internal cluster. Both MLLS and VRLS use identical hyperparameters and training epochs for CIFAR-10 and Fashion MNIST, stopping once the classification loss reaches a predefined threshold on MNIST. We also conduct experiments with 100 and 200 nodes on CIFAR-10, where five nodes are randomly sampled each iteration to simulate more realistic distributed learning. In this case, IW-ERM with true ratios does not act as the upper bound due to the stochastic node sampling. The experiment is run once, and average accuracy across nodes is reported, with label distribution shown in Table 10 in Appendix J. Despite FedBN’s reported slow convergence (Ramezani-Kebrya et al., 2023b), we maintain 15,000 and 10,000 iterations for FedAvg and FedBN on Fashion MNIST and CIFAR-10, respectively, for fair comparison. However, IW-ERM is limited to 5,000 iterations using both true and estimated ratios due to faster convergence.

As shown in Table 2, IW-ERM achieves over 20% higher average accuracy than all baselines on Fashion MNIST, with only a third of the iterations. Notably, even with just 10% of the training data in the first round of global training, the performance remains comparable, demonstrating reduced training complexity. This improvement is attributed to the theoretical benefits of IW-ERM, the robustness of density estimation, and the fact that the aggregation of importance ratios reduces reliance on any single local estimate. Similarly, Table 3 shows that IW-ERM approaches the upper bound on CIFAR-10, outperforming the baselines. Individual node accuracies are detailed in Tables 6-7 in Appendix J. In the 100-node scenario, IW-ERM continues to demonstrate superior performance, requiring only half the iterations, as shown in Table 4. It is important to note that using true ratios does not equate to IW-ERM, given the stochasticity of node selection during training.

7 CONCLUSIONS AND LIMITATIONS

We propose VRLS to address label shift in distributed learning. Paired with IW-ERM, VRLS improves intra- and inter-node label shifts in multi-node settings. Empirically, VRLS consistently outperforms MLLS-based baselines, and IW-ERM with VRLS exceeds all multi-node learning baselines. Theoretical bounds further strengthen our method’s foundation. Future work will explore estimating ratios by relaxing the strict class-conditional assumption and optimizing IW-ERM to reduce time complexity while ensuring scalability and practicality in real-world distributed learning.

ETHICS STATEMENT

No ethical approval was needed as no human subjects were involved. All authors fully support the content and findings.

REPRODUCIBILITY STATEMENT

We ensured reproducibility with publicly available datasets (MNIST, CIFAR-10) and standard models (e.g., ResNet-18). Links to datasets, code, and configurations will be provided upon camera-ready submission. Experiments were run on NVIDIA 3090, A100 GPUs, and Google Colab, with average results and variances reported across multiple trials.

ACKNOWLEDGMENTS

The authors would like to thank Leello Tadesse Dadi and Thomas Pethick for helpful discussions. This work was supported by the Research Council of Norway (RCN) through its Centres of Excellence scheme, Integreat: Norwegian Centre for knowledge-driven machine learning, project number 332645. The work of Changkyu Choi was funded by RCN under grant 309439. The work of Volkan Cevher was supported by Hasler Foundation Program: Hasler Responsible AI (project number 21043), the Army Research Office which was accomplished under Grant Number W911NF-24-1-0048, and the Swiss National Science Foundation (SNSF) under grant number 200021_205011. The computations were performed in part on resources provided by Sigma2 - the National Infrastructure for High-Performance Computing and Data Storage in Norway.

REFERENCES

- Rohit Agrawal and Thibaut Horel. Optimal bounds between f -divergences and integral probability metrics. In *International Conference on Machine Learning (ICML)*, 2020.
- Amr M. Alexandari, Anshul Kundaje, and Avanti Shrikumar. Maximum likelihood with bias-corrected calibration is hard-to-beat at label shift adaptation. In *International Conference on Machine Learning (ICML)*, 2020.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning (ICML)*, 2017.
- Kamyar Azizzadenesheli, Anqi Liu, Fanny Yang, and Animashree Anandkumar. Regularized learning for domain adaptation under label shifts. In *International Conference on Learning Representations (ICLR)*, 2019.
- Arindam Banerjee, Srujana Merugu, Inderjit S. Dhillon, and Joydeep Ghosh. Clustering with bregman divergences. *Journal of Machine Learning Research (JMLR)*, 6(58):1705–1749, 2005.
- Jeremiah Birrell, Paul Dupuis, Markos Katsoulakis, Yannis Pantazis, and Luc Rey-Bellet. (f, Γ) -Divergences: Interpolating between f -divergences and integral probability metrics. *Journal of Machine Learning Research (JMLR)*, 23:1–70, 2022.
- Jonathon Byrd and Zachary C. Lipton. What is the effect of importance weighting in deep learning? In *International Conference on Machine Learning (ICML)*, 2019.
- Artur Back de Luca, Guojun Zhang, Xi Chen, and Yaoliang Yu. Mitigating data heterogeneity in federated learning with data augmentation. *arXiv preprint arXiv:2206.09979*, 2022.
- Jérôme Dedecker, Clémentine Prieur, and Paul Raynaud de Fitte. *Parametrized Kantorovich-Rubinstein theorem and application to the coupling of random variables*. Springer, 2006.
- Fartash Faghri, Iman Tabrizian, Ilia Markov, Dan Alistarh, Daniel M. Roy, and Ali Ramezani-Kebrya. Adaptive gradient quantization for data-parallel SGD. In *Advances in neural information processing systems (NeurIPS)*, 2020.
- Tongtong Fang, Nan Lu, Gang Niu, and Masashi Sugiyama. Rethinking importance weighting for deep learning under distribution shift. *Advances in neural information processing systems (NeurIPS)*, 33:11996–12007, 2020.
- Tongtong Fang, Nan Lu, Gang Niu, and Masashi Sugiyama. Generalizing importance weighting to a universal solver for distribution shift problems. *Advances in neural information processing systems (NeurIPS)*, 36, 2023.

- Saurabh Garg, Yifan Wu, Sivaraman Balakrishnan, and Zachary C. Lipton. A unified view of label shift estimation. In *Advances in neural information processing systems (NeurIPS)*, 2020.
- Saurabh Garg, Sivaraman Balakrishnan, and Zachary C. Lipton. Domain adaptation under open set label shift. *arXiv preprint arXiv:2207.13048*, 2022.
- Saurabh Garg, Nick Erickson, James Sharpnack, Alexander J. Smola, Sivaraman Balakrishnan, and Zachary C. Lipton. Rlsbench: Domain adaptation under relaxed label shift. In *International Conference on Machine Learning (ICML)*, 2023.
- Margalit R. Glasgow, Honglin Yuan, and Tengyu Ma. Sharp bounds for federated averaging (local SGD) and continuous perspective. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2022.
- Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(25):723–773, 2012.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML’17, pp. 1321–1330. JMLR.org, 2017a.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML’17, pp. 1321–1330. JMLR.org, 2017b.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML’17, pp. 1321–1330. JMLR.org, 2017c.
- Jiaxian Guo, Mingming Gong, Tongliang Liu, Kun Zhang, and Dacheng Tao. LTF: A label transformation framework for correcting label shift. In *International Conference on Machine Learning (ICML)*, 2020.
- Sharut Gupta, Kartik Ahuja, Mohammad Havaei, Niladri Chatterjee, and Yoshua Bengio. Fl games: A federated learning framework for distribution shifts. *arXiv preprint arXiv:2205.11101*, 2022.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. Federated learning with compression: Unified analysis and sharp guarantees. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Zhengmian Hu and Heng Huang. Tighter analysis for ProxSkip. In *International Conference on Machine Learning (ICML)*, 2023.
- Jiayuan Huang, Arthur Gretton, Karsten Borgwardt, Bernhard Schölkopf, and Alex Smola. Correcting sample selection bias by unlabeled data. In *Advances in neural information processing systems (NeurIPS)*, 2006.
- Yutao Huang, Lingyang Chu, Zirui Zhou, Lanjun Wang, Jiangchuan Liu, Jian Pei, and Yong Zhang. Personalized cross-silo federated learning on non-iid data. In *AAAI Conference on Artificial Intelligence*, 2021.
- P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. D’Oliveira, S. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, G. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P., M. Mohri, R. Nock, A. Özgür, R. Pagh, M. Raykova, H. Qi, D. Ramage, R. Raskar, D. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.

- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5132–5143. PMLR, 13–18 Jul 2020a.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5132–5143. PMLR, 13–18 Jul 2020b.
- Mikhail Khodak, Maria-Florina Balcan, and Ameet Talwalkar. Adaptive gradient-based meta-learning methods. In *Advances in neural information processing systems (NeurIPS)*, 2019.
- Jinkyu Kim, Geeho Kim, and Bohyung Han. Multi-level branched regularization for federated learning. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 11058–11073. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/kim22a.html>.
- A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- Meelis Kull, Miquel Perello-Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter Flach. *Beyond temperature scaling: obtaining well-calibrated multiclass probabilities with Dirichlet calibration*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- S. Kullback and R. A. Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1):79 – 86, 1951.
- Gil Kur, Eli Putterman, and Alexnader Rakhlin. On the variance, admissibility, and stability of empirical risk minimization. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS ’23, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86:2278–2324, 1998.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In Inderjit S. Dhillon, Dimitris S. Papailiopoulos, and Vivienne Sze (eds.), *Proceedings of Machine Learning and Systems 2020, MLSys 2020, Austin, TX, USA, March 2-4, 2020*. mlsys.org, 2020.
- Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *International Conference on Machine Learning (ICML)*, 2021a.
- Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. FedBN: Federated learning on non-IID features via local batch normalization. In *International Conference on Learning Representations (ICLR)*, 2021b.
- Zachary C. Lipton, Yu-Xiang Wang, and Alexander J. Smola. Detecting and correcting for label shift with black box predictors. In *International Conference on Machine Learning (ICML)*, 2018.
- Zijian Liu, Ta Duy Nguyen, Thien Hang Nguyen, Alina Ene, and Huy Nguyen. High probability convergence of stochastic gradient methods. In *International Conference on Machine Learning (ICML)*, 2023.
- Kangyang Luo, Xiang Li, Yunshi Lan, and Ming Gao. Gradma: A gradient-memory-based accelerated federated learning with alleviated catastrophic forgetting. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3708–3717, 2023. doi: 10.1109/CVPR52729.2023.00361.

