
Adaptive Risk Minimization: A Meta-Learning Approach for Tackling Group Shift

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A fundamental assumption of most machine learning algorithms is that the training
2 and test data are drawn from the same underlying distribution. However, this as-
3 sumption is violated in almost all practical applications: machine learning systems
4 are regularly tested under *distribution shift*, due to temporal correlations, particular
5 end users, or other factors. In this work, we consider the setting where the training
6 data are structured into groups and test time shifts correspond to changes in the
7 group distribution. Prior work has approached this problem by attempting to be
8 robust to all possible test time distributions, which may degrade average perfor-
9 mance. In contrast, we propose to use ideas from meta-learning to learn models
10 that are *adaptable*, such that they can adapt to shift at test time using a batch of
11 unlabeled test points. We acquire such models by learning to adapt to training
12 batches sampled according to different distributions, which simulate structural
13 shifts that may occur at test time. Our primary contribution is to introduce the
14 framework of adaptive risk minimization (ARM), a formalization of this setting
15 that lends itself to meta-learning. We develop meta-learning methods for solving
16 the ARM problem, and compared to a variety of prior methods, these methods
17 provide substantial gains on image classification problems in the presence of shift.

18 1 Introduction

19 The standard assumption in empirical risk minimization (ERM) is that the data distribution at test time
20 will match the training distribution. When this assumption does not hold, the performance of standard
21 ERM methods typically deteriorates rapidly, and this setting is commonly referred to as distribution
22 or dataset *shift* [47, 31]. For instance, we can imagine a handwriting classification system that, after
23 training on a large database of past images, is deployed to specific end users. Some new users have
24 peculiarities in their handwriting, leading to shift in the input distribution. This test scenario must be
25 carefully considered when building machine learning systems for real world applications.

26 Algorithms for handling distribution shift have been studied under a number of frameworks [47].
27 Many of these frameworks aim for *zero shot generalization* to shift, which requires more restrictive
28 but realistic assumptions. For example, one popular assumption is that the training data are provided
29 in *groups* and that distributions at test time will represent either new group distributions or new groups
30 altogether. This assumption is used by, e.g., group distributionally robust optimization (DRO) [24, 53],
31 robust federated learning [44, 32], and domain generalization [4, 17]. Constructing training groups or
32 tasks in practice is generally accomplished by using meta-data, which exists for most commonly used
33 datasets. This assumption allows for more tractable optimization and still permits a wide range of
34 realistic shifts. However, achieving strong zero shot generalization in this setting is still a hard problem.
35 For example, DRO methods, which focus on achieving maximal worst case performance, can often
36 be overly pessimistic and learn models that do not perform well on the actual test distributions [24].

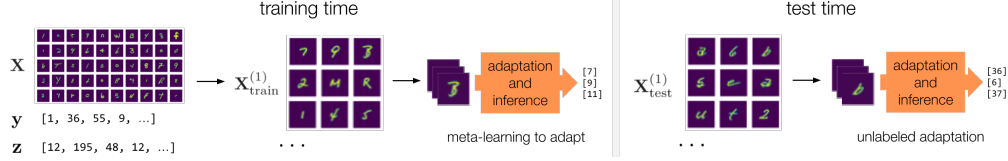


Figure 1: A schematic of the ARM problem setting and approach, described in detail in [Section 3](#). Left: During training, we assume access to labeled data along with group information z , which allows us to construct training distributions that exhibit group distribution shift. For example, a training distribution may place uniform mass on only a single user’s examples. We use these training distributions to learn a model that is adaptable to distribution shift via a form of meta-learning. We detail the specific adaptation procedures (orange box) that we consider in [Section 3](#) and [Figure 2](#). Right: We perform unsupervised adaptation to different test distributions, without requiring zero shot generalization to shift as in prior methods. If the test shifts we observe are similar to those simulated by the training distributions, e.g., we deploy the model to new end users at test time, then we expect that we can effectively adapt to these test distributions for better performance.

In this work, we take a different approach to combating group distribution shift by learning models that are able to deal with shift by *adapting* to the test time distribution. To do so, we assume that we can access a batch of *unlabeled* data points *at test time* – as opposed to individual isolated inputs – which can be used to implicitly infer the test distribution. This assumption is reasonable in many standard supervised learning setups. For example, we do not access single handwritten characters from an end user, but rather collections of characters such as sentences or paragraphs. When combined with the group assumption above, we arrive at a problem setting that is similar to the standard meta-learning setting [\[63\]](#). Meta-learning typically assumes that training data are grouped into tasks and new tasks are encountered at meta-test time, however these new tasks still include labeled examples for adaptation. As illustrated in [Figure 1](#), we instead aim to train a model that uses unlabeled data to adapt to the test distribution, thereby not requiring the model to generalize zero shot to all test distributions as in prior approaches.

The main contribution of this paper is to introduce the framework of adaptive risk minimization (ARM), in which models have the opportunity to adapt to the data distribution at test time based on unlabeled data points. We introduce an algorithm and instantiate a set of methods for solving ARM that, given a set of candidate distribution shifts, meta-learns a model that is adaptable to these shifts. One such method is based on meta-training a model such that simply updating batch normalization statistics [\[25\]](#) provides effective adaptation at test time, and we demonstrate that this simple approach can produce surprisingly strong results. Our experiments demonstrate that the proposed methods are able to outperform prior methods for handling distribution shift in image classification settings exhibiting group shift, including benchmarks for federated learning [\[7\]](#) and testing image classifier robustness [\[19\]](#).

2 Related Work

A number of prior works have studied distributional shift in various forms [\[47\]](#). In this section, we review prior work in robust optimization, meta-learning, and adaptation.

