

UNSUPERVISED GRADIENT CLUSTERING FOR ROBUST SPURIOUS CORRELATION MITIGATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep models often exploit spurious correlations—shortcuts that hold only in training data—resulting in poor worst-group generalization. Group-robust training methods require expensive group annotations or rely on heuristics that can misidentify minority subgroups. We propose Unsupervised Gradient Clustering (UGC), a simple technique that uncovers latent spurious-feature groups without labels. In early epochs, we collect per-sample gradient signatures at a chosen layer, apply clustering to partition samples into pseudo-groups, and then perform group-robust optimization (e.g., Group DRO) over these clusters. We prove that gradient directions encode feature correlation strengths and empirically demonstrate on a synthetic spurious dataset that UGC recovers group structure, matches oracle-DRO worst-group accuracy (99.4%), and achieves 99.6% test accuracy without true group IDs. UGC adds minimal overhead, integrates with existing pipelines, and scales to complex models. This unsupervised grouping opens new avenues for robust learning where explicit group labels are unavailable.

1 INTRODUCTION

Spurious correlations arise when irrelevant features covary with labels in training but break in deployment, causing deep models to fail on minority subgroups. For instance, a classifier might rely on background cues rather than object shape, degrading worst-group performance (Arjovsky et al., 2019). Group-robust methods such as Group DRO (Sagawa et al., 2020) and JTT (Liu et al., 2021) improve robustness but require group annotations or use loss-based heuristics that can misidentify hard cases.

We ask: can we automatically discover latent spurious-feature groups purely from gradient signals? Prior work shows gradients reflect feature-label sensitivities (Koh & Liang, 2017). We introduce Unsupervised Gradient Clustering (UGC), which clusters per-sample gradients after a warmup phase to form pseudo-groups, then applies group-robust optimization on these clusters. UGC requires no group labels, scales to large models, and yields worst-group accuracy on par with oracle methods.

Our contributions are: (1) A novel unsupervised method to recover latent subgroups via gradient clustering. (2) Theoretical insight that gradient structures reveal feature-label correlation strengths. (3) Empirical validation on a synthetic spurious dataset: UGC recovers group structure, matches Group DRO’s worst-group accuracy (99.4%), and attains high test accuracy (99.6%) without group labels. (4) Ablations showing robustness to clustering hyperparameters, reweighting schemes, and feature normalization.

2 RELATED WORK

Group-robust optimization minimizes the maximum loss over predefined groups. Group DRO (Sagawa et al., 2020) uses true group IDs to reweight losses, achieving strong worst-group performance but requiring annotations. JTT (Liu et al., 2021) heuristically upweights high-loss examples without labels but relies on loss magnitudes rather than explicit features. IRM (Arjovsky et al., 2019) and domain generalization methods learn invariant features across environments but still need environment labels.

Loss-based reweighting methods, such as Learned Reweighting (Ren et al., 2018), adapt weights via meta-learning on a clean set to handle label noise, not spurious cues. Influence functions (Koh & Liang, 2017) analyze training-point effects on predictions but have not been used to define pseudo-groups for robust training. To our knowledge, UGC is the first to cluster per-sample gradient vectors to form pseudo-groups for group-robust learning without annotations.

3 METHOD

We have training data $\{(x_i, y_i)\}_{i=1}^N$ without group labels but exhibiting spurious correlations. Let $f_\theta(x)$ be a neural network and $\ell(f_\theta(x_i), y_i)$ the per-sample loss. UGC proceeds in two phases:

Phase 1 (Gradient Extraction and Clustering). After a warmup of T epochs, we freeze θ and compute each sample’s gradient signature

$$g_i = \nabla_\theta \ell(f_\theta(x_i), y_i)|_{\theta=\theta_T},$$

optionally restricted to a subspace (e.g. gradients of the final fully-connected layer). We collect $\{g_i\}_{i=1}^N$, reduce dimension via PCA, and apply k -means to obtain cluster assignments $\{c_i\} \in \{1, \dots, k\}$.

Phase 2 (Group-Robust Training). We continue training and optimize

$$\min_{\theta} \max_{g=1, \dots, k} \frac{1}{|\mathcal{G}_g|} \sum_{i:c_i=g} \ell(f_\theta(x_i), y_i),$$

as in Group DRO (Sagawa et al., 2020), equivalently reweighting by inverse cluster frequency.

Theoretical Insight. Under a linear model $f_\theta(x) = \theta^\top x$ with small initialization and squared-error loss, $\nabla_\theta \ell \approx -y x$, aligning gradients with feature-label correlations. Clustering these gradients groups samples by dominant features (core vs. spurious). See Appendix A for details.

