# Improved Regularization and Robustness for Fine-tuning in Neural Networks

**Dongyue Li** Northeastern University li.dongyu@northeastern.edu Hongyang R. Zhang Northeastern University ho.zhang@northeastern.edu

### Abstract

A widely used algorithm for transfer learning is fine-tuning, where a pre-trained model is fine-tuned on a target task with a small amount of labeled data. When the capacity of the pre-trained model is much larger than the size of the target data set, fine-tuning is prone to overfitting and "memorizing" the training labels. Hence, an important question is to regularize fine-tuning and ensure its robustness to noise. To address this question, we begin by analyzing the generalization properties of fine-tuning. We present a PAC-Bayes generalization bound that depends on the distance traveled in each layer during fine-tuning and the noise stability of the fine-tuned model. We empirically measure these quantities. Based on the analysis, we propose *regularized self-labeling*—the interpolation between regularization and self-labeling methods, including (i) *layer-wise regularization* to constrain the distance traveled in each layer; (ii) self label-correction and label-reweighting to correct mislabeled data points (that the model is confident) and reweight less confident data points. We validate our approach on an extensive collection of image and text data sets using multiple pre-trained model architectures. Our approach improves baseline methods by 1.76% (on average) for seven image classification tasks and 0.75% for a few-shot classification task. When the target data set includes noisy labels, our approach outperforms baseline methods by 3.56% on average in two noisy settings.

# 1 Introduction

Learning from limited labeled data is a fundamental problem in many real-world applications (Ratner et al., 2016, 2017). A common approach to address this problem is fine-tuning a large model that has been pre-trained on publicly available labeled data (He et al., 2019). Since fine-tuning is typically applied to a target task with limited labels, this algorithm is prone to overfitting or "memorization" issues (Tan et al., 2018). These issues worsen when the target task contains noisy labels (Zhang et al., 2016). In this paper, we analyze regularization methods for fine-tuning from both theoretical and empirical perspectives. Based on the analysis, we propose a *regularized self-labeling* approach that improves the generalization and robustness properties of fine-tuning.

Previous works (Li et al., 2018a,b) have proposed regularization methods to constrain the distance between a fine-tuned model and the pre-trained model in the Euclidean norm. Li et al. (2020) provides extensive study to show that the performance of fine-tuning and the benefit of adding regularization depend on the hyperparameter choices. Salman et al. (2020) empirically find that performing adversarial training during the pre-training phase helps learn pre-trained models that transfer better to downstream tasks. The work of Gouk et al. (2021) generalizes the above ideas to various norm choices and finds that projected gradient descent methods perform well for implementing distance-based regularization. Additionally, they derive generalization bounds for fine-tuning using Rademacher complexity. These works focus on settings where there is no label noise in the target data set. When

35th Conference on Neural Information Processing Systems (NeurIPS 2021).

label noise is present, for example, due to applying weak supervision techniques (Ratner et al., 2016), an important question is to design methods that are robust to such noise. The problem of learning from noisy labels has a rich history of study in supervised learning (Natarajan et al., 2013). In contrast, little is known in the transfer learning setting. These considerations motivate us to analyze the generalization and robustness properties of fine-tuning.

In Section 4.1, we begin by conducting a PAC-Bayesian analysis of regularized fine-tuning. This is inspired by recent works that have found PAC-Bayesian analysis correlates with empirical performance better than Rademacher complexity (Jiang et al., 2020). We identify two critical measures for analyzing the generalization performance of fine-tuning. The first measure is the  $\ell_2$  norm of the distance between the pre-trained model (initialization) and the fine-tuned model. The second measure is the perturbed loss of the fine-tuned model, i.e. its loss after the model weights get perturbed by random noise. First, we observe that the fine-tuned weights remain closed to the pre-trained model. Moreover, the top layers travel much further away from the pre-trained model than the bottom layers. Second, we find that fine-tuning from a pre-trained model implies better noise stability than training from a randomly initialized model. In Section 4.2, we evaluate regularized fine-tuning for target tasks with noisy labels. We find that fine-tuning is prone to "memorizing the noisy labels", and regularization helps alleviate such memorization behavior. Moreover, we observe that the neural network has not yet overfitted to the noisy labels during the early phase of fine-tuning. Thus, its prediction could be used to relabel the noisy labels.

We propose an algorithm that incorporates layerwise regularization and self-labeling for improved regularization and robustness based on our results. Figure 1 illustrates the two components. First, we encode layer-wise distance constraints to regularize the model weights at different levels. Compared to (vanilla) fine-tuning, our algorithm reduces the gap between the training and test accuracy, thus alleviating overfitting. Second, we add a self-labeling mechanism that corrects and reweights "noisy labels" based on the neural network's predictions. Figure 1 shows that our algorithm effectively hinders the model from learning the incorrect labels by relabeling them to correct ones.



Figure 1: Red: Layer-wise regularization closes generalization gap. Magenta: Self-labeling relabels noisy data points to their correct label.

In Section 5, we evaluate our proposed algorithm for both transfer learning and few-shot classification tasks with image and text data sets. First, using ResNet-101 (He et al., 2016) as the pre-trained model, our algorithm outperforms previous fine-tuning methods on seven image classification tasks by 1.76% on average and 3.56% when their labels are noisy. Second, we find qualitatively similar results for applying our approach to medical image classification tasks (ChestX-ray14 (Wang et al., 2017; Rajpurkar et al., 2017)) and vision transformers (Dosovitskiy et al., 2020). Finally, we extend our approach to few-shot learning and sentence classification. For these related but different tasks and data modalities, we find an improvement of 0.75% and 0.46% over previous methods, respectively.

In summary, our contributions are threefold. First, we provide a PAC-Bayesian analysis of regularized fine-tuning. Our result implies empirical measures that explain the generalization performance of regularized fine-tuning. Second, we present a regularized self-labeling approach to enhance the generalization and robustness properties of fine-tuning. Third, we validate our approach on an extensive collection of classification tasks and pre-trained model architectures.

# 2 Related work

Fine-tuning is widely used in multi-task and transfer learning, meta-learning, and few-shot learning. Previous works (Li et al., 2018b,a) find that injecting  $\ell_2$  regularization helps improve the performance of fine-tuning. Li et al. (2018b) propose a  $\ell_2$  distance regularization method that penalizes the  $\ell_2$ distance between the fine-tuned weights and the pre-trained weights. Li et al. (2018a) penalize the distance between the *feature maps* as opposed to the layer weights. Chen et al. (2019) design a regularization method that suppresses the spectral components (of the *feature maps*) with small singular values to avoid negative transfer. Gouk et al. (2021) instead encode distance constraints in constrained minimization and use projected gradient descent to ensure the weights are close to the pre-trained model. Salman et al. (2020) show that fine-tuning from adversarially robust pre-trained models outperforms fine-tuning from (standard) pre-trained models.

The robustness of learning algorithms in the presence of label noise has been extensively studied in supervised learning (Natarajan et al., 2013). Three broad ideas for designing robust algorithms include defining novel losses, identifying noisy labels, and regularization methods. Zhang and Sabuncu (2018) design the robust Generalized Cross Entropy (GCE) loss which is a mixture of Cross Entropy (CE) and mean absolute error. The Symmetric Cross Entropy (SCE) (Wang et al., 2019) loss combines reverse cross-entropy with the CE loss. Ma et al. (2020) proposes APL, which normalizes the loss to be robust to noisy labels and combines active and passive loss functions. Thulasidasan et al. (2019) propose DAC, which identifies and suppresses the signals of noisy samples by abstention-based training. Liu et al. (2020) introduce ELR, an early learning regularization approach to mitigate label "memorization". Huang et al. (2020) propose a self-adaptive training method that corrects noisy labels and reweights training data to suppress erroneous signals. These works primarily concern the supervised learning setting. To the best of our knowledge, fine-tuning algorithms under label noise are relatively under-explored in transfer learning. Our approach draws inspiration from the semi-supervised learning literature, which has evaluated pseudo-labeling and self-training approaches given a limited amount of labeled data, and a large amount of unlabeled data (Zou et al., 2019; Xu et al., 2021; Tai et al., 2021). A very recent work considered a sharpness-aware approach and evaluated its performance for fine-tuning from noisy labels (Foret et al., 2021). It would be interesting to see if combining their approach with our ideas could lead to better results.

From a theoretical perspective, the seminal work of Ben-David et al. (2010) considers a setting where labeled data from many source tasks and a target task is available. They show that minimizing a weighted combination of the source and target empirical risks leads to the best result. In the supervised setting, Arora et al. (2019) provide data-dependent generalization bounds for neural networks. Nagarajan and Kolter (2018) and Wei and Ma (2019a,b) provide improved generalization bounds that depend only polynomially on the depth of the neural network characterized by a margin condition. Recent work has found that the PAC-Bayes theory provides generalization measures that correlate with empirical generalization performance better than other alternatives (Jiang et al., 2020). We defer a more extensive review of PAC-Bayesian generalization theory to Section A.

# **3** Preliminaries

**Problem setup.** We begin by formally introducing the setup. Suppose we would like to solve a target task. We have a training data set of size  $n^{(t)}$ . Let  $(x_1^{(t)}, y_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, y_{n^{(t)}}^{(t)})$  be the feature vectors and the labels of the training data set. We assume that every  $x_i^{(t)}$  lies in a *d*-dimensional space denoted by  $\mathcal{X} \subset \mathbb{R}^d$ . Following standard terminologies in statistical learning, we assume all data are drawn from some unknown distribution supported on  $\mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{Y} \subseteq \mathbb{R}$  is the label space. Denote the underlying data distribution as  $\mathcal{P}^{(t)}$ .

For our result in Section 4.1, we consider feedforward neural networks, and our results can also be extended to convolutional neural networks (e.g., using ideas from Long and Sedghi (2020)). Consider an L layer neural network. For each layer i from 1 to L, let  $\psi_i$  be the activation function and let  $W_i$  be the weight matrix at layer i. Given an input z to the *i*-th layer, the output is then denoted as  $\phi_i(z) = \psi_i(W_i z)$ . Thus, the final output of the network is given by

 $f_W(x) = \phi_L \circ \phi_{L-1} \circ \cdots \circ \phi_1(x)$ , for any input  $x \in \mathcal{X}$ ,

where we use  $W = [W_1, \ldots, W_L]$  to include all the parameters of the network for ease of notation. The prediction error of  $f_W$  is measured by a loss function  $\ell(\cdot)$  that is both convex and 1-Lipschitz.

$$\mathcal{L}^{(t)}(f_W) = \mathop{\mathbb{E}}_{(x,y)\sim\mathcal{P}^{(t)}} \left[ \ell(f_W(x), y) \right]. \tag{1}$$

Suppose we have access to a pre-trained source model  $\hat{W}^{(s)}$ . Using  $\hat{W}^{(s)}$  as an initialization, we then fine-tune the layers  $\hat{W}_1^{(s)}, \ldots, \hat{W}_L^{(s)}$  to predict the target labels.

**Applications.** We consider seven image classification data sets described in Table 1. We use ResNets pre-trained on ImageNet as the initialization  $\hat{W}^{(s)}$  (Russakovsky et al., 2015; He et al., 2016). We

Table 1: Basic statistics for seven image classification tasks.

