

Distributionally Robust Post-hoc Classifiers Under Prior Shifts

Jiaheng Wei💜💙, Harikrishna Narasimhan💙, Ehsan Amid💙,
Wen-sheng Chu💙, Yang Liu💜, Abhishek Kumar💙

💙: Google Research

💜: University of California. Santa Cruz

ICLR, 2023

Google Research



UNIVERSITY OF CALIFORNIA
SANTA CRUZ



Brief Overview

Setting

Improve the distribution robustness when there exist prior shifts between training and test datasets, i.e.,

- Shifts in Class Prior $\mathbb{P}(Y = i)$, with label Y ;
- Shifts in Group Prior $\mathbb{P}(G = i)$, with hidden attribute G ;



Brief Overview

One sentence summary

A post-hoc approach that performs scaling adjustments to predictions from a pre-trained model, via minimizing a distributionally robust loss around a target distribution.

Background

Model Robustness under
Distribution Shifts

Distribution Shifts in Class Priors

AKA: Class-Imbalanced Learning

- **Basic Setting:**

For an m -class classification task, denote $\{(x_j, y_j)\}_{j=1}^n$ the training samples drawn from $(X, Y) \sim \mathcal{D}$.

Assume $\pi_i := \mathbb{P}(Y = i)$ is the class prior.

- **Class-level distribution robustness:**

Train on imbalanced π



Aim to perform well on a *target prior distribution* at the test time
(Test on target π)

Distribution Shifts in Group Priors

AKA: Group Distributional Robustness

- **Basic Setting:**

For a m -class classification task, denote $\{(x_j, y_j, a_j)\}_{j=1}^n$ the training samples drawn from $(X, Y, A) \sim \mathcal{D}$.

a_j is the attribute/group information, i.e., male/female.

- **Group-level distribution robustness:**

Train on imbalanced group priors
(without group information in training)