- You-Wei Luo and Chuan-Xian Ren. Generalized label shift correction via minimum uncertainty principle: Theory and algorithm. *ArXiv*, abs/2202.13043, 2022. URL <https://api.semanticscholar.org/CorpusID:247158776>.
- Pranav Mani, Manley Roberts, Saurabh Garg, and Zachary C. Lipton. Unsupervised learning under latent label shift. In *ICML Workshop on Spurious Correlations, Invariance and Stability*, 2022.
- H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera Y. Arcas. Communication-efficient learning of deep networks from decentralized data. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.
- Dexter Neo, Stefan Winkler, and Tsuhan Chen. Maxent loss: Constrained maximum entropy for calibration under out-of-distribution shift. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 21463–21472, 2024.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint arXiv:1701.06548*, 2017.
- Stephan Rabanser, Stephan Günnemann, and Zachary C. Lipton. *Failing Loudly: An Empirical Study of Methods for Detecting Dataset Shift*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- Anichur Rahman, Md Sazzad Hossain, Ghulam Muhammad, Dipanjali Kundu, Tanoy Debnath, Muaz Rahman, Md Saikat Islam Khan, Prayag Tiwari, and Shahab S Band. Federated learning-based ai approaches in smart healthcare: concepts, taxonomies, challenges and open issues. *Cluster computing*, 26(4):2271–2311, 2023.
- Suraj Rajendran, Zhenxing Xu, Weishen Pan, Arnab Ghosh, and Fei Wang. Data heterogeneity in federated learning with electronic health records: Case studies of risk prediction for acute kidney injury and sepsis diseases in critical care. *PLOS Digital Health*, 2(3):e0000117, 2023.
- Ali Ramezani-Kebrya, Fartash Faghri, Ilya Markov, Vitalii Aksenov, Dan Alistarh, and Daniel M. Roy. NUQSGD: Provably communication-efficient data-parallel SGD via nonuniform quantization. *Journal of Machine Learning Research (JMLR)*, 22(114):1–43, 2021.
- Ali Ramezani-Kebrya, Kimon Antonakopoulos, Igor Krawczuk, Justin Deschenaux, and Volkan Cevher. Distributed extra-gradient with optimal complexity and communication guarantees. In *International Conference on Learning Representations (ICLR)*, 2023a.
- Ali Ramezani-Kebrya, Fanghui Liu, Thomas Pethick, Grigorios Chrysos, and Volkan Cevher. Federated learning under covariate shifts with generalization guarantees. *Transactions on Machine Learning Research (TMLR)*, 2023b.
- Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do CIFAR-10 classifiers generalize to CIFAR-10? *arXiv preprint arXiv:1806.00451*, 2018.
- Marco Saerens, Patrice Latinne, and Christine Decaestecker. Adjusting the outputs of a classifier to new a priori probabilities: a simple procedure. In *Neural Computation*, 2002.
- Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90(2):227–244, 2000.
- Avanti Shrikumar, Amr M. Alexandari, and Anshul Kundaje. Adapting to label shift with bias-corrected calibration. *arXiv preprint arXiv:1901.06852v5*, 2019.
- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S. Talwalkar. Federated multi-task learning. In *Advances in neural information processing systems (NeurIPS)*, 2017.
- Masashi Sugiyama, Benjamin Blankertz, Matthias Krauledat, Guido Dornhege, and Klaus-Robert Müller. Importance-weighted cross-validation for covariate shift. In *Joint Pattern Recognition Symposium*, pp. 354–363. Springer, 2006.
- Masashi Sugiyama, Matthias Krauledat, and Klaus-Robert Müller. Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research (JMLR)*, 8(5), 2007.

- Zeyu Sun, Dogyoon Song, and Alfred Hero. Minimum-risk recalibration of classifiers. *Advances in neural information processing systems (NeurIPS)*, 36, 2024.
- Antonio Torralba, Rob Fergus, and William T. Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(11):1958–1970, 2008.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*. Cambridge university press, 2018.
- Cédric Villani. *The Wasserstein distances*. Springer Berlin Heidelberg, 2009.
- Deng-Bao Wang, Lei Feng, and Min-Ling Zhang. Rethinking calibration of deep neural networks: Do not be afraid of overconfidence. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in neural information processing systems (NeurIPS)*, volume 34, pp. 11809–11820. Curran Associates, Inc., 2021.
- Kaibin Wang, Qiang He, Feifei Chen, Chunyang Chen, Faliang Huang, Hai Jin, and Yun Yang. Flexified: Personalized federated learning for edge clients with heterogeneous model architectures. In *Proceedings of the ACM Web Conference 2023*, pp. 2979–2990, 2023.
- Jie Wen, Zhixia Zhang, Yang Lan, Zhihua Cui, Jianghui Cai, and Wensheng Zhang. A survey on federated learning: challenges and applications. *International Journal of Machine Learning and Cybernetics*, 14(2):513–535, 2023.
- Blake Woodworth, Kumar Kshitij Patel, Sebastian Stich, Zhen Dai, Brian Bullins, Brendan McMahan, Ohad Shamir, and Nathan Srebro. Is local SGD better than minibatch SGD? In *International Conference on Machine Learning (ICML)*, 2020.
- Xidong Wu, Feihu Huang, Zhengmian Hu, and Heng Huang. Faster adaptive federated learning. In *AAAI Conference on Artificial Intelligence*, 2023.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Mang Ye, Xiuwen Fang, Bo Du, Pong C Yuen, and Dacheng Tao. Heterogeneous federated learning: State-of-the-art and research challenges. *ACM Computing Surveys*, 56(3):1–44, 2023.
- Mingzhang Yin, Yixin Wang, and David M Blei. Optimization-based causal estimation from heterogeneous environments. *Journal of Machine Learning Research*, 25:1–44, 2024.
- Bianca Zadrozny. Learning and evaluating classifiers under sample selection bias. In *International Conference on Machine Learning (ICML)*, 2004.
- Helen Zhou, Sivaraman Balakrishnan, and Zachary C. Lipton. Domain adaptation under missingness shift. *Artificial Intelligence and Statistics (AISTATS)*, 2023.

The Appendix part is organized as follows:

- All related work are provided in Appendix A.
- Additional details of prior work of BBSE and MLLS are in Appendix B.
- Mathematical proof for label shifts with multiple nodes and IW-ERM is given in Appendix C.
- General algorithmic description is in Appendix D.
- Proof of Theorem 5.1 is in Appendix E.
- Proof of Theorem 5.2 and Convergence-Communication-Privacy guarantees for IW-ERM in Equation (IW-ERM) are provided in Appendix F.
- Complexity analysis is in Appendix G.
- Mathematical notations are summarized in Appendix H.
- Limitations are discussed in Appendix I.
- Additional experiments and experimental details are provided in Appendix J.

A RELATED WORK

In the context of distributed learning with label shifts, importance ratio estimation is tackled either by solving a linear system as in (Lipton et al., 2018; Azizzadenesheli et al., 2019) or by minimizing distribution divergence as in (Garg et al., 2020). In this section, we overview complete related work.

Federated learning (FL). Much of the current research in FL predominantly centers around the minimization of empirical risk, operating under the assumption that each node maintains the same training/test data distribution (Kairouz et al., 2021). Prominent methods in FL include FedAvg (McMahan et al., 2017), FedBN (Li et al., 2021b), FedProx (Li et al., 2020) and SCAFFOLD (Karimireddy et al., 2020a). FedAvg and its variants such as (Huang et al., 2021; Karimireddy et al., 2020b) have been the subject of thorough investigation in optimization literature, exploring facets such as communication efficiency, node participation, and privacy assurance (Faghri et al., 2020; Ramezani-Kebrya et al., 2021; 2023a;b). Subsequent work, such as the study by de Luca et al. (2022), explores Federated Domain Generalization and introduces data augmentation to the training. This model aims to generalize to both in-domain datasets from participating nodes and an out-of-domain dataset from a non-participating node. Additionally, Gupta et al. (2022) introduces FL Games, a game-theoretic framework designed to learn causal features that remain invariant across nodes. This is achieved by employing ensembles over nodes’ historical actions and enhancing local computation, under the assumption of consistent training/test data distribution across nodes. The existing strategies to address statistical heterogeneity across nodes during training primarily rely on heuristic-based personalization methods, which currently lack theoretical backing in statistical learning (Smith et al., 2017; Khodak et al., 2019; Li et al., 2021a). In contrast, we aim to minimize overall test error amid both intra-node and inter-node distribution shifts, a situation frequently observed in real-world scenarios. Techniques ensuring communication efficiency, robustness, and secure aggregations serve as complementary.

Importance ratio estimation Classical Empirical Risk Minimization (ERM) seeks to minimize the expected loss over the training distribution using finite samples. When faced with distribution shifts, the goal shifts to minimizing the expected loss over the target distribution, leading to the development of Importance-Weighted Empirical Risk Minimization (IW-ERM) (Shimodaira, 2000; Sugiyama et al., 2006; Byrd & C. Lipton, 2019; Fang et al., 2020). Shimodaira (2000) established that the IW-ERM estimator is asymptotically unbiased. Moreover, Ramezani-Kebrya et al. (2023b) introduced FTW-ERM, which integrates importance ratio estimation.

Label shift and MLLS family For theoretical analysis, the conditional distribution $p(\mathbf{x}|\mathbf{y})$ is held strictly constant across all distributions (Lipton et al., 2018; Garg et al., 2020; Saelens et al., 2002). Both BBSE (Lipton et al., 2018) and RLLS (Azizzadenesheli et al., 2019) designate a discrete latent space \mathbf{z} and introduce a confusion matrix-based estimation method to compute the ratio w by solving a linear system (Saelens et al., 2002; Lipton et al., 2018). This approach is straightforward and has been proven consistent, even when the predictor is not calibrated. However, its subpar performance is attributed to the information loss inherent in the confusion matrix (Garg et al., 2020).