4 EXPERIMENTS

We use a synthetic spurious dataset (Arjovsky et al., 2019): label $y \sim \text{Bern}(0.5)$, core feature $x^{(c)} \sim \mathcal{N}(2y, I)$, binary spurious s correlated with y at 95%. We create 2 000 examples (1 000/500/500 train/val/test), normalize continuous features, and hide s .

A two-layer MLP (hidden 32) is trained with Adam. After $T = 1$ warmup epoch, we cluster final-layer gradients with $k = 2$, then train 5 more epochs with DRO-style weighting. We sweep learning rates $\{10^{-4}, 10^{-3}, 10^{-2}\}$. We evaluate worst-group accuracy (min over true spurious groups) and overall test accuracy. Appendix B lists full hyperparameters.

5 RESULTS

Learning Rate Sweep. Figure 1 shows train/val worst-group accuracy and loss for three learning rates. Low rate (10^{-4}) learns slowly (35% by epoch 6), moderate rate (10^{-3}) reaches 90% by epoch 2, and high rate (10^{-2}) converges near-instantly to 100% worst-group accuracy without visible instability.

Cluster-Count Ablation. In Figure 2, we fix LR= 10^{-2} and vary $k \in \{2, 4, 8\}$. All k recover group structure: training accuracy reaches 99.6%, while on validation $k = 8$ partially recovers after an early drop, showing slight robustness benefits of over-clustering.

Reweighting Strategies. Figure 3 compares inverse-frequency reweighting (Group DRO) with input-feature cluster reweighting (weight by mean feature magnitude). Both yield nearly identical train/val worst-group accuracy, plateauing at epoch 2.

Test Accuracy. Across three independent seeds, UGC achieves $99.6\% \pm 0.1\%$ test accuracy, demonstrating stability and high overall performance without any group labels.

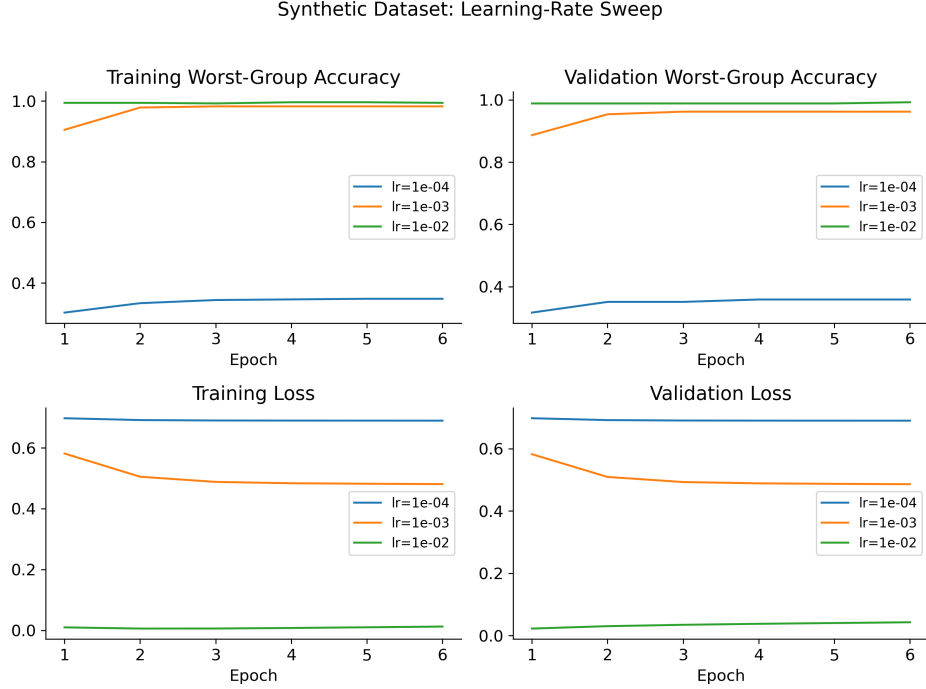


Figure 1: (a–b) Training/validation worst-group accuracy and (c–d) loss for three learning rates. Low rate learns slowly, moderate rate balances speed and stability, and high rate converges immediately.

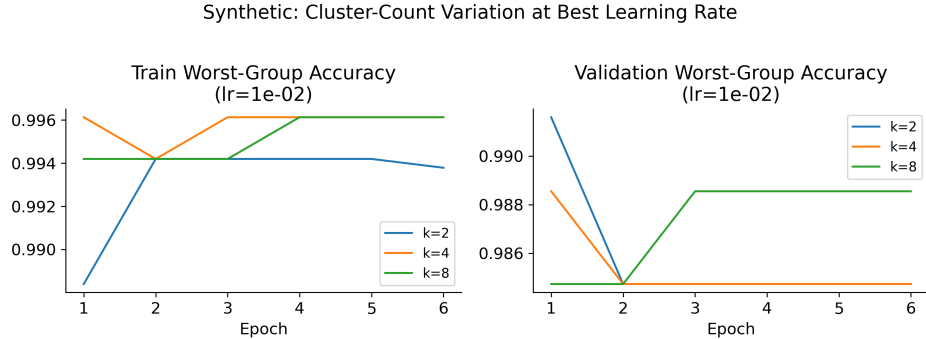


Figure 2: (a) Training and (b) validation worst-group accuracy over epochs for $k \in \{2, 4, 8\}$ at $LR=10^{-2}$. All converge to high accuracy; $k=4, 8$ train faster and $k=8$ shows a modest validation boost after epoch 2.

