

What is Missing in IRM Training and Evaluation? **Challenges and Solutions**

Yihua Zhang¹, Pranay Sharma², Parikshit Ram³, Mingyi Hong⁴, Kush Varshney³, Sijia Liu^{1, 3}

¹Michigan State University, ²Carnegie Mellon University, ³IBM Research, ⁴University of Minnesota, Twin City



What is Missing in IRM Training and Evaluation?

Invariant Risk Minimization

- to acquire environment-agnostic data representations
- to avoid learning **spurious correlations** in the data

Problem Setup^[1]

- ▶ IRM is formulated as a bi-level optimization (BLO) problem: $\min_{\theta \in \mathcal{E}_{tr}} \mathcal{E}_{e \in \mathcal{E}_{tr}} \ell^{(e)}(\mathbf{w}^{*}(\theta) \circ \theta); \ s.t. \ \mathbf{w}^{*}(\theta) \in \min_{\mathbf{w}} \ell^{(e)}(\mathbf{w} \circ \theta), \forall e \in \mathcal{E}_{tr}$
- ▶ IRMv1 simplifies the BLO to a single-level problem: $\min_{\boldsymbol{\theta}} \sum_{e \in \mathcal{E}_{tr}} \left[\ell^{(e)}(\boldsymbol{\theta}) + \gamma \left\| \nabla_{w \mid w=1.0} \ell^{(e)}(w \circ \boldsymbol{\theta}) \right\|_{2}^{2} \right]$

Challenge I: Large-Batch Training

- Large-batch optimization causes **suboptimal** IRM training.
- Our proposal: Small-batch training is effective versus a zoo of large-batch optimization enhancements.

Method

IRMv1

IRMv0

IRM-GAME

REX

BIRM

SPARSEIRM

FISHR

Average



Figure 1. The performance of
three IRM methods (IRMV1,
IRMV0, and REX) vs. batch-
size under Colored-MNIST.

Table 1. Accuracy of IRM methods on Colored-MNIST using the original large-batch implementation (Original), the large-batch optimization-integrated implementations (LSGD/LALR/SAM), and the small-batch training recipe.

Original LSGD LALR

65.82

67.53

67.99

67.85

67.82

65.47

67 59

68 21

67.99

67.93

67.25 67.34 67.63 **68.44**

67.13

65.39

65.69

67.42

67.93

67.72

67.88

67.02

SAM Small

66.23

67.82

68.32

68.13

68.11

68.37

67.73

68.42

68.71

68.81

68.69

Reference [1] Martin Arjovsky et al., Invariant risk minimization. [2] Kartik Ahuja, et al. Invariant risk minimization games.



- \blacktriangleright The evaluation metric adopted by all existing work focuses on a single test environment.
- ▶ IRM performance is sensitive to test environment choices. Singleenvironment evaluation leads to a false sense of invariance.
- > Our proposal: evaluation across multiple environments. Good method should achieve high Avg. Acc. and low Acc. Gap.
- Challenge III: IRM-Game^[2] with Invariant Predictor
- > IRM-Game assigns each environment an individual classifier $w^{(e)}$. The output relies on the ensemble of individual predictors.
- > Our proposal: BLOC-IRM (BLO with Consensus IRM):

minimize $\sum_{e \in \mathcal{E}_{tr}} [\ell^{(e)}(\mathbf{w}^{*}(\theta) \circ \theta) + \gamma \| \nabla_{\mathbf{w}} \ell^{(e)}(\mathbf{w}^{*}(\theta) \circ \theta) \|_{2}^{2}]$

subject to $w^{(e)}(\theta) \in \operatorname{argmin} \ell^{(e)}(w^{(e)} \circ \theta), \forall e \in \mathcal{E}_{tr}; w^*(\theta) = \frac{1}{N} \sum_{e \in \mathcal{E}_{tr}} w^{(e)}(\theta)$

- The lower level ensures (I) per-environment risk minimization and (II) an environment-invariant predictor.
- The upper level minimizes the ERM loss and the regularization term penalizes the lower-level stationarity.
 - This problem can be solved using an ordinary BLO solver.



Figure 3. (Algorithm Comparison) Schematic overview of BLOC-IRM over two training environments (red and green), and its comparison to IRM and IRM-Game.

Dataset	Invariant	Spurious	Env 1	Env 2	
Colored- MNIST	Digit	Color	01	01	
Colored- FMNIST	Object	Color	~ ^	3	
CIFAR- MNIST	CIFAR	MNIST			
Colored- Object	Object	Color	***		
CelebA	Smiling	Hair Color		()	
PACS	Object	Texture	<u>R</u> 20	X 😭	
VLCS	Object	Environment	-	77	

Figure 4. (Dataset Setups) An overview of the 'Invariant' and 'Spurious' features in different datasets used in this work.

