

---

# Informed Learning by Wide Neural Networks: Convergence, Generalization and Sampling Complexity

---

Jianyi Yang<sup>1</sup> Shaolei Ren<sup>1</sup>

## Abstract

By integrating domain knowledge with labeled samples, informed machine learning has been emerging to improve the learning performance for a wide range of applications. Nonetheless, rigorous understanding of the role of injected domain knowledge has been under-explored. In this paper, we consider an informed deep neural network (DNN) with over-parameterization and domain knowledge integrated into its training objective function, and study how and why domain knowledge benefits the performance. Concretely, we quantitatively demonstrate the two benefits of domain knowledge in informed learning — regularizing the label-based supervision and supplementing the labeled samples — and reveal the trade-off between label and knowledge imperfectness in the bound of the population risk. Based on the theoretical analysis, we propose a generalized informed training objective to better exploit the benefits of knowledge and balance the label and knowledge imperfectness, which is validated by the population risk bound. Our analysis on sampling complexity sheds lights on how to choose the hyper-parameters for informed learning, and further justifies the advantages of knowledge informed learning.

## 1. Introduction

The remarkable success of deep neural networks (DNNs), or more generally machine learning, largely relies on the proliferation of data samples with ground-truth labels for supervised learning. Nonetheless, labeled data of high quality can often be very limited and/or extremely expensive to collect in real application domains, including medical

sciences, security-related fields, and specialized engineering areas (von Rueden et al., 2021).

In parallel with the data-driven learning paradigm, domain *knowledge* (which we simply refer to as knowledge) has been utilized to assist with decision making and system designs, with a long history of success. As its name would suggest, domain knowledge is naturally domain-specific and can come from various sources in multiple forms, such as subjective experiences (e.g., medical prognosis), external sources, and scientific laws. For example, partial differential equations are used to govern many flow dynamics in physics, and the Shannon channel capacity is the fundamental principle to guide the design of modern communications systems (Goldsmith, 2005; Willard et al., 2020).

Importantly, domain knowledge has already been, sometimes implicitly, integrated into every stage of the machine learning pipeline, including training data augmentation, hypothesis set selection, model training and hypothesis finalization (more details in Appendix E). For example, differential equations and logic rules from physical sciences and/or common knowledge provide additional constraints or new functional regularization terms for model training (Battaglia et al., 2016; Borghesi et al., 2020; Silvestri et al., 2021; Muralidhar et al., 2018; Xu et al., 2018).

Despite the numerous successful examples (von Rueden et al., 2021; Deng et al., 2020), there still lacks a rigorous understanding of the role of domain knowledge in informed learning. In this paper, we focus on informed DNNs — DNNs with domain knowledge explicitly integrated into the training risk/loss function. Concretely, we consider an over-parameterized DNN with a sufficiently large network width (Neyshabur et al., 2018), and study how domain knowledge affects the DNN from three complementary aspects: convergence, generalization, and sampling complexity.

**Convergence (Theorem 4.1):** We show the convergence of training an informed risk function under milder technical assumptions than the prior works (Section 4.1). More specifically, we show that for inputs within a smooth set (Definition 1), the network outputs converge to the optimal solution jointly determined by all the samples in the set.

**Generalization (Theorems 4.2 and 5.1):** We show in Theo-

---

<sup>1</sup>Department of Electrical and Computer Engineering, University of California, Riverside, CA 92521, United States. Correspondence to: Shaolei Ren <sren@ece.ucr.edu>.

rem 4.2 that the population risk relies on the knowledge imperfectness (Definition 3) as well as knowledge-regularized label imperfectness (Definition 4). Specifically, knowledge has two benefits: regularization for noisy labels and supplementing labels. We propose a generalized informed risk function which disentangles the two effects by introducing another hyper-weight  $\beta$ , followed by the population risk bounds in Theorem 5.1 and Corollary 5.2.

**Sampling Complexity (Corollary 5.3):** By establishing a quantitative equivalence between domain knowledge and labeled samples, we show that domain knowledge (with a reasonable quality) can effectively reduce the number of labeled samples while achieving the same generalization performance, compared to the no-knowledge case.

## 2. Related Work

**Informed Machine Learning.** The broad paradigm of informed machine learning (von Rueden et al., 2021) includes several existing learning frameworks, such as learning using privileged information (LUPI) (Vapnik & Vashist, 2009) where side knowledge is available for labeled samples (Vapnik & Vashist, 2009; Motiian et al., 2016; Sharmanska et al., 2013). Likewise, knowledge distillation (Rahbar et al., 2020; Hinton et al., 2014; Gou et al., 2021; Cho & Hariharan, 2019) transfers prior knowledge from teacher networks to a student network. Some recent studies have also focused on understanding knowledge distillation (Allen-Zhu & Li, 2020). In (Phuong & Lampert, 2019), a generalization bound is derived for knowledge distillation based on linear classifiers and deep linear classifiers, providing insights towards the mechanism of knowledge distillation. The subsequent analysis (Ji & Zhu, 2020; Rahbar et al., 2020) extends to neural networks, showing that the student network may generalize better by exploiting soft labels from the teacher model. Teacher imperfectness is investigated in (Dao et al., 2021), which bounds the learning error and proposes enhanced methods to address imperfect teachers.

Physics-informed neural networks (PINNs) have been recently proposed to solve partial differential equations (PDEs) (Yin et al., 2021; Institute, 2020; Raissi et al., 2017; Baker et al., 2019; Deng et al., 2020; Willard et al., 2020). Besides empirical studies, (Shin et al., 2020) bounds the expected PINN loss, showing that the minimizer of the regularized loss converges to the PDE solution.

More broadly, informed machine learning also includes weakly-supervised learning (Zhou, 2018; Robinson et al., 2020) and few-shot learning (Wang et al., 2020), where knowledge provides weak supervision. Domain-specific constraints (Muralidhar et al., 2018) and semantic information (Xu et al., 2018; Diligenti et al., 2017a) can also be viewed as knowledge injected into training. Our work com-

plements these empirical studies and provides a rigorous understanding of knowledge in a unified framework.

**Over-parameterized neural networks.** Several recent studies (Bahri et al., 2021; Song et al., 2021; Gao et al., 2021; Khanduri et al., 2021; Jacot et al., 2018; Lee et al., 2019; Yang, 2019; Allen-Zhu et al., 2019b; Arora et al., 2019b;a; Cao & Gu, 2019; Allen-Zhu et al., 2019a; Neyshabur et al., 2018) show that over-parameterized neural networks have good convergence and generalization performance. In addition to assuming data separability in a strong sense, another crucial assumption often made in the existing studies is that the network widths increase polynomially with the total number of training samples. In informed DNNs, however, we can have many (unlabeled) training samples fed into the knowledge risk, which hence may not satisfy these assumptions. Thus, we analyze knowledge-informed over-parameterized neural networks under relaxed assumptions (Section 4).

**Regularization.** In the broad context of regularization, (Wei et al., 2019) shows that over-parameterized neural networks with  $l_2$ -regularization can achieve a larger margin and thus better generalization, (Blanc et al., 2020) proves that SGD with label noise is equivalent to an implicit regularization term, while (Wei et al., 2020) shows that the drop-out operation for neural networks has both explicit and implicit regularization effects. These regularizers are usually imposed on the network weights, whereas the knowledge-based regularizer in informed machine learning also incorporates inputs and directly regularizes the network output.

## 3. Informed Neural Network

**Notations:** We use the expression  $[L]$  to denote the set  $\{1, 2, \dots, L\}$  for a positive integer  $L$ . Denote the indicator function as  $\mathbb{1}(x) = 1$  if  $x > 0$ , and  $\mathbb{1}(x) = 0$  otherwise.  $\mathbb{E}$  is the expectation operator and  $\mathbb{P}$  is a probability measure.  $\mathbb{R}^d$  is  $d$ -dimensional real number space.  $\mathcal{N}(x, \sigma^2)$  is the Gaussian distribution with mean  $x$  and variance  $\sigma^2$ . Denote  $|\mathcal{A}|$  as the size of a set  $\mathcal{A}$ . For a vector  $x$ ,  $\|x\|$  is  $l_2$ -norm and  $[x]_j$  is the  $j$ th entry. For a matrix  $\mathbf{X}$ ,  $\|\mathbf{X}\|_2$  represents the spectral norm, and  $\|\mathbf{X}\|$  is the Frobenius norm.  $\mathcal{B}(x, \tau) = \{y \mid \|x - y\| \leq \tau\}$  is the neighborhood domain.

### 3.1. Preliminaries of Neural Networks

Consider a supervised learning task to learn a relationship mapping the input  $x \in \mathcal{X} \subseteq \mathbb{R}^b$  to its output  $y \in \mathcal{Y} \subseteq \mathbb{R}^d$ . The pair of input and output  $(x, y)$  follows a joint distribution  $\mathbb{P}_{XY}$ . More concretely, we consider a fully-connected DNN with an input layer,  $L \geq 1$  hidden layers, and an output layer. Each hidden layer has  $m$  neurons, followed by ReLU activation denoted as  $\sigma(\cdot)$ . Denote  $\mathbf{W}_0 \in \mathbb{R}^{b \times m}$  as the weights for the input layer,  $\mathbf{W}_l \in \mathbb{R}^{m \times m}$  as

the weights for the  $l$ -th layer for  $l \in [L]$ , and  $\mathbf{V} \in \mathbb{R}^{d \times m}$  as the weights for the output layer. We denote the output of the  $l$ -th layer as  $h_l = \sigma(\mathbf{W}_l h_{l-1})$ , for  $l \in [L]$ , where  $h_0$  is the input  $x$ . The output of the neural network can be expressed as  $h_{\mathbf{W}} = \mathbf{V} h_L$ , where  $\mathbf{W} = \{\mathbf{W}_0, \mathbf{W}_1, \dots, \mathbf{W}_L\}$ . Thus, the DNN can be expressed as

$$h_{\mathbf{W}}(x) = \mathbf{V} \sigma(\mathbf{W}_L \sigma(\mathbf{W}_{L-1} \dots \sigma(\mathbf{W}_1 \sigma(\mathbf{W}_0 x))))). \quad (1)$$

Given a DNN  $h_{\mathbf{W}}$ , the risk for a labeled sample  $(x, y)$  is denoted as  $r(h_{\mathbf{W}}(x), y)$ . The goal of the learning task is to learn a DNN that minimizes the population risk:

$$R(h) = \mathbb{E}[r(h(x), y)]. \quad (2)$$

### 3.2. Integration of Knowledge

We consider a commonly-used informed learning method, i.e., integrating knowledge into the neural network during the training stage (von Rueden et al., 2021). During training, a labeled dataset  $S_z = \{(x_1, z_1), \dots, (x_{n_z}, z_{n_z})\}$  with  $n_z$  samples drawn from  $\mathbb{P}_{XZ}$  is provided. We assume  $x_i, i \in [n]$  are drawn from the distribution  $\mathbb{P}_X$ , but the training label  $z_i \in \mathcal{Y}$  may not be the same as the true label  $y_i$  for the input  $x_i$ , because the training label may be of low quality (e.g., corrupted, noisy, and/or quantized) (Cannings et al., 2020; Zhou, 2018). Denote  $h_{\mathbf{W},i} = h_{\mathbf{W}}(x_i)$  as the output of the neural network with respect to the input  $x_i$ . Based on the labeled dataset, the empirical label-based risk can be written as  $\hat{R}_{S_z}(\mathbf{W}) = \frac{1}{n_z} \sum_{S_z} r(h_{\mathbf{W},i}, z_i)$ .

The domain knowledge includes a knowledge-based model  $g(x)$  regarding the input  $x$  and a knowledge-based risk function  $r_K(h_{\mathbf{W}}(x), g(x))$  that relates the DNN's output  $h_{\mathbf{W}}(x)$  to  $g(x)$ . More concrete examples of risk functions for domain knowledge can be found in Appendix F.

For the ease of analysis, we assume that both the risk function  $r$  and the knowledge-based risk function  $r_K$  are Lipschitz continuous, upper bounded, and strongly convex with respect to the network output, and the eigenvalues of their Hessian matrix regarding the network output lie in  $[\rho, 1]$  for  $\rho \in (0, 1]$ . Note that the incorporated domain knowledge may not necessarily be perfect since it can be obtained based on subjective experiences (e.g., medical prognosis) (Muralidhar et al., 2018; Bica et al., 2020), pre-existing machine learning models (Hinton et al., 2014) or theoretical models which itself can deviate from the real physical world (Institute, 2020).

For training, in addition to the labeled dataset  $S_z$ , a dataset  $S_g$  with  $n_g$  unlabeled samples is generated for knowledge-based supervision. Note that  $S_g$  can also include inputs in  $S_z$ , and  $n_g$  can be sufficiently large since unlabeled samples are typically easier to obtain than labeled ones. The training risk of the informed neural network, which we simply refer

to as *informed risk*, is

$$\hat{R}_I(\mathbf{W}) = \frac{1-\lambda}{n_z} \sum_{S_z} r(h_{\mathbf{W},i}, z_i) + \frac{\lambda}{n_g} \sum_{S_g} r_K(h_{\mathbf{W},i}, g_i), \quad (3)$$

where  $\lambda \in [0, 1]$  is a hyper-weight,  $h_{\mathbf{W},i} = h_{\mathbf{W}}(x_i)$ , and  $g_i = g(x_i)$ . Note that Eqn. (3) can also be re-written as

$$\hat{R}_I(\mathbf{W}) = \sum_{S_z \cup S_g} [\mu_i r(h_{\mathbf{W},i}, z_i) + \lambda_i r_K(h_{\mathbf{W},i}, g_i)] \quad (4)$$

with hyper-parameters chosen as  $\mu_i = \frac{1-\lambda}{n_z} \mathbb{1}(x_i \in S_z)$  and  $\lambda_i = \frac{\lambda}{n_g} \mathbb{1}(x_i \in S_g)$ . Eqn. (4) is used for convergence analysis.

