



## Motivation

### Label Shift

In real-world applications, the label distribution  $p(\mathbf{y})$  may change from training and test, while the conditional distribution  $p(\mathbf{x}|\mathbf{y})$  remains unchanged.

### Challenge in Distributed Learning

In multi-node setups, each node may face both intra-node and inter-node label shifts. Standard ERM is insufficient to address this distributional heterogeneity.

### Goal

- Estimate local label shift ratios from few unlabeled test samples, with statistical guarantees.
- Train a global model robust to label shift, in a privacy-preserving distributed framework.

### Key Idea

- VRLS:** Entropy-regularized predictor and density ratio estimation.
- IW-ERM:** Weighted ERM using local ratio estimates to optimize true test risk.

## Problem Setup under Single-Node Case

### Importance Ratio under Label Shift

Given training data  $(\mathbf{x}_i, \mathbf{y}_i) \sim p^{\text{tr}}$  and test data  $(\mathbf{x}_i, \mathbf{y}_i) \sim p^{\text{te}}$ , define the importance ratio as:

$$r(\mathbf{y}) = \frac{p^{\text{te}}(\mathbf{y})}{p^{\text{tr}}(\mathbf{y})}.$$

### Entropy Regularization

Train  $f_{\theta}(\mathbf{x})$  to approximate  $p^{\text{tr}}(\mathbf{y}|\mathbf{x})$  using cross-entropy loss with entropy-based regularization:

$$\Omega(f_{\theta}) = \sum_{c=1}^m \text{softmax}[f_{\theta}(\mathbf{x})]_c \log(\text{softmax}[f_{\theta}(\mathbf{x})]_c).$$

### Why Entropy?

Helps mitigate overconfidence, improving calibration for ratio estimation.

## VRLS Algorithm

**Input:** Training set  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{n_{\text{tr}}}$  and unlabeled test set  $\{\mathbf{x}_j\}_{j=1}^{n_{\text{te}}}$ .

**Train predictor:**

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \left[ \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \ell_{\text{CE}}(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i) + \zeta \Omega(f_{\theta}) \right]$$

**Estimate importance ratio:**

$$\hat{\mathbf{r}} = \arg \max_{\mathbf{r} \in \mathbb{R}_{+}^m} \frac{1}{n_{\text{te}}} \sum_{j=1}^{n_{\text{te}}} \log [f_{\hat{\theta}}(\mathbf{x}_j)^{\top} \mathbf{r}]$$

**Apply  $\hat{\mathbf{r}}$ :** Distributed training with importance weighting.

## Multi-Node Extension

**Setup:** For a system of  $K$  nodes, each node  $k$  has training distribution  $p_k^{\text{tr}}$  and test distribution  $p_k^{\text{te}}$ . A global model  $h_{\mathbf{w}} : \mathcal{X} \rightarrow \mathcal{Y}$  is trained across all nodes. Under label shift:

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

**Local Risk:** Each node  $k$  aims to minimize its true expected test risk:

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

## IW-ERM with VRLS in Distributed Learning

### Phase 1: Local Ratio Estimation (VRLS)

- Each node  $k \in [K]$  estimates its local label shift ratio  $\hat{\mathbf{r}}_k$  by performing the following steps:
  - Train a local predictor  $\hat{f}_{k, \theta}$  with entropy regularization.
  - Estimate the importance ratio  $\hat{\mathbf{r}}_k$  by optimizing the VRLS objective on local test data.

### Phase 2: Ratio Aggregation

- Each node transmits its local label shift estimate  $\hat{\mathbf{r}}_k$  to a central aggregator. The final estimate is:

$$\hat{\mathbf{r}}_k(\mathbf{y}) = \frac{1}{\hat{p}_k^{\text{tr}}(\mathbf{y})} \sum_{j=1}^K \hat{p}_j^{\text{te}}(\mathbf{y}).$$

### Phase 3: Global IW-ERM Training

- Train global model  $h_{\mathbf{w}}$  with reweighted loss that incorporates the estimated label shift ratios:

$$\min_{\mathbf{w}} \sum_{k=1}^K \frac{1}{n_k^{\text{tr}}} \sum_{i=1}^{n_k^{\text{tr}}} \hat{\mathbf{r}}_k(\mathbf{y}_{k,i}) \ell(h_{\mathbf{w}}(\mathbf{x}_{k,i}), \mathbf{y}_{k,i}).$$

**Output:** A global model  $h_{\mathbf{w}}$  adapted to heterogeneous label shifts across distributed nodes.

## Theoretical Highlights

### Ratio Estimation:

Under mild assumptions, the VRLS estimator  $\hat{\mathbf{r}}$  converges to the true optimum  $\mathbf{r}_f^*$  with high probability:

$$\|\hat{\mathbf{r}} - \mathbf{r}_f^*\| \leq O\left(\frac{1}{\sqrt{n_{\text{te}}}}\right) + \text{calib. error}(\theta, \theta^*),$$

due to the MLE structure and bounded softmax outputs.