Robust optimization. DRO methods optimize machine learning systems to be robust to adversarial data distributions, thus optimizing for worst case performance against distribution shift [\[16, 2, 37, 12, 43, 11, 5\]](#). Recent work has shown that these algorithms can be utilized with deep neural networks, with additional care taken for regularization and model capacity [\[53\]](#). Unlike DRO methods, ARM methods do not require the model to generalize zero shot to all test time distribution shifts, but instead trains it to adapt to these shifts.

Also of particular interest are methods for robustness or adaptation to different users [\[21, 8, 27, 13, 35\]](#), a setting commonly referred to as robust or fair federated learning [\[41, 44, 32\]](#). Unlike these works, we consider the federated learning problem setting in which we do not assume access to any labels from any test users, as we partition users into disjoint train and test sets. We argue that this is a realistic setting for many practical machine learning systems – oftentimes, the only available information from the end user is an unlabeled batch of data.

74 **Meta-learning.** Meta-learning [55, 3, 61, 20] has been most extensively studied in the context of
 75 few shot supervised learning methods [54, 63, 48, 14, 58], i.e., *labeled* adaptation. Some other
 76 meta-learning methods adapt using both labeled and unlabeled data, either in the semi supervised
 77 learning setting [49, 67, 33] or the transductive learning setting [38, 1, 23]. These works do not
 78 focus on the same setting of distribution shift and all assume access to labeled data for adaptation.
 79 Prior works in meta-learning for unlabeled adaptation include Yu et al. [66], which adapts a policy
 80 to imitate human demonstrations in the context of robotic learning, and Metz et al. [42], which
 81 meta-learns an update rule for unsupervised representation learning, though they still require labels
 82 to learn a predictive model. We develop a new meta-learning framework for quickly adapting a
 83 predictive model using unlabeled examples, and unlike prior works, we focus on how meta-learning
 84 can address distribution shift.

85 **Adaptation to shift.** Unlabeled adaptation has primarily been studied separately from meta-learning.
 86 Domain adaptation is a prominent framework that assumes access to test examples at training time,
 87 similar to transductive learning [62], thus these methods can only handle a single predefined shift and
 88 do not constitute test time adaptation [10, 65]. Several methods for adaptation at test time have been
 89 developed specifically for dealing with label shift [52, 36, 59]. Other methods adapt using statistics
 90 of the test inputs [34] or optimize self-supervised surrogate losses [60, 64], and these methods have
 91 been shown to perform well across a number of image classification domains. We compare against
 92 these prior methods in [Section 4](#).

93 3 Adaptive Risk Minimization

94 In this section, we first formally describe the ARM problem setting, which builds on the settings
 95 used in prior work for tackling distribution shift. The novel aspect of the ARM setting is that it is
 96 amenable to meta-learning solutions to shift, and we demonstrate this by proposing an objective for
 97 the ARM setting that resembles typical meta-learning objectives. The problem setting and objective
 98 together constitute the ARM problem formulation. We subsequently propose a general algorithm as
 99 well as specific meta-learning approaches for solving the ARM problem.

100 3.1 The ARM problem setting

101 A key goal in machine learning is to develop methods that can go beyond the standard ERM setting
 102 and generalize in the face of distribution shift. Accomplishing this goal necessitates the use of
 103 additional assumptions beyond ERM, and we wish to carefully craft these assumptions such that they
 104 fulfill two properties: they are realistic and applicable to real world problems, and they allow for
 105 powerful and tractable methods. In this work, we choose two assumptions that are well established
 106 in the literature on distribution shift, in order to fulfill the first property, and we develop a novel
 107 meta-learning framework using these assumptions, thus fulfilling the second.

108 The first assumption is that the training data are provided in groups, which, as discussed above, mirrors
 109 analogous assumptions made in group DRO [24], federated learning [41], and meta-learning [63],
 110 among other settings. The second assumption is that we observe batches of test points all together,
 111 rather than one point at a time. Assuming access to multiple test points has been standard in domain
 112 adaptation [10, 65], which makes this assumption for training, as well as recent works studying
 113 test time adaptation [34, 60, 64]. To our knowledge, these assumptions have not been considered
 114 simultaneously in prior work. However, as we detail in this section, it is their conjunction that allows
 115 us to develop meta-learning solutions to shift.

116 In the ARM problem setting, we assume access to a training dataset that consists of N labeled data
 117 points $(\mathbf{x}^{(i)}, y^{(i)}, z^{(i)})$ sampled i.i.d. from the training distribution p . As noted, this differs from
 118 standard supervised learning in that we additionally observe the group $z^{(i)}$ associated with each point,
 119 which is a discrete value $z \in \{1, \dots, S\}$ that can represent tasks, users, or other types of meta-data.
 120 The goal is to learn a model $g(\cdot; \theta) : \mathcal{X} \rightarrow \mathcal{Y}$ that is parameterized by $\theta \in \Theta$ and predicts the output
 121 $y \in \mathcal{Y}$ given the input $\mathbf{x} \in \mathcal{X}$. At test time, we are given batches of K *unlabeled* data points, where
 122 each batch is drawn from a distribution that may differ from both p and the other batch distributions,
 123 and we do not observe either y or z . For example, we can imagine a test scenario that separately
 124 considers each user’s images, as discussed in [Section 1](#).