Consequently, MLLS (Garg et al., 2020) introduces a continuous latent space, resulting in a significant enhancement in estimation performance, especially when combined with a post-hoc calibration method (Shrikumar et al., 2019). It also provides a consistency guarantee with a canonically calibrated predictor. This EM-based MLLS method is both concave and can be solved efficiently.

Discrepancy Measure In information theory and statistics, discrepancy measures play a critical role in quantifying the differences between probability distributions. One such measure is the Bregman Divergence (Banerjee et al., 2005), defined as

$$D_\phi(\mathbf{x}||\mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla \phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle,$$

which encapsulates the difference between the value of a convex function ϕ at two points and the value of the linear approximation of ϕ at one point, leveraging the gradient at another point.

Discrepancy measures are generally categorized into two main families: Integral Probability Metrics (IPMs) and f -divergences. IPMs, including Maximum Mean Discrepancy (Gretton et al., 2012) and Wasserstein distance (Villani, 2009), focus on distribution differences $P - Q$. In contrast, f -divergences, such as KL-divergence (Kullback & Leibler, 1951) and Total Variation distance, operate

on ratios P/Q and do not satisfy the triangular inequality. Interconnections and variations between these families are explored in studies like (f, Γ) -Divergences (Birrell et al., 2022), which interpolate between f -divergences and IPMs, and research outlining optimal bounds between them (Agrawal & Horel, 2020).

MLLS (Garg et al., 2020) employs f -divergence, notably the KL divergence, which is not a metric as it doesn't satisfy the triangular inequality, and requires distribution P to be absolutely continuous with respect to Q . Concerning IPMs, while MMD is reliant on a kernel function, it can suffer from the curse of dimensionality when faced with high-dimensional data. On the other hand, the Wasserstein distance can be reformulated using Kantorovich-Rubinstein duality (Dedecker et al., 2006; Arjovsky et al., 2017) as a maximization problem subject to a Lipschitz constrained function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

B BBSE AND MLLS FAMILY

In this section, we summarize the contributions of BBSE (Lipton et al., 2018) and MLLS (Garg et al., 2020). Our objective is to estimate the ratio $p^{\text{te}}(y)/p^{\text{tr}}(y)$. We consider a scenario with m possible label classes, where $y = c$ for $c \in [m]$. Let $\mathbf{r}^* = [r_1^*, \dots, r_m^*]^\top$ represent the true ratios, with each r_c^* defined as $r_c^* = \frac{p^{\text{te}}(y=c)}{p^{\text{tr}}(y=c)}$ (Garg et al., 2020). We then define a family of distributions over \mathcal{Z} , parameterized by $\mathbf{r} = [r_1, \dots, r_m]^\top \in \mathbb{R}^m$, where r_c is the c -th element of the ratio vector.

$$p_{\mathbf{r}}(\mathbf{z}) := \sum_{c=1}^m p^{\text{te}}(\mathbf{z}|y=c) \cdot p^{\text{tr}}(y=c) \cdot r_c \quad (10)$$

Here, $r_c \geq 0$ for $c \in [m]$ and $\sum_{c=1}^m r_c \cdot p^{\text{tr}}(y=c) = \sum_{c=1}^m p^{\text{te}}(y=c) = 1$ as constraints. When $\mathbf{r} = \mathbf{r}^*$, e.g., $r_c = r_c^*$ for $c \in [m]$, we have $p_{\mathbf{r}}(\mathbf{z}) = p_{\mathbf{r}^*}(\mathbf{z}) = p^{\text{te}}(\mathbf{z})$ (Garg et al., 2020). So our task is to find \mathbf{r} such that

$$\begin{aligned} & \sum_{c=1}^m p^{\text{te}}(\mathbf{z}|y=c) \cdot p^{\text{tr}}(y=c) \cdot r_c \mathbf{x} \\ &= \sum_{c=1}^m p^{\text{tr}}(\mathbf{z}, y=c) \cdot r_c = p^{\text{te}}(\mathbf{z}) \end{aligned} \quad (11)$$

Lipton et al. (2018) introduced Black Box Shift Estimation (BBSE) to address this issue. With a pre-trained classifier f for the classification task, BBSE assumes that the latent space \mathcal{Z} is discrete and defines $p(\mathbf{z}|\mathbf{x}) = \delta_{\arg \max f(\mathbf{x})}$, where the output of $f(\mathbf{x})$ is a probability vector (or a simplex) over m classes. BBSE estimates $p^{\text{te}}(\mathbf{z}|y)$ as a confusion matrix, using both the training and validation data. It calculates $p^{\text{tr}}(y=c)$ from the training set and $p^{\text{te}}(\mathbf{z})$ from the test data. The problem then reduces to solving the following equation:

$$\mathbf{A}\mathbf{w} = \mathbf{B} \quad (12)$$

where $|\mathcal{Z}| = m$, $\mathbf{A} \in \mathbb{R}^{m \times m}$ with $A_{jc} = p^{\text{te}}(\mathbf{z}=j|y=c) \cdot p^{\text{tr}}(y=c)$, and $\mathbf{B} \in \mathbb{R}^m$ with $B_j = p^{\text{te}}(\mathbf{z}=j)$ for $c, j \in [m]$.

The estimation of the confusion matrix in terms of $p^{\text{te}}(\mathbf{z}|y)$ leads to the loss of calibration information (Garg et al., 2020). Furthermore, when defining \mathcal{Z} as a continuous latent space, the confusion matrix becomes intractable since \mathbf{z} has infinitely many values. Therefore, MLLS directly minimizes the divergence between $p^{\text{te}}(\mathbf{z})$ and $p_{\mathbf{r}}(\mathbf{z})$, instead of solving the linear system in Equation (12).

Within the f -divergence family, MLLS seeks to find a weight vector \mathbf{r} by minimizing the KL-divergence $D_{\text{KL}}(p^{\text{te}}(\mathbf{z}), p_{\mathbf{r}}(\mathbf{z})) = \mathbb{E}_{\text{te}} [\log p^{\text{te}}(\mathbf{z})/p_{\mathbf{r}}(\mathbf{z})]$, for $p_{\mathbf{r}}(\mathbf{z})$ defined in Equation (10). Leveraging on the properties of the logarithm, this is equivalent to maximizing the log-likelihood: $\mathbf{r} := \arg \max_{\mathbf{r} \in \mathbb{R}} \mathbb{E}_{\text{te}} [\log p_{\mathbf{r}}(\mathbf{z})]$. Expanding $p_{\mathbf{r}}(\mathbf{z})$, we have

$$\begin{aligned} \mathbb{E}_{\text{te}} [\log p_{\mathbf{r}}(\mathbf{z})] &= \mathbb{E}_{\text{te}} \left[\log \left(\sum_{c=1}^m p^{\text{tr}}(\mathbf{z}, y=c) r_c \right) \right] \\ &= \mathbb{E}_{\text{te}} \left[\log \left(\sum_{c=1}^m p^{\text{tr}}(y=c | \mathbf{z}) r_c \right) + \log p^{\text{tr}}(\mathbf{z}) \right]. \end{aligned} \quad (13)$$

Therefore the unified form of MLLS can be formulated as:

$$\mathbf{r} := \arg \max_{\mathbf{r} \in \mathbb{R}} \mathbb{E}_{\text{te}} \left[\log \left(\sum_{c=1}^m p^{\text{tr}}(y=c | \mathbf{z}) r_c \right) \right]. \quad (14)$$

This is a convex optimization problem and can be solved efficiently using methods such as EM, an analytic approach, and also iterative optimization methods like gradient descent with labeled training data and unlabeled test data. MLLS defines the $p(\mathbf{z}|\mathbf{x})$ as $\delta_{\mathbf{x}}$, plugs in the pre-defined f to approximate $p^{\text{tr}}(y|\mathbf{x})$ and optimizes the following objective:

$$\mathbf{r}_f := \arg \max_{\mathbf{r} \in \mathbb{R}} \ell(\mathbf{r}, f) := \arg \max_{\mathbf{r} \in \mathbb{R}} \mathbb{E}_{\text{te}} [\log(f(\mathbf{x})^T \mathbf{r})] . \quad (15)$$

With the Bias-Corrected Calibration (BCT) (Shrikumar et al., 2019) strategy, they adjust the logits $\hat{f}(\mathbf{x})$ of $f(\mathbf{x})$ element-wise for each class, and the objective becomes:

$$\mathbf{r}_f := \arg \max_{\mathbf{r} \in \mathbb{R}} \ell(\mathbf{r}, f) := \arg \max_{\mathbf{r} \in \mathbb{R}} \mathbb{E}_{\text{te}} [\log(g \circ \hat{f}(\mathbf{x}))^T \mathbf{r}] , \quad (16)$$

where g is a calibration function.

C PROOF OF PROPOSITION 4.1

In the following, we consider four typical scenarios under various distribution shifts described in Table 1 and formulate their IW-ERM with a focus on minimizing R_1 .

C.1 NO INTRA-NODE LABEL SHIFT

For simplicity, we assume that there are only 2 nodes, but our results can be extended to multiple nodes. This scenario assumes $p_k^{\text{tr}}(\mathbf{y}) = p_k^{\text{te}}(\mathbf{y})$ for $k = 1, 2$, but $p_1^{\text{tr}}(\mathbf{y}) \neq p_2^{\text{tr}}(\mathbf{y})$. Node 1 aims to learn $h_{\mathbf{w}}$ assuming $\frac{p_1^{\text{tr}}(\mathbf{y})}{p_2^{\text{tr}}(\mathbf{y})}$ is given. We consider the following IW-ERM that is consistent in minimizing R_1 :

$$\begin{aligned} \min_{h_{\mathbf{w}} \in \mathcal{H}} \frac{1}{n_1^{\text{tr}}} \sum_{i=1}^{n_1^{\text{tr}}} \ell(h_{\mathbf{w}}(\mathbf{x}_{1,i}^{\text{tr}}), \mathbf{y}_{1,i}^{\text{tr}}) \\ + \frac{1}{n_2^{\text{tr}}} \sum_{i=1}^{n_2^{\text{tr}}} \frac{p_1^{\text{tr}}(\mathbf{y}_{2,i}^{\text{tr}})}{p_2^{\text{tr}}(\mathbf{y}_{2,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{2,i}^{\text{tr}}), \mathbf{y}_{2,i}^{\text{tr}}). \end{aligned} \quad (17)$$

Here \mathcal{H} is the hypothesis class of $h_{\mathbf{w}}$. This scenario is referred to as No-LS.

