

Learning Less Generalizable Patterns for Better Test-Time Adaptation

Thomas Duboudin, Emmanuel Dellandréa, Corentin Abgrall, Gilles Hénaff, Liming Chen

► To cite this version:

Thomas Duboudin, Emmanuel Dellandréa, Corentin Abgrall, Gilles Hénaff, Liming Chen. Learning Less Generalizable Patterns for Better Test-Time Adaptation. 18th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISIGRAPP 2023), INSTICC, Feb 2023, Lisbonne, Portugal. 10.5220/0011893800003417. hal-03813534v3

HAL Id: hal-03813534 https://hal.science/hal-03813534v3

Submitted on 23 Feb 2023 $\,$

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Learning Less Generalizable Patterns for Better Test-Time Adaptation

Thomas Duboudin¹ Emmanuel Dellandréa¹ Corentin Abgrall² Gilles Hénaff² Liming Chen¹

Univ. Lyon, Ecole Centrale de Lyon, CNRS, INSA Lyon, Univ. Claude Bernard Lyon 1, Univ. Louis Lumière Lyon 2, LIRIS, UMR5205, 69134 Ecully, France

firstname.lastname@ec-lyon.fr

Abstract

Deep neural networks often fail to generalize outside of their training distribution, particularly when only a single data domain is available during training. While testtime adaptation has yielded encouraging results in this setting, we argue that to reach further improvements, these approaches should be combined with training procedure modifications aiming to learn a more diverse set of patterns. Indeed, test-time adaptation methods usually have to rely on a limited representation because of the shortcut learning phenomenon: only a subset of the available predictive patterns is learned with standard training. In this paper, we first show that the combined use of existing training-time strategies and test-time batch normalization, a simple adaptation method, does not always improve upon the test-time adaptation alone on the PACS benchmark. Furthermore, experiments on Office-Home show that very few trainingtime methods improve upon standard training, with or without test-time batch normalization. Therefore, we propose a novel approach that mitigates the shortcut learning behavior by having an additional classification branch learn less predictive and generalizable patterns. Our experiments show that our method improves upon the state-of-the-art results on both benchmarks and benefits the most to test-time batch normalization.

1. Introduction

Deep neural networks' performance falls sharply when confronted, at test-time, with data coming from a different distribution (or domain) than the training one. A change in lighting, sensor, weather conditions or geographical location can result in a dramatic performance drop [15, 2, 5]. Such environmental changes are commonly encountered when an embedded network is deployed in the wild and Thales LAS France, 78990 Élancourt, France firstname.lastname@fr.thalesgroup.com

exist in such diversity that it is impossible to gather enough data to cover all possible domain shifts. This lack of crossdomain robustness prevents the widespread deployment of deep networks in safety-critical applications. Domain generalization algorithms have been investigated to mitigate the test-time performance drop by modifying the training procedure. Contrary to the domain adaptation research field, no information about the target domain is assumed to be known in domain generalization. Most of the existing works assume to have access to data coming from several identified different domains and try to create a domain invariant representation by finding common predictive patterns [25, 26, 4, 24, 21, 18]. However, such an assumption is quite generous, and in many real-life applications, one does not have access to several data domains but only a single one. As a result, some works study single-source domain generalization [38, 32, 45, 44, 28]. However, most methods were found to perform only marginally better than the standard training procedure when the evaluation is done rigorously on several benchmarks [10, 43]. Another recent paradigm, called test-time adaptation, proposes to use a normally trained network and adapt it with a quick procedure at test-time, using only a batch of unlabeled target samples. This paradigm yielded promising results in the domain generalization setting [41, 40] because they alleviate the main challenges of domain generalization: the lack of information about the target domain and the requirement to be simultaneously robust in advance to every possible shift.

However, test-time adaptation methods suffer from a drawback that limits their adaptation capability, and which can only be corrected at training-time. Indeed, using a standard training procedure, only a subset of predictive patterns is learned, corresponding to the most obvious and efficient ones, while the less predictive patterns are disregarded entirely [33, 14, 29, 13, 31, 2, 8]. This apparent flaw, named shortcut learning, originates from the gradient descent optimization [29] and prevents a test-time method from using all the available patterns. The combination of a training-time patterns diversity-seeking approach with a test-time adaptation method may thus lead to improved results. In this paper, we show that the combined use of test-time batch normalization with the state-of-the-art single-source domain generalization methods does not systematically yield increased results on the PACS benchmark [23] in the single-source setting, despite them being designed to seek normally ignored patterns. Similar experiments on Office-Home [35] yield a similar result, with only a few methods performing better than the standard training procedure.

