
Understanding Post-hoc Adaptation for Improving Subgroup Robustness

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A number of deep learning approaches have recently been proposed to improve
2 model performance on subgroups under-represented in the training set. However,
3 Menon et al. [13] recently showed that, models with poor subgroup performance
4 can still learn representations which contain useful information about these sub-
5 groups. In this work, we explore the representations learned by various approaches
6 to robust learning, finding that different approaches learn very similar represen-
7 tations. We probe a range of post-hoc procedures for making predictions from
8 learned representations, showing that the distribution of the post-hoc validation set
9 is paramount, and that clustering-based methods may be a promising approach.

10 1 Introduction

11 Machine learning systems trained with expected risk minimization (ERM) often struggle to perform
12 well on under-represented subgroups in the training data [2, 5]. For this reason, a number of learning
13 algorithms have been proposed to improve performance across subgroups, many taking advantage of
14 subgroup information [3, 4, 8, 11]. However, it has also been noted that deep models trained with
15 ERM can sometimes learn representations that contain sufficient information to perform well on all
16 subgroups [13], even when the model’s predictions yield large subgroup disparities.

17 In this work, we explore this apparent disconnect between representation learning and prediction in
18 the context of group-robust classification. We consider a three-stage procedure, where a model is
19 trained, then some post-hoc adaptation occurs directly on the learned representations, then tested
20 for performance on various subgroups. We find that ERM-learned representations can be practically
21 identical to those learned by more specialized methods that take advantage of subgroup information.
22 Given this, we explore a range of procedures for post-hoc adaptation of a model by learning a new
23 classifier on the representations directly, using a validation set and no subgroup information. We
24 show that the distribution of the validation set is extremely important to obtaining good subgroup
25 performance, and that clustering methods in representation space may be better than linear classifiers.
26 Our results suggest that learning better representations for subgroup classification is a promising
27 direction, and that post-hoc adaptation can be helpful for improving robustness.

28 2 Background

29 **Notation.** We assume a classification dataset of (input, label) pairs $\{(x_i, y_i)\}_{i=1}^n$, with $x \in \mathcal{X}$, $y \in$
30 $\mathcal{Y} = \{1 \dots Y\}$. We may also have a categorical subgroup variable $\{g_i\}_{i=1}^n$ ($g \in \{1 \dots G\}$): $c_i = g \leftrightarrow$
31 example i is in subgroup g . Subgroup information is assumed available at training/validation time,
32 but not test time. Usually, c will not be distributed uniformly throughout the training set — rather,
33 some subgroups will be smaller, i.e. some values of c will be uncommon. Our goal is to learn a
34 classification function $\bar{f} : \mathcal{X} \rightarrow \mathcal{Y}$ which performs well on all subgroups. The metric by which we

35 will evaluate our success is *worst-group accuracy*, that is, $\min_g \mathbb{E}[\mathbb{1}[\bar{f}(x) = y] | g_i = g]$. In this work,
36 we assume \bar{f} takes the form: $\bar{f} = \operatorname{argmax}_f f(x); f(x) = \sigma(w^\top r(x) + b)$. Here, σ is the softmax
37 function, w is a matrix containing a vector of weights for each class, b is a scalar bias, and r is
38 a function outputting a vector representation of x . The model we use throughout, satisfying this
39 functional form, is a Resnet-50 [7].

40 **Train-Adapt-Test Procedure.** Since we are interested in three aspects of model training (repre-
41 sentation learning, post-hoc adaptation, and subgroup performance on test data), we consider a
42 three-stage procedure. First, we initialize a model and **train** it on a training set with an unbalanced
43 subgroup distribution; here, we learn the prediction function f and, as a byproduct, the representation
44 function r . Next, we focus on r only, and perform post-hoc **adaptation**, using a validation set to
45 learn a simple classifier on top of r . Finally, we **test** our adapted model (the composition of our
46 post-hoc classifier and r) on held-out data, and record the performance on each subgroup.