Datasets	Training	Validation	Test	Classes
Aircrafts (Maji et al., 2013)	3334	3333	3333	100
CUB-200-2011 (Wah et al., 2011)	5395	599	5794	200
Caltech-256 (Griffin et al., 2007)	7680	5120	5120	256
Stanford-Cars (Krause et al., 2013)	7330	814	8441	196
Stanford-Dogs (Khosla et al., 2011)	10800	1200	8580	120
Flowers (Nilsback and Zisserman, 2008)	1020	1020	6149	102
MIT-Indoor (Sharif Razavian et al., 2014)	4824	536	1340	67

perform fine-tuning using the pre-trained network on the above data sets. See Section 5.1 for further description of the training procedure.

**Modeling label noise.** In many settings, fine-tuning is applied to a target task whose labels may contain noise; for example, if the target labels are created using weak supervision techniques (Ratner et al., 2019; Saab et al., 2021). To capture such settings, we denote a noisy data set as  $(x_1^{(t)}, \tilde{y}_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, \tilde{y}_{n^{(t)}}^{(t)})$ , where  $\tilde{y}_i^{(t)}$  is a noisy version of  $y_i^{(t)}$ . We consider two types of label noise: *independent* noise and *correlated* noise. We say that the label noise is *independent* if it is independent of the input feature vector

$$\Pr(\tilde{y}_i^{(t)} = k | y_i^{(t)} = j, x_i^{(t)}) = \Pr(\tilde{y}_i^{(t)} = k | y_i^{(t)} = j) = \eta_{j,k}$$

for some fixed  $\eta_{j,k}$  between 0 and 1. On the other hand, we say that the label noise is *correlated* if it depends on the input feature vector (i.e. the above equation does not hold).

# 4 Our proposed approaches

Given the problem setup described above, next, we study the generalization and robustness properties of regularized fine-tuning methods. First, we present a PAC-Bayes generalization bound for regularized fine-tuning. This result motivates us to evaluate two empirical measures: the fine-tuned distance in each layer and the perturbed loss of the fine-tuned model. Second, we consider fine-tuning from noisy labels. We show that layer-wise regularization prevents the model from memorizing the noisy labels. We then suggest injecting predictions of the model during training, similar to self-training and pseudo-labeling. Finally, we incorporate both components into our *regularized self-labeling* approach, blending the strengths of both to improve the generalization performance and the robustness of fine-tuning.

#### 4.1 Fine-tuning and regularization

We consider the following regularized fine-tuning problem, which constrains the network from traveling too far from the pre-trained initialization  $\hat{W}^{(s)}$  (Li et al., 2018a,b; Gouk et al., 2021).

$$\hat{W} \leftarrow \arg\min \hat{\mathcal{L}}^{(t)}(f_W)$$
s.t.  $\|W_i - \hat{W}_i^{(s)}\|_F \le D_i, \ \forall i = 1, \dots, L.$ 
(2)
(3)

Above,  $D_i$  is a hyperparameter that constrains how far the *i*-th layer  $W_i$  can travel from the pre-trained initialization  $\hat{W}_i^{(s)}$ . Previous works have observed that stronger regularization reduces the generalization gap during fine-tuning (cf. Figure 2). Next, we analyze the generalization error of  $\hat{W}$ , that is,  $\mathcal{L}^{(t)}(f_{\hat{W}}) - \hat{\mathcal{L}}^{(t)}(f_{\hat{W}})$ , where  $\mathcal{L}^{(t)}(f_{\hat{W}})$  is the test loss of  $f_{\hat{W}}$  according to equation (1) and  $\hat{\mathcal{L}}^{(t)}(f_{\hat{W}})$ is the empirical loss of  $f_{\hat{W}}$  on the training data set.

**PAC-Bayesian analysis.** We begin by analyzing the generalization error of  $f_{\hat{W}}$  using PAC-Bayesian



Figure 2: The training and test accuracy of finetuning using early stopping vs. optimizing (2) and (3) on the Indoor data set.



		CUB-200-201	1
$\sigma$	Random	Pre-trained	Adversarial
$10^{-2}$	3.77±0.42	1.45±0.13	$1.76 {\pm} 0.09$
$10^{-3}$	$0.82{\pm}0.07$	$0.62 {\pm} 0.03$	$0.54{\pm}0.03$
$10^{-4}$	$0.81 {\pm} 0.04$	$0.61{\pm}0.03$	$0.61{\pm}0.01$
		Indoor	
$\sigma$	Random	Pre-trained	Adversarial
$10^{-2}$	2.51±0.34	$1.11 \pm 0.09$	0.97±0.07
$10^{-3}$	$0.49 {\pm} 0.09$	$0.36 {\pm} 0.05$	$0.32{\pm}0.04$
$10^{-4}$	$0.44 {\pm} 0.03$	$0.33{\pm}0.02$	$0.30{\pm}0.04$

#### (b) Perturbed loss.

Figure 3: Left: Ratio between the  $\ell_2$ -norm of the fine-tuned distances and the pre-trained network. Right: Perturbed loss of fine-tuned model from random initializations (Random), pre-trained model initializations (Pre-trained), adversarially trained model initializations (Adversarial). Results are shown for the CUB-200-2011 and the Indoor data sets using ResNet-18; Averaged over 10 runs.

tools (McAllester, 1999a,b). This departs from the previous work of Gouk et al. (2021), which bases on their analysis using Rademacher complexity. This change of perspective is inspired by several recent works that have found PAC-Bayesian bounds to better correlate with empirical performance than Rademacher complexity bounds (Jiang et al., 2020). We refer the interested reader to Section A for further references from this line of work. After presenting our result, we provide empirical measures of our result and discuss the practical implication.

**Theorem 4.1** (PAC-Bayes generalization bound for fine-tuning). Suppose for every i = 1, ..., L,  $\|\hat{W}_i^{(s)}\|_2 \leq B_i$ , for a fixed  $B_i > 1$ . Suppose the feature vectors in the domain  $\mathcal{X}$  are all bounded:  $\|x\|_2 \leq C_1$  for every  $x \in \mathcal{X}$ , for some  $C_1 \geq 1$ . Finally, suppose the loss function  $\ell(\cdot)$  is 1-Lipschitz and bounded from above by a fixed constant  $C_2$ . Under these conditions, let  $f_{\hat{W}}$  be the minimizer of regularized fine-tuning, solved from problem (3). Let  $\varepsilon > 0$  be an arbitrary small value. Let H be the maximum over the width over all the L layers and the input dimension d. Then, with probability at least  $1 - 2\delta$  for some small  $\delta > 0$ , the expected loss  $\mathcal{L}^{(t)}(f_{\hat{W}})$  is upper bounded by

$$\mathcal{L}^{(t)}(f_{\hat{W}}) \leq \hat{\mathcal{L}}^{(t)}(f_{\hat{W}}) + \varepsilon + C_2 \sqrt{\frac{\frac{36}{\varepsilon^2} \cdot C_1^2 H \log(4LHC_2) \left(\sum_{i=1}^L \frac{\prod_{j=1}^L (B_j + D_j)}{B_i + D_i}\right)^2 \left(\sum_{i=1}^L D_i^2\right) + 3\ln\frac{n^{(t)}}{\delta} + 8}{n^{(t)}}}.$$
 (4)

*Relation to prior works.* Compared to the result of Gouk et al. (2021), we instead proceed by factoring the generalization error into a noise error (i.e. the  $\varepsilon$  term in equation (4)) and a KL-divergence between  $\hat{W}^{(s)}$  and  $\hat{W}$  (i.e. the final term in equation (4)). The proof of Theorem 4.1 is based on Neyshabur et al. (2018) and is presented in Appendix A. The difference between Theorem 4.1 and the result of Neyshabur et al. (2018) is that our result is stated for the (e.g. cross-entropy) loss function  $\ell(\cdot)$  whereas Neyshabur et al. (2018) states the result using the soft margin loss.

Our result suggests that the *fine-tuned distances*  $\{D_i\}_{i=1}^L$  and the *perturbed loss* (more precisely  $\mathbb{E}_U \left[ \ell(f_{\hat{W}+U}(x), y) \right]$  where every entry of U is drawn from  $\mathcal{N}(0, \sigma^2)$ ) are two important measures for fine-tuning. Next, we empirically evaluate these two measures on real-world data sets.

**Fine-tuned distances**  $\{D_i\}_{i=1}^{L}$ . First, we present our empirical finding for the fine-tuned distances in each layer. We fine-tune a ResNet-18 model (pre-trained on ImageNet) on seven data sets and calculate the Frobenius distance between  $\hat{W}_i^{(s)}$  and  $\hat{W}_i$  for every layer *i*.

Figure 3a shows the ratio of the fine-tuned distances to the pre-trained weights in each layer. First, we observe that the fine-tuned distances are relatively small compared to the pre-trained network. This means that fine-tuning stays within a small local region near the pre-trained model. Second, we find that  $D_i$  varies across layers.  $D_i$  is smaller at lower layers and is larger at higher layers. This observation aligns with the folklore intuition that different layers in Convolutional neural networks play a different role (Guo et al., 2019). Bottom layers extract higher-level representations of features



Figure 4: Training and test accuracy of fine-tuning using early stopping and optimizing equations (5) and (6) on the Indoor data set with different levels of noise rate. Stronger regularization reduces the generalization gap in both settings.

and stay close to the pre-trained model. Top layers extract class-specific features and vary among different tasks. Thus, top layers travel further away from the pre-trained model.

Based on the above empirical finding, we propose layer-wise constraints as regularization. In particular, the constraints for the bottom layers should be small, and the constraints for the top layers should be large. In practice, we find that an exponentially increasing scheme involving a base distance parameter D and a scale factor  $\gamma$  works well. This is summarized in the following regularized fine-tuning problem:

$$\hat{W} \leftarrow \arg\min \hat{\mathcal{L}}^{(t)}(f_W) \tag{5}$$

s.t. 
$$\|W_i - \hat{W}_i^{(s)}\|_E \le D \cdot \gamma^{i-1}, \ \forall i = 1, \dots, L.$$
 (6)

where the scale factor  $\gamma > 1$  according to Figure 3a. The previous work of Gouk et al. (2021) uses the same distance parameter in equation (6), i.e.,  $D_i = D$ ,  $\forall i = 1, \ldots, L$ . In our experiments, we have observed that such constant regularization constraints are only active at the top layers. With layer-wise regularization, we instead extend the constraints to the bottom layers as well. In Table 8 (cf. Section B.3), we compare the value of  $\sum_{i=1}^{L} D_i^2$  between constant regularization and layer-wise regularization. We find that layer-wise regularization leads to a smaller value of  $\sum_{i=1}^{L} D_i^2$ , thus indicating a tighter generalization bound according to Theorem 4.1.

**Perturbed loss**  $\mathbb{E}_{U}\left[\ell(f_{\hat{W}+U}(x), y)\right]$ . Second, we compare the perturbed loss of models finetuned from random, pre-trained, and adversarially robust pre-trained model initializations. The perturbed loss is measured by the average loss on the training set, after perturbing each entry of the layer weights by a Gaussian random variable  $N(0, \sigma^2)$ .

In Table 3b, we report the perturbed losses for both the CUB-200-2011 and the Indoor data set. The results show that models fine-tuned from pre-trained initializations are more stable than models fine-tuned from random initializations.