C.2 LABEL SHIFT ONLY FOR NODE 1

Here we consider label shift only for node 1, i.e., $p_1^{\text{tr}}(\mathbf{y}) \neq p_1^{\text{te}}(\mathbf{y})$ and $p_2^{\text{tr}}(\mathbf{y}) = p_2^{\text{te}}(\mathbf{y})$. We consider the following IW-ERM:

$$\begin{aligned} \min_{h_{\mathbf{w}} \in \mathcal{H}} \frac{1}{n_1^{\text{tr}}} \sum_{i=1}^{n_1^{\text{tr}}} \frac{p_1^{\text{te}}(\mathbf{y}_{1,i}^{\text{tr}})}{p_1^{\text{tr}}(\mathbf{y}_{1,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{1,i}^{\text{tr}}), \mathbf{y}_{1,i}^{\text{tr}}) \\ + \frac{1}{n_2^{\text{tr}}} \sum_{i=1}^{n_2^{\text{tr}}} \frac{p_1^{\text{te}}(\mathbf{y}_{2,i}^{\text{tr}})}{p_2^{\text{tr}}(\mathbf{y}_{2,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{2,i}^{\text{tr}}), \mathbf{y}_{2,i}^{\text{tr}}). \end{aligned} \quad (18)$$

This scenario is referred to as LS on single.

C.3 LABEL SHIFT FOR BOTH NODES

Here we assume $p_1^{\text{tr}}(\mathbf{y}) \neq p_1^{\text{te}}(\mathbf{y})$ and $p_2^{\text{tr}}(\mathbf{y}) \neq p_2^{\text{te}}(\mathbf{y})$, i.e., label shift for both nodes. The corresponding IW-ERM is the same as Eq. Equation (18). This scenario is referred to as LS on both.

Without loss of generality and for simplicity, we set $l = 1$. We consider four typical scenarios under various distribution shifts and formulate their IW-ERM with a focus on minimizing R_1 . The details of these scenarios are summarized in Table 1.

C.4 MULTIPLE NODES

Here we consider a general scenario with K nodes. We assume both intra-node and inter-node label shifts by the following IW-ERM:

$$\min_{h_{\mathbf{w}} \in \mathcal{H}} \sum_{k=1}^K \frac{1}{n_k^{\text{tr}}} \sum_{i=1}^{n_k^{\text{tr}}} \frac{p_k^{\text{te}}(\mathbf{y}_{k,i}^{\text{tr}})}{p_k^{\text{tr}}(\mathbf{y}_{k,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{k,i}^{\text{tr}}), \mathbf{y}_{k,i}^{\text{tr}}), \quad (19)$$

This scenario is referred to as LS on multi.

For the scenario without intra-node label shift, the IW-ERM in Equation (17) can be expressed as

$$\begin{aligned}
& \frac{1}{n_2^{\text{tr}}} \sum_{i=1}^{n_2^{\text{tr}}} \frac{p_1^{\text{tr}}(\mathbf{y}_{2,i}^{\text{tr}})}{p_2^{\text{tr}}(\mathbf{y}_{2,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{2,i}^{\text{tr}}), \mathbf{y}_{2,i}^{\text{tr}}) \\
& \xrightarrow{n_2^{\text{tr}} \rightarrow \infty} \mathbb{E}_{p_2^{\text{tr}}(\mathbf{x}, \mathbf{y})} \left[\frac{p_1^{\text{tr}}(\mathbf{y})}{p_2^{\text{tr}}(\mathbf{y})} \ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y}) \right] \\
& = \int_{\mathcal{Y}} \frac{p_1^{\text{tr}}(\mathbf{y})}{p_2^{\text{tr}}(\mathbf{y})} \mathbb{E}_{p(\mathbf{x}|\mathbf{y})} [\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] p_2^{\text{tr}}(\mathbf{y}) d\mathbf{y} \\
& = \int_{\mathcal{Y}} p_1^{\text{tr}}(\mathbf{y}) \mathbb{E}_{p(\mathbf{x}|\mathbf{y})} [\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] d\mathbf{y} \\
& = \int_{\mathcal{Y}} p_1^{\text{te}}(\mathbf{y}) \mathbb{E}_{p(\mathbf{x}|\mathbf{y})} [\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] d\mathbf{y} \\
& = \mathbb{E}_{p_1^{\text{te}}(\mathbf{x}, \mathbf{y})} [\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] \\
& = R_1(h_{\mathbf{w}}).
\end{aligned} \tag{20}$$

where the second equality holds due to the assumption of the label shift setting and Bayes' theorem: $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y}) \cdot p(\mathbf{y})$, and the fourth equality holds by the assumption that $p_1^{\text{tr}}(\mathbf{y}) = p_1^{\text{te}}(\mathbf{y})$ in the No-LS setting.

For the scenario with label shift only for Node 1 or for both nodes, the IW-ERM in Equation (18) admits

$$\frac{1}{n_2^{\text{tr}}} \sum_{i=1}^{n_2^{\text{tr}}} \frac{p_1^{\text{te}}(\mathbf{y}_{2,i}^{\text{tr}})}{p_2^{\text{tr}}(\mathbf{y}_{2,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{2,i}^{\text{tr}}), \mathbf{y}_{2,i}^{\text{tr}}) \tag{21}$$

$$\xrightarrow{n_2^{\text{tr}} \rightarrow \infty} \mathbb{E}_{p_2^{\text{tr}}(\mathbf{x}, \mathbf{y})} \left[\frac{p_1^{\text{te}}(\mathbf{y})}{p_2^{\text{tr}}(\mathbf{y})} \ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y}) \right] \tag{22}$$

$$= \int_{\mathcal{Y}} \frac{p_1^{\text{te}}(\mathbf{y})}{p_2^{\text{tr}}(\mathbf{y})} \mathbb{E}_{p(\mathbf{x}|\mathbf{y})} [\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] p_2^{\text{tr}}(\mathbf{y}) d\mathbf{y} \tag{23}$$

$$= \int_{\mathcal{Y}} p_1^{\text{te}}(\mathbf{y}) \mathbb{E}_{p(\mathbf{x}|\mathbf{y})} [\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] d\mathbf{y} \tag{24}$$

$$= \mathbb{E}_{p_1^{\text{te}}(\mathbf{x}, \mathbf{y})} [\ell(h_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] \tag{25}$$

$$= R_1(h_{\mathbf{w}}). \tag{26}$$

For multiple nodes, let $k \in [K]$. Similarly, we have

$$\frac{1}{n_k^{\text{tr}}} \sum_{i=1}^{n_k^{\text{tr}}} \frac{p_1^{\text{te}}(\mathbf{y}_{k,i}^{\text{tr}})}{p_k^{\text{tr}}(\mathbf{y}_{k,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{k,i}^{\text{tr}}), \mathbf{y}_{k,i}^{\text{tr}}) \xrightarrow{n_k^{\text{tr}} \rightarrow \infty} R_1(h_{\mathbf{w}}). \tag{27}$$

Then we have

$$\sum_{k=1}^K \frac{1}{n_k^{\text{tr}}} \sum_{i=1}^{n_k^{\text{tr}}} \frac{p_1^{\text{te}}(\mathbf{y}_{k,i}^{\text{tr}})}{p_k^{\text{tr}}(\mathbf{y}_{k,i}^{\text{tr}})} \ell(h_{\mathbf{w}}(\mathbf{x}_{k,i}^{\text{tr}}), \mathbf{y}_{k,i}^{\text{tr}}) \xrightarrow{n_1^{\text{tr}}, \dots, n_K^{\text{tr}} \rightarrow \infty} R_1(h_{\mathbf{w}}). \tag{28}$$

Note that to solve Equation (19), node 1 needs to estimate $\frac{p_1^{\text{te}}(\mathbf{y})}{p_k^{\text{tr}}(\mathbf{y})}$ for all nodes k in Equation (19).

The consistency of Equation (IW-ERM), i.e., convergence in probability, is followed the standard arguments in e.g., (Shimodaira, 2000)[Section 3] and (Sugiyama et al., 2007)[Section 2.2] using the law of large numbers.