47 **Algorithms.** We discuss several training algorithms, described fully in App. A.2. Expected risk
48 minimization (**ERM**) is the usual paradigm for training ML models, where we ignore subgroup
49 information and minimize mean loss on the training set. We also look at two **robust** approaches,
50 which aim to take advantage of subgroup information g . The first is Group Distributionally Robust
51 Optimization (**GDRO**) [19], which aims explicitly to minimize the worst group’s average loss. The
52 second is Invariant Risk Minimization (**IRM**) [1], uses a gradient penalty with the goal of learning a
53 representation such that the same predictive classifier is optimal across subgroups.

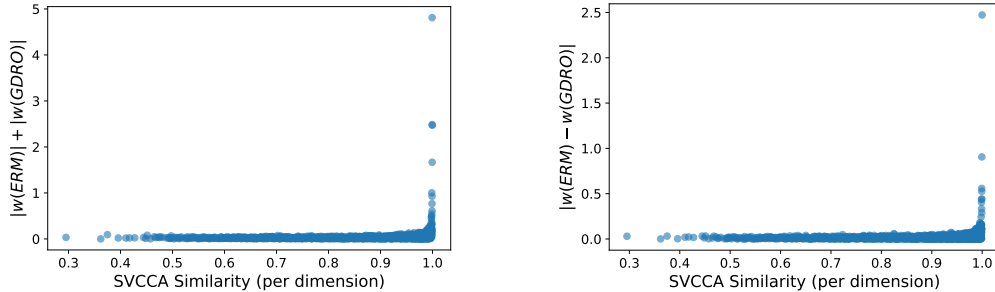
54 **Dataset.** We use the semi-synthetic Waterbirds dataset [19], which is created by pasting pictures
55 of birds from CUB [21] onto backgrounds from Places-365 [22]. The task is to predict whether
56 the bird in the image is a land or water bird; it is confounded by the background, which can be
57 either land or water backgrounds. This yields four subgroups: land birds on land, water birds on
58 water, land birds on water, and water birds on land. In the training set, there is a strong correlation
59 between the bird and background factors: e.g. water birds are usually shown on water backgrounds.
60 The dataset also contains a validation set, which is much more subgroup-balanced than the training
61 set; the authors use the worst-subgroup accuracy on this validation set for early stopping. ERM
62 obtains about 60% worst-group accuracy on Waterbirds, and GDRO/IRM obtain 87-90%. While
63 subgroup classification may be harder on more realistic datasets, Waterbirds is a helpful tool to better
64 understand the properties of various approaches in a research context. Since Waterbirds is small, we
65 always initialize from a model pre-trained on Imagenet.

66 3 Feature Co-Discovery in Robust Learning

67 Deep learning approaches leveraging subgroup information g can improve worst-group performance
68 [8, 9, 11, 12, 20]. We ask here: how do the learned representations reflect these more specialized
69 approaches? Does using subgroup information at training time, or using more specialized algorithms,
70 produce richer or better-separated learned features?

71 Prior literature suggests that differently-performing methods may nonetheless learn similar repre-
72 sentations. For instance, Menon et al. [13] note that ERM features can be used to obtain similar
73 subgroup performance to GDRO by learning a post-hoc linear model on a group-balanced validation
74 set. This suggests that the necessary information needed to improve subgroup performance is already
75 present in the ERM-learned representations, and that it is linearly extractible. In a meta-learning
76 context, Raghu et al. [16] compares the representations learned by a model at its meta-initialization
77 with the representations learned after performing task-specific adaptation. They find evidence of
78 *feature re-use*: that the representations before and after task-specific adaptation are similar, and the
79 meta-initialized model and the task-adapted model differ mostly in their final classification heads.

80 In Fig. 2, we show the similarity of learned representations on Waterbirds across models trained
81 with several loss functions. “None” is an Imagenet pre-trained model (not trained on Waterbirds);
82 “ERM”, “GDRO”, and “IRM” use the methods from Sec. 2, initialized from the “None” model and
83 trained (holding the random seed constant) on Waterbirds. “ERM-2” is the same as “ERM” with
84 a different random seed. To compare the representations, we use SVCCA [17], which determines
85 the most-aligned dimensions between the representations produced by two layers of neurons (here,
86 the final layers of two different models) across some dataset (here, the Waterbirds validation set).



(a) The y-axis shows the sum of the absolute values in from the ERM and GDRO transformed classifiers. We see that the only features which are important for classification in the two models are co-discovered.