*Explaining why adversarial robustness helps fine-tuning.* Recent work (Salman et al., 2020) found that performing adversarial training in the pre-training phase leads to models that transfer better to downstream tasks. More precisely, the adversarial training objective is defined as

$$\tilde{\mathcal{L}}_{adv}^{(s)}(\theta) = \min_{W} \frac{1}{n^{(s)}} \sum_{i=1}^{n^{(s)}} \max_{\|\delta\|_{2} \le \varepsilon} \ell(f_{W}(x_{i}^{(s)} + \delta), y_{i}^{(s)}).$$
(7)

Let  $\hat{W}_{adv}^{(s)}$  be a fine-tuned neural network using a pre-trained model from adversarial training. We measure the perturbed loss of  $\hat{W}_{adv}^{(s)}$  using the pre-trained model provided by Salman et al. (2020) on the CUB-200-2011 and Indoor data set. Table 3b shows that  $f_{\hat{W}_{adv}^{(s)}}$  often incurs lower perturbed losses than models fine-tuned from random and pre-trained initializations. In Section B.3, we additionally observe that the layers of  $\hat{W}_{adv}^{(s)}$  generally have lower Frobenious norms compared to  $\hat{W}^{(s)}$ . Thus, the improved noise stability and the smaller layer-wise norms together imply a tighter generalization bound for adversarial training. We leave a thorough theoretical analysis of adversarial training in the context of fine-tuning to future works.





(b) Average weight injected by self label-reweighting.

Figure 5: Left: Comparing the precision of self-labeling with vs. without regularization. Precision is defined as the number of data points whose labels are corrected (cf. Line 7 in Alg. 1) divided by the number of data points whose labels are changed at every epoch. Right: Comparing the average normalized weight (cf. Line 10 in Alg. 1) of training data points with the correct label vs. data points with an incorrect label.

#### 4.2 Fine-tuning and robustness

Next, we extend our approach to fine-tuning from noisy labels. We compare fine-tuning with and without regularization on a target task with independent label noise to motivate our approach. In Figure 4, we plot the training and test accuracy curves on the Indoor dataset with two noise rates. We observe that the test accuracy increases at first, meaning the network learns from the correct labels during the first few epochs. As the learning process progresses, the test accuracy of fine-tuning decreases, implying the network is "overfitting" to noisy labels. Furthermore, the gap between the training and test accuracy curves is more than 20%. Similar observations have been presented in prior works (Zhang et al., 2016; Li et al., 2018c; Liu et al., 2020), where over-parametrized neural networks can memorize the entire training set. One way to mitigate memorization is to explicitly reduce model capacity (Wu et al., 2020; Yang et al., 2021). Another way is via explicit regularization since the distance constraint significantly reduces the training-test accuracy gaps and improves test accuracy by a significant margin. Based on our empirical observation, the neural network has some discriminating power during the early fine-tuning phase. Thus, we can leverage the model predictions to relabel incorrect labels and reweight data points with incorrect labels. We describe the two components next.

**Self label-correction.** We propose a label correction step to augment the number of correct labels. We leverage the discriminating power of a network during the early phase and correct data points for which the model has high confidence. Concretely, let  $p_i$  be the prediction of the *i*-th data point. Given a confidence threshold  $p_t$  at epoch t, if  $\max(p_i) \ge p_t$  and  $\arg \max(p_i) \ne \tilde{y}_i^{(t)}$ , we identify data point *i* as noisy and relabel it to

$$\tilde{y}_i^{(t)} = \arg\max(\boldsymbol{p_i}).$$

Figure 5a illustrates the precision of this step. We define precision as the number of correctly relabeled data points divided by the total number of data points whose labels are relabeled. We observe that precision is high at the beginning (around 60%) and gradually decreases in later epochs. Moreover, we find precision increases with regularization. This suggests an intricate interaction between regularization and self-labeling during fine-tuning. We elaborate on the role of regularization in label correction by looking at the number of relabeled data points in Section B.3.

**Self label-reweighting.** We incorporate a soft data removal step to prevent the model from overfitting to noisy labels. We identify data points as noisy if their loss values are large and thus down-weight them during training. The model will prioritize data points with correct labels in its gradient by assigning smaller weights to large-loss data points. More precisely, given a data point  $(x_i^{(t)}, \tilde{y}_i^{(t)})$ , we reweight it by  $\omega_i = \exp(-\ell(f_W(x_i^{(t)}), \tilde{y}_i^{(t)})/\tau)$  where  $\tau$  is a temperature parameter. In the implementation, we normalize the weights of every data point in every mini-batch *B*:

$$\hat{\mathcal{L}}_B(f_W) = \frac{1}{\sum_{i \in B} \omega_i} \sum_{i \in B} \omega_i \cdot \ell(f_W(x_i^{(t)}), \tilde{y}_i^{(t)}).$$

#### Algorithm 1 Regularized self-labeling (REGSL)

**Input:** Input nosiy data  $(x_1^{(t)}, \tilde{y}_1^{(t)}), \ldots, (x_{n^{(t)}}^{(t)}, \tilde{y}_{n^{(t)}}^{(t)})$ , network model  $f_W$ , constraint distance D, distance scale factor  $\gamma$ , re-weight start step  $E_r$ , re-weight temperature  $\tau$ , label correction start step  $E_c$ , label correction threhold  $p_t$ , and fine-tuning steps T

1: Initialize model parameters with pre-trained weights  $W = \hat{W}^s$  and the start step t = 0

2: while t < T do Fetch a mini-batch of data  $\{(x_i^{(t)}, \tilde{y}_i^{(t)})\}_{i=1}^m$ 3: for each data point *i* in the batch do 4:

Calculate the predicted probability  $\mathbf{p}_i = \operatorname{softmax}(f_W(x_i^{(t)}))$ if  $t > E_c$  and  $\max(\mathbf{p}_i) > p_t$  and  $\arg \max(\mathbf{p}_i) \neq \tilde{y}_i^{(t)}$  then Correct the label  $\tilde{y}_i^{(t)} = \arg \max(\mathbf{p}_i)$ 5:

- 6:
- 7:
- end if 8:
- 9:
- Calculate loss  $\ell_i = \ell(f_W(x_i^{(t)}), \tilde{y}_i^{(t)})$ if  $t > E_r$  then Calculate the weight  $\omega_i = \exp(-\ell_i/\tau)$  else Set the weight  $\omega_i = 1$ 10:
- 11: end for
- 12:
- Update  $f_W$  by SGD on  $\mathcal{L} = -\frac{1}{\sum_i \omega_i} \sum_i^m \omega_i \ell_i$  and t = t + 1Project the weights  $W_i$  of each layer inside the layer-wise constraint region as in Equation 6 13:

```
14: end while
```

Figure 5b compares the average weight of data points with the correct label to data points with an incorrect label. We find that with our reweighting scheme, the average weight of incorrectly labeled data points is much lower than the average weight of correctly labeled data points. Thus, the reweighting scheme expands the gradient of correctly labeled data points, thus reducing the fraction of noisy data points in the training set.

The pseudo-code for our final approach, which combines both layer-wise regularization and selflabeling, is presented in Algorithm 1.

#### 5 **Experiments**

We evaluate our proposed algorithm in a wide range of tasks and pre-trained networks. First, we show that our algorithm outperforms the baselines by 1.75% on average over seven transfer learning tasks and can generalize to a more distant task of medical images. Second, our algorithm improves the robustness of fine-tuning in the presence of noisy labels compared to previous methods. Our algorithm improves over numerous baselines by 3.56% on average under both independent and correlated label noise. Moreover, our approach can achieve a similar performance boost for finetuning vision transformer models on noisy labeled data. Finally, we show our algorithms improve over fine-tuning baselines by 0.75% on average for the related task of few-shot classification. Our extensive evaluation confirms that our approach applies to a broad range of settings, validating our algorithmic insights. Due to space constraints, we defer the ablation studies to Section B.3. Our code is available at https://github.com/NEU-StatsML-Research/Regularized-Self-Labeling.

#### 5.1 Experimental setup

**Data sets.** First, we evaluate fine-tuning on seven image classification data sets and one medical image data set. The statistics of the seven image data sets are described in Table 1. For the medical imaging task, we consider the ChestX-ray14 data set contains 112120 frontal-view chest X-ray images labeled with 14 different diseases (Wang et al., 2017; Rajpurkar et al., 2017). Second, we evaluate Algorithm 1 under two different kinds of label noise to test the robustness of fine-tuning. We selected the MIT-Indoor data set (Sharif Razavian et al., 2014) as the benchmark data set and randomly flipped the labels of the training samples. The results for the other data sets are similar, which can be found in Section B.2. We consider both independent random noise and correlated noise in our experiments. We generate the independent random noise by flipping the labels uniformly with a given noise rate. We simulate the correlated noise by using the predictions of an auxiliary network as noisy labels. Section B.1 describes the label flipping process. For few-shot image classification, we conduct experiments on the miniImageNet benchmark (Vinyals et al., 2016) following the setup of Tian et al. (2020).

	Aircrafts	CUB-200-2011	Caltech-256	Stanford-Cars	Stanford-Dogs	Flowers	MIT-Indoor
Fine-tuning	73.96±0.34	80.31±0.26	79.41±0.23	89.28±0.14	82.89±0.35	92.92±0.15	$74.78 {\pm} 0.97$
$\ell^2$ -Norm	$74.84 {\pm} 0.15$	$80.80 {\pm} 0.16$	$79.67 \pm 0.31$	$88.96 {\pm} 0.18$	$83.00 \pm 0.10$	$92.76 \pm 0.26$	$76.57 \pm 0.70$
LS	$74.19 \pm 0.41$	$82.22 \pm 0.17$	$81.10 \pm 0.13$	$89.66 {\pm} 0.10$	$84.14 \pm 0.21$	$93.72 {\pm} 0.04$	$76.84 {\pm} 0.31$
$\ell^2$ -SP	$74.09 \pm 0.98$	$81.49 \pm 0.36$	$84.13 \pm 0.09$	$88.96 \pm 0.09$	$88.95 \pm 0.15$	$93.06 {\pm} 0.28$	$78.11 \pm 0.37$
$\ell^2$ -PGM	$74.90 {\pm} 0.26$	$81.23 {\pm} 0.32$	$83.25 {\pm} 0.33$	$88.92 {\pm} 0.39$	$86.48 {\pm} 0.28$	$93.23 {\pm} 0.34$	$77.31 {\pm} 0.30$
REGSL (ours)	$75.32{\pm}0.23$	82.24±0.21	84.90±0.16	89.14±0.22	89.58±0.13	93.82±0.35	79.30±0.31

Table 2: Top-1 test accuracy for fine-tuning ResNet-101 pre-trained on the ILSVRC-2012 subset of ImageNet. Results are averaged over 3 random seeds.

Architectures. We use the ResNet-101 (He et al., 2016) network for transfer learning tasks and ResNet-18 (He et al., 2016) network for ChestX-ray data set. We use ResNet-18 network for the labelnoise experiments and extend the our algorithm to the vision transformer (ViT) model (Dosovitskiy et al., 2020). The ResNet models are pre-trained on ImageNet (Russakovsky et al., 2015) and the ViT model is pre-trained on ImageNet-21k data set (Dosovitskiy et al., 2020). For the few-shot learning tasks, we use ResNet-12 pre-trained on the meta-training data set as in the previous work (Tian et al., 2020). We set four different values of  $D_i$  for the four blocks of the ResNet models in our algorithm. We describe the fine-tuning procedure and the hyperparameters in Section B.1.