3.2 Deriving the ARM objective

We approach the goal of learning adaptable models through the lens of meta-learning. In particular, we define an *adaptation model* as a function $h(\cdot, \cdot; \phi) : \Theta \times \mathcal{X}^K \rightarrow \Theta$, which is parameterized by ϕ . h takes as input the model parameters θ and K unlabeled data points and produces updated parameters θ' after adapting using the K points. h can be initialized as a variety of different adaptation procedures, and we defer this discussion to [subsection 3.3](#). Our goal is to meta-learn ϕ and θ such that h can adapt g using *unlabeled* training data sampled according to a particular group z . Assuming that we will observe batches of data at test time that exhibit a similar type of shift, we can then perform the same procedure for better test performance. This motivates the ARM objective, given by

$$\min_{\theta, \phi} \mathbb{E}_{p_z} \left[\mathbb{E}_{p_{\mathbf{x}_y|z}} \left[\frac{1}{K} \sum_{k=1}^K \ell(g(\mathbf{x}_k; \theta'), y_k) \right] \right], \text{ where } \theta' = h(\theta, \mathbf{x}_1, \dots, \mathbf{x}_K; \phi). \quad (1)$$

A priori, we do not know what $p_{\text{test}}(z)$ will be, i.e., which values of z will be seen at test time. Thus, we draw inspiration from prior work in deep learning that demonstrates that uniformly sampling over a quantity of interest, such as labels or groups, is a strong method for achieving robustness with respect to that quantity [\[57, 6, 53\]](#). We extend this approach to the ARM setting by defining $p(z)$ at training time to place uniform probability mass on each group in the training set, in order to represent all training groups equally.

Standard few shot meta-learning formulations must use disjoint data batches for adaptation and meta-training to avoid label memorization [\[63\]](#). Since labels are not used during adaptation in ARM, we meta-train the adapted model using the same K examples that are used for adaptation. The labels for these examples are used in the meta-training update but not the adaptation itself. Thus, the adaptation matches the ARM setting at meta-test time, in which h adapts the model on the same unlabeled test points that the adapted model then predicts on.

3.3 Optimizing the ARM objective

[Algorithm 1](#) presents a general meta-learning approach for optimizing the ARM objective. In line 5, h outputs updated parameters θ' using an unlabeled batch of data. We assume that h is differentiable with respect to θ and ϕ , thus we can meta-train both θ and ϕ for *post adaptation* performance on a mini batch of data sampled according to a particular group z (line 6). This adaptation is performed using unlabeled data, mimicking the test time procedure detailed in lines 7-8. In practice, we sample mini batches of groups rather than just one group (line 3), to provide a better gradient signal for optimizing ϕ and θ .

We propose two approaches for instantiating the model g and adaptation procedure h in [Algorithm 1](#), which we summarize here and provide full details for in [Appendix B](#). First, we consider a *contextual* approach, shown in [Figure 2](#) (left), in which h summarizes the inputs $\mathbf{x}_1, \dots, \mathbf{x}_K$ into a context \mathbf{c} , which is then used by g as an additional input for predicting on each test point. In this setup, h can learn to provide useful information about the entire batch of K unlabeled data points to g for predicting the correct outputs. In the ARM setup, g is only ever evaluated after adaptation, i.e., with θ' . We can view h as outputting a concatenation of the model parameters and the context $\theta' = [\theta, \mathbf{c}]$.

This approach is inspired by recent works in contextual meta-learning with deep neural networks [\[15, 50\]](#). In line with these works, we propose an ARM-CML implementation of this approach which meta-learns a *context network* $f_{\text{cont}}(\cdot; \phi) : \mathcal{X} \rightarrow \mathbb{R}^D$. Note that f_{cont} is parameterized by ϕ , the parameters of h , as in this method, h has no additional parameters. f_{cont} processes each example \mathbf{x}_k in the mini batch separately to produce $\mathbf{c}_k \in \mathbb{R}^D$ for $k = 1, \dots, K$, where D is a hyperparameter. The average $\mathbf{c} = \frac{1}{K} \sum_{k=1}^K \mathbf{c}_k$ is then used as the context.

Algorithm 1 Meta-Learning for ARM

// Training procedure

Require: # training steps T , batch size K , learning rate η

- 1: **Initialize:** θ, ϕ
- 2: **for** $t = 1, \dots, T$ **do**
- 3: Sample z uniformly from training groups
- 4: Sample $(\mathbf{x}_k, y_k) \sim p(\cdot, \cdot | z)$ for $k = 1, \dots, K$
- 5: $\theta' \leftarrow h(\theta, \mathbf{x}_1, \dots, \mathbf{x}_K; \phi)$
- 6: $(\theta, \phi) \leftarrow (\theta, \phi) - \eta \nabla_{(\theta, \phi)} \sum_{k=1}^K \ell(g(\mathbf{x}_k; \theta'), y_k)$

// Test time adaptation procedure

Require: θ, ϕ , test batch $\mathbf{x}_1, \dots, \mathbf{x}_K$

- 7: $\theta' \leftarrow h(\theta, \mathbf{x}_1, \dots, \mathbf{x}_K; \phi)$
 - 8: $\hat{y}_k \leftarrow g(\mathbf{x}_k; \theta')$ for $k = 1, \dots, K$
-

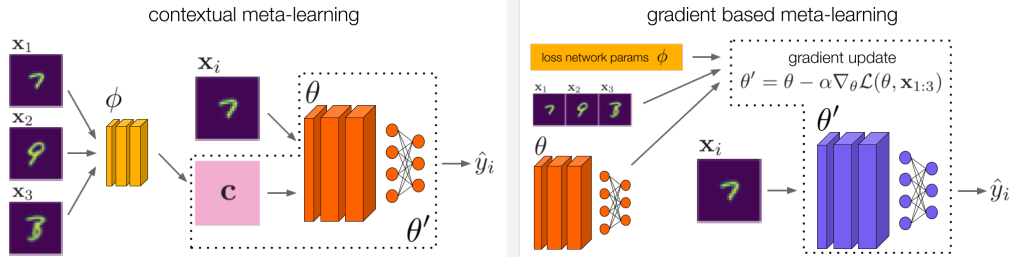


Figure 2: Schematics of the two broad classes of approaches we consider. Left: In the contextual approach, $\mathbf{x}_1, \dots, \mathbf{x}_K$ are summarized into a context \mathbf{c} , and we propose two methods for this summarization, either through a separate context network or using batch normalization activations in the model itself. \mathbf{c} can then be used by the model to infer additional information about the input distribution. Right: In the gradient based approach, an unlabeled loss function \mathcal{L} is used for gradient updates to the model parameters, in order to produce parameters that are specialized to the test inputs and can produce more accurate predictions.