(b) The y-axis shows the absolute difference in the transformed classifier between ERM and GDRO models. We see the features which are treated highly differently between the two models are co-discovered.

Figure 1: Waterbirds feature co-discovery in ERM- and GDRO-learned representations. X-axis shows SVCCA similarity score for each dimension; a feature is co-discovered if it has a score near 1.

87 SVCCA returns a similarity score between 0 and 1 for each pair of (aligned) dimensions describing
 88 how “similar” they are between the two representations; these can be averaged to produce an overall
 89 similarity score which is “a direct multidimensional analogue of Pearson correlation”, and describes
 90 the overall similarity of the two representations. Like Pearson correlation, closer to 1 is more similar.

91 Fig. 2 shows that ERM, GDRO and IRM all
 92 learn very similar representations on Waterbirds,
 93 with SVCCA overall similarities of ~ 0.9 . This
 94 is comparable to the similarity between two dif-
 95 ferent seeds of ERM, suggesting that the impact
 96 of these different algorithms on the represen-
 97 tation is minimal here, despite methodological
 98 differences and the availability of subgroup in-
 99 formation g to GDRO and IRM. Also, these
 100 methods are more similar to each other (~ 0.9)
 101 than to their initialization (~ 0.7). This ex-
 102 plains the result from Menon et al. [13]: we
 103 can match GDRO performance using a post-
 104 hoc linear model on ERM representations be-
 105 cause GDRO representations and ERM repre-
 106 sentations are roughly equal. Contrasted with
 107 Raghu et al. [16], we show that two models, us-
 108 ing different learning algorithms, learned similar features to each other, but fairly different features
 109 from the initialized model; we call this *feature co-discovery* (rather than re-use [16]).

	ERM	GDRO	IRM	None	ERM-2
ERM	1	0.89	0.91	0.71	0.93
GDRO	0.89	1	0.97	0.72	0.88
IRM	0.91	0.97	1	0.73	0.9
None	0.71	0.72	0.73	1	0.71
ERM-2	0.93	0.88	0.9	0.71	1

Figure 2: SVCCA overall similarities between learned representations on Waterbirds. None is an Imagenet pre-trained model; ERM/GDRO/IRM are as in Sec. 2, with the same seed, using None as an initialization; ERM-2 is like ERM but with a different seed. 1 is perfect similarity.

110 We explore feature co-discovery between ERM and GDRO in Fig 1 (with similar ERM/IRM results
 111 in App. B). We use SVCCA to transform the representations from ERM and GDRO into their most
 112 similar directions, and obtain an equivalent linear classifier in transformed space for each method,
 113 such that the transformed representations and the transformed classifier together output the same
 114 logits as the non-transformed model. In Fig 1a, we plot the SVCCA similarity scores for each
 115 dimension on the x-axis. On the y-axis, we plot the sum of the absolute values of the transformed
 116 classifier weights from ERM and GDRO. We observe that all dimensions which are important for
 117 classification in either model are co-discovered (i.e. have high similarity score). In Fig 1b, we instead
 118 plot the absolute difference between the ERM and GDRO transformed classifiers on the y-axis. We
 119 observe that all dimensions where these two methods *differ* for classification are also co-discovered.
 120 Both of these plots suggest that the improvement in subgroup disparities on Waterbirds shown by
 121 GDRO is due to, not improvements in learned features, but a classification layer which weights the
 122 same features differently. See more experimental details for this section and the next in App. A.