**Baselines.** For the transfer learning tasks, we use the Frobenius norm for distance regularization. We present ablation studies to compare the Frobenius norm and other norms in Section B.3. we compare our algorithm with fine-tuning with early stop (Fine-tuning), fine-tuning with weight decay ( $\ell^2$ -Norm), label smoothing (LS) method formulated in Müller et al. (2019),  $\ell^2$ -SP (Li et al., 2018b), and  $\ell^2$ -PGM (Gouk et al., 2021). For testing the robustness of our algorithm, in addition to the baselines described above, we adopt several baselines that have been proposed in the supervised learning setting, including GCE (Zhang and Sabuncu, 2018), SCE (Wang et al., 2019), DAC (Thulasidasan et al., 2019), APL (Ma et al., 2020), ELR (Liu et al., 2020) and self-adaptive training (SAT) (Huang et al., 2020). We compare with the same baselines to transfer learning in few-shot classification tasks. We describe the implementation and hyper-parameters of baselines in Section B.1.

#### 5.2 Experimental results

Improved regularization for transfer learning. We report the test accuracy of fine-tuning ResNet-101 on seven data sets in Table 2. We observe that our method performs the best among six regularized fine-tuning methods on average. Our proposed layer-wise regularization method provides an improvement of 1.76% on average over distance-based regularization (Gouk et al., 2021). In particular, our approach outperforms  $\ell^2$ -PGM by  $2 \sim 3\%$  on both Stanford-Dogs and MIT-Indoor data sets. This suggests that adding appropriate distance constraints for all the layers is better than applying only one constraint to the top layers. In Table 10 (cf. Section B.3), we note that using the MARS norm of Gouk et al. (2021) yields comparable results to using the  $\ell_2$  norm with our approach.

Next, we apply our approach to medical image classification, a more distant transfer task (relative to the source data set ImageNet). We report the mean AUROC (averaged over predicting all 14 labels) on the ChestX-ray14 data set in Table 6a (cf. Section B.2). With layer-wise regularization, we see an 0.39% improvement over the baseline methods. Finally, we extend our approach to text classification. The results are consistent with the above and can be found in Table 7 (cf. Section B.2).

**Improved robustness under label noise.** Next, we report the test accuracy of our approach on the Indoor data set with independent and correlated noise in Table 3. We find that our approach consistently outperforms baseline methods by  $1 \sim 3\%$  for various settings involving label noise. First, our method (REGSL) improves the performance by over 4% on average compared to distance-based regularization (Gouk et al., 2021). This result implies our method is more robust to label noise in the training labels. Second, our method outperforms previous supervised training methods by **3.56**% on average. This result suggests that regularization is critical for fine-tuning. Taken together, our results suggest that regularization and self-labeling complement each other during fine-tuning. This is reinforced by our ablation study in Section B.3; We study the influence of each component of our algorithm and find that removing any component degrades performance.

We further apply our approach to fine-tuning pre-trained ViT (Dosovitskiy et al., 2020) from noisy labels, under the same setting as Table 3. Table 6b shows the result (cf. Section B.2). First, we find that

			correlated noise			
Data sets	Methods	20%	40%	60%	80%	25.18%
	Fine-tuning	$65.02 {\pm} 0.39$	57.49±0.39	44.60±0.95	27.09±0.19	67.49±0.74
	LS	$67.04 {\pm} 0.58$	$58.98 {\pm} 0.57$	$48.56 {\pm} 0.53$	$25.82{\pm}0.90$	$68.06 {\pm} 0.76$
	$\ell^2$ -PGM	$69.45 {\pm} 0.19$	$62.74 {\pm} 0.60$	$51.24 {\pm} 0.15$	$30.15 {\pm} 0.94$	$68.86 {\pm} 0.27$
	GCE	$70.45 {\pm} 0.40$	64.73±0.21	$54.10 {\pm} 0.16$	$29.58 {\pm} 0.26$	$68.88 {\pm} 0.18$
Indoor	SCE	$69.43 {\pm} 0.20$	$64.68 {\pm} 0.45$	$55.07 {\pm} 0.52$	$29.85 {\pm} 0.44$	$69.00 \pm 1.22$
muoor	DAC	$64.45 {\pm} 0.31$	$59.73 {\pm} 0.27$	$47.44 {\pm} 0.09$	$26.69 {\pm} 0.34$	$68.06 \pm 1.32$
	APL	$70.05 {\pm} 0.41$	$66.22 {\pm} 0.10$	$52.51 {\pm} 0.66$	$30.90 {\pm} 0.37$	68.31±0.77
	SAT	$68.98 {\pm} 0.63$	$63.43 {\pm} 0.72$	$52.84 {\pm} 0.38$	$29.60 {\pm} 0.53$	$67.06 {\pm} 0.55$
	ELR	$71.43 {\pm} 0.80$	$66.34{\pm}0.48$	$55.22 \pm 0.73$	$31.24{\pm}0.19$	$69.38 {\pm} 0.49$
	REGSL (ours)	72.51±0.46	68.13±0.16	57.59±0.55	34.08±0.79	70.12±0.83

Table 3: Top-1 test accuracy of fine-tuning ResNet-18 on the Indoor data set with various settings of noisy labels in the training set. Results are averaged over 3 random seeds.

our approach improves upon the best regularization methods by 1.17% averaged over two settings. Second, we find that our approach also improves upon self-labeling by 13.57% averaged over two settings. These two results again highlight that both regularization and self-labeling contribute to the final result. While regularization prevents the model from overfitting to the random labels, self-labeling injects the belief of the fine-tuned model into the noisy data set.

**Extension to few-shot classification.** We extend our approach to a few-shot image classification task. We compare our approach to the baseline regularization methods. In Table 4, we report the average accuracy of 600 sampled tasks from the meta-test split of the miniImageNet benchmark (Vinyals et al., 2016). We find that our layer-wise regularization method achieves 0.75% improvement on average compared to previous regularization methods. Additionally, regularization methods generally improve the performance of fine-tuning by over 1%.

Table 4: Test accuracy with 95% confidence interval for fine-tuning ResNet-12 over 600 meta-test splits of miniImageNet.

Methods	miniIma	miniImageNet				
1010thotas	5-way-1-shot	5-way-5-shot				
Fine-tuning	$60.56 \pm 0.78$	$76.42\pm0.36$				
$\ell^2$ -Norm	$60.93 \pm 0.81$	$76.57\pm0.55$				
LS	$61.31\pm0.77$	$76.73\pm0.62$				
$\ell^2$ -SP	$61.48 \pm 0.76$	$77.02\pm0.62$				
$\ell^2$ -PGM	$61.35\pm0.83$	$77.33\pm0.59$				
REGSL (ours)	$\textbf{61.71} \pm \textbf{0.77}$	$\textbf{78.03} \pm \textbf{0.54}$				

Ablation studies. We study the influence of removing each component from our algorithm. Results

shown in Table 11 (cf. Section B.3) suggest that they all degrade performance. Thus, all three components in Algorithm 1 contribute to the final performance. The importance of each component depends on the noise rate. In particular, if the noise rate is higher than 40%, label-correction is the most critical component. Other ablation studies we present include: comparing the  $\ell_2$  norm and the MARS norm in our algorithm and comparing the distance parameters between layer-wise and constant regularization. We leave the details to Section B.3.

# 6 Conclusion

This paper studied regularization methods for fine-tuning as well as their robustness properties. We investigated the generalization error of fine-tuning using PAC-Bayesian techniques. This leads to two empirical measures, which we empirically computed. The analysis inspired us to consider layer-wise regularization for fine-tuning. This approach performs well on both pre-trained ResNets and ViTs; we have found the fine-tuned distance using vanilla fine-tuning and encoded the distance patterns in layer-wise regularization. We then evaluated the performance of regularized fine-tuning from noisy labels. We proposed a self label-correction and label-reweighting approach in the noisy setting. We found that regularization and self-labeling complement each other in our experiments, leading to significant improvement over previous methods, which either use regularization methods or use self-labeling methods but not both. We discuss the limitation of our work in Section C.

**Negative societal impacts.** Our work concerns the theory and the empirical performance of finetuning for transfer learning. Due to its technical nature, there is very little negative societal impact. Since transfer learning is commonly used in practice, one avenue that deserves attention for these applications is to carefully evaluate bias or fairness metrics before applying our proposed approaches. **Acknowledgment.** Thanks to the anonymous reviewers and the area chairs for carefully reading through our work and providing constructive feedback. Thanks to the program committee chairs for communicating with us regarding an omitted factor in equation (14), which we have added in the current version. HZ is grateful to Sen Wu and Chris Ré for stimulating discussions. HZ is supported by Northeastern University's Khoury seed/proof-of-concept grant.

### References

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A nextgeneration hyperparameter optimization framework. In *International Conference on Knowledge Discovery and Data Mining*, 2019.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*. PMLR, 2019.
- Peter Bartlett, Dylan J Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. *arXiv preprint arXiv:1706.08498*, 2017.
- Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *The Journal of Machine Learning Research*, 20(1):2285–2301, 2019.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 2010.
- Xinyang Chen, Sinan Wang, Bo Fu, Mingsheng Long, and Jianmin Wang. Catastrophic forgetting meets negative transfer: Batch spectral shrinkage for safe transfer learning. 2019.
- Guneet S Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. arXiv preprint arXiv:1909.02729, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint arXiv:1703.11008*, 2017.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *ICLR*, 2021.
- Henry Gouk, Timothy M Hospedales, and Massimiliano Pontil. Distance-based regularisation of deep networks for fine-tuning. *ICLR*, 2021.
- Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007.
- Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. Spottune: transfer learning through adaptive fine-tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Conference on computer vision and pattern recognition, pages 770–778, 2016.
- Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *International Conference on Computer Vision*, 2019.
- Minqing Hu and Bing Liu. Mining and summarizing customer reviews. In ACM SIGKDD international conference on Knowledge discovery and data mining, 2004.
- Lang Huang, Chao Zhang, and Hongyang Zhang. Self-adaptive training: beyond empirical risk minimization. Advances in Neural Information Processing Systems, 2020.

- Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *ICLR*, 2020.
- Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Li Fei-Fei. Novel dataset for fine-grained image categorization. In *Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition*, 2011.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Workshop on 3D Representation and Recognition*, 2013.
- Hao Li, Pratik Chaudhari, Hao Yang, Michael Lam, Avinash Ravichandran, Rahul Bhotika, and Stefano Soatto. Rethinking the hyperparameters for fine-tuning. *arXiv preprint arXiv:2002.11770*, 2020.
- Xin Li and Dan Roth. Learning question classifiers. In International conference on Computational linguistics, 2002.
- Xingjian Li, Haoyi Xiong, Hanchao Wang, Yuxuan Rao, Liping Liu, and Jun Huan. Delta: Deep learning transfer using feature map with attention for convolutional networks. In *International Conference on Learning Representations*, 2018a.
- Xuhong Li, Yves Grandvalet, and Franck Davoine. Explicit inductive bias for transfer learning with convolutional networks. *arXiv preprint arXiv:1802.01483*, 2018b.
- Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47. PMLR, 2018c.
- Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning regularization prevents memorization of noisy labels. *arXiv preprint arXiv:2007.00151*, 2020.
- Philip M Long and Hanie Sedghi. Generalization bounds for deep convolutional neural networks. *ICLR*, 2020.
- Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah Erfani, and James Bailey. Normalized loss functions for deep learning with noisy labels. In *International Conference on Machine Learning*, 2020.
- Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- David A McAllester. Pac-bayesian model averaging. In *Proceedings of the twelfth annual conference* on Computational learning theory, pages 164–170, 1999a.
- David A McAllester. Some pac-bayesian theorems. Machine Learning, 37(3):355–363, 1999b.
- Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help? *arXiv* preprint arXiv:1906.02629, 2019.
- Vaishnavh Nagarajan and Zico Kolter. Deterministic pac-bayesian generalization bounds for deep networks via generalizing noise-resilience. In *International Conference on Learning Representations*, 2018.
- Nagarajan Natarajan, Inderjit S Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Learning with noisy labels. In *Conference on Neural Information Processing Systems*, 2013.
- Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A pac-bayesian approach to spectrallynormalized margin bounds for neural networks. *ICLR*, 2018.
- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics & Image Processing*, 2008.

- Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Annual meeting on Association for Computational Linguistics*, 2004.
- Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Annual meeting on association for computational linguistics*, 2005.
- Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, et al. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. arXiv preprint arXiv:1711.05225, 2017.
- Alexander Ratner, Christopher De Sa, Sen Wu, Daniel Selsam, and Christopher Ré. Data programming: Creating large training sets, quickly. Advances in neural information processing systems, 2016.
- Alexander Ratner, Stephen H Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. Snorkel: Rapid training data creation with weak supervision. In *International Conference on Very Large Data Bases*. NIH Public Access, 2017.
- Alexander Ratner, Braden Hancock, Jared Dunnmon, Frederic Sala, Shreyash Pandey, and Christopher Ré. Training complex models with multi-task weak supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4763–4771, 2019.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 2015.
- Khaled Saab, Sarah M Hooper, Nimit S Sohoni, Jupinder Parmar, Brian Pogatchnik, Sen Wu, Jared A Dunnmon, Hongyang R Zhang, Daniel Rubin, and Christopher Ré. Observational supervision for medical image classification using gaze data. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 603–614. Springer, 2021.
- Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? *arXiv preprint arXiv:2007.08489*, 2020.
- Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Conference on computer vision and pattern recognition workshops*, 2014.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Conference on empirical methods in natural language processing*, 2013.
- Kai Sheng Tai, Peter Bailis, and Gregory Valiant. Sinkhorn label allocation: Semi-supervised classification via annealed self-training. *arXiv preprint arXiv:2102.08622*, 2021.
- Chuanqi Tan, Fuchun Sun, Tao Kong, Wenchang Zhang, Chao Yang, and Chunfang Liu. A survey on deep transfer learning. In *International conference on artificial neural networks*, 2018.
- Sunil Thulasidasan, Tanmoy Bhattacharya, Jeff Bilmes, Gopinath Chennupati, and Jamal Mohd-Yusof. Combating label noise in deep learning using abstention. arXiv preprint arXiv:1905.10964, 2019.
- Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? *arXiv preprint arXiv:2003.11539*, 2020.
- Joel A Tropp. An introduction to matrix concentration inequalities. *arXiv preprint arXiv:1501.01571*, 2015.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. *arXiv preprint arXiv:1606.04080*, 2016.
- C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.

- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Conference on computer vision and pattern recognition*, 2017.
- Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In *International Conference on Computer Vision*, pages 322–330, 2019.
- Colin Wei and Tengyu Ma. Data-dependent sample complexity of deep neural networks via lipschitz augmentation. *arXiv preprint arXiv:1905.03684*, 2019a.
- Colin Wei and Tengyu Ma. Improved sample complexities for deep networks and robust classification via an all-layer margin. *arXiv preprint arXiv:1910.04284*, 2019b.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 2005.
- Sen Wu, Hongyang R Zhang, and Christopher Ré. Understanding and improving information transfer in multi-task learning. *ICLR*, 2020.
- Yi Xu, Lei Shang, Jinxing Ye, Qi Qian, Yu-Feng Li, Baigui Sun, Hao Li, and Rong Jin. Dash: Semi-supervised learning with dynamic thresholding. In *International Conference on Machine Learning*, pages 11525–11536. PMLR, 2021.
- Fan Yang, Hongyang R Zhang, Sen Wu, Weijie J Su, and Christopher Ré. Analysis of information transfer from heterogeneous sources via precise high-dimensional asymptotics. *arXiv preprint arXiv:2010.11750*, 2021.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- Zhilu Zhang and Mert R Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *arXiv preprint arXiv:1805.07836*, 2018.
- Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P Adams, and Peter Orbanz. Non-vacuous generalization bounds at the imagenet scale: a pac-bayesian compression approach. *arXiv preprint arXiv:1804.05862*, 2018.
- Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. Confidence regularized self-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5982–5991, 2019.

# A Proof of Theorem 4.1

**Background.** PAC-Bayesian generalization theory offers a nice way to blend data-dependent properties such as noise stability and sharpness into generalization bounds. Recent works (e.g., Bartlett et al. (2017); Neyshabur et al. (2018)) introduced generalization bounds for multi-layer neural networks in an attempt to explain why neural network models generalize well despite having more trainable parameters than the number of training examples. On the one hand, the VC dimension of a neural network is known to be roughly equal to its number of parameters (Bartlett et al., 2019). On the other hand, the number of parameters is not a good capacity measure for neural nets, as evidenced by the popular work of Zhang et al. (2016). The bounds of Bartlett et al. (2017) (and follow-up works) provide a more meaningful notion of "capacity" compared to VC-dimension.

While these bounds constitute an improvement compared to classical learning theory, it is unclear if these bounds are tight or non-vacuous. One response to this criticism is a computational framework from Dziugaite and Roy (2017). That work shows that directly optimizing the PAC-Bayes bound leads to a much smaller bound and low test error simultaneously (see also Zhou et al. (2018) for a large-scale study). The recent work of Jiang et al. (2020) further compared different complexity notions and noted that the ones given by PAC-Bayes tools correlate better with empirical performance. In particular, we will use the following classical result (McAllester, 1999a,b).

**Theorem A.1** (Theorem 1 in McAllester (1999a)). Let  $\mathcal{H}$  be some hypothesis class. Let P be a prior distribution on  $\mathcal{H}$  that is independent of the training set. Let  $Q_S$  be a posterior distribution on  $\mathcal{H}$  that may depend on the training set S. Suppose the loss function is bounded from above by C. With probability  $1 - \delta$  over the randomness of the training set, the following holds

$$\mathop{\mathbb{E}}_{h \sim Q_S} \left[ \mathcal{L}(h) \right] \le \mathop{\mathbb{E}}_{h \sim Q_S} \left[ \hat{\mathcal{L}}(h) \right] + C \sqrt{\frac{\operatorname{KL}(Q_S \mid\mid P) + 3\ln\frac{n}{\delta} + 8}{n}}.$$
(8)

We remark that the original statement in McAllester (1999a) requires the loss function is bounded between 0 and 1. The above statement modifies the original statement and instead applies to a loss function bounded between 0 and C, for some fixed constant C > 0. This is achieved by rescaling the loss by 1/C, leading the C factor in the right hand side of equation (8). To invoke the above result in our setting, we set the prior distribution  $P = \mathcal{N}(\hat{W}^{(s)}, \sigma^2 \operatorname{Id})$ , where  $\hat{W}^{(s)}$  are the weights of the pre-trained nework. The posterior distribution  $Q_S$  is centered at the fine-tuned model as  $\mathcal{N}(\hat{W}, \sigma^2 \operatorname{Id})$ . Based on the above result, we present the proof of Theorem 4.1.

**Proof of Theorem 4.1.** First, we show that the KL divergence between P and  $Q_S$  is equal to  $\frac{1}{2\sigma^2} \|\hat{W} - \hat{W}^{(s)}\|^2$ . We expand the definition using the density of multivariate normal distributions.

$$\begin{aligned} \text{KL}(P \mid\mid Q_S) &= \mathop{\mathbb{E}}_{W \sim \mathcal{P}} \left[ \log \left( \frac{\Pr(W \sim P)}{\Pr(W \sim Q_S)} \right) \right] \\ &= \mathop{\mathbb{E}}_{W \sim \mathcal{P}} \left[ \log \frac{\exp(-\frac{1}{2\sigma^2} \|W - \hat{W}^{(s)}\|^2)}{\exp(-\frac{1}{2\sigma^2} \|W - \hat{W}\|^2)} \right] \\ &= -\frac{1}{2\sigma^2} \mathop{\mathbb{E}}_{W \sim \mathcal{P}} \left[ \|W - \hat{W}^{(s)}\|^2 - \|W - \hat{W}\|^2 \right] \\ &= \frac{1}{2\sigma^2} \mathop{\mathbb{E}}_{W \sim \mathcal{P}} \left[ \langle \hat{W}^{(s)} - \hat{W}, 2W - \hat{W}^{(s)} - \hat{W} \rangle \right] \\ &= \frac{1}{2\sigma^2} \|\hat{W} - \hat{W}^{(s)}\|_F^2 \leq \frac{\sum_{i=1}^L D_i^2}{2\sigma^2}, \end{aligned}$$

where the last line uses the fact that  $\|\hat{W}_i - \hat{W}_i^{(s)}\|_F \leq D_i$ , for all  $1 \leq i \leq L$ .

Next, let the dimension of  $W_i$  be  $d_{i-1}$  times  $d_i$ , for every i = 1, 2, ..., L. In particular,  $d_0$  is equal to the dimension of the input data points. Let  $H = \max_{0 \le i \le L} d_i$  be the maximum matrix dimension size across all the *L* layers. Let  $e = \sigma \sqrt{2H \log(2L \cdot H/\delta)}$ . We show that for any  $\delta > 0$  and any *L*-layer feedforward neural network parameterized by  $W = [W_1, W_2, ..., W_L]$ , with probability at least  $1 - \delta$ , the perturbation  $Q_S$  increases the loss function  $\hat{\mathcal{L}}(h)$  by at most

$$C_1 \cdot e \cdot \left( \sum_{i=1}^{L} \frac{\prod_{j=1}^{L} \left( \|W_j\|_2 + e \right)}{\|W_i\|_2 + e} \right).$$
(9)

Consider some data point (x, y) evaluated at a neural network function  $f_W$ . Let  $U_i$  be a Gaussian perturbation matrix on  $W_i$  with mean zero and entrywise variance  $\sigma^2$ , for i = 1, ..., L. Denote by  $U = [U_1, U_2, ..., U_L]$ . Since  $\ell(x, y)$  is 1-Lipschitz over x, we get

$$|\ell(f_{W+U}(x), y) - \ell(f_W(x), y)| \le ||f_{W+U}(x) - f_W(x)||_2.$$
(10)

Since  $U_i \in \mathbb{R}^{d_{i-1} \times d_i}$  is a random matrix with i.i.d. entries sampled from the standard Gaussian distribution, we can upper bound the operator norm of  $U_i$  by applying well-known concentration results. In particular, by equation (4.1.8) and section (4.2.2) of Tropp (2015), we get

$$\Pr\left[\frac{1}{\sigma}\|U_i\|_2 \ge t\right] \le (d_{i-1} + d_i) \cdot \exp\left(-\frac{-t^2}{2\max(d_{i-1}, d_i)}\right).$$
(11)

Thus, for all i = 1, 2, ..., L, with probability at most  $\delta/L$ , we have that

$$\frac{1}{\sigma} \|U_i\|_2 \le \sqrt{2H \log(2L \cdot H/\delta)}.$$
(12)

In particular, the right hand side above is obtained by setting  $2H \exp\left(-\frac{t^2}{2H}\right) = \delta/L$  and solve for  $t = \sqrt{2H \log(2L \cdot H/\delta)}$ . Since the right hand side of equation (11) is at most  $2H \exp\left(-\frac{t^2}{2H}\right)$ , by union bound, we conclude that with probability at most  $1 - \delta$ , for all i = 1, 2, ..., L, the operator norm of  $U_i$  is at most  $\sigma \cdot \sqrt{2H \log(2L \cdot H/\delta)}$ .