To train the informed DNN, we consider a gradient descent approach in Algorithm 1 shown in Appendix A. This training approach has also been commonly considered in the literature (Allen-Zhu et al., 2019b; Zou & Gu, 2019; Du et al., 2019) for theoretical analysis of standard DNNs without domain knowledge. For the sake of analysis, we also define a hypothesis space  $\mathcal{H} = \{h_{\mathbf{W}} \mid \mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)\}$  where  $\mathbf{W}^{(0)}$  is the initialized weight and  $\tau$  is the maximum distance between the weights in gradient descent and the initialized weights. We denote  $h_l^{(0)}(x), l \in [L]$  as the output of the  $l$ -th layer for an input  $x$  at initialization.

*Remark 1.* The considered informed learning is relevant to several other frameworks. For example, it can model weakly-supervised learning (Zhou, 2018; Wang et al., 2020) with a few (possibly imperfectly) labeled samples as well as other weak supervision signals (i.e., knowledge). Besides, by viewing  $\{z_i\}$  as hard labels and the knowledge-based model  $g(x)$  as soft labels provided by a teacher model, the informed learning captures knowledge distillation (Hinton et al., 2014; Phuong & Lampert, 2019; Rahbar et al., 2020). Thus, our work can complement the existing analysis for the aforementioned learning frameworks from a different and more unified perspective. Additionally, PAC-Bayesian learning optimizes the PAC-Bayesian bound which is a trade-off between the empirical error and a regularization term based on a prior distribution given by knowledge (Guedj, 2019; Amit & Meir, 2018; Germain et al., 2016). But, different from PAC-Bayesian learning which considers random hypothesis, we analyze an over-parameterized neural network with a predetermined architecture.

## 4. Effects of Domain Knowledge

### 4.1. Convergence

Since the domain knowledge is integrated into a neural network during training, it is important to analyze the convergence to understand how the label and knowledge supervision jointly determine the network output. While convergence based on gradient descent for over-parameterized

neural networks has been studied extensively (Bahri et al., 2021; Allen-Zhu et al., 2019b; Zou & Gu, 2019; Arora et al., 2019a; Du et al., 2019), the current analysis is *not* suitable to study the convergence of informed over-parameterized neural networks. The reasons are summarized as follows.

• **Inapplicable for multiple supervisions.** Typically, assuming one unique label for each distinct training sample and a large enough network width, the prior studies show that the neural network can fit to the labels, i.e., the network output for each training input converges to the corresponding label (Zhang et al., 2021; Arora et al., 2019a; Zou & Gu, 2019; Oymak & Soltanolkotabi, 2020). But, in our case, one training input can have multiple supervisions from both label and knowledge with possibly different forms of risks. Thus, the network output for an input may not be necessarily determined by a unique label. The convergence of knowledge distillation supervised by both hard and soft labels is studied by (Rahbar et al., 2020), but only the quadratic risk and shallow networks are considered.

• **Strong data separability assumption.** Some prior studies require a lower-bounded distance of any two samples (Allen-Zhu et al., 2019b; Zou & Gu, 2019; Du et al., 2019), but this may not be satisfied for an informed DNN because the input samples for label-based and knowledge-based risks can be very close or even the same. Other studies assume data separability by a neural tangent model (Chen et al., 2021b; Ji & Telgarsky, 2020; Cao & Gu, 2020; Nitanda et al., 2019), but data separability by a neural tangent model is not well defined for training with multiple supervisions in informed DNNs.

To address these challenges, we provide convergence analysis for informed over-parameterized neural networks based on a new data separability assumption of smooth sets. The construction of smooth sets approximates the space  $\mathcal{X}$  with discrete pieces, each containing samples that jointly satisfy the smooth properties. The smooth sets are formally defined below, followed by the data separability assumption.

**Definition 1** (Smooth sets). Given  $\phi > 0$ , construct a  $\phi$ -net (Clarkson, 2006)  $\mathcal{X}_\phi = \{x'_k, k \in [N], x'_k \in \mathcal{X}\}$  with  $N \sim O(1/\phi^b)$  such that  $\forall x'_i, x'_j \in \mathcal{X}_\phi$  and  $x'_i \neq x'_j$ ,  $\|x'_i - x'_j\| \geq \phi$  holds, and  $\forall x_i \in S_z \cup S_g$ , there exists at least one  $x'_k \in \mathcal{X}_\phi$  satisfying  $\|x_i - x'_k\| \leq \phi$ . Each input  $x'_k \in \mathcal{X}_\phi$ , referred to as a representative input, determines a smooth set  $\mathcal{C}_{\phi,k} = \{x \in \mathcal{X} \mid \|x - x'_k\| \leq \phi, \|x - x'_j\| \geq \phi/2, \forall j \neq k, x'_k, x'_j \in \mathcal{X}_\phi\}$ . The index set of training samples within the  $k$ th smooth set is  $\mathcal{I}_{\phi,k} = \{i \mid x_i \in S_z \cup S_g, x_i \in \mathcal{C}_{\phi,k}\}, k \in [N]$ .

**Assumption 1** (Data separability by smooth sets). For each smooth set  $k$  with representative sample  $x'_k$ , there exists a non-empty subset of neuron indices  $\mathcal{G}_{k,\alpha} \in [m]$  with size  $|\mathcal{G}_{k,\alpha}| = \alpha m, \alpha \in (0, 1]$  such that at initialization,  $\forall i \in \mathcal{I}_{\phi,k}, \forall j \in \mathcal{G}_{k,\alpha}, \mathbb{1} \left( \left[ h_L^{(0)}(x_i) \right]_j \geq 0 \right) =$

$\mathbb{1} \left( \left[ h_L^{(0)}(x'_k) \right]_j \geq 0 \right)$ , and  $\forall j \notin \mathcal{G}_{k,\alpha}$ , the pre-activation of the  $L$ -th layer  $\left| \left[ \mathbf{W}_L^{(0)} h_{L-1}^{(0)}(x_i) \right]_j \right| \geq \frac{3\sqrt{2\pi}\phi^{b+1}}{16\sqrt{m}}$ .

Instead of requiring a lower-bounded distance of any two training samples, the data separability assumption requires that, at initialization, for samples in one smooth set, the outputs of the last hidden layer either have the same signs as those of the representative sample, or their absolute values are larger than a very small threshold. Thus, this data separability assumption is set-wise and addresses the cases where two training inputs are very close or the same, and hence is milder than the one in existing studies (e.g., (Allen-Zhu et al., 2019b)). The parameter  $\alpha$  indicates slackness: with larger  $\alpha$ , more neurons have the same signs. Actually, data separation by smooth set with  $\phi > 0$  in Assumption 1 always exists: when  $\phi$  is small enough such that only one inputs or several same inputs are included in a smooth set, Assumption 1 is satisfied with  $\alpha = 1$ . Even in this worst case, our assumption is still milder than the data separability assumption considered in (Allen-Zhu et al., 2019b; Zou & Gu, 2019) that excludes the existence of two training samples with the same inputs but different supervisions.

With the data-separability assumption by smooth sets, we are ready to show the labels and knowledge jointly determine the network output for training inputs. We introduce the notation *effective label*, as formally defined below.

**Definition 2** (Effective label). For the  $k$ -th smooth set, define the effective label as  $y_{\text{eff},k} = \arg \min_h \sum_{i \in \mathcal{I}_{\phi,k}} \{\mu_i r(h, z_i) + \lambda_i r_K(h, g_i)\}$  with  $\mu_i, \lambda_i$  defined in Eqn. (3) and  $h$  in the space of network output, and the effective optimal risk as  $r_{\text{eff},k} = \sum_{i \in \mathcal{I}_{\phi,k}} \{\mu_i r(y_{\text{eff},k}, z_i) + \lambda_i r_K(y_{\text{eff},k}, g_i)\}$ .

Next, we show the convergence analysis. Note that the proof based on the data separability by smooth sets (Assumption 1) invalidates the proofs in previous studies, and we need new lemmas that lead to novel convergence to effective labels in Definition 2. In particular, in Lemma B.1, to approximate the outputs in the smooth set  $k$  by the output of the representative input  $x'_k$ , we need to bound the difference of the outputs with respect to  $x'_k$  and an input in the smooth set  $k$ . Also, based on Assumption 1, we derive in Lemma B.4 the gradient lower bound which relies on the number of smooth sets  $N$  instead of the sample size  $n_z + n_g$  in the previous analysis. This makes the network width  $m$  in our analysis directly rely on the smooth set size  $\phi$ . Moreover, in Lemma B.5, we prove based on the definition of smooth sets that the first-order approximation error of the total informed risk depends on the difference between the risk and effective risk in Definition 2. This is important to prove the convergence to the effective labels. The details of the convergence analysis are deferred to Appendix B.4.

**Theorem 4.1.** *Assume that the network width satisfies  $m \geq \Omega(\phi^{-11b-4} L^{15} d \rho^{-4} \bar{\lambda}^{-4} \alpha^{-4} \log^3(m))$ , and the step size is set as  $\eta = O(\frac{d}{L^2 m})$ . With Assumptions 1 satisfied, for any  $\epsilon > 0$  and  $\phi \leq \tilde{O}(\epsilon L^{-9/2} \log^{-3}(m))$ , we have with probability at least  $1 - O(\phi)$ , by gradient descent after  $T = O\left(\frac{L^2}{\phi^{1+2b} \rho \lambda \alpha} \log(\epsilon^{-1} \log(\phi^{-1}))\right)$  steps, the informed risk in Eqn. (4) is bounded as:  $\hat{R}_1(\mathbf{W}^{(T)}) - \hat{R}_{\text{eff}} \leq O(\epsilon)$ , where  $\hat{R}_{\text{eff}} = \sum_{k=1}^N r_{\text{eff},k}$ ,  $\bar{\lambda} = \Omega(\min(1 - \lambda, \lambda) \mathbf{1}(\lambda \in (0, 1)) + \mathbf{1}(\lambda \in \{0, 1\}))$ . Also, the DNN outputs satisfy:*

$$\sum_{S_z \cup S_g} (\mu_i + \lambda_i) \|h_{\mathbf{W}^{(T)}}(x_i) - y_{\text{eff},k(x_i)}\|^2 \leq O(\epsilon),$$

where  $k(x_i)$  is the index of the smooth set that includes  $x_i$ ,  $\mu_i = \frac{1-\lambda}{n_z} \mathbf{1}(x_i \in S_z)$  and  $\lambda_i = \frac{\lambda}{n_g} \mathbf{1}(x_i \in S_g)$ .

**Remark 2.** The convergence analysis in Theorem 4.1 addresses the limitations mentioned at the beginning of this section. First, instead of fitting a unique label for each input, the informed neural network with multiple supervisions converges to effective labels. Second, the data separability assumption is enough for convergence analysis of informed neural networks. Another observation is that with smaller  $\phi$  and smaller  $\alpha$ , Assumption 1 becomes milder, but a larger network width and more training steps are needed to guarantee convergence.

Additionally, different from previous convergence analysis where the width  $m$  increases directly with the sample size, the network width  $m$  in our analysis depends on the smooth set size  $\phi$  and is non-decreasing with sample size (i.e.,  $m$  may not always increase with the sample size). To see this, given a construction of smooth sets by size  $\phi$  that meets Assumption 1, if we continue to add (either labeled or knowledge-supervised) training samples that lie in the existing smooth sets and satisfy Assumption 1, the width  $m$  remains the same, and smaller  $\phi$  (larger  $m$ ) is needed to guarantee the convergence only when the added samples violate Assumption 1 under the current  $\phi$ . The large network width needed for analysis is due to the limitation of over-parameterization techniques, while in practice a much smaller network width is enough. Albeit beyond the scope of our study, addressing the gap between theory and practice is clearly important and still active research in the community (Bahri et al., 2021).

**Remark 3.** We can get more insights about the effects of labels and knowledge from the conclusion that the network outputs converge to the corresponding effective labels in Definition 2. On the one hand, if knowledge is applied to the samples within the same smooth sets as labeled samples, knowledge-based supervision and label-based supervision jointly determine the network output together: knowledge serves as a regularization for labels in this case. On the other hand, if a smooth set only contains knowledge-supervised

samples, the network output is determined solely by knowledge: knowledge supplements labeled samples (albeit possibly imperfectly) to provide additional supervision.

## 4.2. Generalization

We now formally analyze how the domain knowledge affects the generalization performance. From our convergence analysis, there are two different effects of knowledge (Remark 3). We characterize the two effects by formally defining knowledge imperfectness and knowledge-regularized label imperfectness. Before this, we list some notations for further analysis. Given a  $\phi$ -net  $\mathcal{X}_\phi$  (Definition 1),  $\mathcal{U}_\phi(S_z) = \{k \in [N] \mid \exists x \in S_z, x \in \mathcal{C}_{\phi,k}\}$  is the index collection of smooth sets that contain at least one labeled sample, and  $\mathcal{X}_\phi(S_z) = \bigcup_{k \in \mathcal{U}_\phi(S_z)} \mathcal{C}_{\phi,k}$  is the region covered by the smooth sets in  $\mathcal{U}_\phi(S_z)$ .  $S'_g = S_g \cap \mathcal{X}_\phi(S_z)$  is the knowledge supervised dataset with samples share the common smooth sets with labeled samples in  $S_z$  while the samples in  $S''_g = S_g \setminus S'_g$  lie in smooth sets without labeled samples. Denote  $n'_g = |S'_g|$  and  $n''_g = |S''_g|$ .