We thus propose a new method, called L2GP, which encourages a network to learn new predictive patterns rather than exploiting and refining already learned ones. To find such patterns, we propose to look for predictive patterns that are less generalizable than the naturally learned ones. These less generalizable patterns match the ones normally ignored because of the simplicity bias of deep networks that promotes the learning of a representation with a high generalization capability [19, 7]. Our method requires two classifiers added to a features extractor. They are trained asymmetrically: one is trained normally (with the standard cross-entropy classification loss only), and the other with both a cross-entropy loss and an additional shortcut avoidance loss. This loss slightly encourages memorization rather than generalization by learning batch-specific patterns, i.e. patterns that lower the loss on the running batch but with a limited effect on the other batches of data. The features extractor is trained with respect to both classification branches.

To summarize, our contribution is threefold:

- To the best of our knowledge, we are the first to investigate the effect of training-time single-source methods on a test-time adaptation strategy. We show that it usually does not increase performance and can even have an adverse effect.
- We apply, for the first time, several state-of-the-art single-source domain generalization algorithms on the more challenging and rarely used Office-Home benchmark and showed that very few yield a robust cross-domain representation.
- We propose an original algorithm to learn a larger than usual subset of predictive features. We show that it yields results competitive or over the state-of-the-art with the combination of test-time batch normalization.



Figure 1. Schema of our bi-headed architecture. The naming convention is the same as the one used in algorithm 1.

2. Related Works

2.1. Single-Source Domain Generalization

Most domain generalization algorithms require several identified domains to enforce some level of distributional invariance. Because this is an unrealistic hypothesis in some situations (such as in healthcare or defense-related tasks), methods were developed to deal with a domain shift issue with only one single domain available during training. Some of them rely on a domain shift invariance hypothesis. A commonly used invariance hypothesis is the texture shift hypothesis. Indeed, many domain shifts are primarily textures shifts, and using style-transfer-based data augmentation will improve the generalization. It can be done explicitly by training a model on stylized images [38, 20] or implicitly in the internal representation of the network [44, 28]. Such methods are limited to situations where it is indeed a shift of the hypothesized nature that is encountered. Others wish to learn a larger set of predictive patterns to make the network more robust should one or several training-time predictive patterns be missing at test-time. Volpi et al. [36] and Zhang et al. [45] propose to incrementally add adversarial images crafted to maximize the classification error of the network to the training dataset. These images no longer contain the original obvious predictive patterns, which then forces the learning of new patterns. These strategies are inspired by adversarial training methods [17, 22] that were originally designed to improve adversarial robustness in deep networks. Wang et al. [38] used a similar approach in an online fashion, without the impractical ever-growing training dataset, and combined it with a style augmentation approach. Huang et al. [18] and Shi et al. [32] used a dropout-based [34] strategy to prevent the network from relying only on the most predictive patterns by muting the most useful channels or mitigating the texture bias. These methods were evaluated in the singlesource setting on several benchmarks, including the very common PACS dataset.

Algorithm 1: Learning Less Generalizable Patterns (L2GP)

- 1 Method specific hyper-parameters:
- 2 weight for the shortcut avoidance loss α
- 3 step size used for the gradient perturbation lr_+
- 4 Networks:
- 5 features extractor f, and its weights W (ResNet18 without its last linear layer)
- first classifier c_1 (single linear layer) 6
- 7 - second classifier c_2 (single linear layer)
- while training is not over do 8
- sample 2 batches of data $\{(x_i, y_i), i = 0...N 1\}, \{(\tilde{x}_i, \tilde{y}_i), i = 0...N 1\}$ 9
- calculate the cross-entropy loss \mathcal{L} on the first batch for both branches on the original weights W: 10
- 11
- $$\begin{split} \mathcal{L}(f,c_1) &= \frac{1}{N} \sum_i \mathcal{L}[c_1(f(W,x_i)),y_i] \\ \mathcal{L}(f,c_2) &= \frac{1}{N} \sum_i \mathcal{L}[c_2(f(W,x_i)),y_i] \end{split}$$
 12
- calculate the gradient of the cross-entropy loss \mathcal{L} w.r.t W on the first batch: 13
- 14
- $\nabla_W \mathcal{L} = \nabla_W \frac{1}{N} \sum_i \mathcal{L}[c_2(f(W, x_i)), y_i]$ add the perturbation to the running weight W, and track this addition in the computational graph: 15
- $W_+ = W + lr_+ \nabla_W \mathcal{L}$ 16
- calculate the shortcut avoidance loss on the second batch: 17
- 18
- $\begin{aligned} \mathcal{L}_{sa}(f,c_2) &= \frac{1}{N} \sum_i ||c_2(f(W,\tilde{x}_i)) c_2(f(W_+,\tilde{x}_i))||_1 \\ \text{update all networks to minimize } \mathcal{L}_{total}(f,c_1,c_2) &= \frac{1}{2} (\mathcal{L}(f,c_1) + \mathcal{L}(f,c_2)) + \alpha \mathcal{L}_{sa}(f,c_2) \end{aligned}$ 19
- 20 end