D ALGORITHMIC DESCRIPTION

```

1
2 # Split the training dataset on each node
3 trainsets = target_shift.split_dataset(trainset.data, trainset.targets,
    node_label_dist_train, transform=transform_train)
4
5 # Split the test dataset on each node
6 testsets = target_shift.split_dataset(testset.data, testset.targets,
    node_label_dist_test, transform=transform_test)
7
8 # Initialize K local models (nets) for each node
9 nets = [initialize_model() for _ in range(node_num)]
10
11 # Initialize the estimator for each local model
12 estimators = [LS_RatioModel(nets[k]) for k in range(node_num)]
13
14 # Initialize tensors to store the estimated ratios, values, and marginal
    values for each pair of nodes.
15 estimated_ratios = torch.zeros(node_num, node_num, nclass)
16 estimated_values = torch.zeros(node_num, node_num, nclass)
17 marginal_values = torch.zeros(node_num, nclass)
18
19 # Phase 1: Compute the estimated ratios for each node pair (k, j)
20 for k in range(node_num):
21     for j in range(node_num):
22         # Perform test on node k using node j's testset
23         estimated_ratios[k, j] = estimators[k](testsets[j].data.cpu().
            numpy())
24
25 # Phase 2: Compute the marginal values on each node's training set
26 for i, trainset in enumerate(trainsets):
27     marginal_values[i] = marginal(trainset.targets)
28
29 # Phase 3: Compute the final estimated values for each node
30 for k in range(node_num):
31     for j in range(node_num):
32         estimated_values[k, j] = marginal_values[j] * estimated_ratios[k,
            j]
33
34 # Aggregate the estimated values across nodes
35 aggregated_values = torch.sum(estimated_values, dim=1)
36
37 # Compute the final ratios for each node
38 ratios = (aggregated_values / marginal_values).to(args.device)

```

Listing 1: Our VRLS in distributed learning. It is the implementation of Algorithm 2

E PROOF OF THEOREM 5.1

Proof. Let $H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) = -\log(f(\mathbf{x}, \boldsymbol{\theta})^\top \mathbf{r})$. From the strong convexity in Lemma E.7, we have that

$$\|\hat{\mathbf{r}}_{n^{\text{te}}} - \mathbf{r}_{f^*}\|_2^2 \leq \frac{2}{\mu p_{\min}} (\mathcal{L}_{\boldsymbol{\theta}^*}(\hat{\mathbf{r}}_{n^{\text{te}}}) - \mathcal{L}_{\boldsymbol{\theta}^*}(\mathbf{r}_{f^*})) \quad (29)$$

Now focusing on the term on the right-hand side, we find by invoking Lemma E.4 that

$$\begin{aligned} & \mathcal{L}_{\boldsymbol{\theta}^*}(\hat{\mathbf{r}}_{n^{\text{te}}}) - \mathcal{L}_{\boldsymbol{\theta}^*}(\mathbf{r}_{f^*}) \\ & \leq \mathbb{E} \left[H(\hat{\mathbf{r}}_{n^{\text{te}}}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}) \right] - \mathbb{E} \left[H(\mathbf{r}_{f^*}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}) \right] + 2L \mathbb{E} \left[\|\hat{\boldsymbol{\theta}}_{n^{\text{tr}}} - \boldsymbol{\theta}^*\|_2 \right] \\ & = \mathbb{E} \left[H(\hat{\mathbf{r}}_{n^{\text{te}}}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}) \right] - \frac{1}{n^{\text{te}}} \sum_{j=1}^{n^{\text{te}}} H(\hat{\mathbf{r}}_{n^{\text{te}}}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}_j) + \frac{1}{n^{\text{te}}} \sum_{j=1}^{n^{\text{te}}} H(\hat{\mathbf{r}}_{n^{\text{te}}}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}_j) \\ & \quad - \mathbb{E} \left[H(\mathbf{r}_{f^*}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}) \right] + 2L \mathbb{E} \left[\|\hat{\boldsymbol{\theta}}_{n^{\text{tr}}} - \boldsymbol{\theta}^*\|_2 \right] \\ & \leq \mathbb{E} \left[H(\hat{\mathbf{r}}_{n^{\text{te}}}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}) \right] - \frac{1}{n^{\text{te}}} \sum_{j=1}^{n^{\text{te}}} H(\hat{\mathbf{r}}_{n^{\text{te}}}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}_j) + \frac{1}{n^{\text{te}}} \sum_{j=1}^{n^{\text{te}}} H(\mathbf{r}_{f^*}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}_j) \\ & \quad - \mathbb{E} \left[H(\mathbf{r}_{f^*}, \hat{\boldsymbol{\theta}}_{n^{\text{tr}}}, \mathbf{x}) \right] + 2L \mathbb{E} \left[\|\hat{\boldsymbol{\theta}}_{n^{\text{tr}}} - \boldsymbol{\theta}^*\|_2 \right], \end{aligned} \quad (30)$$

where in the last inequality we used the fact that $\hat{\mathbf{r}}_n$ is a minimizer of $\mathbf{r} \mapsto \frac{1}{n} \sum_{j=1}^n H(\mathbf{r}, \hat{\boldsymbol{\theta}}_t, \mathbf{x}_j)$. Finally by using Lemma E.5 and Lemma E.6 with $\delta/2$ each, we have that with probability $1 - \delta$,

$$\mathcal{L}_{\boldsymbol{\theta}^*}(\hat{\mathbf{r}}_{n^{\text{te}}}) - \mathcal{L}_{\boldsymbol{\theta}^*}(\mathbf{r}_{f^*}) \leq \frac{4}{\sqrt{n^{\text{te}}}} \text{Rad}(\mathcal{F}) + 2L \mathbb{E} \left[\|\hat{\boldsymbol{\theta}}_{n^{\text{tr}}} - \boldsymbol{\theta}^*\|_2 \right] + 4B \sqrt{\frac{\log(4/\delta)}{n^{\text{te}}}} \quad (31)$$

Plugging this back into Equation (29), we have that

$$\|\hat{\mathbf{r}}_{n^{\text{te}}} - \mathbf{r}_{f^*}\|_2^2 \leq \frac{2}{\mu p_{\min}} \left(\frac{4}{\sqrt{n^{\text{te}}}} \text{Rad}(\mathcal{F}) + 4B \sqrt{\frac{\log(4/\delta)}{n^{\text{te}}}} \right) + \frac{4L}{\mu p_{\min}} \mathbb{E} \left[\|\hat{\boldsymbol{\theta}}_{n^{\text{tr}}} - \boldsymbol{\theta}^*\|_2 \right]. \quad (32)$$

□

Lemma E.1. For any $\mathbf{r} \in \mathbb{R}_+^m$, $\boldsymbol{\theta} \in \Theta$, $\mathbf{x} \in \mathcal{X}$, we have that

$$\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}) \leq \frac{1}{p_{\min}}.$$

Proof. Applying Hölder's inequality we have that

$$\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}) \leq \|\mathbf{r}\|_\infty \|f(\mathbf{x}, \boldsymbol{\theta})\|_1 = \|\mathbf{r}\|_\infty.$$

Moreover, since $\mathbf{r} \in \mathbb{R}_+^m$, we have that $\sum_y r_y p_{tr}(y) = 1$. This implies that $\|\mathbf{r}\|_\infty \leq \frac{1}{p_{\min}}$, which yields the result. □

Lemma E.2 (Implication of Assumption 5.1). Under Assumption 5.1, there exists $B > 0$ such that for any $\mathbf{r} \in \mathbb{R}_+^m$, $\boldsymbol{\theta} \in \Theta$, $\mathbf{x} \in \mathcal{X}$,

$$|\log(\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}))| \leq B.$$

Proof. Since $\mathbf{r} \in \mathbb{R}_+^m$, it has at least one non-zero coordinate and $f(\mathbf{x}, \boldsymbol{\theta})$ is the output of a softmax layer so all of its coordinates are non-zero. Consequently,

$$\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}) > 0$$

So by Assumption 5.1, the function $(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) \mapsto \log(\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}))$ is defined and continuous over a compact set, so there exists a constant B giving us the result. □

Lemma E.3 (Population Strong Convexity). *Let $H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) = -\log(\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}))$. Under Assumption 5.2, the function*

$$\mathcal{L}_{\boldsymbol{\theta}^*} : \mathbf{r} \mapsto \mathbb{E} \left[H(\mathbf{r}, \boldsymbol{\theta}^*, \mathbf{x}) \right]$$

is μp_{\min} -strongly convex.

Proof. We first compute the Hessian of \mathcal{L} to find that

$$\nabla^2 \mathcal{L}(\mathbf{r}) = \mathbb{E} \left[\frac{1}{(\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}^*))^2} f(\mathbf{x}, \boldsymbol{\theta}^*) f(\mathbf{x}, \boldsymbol{\theta}^*)^\top \right].$$

Since by Lemma E.1, we have that $\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}^*) \leq p_{\min}^{-1}$, we conclude that

$$\nabla^2 \mathcal{L}(\mathbf{r}) \succeq p_{\min} \mathbb{E} \left[f(\mathbf{x}, \boldsymbol{\theta}^*) f(\mathbf{x}, \boldsymbol{\theta}^*)^\top \right] \succeq \mu p_{\min} \mathbf{I}_m.$$

□

Lemma E.4 (Lipschitz Parametrization). *Let $H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) = -\log(f(\mathbf{x}, \boldsymbol{\theta})^\top \mathbf{r})$. There exists $L > 0$ such that for any $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta$, and $\mathbf{r} \in \mathbb{R}_+^m$, we have that*

$$|H(\mathbf{r}, \boldsymbol{\theta}_1, \mathbf{x}) - H(\mathbf{r}, \boldsymbol{\theta}_2, \mathbf{x})| \leq L \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2.$$

Proof. The gradient of H with respect to $\boldsymbol{\theta}$ is given by

$$\nabla_{\boldsymbol{\theta}} H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) = -\frac{1}{f(\mathbf{x}, \boldsymbol{\theta})^\top \mathbf{r}} \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta})$$

Reasoning like in Lemma E.1, we know that $\frac{1}{f(\mathbf{x}, \boldsymbol{\theta})^\top \mathbf{r}}$ is defined and continuous over the compact set of its parameters, we also know that f is a neural network parametrized by $\boldsymbol{\theta}$, hence $\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta})$ is bounded when $\boldsymbol{\theta}$ and \mathbf{x} are bounded. Consequently, under Assumption 5.1, there exists a constant $L > 0$ such that

$$\|\nabla_{\boldsymbol{\theta}} H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x})\|_2 \leq L.$$

□

Lemma E.5 (Uniform Bound 1). *Let $\delta \in (0, 1)$, with probability $1 - \delta$, we have that*

$$\begin{aligned} & \mathbb{E} \left[H(\hat{\mathbf{r}}_n, \hat{\boldsymbol{\theta}}_t, \mathbf{x}) \right] - \frac{1}{n} \sum_{j=1}^n H(\hat{\mathbf{r}}_n, \hat{\boldsymbol{\theta}}_t, \mathbf{x}_j) \\ & \leq \frac{2}{\sqrt{n}} \text{Rad}(\mathcal{F}) + 2B \sqrt{\frac{\log(4/\delta)}{n}}. \end{aligned} \tag{33}$$