Let  $z_i$  be the input to the *i*-th layer of  $f_W$ . Let  $\tilde{z}_i$  be the input to the *i*-th layer of  $f_{W+U}$ . We can expand the right hand side of equation (10) as

$$\begin{aligned} & \|\psi_L((W_L + U_L)\tilde{z}_L) - \psi_L(W_L z_L)\|_2 \tag{13} \\ & \leq \|(W_L + U_L)\tilde{z}_L - W_L z_L\|_2 \qquad (\text{since } \psi_L(\cdot) \text{ is 1-Lipschitz}) \\ & \leq \|W_L(\tilde{z}_L - z_L)\|_2 + \|U_L\tilde{z}_L\|_2 \qquad (\text{by triangle inequality}) \\ & \leq \|W_L\|_2 \cdot \|\tilde{z}_L - z_L\|_2 + \|U_L\tilde{z}_L\|_2 \qquad (\|Wx\|_2 \leq \|W\|_2 \cdot \|x\|_2 \text{ for any vector } x) \\ & \leq \|W_L\|_2 \cdot \|\tilde{z}_L - z_L\|_2 + \sigma \cdot \sqrt{2H\log(2L \cdot H/\delta)} \cdot \|\tilde{z}_L\|_2, \qquad (14) \end{aligned}$$

where the last step is by equation (12). Recall from Section 3 that the *i*-th layer takes an input z and outputs  $\psi_i(W_i z)$ , for every i = 1, ..., L. Because we have assumed that  $\psi_i(\cdot)$  is 1-Lipschitz, for all i = 1, ..., L - 1, we can bound

$$\begin{split} \|\tilde{z}_{L}\|_{2} &= \phi_{L-1} \circ \phi_{L-2} \cdots \circ \phi_{1}(x) \leq \Big(\prod_{i=1}^{L-1} \|W_{i} + U_{i}\|_{2}\Big) \cdot \|x\|_{2} \\ &\leq \Big(\prod_{i=1}^{L-1} \|W_{i} + U_{i}\|_{2}\Big) \cdot C_{1} \quad \text{(since } \|x\|_{2} \leq C_{1} \text{ for } x \in \mathcal{X}) \\ &\leq \Big(\prod_{i=1}^{L-1} (\|W_{i}\|_{2} + e)\Big) \cdot C_{1}. \quad \text{(by equation (12))} \end{split}$$

Applying the above to equation (14), we have shown

$$\|f_{W+U}(x) - f_W(x)\|_2 \le \|W_L\|_2 \cdot \|\tilde{z}_L - z_L\|_2 + \sigma \cdot \sqrt{2H\log(2L \cdot H/\delta)} \cdot \left(\prod_{i=1}^{L-1} (\|W_i\|_2 + e)\right) \cdot C_1.$$
(15)

Next, we can expand the difference between  $\tilde{z}_L = \psi_{L-1}((W_{L-1} + U_{L-1})\tilde{z}_{L-1})$  and  $z_L = \psi_{L-1}(W_{L-1}z_{L-1})$  similar to equation (13). In general, for any i = 1, 2, ..., L, we get

$$\|\tilde{z}_{i} - z_{i}\|_{2} \le \|W_{i-1}\|_{2} \cdot \|\tilde{z}_{i-1} - z_{i-1}\|_{2} + \sigma \cdot \sqrt{2H\log(2L \cdot H/\delta)} \cdot \left(\prod_{i=1}^{i-2} (\|W_{i}\|_{2} + e)\right) \cdot C_{1}.$$
 (16)

By repeatedly applying equation (16) together with equation (15) (relaxing  $||W_i||_2$  to  $||W_i||_2 + e$ ), we can show that with probability at least  $1 - \delta$ ,

$$\|f_{W+U}(x) - f_W(x)\|_2 \le \sigma \cdot C_1 \sqrt{2H \log(2L \cdot H/\delta)} \cdot \Big(\sum_{i=1}^L \frac{\prod_{j=1}^L (\|W_j\|_2 + e)}{\|W_i\|_2 + e}\Big).$$

Combined with equation (10), we have shown that equation (9) holds. Finally, with equation (9) and the fact that the loss function  $\ell(\cdot)$  is bounded from above by some fixed value  $C_2$ , we can bound the expectation of  $\hat{\mathcal{L}}(h)$  for some  $h \sim Q_S$  as

$$\mathbb{E}_{h\sim Q_S}\left[\hat{\mathcal{L}}(h)\right] \leq \hat{\mathcal{L}}(f_{\hat{W}}) + (1-\delta) \cdot \sigma \cdot C_1 \sqrt{2H \log\left(\frac{2L \cdot H}{\delta}\right)} \cdot \left(\sum_{i=1}^L \frac{\prod_{j=1}^L (\|\hat{W}_j\|_2 + e)}{\|\hat{W}_i\|_2 + e}\right) + \delta \cdot C_2, \quad (17)$$

Note that the above inequality applies for any (small) value of  $\delta$ . Thus, we can minimize the right hand side by appropriately setting  $\delta$ . For our purpose, we will show that there exists some  $\delta$  such that the right hand side of equation (17) (leaving out  $\hat{\mathcal{L}}(f_{\hat{W}})$ ) is at most

$$2\sigma C_1 \Big( \sum_{i=1}^{L} \frac{\prod_{j=1}^{L} (\|\hat{W}_j\|_2 + e)}{\|\hat{W}_i\|_2 + e} \Big) \cdot \sqrt{2H \log \left(4L \cdot H \cdot C_2\right)}.$$
(18)

To see that the above is true, we will abstract away the precise scalars and instead write the right hand side of equation (17) as  $g(\delta) = (1 - \delta)A_1\sqrt{\log(A_2/\delta)} + C_2\delta$ , for some fixed  $A_1$  and  $A_2$  that does not depend on  $\delta$ . We note that subject to

$$C_2\delta \le (1-\delta)A_1\sqrt{\log(A_2/\delta)},\tag{19}$$

 $g(\delta)$  is at most

$$g(\delta) \le 2(1-\delta)A_1\sqrt{\log(A_2/\delta)} \le 2A_1\sqrt{\log\left(\frac{A_2}{\delta}\right)}.$$
(20)

Therefore, we only need to lower bound  $\delta$  based on the constraint (19). In particular, the largest possible  $\delta^*$  under constraint (19) is achieved when both sides equal:

$$C_2 \delta^* = (1 - \delta^*) A_1 \sqrt{\log(A_2/\delta^*)} \ge \frac{A_1}{2} \sqrt{\log(A_2)},$$

which implies that  $\delta^* \ge \frac{A_1}{2C_2}\sqrt{\log(A_2)}$ . Thus, we have shown that  $g(\delta^*)$  is less than equation (18), using equation (20) and the lower bound on  $\delta^*$ . Additionally,  $A_1$  is on the order of a fixed constant based on  $\sigma$  defined below.

Using a similar argument since our bound on the perturbed loss holds point-wise for every  $x \in \mathcal{X}$ , we can likewise prove that with probability  $1 - \delta$ ,

$$\mathbb{E}_{h \sim Q_S} \left[ \mathcal{L}(h) \right] \ge \mathcal{L}(f_{\hat{W}}) - 2\sigma C_1 \Big( \sum_{i=1}^L \frac{\prod_{j=1}^L (\|\hat{W}_j\|_2 + e)}{\|\hat{W}_i\|_2 + e} \Big) \cdot \sqrt{2H \log(4L \cdot H \cdot C_2)}$$

By applying equation (18) and the above to equation (8) in Theorem A.1, we thus conclude that for any  $\sigma$ , with probability at least  $1 - 2\delta$  over the training data set, the following holds

$$\mathcal{L}(f_{\hat{W}}) \leq \hat{\mathcal{L}}(f_{\hat{W}}) + 4\sigma \cdot C_1 \cdot \sqrt{2H \log(4L \cdot H \cdot C_2)} \left( \sum_{i=1}^{L} \frac{\prod_{j=1}^{L} (\|\hat{W}_j\|_{2}+e)}{\|\hat{W}_i\|_{2}+e} \right) + C_2 \sqrt{\frac{\sum_{i=1}^{L} \frac{D_i^2}{2\sigma^2} + 3\ln\frac{n(t)}{\delta} + 8}{n^{(t)}}}.$$
 (21)

Since  $\|\hat{W}_i^{(s)}\|_2 \leq B_i$  and  $\|\hat{W}_i - \hat{W}_i^{(s)}\|_F \leq D_i$ , for any i = 1, ..., L, we have that  $\|\hat{W}_i\|_2 \leq B_i + D_i$ . Thus, we can upper bound the above by replacing every  $\|\hat{W}_i\|_2$  with  $B_i + D_i$ .

By setting  $\sigma = \frac{\varepsilon}{6\sigma C_1 \cdot \alpha \sqrt{2H \log(4LHC_2)}}$  where  $\alpha = \left(\sum_{i=1}^L \frac{\prod_{j=1}^L (B_j + D_j)}{B_i + D_i}\right)$ , we get  $e \leq \frac{1}{6\alpha}$  since  $C_1 \geq 1$ . To finish the proof, we will show

$$\sum_{i=1}^{L} \frac{\prod_{j=1}^{L} (B_j + D_j + e)}{B_i + D_i + e} \le 3\alpha/2.$$

In particular, we will show that for every  $k = 1, 2, ..., L, e \leq \frac{B_k + D_k}{6L}$ . To see this, recall that  $B_k \geq 1$  and  $D_k \geq 0$ . Therefore,  $\alpha \geq L$  for every k and  $e \leq \frac{1}{6\alpha} \leq \frac{B_k + D_k}{2}$ . This implies

$$\sum_{i=1}^{L} \frac{\prod_{j=1}^{L} (B_j + D_j + e)}{(B_i + D_i + e)} \le \left(1 + \frac{1}{6L}\right)^{L-1} \sum_{i=1}^{L} \frac{\prod_{j=1}^{L} (B_j + D_j)}{B_i + D_i} \le \frac{3}{2}\alpha$$