21 At test-time: use $c_1 \circ f$ (discard c_2) combined with test-time batch normalization

2.2. Test-Time Adaptation

Test-time adaption has emerged as a promising paradigm to deal with domain shifts. Waiting to gather information about the target domain, in the shape of an unlabeled batch of samples (or even a single sample), alleviates the main drawbacks of training-time domain generalization methods: the lack of information about the target domain, and the necessity to simultaneously adapt to all possible shifts. The simplest test-time adaptation strategy consists of replacing the training-time statistics in the batch normalization layers with the running test batch statistics. It is now a mandatory algorithm block for almost all methods [27, 3, 41, 16, 30]. This strategy was originally designed to deal with test-time image corruptions but proved to be efficient in a more general domain shift setting [41, 40]. In a situation where samples of a test batch cannot be assumed to come from the same distribution, workarounds requiring a single sample were developed by mixing test-time and training-time statistics [41, 40, 16, 30], or by using data augmentation [16]. Some solutions, such as the work of Yang *et al.* [40] or Wang et al. [37], further rely on test-time entropy minimization to remove inconsistent features from the prediction. Finally, Zhang et al. [42] quickly adapt a network to make consistent predictions between different augmentations of the same test sample. All these strategies rely on a model trained with the standard training procedure.

3. Method

Our approach requires two classification layers plugged after the same features extractor: one will be tasked with learning the patterns that are normally learned (as they are not necessarily spurious and, therefore, should not be systematically ignored), and the other the normally "hidden" ones. This lightweight modification of the standard architecture, illustrated in figure 1, is compatible with many networks and tasks. The primary branch, consisting in the features extractor and the primary classifier, is trained to minimize the usual cross-entropy loss (algo. 1, lines 11). The secondary one is trained to minimize the cross-entropy loss (algo. 1, line 12) alongside a novel shortcut avoidance loss. The complete procedure is available in algorithm 1.

If we are able to update a model in a direction that lowers the loss value on a certain batch of data, but does not produce a similar decrease on another batch of the same distribution, it means that the patterns learned are both predictive as they lower the loss and generalize poorly, *i.e.* they are less predictive. These are precisely the patterns we are looking for. Our shortcut avoidance loss follows this idea. We first compute a new set of weights for the secondary branch by applying a single cross-entropy gradient ascent step to the branch weights (algo. 1, lines 13-16). The gradient is computed on the original running batch, already used for the cross-entropy losses. We, then,

	without TTBN		with TTBN					
Method	Avg. Val. Acc.	Avg. Test Acc.	Avg. Val. Acc.	Avg. Test Acc.				
PACS dataset								
ERM	96.8 ± 0.4	52.0 ± 1.9	97.4 ± 0.3	66.1 ± 1.1				
RSC [18]	97.7 ± 0.4	54.3 ± 1.8	97.2 ± 0.2	58.7 ± 1.6				
InfoDrop [32]	96.6 ± 0.3	53.4 ± 2.0	95.9 ± 0.3	65.5 ± 1.0				
ADA [36]	96.9 ± 0.8	55.9 ± 2.9	96.6 ± 1.1	66.5 ± 1.2				
ME-ADA [45]	96.7 ± 1.3	54.7 ± 3.1	96.5 ± 0.9	66.7 ± 2.0				
EFDM [44]	96.9 ± 0.5	59.6 ± 2.3	97.5 ± 0.5	71.3 ± 1.0				
SagNet [28]	97.2 ± 0.7	57.9 ± 2.9	97.8 ± 0.7	62.4 ± 1.8				
L.t.D [38]	97.9 ± 1.0	59.9 ± 2.7	97.6 ± 0.7	66.3 ± 1.5				
Spectral Decoupling [29]	95.9 ± 0.4	52.9 ± 2.6	96.2 ± 0.7	66.7 ± 1.1				
L2GP (ours)	98.6 ± 0.2	56.1 ± 2.7	96.4 ± 0.3	71.3 ± 0.6				
Office-Home dataset								
ERM	82.0 ± 0.8	52.0 ± 0.8	81.6 ± 1.1	52.6 ± 0.6				
RSC [18]	80.9 ± 0.4	49.2 ± 0.7	80.2 ± 0.5	48.9 ± 0.7				
InfoDrop [32]	76.4 ± 0.8	45.9 ± 0.5	77.1 ± 0.7	46.4 ± 0.6				
ADA [36]	81.2 ± 2.6	50.4 ± 0.9	80.3 ± 2.0	50.0 ± 0.7				
ME-ADA [45]	78.9 ± 1.4	49.8 ± 0.6	81.4 ± 1.2	50.0 ± 0.7				
EFDM [44]	82.9 ± 0.5	52.8 ± 0.6	83.3 ± 1.0	53.3 ± 0.5				
SagNet [28]	81.5 ± 1.5	51.9 ± 0.7	81.1 ± 1.1	51.8 ± 0.9				
L.t.D [38]	81.0 ± 1.2	50.9 ± 0.7	81.7 ± 2.7	51.2 ± 0.8				
Spectral Decoupling [29]	83.8 ± 0.7	52.5 ± 0.5	82.5 ± 0.6	53.2 ± 0.3				
L2GP (ours)	84.0 ± 0.6	53.4 ± 0.6	83.8 ± 0.5	54.5 ± 0.3				