Proof. Let $\delta \in (0, 1)$. Since $\hat{\mathbf{r}}_n$ is learned from the samples \mathbf{x}_j , we do not have independence, which would have allowed us to apply a concentration inequality. Hence, we derive a uniform bound as follows. We begin by observing that:

$$\begin{aligned} & \mathbb{E} \left[H(\hat{\mathbf{r}}_n, \hat{\boldsymbol{\theta}}_t, \mathbf{x}) \right] - \frac{1}{n} \sum_{j=1}^n H(\hat{\mathbf{r}}_n, \hat{\boldsymbol{\theta}}_t, \mathbf{x}_j) \\ & \leq \sup_{\mathbf{r}, \boldsymbol{\theta}} \left(\mathbb{E} \left[H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) \right] - \frac{1}{n} \sum_{j=1}^n H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}_j) \right) \end{aligned}$$

Now since Lemma E.2 holds, we can apply McDiarmid's Inequality to get that with probability $1 - \delta$, we have:

$$\begin{aligned} & \sup_{\mathbf{r}, \boldsymbol{\theta}} \left(\mathbb{E} \left[H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) \right] - \frac{1}{n} \sum_{j=1}^n H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}_j) \right) \\ & \leq \mathbb{E} \left[\sup_{\mathbf{r}, \boldsymbol{\theta}} \left(\mathbb{E} \left[H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}) \right] - \frac{1}{n} \sum_{j=1}^n H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}_j) \right) \right] + 2B \sqrt{\frac{\log(2/\delta)}{n}} \end{aligned}$$

The expectation of the supremum on the right-hand side can be bounded by the Rademacher complexity of $\mathcal{F} := \{\mathbf{x} \mapsto \mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}), (\mathbf{r}, \boldsymbol{\theta}) \in \mathbb{R}_+^m \times \Theta\}$, and we obtain:

$$\begin{aligned} & \sup_{\mathbf{r}, \boldsymbol{\theta}} \left(\mathbb{E}[H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x})] - \frac{1}{n} \sum_{j=1}^n H(\mathbf{r}, \boldsymbol{\theta}, \mathbf{x}_j) \right) \\ & \leq \frac{2}{\sqrt{n}} \text{Rad}(\mathcal{F}) + 2B \sqrt{\frac{\log(2/\delta)}{n}}. \end{aligned} \quad (34)$$

□

Lemma E.6 (Uniform Bound 2). *Let $\delta \in (0, 1)$, with probability $1 - \delta$, we have that*

$$\begin{aligned} & \mathbb{E} \left[H(\mathbf{r}_{f^*}, \hat{\boldsymbol{\theta}}_t, \mathbf{x}) \right] - \frac{1}{n} \sum_{j=1}^n H(\mathbf{r}_{f^*}, \hat{\boldsymbol{\theta}}_t, \mathbf{x}_j) \\ & \leq \frac{2}{\sqrt{n}} \text{Rad}(\mathcal{F}) + 2B \sqrt{\frac{\log(2/\delta)}{n}}. \end{aligned} \quad (35)$$

Proof. The proof is identical to that of Lemma E.5. □

Lemma E.7 (Strong Convexity of Population Loss). *Let $\mathcal{L}(\mathbf{r}, \boldsymbol{\theta})$ be the population loss as defined in Lemma E.7. We establish that $\mathcal{L}(\mathbf{r}, \boldsymbol{\theta})$ is μp_{\min} -strongly convex under the assumptions of calibration (Assumption 5.2).*

Proof. We compute the Hessian of the population loss \mathcal{L} as in Lemma E.7, obtaining that:

$$\nabla^2 \mathcal{L}(\mathbf{r}) = \mathbb{E} \left[\frac{1}{(\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}))^2} f(\mathbf{x}, \boldsymbol{\theta}) f(\mathbf{x}, \boldsymbol{\theta})^\top \right].$$

From Lemma E.1, we have that $\mathbf{r}^\top f(\mathbf{x}, \boldsymbol{\theta}) \leq p_{\min}^{-1}$. Therefore, we conclude:

$$\nabla^2 \mathcal{L}(\mathbf{r}) \succeq p_{\min} \mathbb{E} \left[f(\mathbf{x}, \boldsymbol{\theta}) f(\mathbf{x}, \boldsymbol{\theta})^\top \right] \succeq \mu p_{\min} \mathbf{I}_m.$$

□

Lemma E.8 (Bound on Empirical Loss). *Under Assumption 5.1, the empirical loss $\mathcal{L}_{n^{te}}(\mathbf{r}, \hat{\boldsymbol{\theta}}_{n^{tr}})$ satisfies the following concentration bound:*

$$\mathbb{P} \left(\sup_{\mathbf{r} \in \mathbb{R}_+^m} \left| \mathcal{L}_{n^{te}}(\mathbf{r}, \hat{\boldsymbol{\theta}}_{n^{tr}}) - \mathcal{L}(\mathbf{r}, \hat{\boldsymbol{\theta}}_{n^{tr}}) \right| > \epsilon \right) \leq 2 \exp(-cn^{te}\epsilon^2).$$

Proof. This result follows from standard concentration inequalities, such as McDiarmid's inequality, together with the Lipschitz continuity of the loss function \mathcal{L} with respect to the samples. □

F PROOF OF THEOREM 5.2 AND CONVERGENCE-COMMUNICATION GUARANTEES FOR IW-ERM WITH VRLS

We now establish convergence rates for IW-ERM with VRLS and show our proposed importance weighting achieves *the same rates* with the data-dependent *constant terms* increase linearly with $\max_{y \in \mathcal{Y}} \sup_f r_f(y) = r_{\max}$ under negligible communication overhead over the baseline ERM-solvers without importance weighting. In Appendix F, we establish tight convergence rates and communication guarantees for IW-ERM with VRLS in a broad range of importance optimization settings including convex optimization, second-order differentiability, composite optimization with proximal operator, optimization with adaptive step-sizes, and nonconvex optimization, along the lines of e.g., (Woodworth et al., 2020; Haddadpour et al., 2021; Glasgow et al., 2022; Liu et al., 2023; Hu & Huang, 2023; Wu et al., 2023; Liu et al., 2023).

By estimating the ratios locally and absorbing into local losses, we note that the properties of the modified local loss w.r.t. the neural network parameters \mathbf{w} , e.g., convexity and smoothness, do not change. The data-dependent parameters such as Lipschitz and smoothness constants for $\ell \circ h_{\mathbf{w}}$ w.r.t. \mathbf{w} are scaled linearly by r_{\max} . Our method of importance ratio estimation trains the pre-defined predictor *exclusively using local training data*, which implies IW-ERM with VRLS achieves the same privacy guarantees as the baseline ERM-solvers without importance weighting. For ratio estimation, the communication between clients involves only the estimated marginal label distribution, instead of data, ensuring negligible communication overhead. Given the size of variables to represent marginal distributions, which is by orders of magnitude smaller than the number of parameters of the underlying neural networks for training and the fact that ratio estimation involves only one round of communication, the overall communication overhead for ratio estimation is masked by the communication costs of model training. The communication costs for IW-ERM with VRLS over the course of optimization are exactly the same as those of the baseline ERM-solvers without importance weighting. All in all, importance weighting does not negatively impact communication guarantees throughout the course of optimization, which proves Theorem 5.2.

In the following, we establish tight convergence rates and communication guarantees for IW-ERM with VRLS in a broad range of importance optimization settings including convex optimization, second-order differentiability, composite optimization with proximal operator, optimization with adaptive step-sizes, and nonconvex optimization.

For convex and second-order Differentiable optimization, we establish a lower bound on the convergence rates for IW-ERM in with VRLS and local updating along the lines of e.g., (Glasgow et al., 2022, Theorem 3.1).

Assumption F.1 (PL with Compression). 1) The $\ell(h_{\mathbf{w}}(\mathbf{x}), y)$ is β -smoothness and convex w.r.t. \mathbf{w} for any (\mathbf{x}, y) and satisfies Polyak-Łojasiewicz (PL) condition (there exists $\alpha_{\ell} > 0$ such that, for all $\mathbf{w} \in \mathcal{W}$, we have $\ell(h_{\mathbf{w}}) \leq \|\nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})\|_2^2 / (2\alpha_{\ell})$); 2) The compression scheme \mathcal{Q} is unbiased with bounded variance, i.e., $\mathbb{E}[\mathcal{Q}(\mathbf{x})] = \mathbf{x}$ and $\mathbb{E}[\|\mathcal{Q}(\mathbf{x}) - \mathbf{x}\|_2^2] \leq q\|\mathbf{x}\|_2^2$; 3) The stochastic gradient $\mathbf{g}(\mathbf{w}) = \tilde{\nabla}_{\mathbf{w}} \ell(h_{\mathbf{w}})$ is unbiased, i.e., $\mathbb{E}[\mathbf{g}(\mathbf{w})] = \nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})$ for any $\mathbf{w} \in \mathcal{W}$ with bounded variance $\mathbb{E}[\|\mathbf{g}(\mathbf{w}) - \nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})\|_2^2]$.

For nonconvex optimization with PL condition and communication compression, we establish convergence and communication guarantees for IW-ERM with VRLS, compression, and local updating along the lines of e.g., (Haddadpour et al., 2021, Theorem 5.1).