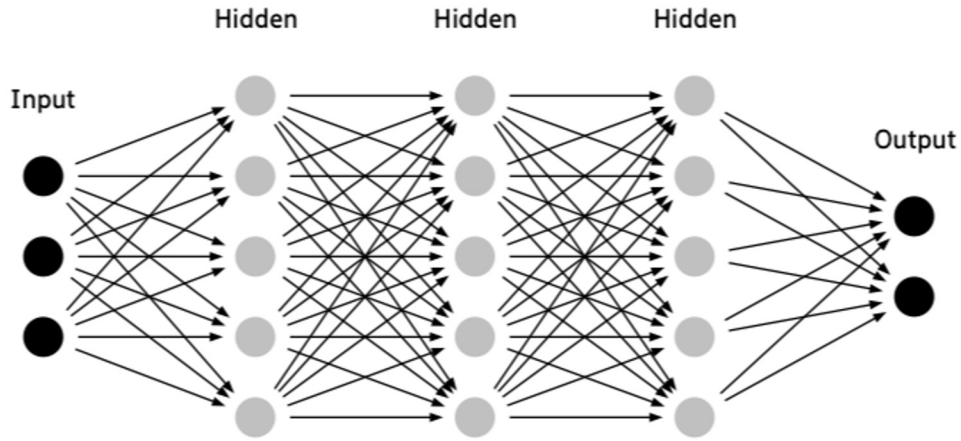


Aim to perform well on a *target group prior* during the test time

Motivations

Deep neural nets
capture core features well

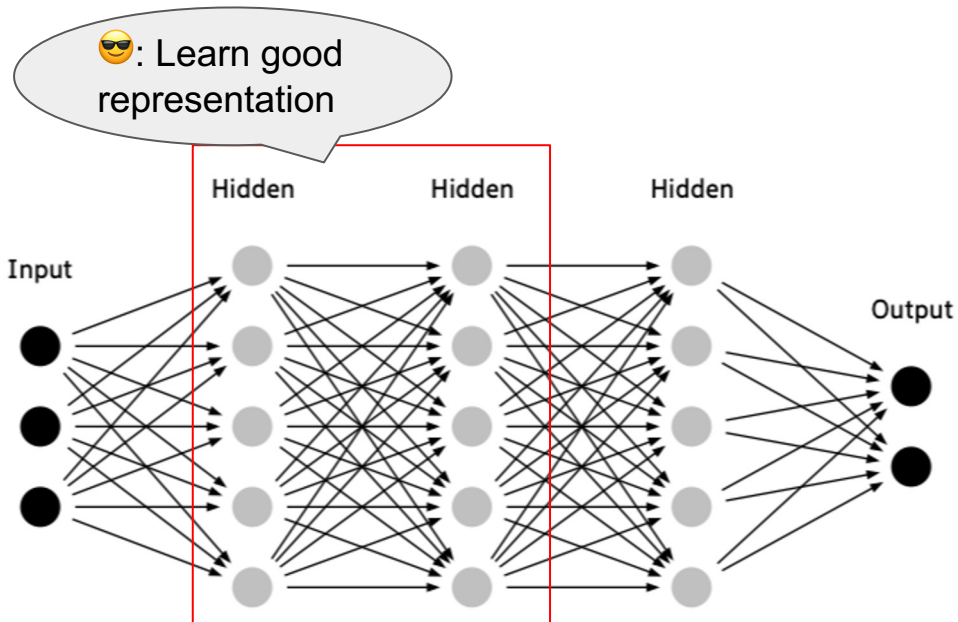
Motivation 1



Deep neural networks could learn the core features sufficiently well, even if they appear to perform poorly on minority classes/groups.

[1] Last layer re-training is sufficient for robustness to spurious correlations.
[ICML 2022 workshop]

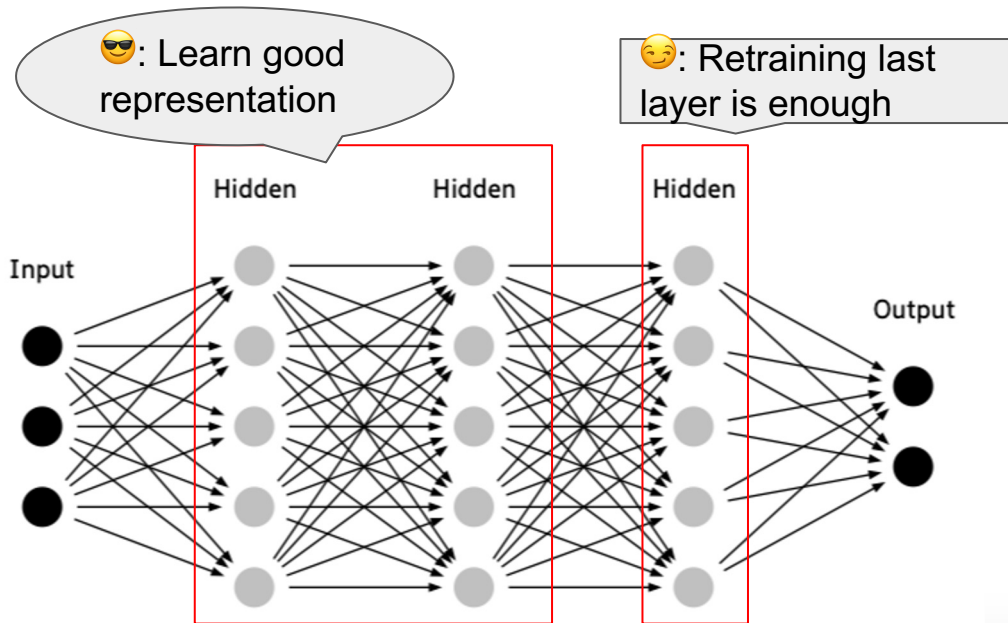
Motivation 1



Deep neural networks could learn the core features sufficiently well, even if they appear to perform poorly on minority classes/groups.

[1] Last layer re-training is sufficient for robustness to spurious correlations.
[ICML 2022 workshop]

Motivation 1



Deep neural networks could learn the core features sufficiently well, even if they appear to perform poorly on minority classes/groups.

[1] Last layer re-training is sufficient for robustness to spurious correlations.
[ICML 2022 workshop]

A spectrum of controlled distribution shifts:

Distribution Robust Evaluation (DRE metric)

$$\text{DRE}(D, \delta): \quad \min_g \sum_i g_i \text{Acc}_i, \quad \text{s.t.} \quad \sum_{i \in [m]} g_i = 1, g_i \geq 0, D(g, u) \leq \delta.$$

Minimize the weighted sum of per-class/group accuracy Acc_i , where:

D : the divergence; u : a *target distribution*; δ : the perturbation;

Special cases

$\delta = 0$ evaluates w.r.t. any target distribution u ;

$\delta = \infty$ returns the worst class/group accuracy.