Table 1. Performances of our approach and comparison with the state-of-the-art.

compare the predictions of the secondary branch with the current weights and the computed altered weights (algo. 1, lines 17-18). This difference in predictions constitutes our shortcut avoidance loss.

Our approach requires the sampling of two batches of data simultaneously because the shortcut avoidance loss is computed on a batch of data different from the one used to compute the applied gradient. As the features learned in the applied gradient generalize from one batch of data to the other, the altered weights' predictions are a lot less accurate than the running weights' predictions (cross-entropy gradient ascent). As a result, these predictions differ greatly. By training the secondary branch to minimize the gap between both predictions, we are pushing the weights toward an area in which the applied gradient does not change the network's secondary output. This would mean that the patterns extracted for the second batch are different from the ones learned in the applied gradient. By adding the cross-entropy loss to the training procedure, we are driving the network to learn weights that are both predictive for the running classification batch but that have a low effect on the predictions of another batch and are, hence, less predictive. Note that the running network's weights are optimized with regard to both sides of the

shortcut avoidance loss. The addition of the gradient must thus be tracked in the computational graph. This is akin to the MAML [6] meta-learning framework in which the starting point of a few optimization steps is itself optimized.

During the evaluation, only the first classifier is used, and the secondary one can be discarded. Indeed, the first classifier uses every available feature at its disposal, including those learned by the secondary branch, while the secondary branch only favors less simple features. Furthermore, we use test-time batch normalization (abbreviated as TTBN). This method has been chosen because of its simplicity and its wide range of applicability. We do not use the usual exponential average training mean and standard deviation (computed during training) in the batch normalization layers. Instead, we first calculate the statistics on the running test batch and use them to update an exponential average of the test statistics, as in [27, 3, 41, 16, 30], before using this estimate to normalize the features. A correct target statistics approximation can be reached only if all samples encountered at test-time come from the same data distribution. This is a realistic scenario for applications like autonomous driving, in which the data distribution is not expected to change over the course of a few consecutive images. Several methods [41, 16] provide ways to circumvent this issue

	without	t TTBN	with TTBN							
Ablation	Avg. Val. Acc.	Avg. Test Acc.	Avg. Val. Acc.	Avg. Test Acc.						
PACS dataset										
Double branch only (A)	96.8 ± 0.6	53.4 ± 2.6	96.4 ± 0.3	67.4 ± 0.8						
Detached loss term (B)	97.5 ± 0.1	52.6 ± 2.3	97.3 ± 0.3	68.2 ± 1.3						
Secondary prediction branch (C)	98.0 ± 0.1	53.4 ± 2.8	96.9 ± 0.2	70.1 ± 0.4						
Single branch (D)	92.8 ± 1.1	46.4 ± 4.9	93.0 ± 0.9	51.2 ± 5.1						
Complete method	98.6 ± 0.2	56.1 ± 2.7	96.4 ± 0.3	71.3 ± 0.6						
Office-Home dataset										
Double branch only (A)	82.7 ± 0.4	52.8 ± 0.5	82.6 ± 0.3	53.5 ± 0.4						
Detached loss term (B)	83.5 ± 0.7	52.7 ± 0.6	82.3 ± 0.6	54.0 ± 0.6						
Secondary prediction branch (C)	81.3 ± 0.4	53.9 ± 0.7	83.8 ± 0.6	54.8 ± 0.5						
Single branch (D)	82.6 ± 0.7	53.7 ± 0.4	82.0 ± 0.5	54.3 ± 0.5						
Complete method	84.0 ± 0.6	53.4 ± 0.6	83.8 ± 0.5	54.5 ± 0.3						

Table 2. Ablation study.

if needed.

4. Experiments and results

4.1. Baselines for comparison and experimental setup

We compare our approach with the standard training procedure (expected risk minimization, abbreviated ERM), with several methods designed for single-source domain generalization [38, 44, 28, 36, 45, 32], with Spectral Decoupling [29], a method designed to reduce the shortcutlearning phenomenon in deep networks, and with RSC [18], and InfoDrop [32], that are domain generalization algorithms which do not explicitly require several training domains. These baselines were selected because they yield state-of-the-art results, are representative of the main ideas in the single-source domain generalization research community, and because they have a publicly available implementation. This was a necessity as the original works' results were given without any test-time adaptation, and trained models were not provided. Our experiments are conducted on the PACS (7 classes, 4 domains, around 10k images in total), and the Office-Home (65 classes, 4 domains, around 15k images in total) benchmarks. PACS has already been used in the single-source setting in several works, but not Office-Home.