**Definition 3** (Knowledge imperfectness). Let  $h_K^* = \min_h \frac{1}{n''_g} \sum_{S''_g} [r_K(h(x_i), g(x_i))]$  be the optimal hypothesis for the knowledge-based risk on the dataset  $S''_g$ . The imperfectness of domain knowledge  $K$  applied to the dataset  $S''_g$  is defined as  $\hat{Q}_{K,S''_g} = \frac{1}{n''_g} \sum_{x_i \in S''_g} r(h_K^*(x_i), y_i)$  where  $y_i$  is the true label of  $x_i$ . Correspondingly, let  $\bar{h}_K^* = \min_h \mathbb{E}[r_K(h(x), g(x))]$  be the optimal hypothesis for the expected knowledge-based risk, and the expected imperfectness of domain knowledge  $K$  is defined as  $Q_K = \mathbb{E}[r(\bar{h}_K^*(x), y)]$ .

The (empirical or expected) knowledge imperfectness is defined as the risk under the hypothesis optimally learned by knowledge-based supervision. Thus, it measures the extent to which the domain knowledge is inconsistent with the true labels, measured in terms of the risk over the hypothesis set  $\mathcal{H}$ . Besides knowledge-based supervision, the network outputs for some smooth sets that contain both samples for knowledge risks and labeled samples are jointly determined by label-based and knowledge-based supervisions. Thus, we define knowledge-regularized label imperfectness below.

**Definition 4** (Knowledge-regularized label imperfectness). Let  $h_{R,\beta}^* = \arg \min_h \frac{1-\beta}{n_z} \sum_{S_z} r(h(x_i), z_i) + \frac{\beta}{n'_g} \sum_{S'_g} r_K(h(x_i), g(x_i))$  be the optimal hypothesis for the knowledge-regularized risk and  $\beta \in [0, 1]$ . The knowledge-regularized label imperfectness is  $\hat{Q}_{R,S_z,S'_g}(\beta) = \frac{1}{n_z} \sum_{S_z} r(h_{R,\beta}^*(x_i), y_i)$ , where  $y_i$  is the true label regarding  $x_i$ . Correspondingly, with  $\bar{h}_{R,\beta}^* = \arg \min_h \mathbb{E}[\frac{1-\beta}{n_z} \sum_{S_z} r(h(x_i), z_i) + \frac{\beta}{n'_g} \sum_{S'_g} r_K(h(x_i), g(x_i))]$  being the optimal hypothesis for the regularized risk, the expected knowledge regularized

label imperfectness is  $Q_R(\beta) = \mathbb{E} \left[ r(\bar{h}_{R,\beta}^*(x), y) \right]$ .

Like knowledge imperfectness, knowledge-regularized label imperfectness indicates the risk of the hypothesis optimally learned by joint supervision from labels and knowledge. We see that when  $\beta = 0$ ,  $\widehat{Q}_R(0)$  (or  $Q_R(0)$ ) is the imperfectness of pure label-based supervision. Thus, the gain due to knowledge is  $\Delta\widehat{Q}_{R,\beta} = \widehat{Q}_R(0) - \widehat{Q}_R(\beta)$  (or  $\Delta Q_{R,\beta} = Q_R(0) - Q_R(\beta)$  for the expected version). We show in the following theorem how the two types of imperfectness affect the population risk trained on the informed risk in Eqn. (3). The details are deferred to Appendix B.5.1.

**Theorem 4.2.** *With  $\mathbf{W}^{(T)}$  trained on Eqn. (3),  $\phi \leq \bar{O}(\epsilon^2 L^{-9/2} \log^{-3}(m))$ ,  $\phi \leq (\sqrt{\epsilon}/n_z)^{1/b}$ , and other assumptions the same as Theorem 4.1, with probability at least  $1 - O(\phi) - \delta$ ,  $\delta \in (0, 1)$ , the population risk satisfies*

$$R(h_{\mathbf{W}^{(T)}}) \leq O(\sqrt{\epsilon}) + (1 - \lambda)\widehat{Q}_{R,S_z,S'_g}(\beta_\lambda) + \lambda\widehat{Q}_{K,S''_g} + O\left(\Phi + \sqrt{\log(1/\delta)}\right) \left( \frac{1 - \lambda}{\sqrt{n_z}} + \frac{\lambda}{\sqrt{n_g}} \right),$$

where  $\beta_\lambda = \frac{\lambda n'_g}{(1 - \lambda)n_g + \lambda n'_g}$ ,  $\widehat{Q}_{R,S_z,S'_g}(\beta_\lambda)$  is the knowledge-regularized label imperfectness in Definition 4 and  $\widehat{Q}_{K,S''_g}$  is the knowledge imperfectness in Definition 3 applied to  $S''_g$ , and  $\Phi = O(4^L L^{3/2} m^{1/2} \phi^{-b-1/2} d \rho^{-1/2} \bar{\lambda}^{-1/2} \alpha^{-1/2})$ .

*Remark 4.* Theorem 4.2 shows that by training on the informed risk (3), knowledge affects the generation performance in the following two ways.

- **Knowledge for regularization.** When knowledge is applied to sample inputs inside the same smooth sets as labeled samples, it serves as an explicit regularization for label-based supervision, possibly reducing the label imperfectness from  $\widehat{Q}_{R,S_z,S'_g}(0)$  to  $\widehat{Q}_{R,S_z,S'_g}(\beta_\lambda)$ .

- **Knowledge for supplementing labels.** The generalization error is in the order of  $O\left(\frac{1 - \lambda}{\sqrt{n_z}} + \frac{\lambda}{\sqrt{n_g}}\right)$ . When no knowledge is used ( $\lambda = 0$ ), the order is as large as  $O\left(\frac{1}{\sqrt{n_z}}\right)$ . If knowledge is applied ( $\lambda > 0$ ), then the generalization error decreases with the increasing of knowledge-supervised sample size  $n_g$ . Thus, when knowledge is applied to smooth sets without labeled samples, it serves as an (possibly imperfect) supplement for labels, while introducing knowledge imperfectness  $\widehat{Q}_{K,S''_g}$ .

The hyper-parameter  $\lambda$  can be used to balance the introduced imperfectness and generalization error from label and knowledge supervision. However, by the risk bound, it is hard to use one hyper-parameter  $\lambda$  to control the two effects of knowledge, which will be further discussed in the next section.

## 5. A Generalized Training Objective

In the informed risk in Eqn. (3), only one hyper-weight  $\lambda$  is present, controlling the two different effects of knowledge (Remark 4). To better reap the benefits of knowledge, we consider a generalized informed risk in Eqn.(5) by introducing another hyper-weight  $\beta$ , which introduces more flexibility to govern the roles of domain knowledge.

$$\begin{aligned} \hat{R}_{I,G}(\mathbf{W}) = & \frac{(1 - \lambda)(1 - \beta)}{n_z} \sum_{S_z} r(h_{\mathbf{W},i}, z_i) + \\ & \frac{(1 - \lambda)\beta}{n'_g} \sum_{S'_g} r_K(h_{\mathbf{W},i}, g_i) + \frac{\lambda}{n''_g} \sum_{S''_g} r_K(h_{\mathbf{W},i}, g_i), \end{aligned} \quad (5)$$

where  $\beta, \lambda \in [0, 1]$ ,  $h_{\mathbf{W},i} = h_{\mathbf{W}}(x_i)$ ,  $g_i = g(x_i)$ .

In Eqn. (5), the two hyper-parameters  $\lambda$  and  $\beta$  can jointly control the knowledge effects (and the introduced imperfectness) when knowledge is applied. The hyperparameter  $\beta$  is used to control the knowledge regularization strength. By Remark 4, knowledge-supervised samples in  $S'_g$  serve as an explicit regularization for label-based supervision while introducing knowledge-regularized label imperfectness  $Q_R(\beta)$ . Thus, when  $\beta$  is larger, more effects from  $S'_g$  are incorporated and the regularization effect from knowledge is stronger. Also, we use  $\lambda$  to adjust the effect of supplementing labels and the introduction of  $Q_K$ . By Remark 4,  $S''_g$  serves as a supplement for labels while introducing the knowledge imperfectness  $Q_K$ . Thus, with larger  $\lambda$ , more effects from  $S''_g$  are incorporated, which means we incorporate more effects of data supplement from knowledge and also knowledge imperfectness  $Q_K$  but less effect of knowledge regularization and knowledge-regularized label imperfectness  $Q_R$ . The benefit of the training objective in Eqn. (5) will be explained formally in Theorem 5.1 and Corollary 5.2.

Compared with the objective in Eqn. (3) with only one hyper-parameter  $\lambda$ , Eqn. (5) introduces another hyper-parameter  $\beta$  to independently adjust the degree of the knowledge regularization, making Eqn. (5) more general and flexible. To train on Eqn. (5), we need to separate dataset for knowledge supervision into two datasets  $S'_g$  and  $S''_g$  based on whether an input is close to a labeled input and assign different hyper-weights to them. The knowledge-based dataset separation is determined by  $\phi$  in Definition 1. Specifically, when the network width goes to infinity ( $\phi$  goes to zero),  $S'_g$  shares the same inputs as  $S_z$ , but  $S'_g$  and  $S_z$  are supervised by knowledge and labels, respectively. We have  $S''_g = S_g \setminus S'_g = S_g \setminus S_z$  which supplements the labels as shown in Remark 4. Note that when the knowledge is perfect and knowledge-supervised samples are sufficient, we do not need labeled samples, i.e.,  $S_z = \emptyset$  and we set  $\lambda = 1, \beta = 1$ . Then, we have  $S''_g = S_g$  and Eqn. (5) becomes a purely knowledge-based risk. When no knowledge

is applied, we set  $\lambda = 0, \beta = 0$ , and Eqn. (5) becomes a purely label-based risk. In general cases when labels and knowledge are both used, hyper-parameters  $\lambda$  and  $\beta$  are used to control the effects of knowledge.

### 5.1. Population Risk

Note that Eqn. (5) can also be written as the form of Eqn. (4) with hyper-parameters chosen as  $\mu_i = \frac{(1-\lambda)(1-\beta)}{n_z} \mathbb{1}(x_i \in S_z)$  and  $\lambda_i = \frac{(1-\lambda)\beta}{n'_g} \mathbb{1}(x_i \in S'_g) + \frac{\lambda}{n''_g} \mathbb{1}(x_i \in S''_g)$ , so Theorem 4.1 for convergence still holds. Next, we bound the population risk based on the generalized informed risk. The details are given in Appendix B.5.2.

**Theorem 5.1.** *Assume that  $\mathbf{W}^{(T)}$  trained on Eqn. (5) and other assumptions are the same with those of Theorem 4.1, setting  $\phi : \phi \leq \tilde{O}(\epsilon^2 L^{-9/2} \log^{-3}(m))$  and  $\phi \leq (\sqrt{\epsilon}/n_z)^{1/b}$ , with probability at least  $1 - O(\phi) - \delta, \delta \in (0, 1)$ , the population risk satisfies*

$$R(h_{\mathbf{W}^{(T)}}) \leq O(\sqrt{\epsilon}) + (1-\lambda)\widehat{Q}_{R,S_z,S'_g}(\beta) + \lambda\widehat{Q}_{K,S''_g} + O\left(\Phi + \sqrt{\log(1/\delta)}\right) \left(\frac{1-\lambda}{\sqrt{n_z}} + \frac{\lambda}{\sqrt{n''_g}}\right),$$

where  $\beta$  and  $\lambda$  are trade-off hyper-parameters in Eqn. (5)

Additionally, to obtain more insights for sampling complexity, we further bound the population risk in terms of expected imperfectness, at the expense of some tightness. The proof details are deferred to Appendix B.5.3.

**Corollary 5.2.** *With the same assumptions as in Theorem 5.1, with probability at least  $1 - O(\phi) - \delta, \delta \in (0, 1)$ , the population risk satisfies*

$$R(h_{\mathbf{W}^{(T)}}) \leq O(\sqrt{\epsilon}) + (1-\lambda)Q_R(\beta) + \lambda Q_K + O\left(\Phi + \log^{1/4}(1/\delta)\right) \sqrt{\frac{1-\lambda}{\sqrt{n_z}} + \frac{\lambda}{\sqrt{n''_g}}},$$

where  $Q_R(\beta)$  is the expected knowledge-regularized label imperfectness in Definition 4,  $Q_K$  is the expected knowledge imperfectness in Definition 3.

**Remark 5.** Theorem 5.1 and Corollary 5.2 show that by training on the generalized informed risk in Eqn. (5), label and knowledge supervision jointly affect the population risk while introducing a combination of knowledge-regularized label imperfectness  $Q_R(\beta)$  and knowledge imperfectness  $Q_K$ . The effect of knowledge regularization is controlled by  $\beta$  and the trade-off between the two imperfectness terms and the trade-off between the two generalization errors  $\frac{1-\lambda}{\sqrt{n_z}}$  and  $\frac{\lambda}{\sqrt{n''_g}}$  are both controlled by  $\lambda$ . Thus, this gives us more flexibility to adjust how much domain knowledge is incorporated when it plays different roles in informed learning as discussed in Remark 4. Also, as shown by

the population risk bounds, we can tune the two hyper-parameters separately — we can first tune  $\beta$  to minimize  $Q_R(\beta)$ , and then tune  $\lambda$  to balance  $Q_R(\beta)$  and  $Q_K$ , and also balance the generalization errors due to sizes of datasets.