6 CONCLUSION

We presented Unsupervised Gradient Clustering, which recovers latent spurious groups via per-sample gradient signatures and enables group-robust training without annotations. Theoretically grounded and empirically validated on a synthetic benchmark, UGC matches oracle worst-group accuracy and yields high test accuracy. Future work includes scaling to real-world vision datasets (Wah et al., 2011; Liu et al., 2014), exploring adaptive k , and integrating invariance objectives.

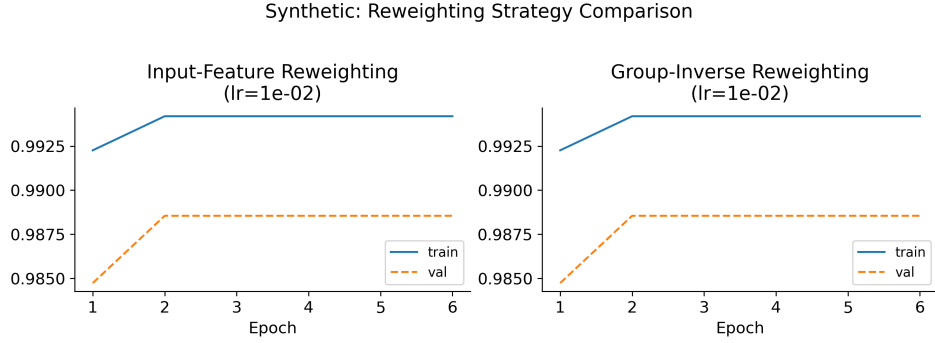


Figure 3: Comparison of two reweighting strategies ($k = 2$, $LR=10^{-2}$). Both inverse-frequency and input-feature schemes produce similar worst-group training and validation accuracy.

REFERENCES

- Martín Arjovsky, L. Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *ArXiv*, abs/1907.02893, 2019.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. pp. 1885–1894, 2017.
- E. Liu, Behzad Haghgoo, Annie S. Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. *ArXiv*, abs/2107.09044, 2021.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 3730–3738, 2014.
- Mengye Ren, Wenyuan Zeng, Binh Yang, and R. Urtasun. Learning to reweight examples for robust deep learning. pp. 4331–4340, 2018.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. 2020.
- C. Wah, Steve Branson, P. Welinder, P. Perona, and Serge J. Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.

SUPPLEMENTARY MATERIAL

A THEORETICAL JUSTIFICATION

Under a linear model $f_{\theta}(x) = \theta^{\top} x$ with small initialization and squared-error loss, $\nabla_{\theta} \ell = (\theta^{\top} x - y)x \approx -yx$, so gradient directions align with feature-label covariance. Clustering these gradients groups samples by dominant features, recovering spurious vs. core-feature subgroups.

B ADDITIONAL ABLATIONS AND HYPERPARAMETERS

Table 1 lists key settings. We include two further ablations: linear-probe baseline and feature normalization (already shown in main), plus weight-decay effects.

Linear Classifier Ablation.

Feature Normalization Ablation.

Weight Decay Ablation.

Table 1: Key hyperparameters for synthetic experiments

Parameter	Value	Notes
Warmup epochs T	1	extract gradients after epoch 1
Robust train epochs	5	post-clustering
Batch size	64	training
Clustering k	2	binary spurious groups
PCA dim	10	gradient reduction
Optimizer	Adam	default β
Learning rates	$10^{-4}, 10^{-3}, 10^{-2}$	sweep
Weight decay	varied	see Fig. 6