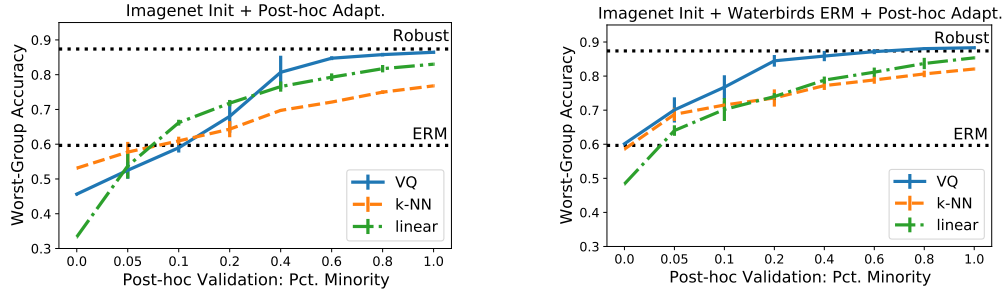


Figure 3: Waterbirds post-hoc classification results. On the left, we use the representations obtained by a model pre-trained on Imagenet; on the right, we use the representations learned by ERM. On the x-axis, we show (number of minority examples) / (number of majority examples) in the post-hoc validation set. The y-axis shows the worst-group accuracy from the 4 subgroups on the test set. “Robust” shows the performance of GDRO/IRM, and “ERM” shows the performance of ERM; both are fully trained, starting from a pretrained Imagenet initialization. Average over 3 seeds shown.

123 4 Experiments: Post-hoc Classification

124 In Sec. 3, we show representations learned by ERM are similar to those learned by more specialized
 125 algorithms, including those using group information g . This supports the findings of Menon et al.
 126 [13], who show that by doing post-hoc logistic regression with a subgroup-balanced dataset, one
 127 can match the subgroup performance of GDRO. Here, we attempt to empirically disentangle two
 128 aspects of this finding, to determine which factor of the post-hoc learning procedure is important for
 129 subgroup performance. First is the post-hoc classification algorithm used: training a post-hoc linear
 130 model might induce larger disparities than a more flexible model would. Second is the distribution of
 131 the post-hoc validation set. As discussed previously, the given Waterbirds validation set is subgroup-
 132 balanced ($P(\text{background} \mid \text{label}) = 0.5$) — this is very different from the training set, and may encode
 133 a lot of information about the subgroups, even if only used post-hoc.

134 The three adaptation algorithms we look at are logistic regression, k-NN and vector quantization
 135 (VQ). In each, we train the post-hoc model on the representations in the validation set, and test them
 136 on representations from a second held-out set (the test set). For our VQ classifier: we fit a k-means
 137 model to each of the two classes on the validation set, returning two sets of k centroids, and use these
 138 2k centroids in a 1-NN classifier for a given test point. To probe the importance of the validation
 139 distribution, we keep only $p\%$ of the minority examples in the validation set (minority examples have
 140 e.g. water bird on land background), and vary p . At $p = 1$, we have the original validation set — for
 141 each y , the number of land and water backgrounds are the same. At $p = 0$, we have removed all
 142 minority examples from the validation set. At $p \approx 0.05$, we have a validation set which is distributed
 143 similarly to the training set.

144 In Fig. 3, we show the effects of these two factors for the representations at the Imagenet pre-trained
 145 initialization, and the representations after ERM training on Waterbirds. We note that the distribution
 146 of the validation set is important, with a large difference in worst-group accuracy between the training
 147 set’s proportion ($p = 0.05$) and the given validation set ($p = 1$). Indeed, post-hoc adaptation on the
 148 given validation set, using only a pre-trained model (which never sees any Waterbirds training data)
 149 matches the performance of specialized robust methods, which trains on the full Waterbirds training
 150 set as well as subgroup information, and might be considered an upper bound for post-hoc approaches.
 151 Secondly, the difference between the methods is smaller but noteworthy, and the difference between
 152 adaptation and the original model (horizontal line labelled “ERM”) is large. Using the current
 153 validation set, VQ achieves the best performance, and this advantage is robust to perturbations in p :
 154 e.g. even at $p = 0.05$, VQ improves significantly over both the original model and linear adaptation.

155 **Conclusion.** We draw attention to the similarity of representations learned by robust approaches to
 156 those learned by ERM — this suggests that robust approaches which improve representation learning
 157 are potentially promising. On the other hand, the utility of post-hoc adaptation here stresses the
 158 richness of ERM representations, and it may be better to find methods which harness that richness
 159 through adaptation, rather than learning different ones.