Summary: DRE metric measures the worst expected accuracy in a δ -radius ball around the target distribution u .

Goal:

😎: Improve the
distributional robustness

$$\begin{aligned} \min_{\theta} \quad & \max_g \sum_{i \in [m]} g_i \mathbb{P}_{(X,Y=i) \sim \mathcal{D}}(h_{\theta}(X) \neq Y), \\ \text{s. t.} \quad & \sum_{i \in [m]} g_i = 1, g_i \geq 0, D(g, u) \leq \delta. \end{aligned}$$

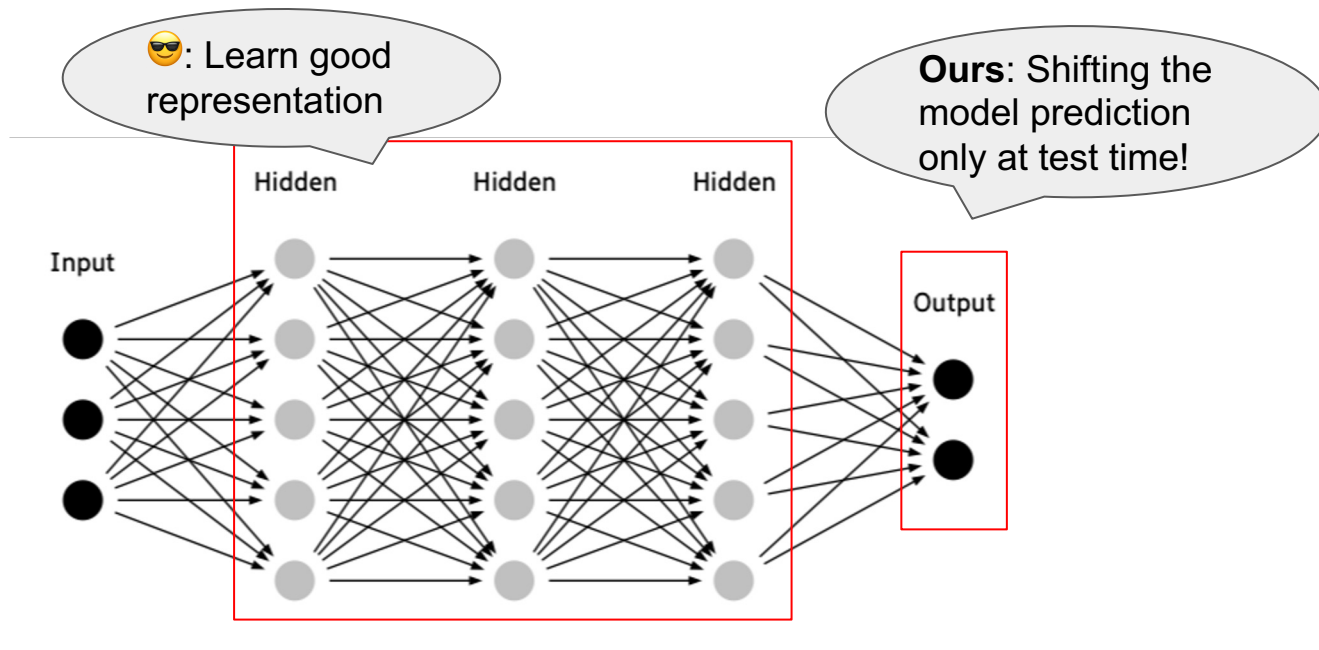
h_{θ} denotes the deep neural nets.

🤔: Can we optimize the model performance under the controlled distribution shifts, by **only** shifting the model prediction at the **test time**?