### 5.2. Sampling Complexity

We discuss the choices of hyper-parameters  $\beta$  and  $\lambda$  in different cases to guarantee a small population risk, and give the sampling complexity in each case, whose details are deferred to Appendix B.5.4.

**Corollary 5.3 (Sampling Complexity).** *With the same set of assumptions as in Corollary 5.2 and setting  $\beta^* = \arg \min_{\beta \in [0,1]} Q_R(\beta)$ , with probability at least  $1 - O(\phi) - \delta, \delta \in (0, 1)$ , to guarantee a population risk no larger than  $\sqrt{\epsilon}$ , we have the following cases:*

- If  $Q_K \leq \sqrt{\epsilon}$ , set  $\lambda = 1$ , the sampling complexity for labels is  $n_z = 0$  and the sampling complexity for knowledge-supervision is  $n_g \sim O(1/(\epsilon^2 - \epsilon^3))$ .
- If  $Q_K > \sqrt{\epsilon}$  and  $\frac{\sqrt{\epsilon}}{Q_K} + \frac{\sqrt{\epsilon}}{Q_R(\beta^*)} \geq 1$ , set  $\lambda = \frac{\sqrt{\epsilon}}{Q_K}$ , the sampling complexity for labels is  $n_z \sim O\left((1/\epsilon - 1/(\sqrt{\epsilon}Q_K))^2\right)$  and the sampling complexity for knowledge-supervision is  $n_g \sim O(1/((\epsilon - \epsilon^2)Q_K^2))$ .
- If  $\frac{\sqrt{\epsilon}}{Q_K} + \frac{\sqrt{\epsilon}}{Q_R(\beta^*)} < 1$ , a population risk as low as  $\sqrt{\epsilon}$  cannot be achieved no matter what  $\lambda$  is and how many samples are used.

**Remark 6.** In practice, unlabeled samples are typically cheaper to obtain than labeled samples. If  $Q_K \leq \sqrt{\epsilon}$ , the domain knowledge is good enough for supervision, and thus we can perform purely knowledge-based training without any labeled samples and guarantee a population risk no larger than  $\sqrt{\epsilon}$  with  $n''_g \sim O(1/\epsilon^2)$ , and hence  $n_g \sim O(1/(\epsilon^2 - \epsilon^3))$ . When the knowledge imperfectness  $Q_K > \sqrt{\epsilon}$ , we discuss the following two cases. First, if  $\frac{\sqrt{\epsilon}}{Q_K} + \frac{\sqrt{\epsilon}}{Q_R(\beta^*)} \geq 1$ , we can choose  $\lambda$  from  $\left[1 - \frac{\sqrt{\epsilon}}{Q_R(\beta^*)}, \frac{\sqrt{\epsilon}}{Q_K}\right]$  to control the risk from knowledge and label imperfectness as low as  $\sqrt{\epsilon}$ . We thus choose the largest  $\lambda = \frac{\sqrt{\epsilon}}{Q_K}$  to reduce the label sampling complexity. In this case, knowledge is not good enough, but label imperfectness is not too large. Thus, we can guarantee a population risk no larger than  $\sqrt{\epsilon}$  with labeled samples  $n_z \sim O\left((1/\epsilon - 1/(\sqrt{\epsilon}Q_K))^2\right)$  and knowledge supervised samples  $n_g \sim O(1/((\epsilon - \epsilon^2)Q_K^2))$ . Finally, if  $\frac{\sqrt{\epsilon}}{Q_K} + \frac{\sqrt{\epsilon}}{Q_R(\beta^*)} < 1$ , we cannot guarantee a population risk less than  $\sqrt{\epsilon}$  no matter what  $\lambda$  is and how many samples are used since the neither knowledge nor labels are of high enough quality.

In summary, the extreme cases are: Case (a) where the knowledge supervision alone is nearly perfect, and Case (c) where the knowledge and labels are both of low quality. Usually, we are in Case (b) where knowledge is imperfect but labels (after knowledge regularization) are good enough. In contrast, DNNs without using domain knowledge requires the label imperfectness  $Q_{R,0}$  not to exceed  $\sqrt{\epsilon}$ ; otherwise, the population risk cannot be guaranteed to be no greater than  $\sqrt{\epsilon}$ . The informed DNNs relaxes this requirement by requiring  $\frac{\sqrt{\epsilon}}{Q_K} + \frac{\sqrt{\epsilon}}{Q_{R(\beta^*)}} \geq 1$ . In addition, the incorporation of domain knowledge reduces the labeled sampling complexity from  $n_z \sim O(\frac{1}{\epsilon^2})$  in the traditional no-knowledge setting to  $n_z \sim O\left(\frac{1}{\epsilon} - \frac{1}{(\sqrt{\epsilon}Q_K)^2}\right)$ . In other words, the incorporation of knowledge is equivalent to  $O(\frac{2}{\epsilon^{3/2}Q_K} - \frac{1}{\epsilon Q_K^2})$  labeled samples, establishing a quantitative comparison between knowledge supervision and labeled samples.

## 6. Further Discussions

**Summary of analysis.** The convergence analysis in Theorem 4.1 introduces the concept of smooth sets and explains how the neural network output behaves by training on an informed risk. The generalization analysis in Theorem 4.2 explicitly shows the two different effects the domain knowledge has on the population risk (i.e., regularizing labels and supplementing labels). Based on this observation, we propose a generalized informed risk in Eqn. 5 to get more flexibility to control the two effects of knowledge, which is validated by Theorem 5.1 and its Corollary 5.2. Finally, the sampling complexity in Corollary 5.3 shows the effects of joint knowledge and label supervision in a quantitative way.

**Understanding knowledge distillation from the perspective of informed learning.** Knowledge distillation is extremely useful in practice (e.g., for model compression (Hinton et al., 2014)). Here, we show how our analysis complement the existing understanding of knowledge distillation (Hinton et al., 2014; Phuong & Lampert, 2019; Rahbar et al., 2020; Dao et al., 2021; Ji & Zhu, 2020) from the perspective of hard label and teacher’s knowledge imperfectness. In our formulation, hard labels are  $\{z_i\}$  in the labeled dataset, whose imperfectness (non-softness) is measured by  $Q_R(0)$ . In Theorems 4.2, 5.1, and Corollary 5.2, by viewing the teacher model  $g(x)$  as domain knowledge, we show the teacher benefits the student training by providing a regularization gain  $\Delta Q_{R,\beta}$ , and reducing the sampling complexity of hard labels by Corollary 5.3. The knowledge-regularized label imperfectness  $Q_{R,\beta}$  can be less than pure label imperfectness  $Q_R(0)$  because the soft label can smooth the network output within each smooth set. But, given the teacher (knowledge) imperfectness  $Q_K$ , there exists a trade-off between hard label and teacher supervision.

Importantly, our results are in line with the observations and also complement the analysis in (Ji & Zhu, 2020). Specifically, (Ji & Zhu, 2020) uses NTK to show that the soft labels provided by a teacher model (knowledge) are easier to learn than hard labels while hard labels can correct imperfect teachers pointwise, exhibiting a trade-off between hard labels and the imperfect teacher. We define the hard label and teacher (knowledge) imperfectness, and show that for a neural network with finite width, hard labels and teacher’s knowledge compensate for each other within each smooth set. In consistency with our results, (Rahbar et al., 2020) based on NTK also presents a trade-off between labels and the imperfect teacher. The teacher model imperfectness is also observed by (Dao et al., 2021) which measures the teacher imperfectness by the squared norm of the difference of the soft label and the true Bayesian class probability. Note, however, that our analysis *cannot* adequately explain the benefit of knowledge distillation for the perspective of feature learning due to the inherent limitations of over-parameterization techniques, which are further discussed in (Allen-Zhu & Li, 2020).

## 7. Numerical Results

### 7.1. Problem Setup

We consider an informed DNN with domain knowledge in the form of constraints to learn a Bohachevsky function. The learning task is to learn a relationship  $y(x)$ . The learner is provided with a dataset with labeled samples  $S_z = \{(x_i, z_i), i \in [n_z]\}$ , having possibly noisy labels  $z_i = y(x_i) + n_i, n_i \sim \mathcal{N}(0, \sigma_z^2)$ , and an unlabeled dataset  $S_g = \{(x_i), i \in [n_g]\}$ . Additionally, the learner is informed with the constraint knowledge, which includes an upper bound  $g_{ub}(x)$  and a lower bound  $g_{lb}(x)$  on the true label corresponding to input  $x$ , i.e.  $g_{lb}(x) \leq y(x) \leq g_{ub}(x)$ . A neural network  $h_W(x)$  is used for learning and the metric of interest is the mean square error (MSE) of the network output  $h_W(x)$  with respect to the true label  $y(x)$  on a test dataset  $S_t$ , which is expressed as  $\hat{R}_{S_t}(h_W) = \frac{1}{2|S_t|} \sum_{(x_i, y_i) \in S_t} (h_W(x_i) - y_i)^2$ . Assume that the relationship to be learned is governed by a multi-dimensional Bohachevsky function  $y(x) = x\mathbf{A}\mathbf{A}^\top x^\top - c \cos(a^\top x) + c$ , where  $\mathbf{A}$  is a  $b \times b$  matrix,  $a$  is a  $b$ -dimensional vector and  $c$  is a constant. The constraint knowledge includes an upper bound model  $g_{ub}(x) = x\mathbf{A}\mathbf{A}^\top x^\top + ub$  with  $ub \geq 2c$ , and a lower bound model  $g_{lb}(x) = x\mathbf{A}\mathbf{A}^\top x^\top + lb$  with  $lb \leq 0$ . While it is not strongly convex and hence deviates from the assumptions in our theoretical analysis, we use ReLU as the knowledge-based risk function, i.e., the knowledge-based risk is written as  $r_K(h_W(x)) = \text{relu}(h_W(x) - g_{ub}(x)) + \text{relu}(g_{lb}(x) - h_W(x))$ . If  $ub - lb$  is larger, the uncertainty of the label given the knowledge is larger — the knowledge imperfectness is higher. We choose  $(lb, ub)$  as  $(0, 0.6)$  and



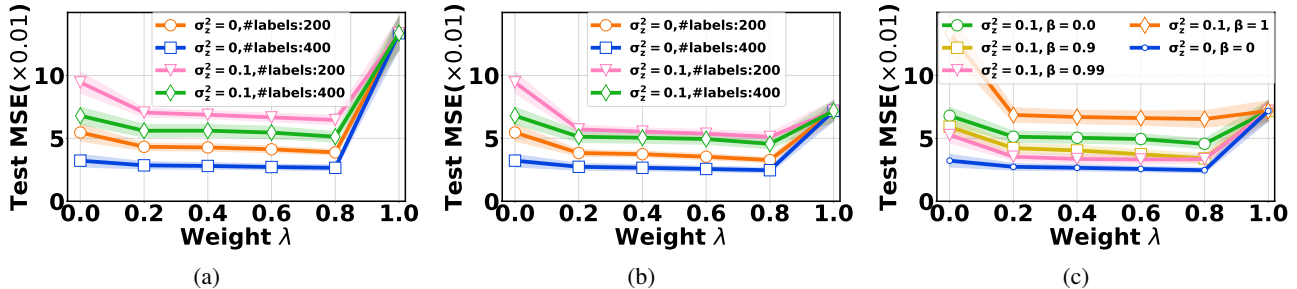


Figure 1. Test MSE under different hyper-parameters.  $\sigma_z^2 = 0$  means using perfect labels;  $\sigma_z^2 = 0.1$  means using imperfect labels with noise variance 0.1; knowledge imperfectness is determined by  $ub$  and  $lb$  in the problem setting. (a) Training on the standard informed objective Eqn. (3) using knowledge with high imperfectness; (b) Training on the standard informed objective Eqn. (3) using knowledge with low imperfectness; (c) Training on the generalized informed objective Eqn. (5) using knowledge with low imperfectness and 400 labels.

(0, 0.8) respectively to show the performances under low and high knowledge imperfectness. More details of the setup are in Appendix G.1.

## 7.2. Results

The curves of test MSE with different knowledge and label settings are shown in Fig. 1. In all the three figures, test MSE in  $\lambda = 0$  approximately measures knowledge-regularized label imperfectness in Definition 4, while test MSE in  $\lambda = 1$  approximately measures knowledge imperfectness in Definition 3. We first use the training objective Eqn.(3) in Fig. 1(a) and Fig. 1(b) to show the effect of adjusting  $\lambda$ , which controls the knowledge effects (see Remark 4). From both Fig. 1(a) and Fig. 1(b), we see that the test MSE is smaller when there are more labeled samples and when label noise variance is lower. Importantly, domain knowledge helps reduce the MSE compared with pure label-supervised learning, especially for the cases with fewer labels and high label noise variance. Also, by comparing Fig. 1(a) and Fig. 1(b), we can find that the test MSE is lower when the knowledge imperfectness is lower. Additionally, Fig. 1(c) gives the test MSEs training on the generalized objective (5) under different  $\beta$  when the labeled dataset size is 400, showing that the test risk can be reduced by adjusting  $\beta$  which controls the knowledge regularization effect (see Remark 5). We can find that by properly adjusting  $\beta$ , the test MSEs under label noise are very close to that without label noise (the blue line). When  $\beta = 1$ , the test MSE is the highest since no labeled data is used to provide supervision.

More results, including another application of learning to manage wireless spectrum, are available in Appendix G.2.