Prior works outside of meta-learning have also investigated ways of conditioning predictions on a batch of data. One prominent technique, assuming that the model g is parameterized by a deep neural network with batch normalization layers [25], is to compute the normalization statistics for these layers using batches of test inputs, rather than the standard test time procedure of using the running statistics computed over the course of training. Several works have demonstrated the empirical effectiveness of this simple strategy, e.g., Li et al. [34], Kaku et al. [28], Nado et al. [45], Schneider et al. [56]. One advantage of this method’s simplicity is that it is easy to translate into the ARM setting, in order to arrive at a meta-learning version of this method which we call ARM-BN. The key difference between ARM-BN and BN is that, in ARM-BN, the model is trained to adapt using batches of training points sampled from the same group, following Algorithm 1. We can interpret this method through the contextual approach described above: if we view the running statistics used by standard BN as learned parameters of the model, then h replaces these parameters with statistics computed on the batch of inputs, which then serves as the context \mathbf{c} . In ARM-BN, the model is meta-trained to make effective use of this adaptation procedure, thus leading to more effective adaptation at test time. We provide complete details on ARM-BN in Appendix B.

As shown in Figure 2 (right), a distinct approach draws inspiration from gradient based meta-learning, where the goal is to learn parameters θ that are amenable to gradient updates on a loss function in order to quickly adapt to a new problem [14]. In other words, h produces $\theta' = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathbf{x}_1, \dots, \mathbf{x}_K)$, where α is a hyperparameter. Note that the loss function \mathcal{L} used in the gradient updates may be different from the original supervised loss function ℓ . In particular, in the setting of unlabeled adaptation, \mathcal{L} must be defined such that it operates on only the inputs \mathbf{x} , rather than the input output pairs that ℓ receives. Akin to Yu et al. [66], we propose a learned loss (ARM-LL) method that learns to modulate the output features of the model g . In our implementation, we assume that g produces output features $\mathbf{o} \in \mathbb{R}^{|\mathcal{Y}|}$ that are used as logits when making predictions. With this assumption, we define \mathcal{L} to be the composition of g and a loss network $f_{\text{loss}}(\cdot; \phi) : \mathbb{R}^{|\mathcal{Y}|} \rightarrow \mathbb{R}$, which takes in the output features from g and produces a scalar. Note that, similar to the CML method, f_{loss} is parameterized by ϕ , as h has no additional parameters. The ℓ_2 -norm of these scalars across the batch of test inputs is used as the loss for updating the model parameter θ . In other words,

$$\mathcal{L}(\theta, \mathbf{x}_1, \dots, \mathbf{x}_K) = \|\mathbf{v}\|_2, \text{ where } \mathbf{v} = [f_{\text{loss}}(g(\mathbf{x}_1; \theta); \phi), \dots, f_{\text{loss}}(g(\mathbf{x}_K; \theta); \phi)].$$

Before adaptation, the output features \mathbf{o} from g need not be suitable logits for prediction, as g is not evaluated for predictive performance using the unadapted parameters θ . Instead, \mathbf{o} may represent, for example, general features of the input \mathbf{x} . These features can then be used by f_{loss} in order to provide a gradient signal that adapts g to output accurate logits.

4 Experiments

Our experiments are designed to answer the following questions:

- (1) Do methods for adaptive risk minimization learn models that can adapt to group shift?
- (2) How do these methods compare to prior methods for robustness and adaptation?
- (3) Can we loosen the assumptions of accessing groups, at training time, and batches, at test time?

4.1 Evaluation domains and protocol

We evaluate on four image classification benchmarks, which span a range of problem settings including federated learning and robustness, demonstrating the general applicability of the proposed methods. Experimental details are provided in full in [Appendix C](#).

Rotated MNIST. We study a modified version of MNIST where images are rotated in 10 degree increments, from 0 to 130 degrees. We use only 108 training data points for each of the 2 smallest groups (120 and 130 degrees), and 324 points for rotations 90 to 110, whereas the overall training set contains 32292 points. At test time, we generate images from the MNIST test set with a certain rotation, and we consider each method’s worst case and average accuracy across groups.

Federated Extended MNIST (FEMNIST). The extended MNIST (EMNIST) dataset [9] consists of images of handwritten uppercase and lowercase letters, in addition to digits. FEMNIST [7] is a version of EMNIST that associates each handwritten character with its user. We measure each method’s worst case and average accuracy across 35 test users, which are held out and thus *disjoint* from the training users.

Corrupted image datasets. We evaluate the proposed methods and all comparisons on modified versions of CIFAR-10-C and Tiny ImageNet-C [19], which augment the CIFAR-10 [30] and Tiny ImageNet datasets, respectively, with common image corruptions that vary in type and severity. We modify the protocol from Hendrycks and Dietterich [19] to fit into the ARM problem setting by using a set of 56 corruptions for the training data, and we define each corruption to be a group. We use a disjoint set of 22 corruptions for the test data, and we measure worst case and average accuracy across the test corruptions.