Distributional ROburst PoSt-hoc Approach

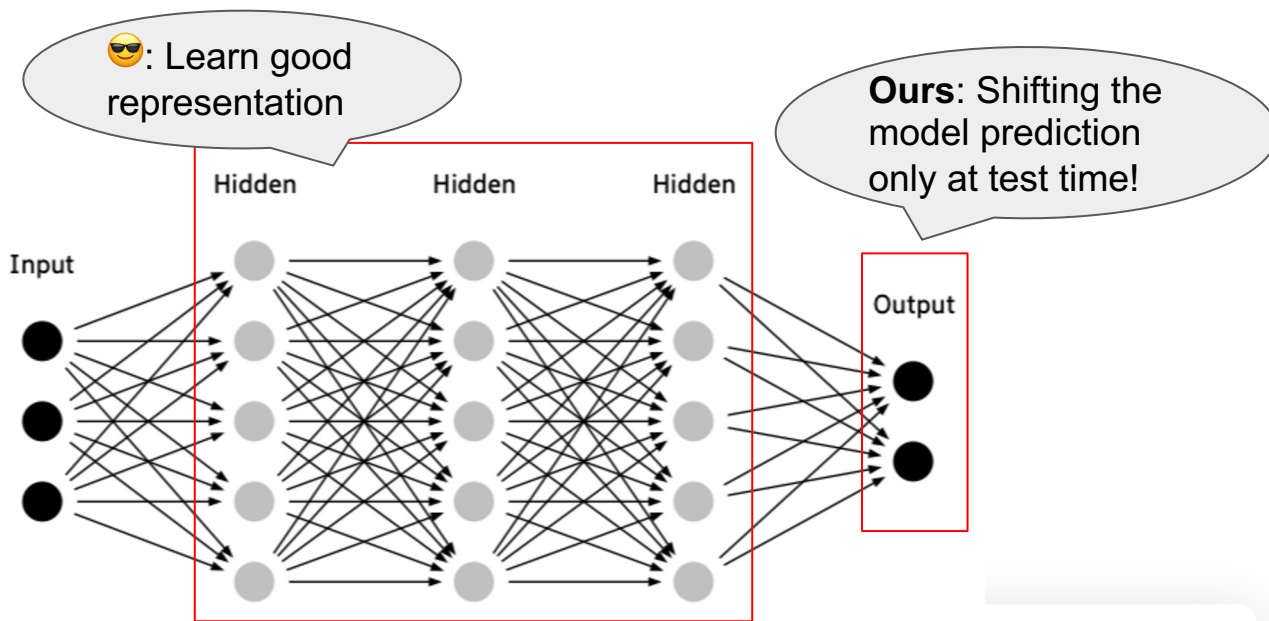
Scale the model prediction
at test time (DROPS)

DROPS



👍: Test time scaling helps with improving the robustness under any **DRE** metric, efficiently.

DROPS



The intuition

1. **Down-scale:**
The model predicted probability of majority training classes at the test time.
1. **Up-scale:**
The model predicted probability of minority training classes at the test time.

How to scale model predictions

Class-imbalanced setting

$$\min_{f: \mathcal{X} \rightarrow \Delta_m} \max_{g \in \mathbb{G}(\delta)} \sum_{i=1}^m g_i \ell_i(f). \quad (1)$$

Minimize the sum of worst-case re-weighted per-class loss

Clarifications

$\ell_i(f)$ is the expected loss for each class,
 $\mathbb{G}(\delta) = \{g \in \Delta_m \mid D(g, u) \leq \delta\}$, where:
 $D: \Delta_m \times \Delta_m \rightarrow \mathbb{R}_+$ is the divergence;
 u is a *target distribution*, δ is the perturbation.

Generates optimal class-weights g_i^*

How to scale model predictions

Class-imbalanced setting

$$\mathbf{DROPS:} \quad h(x) \in \operatorname{argmax}_{i \in [m]} \frac{g_i^*}{\pi_i} \cdot \hat{\eta}_i(x). \quad (2)$$

Clarifications

$\hat{\eta}_i(x)$ is the predicted probability of sample x belonging to class i ;

Insights

Model prediction of class i is upscaled by **DROPS** if:

(1) A large weight g_i^* is assigned; or (2) Class i has a small prior π_i .

An Empirical Sketch (DROPS)

Solve the Lagrangian form of Eqn. (1) via a validation set for a number of iterations.

At iteration t , do:

- **Step 1:** updating the Lagrangian multiplier $\lambda^{(t)}$;
- **Step 2:** updating the weights $g^{(t)}$;
- **Step 3:** scaling the predictions.