## 8. Conclusion

In this paper, we consider an informed DNN with domain knowledge integrated with its training risk function. We quantitatively demonstrate that domain knowledge can improve the generalization performance and reduce the sam-

pling complexity, while also impacting the point to which the network output converges. Our analysis also reveals that knowledge affects the generalization performance in two ways: regularizing the label supervision, and supplementing the labeled samples. Finally, we discuss how an informed DNN relates to other learning frameworks.

## Acknowledgment

This work was supported in part by the U.S. NSF under CNS-1910208.

## References

- Allen-Zhu, Z. and Li, Y. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- Allen-Zhu, Z., Li, Y., and Liang, Y. Learning and generalization in overparameterized neural networks, going beyond two layers. *NeurIPS*, 2019a.
- Allen-Zhu, Z., Li, Y., and Song, Z. A convergence theory for deep learning via over-parameterization. In *ICML*, 2019b.
- Amit, R. and Meir, R. Meta-learning by adjusting priors based on extended pac-bayes theory. In *ICML*, pp. 205–214, 2018.
- Arora, S., Du, S., Hu, W., Li, Z., and Wang, R. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *ICML*, 2019a.
- Arora, S., Du, S. S., Hu, W., Li, Z., Salakhutdinov, R., and Wang, R. On exact computation with an infinitely wide neural net. *NeurIPS*, 2019b.
- Bahri, Y., Gu, Q., Karbasi, A., and Sedghi, H. Over-parameterization: Pitfalls and opportunities. In *ICML Workshop*, 2021. URL

<https://icml.cc/Conferences/2021/ScheduleMultitrack?event=8357>.

- Baker, N., Alexander, F., Bremer, T., Hagberg, A., Kevrekidis, Y., Najm, H., Parashar, M., Patra, A., Sethian, J., Wild, S., et al. Workshop report on basic research needs for scientific machine learning: Core technologies for artificial intelligence. Technical report, USDOE Office of Science (SC), Washington, DC (United States), 2019.
- Bamler, R., Salehi, F., and Mandt, S. Augmenting and tuning knowledge graph embeddings. In *UAI*, pp. 508–518. PMLR, 2020.
- Bardenet, R., Brendel, M., Kégl, B., and Sebag, M. Collaborative hyperparameter tuning. In *ICML*, 2013.
- Bartlett, P., Foster, D. J., and Telgarsky, M. Spectrally-normalized margin bounds for neural networks. *NeurIPS*, 2017.
- Battaglia, P., Pascanu, R., Lai, M., Rezende, D. J., et al. Interaction networks for learning about objects, relations and physics. In *NeurIPS*, pp. 4502–4510, 2016.
- Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- Beck, C., Weinan, E., and Jentzen, A. Machine learning approximation algorithms for high-dimensional fully nonlinear partial differential equations and second-order backward stochastic differential equations. *Journal of Nonlinear Science*, 29(4):1563–1619, 2019.
- Benaim, S. and Wolf, L. One-shot unsupervised cross domain translation. In *NeurIPS*, pp. 2104–2114, 2018.
- Bica, I., Alaa, A. M., Jordon, J., and van der Schaar, M. Estimating counterfactual treatment outcomes over time through adversarially balanced representations. In *ICLR*, 2020. URL <https://openreview.net/forum?id=BJg866NFvB>.
- Blanc, G., Gupta, N., Valiant, G., and Valiant, P. Implicit regularization for deep neural networks driven by an ornstein-uhlenbeck like process. In *Conference on learning theory*, pp. 483–513. PMLR, 2020.
- Borghesi, A., Baldo, F., and Milano, M. Improving deep learning models via constraint-based domain knowledge: a brief survey. *arXiv preprint arXiv:2005.10691*, 2020.
- Cannings, T. I., Fan, Y., and Samworth, R. J. Classification with imperfect training labels. *Biometrika*, 107(2):311–330, 2020.
- Cao, Y. and Gu, Q. Generalization bounds of stochastic gradient descent for wide and deep neural networks. *NeurIPS*, 2019.
- Cao, Y. and Gu, Q. Generalization error bounds of gradient descent for learning over-parameterized deep relu networks. In *AAAI*, 2020.
- Chen, Y., Gao, R., Liu, F., and Zhao, D. Modulenet: Knowledge-inherited neural architecture search. *IEEE Transactions on Cybernetics*, 2021a.
- Chen, Z., Cao, Y., Zou, D., and Gu, Q. How much over-parameterization is sufficient to learn deep relu networks? *ICLR*, 2021b.
- Chiang, M., Hande, P., and Lan, T. *Power control in wireless cellular networks*. Now Publishers Inc, 2008.
- Cho, J. H. and Hariharan, B. On the efficacy of knowledge distillation. In *ICCV*, pp. 4794–4802, 2019.
- Clarkson, K. L. Building triangulations using  $\epsilon$ -nets. In *STOC*, 2006.
- Cui, W., Shen, K., and Yu, W. Spatial deep learning for wireless scheduling. *IEEE Journal on Selected Areas in Communications*, 37(6):1248–1261, 2019.
- Dao, T., Kamath, G. M., Syrgkanis, V., and Mackey, L. Knowledge distillation as semiparametric inference. *ICLR*, 2021.
- Deng, C., Ji, X., Rainey, C., Zhang, J., and Lu, W. Integrating machine learning with human knowledge. *iScience*, 23(11):101656, 2020.
- Diligenti, M., Gori, M., and Sacca, C. Semantic-based regularization for learning and inference. *Artificial Intelligence*, 244:143–165, 2017a.
- Diligenti, M., Roychowdhury, S., and Gori, M. Integrating prior knowledge into deep learning. In *ICMLA*, pp. 920–923, 2017b.
- Du, S., Lee, J., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. In *ICML*, pp. 1675–1685, 2019.
- Fang, Y., Kuan, K., Lin, J., Tan, C., and Chandrasekhar, V. Object detection meets knowledge graphs.(2017). In *IJCAI*, 2017.
- Furlanello, T., Lipton, Z., Tschannen, M., Itti, L., and Anandkumar, A. Born again neural networks. In *ICML*, pp. 1607–1616, 2018.
- Gao, H., Shou, Z., Zareian, A., Zhang, H., and Chang, S.-F. Low-shot learning via covariance-preserving adversarial augmentation networks. In *NeurIPS*, pp. 975–985, 2018a.

- Gao, T., Liu, H., Liu, J., Rajan, H., and Gao, H. A global convergence theory for deep relu implicit networks via over-parameterization. In *ICML*, 2021.
- Gao, Y., Xu, H., Lin, J., Yu, F., Levine, S., and Darrell, T. Reinforcement learning from imperfect demonstrations. *ICLR Workshop*, 2018b.
- Garcia, V. and Bruna, J. Few-shot learning with graph neural networks. *ICLR*, 2018.
- Germain, P., Bach, F., Lacoste, A., and Lacoste-Julien, S. Pac-bayesian theory meets bayesian inference. *NeurIPS*, 2016.
- Goldsmith, A. *Wireless Communications*. Cambridge University Press, 2005.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- Gore, A. D. and Karandikar, A. Link scheduling algorithms for wireless mesh networks. *IEEE Communications Surveys & Tutorials*, 13(2):258–273, 2010.
- Gou, J., Yu, B., Maybank, S. J., and Tao, D. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.
- Guedj, B. A primer on pac-bayesian learning. *Proceedings of the 2nd congress of the Société Mathématique de France*, pp. 391–414, 2019.
- Hester, T., Vecerik, M., Pietquin, O., Lanctot, M., Schaul, T., Piot, B., Horgan, D., Quan, J., Sendonaris, A., Dulac-Arnold, G., et al. Deep q-learning from demonstrations. *AAAI*, 2018.
- Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. *Neurips Deep Learning Workshop*, 2014.
- Hong, M. and Luo, Z.-Q. Signal processing and optimal resource allocation for the interference channel. In *Academic Press Library in Signal Processing*, volume 2, pp. 409–469. 2014.
- Humbird, K. D., Peterson, J. L., and McClarren, R. G. Deep neural network initialization with decision trees. *IEEE transactions on neural networks and learning systems*, 30(5):1286–1295, 2018.
- Husken, M. and Goerick, C. Fast learning for problem classes using knowledge based network initialization. In *IEEE IJCNN*, 2000.
- Institute, T. A. T. Physics-informed machine learning, 2020. <https://www.turing.ac.uk/research/theory-and-method-challenge-fortnights/physics-informed-machine-learning>.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. *NeurIPS*, 2018.
- Ji, G. and Zhu, Z. Knowledge distillation in wide neural networks: Risk bound, data efficiency and imperfect teacher. *NeurIPS*, 2020.
- Ji, Z. and Telgarsky, M. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. *ICLR*, 2020.
- Karpatne, A., Watkins, W., Read, J., and Kumar, V. Physics-guided neural networks (pgnn): An application in lake temperature modeling. *arXiv preprint arXiv:1710.11431*, 2017.
- Khanduri, P., Yang, H., Hong, M., Liu, J., Wai, H. T., and Liu, S. Decentralized learning for overparameterized problems: A multi-agent kernel approximation approach. In *International Conference on Learning Representations*, 2021.
- Khoo, Y., Lu, J., and Ying, L. Solving parametric pde problems with artificial neural networks. *European Journal of Applied Mathematics*, 32(3):421–435, 2021.
- Klautau, A., Batista, P., González-Prelcic, N., Wang, Y., and Heath, R. W. 5g mimo data for machine learning: Application to beam-selection using deep learning. In *ITA*, pp. 1–9, 2018.
- Kurata, G., Xiang, B., and Zhou, B. Improved neural network-based multi-label classification with better initialization leveraging label co-occurrence. In *NAACL: Human Language Technologies*, pp. 521–526, 2016.
- Lee, J., Xiao, L., Schoenholz, S. S., Bahri, Y., Novak, R., Sohl-Dickstein, J., and Pennington, J. Wide neural networks of any depth evolve as linear models under gradient descent. *NeurIPS*, 2019.
- Liang, F., Shen, C., Yu, W., and Wu, F. Towards optimal power control via ensembling deep neural networks. *IEEE Transactions on Communications*, 68(3):1760–1776, 2019.
- Lu, L., Meng, X., Mao, Z., and Karniadakis, G. E. Deepxde: A deep learning library for solving differential equations. *SIAM Review*, 63(1):208–228, 2021.
- Maher, M. and Sakr, S. Smartml: A meta learning-based framework for automated selection and hyperparameter tuning for machine learning algorithms. In *EDBT*, 2019.

- Marino, K., Salakhutdinov, R., and Gupta, A. The more you know: Using knowledge graphs for image classification. *CVPR*, 2016.
- Motiiian, S., Piccirilli, M., Adjeroh, D. A., and Doretto, G. Information bottleneck learning using privileged information for visual recognition. In *CVPR*, June 2016.
- Muralidhar, N., Islam, M. R., Marwah, M., Karpatne, A., and Ramakrishnan, N. Incorporating prior domain knowledge into deep neural networks. In *IEEE Big Data*, pp. 36–45, 2018.
- Neyshabur, B., Li, Z., Bhojanapalli, S., LeCun, Y., and Srebro, N. The role of over-parametrization in generalization of neural networks. In *ICLR*, 2018.
- Nitanda, A., Chinot, G., and Suzuki, T. Gradient descent can learn less over-parameterized two-layer neural networks on classification problems. *arXiv preprint arXiv:1905.09870*, 2019.
- Oymak, S. and Soltanolkotabi, M. Toward moderate overparameterization: Global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 1(1):84–105, 2020.
- Pfrommer, J., Zimmerling, C., Liu, J., Kärger, L., Henning, F., and Beyerer, J. Optimisation of manufacturing process parameters using deep neural networks as surrogate models. *Procedia CiRP*, 72:426–431, 2018.
- Phuong, M. and Lampert, C. Towards understanding knowledge distillation. In *ICML*, pp. 5142–5151, 2019.
- Rahbar, A., Panahi, A., Bhattacharyya, C., Dubhashi, D., and Chehreghani, M. H. On the unreasonable effectiveness of knowledge distillation: Analysis in the kernel regime. *arXiv preprint arXiv:2003.13438*, 2020.
- Raissi, M., Perdikaris, P., and Karniadakis, G. E. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10561*, 2017.
- Raissi, M., Perdikaris, P., and Karniadakis, G. E. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- Ramsey, C. L. and Grefenstette, J. J. Case-based initialization of genetic algorithms. In *ICGA*, pp. 84–91, 1993.
- Robinson, J., Jegelka, S., and Sra, S. Strength from weakness: Fast learning using weak supervision. In *ICML*, pp. 8127–8136, 2020.
- Sanayei, S. and Nosratinia, A. Antenna selection in mimo systems. *IEEE Communications magazine*, 42(10):68–73, 2004.
- Sharmanska, V., Quadrianto, N., and Lampert, C. H. Learning to rank using privileged information. In *ICCV*, 2013.
- Shin, Y., Darbon, J., and Karniadakis, G. E. On the convergence and generalization of physics informed neural networks. *arXiv preprint arXiv:2004.01806*, 2020.
- Silvestri, M., Lombardi, M., and Milano, M. Injecting domain knowledge in neural networks: a controlled experiment on a constrained problem. pp. 266–282, 2021.
- Song, C., Ramezani-Kebrya, A., Pethick, T., Eftekhari, A., and Cevher, V. Subquadratic overparameterization for shallow neural networks. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *NeurIPS*, 2021.
- Sun, H., Chen, X., Shi, Q., Hong, M., Fu, X., and Sidiropoulos, N. D. Learning to optimize: Training deep neural networks for interference management. *IEEE Transactions on Signal Processing*, 66(20):5438–5453, 2018.
- Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P. H., and Hospedales, T. M. Learning to compare: Relation network for few-shot learning. In *CVPR*, pp. 1199–1208, 2018.
- Towell, G. G. and Shavlik, J. W. Knowledge-based artificial neural networks. *Artificial intelligence*, 70(1-2):119–165, 1994.
- Van Rijn, J. N. and Hutter, F. Hyperparameter importance across datasets. In *ACM SIGKDD*, pp. 2367–2376, 2018.
- Vapnik, V. and Vashist, A. A new learning paradigm: Learning using privileged information. *Neural networks*, 22(5-6):544–557, 2009.
- von Rueden, L., Mayer, S., Beckh, K., Georgiev, B., Giesselbach, S., Heese, R., Kirsch, B., Walczak, M., Pfrommer, J., Pick, A., et al. Informed machine learning—a taxonomy and survey of integrating prior knowledge into learning systems. *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- Wang, Y., Yao, Q., Kwok, J. T., and Ni, L. M. Generalizing from a few examples: A survey on few-shot learning. *ACM Computing Surveys*, 53(3):1–34, 2020.
- Wei, C., Lee, J. D., Liu, Q., and Ma, T. Regularization matters: Generalization and optimization of neural nets vs their induced kernel. In *NeurIPS*, pp. 9712–9724, 2019.
- Wei, C., Kakade, S., and Ma, T. The implicit and explicit regularization effects of dropout. In *ICML*, 2020.