Detecate	Mathada		independe	ent noise	
Datasets	Methous	20%	40%	60%	80%
	Fine-tuning	$68.28 \pm 0.34$	$56.88 \pm 0.38$	$39.54 \pm 0.42$	$16.21\pm0.38$
	LS	$69.89 \pm 0.85$	$59.18\pm0.43$	$41.58\pm0.45$	$16.96\pm0.45$
	$\ell^2$ -PGM	$69.71\pm0.31$	$58.91 \pm 0.42$	$41.52\pm0.74$	$16.76\pm0.41$
CUB-200-2011	GCE	$69.54 \pm 0.25$	$60.15\pm0.93$	$41.84\pm0.47$	$17.77\pm0.32$
	SCE	$70.68\pm0.67$	$61.33\pm0.57$	$43.67\pm0.37$	$17.62\pm0.19$
	DAC	$68.58 \pm 0.25$	$57.37 \pm 0.29$	$40.00\pm0.37$	$15.92\pm0.83$
	APL	$70.59\pm0.18$	$59.68 \pm 0.69$	$40.46\pm0.24$	$14.69\pm0.36$
	SAT	$68.69 \pm 0.38$	$57.34 \pm 0.18$	$38.75\pm0.40$	$15.14\pm0.20$
	ELR	$69.92\pm0.14$	$58.69 \pm 0.02$	$40.76\pm0.68$	$16.68\pm0.79$
	REGSL (ours)	$\textbf{71.76} \pm \textbf{0.49}$	$\textbf{62.79} \pm \textbf{0.23}$	$\textbf{45.85} \pm \textbf{0.66}$	$\textbf{17.88} \pm \textbf{0.25}$
	Fine-tuning	$83.13\pm0.15$	$72.23\pm0.40$	$55.27\pm0.32$	$29.35\pm0.74$
	LS	$83.62\pm0.30$	$72.35\pm0.47$	$54.23\pm0.24$	$28.60\pm0.45$
	$\ell^2$ -PGM	$83.45\pm0.70$	$73.24\pm0.20$	$56.51\pm0.40$	$31.04\pm0.88$
	GCE	$83.09\pm0.22$	$71.74\pm0.65$	$54.73\pm0.90$	$28.38 \pm 0.95$
Flowers	SCE	$83.45\pm0.14$	$73.11\pm0.05$	$55.96 \pm 0.97$	$29.22\pm0.07$
Flowers	DAC	$83.40\pm0.27$	$72.73\pm0.34$	$55.27\pm0.40$	$28.95\pm0.89$
	APL	$83.42\pm0.11$	$72.06\pm0.23$	$54.96\pm0.36$	$28.86\pm0.64$
	SAT	$82.49\pm0.25$	$72.52\pm0.28$	$55.10\pm0.50$	$27.82\pm0.38$
	ELR	$83.38\pm0.38$	$72.53\pm0.48$	$55.74\pm0.42$	$30.17\pm0.26$
	REGSL (ours)	$\textbf{83.79} \pm \textbf{0.39}$	$\textbf{73.32} \pm \textbf{0.66}$	$\textbf{57.52} \pm \textbf{0.15}$	$\textbf{31.09} \pm \textbf{0.25}$

Table 5: Top-1 test accuracy on CUB-200-2011 and Flowers dataset with independent label noise injected in the training set. Results are averaged over 3 random seeds.

Thus, plugging in the value of  $\sigma$  above, we conclude that the second term in equation (21) is at most  $\varepsilon$ . By applying the value of  $\sigma$  to the third term in equation (21), we conclude

$$\mathcal{L}(f_{\hat{W}}) \leq \hat{\mathcal{L}}(f_{\hat{W}}) + \varepsilon + C_2 \sqrt{\frac{\frac{36}{\varepsilon^2} \cdot C_1^2 H \log(4LHC_2) \left(\sum_{i=1}^L \frac{\prod_{j=1}^L (B_j + D_j)}{B_i + D_i}\right)^2 \left(\sum_{i=1}^L D_i^2\right) + 3\ln\frac{n^{(t)}}{\delta} + 8}{n^{(t)}}}$$

Thus, we have proved equation (4) is true. The proof is complete.

*Implication.* Our result implies a tradeoff in setting the distance parameter  $D_i$ . While a larger  $D_i$  increases the network's "capacity", the generalization bound gets worse as a result.

### **B** Extended experimental setup and results

#### **B.1** Additional details of the experimental setup

We extend the details of our experiment setup. First, we introduce the data sets in our experiments and describe the pre-processing steps. Second, we present the model architectures and the fine-tuning procedures. Third, we describe the baselines in detail. Finally, we describe the implementation and the hyperparameters for our algorithm and the baselines.

**Data sets.** We evaluate fine-tuning on seven image classification data sets covering multiple applications, including fine-grained object recognition, scene recognition, and general object recognition. We divide each dataset into the training set, the validation set, and the test set for data split. We adopt the standard splitting of data given in the Aircrafts dataset (Maji et al., 2013). For Caltech-256 dataset (Griffin et al., 2007), we use the setting in Li et al. (2018b) that randomly samples 30, 20, 20 images of each class for training, validation, and test set, respectively. For other data sets, we split 10% of the training set as the validation set and use the standard test set. We describe the statistics of the seven data sets in Table 1. When fine-tuning on the seven data sets, we preserve the original pixel and resize the shorter side to 256 pixels. The image samples are normalized with the mean and std values over ImageNet data (Russakovsky et al., 2015). Moreover, we apply commonly used

Table 6:	Left: M	lean AUF	ROC of a	fine-tunin	g ResN	et-18 o	n the (	ChestX	-ray14	data s	et. Rigl	nt: 🛛	Гор-1
accuracy	of fine-	-tuning V	'iT on th	e indoor	lata set.	Result	s are a	average	d over	3 rand	lom see	ds.	

(a) The Chest	X-ray14 data set.	(b)	(b) The indoor data set.			
Methods Mean AUROC		Methods	Independent Noise			
Fine-tuning	$0.8159 \pm 0.0667$		40%	80%		
$\ell^2$ -Norm LS $\ell^2$ -SP $\ell^2$ PCM	$\begin{array}{c} 0.8198 \pm 0.0644 \\ 0.7885 \pm 0.0578 \\ 0.8231 \pm 0.0658 \\ 0.8235 \pm 0.0626 \end{array}$	Fine-tuning Regularization Self-labeling	$\begin{array}{c} 75.87 \pm 1.15 \\ 82.66 \pm 0.44 \\ 79.13 \pm 0.37 \end{array}$	$\begin{array}{c} 32.06 \pm 3.17 \\ 63.01 \pm 0.83 \\ 41.69 \pm 0.16 \end{array}$		
$\frac{\ell \text{ -FGM}}{\text{REGSL (ours)}}$	$0.8233 \pm 0.0030$ $0.8274 \pm 0.0654$	REGSL (ours)	$\textbf{83.48} \pm \textbf{0.29}$	$64.50 \pm 0.53$		

data augmentation methods on the training samples, including random scale cropping and random flipping. The images are resized to  $224 \times 224$  as the input for the model. The training and batch size are both 16.

To test the robustness of fine-tuning, we evaluate Algorithm 1 under two different settings with label noise. We select MIT-Indoor dataset (Sharif Razavian et al., 2014) as the benchmark dataset and randomly flipped the labels of the training samples. The results for the other datasets are similar, which can be found in Appendix B.2. We consider two scenarios of label noise in our experiments, independent random noise and correlated noise. Random label noise is generated by flipping the labels of a given proportion of training samples to other classes uniformly. For the correlated noise setting, the label noise is dependent on the sample. We simulate the correlated noisy label by training an auxiliary network on a held-out dataset to a certain accuracy. We then use the prediction of the auxiliary network as noisy labels. For few-shot learning, we conduct experiments on a few-shot image recognition benchmark, miniImageNet (Vinyals et al., 2016) and follow the meta training and meta test setup in Tian et al. (2020).

Architectures. For image data sets, we use ResNet-101 (He et al., 2016) network which is pre-trained on ImageNet dataset (Russakovsky et al., 2015). For the label-noise experiments, we mainly use ResNet-18 (He et al., 2016) network. In the transfer learning and label-noise experiments, we finetune the model with Adam optimizer (Kingma and Ba, 2014) with an initial learning rate 0.0001 for 30 epochs and decay the learning rate by 0.1 every 10 epochs. We report the average Top-1 accuracy on the test set of 3 random seeds. For the few-shot learning tasks, we use ResNet-12 pre-trained on the meta-training dataset as in the previous work (Tian et al., 2020). To fine-tune on the meta-test sets, we use Adam optimizer (Kingma and Ba, 2014) with an initial learning rate  $5e^{-5}$  and fine-tune the model on the training set for 25 epochs. We report the average classification accuracies on 600 sampled tasks from the meta-test split.

**Baselines.** For the transfer learning tasks, we focus on using the Frobenius norm for distance regularization. We present ablation studies to compare the Frobenius norm and other norms in Appendix B.2. we compare our algorithm with fine-tuning with early stop (Fine-tuning), fine-tuning with weight decay ( $\ell^2$ -Norm), label smoothing (LS) method formulated in (Müller et al., 2019),  $\ell^2$ -SP (Li et al., 2018b), and  $\ell^2$ -PGM (Gouk et al., 2021). For testing the robustness of our algorithm, in addition to the baselines described above, we adopt several baselines that have been proposed in the supervised learning setting, including GCE (Zhang and Sabuncu, 2018), SCE (Wang et al., 2019), DAC (Thulasidasan et al., 2019), APL (Ma et al., 2020), ELR (Liu et al., 2020) and self-adaptive training (SAT) (Huang et al., 2020). We compare with the same fine-tuning baselines in transfer learning for the few-shot classification tasks.

**Implementation and hyperparameters.** For all baselines and data sets, all regularization hyperparameters are searched on the validation dataset by the Optuna optimization framework (Akiba et al., 2019). In our proposed algorithm, we search the constraint distance D in [0.05, 10] and the distance scale factor  $\gamma$  in [1, 5] by sampling. For setting layer-wise constraints, we set four different constraints for the four blocks of the ResNet. Concretely, ResNet-101 has four blocks, with 10, 13, 67, and 10 convolutional layers for each block, respectively. For every block, we set the same constraint value for all the layers within the block. Therefore, there are four distance values  $(D, D \cdot \gamma, D \cdot \gamma^2, D \cdot \gamma^3)$  each block. Similar constraints are set for ResNet-18 and ResNet-12 models in the experiments. The search space for other hyperparameters is shown as follows. Re-weight start step  $E_r$  is searched in  $\{3, 5, 8, 10, 13\}$ . Re-weight temperature factor  $\gamma$  is searched in  $\{3.0, 2.0, 1.5, 1.0\}$ . Label correction

Noise Rate 0%	MR	CR	MPQA	SUBJ	TREC
Fine-tuning $\ell^2$ -PGM REGSL (ours)	83.37±0.70 84.16±0.41 <b>84.20±0.47</b>	83.29±0.80 83.87±0.66 <b>84.35±0.60</b>	87.56±0.70 87.77±0.62 <b>87.95±0.65</b>	93.14±0.42 93.16±0.21 <b>93.50±0.17</b>	83.28±0.86 <b>84.48±0.52</b> 83.73±0.75
Noise Rate 40%	MR	CR	MPQA	SUBJ	TREC
Fine-tuning $\ell^2$ -PGM REGSL (ours)	82.91±0.25 83.45±0.28 <b>83.54±0.40</b>	77.67±1.11 79.47±0.37 <b>79.89</b> ± <b>0.51</b>	82.85±0.98 84.03±0.41 <b>84.15±0.99</b>	90.78±0.88 72.14±0.21 <b>72.48±0.45</b>	$70.80{\pm}0.22 \\73.48{\pm}0.90 \\\textbf{74.80{\pm}0.87}$

Table 7: Top-1 accuracy of fine-tuning a three-layer feed forward network pre-trained on SST with no label noise and with independent label noise. Results are averaged over 5 random seeds.

Table 8: Comparing  $\sum_{i=1}^{L} D_i^2$  between layer-wise and constant regularization.