4.2 Comparisons and ablations

We compare the ARM methods against several prior methods designed for robustness and adaptation. We summarize the comparisons here and again provide additional details in [Appendix C](#).

Group robustness. Sagawa et al. [53] recently proposed a state-of-the-art method for group robustness, and we refer to this approach as distributionally robust neural networks (DRNN). Their work also evaluates a strong upweighting (UW) baseline that samples uniformly from each group, and so we also evaluate this approach in our experiments. Note that, for CIFAR-10-C and Tiny ImageNet-C, UW is equivalent to ERM, as the groups all have an equal number of data points.

Test time adaptation. We evaluate the general approach of using test batches to compute batch normalization (BN) statistics, which has been proposed in several prior works [34, 28, 45, 56]. We also compare to test time training (TTT) [60], which adapts the model at test time using a self-supervised rotation prediction loss. These methods have previously achieved strong results, even without meta-learning, due to their favorable inductive biases for tasks such as image classification [60].

Robustness methods assume access to training groups but not test batches, whereas adaptation methods assume the opposite. Thus, at a high level, we can view the comparisons to these two broad classes of methods as evaluating the importance of each of these assumptions. We also conduct experiments in [subsection 4.4](#) in which we test ARM methods under looser assumptions.

Ablations. We also include ablations of the ARM-CML and ARM-LL methods, which sample mini-batches of unlabeled examples uniformly from all groups, rather than sampling from a single group to induce distribution shift. These “context ablation” and “learned loss ablation” are similar to test time adaptation methods in that they do not assume access to training groups. However, these methods lack the inductive bias of BN and TTT, as they instead use learned context and loss networks. These ablations validate the importance of adapting to a specific group.

4.3 Quantitative evaluation and comparisons

In [Table 1](#), we summarize the results. From these results, we highlight several key takeaways:

ARM methods consistently improve robustness and performance. Across all of our experiments, ARM methods significantly increase both worst case and average accuracy compared to all other methods. ARM-BN in particular achieves the best performance on most domains, demonstrating the effectiveness of using meta-training to improve an already strong inductive bias that empirically

Method	MNIST		FEMNIST		CIFAR-10-C		Tiny ImageNet-C	
	WC	Avg	WC	Avg	WC	Avg	WC	Avg
ERM	74.3 \pm 1.7	93.6 \pm 0.4	62.9 \pm 1.9	80.1 \pm 0.9	49.6 \pm 0.1	69.8 \pm 0.4	19.3 \pm 0.5	41.4 \pm 0.2
UW*	80.2 \pm 0.1	94.8 \pm 0.2	61.8 \pm 0.9	80.1 \pm 0.3	—	—	—	—
DRNN	79.3 \pm 1.1	94.8 \pm 0.1	58.1 \pm 0.7	74.4 \pm 0.8	44.5 \pm 0.5	70.7 \pm 0.6	19.9 \pm 0.3	41.6 \pm 0.2
BN	75.1 \pm 0.2	93.9 \pm 0.1	66.9 \pm 0.8	81.1 \pm 0.3	62.5 \pm 0.2	79.4 \pm 0.3	23.9 \pm 0.2	42.8 \pm 0.2
TTT	81.1 \pm 0.3	95.4 \pm 0.1	64.1 \pm 0.2	83.4 \pm 0.1	66.6 \pm 0.6	75.6 \pm 0.8	19.7 \pm 0.4	41.4 \pm 0.3
CML ablation	78.2 \pm 0.6	94.2 \pm 0.1	64.4 \pm 0.7	81.5 \pm 0.7	47.8 \pm 0.1	68.2 \pm 0.1	19.6 \pm 0.4	42.3 \pm 0.2
LL ablation	82.4 \pm 0.3	94.8 \pm 0.2	61.9 \pm 0.2	79.3 \pm 0.6	61.5 \pm 0.2	68.3 \pm 0.5	25.8 \pm 0.4	41.7 \pm 0.1
ARM-CML	88.7 \pm 0.6	96.7 \pm 0.1	67.8 \pm 1.3	85.7 \pm 0.3	67.7 \pm 0.5	79.2 \pm 0.3	21.4 \pm 0.2	43.3 \pm 0.4
ARM-BN	82.8 \pm 0.4	95.3 \pm 0.1	72.6 \pm 0.3	85.7 \pm 0.1	71.1 \pm 0.1	80.9 \pm 0.2	27.7 \pm 0.2	44.9 \pm 0.2
ARM-LL	87.2 \pm 0.5	96.3 \pm 0.2	69.6 \pm 2.1	85.6 \pm 0.5	66.9 \pm 0.2	75.7 \pm 0.3	27.1 \pm 0.3	44.2 \pm 0.4

Table 1: Worst case (WC) and average (Avg) top 1 accuracy on rotated MNIST, FEMNIST, CIFAR-10-C, and Tiny ImageNet-C across all methods, where means and standard errors are reported across three separate runs of each method. ARM methods consistently achieve greater robustness, measured by WC, and Avg performance compared to prior methods. *The UW baseline is equivalent to ERM for CIFAR-10-C and Tiny ImageNet-C.

works well for image classification. ARM-CML and ARM-LL also generally improve upon the other methods for almost all metrics, and we suspect that these more expressive methods could perform better than ARM-BN for other modalities such as natural language and video [66].