- Willard, J., Jia, X., Xu, S., Steinbach, M. S., and Kumar, V. Integrating physics-based modeling with machine learning: A survey. *arXiv preprint arXiv:2003.04919*, 2020.
- Xu, J., Zhang, Z., Friedman, T., Liang, Y., and Broeck, G. A semantic loss function for deep learning with symbolic knowledge. In *ICML*, pp. 5502–5511, 2018.
- Yang, G. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*, 2019.
- Yin, Y., Le Guen, V., Dona, J., de Bézenac, E., Ayed, I., Thome, N., and Gallinari, P. Augmenting physical models with deep networks for complex dynamics forecasting. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124012, 2021.
- Zappone, A., Di Renzo, M., and Debbah, M. Wireless networks design in the era of deep learning: Model-based, ai-based, or both? *IEEE Transactions on Communications*, 67(10):7331–7376, 2019. doi: 10.1109/TCOMM.2019.2924010.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Zhang, Y., Tang, H., and Jia, K. Fine-grained visual categorization using meta-learning optimization with sample selection of auxiliary data. In *ECCV*, pp. 233–248, 2018.
- Zhou, Z.-H. A brief introduction to weakly supervised learning. *National science review*, 5(1):44–53, 2018.
- Zou, D. and Gu, Q. An improved analysis of training over-parameterized deep neural networks. In *NeurIPS*, 2019.
- Zou, D., Cao, Y., Zhou, D., and Gu, Q. Gradient descent optimizes over-parameterized deep relu networks. *Machine Learning*, 109(3):467–492, 2020.

## Appendix

### A. Training Algorithm

To train the knowledge-informed DNN, we consider a gradient descent approach in Algorithm 1. This training approach has also been commonly considered in the literature (Allen-Zhu et al., 2019b; Zou & Gu, 2019; Du et al., 2019) for theoretical analysis of standard DNNs without domain knowledge.

---

**Algorithm 1** Informed Neural Network Training by Gradient Descent
 

---

**Initialization:** Initialize each entry of weights  $\mathbf{W}_0^{(0)}, \mathbf{W}_l^{(0)}, l \in [L]$  independently by  $\mathcal{N}(0, \frac{2}{m})$  and each entry of  $\mathbf{V}^{(0)}$  independently by  $\mathcal{N}(0, \frac{1}{d})$ .  
**for**  $t = 0, \dots, T - 1$  **do**  
     Update the weights as  $\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} + \eta \nabla_{\mathbf{W}} \hat{R}_{\text{I}}(\mathbf{W}^{(t)})$ .  
**end for**  
**Output:**  $\mathbf{W}^{(T)}$ .

---

### B. Notations, Key Lemmas and Proofs of Main Results in Section 4 and Section 5

#### B.1. Further Notations

Before the proofs, we list some additional notations as below. Denote  $n' = |S_z| + |S_g| = n_z + n_g$ . We assign the samples in the dataset  $S_z$  with indices from 1 to  $n_z$  and the samples in the dataset  $S_g$  with indices from  $n_z + 1$  to  $n'$ . The informed risk of an informed DNN in Eqns. (3),(5) can be re-written as

$$\hat{R}_{\text{I}}(\mathbf{W}) = \sum_{i=1}^{n'} [\mu_i r(h_{\mathbf{W}}(x_i), z_i) + \lambda_i r_{\text{K}}(h_{\mathbf{W}}(x_i), g(x_i))], \quad (6)$$

where  $\sum_{i=1}^{n'} (\mu_i + \lambda_i) = 1$ . Thus, in Eqn. (3), we have  $\mu_i = \frac{1-\lambda}{n_z} \mathbb{1}(x_i \in S_z)$  and  $\lambda_i = \frac{\lambda}{n_g} \mathbb{1}(x_i \in S_g)$ ; in Eqn. (5), we have  $\mu_i = \frac{(1-\lambda)(1-\beta)}{n_z} \mathbb{1}(x_i \in S_z)$  and  $\lambda_i = \frac{(1-\lambda)\beta}{n'_g} \mathbb{1}(x_i \in S'_g) + \frac{\lambda}{n'_g} \mathbb{1}(x_i \in S''_g)$ . We prove convergence for the above three risks.

For any input  $x_i, i \in [n']$ , we denote the DNN output with respect to weight  $\mathbf{W}$  as  $h_{\mathbf{W},i} = h_{\mathbf{W}}(x_i)$ . To express the output of the ReLU activation of the  $l$ -th layer for an input sample  $x_i$ , for  $l \in [L]$  and  $i \in [n]$ , we denote a diagonal matrix  $\mathbf{D}_{l,i}$  with its  $j$ -th (for  $j \in [m]$ ) diagonal entry as  $\mathbb{1}([\mathbf{W}_l h_{l-1}]_j \geq 0)$ . Thus, given the input  $x_i$ , the DNN output can be expressed as

$$h_{\mathbf{W},i} = \mathbf{V} \mathbf{D}_{L,i} \mathbf{W}_L \mathbf{D}_{L-2,i} \cdots \mathbf{D}_{0,i} \mathbf{W}_0 x_i. \quad (7)$$

Also, we denote the informed risk for hypothesis  $h \in \mathcal{H}$  and input  $x_i$  as

$$r_{\text{I},i} = \mu_i r(h(x_i), z_i) + \lambda_i r_{\text{K}}(h(x_i), g(x_i)) \quad (8)$$

The gradient of informed risk with respect to the hypothesis output is

$$u_i(h(x_i)) = \nabla_h r_{\text{I},i}(h(x_i)) = \mu_i \nabla_h r(h(x_i), z_i) + \lambda_i \nabla_h r_{\text{K}}(h(x_i), g(x_i)). \quad (9)$$

After constructing the smooth sets, denote for the  $k$ th smooth set, the sum of indices as  $M_k = \sum_{\mathcal{I}_{\phi,k}} (\mu_i + \lambda_i)$ . Denote the sum risk of the  $k$ th smooth set for hypothesis  $h \in \mathcal{H}$  as

$$\bar{r}_{\text{I},k}(h(x_i)) = \sum_{i \in \mathcal{I}_{\phi,k}} r_{\text{I},i}(h(x_i)). \quad (10)$$

Thus, the effective label given in Definition 2 is written as  $y_{\text{eff},k} = \arg \min_h \bar{r}_{\text{I},k}(h)$  with  $h$  in the space of network output, and the optimal effective risk is written as  $r_{\text{eff},k} = \bar{r}_{\text{I},k}(y_{\text{eff},k})$ .

We then give some key technical lemmas which are the foundations for our further analysis. The proofs for these lemmas are shown in Appendix C.

## B.2. Forward Perturbation Regarding Inputs

The forward perturbation for weights in the weight update range is proved in (Allen-Zhu et al., 2019b). However, to characterize the smooth sets, it is important to prove forward perturbation for inputs in a smooth set, which is given as follows.

**Lemma B.1.** *For any  $i \in \mathcal{I}_{\phi,k}$ ,  $k \in [N]$ , let  $h_{l,k} = h_l(x'_k)$ ,  $h_{l,i} = h_l(x_i)$ , and  $f_{l,k} = \mathbf{W}_l h_{l-1}(x'_k)$ ,  $f_{l,i} = \mathbf{W}_l h_{l-1}(x_i)$  and denote  $\mathbf{D}'_{l,i,k} \in \mathbb{R}^{m \times m}$  as the diagonal matrix with  $[\mathbf{D}'_{l,i,k}]_{j,j} = \mathbf{1}([f_{l,i}]_j \geq 0) - \mathbf{1}([f_{l,k}]_j \geq 0)$ . Assuming  $\phi \leq O(L^{-9/2} \log^{-3}(m) \log^{-3/4}(1/\phi))$ , we have with probability at least  $1 - \phi$  over the randomness of  $\mathbf{W}^{(0)}$ ,*

(a) *At initialization,  $\|\mathbf{D}'_{l,i,k}\|_0 \leq O(m\phi^{2/3} L \log^{1/2}(1/\phi))$*

(b) *For  $\mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$  with  $\tau \leq O(\phi^{3/2})$  we have  $\|h_{l,i} - h_{l,k}\| \leq O(L^{5/2} \phi \sqrt{\log(m) \log(1/\phi)})$  and  $\|f_{l,i} - f_{l,k}\| \leq O(L^{5/2} \phi \sqrt{\log(m) \log(1/\phi)})$ .*

The proof of Lemma B.1 is given in Section C.1. The forward perturbation regarding inputs indicates the smoothness property of neural networks with respect to inputs. For compactness, we absorb the logarithmically increasing terms into  $\tilde{O}$  and denote  $\tilde{O}(L^{5/2} \phi \log^{1/2}(m)) = O(L^{5/2} \phi \sqrt{\log(m) \log(1/\phi)})$ ,  $\tilde{O}(L^{-9/2} \log^{-3}(m)) = O(L^{-9/2} \log^{-3}(m) \log^{-3/4}(1/\phi))$  in the following analysis.

## B.3. Properties of Strong Convexity

Since our analysis is based on strongly convex risk functions, we give some key properties of strongly convex functions.

**Lemma B.2** (Properties of Strong Convexity). *If a strongly convex function  $r(h)$  has a minimum value of  $r(h^*) = r_{\min}$  and the eigenvalues of its Hessian matrix lie in  $[\rho, 1]$ , then we have  $\|\nabla r(h)\|^2 \leq 2(r(h) - r_{\min})$ ,  $\|\nabla r(h)\|^2 \geq 2\rho(r(h) - r_{\min})$  and  $\|h^* - h\| \leq \frac{2}{\rho} \|\nabla r(h)\|$ .*

**Lemma B.3.** *If the risk functions  $r$  and  $r_K$  are strongly convex with their eigenvalues of Hessian matrices in  $[\rho, 1]$ , then we have for hypothesis  $h \in \mathcal{H}$ , if  $\|h(x_i) - h(x'_k)\| \leq \tilde{O}(L^{5/2} \phi \log^{1/2}(m))$  for  $i \in \mathcal{I}_{\phi,k}$ ,  $k \in [N]$ , the sum risk gradient for a smooth set with  $\phi \leq \tilde{O}(L^{-9/2} \log^{-3}(m))$  with respect to  $h$  satisfies,*

$$\begin{aligned} \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i(h(x_i)) \right\|^2 &\geq 2M_k \rho (\bar{r}_{\mathcal{I},k} - r_{\text{eff},k}) - M_k^2 \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \\ \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i(h(x_i)) \right\|^2 &\leq 2M_k (\bar{r}_{\mathcal{I},k} - r_{\text{eff},k}) + M_k^2 \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \\ \sum_{i \in \mathcal{I}_{\phi,k}} (\mu_i + \lambda_i) \|h(x_i) - y_{\text{eff},k}\|^2 &\leq \frac{1}{\rho^2} O\left(\bar{r}_{\mathcal{I},k} - r_{\text{eff},k} + M_k \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right), \end{aligned}$$

where  $u_i$  is defined in Eqn. (9) and  $M_k = \sum_{\mathcal{I}_{\phi,k}} (\mu_i + \lambda_i)$ .

Lemma B.2 and B.3 are proved in Section C.2.

## B.4. Proof of Theorem 4.1.

In this section, we prove the convergence for informed risks in Eqn. (3), Eqn. (5). First, the gradient lower bound, semi-smoothness of the risk function, and initialized risk bound are proved.

**Lemma B.4** (Gradient Lower Bound). *For any  $\mathbf{W} : \|\mathbf{W} - \mathbf{W}^{(0)}\| \leq \tau$  with  $\tau = O(N^{-9/2} \phi^{3/2} \rho^{3/2} \bar{\lambda}^{3/2} \alpha^{3/2} L^{-15/2} \log^{-3/2}(m))$  and  $\phi \leq \tilde{O}(L^{-9/2} \log^{-3}(m))$ , with Assumption 1 satisfied, we have with probability at least  $1 - O(\phi)$  over the randomness of  $\mathbf{W}^{(0)}$ , the gradient of label-based data risk satisfies*

$$\left\| \nabla_{\mathbf{W}} \hat{R}_1(\mathbf{W}) \right\|_F^2 \geq \Omega\left(\frac{\alpha m \phi \rho \bar{\lambda}}{dN^2}\right) \left(\hat{R}_1(\mathbf{W}) - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right).$$

where  $\hat{R}_{\text{eff}} = \sum_{k=1}^N r_{\text{eff},k}$ , and  $\bar{\lambda}$  is a parameter with lower bound  $\Omega(\min(1 - \lambda, \lambda) \mathbf{1}(\lambda \in (0, 1)) + \mathbf{1}(\lambda \in \{0, 1\}))$ .