Experiment Results

Algorithm	COLORED Avg Acc (%	о-Овјест 6) Acc Gap	CIFAR-] Avg Acc (9	MNIST 6) Acc Gap	CEL Avg Acc (%	EBA b) Acc Gap	VL Avg Acc (9	CS 6) Acc Gap	PA Avg Acc (%	CS 6) Acc Gap
ERM	41.11 ± 1.44	$86.43{\pm}2.89$	$40.39{\pm}1.32$	85.53±2.33	$72.38{\pm}0.29$	$10.73{\pm}0.36$	$63.23{\pm}0.23$	$12.39{\pm}0.35$	$69.95{\scriptstyle\pm0.35}$	$14.32{\pm}0.75$
IRMv1	64.42 ± 0.21	$4.18 {\pm} 0.29$	61.49 ± 0.29	7.17±0.33	$72.49 {\pm} 0.38$	$10.15 {\pm} 0.27$	62.72 ± 0.29	$12.74 {\pm} 0.27$	68.93±0.33	$14.99 {\pm} 0.51$
IRMv0	$62.39 {\pm} 0.25$	5.36 ± 0.31	$60.14 {\pm} 0.18$	$8.83 {\pm} 0.39$	72.42 ± 0.35	$10.43{\pm}0.38$	$62.59 {\pm} 0.32$	$12.99 {\pm} 0.36$	68.72 ± 0.29	$15.29 {\pm} 0.71$
IRM-GAME	$62.88 {\pm} 0.34$	5.59 ± 0.28	$60.44 {\pm} 0.31$	6.72 ± 0.41	$72.18 {\pm} 0.44$	$12.32 {\pm} 0.41$	62.31 ± 0.38	$13.37 {\pm} 0.62$	68.12 ± 0.22	$15.77 {\pm} 0.66$
REX	63.37 ± 0.35	5.42 ± 0.31	62.32 ± 0.24	5.55 ± 0.32	72.34 ± 0.26	10.31 ± 0.23	63.19 ± 0.31	12.87 ± 0.31	69.43 ± 0.34	15.31 ± 0.67
BIRM	65.11 ± 0.27	3.31±0.22	62.99 ± 0.35	5.23 ± 0.36	$72.93 {\pm} 0.28$	$9.92 {\pm} 0.33$	$63.33 {\pm} 0.40$	$12.13 {\pm} 0.23$	69.34 ± 0.25	$15.76 {\pm} 0.49$
SPARSEIRM	64.97 ± 0.39	3.97 ± 0.25	62.16 ± 0.29	4.14±0.31	72.42 ± 0.33	$9.79 {\pm} 0.21$	$62.86 {\pm} 0.26$	12.79 ± 0.35	69.52 ± 0.39	$15.81 {\pm} 0.82$
FISHR	64.07 ± 0.23	4.41 ± 0.29	61.79 ± 0.25	5.55 ± 0.21	72.89 ± 0.25	$9.42 {\pm} 0.32$	63.44 ± 0.37	11.93 ± 0.42	70.21 ± 0.22	14.52±0.43
BLOC-IRM	65.97±0.33	$4.10{\pm}0.36$	63.69 ±0.32	$4.89{\pm}0.36$	$\textbf{73.35}{\pm}0.32$	$\textbf{8.79}{\pm}0.21$	63.62±0.35	$11.55{\pm}0.32$	70.31±0.21	$14.73{\pm}0.65$

Table 2. (Main Results) IRM performance comparison between BLOC-IRM and other baselines.



Environment	$p_{tr} \in \{0\}$.1, 0.15	$p_{tr} \in \{0.1, 0.15, 0.2\}$		
Metrics (%)	Avg Acc	Acc Gap	Avg Acc	Acc Gap	
Optimum	75.00	0.00	75.00	0.00	
GRAYSCALE	73.82 ± 0.11	0.37 ± 0.05	$73.97 {\pm} 0.14$	$0.29 {\pm} 0.08$	
ERM	49.21±0.79	$91.88{\pm}3.31$	$49.03{\scriptstyle\pm0.93}$	$92.17{\pm}3.04$	
IRMv1	67.36±0.31	2.77 ± 0.15	$67.11 {\pm} 0.34$	$2.42 {\pm} 0.12$	
IRMv0	67.01 ± 0.42	2.85 ± 0.18	66.71 ± 0.42	2.36 ± 0.19	
IRM-GAME	66.39 ± 0.72	4.47 ± 0.61	65.93 ± 0.53	4.25 ± 0.84	
REX	66.82 ± 0.44	2.59 ± 0.11	67.14 ± 0.38	2.16 ± 0.13	
BIRM	67.35 ± 0.39	2.65 ± 0.10	$68.05 {\pm} 0.43$	1.99±0.07	
SPARSEIRM	67.12 ± 0.53	2.33 ± 0.18	67.72 ± 0.41	2.11 ± 0.19	
FISHR	67.22 ± 0.43	2.44 ± 0.15	$67.32 {\pm} 0.39$	$2.59 {\pm} 0.15$	
BLO-IRM	68.72±0.41	2.19 ±0.15	68.89 ±0.31	$2.39{\pm}0.09$	

Figure 5. Ablation study on model size on the dataset Colored-MNIST

Table 3. Ablation study on different training environments on the dataset Colored-MNIST.

100 80 60 Test Env 20 ____ IRM-Game Fishr 0 0.0 0.2 0.4 0.6 0.8 1.0

Data Environment (β)

 β Dist. I+ $(1 - \beta)$ Dist. II

Training Env. I: $\beta = 0.1$

Training Env. II: $\beta = 0.2$

Figure 2. Performance comparison of different IRM methods under diverse test environments. Existing methods only evaluate at $\beta = 0.9$.

Paper