

DATA W ULATIONS

Setting:

- Training data: Consists of imbalanced subgroups.
- Majority groups: Spurious features correlated with class labels.
- Minority groups: Only core features are predictive of class label. Worst group



Problem:

Core feature

- 1. ERM classifier predictions are incorrect on underrepresented subgroups.
- 2. Problem can be mitigated given group labels. But they are difficult to collect.

PREVIOUS APPROACHES

Use **train + validation** data with group annotations:

- Drastically improve worst group accuracy
- Require large amounts of annotated data

Use **validation** data with group annotations:

- Similar worst-group accuracy, at a reduced labeling cost
- Sometimes impossible to collect group-annotations (e.g. ethnicity, sexual orientation etc.)

Can we improve worst-group accuracy when **no group labels** are available?

YES! We perform model selection without group labels!

Boosting subgroup performance without any group annotations

Vincent Bardenhagen, Alexandru Ţifrea, Fanny Yang Department of Computer Science, ETH Zurich



DEBIASED CLA

First stage: Train biased predictor \hat{f}_{t_1,θ_1} using regularization t_1 and optimal hyperparameters θ_1 . Annotate samples in the error set of \hat{f}_{t_1,θ_1} as minority (g = 1). $\bar{S}(\hat{f}_{t_1,\theta_1}) = \{ (x_i, y_i, g_i) : (x_i, y_i) \in S, g_i = \mathbb{1}[\hat{f}_{t_1,\theta_1}(x_i) \neq y_i] \}$

Second stage: Train unbiased predictor using IW/GDRO:

 $\hat{f}_{t_1,\theta_1,t_2,\theta_2} = \arg\min \mathcal{L}_{\mathrm{IW}}(f,\bar{S}(\hat{f}_{t_1,\theta_1}))$ $f \in \mathcal{F}(t_1, \theta_1, t_2, \theta_2)$

Prior work: Hyperparameter tuning **requires group labels!** Select $t_1^*, \theta_1^*, t_2^*, \theta_2^*$ using WgAcc on group-annotated $\overline{V}_{\text{oracle}}$.

MODEL SELECTION WITHOUT GROUP LABELS

- . Early-stopping after one epoch: $t_1^* = 1$
- 2. Optimize AvgAcc on *V* to increase bias:

 $\theta_1^* \in \arg\max\operatorname{AvgAcc}(\hat{f}_{t_1,\theta_1}, V)$

3. Optimize WgAcc wrt estimated group labels $\overline{V}(\hat{f}_i)$:

EXPERIMENTAL RESULTS

		Corrupt-MNIST	Waterbirds	Cele
No group	ERM	71.2	74.9	60
labels	Ours	96.5	78.5	78
Val group	ERM WG	79.8	86.7	77
labels	JTT	91.3	86.7	81
Train & val group labels	GDRO	93.1	89.4	92

Our procedure significantly outperforms the ERM baseline

similar performance to approaches that use group labels.