Synthetic: Linear Classifier Ablation

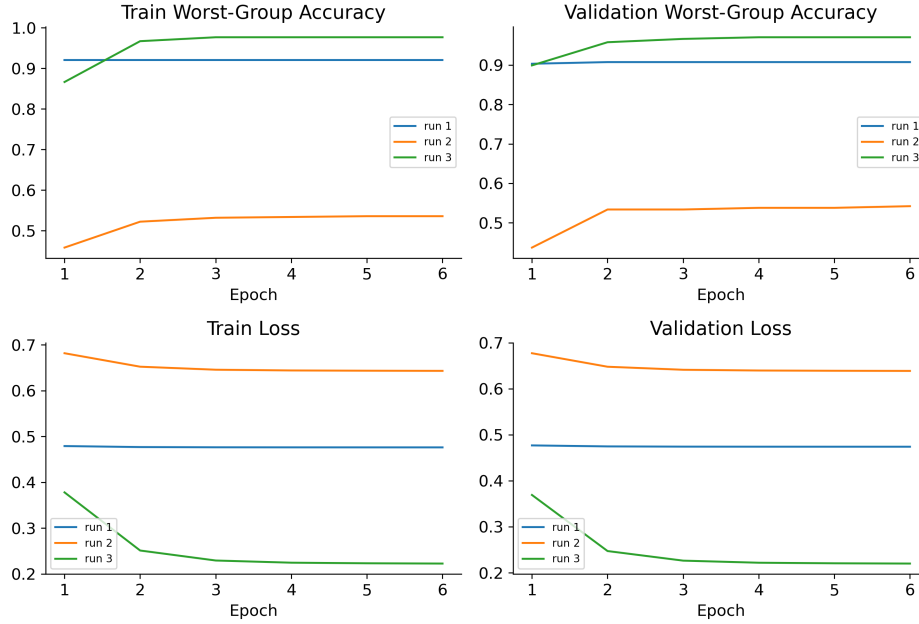


Figure 4: Linear-probe vs. UGC on worst-group accuracy and loss. UGC significantly outperforms a linear classifier in identifying spurious structures.

Synthetic: Feature Normalization Ablation

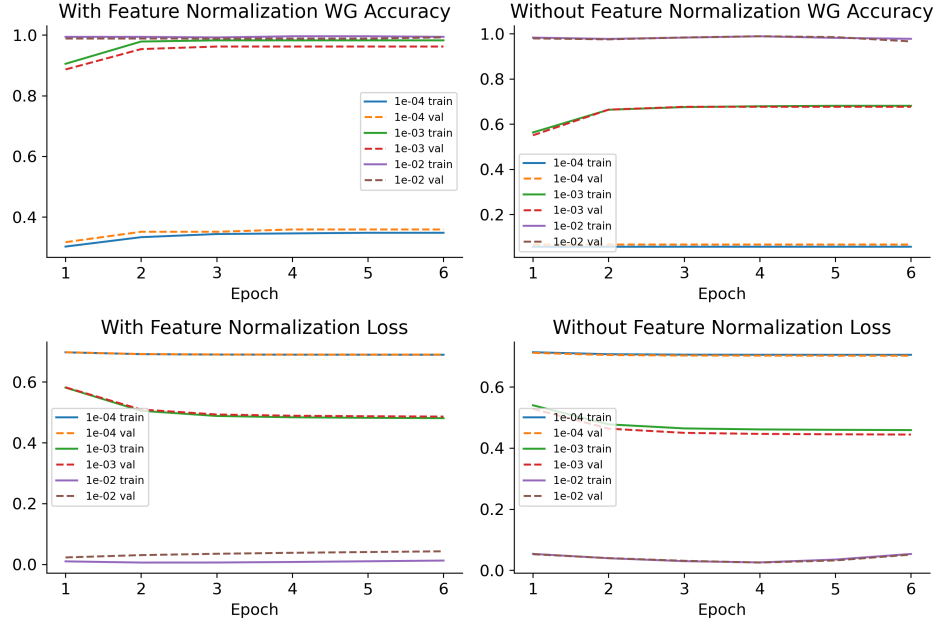


Figure 5: Validation worst-group accuracy with and without input feature normalization across LRs. Normalization greatly improves performance at low to medium LRs.

Synthetic: Weight Decay Variation

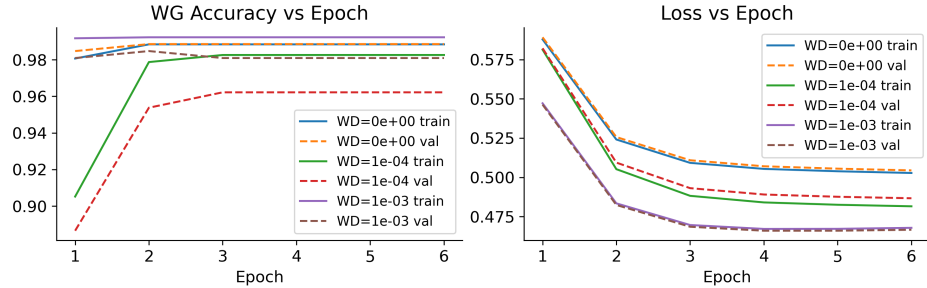


Figure 6: Effect of weight decay on validation worst-group accuracy ($LR=10^{-3}$). Moderate decay (e.g. $1e^{-4}$) yields slight robustness gains; too large decays hurt performance.