Experiments

On class-imbalanced learning
& Group distributional robustness

Experiments on Class-Imbalanced Learning

Synthetic class-imbalanced CIFAR-10 & CIFAR-100

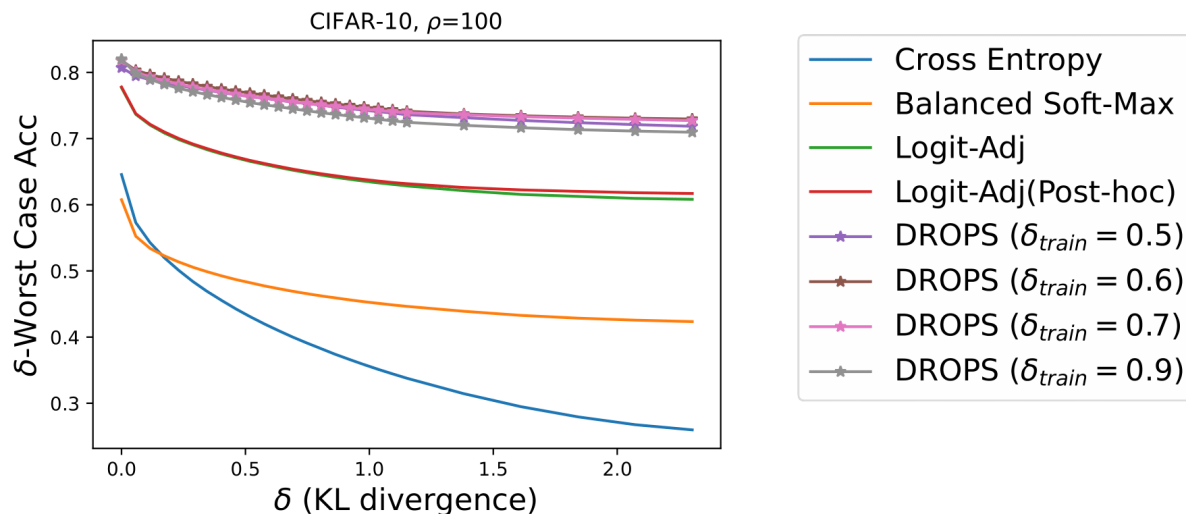
- Simulation: with the increasing of class index, # of selected samples per class has an exponential decay.
Imbalance ratio: $\rho = \frac{\max_i \pi_i}{\min_i \pi_i}$.
- Train on imbalanced dataset; validate and test on balanced datasets.
- For DROPS, get g_i^* under different perturbation level δ_{train} , under D_{KL} .
- Performance evaluation: for both D_{KL} and $D_{\text{R-KL}}$, report the model performance in the metric $\text{DRE}(D, \delta_{\text{test}})$, with a list of δ_{test} .

Experiments on Class-Imbalanced Learning

CIFAR-10 (imbalance ratio $\rho = 100$) – under $DRE(D_{KL}, \delta_{test})$

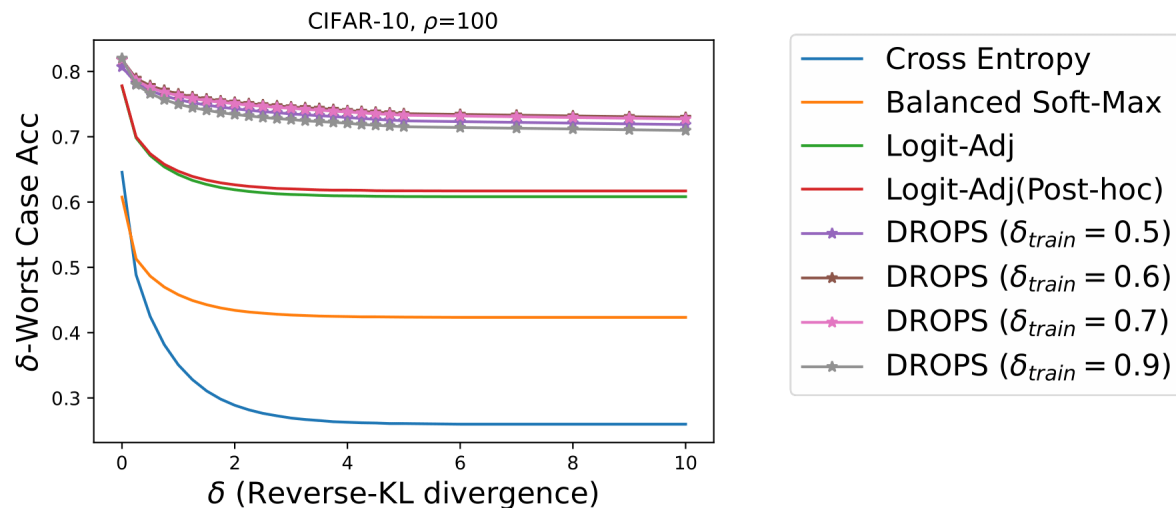
X axis: divergence from the target distribution δ_{test} ;