Theorem F.1 (Convergence and Communication Bounds for Nonconvex Optimization with PL). Let κ denote the condition number, τ denote the number of local steps, R denote the number of communication rounds, and $\max_{y \in \mathcal{Y}} \sup_f r_f(y) = r_{\max}$. Under Assumption F.1, suppose Algorithm 2 with τ local updates and communication compression (Haddadpour et al., 2021, Algorithm 1) is run for $T = \tau R$ total stochastic gradients per node with fixed step-sizes $\eta = 1/(2r_{\max}\beta\gamma\tau(q/K + 1))$ and $\gamma \geq K$. Then we have $\mathbb{E}[\ell(h_{\mathbf{w}_T}) - \ell(h_{\mathbf{w}^*})] \leq \epsilon$ by setting

$$R \lesssim \left(\frac{q}{K} + 1\right) \kappa \log\left(\frac{1}{\epsilon}\right) \quad \text{and} \quad \tau \lesssim \left(\frac{q + 1}{K(q/K + 1)\epsilon}\right). \quad (36)$$

Assumption F.2 (Nonconvex Optimization with Adaptive Step-sizes). 1) The $\ell \circ h_{\mathbf{w}}$ is β -smoothness with bounded gradients; 2) The stochastic gradients $\mathbf{g}(\mathbf{w}) = \tilde{\nabla}_{\mathbf{w}} \ell(h_{\mathbf{w}})$ is unbiased with bounded variance $\mathbb{E}[\|\mathbf{g}(\mathbf{w}) - \nabla_{\mathbf{w}} \ell(h_{\mathbf{w}})\|_2^2]$; 3) Adaptive matrices A_t constructed as in (Wu et al., 2023, Algorithm 2) are diagonal and the minimum eigenvalues satisfy $\lambda_{\min}(A_t) \geq \rho > 0$ for some $\rho \in \mathbb{R}_+$.

For nonconvex optimization with adaptive step-sizes, we establish convergence and communication guarantees for IW-ERM with VRLS and local updating along the lines of e.g., (Wu et al., 2023, Theorem 2).

Theorem F.2 (Convergence and Communication Guarantees for Nonconvex Optimization with Adaptive Step-sizes). *Let τ denote the number of local steps, R denote the number of communication rounds, and $\max_{y \in \mathcal{Y}} \sup_f r_f(y) = r_{\max}$. Under Assumption F.2, suppose Algorithm 2 with τ local updates is run for $T = \tau R$ total stochastic gradients per node with an adaptive step-size similar to (Wu et al., 2023, Algorithm 2). Then we $\mathbb{E}[\|\nabla_{\mathbf{w}} \ell(h_{\mathbf{w}_T})\|_2] \leq \epsilon$ by setting:*

$$T \lesssim \frac{r_{\max}}{K\epsilon^3} \quad \text{and} \quad R \lesssim \frac{r_{\max}}{\epsilon^2}. \quad (37)$$

Assumption F.3 (Composite Optimization with Proximal Operator). *1) The $\ell \circ h_{\mathbf{w}}$ is smooth and strongly convex with condition number κ ; 2) The stochastic gradients $\mathbf{g}(\mathbf{w}) = \tilde{\nabla}_{\mathbf{w}} \ell(h_{\mathbf{w}})$ is unbiased.*

For composite optimization with strongly convex and smooth functions and proximal operator, we establish an upper bound on oracle complexity to achieve ϵ error on the Lyapunov function defined as in (Hu & Huang, 2023, Section 4) for Gradient Flow-type transformation of IW-ERM with VRLS in the limit of infinitesimal step-size.

Theorem F.3 (Oracle Complexity of Proximal Operator for Composite Optimization). *Let κ denote the condition number. Under Assumption F.3, suppose Gradient Flow-type transformation of Algorithm 2 with VRLS and Proximal Operator evolves in the limit of infinitesimal step-size (Hu & Huang, 2023, Algorithm 3). Then it achieves $\mathcal{O}(r_{\max} \sqrt{\kappa} \log(1/\epsilon))$ Proximal Operator Complexity.*

G COMPLEXITY ANALYSIS

In our algorithm, the ratio estimation is performed once in parallel before the IW-ERM step.

In the experiments, we used a simple network to estimate the ratios in advance, which required significantly less computational effort compared to training the global model. Although IW-ERM with VRLS introduces additional computational complexity compared to the baseline FedAvg, it results in substantial improvements in overall generalization, particularly under challenging label shift conditions.

H MATHEMATICAL NOTATIONS

In this appendix, we provide a summary of mathematical notations used in this paper in Table 5:

Table 5: Math Symbols

Math Symbol	Definition
\mathcal{X}	Compact metric space for features
\mathcal{Y}	Discrete label space with $ \mathcal{Y} = m$
K	Number of clients in an FL setting
\mathcal{S}_k	All samples in the training set of client k
h_w	Hypothesis function $h_w : \mathcal{X} \rightarrow \mathcal{Y}$
\mathcal{H}	Hypothesis class for h_w
\mathcal{Z}	Mapping space from \mathcal{X} , which can be discrete or continuous

I LIMITATIONS

The distribution shifts observed in real-world data are often not fully captured by the label shift or relaxed distribution shift assumptions. In our experiments, we applied mild test data augmentation to approximate the relaxed label shift and manage ratio estimation errors for both the baselines and our method. However, the label shift assumption remains overly restrictive, and the relaxed label shift lacks robust empirical validation in practical scenarios.

Additionally, IW-ERM’s parameter estimation relies on local predictors at each client, which limits its scalability. In practice, a simpler global predictor could be sufficient for parameter estimation and IW-ERM training. Future research could explore VRLS variants capable of effectively handling more complex distribution shifts in challenging datasets, such as CIFAR-10.1 (Recht et al., 2018; Torralba et al., 2008), as suggested in (Garg et al., 2023).

J EXPERIMENTAL DETAILS AND ADDITIONAL EXPERIMENTS

In this section, we provide experimental details and additional experiments. In particular, we validate our theory on multiple clients in a federated setting and show that our IW-ERM outperforms FedAvg and FedBN baselines *under drastic and challenging label shifts*.

J.1 EXPERIMENTAL DETAILS

In single-client experiments, a simple MLP without dropout is used as the predictor for MNIST, and ResNet-18 for CIFAR-10.

For experiments in a federated learning setting, both MNIST (LeCun et al., 1998) and Fashion MNIST (Xiao et al., 2017) datasets are employed, each containing 60,000 training samples and 10,000 test samples, with each sample being a 28 by 28 pixel grayscale image. The CIFAR-10 dataset (Krizhevsky) comprises 60,000 colored images, sized 32 by 32 pixels, spread across 10 classes with 6,000 images per class; it is divided into 50,000 training images and 10,000 test images. In this setting, the objective is to minimize the cross-entropy loss. Stochastic gradients for each client are calculated with a batch size of 64 and aggregated on the server using the Adam optimizer. LeNet is used for experiments on MNIST and Fashion MNIST with a learning rate of 0.001 and a weight decay of 1×10^{-6} . For CIFAR-10, ResNet-18 is employed with a learning rate of 0.0001 and a weight decay of 0.0001. Three independent runs are implemented for 5-client experiments on Fashion MNIST and CIFAR-10, while for 10 clients, one run is conducted on CIFAR-10. The regularization coefficient ζ in Equation (4) is set to 1 for all experiments. All experiments are performed using a single GPU on an internal cluster and Colab.

Importantly, the training of the predictor for ratio estimation on both the baseline MLLS and our VRLS is executed with identical hyperparameters and epochs for CIFAR-10 and Fashion MNIST. The training is halted once the classification loss reaches a predefined threshold on MNIST.

J.2 RELAXED LABEL SHIFT EXPERIMENTS

In conventional label shift, it is assumed that $p(\mathbf{x} | y)$ remains unchanged across training and test data. However, this assumption is often too strong for real-world applications, such as in healthcare, where different hospitals may use varying equipment, leading to shifts in $p(\mathbf{x} | y)$ even with the same labels (Rajendran et al., 2023). Relaxed label shift loosens this assumption by allowing small changes in the conditional distribution (Garg et al., 2023; Luo & Ren, 2022).

To formalize this, we use the distributional distance \mathcal{D} and a relaxation parameter $\epsilon > 0$, as defined by Garg et al. (2023): $\max_y \mathcal{D}(p_{\text{tr}}(\mathbf{x} | y), p_{\text{te}}(\mathbf{x} | y)) \leq \epsilon$. This allows for slight differences in feature distributions between training and testing, capturing a more realistic scenario where the conditional distribution is not strictly invariant.

In our case, visual inspection suggests that the differences between temporally distinct datasets, such as CIFAR-10 and CIFAR-10.1_v6 (Torralba et al., 2008; Recht et al., 2018), may not meet the assumption of a small ϵ . To address this, we instead simulate controlled shifts using test data augmentation, allowing us to regulate the degree of relaxation, following the approach outlined in Garg et al. (2023).

J.3 ADDITIONAL EXPERIMENTS

In this section, we provide supplementary results, visualizations of accuracy across clients and tables showing dataset distribution in FL setting and relaxed label shift.

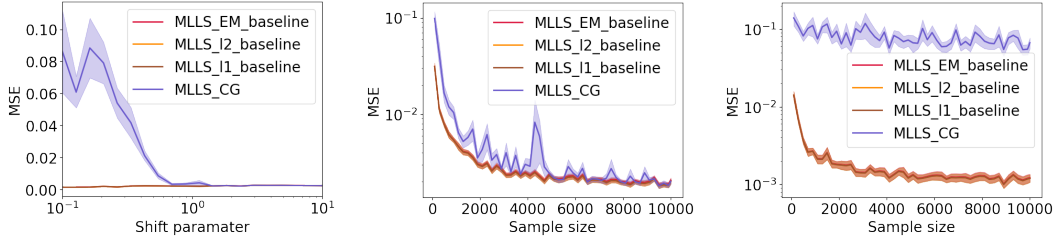


Figure 3: MSE analysis on MNIST for MLLS baselines. **Left:** Performance evaluation across various alpha values, comparing different methods: MLLS_EM, MLLS_L1, MLLS_L2, and MLLS_CG. MLLS_L1 and MLLS_L2 utilize convex optimization with L_1 and L_2 regularization for estimating our limited test sample problem, respectively, and are solved directly with a convex solver. In contrast, MLLS_CG uses conjugate gradient descent and MLLS_EM solves this convex optimization problem with EM algorithm. Both the EM and convex optimization methods (MLLS_L1, MLLS_L2) demonstrate superior and more consistent performance, especially under severe label shift conditions, when compared to MLLS_CG. **Middle:** At an alpha value of 1.0, the MSE analysis shows comparable performance across most methods, with the exception of MLLS_CG, which lags behind. **Right:** For alpha=0.1, MLLS_CG performs significantly worse than the EM and convex optimization methods, consistent with the trends observed in the left plot.