The proof of Lemma B.4 can be found in Section C.3.

**Lemma B.5.** For any  $\mathbf{W}, \mathbf{W}' \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ ,  $\tau \in \left[ \Omega(d^{3/2}m^{-3/2}L^{-5/2}\log^{-3/2}(m)), O(L^{-9/2}[\log^{-3}(m)]) \right]$  and  $\phi \leq \tilde{O}(L^{-9/2}\log^{-3}(m))$ , with probability at least  $1 - O(\phi)$  over the randomness of  $\mathbf{W}^{(0)}$ , we have

$$\begin{aligned} \hat{R}_I(\mathbf{W}') &\leq \hat{R}_I(\mathbf{W}) + \left\langle \nabla_{\mathbf{W}} \hat{R}_I(\mathbf{W}), \mathbf{W}' - \mathbf{W} \right\rangle + O(L^2m/d) \|\widehat{\mathbf{W}}\|^2 \\ &\quad + \left( \sqrt{\left( \hat{R}_I(\mathbf{W}) - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2}\phi\log^{1/2}(m)) \right)} \right) O\left( N^{1/2}\tau^{1/3}L^{5/2}\sqrt{m\log(m)}d^{-1/2} \right) \|\widehat{\mathbf{W}}\| \end{aligned}$$

The proof of Lemma B.5 can be found in Section C.4.

**Lemma B.6.** If  $m \geq \Omega(L\log(NL\phi^{-1}))$  and  $\phi \leq \tilde{O}(L^{-9/2}\log^{-3}(m))$ , with probability at least  $1 - O(\phi)$  over the randomness of  $\mathbf{W}^{(0)}$ , at initialization, we have for any  $x_i, i \in [n']$ ,

$$\begin{aligned} \|h_{\mathbf{W}^{(0)}, i}\| &\leq O\left(\log^{1/2}(1/\phi)\right), \\ \text{and } \hat{R}_I(\mathbf{W}^{(0)}) - \hat{R}_{\text{eff}} &\leq O\left(\log^{1/2}(1/\phi)\right). \end{aligned}$$

The proof of Lemma B.6 can be found in Section C.5.

### Proof of Theorem 4.1.

*Proof.* **Convergence of the informed risk.** We first assume  $\tau = \frac{\Gamma}{\sqrt{m}}$  with  $\Gamma = Nd^{1/2}\phi^{-1/2}\rho^{-1/2}\bar{\lambda}^{-1/2}\alpha^{-1/2}$ . Hence, with the choice of  $m$ , we have  $\tau = O(N^{-9/2}\phi^{3/2}\rho^{3/2}\bar{\lambda}^{3/2}\alpha^{3/2}L^{-15/2}\log^{-3/2}(m))$ . We get the recursion inequality based on gradient descent. By the weight update rule of gradient descent, we have  $\mathbf{W}^{(t)} - \mathbf{W}^{(t-1)} = -\eta \nabla \hat{R}_I(\mathbf{W}^{(t-1)})$ . Let  $\Psi = \tilde{O}(L^{5/2}\phi\log^{1/2}(m))$ . By Lemma B.5, we have

$$\begin{aligned} &\hat{R}_I(\mathbf{W}^{(t+1)}) - \hat{R}_{\text{eff}} - \Psi \\ &\leq \hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi - (\eta - O(\eta^2L^2m/d)) \left\| \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\|^2 \\ &\quad + \eta \sqrt{2N \left( \hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi \right)} O\left( \tau^{1/3}L^{5/2}\sqrt{m\log(m)}d^{-1/2} \right) \left\| \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\| \\ &\leq \hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi - \Omega(\eta) \left\| \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\|^2 \\ &\quad + \eta\Omega \left( N^{3/2}\tau^{1/3}L^{5/2}\log^{1/2}(m)\phi^{-1/2}\rho^{-1/2}\bar{\lambda}^{-1/2}\alpha^{-1/2} \right) \left\| \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\|^2 \\ &\leq \hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi - \Omega(\eta) \left\| \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\|^2, \end{aligned} \tag{11}$$

where the second inequality holds by the choice of  $\eta = O(\frac{d}{L^2m})$  such that  $O(\eta L^2m/d) = O(1)$  and the gradient lower bound in Lemma B.4, and the last inequality holds by the choice of  $m \geq \Omega(N^{11}L^{15}d\phi^{-4}\rho^{-4}\bar{\lambda}^{-4}\alpha^{-4}\log^3(m))$  such that  $\Omega(N^{3/2}\tau^{1/3}L^{5/2}\log^{1/2}(m)\phi^{-1/2}\rho^{-1/2}\bar{\lambda}^{-1/2}\alpha^{-1/2}) \leq O(1)$ .

Further, by Lemma B.4, we have

$$\hat{R}_I(\mathbf{W}^{(t+1)}) - \hat{R}_{\text{eff}} - \Psi \leq \left( 1 - \Omega\left(\frac{\eta\alpha m\phi\rho\bar{\lambda}}{dN^2}\right) \right) \left( \hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi \right).$$

Based on the iteration of the recursion inequality, with probability at least  $1 - O(\phi)$ , we have

$$\begin{aligned} &\hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi \\ &\leq \left( 1 - \Omega\left(\frac{\eta\alpha m\phi\rho\bar{\lambda}}{dN^2}\right) \right)^t \left( \hat{R}_I(\mathbf{W}^{(0)}) - \hat{R}_{\text{eff}} - \Psi \right) \\ &\leq \left( 1 - \Omega\left(\frac{\eta\alpha m\phi\rho\bar{\lambda}}{dN^2}\right) \right)^t O(\log^{1/2}(1/\phi)), \end{aligned}$$



where the last inequality comes from Lemma B.6. Then, by taking logarithm, we get

$$\begin{aligned} \ln \left( \left( \hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi \right) \right) &\leq t \ln \left( 1 - \Omega \left( \frac{\eta \alpha m \phi \rho \bar{\lambda}}{dN^2} \right) \right) + \frac{1}{2} \ln \log(1/\phi) \\ &\leq -t \Omega \left( \frac{\eta \alpha m \phi \rho \bar{\lambda}}{dN^2} \right) + \frac{1}{2} \ln \log(1/\phi). \end{aligned}$$

Since  $\eta = O(\frac{d}{L^2 m})$ , after  $T = O\left(\frac{L^2 N^2}{\phi \rho \lambda \alpha} \ln(\epsilon^{-1} \log(\phi^{-1}))\right)$  iterations, for any  $\epsilon > 0$ , we have

$$\hat{R}_I(\mathbf{W}^{(T)}) - \hat{R}_{\text{eff}} \leq O(L^{5/2} \phi \log^{1/2}(1/\phi) \log^{1/2} m) + \epsilon.$$

By setting  $\phi$  as  $\phi \log^{1/2}(1/\phi) \leq \epsilon L^{-5/2} \log^{-1/2} m$  (which satisfies the assumption of  $\phi$  in Theorem 4.1), we can bound  $\hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}}$  by a small positive quantity  $\epsilon$ .

**Verify the weight update range.** Now, we verify that the assumption  $\|\mathbf{W}^{(t)} - \mathbf{W}^{(0)}\| \leq \frac{\Gamma}{\sqrt{m}}$  holds. Denote  $\bar{R}_I(\mathbf{W}^{(t)}) = \hat{R}_I(\mathbf{W}^{(t)}) - \hat{R}_{\text{eff}} - \Psi$ . By Eqn. (11), we have

$$\bar{R}_I(\mathbf{W}^{(t+1)}) - \bar{R}_I(\mathbf{W}^{(t)}) \leq -\Omega(\eta) \left\| \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\|^2.$$

Then, we have

$$\begin{aligned} \sqrt{\bar{R}_I(\mathbf{W}^{(t+1)})} - \sqrt{\bar{R}_I(\mathbf{W}^{(t)})} &= \frac{\bar{R}_I(\mathbf{W}^{(t+1)}) - \bar{R}_I(\mathbf{W}^{(t)})}{\sqrt{\bar{R}_I(\mathbf{W}^{(t+1)})} + \sqrt{\bar{R}_I(\mathbf{W}^{(t)})}} \\ &\leq \frac{-\Omega(\eta) \left\| \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\|^2}{2\sqrt{\bar{R}_I(\mathbf{W}^{(t)})}} \leq -O\left(\frac{m^{1/2} \phi^{1/2} \rho^{1/2} \bar{\lambda}^{1/2} \alpha^{1/2}}{d^{1/2} N}\right) \left\| \eta \nabla \hat{R}_I(\mathbf{W}^{(t)}) \right\| \end{aligned}$$

where the last inequality follows from Lemma B.4.

By the triangle inequality, for any  $t \in [T]$ , we have

$$\begin{aligned} \left\| \mathbf{W}^{(t)} - \mathbf{W}^{(0)} \right\| &\leq \sum_{s=0}^t \left\| \eta \nabla \hat{R}_I(\mathbf{W}^{(s)}) \right\| \\ &\leq O\left(\frac{d^{1/2} N}{m^{1/2} \phi^{1/2} \rho^{1/2} \bar{\lambda}^{1/2} \alpha^{1/2}}\right) \sqrt{\bar{R}_I(\mathbf{W}^{(0)})} \leq O\left(\frac{\Gamma}{\sqrt{m}}\right), \end{aligned} \tag{12}$$

where  $\Gamma = Nd^{1/2} \phi^{-1/2} \rho^{-1/2} \bar{\lambda}^{-1/2} \alpha^{-1/2}$ . Hence, with the choice of  $m$ , we have  $\tau = O(N^{-9/2} \phi^{3/2} \rho^{3/2} \bar{\lambda}^{3/2} \alpha^{3/2} L^{-15/2} \log^{-3/2}(m))$ .

**Convergence of network output.** By Lemma B.3, we have

$$\sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi, k}} (\mu_i + \lambda_i) \left\| h_{\mathbf{W}^{(T)}}(x_i) - y_{\text{eff}, k} \right\|^2 \leq \frac{1}{\rho^2} O\left(\left(\hat{R}_I(\mathbf{W}^{(T)}) - \hat{R}_{\text{eff}}\right) + \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right) \leq O(\epsilon).$$

Denoting  $k(x_i)$  as the index of the cell containing  $x_i$  and rearranging the above summation, we have

$$\sum_{x_i \in S_z} \mu_i \left\| h_{\mathbf{W}^{(T)}}(x_i) - y_{\text{eff}, k(x_i)} \right\|^2 + \sum_{x_j \in S_g} \lambda_j \left\| h_{\mathbf{W}^{(T)}}(x_j) - y_{\text{eff}, k(x_j)} \right\|^2 \leq O(\epsilon).$$

□

## B.5. Proof of Generalization

In this section, we prove the generalization bound based on Rademacher complexity. We first present the bound of Rademacher complexity for neural networks.

**Lemma B.7.** [Theorem 3.3 in (Bartlett et al., 2017), Lemma A.3 in (Chen et al., 2021b)] If risk functions are 1-Lipschitz continuous, with probability at least  $1 - L \exp(-\Omega(m))$ , the Rademacher complexity  $\mathfrak{R}_S(\mathcal{F})$  for the risk set  $\mathcal{F} = \{r(h_{\mathbf{W}}(x), y) : (x, y) \in \mathcal{X} \times \mathcal{Y}, \|\mathbf{W} - \mathbf{W}^{(0)}\| \leq \tau\}$ ,  $\tau = \frac{\Gamma}{\sqrt{m}}$  with  $\Gamma = Nd^{1/2}\phi^{-1/2}\rho^{-1/2}\bar{\lambda}^{-1/2}\alpha^{-1/2}$  given a dataset  $S$  of  $n$  samples is bounded as

$$\mathfrak{R}_S(\mathcal{F}) \leq \Phi/\sqrt{n}, \quad (13)$$

where  $\Phi = O(4^L L^{3/2} m^{1/2} \phi^{-b-1/2} d \rho^{-1/2} \bar{\lambda}^{-1/2} \alpha^{-1/2})$ .

Then, we need to bound the error between effective labels in Definition 2 and the output of optimal hypothesis in Definitions 3 and 4.

**Lemma B.8.** Consistent with Definition 2, assume that for any smooth set  $k \in \mathcal{U}_\phi(S_z)$  (containing at least one labeled sample),  $y_{\text{eff},k}$  equivalently minimizes  $\sum_{i \in I_{\phi,k}} \frac{1-\beta}{n_z} \mathbb{1}(x_i \in S_z) r(h, z_i) + \frac{\beta}{n'_g} \mathbb{1}(x_i \in S'_g) r_{\mathbf{K}}(h, g_i)$ , and for any smooth set  $k \in [N] \setminus \mathcal{U}_\phi(S_z)$  (not containing labeled sample),  $y_{\text{eff},k}$  equivalently minimizes  $\sum_{i \in I_{\phi,k}} \frac{1}{n'_g} r_{\mathbf{K}}(h, g_i)$ .