### Distributed Convergence:

Let  $h(T)$  denote the baseline convergence rate (e.g.,  $O(\frac{1}{\sqrt{T}})$  for SGD). IW-ERM with VRLS satisfies:

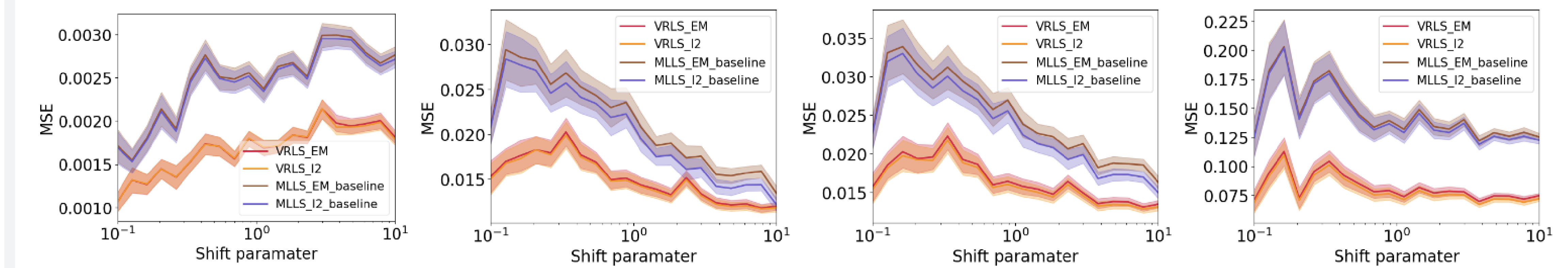
$$\ell(h_{\mathbf{w}}) - \ell^* \leq O(R_{\text{max}}) \cdot h(T).$$

Communication overhead is minimal: each node transmits only a one-shot  $\hat{\mathbf{r}}_k$ , preserving privacy and incurring no extra per-iteration cost.

## Experiments and Results

### Single-Node Ratio Estimation:

- Datasets: MNIST and CIFAR-10, with synthetic label shift severity  $\alpha$ .
- VRLS achieves lower MSE compared to EM-based baseline.



MNIST

CIFAR-10

CIFAR-10, Relaxed

CIFAR-10, Relax-m

### Multi-Node Setup:

- Up to 200 nodes, each with local distributions  $p_k^{\text{tr}}$  and  $p_k^{\text{te}}$ .
- Our IW-ERM with VRLS outperforms FedAvg, FedProx, etc.
- Each node transmits its local  $\hat{\mathbf{r}}_k$  once to the server; no raw data is shared.

### Accuracy across 100 Nodes on Fashion-MNIST:

Method	Accuracy
IW-ERM	<b>0.7520 ± 0.0209</b>
IW-ERM (small)	0.7376 ± 0.0099
FedAvg	0.5472 ± 0.0297
FedBN	0.5359 ± 0.0306
FedProx	0.5606 ± 0.0070
SCAFFOLD	0.5774 ± 0.0036
Upper Bound	0.8273 ± 0.0041

### Accuracy across varying node numbers on CIFAR-10:

Nodes	IW-ERM	FedAvg	FedBN
100	<b>0.5354</b>	0.3915	0.1537
200	<b>0.6216</b>	0.5942	0.1753

## Future Directions

- Beyond  $p(\mathbf{x}|\mathbf{y})$  invariance: handling mixed or relaxed shift conditions.
- Improving ratio estimation: tighter generalization bounds and faster statistical convergence.
- Scaling VRLS to large-scale and privacy-preserving distributed systems.

## Acknowledgments

Experiments were done on UiO Sigma2 resources. This work was supported by the Research Council of Norway through Integreat, Norwegian Centre for Knowledge-driven Machine Learning (332645) and Visual Intelligence, Norwegian Centre for Research-based Innovation (309439). The work of Volkan Cevher was supported by the Hasler Foundation Program: Hasler Responsible AI (21043), the Army Research Office (W911NF-24-1-0048), and the Swiss National Science Foundation (200021-205011).