References

- 160
- 161 [1] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. Invariant risk minimization. *arXiv*
162 *preprint arXiv:1907.02893*, 2019.
- 163 [2] J. Buolamwini and T. Gebru. Gender shades: Intersectional accuracy disparities in commercial
164 gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91.
165 PMLR, 2018.
- 166 [3] E. Creager, J.-H. Jacobsen, and R. Zemel. Environment inference for invariant learning. In
167 *ICML Workshop on Uncertainty and Robustness*, 2020.
- 168 [4] N. Dagaev, B. D. Roads, X. Luo, D. N. Barry, K. R. Patil, and B. C. Love. A too-good-to-be-true
169 prior to reduce shortcut reliance. *arXiv preprint arXiv:2102.06406*, 2021.
- 170 [5] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann.
171 Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- 172 [6] I. Gulrajani and D. Lopez-Paz. In search of lost domain generalization. *arXiv preprint*
173 *arXiv:2007.01434*, 2020.
- 174 [7] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In
175 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–
176 778, 2016.
- 177 [8] C. Heinze-Deml and N. Meinshausen. Conditional variance penalties and domain shift robust-
178 ness. *Machine Learning*, 110(2):303–348, 2021.
- 179 [9] W. Hu, G. Niu, I. Sato, and M. Sugiyama. Does distributionally robust supervised learning
180 give robust classifiers? In *International Conference on Machine Learning*, pages 2029–2037.
181 PMLR, 2018.
- 182 [10] J. Johnson, M. Douze, and H. Jégou. Billion-scale similarity search with gpus. *arXiv preprint*
183 *arXiv:1702.08734*, 2017.
- 184 [11] B. Kim, H. Kim, K. Kim, S. Kim, and J. Kim. Learning not to learn: Training deep neural
185 networks with biased data. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
186 *and Pattern Recognition*, pages 9012–9020, 2019.
- 187 [12] D. Krueger, E. Caballero, J.-H. Jacobsen, A. Zhang, J. Binas, D. Zhang, R. L. Priol, and
188 A. Courville. Out-of-distribution generalization via risk extrapolation (rex). *arXiv preprint*
189 *arXiv:2003.00688*, 2020.
- 190 [13] A. K. Menon, A. S. Rawat, and S. Kumar. Overparameterisation and worst-case generalisation:
191 friend or foe? In *International Conference on Learning Representations*, 2020.
- 192 [14] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
193 N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning
194 library. *Advances in neural information processing systems*, 32:8026–8037, 2019.
- 195 [15] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
196 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
197 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
198 *Learning Research*, 12:2825–2830, 2011.
- 199 [16] A. Raghu, M. Raghu, S. Bengio, and O. Vinyals. Rapid learning or feature reuse? towards
200 understanding the effectiveness of maml. *arXiv preprint arXiv:1909.09157*, 2019.
- 201 [17] M. Raghu, J. Gilmer, J. Yosinski, and J. Sohl-Dickstein. Svcca: Singular vector canonical corre-
202 lation analysis for deep learning dynamics and interpretability. *arXiv preprint arXiv:1706.05806*,
203 2017.
- 204 [18] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy,
205 A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International*
206 *journal of computer vision*, 115(3):211–252, 2015.

- 207 [19] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Distributionally robust neural networks
 208 for group shifts: On the importance of regularization for worst-case generalization. *arXiv*
 209 *preprint arXiv:1911.08731*, 2019.
- 210 [20] M. Srivastava, T. Hashimoto, and P. Liang. Robustness to spurious correlations via human
 211 annotations. In *International Conference on Machine Learning*, pages 9109–9119. PMLR,
 212 2020.
- 213 [21] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011
 214 dataset. 2011.
- 215 [22] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba. Places: A 10 million image
 216 database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*,
 217 40(6):1452–1464, 2017.

218 A Experimental Details

219 A.1 Model Training

220 We train (where not otherwise indicated) using SGD with learning rate of 1e-4, momentum of 0.9,
 221 L2 weight regularization of 1e-4, and batch size 128. We used early stopping with patience 20
 222 for all models, using an early stopping metric of worst-group accuracy on the validation set for all
 223 models except those for the coloured lines in Fig. 3, which use reweighted validation loss as an early
 224 stopping metric (i.e. validation loss where the subgroup losses are reweighted to match the training
 225 distribution’s subgroup distribution). We used $K = 5$ for GDRO and $\lambda = 3$ for IRM. We train all
 226 models in Pytorch [14] using their Imagenet-pretrained initialization [18].