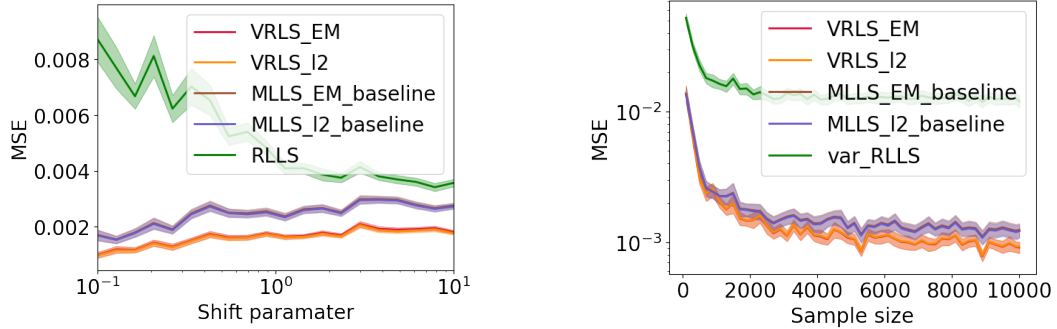


Figure 4: In our detailed analysis with the MNIST dataset, we conduct a thorough comparison of VRLS alongside MLLS (Garg et al., 2020), EM (Saerens et al., 2002), and also RLLS (Azizzadenesheli et al., 2019).

Table 6: LeNet on Fashion MNIST with label shift across 5 clients. 15,000 iterations for FedAvg and FedBN; 5,000 for Upper Bound (FTW-ERM) using true ratios and our IW-ERM. To mention, to train our predictor, we use a simplest MLP and employ linear kernel.

FMNIST	Our IW-ERM	FedAvg	FedBN	Upper Bound
Avg. accuracy	0.7520 ± 0.0209	0.5472 ± 0.0297	0.5359 ± 0.0306	0.8273 ± 0.0041
Client 1 accuracy	0.7162 ± 0.0059	0.3616 ± 0.0527	0.3261 ± 0.0296	0.8590 ± 0.0062
Client 2 accuracy	0.9266 ± 0.0125	0.9060 ± 0.0157	0.9035 ± 0.0162	0.9357 ± 0.0037
Client 3 accuracy	0.6724 ± 0.0467	0.3279 ± 0.0353	0.3612 ± 0.0814	0.7896 ± 0.0109
Client 4 accuracy	0.7979 ± 0.0448	0.6858 ± 0.0105	0.6654 ± 0.0121	0.8098 ± 0.0112
Client 5 accuracy	0.6468 ± 0.0248	0.4548 ± 0.0655	0.4234 ± 0.0387	0.7426 ± 0.0257

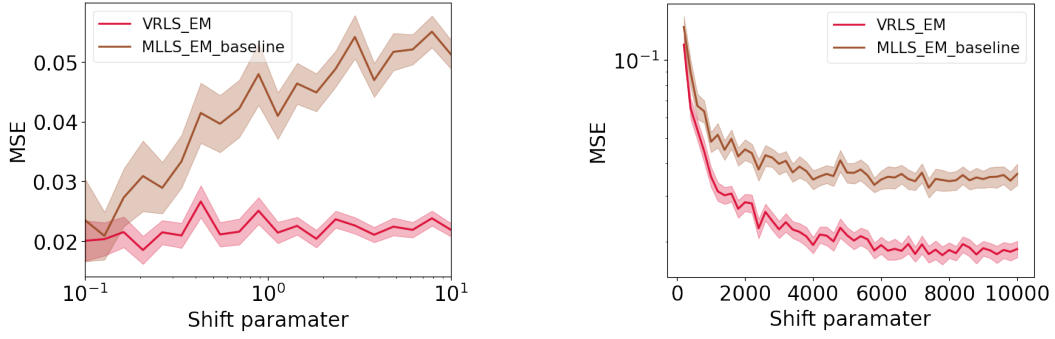


Figure 5: In this experiment with Fashion MNIST, a simple MLP with dropout were employed.

Table 7: ResNet-18 on CIFAR-10 with label shift across 5 clients. For fair comparison, we run 5,000 iterations for our method and Upper Bound, while 10000 for FedAvg and FedBN.

CIFAR-10	Our IW-ERM	FedAvg	FedBN	Upper Bound
Avg. accuracy	0.5640 \pm 0.0241	0.4515 \pm 0.0148	0.4263 \pm 0.0975	0.5790 \pm 0.0103
Client 1 accuracy	0.6410 \pm 0.0924	0.5405 \pm 0.1845	0.5321 \pm 0.0620	0.7462 \pm 0.0339
Client 2 accuracy	0.8434 \pm 0.0359	0.3753 \pm 0.0828	0.4656 \pm 0.2158	0.7509 \pm 0.0534
Client 3 accuracy	0.4591 \pm 0.1131	0.3973 \pm 0.1333	0.2838 \pm 0.1055	0.5845 \pm 0.0854
Client 4 accuracy	0.4751 \pm 0.1241	0.5007 \pm 0.1303	0.5256 \pm 0.1932	0.3507 \pm 0.0578
Client 5 accuracy	0.4013 \pm 0.0430	0.4429 \pm 0.1195	0.5603 \pm 0.1581	0.4627 \pm 0.0456

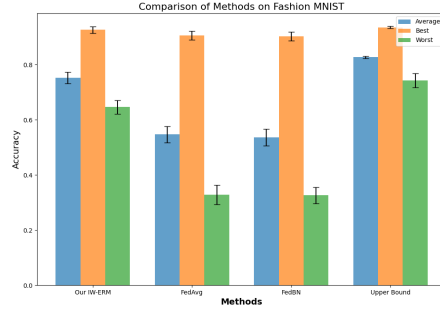


Figure 6: The average, best-client, and worst-client accuracy, along with their standard deviations, are derived from Table 6. Our method exhibits the lowest standard deviation, showcasing the most robust accuracy amongst the compared methods.

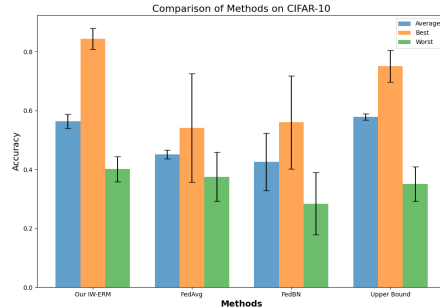


Figure 7: The average, best-client, and worst-client accuracy, along with their standard deviations, are derived from Table 7.

Table 8: Label distribution on Fasion MNIST with 5 nodes, with the majority of classes possessing a limited number of training and test images across each node.

		Class									
		0	1	2	3	4	5	6	7	8	9
Node 1	Train	34	34	34	34	34	5862	34	34	34	34
	Test	977	5	5	5	5	5	5	5	5	5
Node 2	Train	34	34	34	34	34	34	5862	34	34	34
	Test	5	977	5	5	5	5	5	5	5	5
Node 3	Train	34	34	34	34	34	34	34	5862	34	34
	Test	5	5	977	5	5	5	5	5	5	5
Node 4	Train	34	34	34	34	34	34	34	34	5862	34
	Test	5	5	5	977	5	5	5	5	5	5
Node 5	Train	34	34	34	34	34	34	34	34	34	5862
	Test	5	5	5	5	977	5	5	5	5	5

Table 9: Label distribution on CIFAR-10 with 5 clients, with the majority of classes possessing a limited number of training and test images across each client.

		Class									
		0	1	2	3	4	5	6	7	8	9
Node 1	Train	34	34	34	34	34	5862	34	34	34	34
	Test	977	5	5	5	5	5	5	5	5	5
Node 2	Train	34	34	34	34	34	34	5862	34	34	34
	Test	5	977	5	5	5	5	5	5	5	5
Node 3	Train	34	34	34	34	34	34	34	5862	34	34
	Test	5	5	977	5	5	5	5	5	5	5
Node 4	Train	34	34	34	34	34	34	34	34	5862	34
	Test	5	5	5	977	5	5	5	5	5	5
Node 5	Train	34	34	34	34	34	34	34	34	34	5862
	Test	5	5	5	5	977	5	5	5	5	5

Table 10: Label distribution on CIFAR-10 with 100 clients, wherein groups of 10 clients share the same distribution and ratios. The majority of classes possess a limited quantity of training and test images on each client.

		Class				
		0	1	2	3	4
Client 1-10	Train	95/100	5/9	5/9	5/9	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 11-20	Train	5/9	95/100	5/9	5/9	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 21-30	Train	5/9	5/9	95/100	5/9	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 31-40	Train	5/9	5/9	5/9	95/100	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 41-50	Train	5/9	5/9	5/9	5/9	95/100
	Test	5/9	5/9	5/9	5/9	5/9
Client 51-60	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	5/9	5/9	5/9	95/100
Client 61-70	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	5/9	5/9	95/100	5/9
Client 71-80	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	5/9	95/100	5/9	5/9
Client 81-90	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	95/100	5/9	5/9	5/9
Client 91-100	Train	5/9	5/9	5/9	5/9	5/9
	Test	95/100	5/9	5/9	5/9	5/9

		Class				
		5	6	7	8	9
Client 1-10	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	5/9	5/9	5/9	95/100
Client 11-20	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	5/9	5/9	95/100	5/9
Client 21-30	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	5/9	95/100	5/9	5/9
Client 31-40	Train	5/9	5/9	5/9	5/9	5/9
	Test	5/9	95/100	5/9	5/9	5/9
Client 41-50	Train	5/9	5/9	5/9	5/9	5/9
	Test	95/100	5/9	5/9	5/9	5/9
Client 51-60	Train	95/100	5/9	5/9	5/9	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 61-70	Train	5/9	95/100	5/9	5/9	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 71-80	Train	5/9	5/9	95/100	5/9	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 81-90	Train	5/9	5/9	5/9	95/100	5/9
	Test	5/9	5/9	5/9	5/9	5/9
Client 91-100	Train	5/9	5/9	5/9	5/9	95/100
	Test	5/9	5/9	5/9	5/9	5/9