For a classification task, using a ResNet [12], our architectural changes break down to adding a single fully connected layer after the average pooling layer, next to the original last classification layer. To avoid a target domain information leak, the models selected for the test are those with the best validation accuracy. Furthermore, we chose to use the same common hyper-parameters for all baselines to precisely measure the effect of the training procedure modifications rather than the influence of a perhaps better than usual hyper-parameter. This change of hyper-parameters and differences in the model selection process are responsible for some inconsistencies between the results reported in the original works and ours (such as with SagNet [28]: 61.9% average accuracy on PACS in the original work, 57.9 in our own). Further experimental details, including common hyper-parameters and hyper-parameters selected for our approach and the comparison baselines, are available in the supplementary material.

4.2. Results and analysis

Our main results are available in table 1. The reported results are the mean, over the 12 distinct pairs of training and test domains, of the averages and standard deviations, over 3 runs, of the validation and test accuracies. More details about the precise calculation process are given in the supplementary material. Used alongside test-time batch normalization, our method reaches a performance similar to that of EFDM [44] on the PACS datasets but exceeds it on the Office-Home datasets. When test-time batch normalization is not used, our method remains state-of-the-art on the Office-Home dataset but falls behind the style-transfer-based methods on the PACS dataset by a noticeable margin. Besides, our approach also benefits the accuracy on the validation sets.

We observe a completely different behavior between experiments on PACS and Office-Home. While all the existing methods improve upon the standard training procedure (ERM) on PACS, only EFDM, spectral decoupling [29], and our method yield better results on Office-Home. Likewise, while always positive, the effect of the test-time batch normalization is much more noticeable on PACS than on Office-Home. Furthermore, it is interesting to notice that



Figure 2. Mean absolute difference for ERM and our approach.

the performance gain due to the test-time batch normalization is highly dependant on the training-time method used. Indeed, the gain is the highest when our approach or ERM is used and only reaches a result closely similar to ERM or below in most of the other cases. We hypothesize that the domain shifts of the PACS datasets are mostly textures shifts, while they are not for the Office-Home datasets. This would explain why test-time batch normalization yields a large improvement on the PACS benchmark: the simple use of test-time statistics, that encode textures [3], is enough to significantly bridge the domain gap. It would also explain why the methods reaching the highest results [44, 28, 38] in the usual setting (without test-time batch normalization) are all style-transfer-based methods. As our approach is not related to style transfer in any way, we are able to reach a higher accuracy on Office-Home than other existing works. Regarding the effect of different training-time methods, we hypothesize that the magnitude of the gain is related to whether the method is really learning a more diverse set of patterns or rather only weighting differently patterns that would also be learned naturally. This would explain why several methods that improve upon ERM without test-time batch normalization only perform precisely as well once it is used. Style-transfer-based methods, for instance, essentially grant a higher importance to shape-based patterns rather than texture-based patterns but not necessarily learn new patterns.

We also conducted an extensive ablation study to

understand and demonstrate the necessity of our choices. As a sanity check, we first study the $\alpha = 0$ situation: a single features extractor on which two classification layers are plugged in, trained only with the cross-entropy on the same batch at each iteration for both branches (line A in the table 2). The differences in initialization of the classifiers may have an implicit ensembling effect, as in MIMO [11], which could lead to a better out-of-distribution generalization without the need for the shortcut avoidance loss. This experiment yields a small increase of performance on both benchmarks, but it remains far below our approach, whose gain is, therefore, not coming from an implicit ensembling mechanism. We also study the effect of detaching from the computational graph the $c_2(f(W, \tilde{x}_i))$ term (not optimizing the features extractor with respect to this part of the loss) in the shortcut avoidance loss (line B), as this could lead to a substantial improvement in memory consumption, and as the simultaneous optimization on both terms in not needed per se to decrease the generalization ability of the network. This experiment shows a decreased performance as well. The detachment most likely only results in a slower learning as the constraint's gradient pushes in the reverse direction of the classification loss gradient. This behavior is prevented when the features extractor is optimized with regard to both terms of the regularization: pushing in the reverse direction of the classification gradient will only slide the difference in the parameter space but not shorten the gap. Then, to show that the performance gain is effectively linked to a mitigation of the shortcut learning phenomenon, we conduct two experiments. Firstly, we study the impact of using the secondary prediction branch at test-time rather than the primary one (line C). This experiment results in performances fairly similar to the first branch, only lower in validation. This was to be expected as the secondary branch is precisely trained so that it generalizes less on the training domain. Secondly, we study the effect of applying our shortcut avoidance loss on an architecture without the added secondary branch (line D). The shortcut avoidance loss is thus applied to the original classifier. The results show a dramatic drop in accuracy on the PACS dataset but not on the Office-Home dataset. This difference is most likely due to the higher diversity in Office-Home, which prevents the original patterns from being ignored.