227 A.2 Algorithms

228 We define here a number of different supervised classification algorithms of interest. We let ℓ be
 229 example-wise cross-entropy. Let the number of examples in a group g be n_g , and let the following
 230 shorthand describe the average loss for a group g : $\ell_g(f) = \frac{1}{n_g} \sum_{i=1}^n \ell(f(x_i), y_i) \mathbb{1}\{g_i = g\}$.

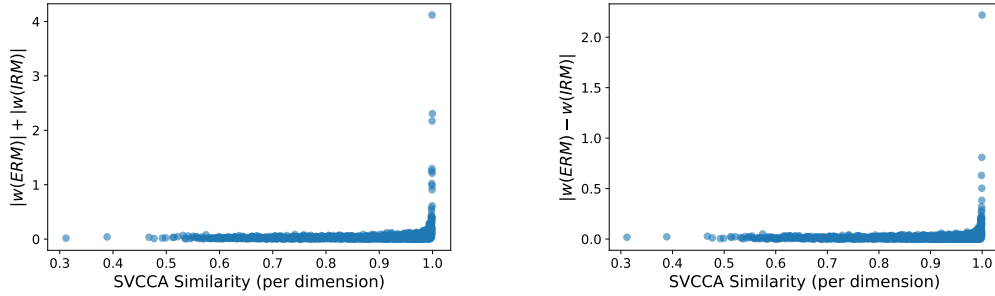
231 **Expected Risk Minimization** Expected risk minimization (ERM) is the usual paradigm for
 232 training ML models. In ERM, we minimize the mean loss ℓ on the training set: $\mathcal{L}_{ERM}(f) =$
 233 $\frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$.

234 **Group DRO.** Group Distributionally Robust Optimization (GDRO) [19] was proposed specifically
 235 for subgroup-based learning, and takes advantage of the subgroup information g . This loss aims
 236 to ensure that no group’s loss is that bad, and the group adjustment term (with hyperparameter
 237 K) ensures greater focus on smaller groups, which may otherwise be ignored: $\mathcal{L}_{GDRO}(f) =$
 238 $\max_{g \in \{1 \dots G\}} \left\{ \ell_g(f) + \frac{K}{\sqrt{n_g}} \right\}$.

239 **IRM.** Invariant risk minimization (IRM) [1] is another method that uses subgroup information.
 240 The intuition for this is somewhat involved [1]; the overarching motivation is that each environment
 241 should learn a representation such that the same predictive classifier is optimal across environments.
 242 The second term below is a gradient penalty on the output of f , w is a constant multiplier on the
 243 output of f , and λ is a hyperparameter: $\mathcal{L}_{IRM}(f) = \sum_{g=1}^G \ell_g(f) + \lambda \|\nabla_{w|w=1} \ell_g(w \cdot f)\|$

244 A.3 Post-hoc Training

245 We used the implementations of k-NN and k-Means implemented in the `faiss` package [10]. We
 246 used the implementation of logistic regression implemented in the `scikit-learn` package [15]
 247 with the “lbfgs” solver. For each method, we searched over 5 hyperparameter values and chose the
 248 best one (by worst-group error on the test set) for each value reported, as suggested by Gulrajani
 249 and Lopez-Paz [6]. For VQ and k-NN, we loop over the values of $k = 1, 2, 4, 8, 16$. For logistic
 250 regression, we loop over the L2-regularization value $C = 0.1, 0.2, 0.5, 1.0, 2.0$.



(a) The y-axis shows the sum of the absolute values in from the ERM and IRM transformed classifiers. We see that the only features which are important for classification in the two models are co-discovered.

(b) The y-axis shows the absolute difference in the transformed classifier between ERM and IRM models. We see the features which are treated highly differently between the two models are co-discovered.

Figure 4: Waterbirds feature co-discovery analysis on ERM- and IRM-learned representations. The x-axis shows the SVCCA similarity score for each dimension in both plots; a feature is co-discovered if it has a score near 1.

251 B Feature Co-Discovery in ERM and IRM

252 In Fig. 1, we show evidence of feature co-discovery between ERM and GDRO. Since GDRO and
 253 IRM have very similar representations (as shown in Fig. 2), we would expect a similar pattern to
 254 hold in IRM’s representations. For completeness, we show the analogous results for ERM and IRM
 255 in Fig. 4.