Robustness methods suffer from pessimism, which hurts their performance. DRNN generally results in worse average case and, surprisingly, worst case performance, which we hypothesize may be due to optimization difficulties or overfitting to the training groups. In particular, methods such as DRNN were originally evaluated in settings where the training and test groups were semantically the same [53], whereas our FEMNIST setup tests on held out users and our CIFAR-10-C and Tiny ImageNet-C setups test on held out corruptions. Indeed, for FEMNIST, we also test q -FedAvg [32], a state-of-the-art method for fair federated learning. q -FedAvg was also originally evaluated with the same users at training and test time, and in our setup, this method also performs poorly, achieving 58.2 ± 1.0 worst case and 80.8 ± 0.3 average accuracy.

With an appropriate inductive bias, test time adaptation methods perform well. In our evaluation, one of the strongest prior methods is the simple BN method. This method performs well across all metrics, though it typically still lags behind ARM methods and ARM-BN in particular. As discussed above, we believe that this adaptation procedure performs well as it constitutes an inductive bias that is well suited for image classification. TTT offers additional support for this hypothesis: this method also works well across most metrics, in line with previous results on the original versions of the corrupted image benchmarks [60], but works notably well for rotated MNIST. This may be because the inductive bias associated with the auxiliary task of rotating images allows the classifier to specifically be more robust to rotation shift.

In summary, in our experiments, we observe poor performance from robustness methods, varying performance from adaptation methods, and the strongest performance from ARM methods. As ARM methods make the strongest assumptions, we next discuss how these assumptions may be loosened.

4.4 Loosening the training group and test batch assumption

We present a preliminary investigation of the feasibility and effectiveness of ARM methods without the training group and test batch assumptions. In these experiments, we focus on the rotated MNIST domain and ARM-CML, which achieved the best performance on this domain.

Unknown groups. In the case of unknown groups, one option is to use unsupervised learning to discover group structure in the training data. To test this option, we train a variational autoencoder (VAE) [29, 51] with discrete latent variables [26, 40] using the training dataset. Specifically, we define the latent variable, which we denote as c to differentiate from the group z , to be Categorical with 12 possible discrete values, which we purposefully choose to be smaller than the number of rotations. The VAE is not given any information about the ground truth z ; however, we encode the notion that c is independent of y by conditioning the decoder on the label. We use the inference network of the VAE to assign groups to the training data, and we run ARM-CML using these learned groups. In Table 2, we see that ARM-CML in this setting outperforms ERM and is competitive with TTT on the ground truth groups, which as discussed earlier encodes a strong inductive bias for solving this task. Figure 3 (left) visualizes samples from the VAE for different values of y and c .

Method	WC	Avg
ERM	74.3 ± 1.7	93.6 ± 0.4
TTT	81.1 ± 0.3	95.4 ± 0.1
ARM-CML	81.7 ± 0.3	95.2 ± 0.3

Table 2: Using learned groups, ARM-CML outperforms ERM and matches the performance of TTT on rotated MNIST. This result may be improved by techniques for learning more diverse groups for meta-training.

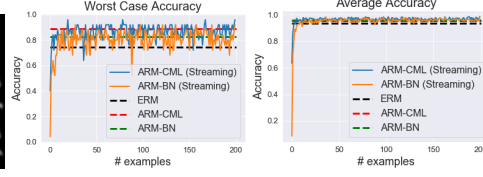


Figure 3: Left: Visualizing VAE samples conditioned on different values of y (x axis) and c (y axis). The VAE learns to use c to represent rotation information. Right: In the streaming setting, ARM methods reach strong performance after fewer than 10 data points, despite meta-training with batch sizes of 50.

This result suggests that, when group information is not provided, a viable approach is to learn groups for ARM methods. Discovering disentangled factors of variation without supervision is, in the most general sense, an impossible problem [39]. However, when combined with meta-learning, the learned groups need not perfectly reflect the test time distributions; rather, the groups should cover many different distributions to allow for meta-training the model such that it can adapt to new test distributions. This advantage was noted by Hsu et al. [22], who show that even simple techniques such as overcomplete clustering can be effective for defining meta-training tasks. Incorporating techniques from this prior work is a promising direction for building on our results.

Streaming test points. When we cannot access a batch of test points all at once, we can augment the proposed ARM methods to be sequential. ARM-CML and ARM-BN can update their average context and normalization statistics, respectively, after observing each new test point, and ARM-LL can perform small gradient updates on each point. In Figure 3 (right), we visualize the performance of ARM methods, using models that were meta-trained with batch sizes of 50 but evaluated in this streaming setting. We see that both ARM-CML and ARM-BN are able to achieve their original average accuracy within observing a few data points, well before the training batch size of 50. We describe in detail how each ARM method can be applied to the streaming setting in Appendix B. Next, we qualitatively analyze why ARM-CML achieves better performance compared to ERM in the case of FEMNIST.

4.5 Qualitative analysis and observations

In Figure 4, we present an example of how ARM methods can improve test accuracy by adapting to specific users. We visualize a batch of 50 examples from a randomly sampled FEMNIST test user, and we highlight an ambiguous example. Models trained via ERM and ARM-CML, when only given a batch size of 2 as shown by the black dashed box, incorrectly classify this example as “2”. However, when given access to the entire batch of 50 images, which contain examples of class “2” and “a” from this user, the ARM-CML trained model successfully adapts this prediction to instead output “a”, which is the correct label. In general, we find that most examples of adaptation in FEMNIST occur for similarly ambiguous examples, e.g., “l” versus “I”, though not all examples were interpretable.