To further show the effect of our loss, we track during training a measure of the diversity of the learned patterns for both our approach and ERM. Inspired by [1], we use the mean absolute difference (MAD) between normalized convolutional filters f (or neurons for fully connected layers) of a certain layer, computed over all layers L of size N_L and training domains D, for an epoch t, following the equation 1. The results are available in figure 2 and show a systematic increase in the diversity of the learned patterns for our approach compared to ERM, for both benchmarks. Finally, as the tuning of hyper-parameters in the domain generalization setting is a critical issue, we conduct a broad hyperparameters sensitivity analysis, available in the supplementary material in table 3. Our study shows a relatively low sensitivity and a large match between hyper-parameters fit for all training-test pairs of PACS and Office-Home.

$$MAD(t) = \sum_{D} \sum_{L} \frac{1}{N_L^2} \sum_{i,j} ||f_{t,D,L,i} - f_{t,D,L,j}||_1$$
(1)

5. Conclusion

In this paper, we investigated the behavior of different single-source methods when used in conjunction with testtime batch normalization on the PACS and Office-Home benchmarks. We showed that test-time batch normalization always has a positive, yet highly variable, influence and that, most of the time, the addition of a training-time method is superfluous. We hypothesized that this lack of additional performance was linked to the selection behavior of some algorithms, which still learn the same subset of patterns as the standard training, but weigh them differently. We thus proposed a novel approach learning normally "hidden" patterns by looking for predictive patterns that generalize less. We showed that it yielded state-of-the-art results on both benchmarks and benefits the most to test-time batch normalization. Future work will be dedicated to a better understanding of the origin of this test-time batch normalization variability and to experiments with our method on the DomainBed [10] benchmark.

Acknowledgement

This work was in part supported by the 4D Vision project funded by the Partner University Fund (PUF), a FACE program, as well as the French Research Agency, l'Agence Nationale de Recherche (ANR), through the projects Learn Real (ANR-18-CHR3-0002-01), Chiron (ANR-20-IADJ-0001-01), Aristotle (ANR-21-FAI1-0009-01), and the joint support of the French national program of investment of the future and the regions through the PSPC FAIR Waste project.

References

- Babajide O. Ayinde, Tamer Inanc, and Jacek M. Zurada. Regularizing deep neural networks by enhancing diversity in feature extraction. *IEEE Transactions on Neural Networks* and Learning Systems, 2019.
- [2] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *IEEE/CVF European conference on computer vision*, 2018.
- [3] Philipp Benz, Chaoning Zhang, Adil Karjauv, and In So Kweon. Revisiting batch normalization for improving corruption robustness. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021.
- [4] Fabio M Carlucci, Antonio D'Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. Domain generalization by solving jigsaw puzzles. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [5] Alex J DeGrave, Joseph D Janizek, and Su-In Lee. Ai for radiographic covid-19 detection selects shortcuts over signal. *Nature Machine Intelligence*, 2021.
- [6] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Modelagnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, 2017.
- [7] Tomer Galanti and Tomaso Poggio. Sgd noise and implicit low-rank bias in deep neural networks. arXiv preprint arXiv:2206.05794, 2022.
- [8] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2020.
- [9] Tejas Gokhale, Rushil Anirudh, Jayaraman J Thiagarajan, Bhavya Kailkhura, Chitta Baral, and Yezhou Yang. Improving diversity with adversarially learned transformations for domain generalization. arXiv preprint arXiv:2206.07736, 2022.
- [10] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International Conference on Learning Representations*, 2021.
- [11] Marton Havasi, Rodolphe Jenatton, Stanislav Fort, Jeremiah Zhe Liu, Jasper Snoek, Balaji Lakshminarayanan, Andrew Mingbo Dai, and Dustin Tran. Training independent subnetworks for robust prediction. In *International Conference on Learning Representations*, 2021.

- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016.
- [13] Katherine Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture bias in convolutional neural networks. *Advances in Neural Information Processing Systems*, 2020.
- [14] Katherine Hermann and Andrew Lampinen. What shapes feature representations? exploring datasets, architectures, and training. Advances in Neural Information Processing Systems, 2020.
- [15] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International Conference on Machine Learning*, 2018.
- [16] Xuefeng Hu, Gokhan Uzunbas, Sirius Chen, Rui Wang, Ashish Shah, Ram Nevatia, and Ser-Nam Lim. Mixnorm: Test-time adaptation through online normalization estimation. arXiv preprint arXiv:2110.11478, 2021.
- [17] Ruitong Huang, Bing Xu, Dale Schuurmans, and Csaba Szepesvári. Learning with a strong adversary. arXiv preprint arXiv:1511.03034, 2015.
- [18] Zeyi Huang, Haohan Wang, Eric P. Xing, and Dong Huang. Self-challenging improves cross-domain generalization. In IEEE/CVF European Conference on Computer Vision, 2020.
- [19] Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. The lowrank simplicity bias in deep networks. arXiv preprint arXiv:2103.10427, 2021.
- [20] Philip TG Jackson, Amir Atapour Abarghouei, Stephen Bonner, Toby P Breckon, and Boguslaw Obara. Style augmentation: data augmentation via style randomization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [21] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). arXiv preprint arXiv:2003.00688, 2020.
- [22] Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *International Conference on Learning Representations*, 2016.
- [23] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M. Hospedales. Deeper, broader and artier domain generalization. In *IEEE/CVF International Conference on Computer Vision*, 2017.
- [24] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Learning to generalize: Meta-learning for domain generalization. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [25] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.
- [26] Daniel Moyer, Shuyang Gao, Rob Brekelmans, Aram Galstyan, and Greg Ver Steeg. Invariant representations without adversarial training. In Advances in Neural Information Processing Systems, 2018.

- [27] Zachary Nado, Shreyas Padhy, D Sculley, Alexander D'Amour, Balaji Lakshminarayanan, and Jasper Snoek. Evaluating prediction-time batch normalization for robustness under covariate shift. arXiv preprint arXiv:2006.10963, 2020.
- [28] Hyeonseob Nam, HyunJae Lee, Jongchan Park, Wonjun Yoon, and Donggeun Yoo. Reducing domain gap by reducing style bias. In *IEEE/CVF Conference on Computer Vision* and Pattern Recognition, 2021.
- [29] Mohammad Pezeshki, Sékou-Oumar Kaba, Yoshua Bengio, Aaron Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. arXiv preprint arXiv:2011.09468, 2020.
- [30] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. Improving robustness against common corruptions by covariate shift adaptation. Advances in Neural Information Processing Systems, 2020.
- [31] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. In *Advances in Neural Information Processing Systems*, 2020.
- [32] Baifeng Shi, Dinghuai Zhang, Qi Dai, Zhanxing Zhu, Yadong Mu, and Jingdong Wang. Informative dropout for robust representation learning: A shape-bias perspective. In *International Conference on Machine Learning*, 2020.
- [33] Sahil Singla, Besmira Nushi, Shital Shah, Ece Kamar, and Eric Horvitz. Understanding failures of deep networks via robust feature extraction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [34] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 2014.
- [35] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *IEEE/CVF Conference* on Computer Vision and Pattern Recognition, 2017.
- [36] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John C Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. In Advances in Neural Information Processing Systems, 2018.
- [37] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *International Conference on Learning Representations*, 2021.
- [38] Zijian Wang, Yadan Luo, Ruihong Qiu, Zi Huang, and Mahsa Baktashmotlagh. Learning to diversify for single domain generalization. In *IEEE/CVF International Conference* on Computer Vision, 2021.
- [39] Zhenlin Xu, Deyi Liu, Junlin Yang, Colin Raffel, and Marc Niethammer. Robust and generalizable visual representation learning via random convolutions. In *International Conference on Learning Representations*, 2021.
- [40] Tao Yang, Shenglong Zhou, Yuwang Wang, Yan Lu, and Nanning Zheng. Test-time batch normalization. arXiv preprint arXiv:2205.10210, 2022.

- [41] Fuming You, Jingjing Li, and Zhou Zhao. Test-time batch statistics calibration for covariate shift. *arXiv preprint arXiv:2110.04065*, 2021.
- [42] Marvin Mengxin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via adaptation and augmentation. In *NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications*, 2021.
- [43] Xingxuan Zhang, Linjun Zhou, Renzhe Xu, Peng Cui, Zheyan Shen, and Haoxin Liu. Nico++: Towards better benchmarking for domain generalization. arXiv preprint arXiv:2204.08040, 2022.
- [44] Yabin Zhang, Minghan Li, Ruihuang Li, Kui Jia, and Lei Zhang. Exact feature distribution matching for arbitrary style transfer and domain generalization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [45] Long Zhao, Ting Liu, Xi Peng, and Dimitris Metaxas. Maximum-entropy adversarial data augmentation for improved generalization and robustness. *Advances in Neural Information Processing Systems*, 2020.

Supplementary Material

Results details

The results were obtained as follows:

- For the 12 distinct pairs of training and test domains, we calculate the average and the standard deviation of the validation and test accuracies over 3 runs with different random seeds (because the effect of the network's initialization on the test accuracy is greater than usual in a test-time domain shift situation).
- The reported numbers are the non-weighted mean over all distinct pairs of the average accuracies per trainingtest pair previously computed ± the mean over all distinct pairs of the pairwise standard deviation (as we are interested in the randomness of the initialization rather than the variation of accuracies between training-test pairs).