Y-axis: $DRE(D, \delta_{test})$; **Higher curve** \rightarrow **More robustness**



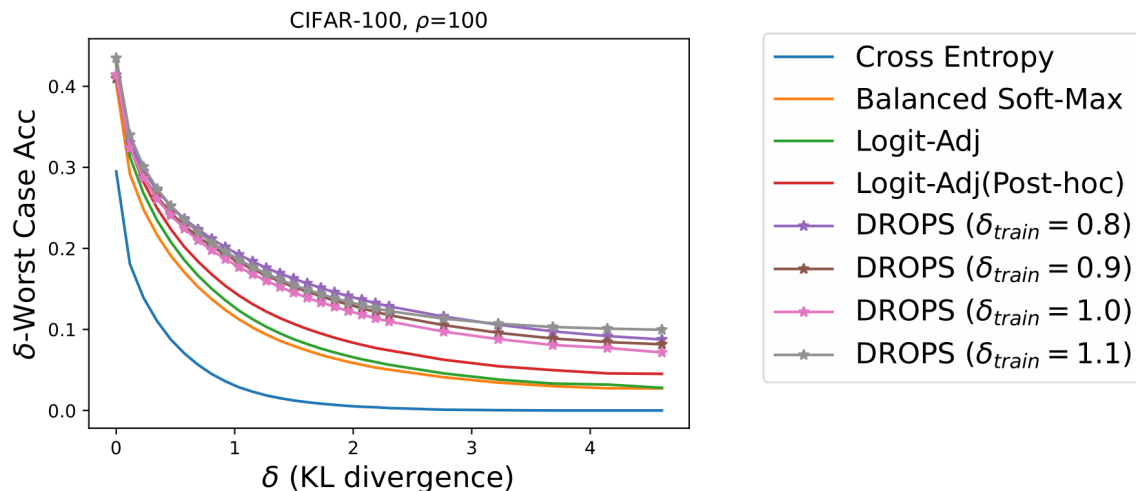
Experiments on Class-Imbalanced Learning

CIFAR-10 (imbalance ratio $\rho = 100$) – under DRE($D_{R-KL}, \delta_{\text{test}}$)



Experiments on Class-Imbalanced Learning

CIFAR-100 (imbalance ratio $\rho = 100$) – under $\text{DRE}(D_{\text{KL}}, \delta_{\text{test}})$



Experiments on Group Distributional Robustness (DRO)

Table 2: Performance comparisons on Waterbirds and CelebA: averaged accuracy with training prior weight and uniform weight, δ -worst accuracy for several δ perturbations, and worst group-level test accuracy are reported. The *Group Info* column indicates whether the group labels are available to the methods on train/validation sets. The symbol $\checkmark\checkmark$ means the method also re-trains the last layer on the validation set (w/o needing the group information for training).