	Aircrafts	CUB-200-2011	Caltech-256	Stanford-Cars	Stanford-Dogs	Flowers	MIT-Indoor
Layer-wise	250.20	840.09	90.42	4150.02	20.72	252.45	57.16
Constant	2465.05	1225.69	195.72	2480.08	211.78	3151.28	295.49

start step  $E_c$  is searched in {5, 8, 10, 13, 15}. Label correction threshold p is set as 0.90 in all experiments. The validation set sizes used for evaluating these hyperparameters choices range from 536 to 5120 (cf. Table 1). For the results in Table 2, we search for 20 different trials of the hyperparameters for the proposed algorithm and the baselines. For the results in Table 3, we search 20 times on the self-labeling parameters and 20 times on the regularization parameters. For the few-show learning experiments in Table 4, we use the validation set of Mini-ImageNet, following the training procedure of Dhillon et al. (2019) and search for 20 different trials.

We report the results for baselines by running the official open-sourced implementations. We describe the hyperparameters space for baselines as follows. For GCE (Zhang and Sabuncu, 2018), we search the factor q in their proposed truncated loss  $\mathcal{L}_q$  in {0.4, 0.8, 1.0} and set the factor k as 0.5. We also search the start pruning epoch in {3, 5, 8, 10, 13}. For SCE (Wang et al., 2019), we search the  $\alpha$  and A in their proposed symmetric cross entropy loss. The factor  $\alpha$  is searched in {0.01, 0.10, 0.50, 1.00}, and the factor A is searched in {-2, -4, -6, -8}. For DAC (Thulasidasan et al., 2019), we search the start abstention epoch in {3, 5, 8, 10, 15} and set the other hyperparameters as the same in their paper. For APL (Ma et al., 2020), we choose the active loss as normalized cross entropy loss and the passive loss as reversed cross entropy loss. We search the loss factor  $\alpha$  in {1, 10, 100} and  $\beta$  in {0.1, 1.0, 10}. For ELR (Liu et al., 2020), we search the momuntum factor  $\beta$  in {0.5, 0.7, 0.9, 0.99} and the weight factor  $\lambda$  in {0.05, 0.3, 0.5, 0.7}. For SAT (Huang et al., 2020), the start epoch is searched in {3, 5, 8, 10, 13}, and the momuntum is searched in {0.6, 0.8, 0.9, 0.99}. We use the same number of trials in tuning hyperparameters for baselines.

#### **B.2** Extended experimental results

We present additional results. First, we apply our algorithm to another data set with label noise (same label noise as Table 3). Second, we apply our algorithm to sentence classification tasks.

Additional results under label noise. We report the test accuracy results of our algorithm on CUB-200-2011 and Flowers data sets with independent label noise in Table 5. From the table, we show that our algorithm still outperforms previous methods by a significant margin, which aligns with our results of the MIT-Indoor data set in Table 3.

**Results on ChestX-ray14.** We apply our approach to medical image classification, which is a more distant transfer task (relative to the source data set ImageNet). We report the mean AUROC (averaged over predicting all 14 labels) on the ChestX-ray14 data set in Table 6a. With layer-wise regularization, we see **0.39**% improvement compared to the baseline methods.

**Results on ViT.** We apply our approach to fine-tuning pre-trained ViT (Dosovitskiy et al., 2020) from noisy labels, under the same setting as Table 3). Table 6b shows the result. First, we find that our approach improves upon the best regularization methods by 1.17% averaged over two settings. Second, we find that our approach also improves upon self-labeling by 13.57% averaged over two



Figure 6: Frobenious norms of each layer in ResNet-18 pre-trained with standard and adversarial training methods from different noise levels ( $\epsilon = 0.01, 0.25, 5$ ). Higher  $\epsilon$  implies lower norms.

settings. These two results again highlight that both regularization and self-labeling contribute to the final result. While regularization prevents the model from over-fitting to the random labels, self-labeling injects the belief of the fine-tuned model into the noisy data set.

**Extension to text classification.** We apply our algorithm in text data domains. We conduct an experiment on sentiment classification tasks using a three-layer feedforward neural network (or multi-layer perceptron). We considered six text classification data sets: SST (Socher et al., 2013), MR (Pang and Lee, 2005), CR (Hu and Liu, 2004), MPQA (Wiebe et al., 2005), SUBJ (Pang and Lee, 2004), and TREC (Li and Roth, 2002). We use SST as the source task and one of the other tasks as the target task. We compared our proposed algorithm with fine-tuning and L2-PGM (Gouk et al., 2021). In Table 7, we report the results evaluated under a setting with no label noise and a setting with independent label noise (same setting in Table 3). The results show our proposed algorithm can outperform these baseline methods under both settings.

#### **B.3** Ablation studies

We describe several ablation studies. First, we compare the performance between using the Frobenius norm and the MARS norm proposed in Gouk et al. (2021). Second, we show the norm values of the adversarial robust pre-trained models in each layer and explain their improvement for fine-tuning on downstream tasks. Third, we compare our proposed layer-wise constraints and a constant and discuss their correlation to the generalization performance in Theorem 4.1. Fourth, we analyze the role of regularization in our proposed self label-correction method. Finally, we conduct an ablation study to study the influence of three components in our algorithm.

**Comparing layer-wise and constant regularization:**  $\sum_{i=1}^{L} D_i^2$ . To show the bound in Proposition 4.1 is nonvacuous, We compare the value of  $\sum_{i=1}^{L} D_i^2$  from Equation 4 under both layer-wise constraints and constant regularization in Table 8. The values for every  $D_i$  are from the parameters set in Table 2 for these data sets. We can see that layer-wise constraints incur a smaller value of  $\sum_{i=1}^{L} D_i^2$  compared to constant regularization, which correlates with their better generalization performance.

**Comparing norms of adversarial robust pre-trained models.** We empirically observe that adversarially robust pre-trained models have lower Frobenius norms over each layer shown in Figure 6. The lower norms of robust models would induce smaller generation error, which explains their improvement in fine-tuning downstream tasks.

**Comparing label-correction precision with and without regularization.** We observe from Figure 5a that combining self-labeling and regularization is more effective than using just self-labeling. This is evidenced by the gap between correction with regularization and without regularization in the figure. In Table 9, we show the denominator of relabeling precision (i.e., the overall number of data points that are relabeled in line 7 of Algorithm 1) in the correction process. We find that the denominator is different between with and without regularization. The denominator is significantly smaller with regularization than without regularization, indicating that regularization prevents the model from over-fitting the noisy labels. The denominator decreases only with regularization. Without

Table 9: Denominator (total number of relabeled data points) and numerator (the number of correctly relabeled data) for calculating precision of label-correction in Figure 5a.

			0				
1	4	7	10	13	16	19	21
28	32	32	29	26	24	21	19
15	17	16	13	12	10	9	8
56	66	64	64	62	58	59	56
34	33	28	26	21	18	17	14
	1 28 15 56 34	1         4           28         32           15         17           56         66           34         33	1         4         7           28         32         32           15         17         16           56         66         64           34         33         28	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1       4       7       10       13       16         28       32       32       29       26       24         15       17       16       13       12       10         56       66       64       64       62       58         34       33       28       26       21       18	1       4       7       10       13       16       19         28       32       32       29       26       24       21         15       17       16       13       12       10       9         56       66       64       64       62       58       59         34       33       28       26       21       18       17

Table 10: Top-1 test accuracy of using different norms for fine-tuning ResNet-101 pre-trained on the ILSVRC-2012 subset of ImageNet. Results are averaged over 3 random seeds.

	Aircrafts	CUB-200-2011	Caltech-256	Stanford-Cars	Stanford-Dogs	Flowers	MIT-Indoor		
PGM (MARS)	$74.34 {\pm} 0.15$	$81.14 {\pm} 0.16$	$83.28 {\pm} 0.31$	$89.01 {\pm} 0.18$	$87.90 {\pm} 0.10$	93.45±0.23	$78.48 {\pm} 0.29$		
$PGM(\ell^2)$	$74.90 {\pm} 0.26$	$81.23 \pm 0.32$	$83.25 \pm 0.33$	$88.92 \pm 0.39$	$86.48 \pm 0.28$	$93.23 \pm 0.34$	$77.31 \pm 0.30$		
Ours (MARS)	$75.10 \pm 0.32$	$82.03 \pm 0.32$	$85.33 \pm 0.33$	$89.29 \pm 0.25$	$90.94 \pm 0.28$	$93.67 \pm 0.34$	$79.40 {\pm} 0.30$		
Ours $(\ell^2)$	$75.32 {\pm} 0.23$	$82.24 \pm 0.21$	$84.90 {\pm} 0.16$	$89.14 {\pm} 0.22$	$89.58 {\pm} 0.13$	$93.82 {\pm} 0.35$	$79.30 {\pm} 0.31$		

Table 11: Removing any component in our algorithm leads to worst performance. Results are on the Indoor data set with independent label noise. Results are averaged over 3 random seeds.

Indoor	20%	independ 40%	lent noise 60%	80%	correlated noise 25.18%
REGSL (ours)	72.51±0.46	68.13±0.16	57.59±0.55	34.08±0.79	70.12±0.83
w/o regularization	$71.94{\pm}0.43$	$67.84{\pm}0.38$	$57.24 \pm 0.43$	$33.78 {\pm} 0.30$	$69.43 {\pm} 0.36$
w/o label correction	$70.92 {\pm} 0.41$	$59.10 {\pm} 0.24$	$47.81 {\pm} 0.35$	$28.42 \pm 0.46$	$69.78 {\pm} 0.34$
w/o label removal	$70.32 {\pm} 0.65$	$66.57 {\pm} 0.76$	$55.37 {\pm} 0.28$	$29.43 {\pm} 0.88$	$67.96 {\pm} 0.49$
w/o self-labeling	$70.23 {\pm} 0.25$	$64.40 {\pm} 0.58$	$54.20{\pm}0.68$	$32.54{\pm}0.43$	$69.05 {\pm} 0.09$

regularization, the denominator remains roughly the same. Additionally, the number of incorrectly labeled data points by the self-labeling is much higher (e.g., 42 vs. 11 at epoch 21).

**Comparing**  $\ell_2$  **norm and the MARS norm.** We report the results of using different norms for constraints in our algorithm in Table 10. We compare the performance of the Frobenius norm ( $\ell^2$ ), and the MARS norm proposed in (Gouk et al., 2021). From the table, we show that the results of the two norms are similar. In the paper, we focus on the Frobenius norm for our discussion.

**Influence of different components in** REGSL. We study the influence of each component of our algorithm: layer-wise constraint, label correction, and label removal. We remove these components, respectively, and run the same experiments on the MIT-Indoor dataset with different kinds of label noise. Furthermore, we also include a row of results using only layer-wise constraints without self-labeling (containing label correction and removal). As shown in Table 11, removing any component from our algorithm could hurt the performance. This suggests that only incorporating these components can prevent both over-fitting and label memorization of models. In addition, We can see from Table 11 that when the noise rate is 20%, the self-labeling part (including label correction and removal) is more critical than regularization. When the noise rate is 40% or higher, the label correction part is the most important.

# C Limitation

Our work focuses on the regularization and robustness properties of fine-tuning. We mainly focus on feedforward neural networks, including convolutional neural networks. Thus, we have not considered Other settings such as recurrent neural networks in this work. We mainly focus on independent label noise for evaluating the robustness. Other types of label noise, for example, using weak supervision, are not well-understood. This is left for future work. Finally, our treatment of label noise using self-labeling is based on empirical heuristics. It would be fascinating to come up with a more principled approach to tackle label noise.

# Checklist

- 1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes]
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes]
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
  - (b) Did you include complete proofs of all theoretical results? [Yes]
- 3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes]
  - (b) Did you mention the license of the assets? [Yes]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
- 5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]