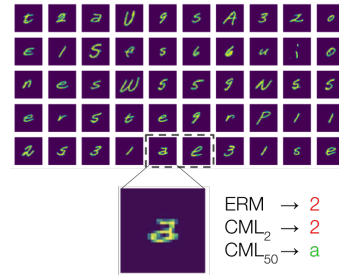


Figure 4: Visualizing one batch of 50 images from a FEMNIST test user. The ARM-CML method, using the entire batch, is able to successfully adapt to output the correct label “a” on the ambiguous example, shown enlarged, whereas ERM incorrectly outputs “2”.

5 Discussion and Future Work

We presented adaptive risk minimization (ARM), a problem formulation for learning models that can robustly adapt in the face of group distribution shift at test time using only a batch of unlabeled test examples. We devised an algorithm and a set of methods for optimizing the ARM objective that meta-learns models that are adaptable to different distributions of training data. Empirically, we observed that ARM methods consistently improve performance in terms of both average and worst case metrics, as compared to a number of prior approaches for handling shift. Two exciting directions for future work are to further explore the unknown groups setting, potentially drawing inspiration from Hsu et al. [22] as discussed, and to develop more sophisticated ARM approaches.

References

- [1] A. Antoniou and A. Storkey. Learning to learn via self-critique. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [2] A. Ben-Tal, D. den Hertog, A. De Waegenaere, B. Melenberg, and G. Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 2013.
- [3] S. Bengio, Y. Bengio, J. Cloutier, and J. Gecsei. On the optimization of a synaptic learning rule. In *Optimality in Artificial and Biological Neural Networks*, 1992.
- [4] G. Blanchard, G. Lee, and C. Scott. Generalizing from several related classification tasks to a new unlabeled sample. In *Advances in Neural Information Processing Systems (NIPS)*, 2011.
- [5] J. Blanchet, Y. Kang, and K. Murthy. Robust Wasserstein profile inference and applications to machine learning. *arXiv preprint arXiv:1610.05627*, 2016.
- [6] M. Buda, A. Maki, and M. Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 2018.
- [7] S. Caldas, S. Duddu, P. Wu, T. Li, J. Konečný, H. McMahan, V. Smith, and A. Talwalkar. LEAF: A benchmark for federated settings. In *Workshop on Federated Learning for Data Privacy and Confidentiality*, 2019.
- [8] F. Chen, M. Luo, Z. Dong, Z. Li, and X. He. Federated meta-learning with fast convergence and efficient communication. *arXiv preprint arXiv:1802.07876*, 2018.
- [9] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik. EMNIST: An extension of MNIST to handwritten letters. *arXiv preprint arXiv:1702.05373*, 2017.
- [10] G. Csurka. Domain adaptation for visual applications: A comprehensive survey. *arXiv preprint arXiv:1702.05374*, 2017.
- [11] J. Duchi, P. Glynn, and H. Namkoong. Statistics of robust optimization: A generalized empirical likelihood approach. *arXiv preprint arXiv:1610.03425*, 2016.
- [12] P. Esfahani and D. Kuhn. Data-driven distributionally robust optimization using the Wasserstein metric: Performance guarantees and tractable reformulations. *arXiv preprint arXiv:1505.05116*, 2015.
- [13] A. Fallah, A. Mokhtari, and A. Ozdaglar. Personalized federated learning: A meta-learning approach. *arXiv preprint arXiv:2002.07948*, 2020.
- [14] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, 2017.
- [15] M. Garnelo, D. Rosenbaum, C. Maddison, T. Ramalho, D. Saxton, M. Shanahan, Y. Teh, D. Rezende, and S. Eslami. Conditional neural processes. In *International Conference on Machine Learning (ICML)*, 2018.
- [16] A. Globerson and S. Roweis. Nightmare at test time: Robust learning by feature deletion. In *International Conference on Machine Learning (ICML)*, 2006.
- [17] I. Gulrajani and D. Lopez-Paz. In search of lost domain generalization. *arXiv preprint arXiv:2007.01434*, 2020.
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [19] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations (ICLR)*, 2019.
- [20] S. Hochreiter, A. Younger, and P. Conwell. Learning to learn using gradient descent. In *International Conference on Artificial Neural Networks (ICANN)*, 2001.

- [21] S. Horiguchi, S. Amano, M. Ogawa, and K. Aizawa. Personalized classifier for food image recognition. *IEEE Transactions on Multimedia*, 2018.
- [22] K. Hsu, S. Levine, and C. Finn. Unsupervised learning via meta-learning. In *International Conference on Learning Representations (ICLR)*, 2019.
- [23] S. Hu, P. Moreno, Y. Xiao, X. Shen, G. Obozinski, N. Lawrence, and A. Damianou. Empirical Bayes transductive meta-learning with synthetic gradients. In *International Conference on Learning Representations (ICLR)*, 2020.
- [24] W. Hu, G. Niu, I. Sato, and M. Sugiyama. Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning (ICML)*, 2018.
- [25] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, 2015.
- [26] E. Jang, S. Gu, and B. Poole. Categorical reparameterization with Gumbel-softmax. In *International Conference on Learning Representations (ICLR)*, 2017.
- [27] Y. Jiang, J. Konečný, K. Rush, and S. Kannan. Improving federated learning personalization via model agnostic meta learning. *arXiv preprint arXiv:1909.12488*, 2019.
- [28] A. Kaku, S. Mohan, A. Parnandi, H. Schambra, and C. Fernandez-Granda. Be like water: Robustness to extraneous variables via adaptive feature normalization. *arXiv preprint arXiv:2002.04019*, 2020.
- [29] D. Kingma and M. Welling. Auto-encoding variational Bayes. In *International Conference on Learning Representations (ICLR)*, 2014.
- [30] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [31] D. Lazer, R. Kennedy, G. King, and A. Vespignani. The parable of Google flu: Traps in big data analysis. *Science*, 2014.
- [32] T. Li, M. Sanjabi, A. Beirami, and V. Smith. Fair resource allocation in federated learning. In *International Conference on Learning Representations (ICLR)*, 2020.
- [33] X. Li, Q. Sun, Y. Liu, Q. Zhou, S. Zheng, T. Chua, and B. Schiele. Learning to self-train for semi-supervised few-shot classification. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [34] Y. Li, N. Wang, J. Shi, J. Liu, and X. Hou. Revisiting batch normalization for practical domain adaptation. In *International Conference on Learning Representations Workshop (ICLRW)*, 2017.
- [35] S. Lin, Y. Guang, and J. Zhang. Real-time edge intelligence in the making: A collaborative learning framework via federated meta-learning. *arXiv preprint arXiv:2001.03229*, 2020.
- [36] Z. Lipton, Y. Wang, and A. Smola. Detecting and correcting for label shift with black box predictors. In *International Conference on Machine Learning (ICML)*, 2018.
- [37] A. Liu and B. Ziebart. Robust classification under sample selection bias. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- [38] Y. Liu, J. Lee, M. Park, S. Kim, E. Yang, S. Hwang, and Y. Yang. Learning to propagate labels: Transductive propagation network for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2019.
- [39] F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Schölkopf, and O. Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *International Conference on Machine Learning (ICML)*, 2019.
- [40] C. Maddison, A. Mnih, and Y. Teh. The Concrete distribution: A continuous relaxation of discrete random variables. In *International Conference on Learning Representations (ICLR)*, 2017.