Hyper-parameters details

Data: for all the methods and benchmarks, we use the data augmentation described in [18] (random resized crops, color jitter, random horizontal flips, random grayscale). For a particular domain used in training, 90% of the dataset is used for training and the remaining 10% for validation. The test set is obtained using another domain dataset entirely.

Common hyper-parameters: experiments were conducted with a ResNet18 [12] trained for 100 epochs, with the stochastic gradient descent, a learning rate of 1e - 3, a batch size of 64, a weight decay of 1e - 5, and a Nesterov momentum of 0.9. After 80 epochs, the learning rate is divided by 10. The exponential average momentum used in

the batch normalization layers at test-time is set to 0.1.

L2GP (ours): the gradient ascent learning rate lr_+ is set to 1.0 and the α weight for the shortcut avoidance loss to 1.0 as well, for all the experiments, that is, for all the training-test pairs on both the PACS and the Office-Home datasets. These hyper-parameters were first set arbitrarily to plausible values and then confirmed to be effective on the PACS benchmark by looking at target performance. They were finally reused as is on the Office-Home benchmark. This hyper-parameters selection strategy may seem sub-optimal but is, in fact, more and more used in domain generalization problems [9, 39]: a method requiring a new and careful hyper-parameters setting for each new dataset encountered is impractical, even more so when the target data distribution is unknown and cannot thus be used to help the setting.

Comparison baselines specifics hyper-parameters are detailed below. For the experiments on the PACS datasets, on which most of the baselines were tested, we use the same hyper-parameters as in the original works. For the Office-Home datasets, we used the hyper-parameters of the multi-source setting if available. If the methods did not have quantitative hyper-parameters, such as EFDM [44] with the choice of mixing-layers depths, we used the ones proposed for the PACS experiments for the ones on Office-Home. Likewise, if no rigorous hyper-parameters setting strategy was detailed in the original work, we used the PACS hyper-parameters for experiments on Office-Home. Finally, for the Spectral Decoupling work that was never evaluated on neither PACS nor Office-Home, we conducted a simple hyper-parameters search using a single training-test domains pair, and transferred them as is to the other pairs with the same training domain.

RSC: the percentage of channels (or spatial cross-channel locations) to be dropped is initialized at 30% and is increased every 10 epochs linearly to reach 90% for the last ten. Spatial cross-channel locations dropout and channel all-locations dropout are applied in a mutually exclusive way with the same probability. All samples in a batch are subject to dropout.

InfoDrop: half the layers are subjected to the info-dropout. The dropout rate is set to 1.5, the temperature to 0.1, the bandwidth to 1.0, and the radius to 3.

ADA: the number of adversarial gradient ascent steps is set to 25, and the learning rate for the adversarial gradient ascent steps is set to 50. The γ and η factors are respectively set to 10.0 and 50.0. Adversarial images are added to the training set every 10 epoch.

Avg. test Acc. on PACS - Avg. test Acc. on Office-Home										
$lr_+\downarrow/\alpha \rightarrow$	10^{-3}	10^{-2}	0.1	1.0	10.0	100.0				
10^{-3}	66.9 - 53.7	66.8 - 53.1	67.7 - 53.2	67.0 - 53.5	67.4 - 53.2	68.5 - 53.7				
10^{-2}	67.8 - 53.2	67.8 - 53.1	67.6 - 53.3	67.7 - 52.4	68.6 - 53.9	70.6 - 53.2				
0.1	67.8 - 53.0	67.5 - 53.2	67.4 - 53.3	69.5 - 53.8	71.3 - 54.7	69.2 - 51.9				
1.0	67.1 - 53.3	68.0 - 53.4	69.0 - 53.8	71.3 - 54.4	70.3 - 52.6	20.1 - 49.9				
10.0	67.8 - 52.9	67.2 - 53.4	67.4 - 53.3	66.0 - 53.9	54.4 - 51.9	15.0 - 5.2				
100.0	66.2 - 53.2	67.9 - 53.4	67.4 - 53.4	67.8 - 53.3	60.5 - 53.2	14.5 - 2.0				

Table 3. Broad hyper-parameters sensitivity analysis.

ME-ADA: The same hyper-parameters as the ones above are used.

EFDM: the EFDMix layers are inserted after the first 3 residual blocks in the ResNet architecture.

SagNet: The randomization stage and the adversarial weight of SagNets are fixed to 3 and 0.1 for all experiments, as in the original work. A gradient clipping to 0.1 is applied to the adversarial loss.

L.t.D: α_1 and α_2 weights for the additional losses were set to 1.0, β to 0.1, for all experiments.

Spectral Decoupling: the weight of the spectral decoupling constraint (an L2-norm on the network's output) is set to 0.001 for experiments on Office-Home Experiments, and to 0.01 for experiments on PACS.