Method	Group Info		(Train prior)		Waterbirds ($r \rightarrow$ uniform prior)				
	Train	Val	Averaged	Averaged	$\delta = 0.05$ -Worst	$\delta = 0.1$ -Worst	$\delta = 0.2$ -Worst	$\delta = 0.5$ -Worst	Worst
ERM	\times	\checkmark	98.08 \pm 0.20	88.09 \pm 0.90	84.31 \pm 1.24	82.71 \pm 1.41	80.51 \pm 1.68	76.65 \pm 2.25	70.65 \pm 3.32
JTT	\times	\checkmark	93.05 \pm 0.36	89.56 \pm 0.69	88.59 \pm 0.74	88.24 \pm 0.75	87.81 \pm 0.77	87.13 \pm 0.85	86.18 \pm 1.08
DROPS	\times	\checkmark	97.95 \pm 0.16	88.55 \pm 1.15	85.50 \pm 1.51	84.33 \pm 1.64	82.81 \pm 1.79	80.41 \pm 2.00	77.14 \pm 2.15
G-DRO	\checkmark	\checkmark	93.03 \pm 0.34	91.67 \pm 0.22	91.23 \pm 0.33	91.06 \pm 0.34	90.83 \pm 0.40	90.44 \pm 0.52	89.85 \pm 0.73
SUBG	\checkmark	\checkmark	91.97 \pm 0.50	90.05 \pm 0.44	89.46 \pm 0.40	89.24 \pm 0.39	88.98 \pm 0.40	88.59 \pm 0.49	88.12 \pm 0.76
DFR _{Tr} ^{Tr}	\checkmark	\checkmark	95.83 \pm 0.94	93.45 \pm 0.49	92.77 \pm 0.48	92.51 \pm 0.51	92.16 \pm 0.55	91.58 \pm 0.67	90.72 \pm 0.91
DFR _{Tr} ^{Val}	\times	$\checkmark\checkmark$	93.17 \pm 1.30	93.29 \pm 0.80	92.98 \pm 0.84	92.86 \pm 0.85	92.70 \pm 0.87	92.42 \pm 0.88	92.01 \pm 0.88
DROPS*	\times	$\checkmark\checkmark$	93.01 \pm 1.32	93.42 \pm 0.61	93.08 \pm 0.81	92.98 \pm 0.90	92.76 \pm 1.02	92.45 \pm 1.23	91.99 \pm 1.56
Method	Group Info		(Train prior)		CelebA ($r \rightarrow$ uniform prior)				
	Train	Val	Averaged	Averaged	$\delta = 0.05$ -Worst	$\delta = 0.1$ -Worst	$\delta = 0.2$ -Worst	$\delta = 0.5$ -Worst	Worst
ERM	\times	\checkmark	95.33 \pm 0.12	81.18 \pm 1.60	73.78 \pm 2.36	70.53 \pm 2.69	65.97 \pm 3.15	57.82 \pm 3.98	44.89 \pm 5.30
JTT	\times	\checkmark	87.78 \pm 0.73	85.04 \pm 1.05	83.34 \pm 1.35	82.89 \pm 1.49	82.35 \pm 1.71	81.51 \pm 2.12	79.71 \pm 2.85
DROPS	\times	\checkmark	90.67 \pm 0.76	89.05 \pm 0.88	86.80 \pm 1.34	85.89 \pm 1.55	84.67 \pm 1.88	82.59 \pm 2.50	79.44 \pm 3.58
G-DRO	\checkmark	\checkmark	92.59 \pm 0.87	90.95 \pm 0.52	90.18 \pm 0.46	89.89 \pm 0.43	89.53 \pm 0.37	88.97 \pm 0.26	88.21 \pm 0.17
SUBG	\checkmark	\checkmark	91.09 \pm 0.62	88.75 \pm 0.34	87.69 \pm 0.31	87.27 \pm 0.41	86.72 \pm 0.73	85.80 \pm 0.29	84.45 \pm 0.28
DFR _{Tr} ^{Tr}	\checkmark	\checkmark	90.36 \pm 0.64	89.25 \pm 0.73	87.30 \pm 0.97	86.53 \pm 1.09	85.50 \pm 1.26	83.77 \pm 1.63	81.22 \pm 2.31
DFR _{Tr} ^{Val}	\times	$\checkmark\checkmark$	90.90 \pm 0.79	91.13 \pm 0.40	90.32 \pm 0.54	90.02 \pm 0.58	89.64 \pm 0.65	89.05 \pm 0.77	88.25 \pm 1.03
DROPS*	\times	$\checkmark\checkmark$	95.69 \pm 0.41	93.59 \pm 0.36	92.64 \pm 0.46	92.28 \pm 0.51	91.82 \pm 0.58	91.10 \pm 0.68	90.19 \pm 0.81

Theoretical Results

- Theorem 1 and Lemma 2: the optimal scaling of DROPS;
- Section 4.2: empirical implementation of DROPS;
- Theorem 3: convergence analysis of DROPS.

Empirical Results

- Table1: Experiments results of class-imbalanced CIFAR datasets;
- Table2: Experiments of group distributional robustness on Waterbirds, CelebA.

Thanks for watching!

Paper



Code