- [41] H. McMahan, E. Moore, D. Ramage, S. Hampson, and B. Arcas. Communication-efficient learning of deep networks from decentralized data. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.
- [42] L. Metz, N. Maheswaranathan, B. Cheung, and J. Sohl-Dickstein. Meta-learning update rules for unsupervised representation learning. In *International Conference on Learning Representations (ICLR)*, 2019.
- [43] T. Miyato, S. Maeda, M. Koyama, K. Nakae, and S. Ishii. Distributional smoothing with virtual adversarial training. *arXiv preprint arXiv:1507.00677*, 2015.
- [44] M. Mohri, G. Sivek, and A. Suresh. Agnostic federated learning. In *International Conference on Machine Learning (ICML)*, 2019.
- [45] Z. Nado, S. Padhy, D. Sculley, A. D’Amour, B. Lakshminarayanan, and J. Snoek. Evaluating prediction-time batch normalization for robustness under covariate shift. *arXiv preprint arXiv:2006.10963*, 2020.
- [46] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [47] J. Quiñero Candela, M. Sugiyama, A. Schwaighofer, and N. Lawrence. *Dataset Shift in Machine Learning*. The MIT Press, 2009.
- [48] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2017.
- [49] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. Tenenbaum, H. Larochelle, and R. Zemel. Meta-learning for semi-supervised few-shot classification. In *International Conference on Learning Representations (ICLR)*, 2018.
- [50] J. Requeima, J. Gordon, J. Bronskill, S. Nowozin, and R. Turner. Fast and flexible multi-task classification using conditional neural adaptive processes. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [51] D. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning (ICML)*, 2014.
- [52] A. Royer and C. Lampert. Classifier adaptation at prediction time. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [53] S. Sagawa, P. Koh, T. Hashimoto, and P. Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations (ICLR)*, 2020.
- [54] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap. Meta-learning with memory-augmented neural networks. In *International Conference on Machine Learning (ICML)*, 2016.
- [55] J. Schmidhuber. Evolutionary principles in self-referential learning. *Diploma thesis, Institut f. Informatik, Tech. Univ. Munich*, 1987.
- [56] S. Schneider, E. Rusak, L. Eck, O. Bringmann, W. Brendel, and M. Bethge. Improving robustness against common corruptions by covariate shift adaptation. *arXiv preprint arXiv:2006.16971*, 2020.
- [57] L. Shen, Z. Lin, and Q. Huang. Relay backpropagation for effective learning of deep convolutional neural networks. In *European Conference on Computer Vision (ECCV)*, 2016.
- [58] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.

- 481 [59] M. Sulc and J. Matas. Improving CNN classifiers by estimating test-time priors. In *IEEE*
482 *International Conference on Computer Vision (ICCV)*, 2019.
- 483 [60] Y. Sun, X. Wang, Z. Liu, J. Miller, A. Efros, and M. Hardt. Test-time training with self-
484 supervision for generalization under distribution shifts. In *International Conference on Machine*
485 *Learning (ICML)*, 2020.
- 486 [61] S. Thrun and L. Pratt. *Learning to Learn*. Springer Science & Business Media, 1998.
- 487 [62] V. Vapnik. *Statistical Learning Theory*. Wiley New York, 1998.
- 488 [63] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra. Matching networks for
489 one shot learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- 490 [64] D. Wang, E. Shelhamer, S. Liu, B. Olshausen, and T. Darrell. Fully test-time adaptation by
491 entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.
- 492 [65] G. Wilson and D. Cook. A survey of unsupervised deep domain adaptation. *ACM Transactions*
493 *on Intelligent Systems and Technology (TIST)*, 2020.
- 494 [66] T. Yu, C. Finn, A. Xie, S. Dasari, P. Abbeel, and S. Levine. One-shot imitation from observing
495 humans via domain-adaptive meta-learning. In *Robotics: Science and Systems (RSS)*, 2018.
- 496 [67] R. Zhang, T. Che, Z. Ghahramani, Y. Bengio, and Y. Song. MetaGAN: An adversarial approach
497 to few-shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.