(a) Letting  $h_{\mathbf{K}}^*$  and  $h_{\mathbf{R},\beta}^*$  be the optimal hypothesis for empirical risks in Definitions 3 and 4, respectively, we have with probability at least  $1 - O(\phi)$  over the randomness of  $\mathbf{W}^{(0)}$ ,

$$\frac{1}{n''_g} \sum_{S''_g} \|h_{\mathbf{K},i}^* - y_{\text{eff},k(x_i)}\|^2 \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right),$$

and

$$\frac{1-\beta}{n_z} \sum_{S_z} \|h_{\mathbf{R},\beta,i}^* - y_{\text{eff},k(x_i)}\|^2 + \frac{\beta}{n'_g} \sum_{S'_g} \|h_{\mathbf{R},\beta,i}^* - y_{\text{eff},k(x_i)}\|^2 \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right),$$

where  $\tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) = O(L^{5/4} \phi^{1/2} \log^{1/4}(1/\phi) \log^{1/4} m)$ .

(b) Letting  $\bar{h}_{\mathbf{K}}^*$  and  $\bar{h}_{\mathbf{R},\beta}^*$  be the optimal hypothesis for the expected risks in Definitions 3 and 4, respectively, we have with probability at least  $1 - O(\phi) - \delta$ ,

$$\frac{1}{n''_g} \sum_{S''_g} \|\bar{h}_{\mathbf{K},i}^* - y_{\text{eff},k(x_i)}\|^2 \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n''_g}}\right),$$

and

$$\frac{1-\beta}{n_z} \sum_{S_z} \|\bar{h}_{\mathbf{R},\beta,i}^* - y_{\text{eff},k(x_i)}\|^2 + \frac{\beta}{n'_g} \sum_{S'_g} \|\bar{h}_{\mathbf{R},\beta,i}^* - y_{\text{eff},k(x_i)}\|^2 \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n_z}}\right),$$

where  $\tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) = O(L^{5/4} \phi^{1/2} \log^{1/4}(1/\phi) \log^{1/4} m)$ .

Proof of Lemma B.8 is given in Section C.6.

### B.5.1. PROOF OF THEOREM 4.2

*Proof.* By generalization bound with Rademacher complexity and Lemma B.7, the population risk is bounded with probability at least  $1 - \delta$ ,  $\delta \in (0, 1)$  as

$$\begin{aligned} & R(\mathbf{W}^{(T)}) \\ &= (1-\lambda)R(\mathbf{W}^{(T)}) + \lambda R(\mathbf{W}^{(T)}) \\ &\leq \frac{1-\lambda}{n_z} \sum_{S_z} r(h_{\mathbf{W}^{(T)},i}, y_i) + \frac{\lambda}{n'_g} \sum_{S'_g} r(h_{\mathbf{W}^{(T)},i}, y_i) + O\left(\Phi + \sqrt{\log(1/\delta)}\right) \left( (1-\lambda)\sqrt{\frac{1}{n_z}} + \lambda\sqrt{\frac{1}{n'_g}} \right) \end{aligned} \quad (14)$$

For the empirical risk, we have

$$\begin{aligned}
 & \frac{1-\lambda}{n_z} \sum_{S_z} r(h_{\mathbf{W}^{(T)},i}, y_i) + \frac{\lambda}{n_g} \sum_{S_g} r(h_{\mathbf{W}^{(T)},i}, y_i) \\
 & \leq \frac{1-\lambda}{n_z} \sum_{S_z} (r(y_{\text{eff},k(x_i)}, y_i) + \|h_{\mathbf{W}^{(T)},i} - y_{\text{eff},k(x_i)}\|) + \frac{\lambda}{n_g} \sum_{S_g} [r(y_{\text{eff},k(x_i)}, y_i) + \|h_{\mathbf{W}^{(T)},i} - y_{\text{eff},k(x_i)}\|] \quad (15) \\
 & \leq \sqrt{\epsilon} + \frac{1-\lambda}{n_z} \sum_{S_z} r(y_{\text{eff},k(x_i)}, y_i) + \frac{\lambda}{n_g} \sum_{S_g} r(y_{\text{eff},k(x_i)}, y_i)
 \end{aligned}$$

where the first inequality holds because of the 1-Lipschitz of risk functions such that  $r(h_{\mathbf{W}^{(t)},i}, y_i) - r(y_{\text{eff},k(x_i)}, y_i) \leq \|h_{\mathbf{W}^{(t)},i} - y_{\text{eff},k(x_i)}\|$ , and the second inequality follows from the convergence of network output in Theorem 4.1.

Since with  $\beta_\lambda = \frac{\lambda n'_g}{(1-\lambda)n_g + \lambda n'_g}$ , the training objective in Eqn.(3) can also be written as  $\hat{R}_1(\mathbf{W}) = \left(1 - \lambda + \frac{\lambda n'_g}{n_g}\right) \left(\frac{1-\beta_\lambda}{n_z} \sum_{S_z} r(h_{\mathbf{W},i}, z_i) + \frac{\beta_\lambda}{n'_g} \sum_{S'_g} r_K(h_{\mathbf{W},i}, g_i)\right) + \frac{\lambda n''_g}{n_g} \sum_{S''_g} r_K(h_{\mathbf{W},i}, g_i)$ , for any smooth set  $k \in \mathcal{U}_\phi(S_z)$  (containing at least one labeled sample),  $y_{\text{eff},k}$  equivalently minimizes  $\sum_{i \in I_{\phi,k}} \frac{1-\beta_\lambda}{n_z} \mathbb{1}(x_i \in S_z) r(h, z_i) + \frac{\beta_\lambda}{n'_g} \mathbb{1}(x_i \in S'_g) r_K(h, g_i)$  by Definition 2. Thus, the bounds of the differences between optimal hypothesis and effective labels in Lemma B.8 hold for  $\beta = \beta_\lambda$  and  $h_{\mathbf{R},\beta_\lambda}^*$ . Next, we can bound the total effective risk in terms of label and knowledge imperfectness, with probability at least  $1 - \tilde{O}(\phi)$ ,

$$\begin{aligned}
 & \frac{1-\lambda}{n_z} \sum_{S_z} r(y_{\text{eff},k(x_i)}, y_i) + \frac{\lambda}{n_g} \sum_{S_g} r(y_{\text{eff},k(x_i)}, y_i) \\
 & = \frac{1-\lambda}{n_z} \sum_{S_z} r(y_{\text{eff},k(x_i)}, y_i) + \frac{\lambda}{n_g} \sum_{S'_g} r(y_{\text{eff},k(x_i)}, y_i) + \frac{\lambda}{n_g} \sum_{S''_g} r(y_{\text{eff},k(x_i)}, y_i) \\
 & \leq \frac{1-\lambda}{n_z} \sum_{S_z} r(h_{\mathbf{R},\beta_\lambda,i}^*, y_i) + \frac{\lambda}{n_g} \sum_{S'_g} r(h_{\mathbf{R},\beta_\lambda,i}^*, y_i) + \frac{\lambda}{n_g} \sum_{S''_g} r(h_{\mathbf{K},i}, y_i) \quad (16) \\
 & \quad + \frac{1-\lambda}{n_z} \sum_{S_z} \|h_{\mathbf{R},\beta_\lambda,i}^* - y_{\text{eff},k(x_i)}\| + \frac{\lambda}{n_g} \sum_{S'_g} \|h_{\mathbf{R},\beta_\lambda,i}^* - y_{\text{eff},k(x_i)}\| + \frac{\lambda}{n_g} \sum_{S''_g} \|h_{\mathbf{K},i}^* - y_{\text{eff},k(x_i)}\| \\
 & \leq (1-\lambda) \hat{Q}_{\mathbf{R},S_z,S'_g}(\beta_\lambda) + \lambda \hat{Q}_{\mathbf{K},S''_g} + O(\sqrt{\epsilon})
 \end{aligned}$$

where  $\beta_\lambda = \frac{\lambda n'_g}{(1-\lambda)n_g + \lambda n'_g}$ , the first inequality comes from the Lipschitz continuity of risk functions, and the last inequality holds by Lemma B.8 and the assumption  $\phi \leq \tilde{O}(\epsilon^2 L^{-9/2} \log^{-3}(m))$  such that  $\frac{1-\lambda}{n_z} \sum_{S_z} \|h_{\mathbf{R},\beta_\lambda,i}^* - y_{\text{eff},k(x_i)}\| + \frac{\lambda}{n_g} \sum_{S'_g} \|h_{\mathbf{R},\beta_\lambda,i}^* - y_{\text{eff},k(x_i)}\| \leq \left(1 - \lambda + \frac{\lambda n'_g}{n_g}\right) \sqrt{\frac{1-\beta_\lambda}{n_z} \sum_{S_z} \|h_{\mathbf{R},\beta_\lambda,i}^* - y_{\text{eff},k(x_i)}\|^2 + \frac{\beta_\lambda}{n'_g} \sum_{S'_g} \|h_{\mathbf{R},\beta_\lambda,i}^* - y_{\text{eff},k(x_i)}\|^2} \leq O(\sqrt{\epsilon})$  and  $\frac{\lambda}{n_g} \sum_{S'_g} \|h_{\mathbf{R},\beta_\lambda,i}^* - y_{\text{eff},k(x_i)}\| \leq \frac{\lambda n''_g}{n_g} \sqrt{\frac{1}{n''_g} \sum_{S''_g} \|h_{\mathbf{K},i}^* - y_{\text{eff},k(x_i)}\|^2} \leq O(\sqrt{\epsilon})$ . In the last inequality of (16), we absorb  $\frac{\lambda}{n_g} \sum_{S'_g} r(h_{\mathbf{R},\beta_\lambda,i}^*, y_i)$  into  $O(\sqrt{\epsilon})$  because the risk functions are upper bounded and  $S'_g$  is the set of samples sharing the same smooth sets with  $S_z$ , and so  $\frac{\lambda}{n_g} \sum_{S'_g} r(h_{\mathbf{R},\beta_\lambda,i}^*, y_i) \leq O(\frac{n'_g}{n_g}) \leq O(n_z \phi^b) \leq O(\sqrt{\epsilon})$ .

Substituting Eqn.(16) and (15) into Eqn. (14), the population risk is bounded with probability at least  $1 - O(\phi) - \delta$ ,  $\delta \in (0, 1)$  as

$$\begin{aligned}
 R(\mathbf{W}^{(T)}) & = (1-\lambda)R(\mathbf{W}^{(T)}) + \lambda R(\mathbf{W}^{(T)}) \\
 & \leq \sqrt{\epsilon} + (1-\lambda) \hat{Q}_{\mathbf{R},S_z,S'_g}(\beta_\lambda) + \lambda \hat{Q}_{\mathbf{K},S''_g} + O\left(\Phi + \sqrt{\log(1/\delta)}\right) \left( (1-\lambda) \sqrt{\frac{1}{n_z}} + \lambda \sqrt{\frac{1}{n_g}} \right).
 \end{aligned}$$

□

## B.5.2. PROOF OF THEOREM 5.1

*Proof.* Based on the construction of smooth sets in Definition 1, denote  $\mathcal{X}' = \mathcal{X}_\phi(S_z) = \bigcup_{k \in \mathcal{U}_\phi} (S_z) \mathcal{C}_{\phi,k}$  as the region covered by the smooth sets containing at least one sample in  $S_z$ , and let  $\mathcal{X}'' = \mathcal{X}'/\mathcal{X}'$ . Let  $\mathbb{P}_{\mathcal{X}'} = \int_{\mathcal{X}', \mathcal{Y}} p(y | x) dy p(x) dx$  and  $\mathbb{P}_{\mathcal{X}''} = \int_{\mathcal{X}'', \mathcal{Y}} p(y | x) dy p(x) dx$  where  $p(x)$  and  $p(y | x)$  are probability densities. Then we have  $P_{\mathcal{X}'} \leq O(n_z/N) = O(n_z \phi^b) = O(\sqrt{\epsilon})$  by the assumption that  $\phi \leq (\sqrt{\epsilon}/n_z)^{1/b}$ , and  $\mathbb{E}[r(h_{\mathbf{W}^{(T)}}(x), y)] = \mathbb{E}_{\mathbb{P}_{\mathcal{X}'}}[r(h_{\mathbf{W}^{(T)}}(x), y)] \mathbb{P}_{\mathcal{X}'} + \mathbb{E}_{\mathbb{P}_{\mathcal{X}''}}[r(h_{\mathbf{W}^{(T)}}(x), y)] \mathbb{P}_{\mathcal{X}''}$ . By generalization bound with Rademacher complexity and Lemma B.7, the population risk is bounded with probability at least  $1 - \delta$ ,  $\delta \in (0, 1)$  as

$$\begin{aligned} R(\mathbf{W}^{(T)}) &= (1 - \lambda) \mathbb{E}[r(h_{\mathbf{W}^{(T)}}(x), y)] + \lambda \mathbb{E}[r(h_{\mathbf{W}^{(T)}}(x), y)] \\ &= (1 - \lambda) \mathbb{E}[r(h_{\mathbf{W}^{(T)}}(x), y)] + \lambda \mathbb{E}_{\mathbb{P}_{\mathcal{X}''}}[r(h_{\mathbf{W}^{(T)}}(x), y)] + \lambda (\mathbb{E}_{\mathbb{P}_{\mathcal{X}'}}[r(h_{\mathbf{W}^{(T)}}(x), y)] - \mathbb{E}_{\mathbb{P}_{\mathcal{X}''}}[r(h_{\mathbf{W}^{(T)}}(x), y)]) \mathbb{P}_{\mathcal{X}'} \\ &\leq \frac{(1 - \lambda)}{n_z} \sum_{S_z} r(h_{\mathbf{W}^{(T)}, i}, y_i) + \frac{\lambda}{n_g''} \sum_{S_g''} r(h_{\mathbf{W}^{(T)}, i}, y_i) + \lambda O(\sqrt{\epsilon}) + O\left(\Phi + \sqrt{\log(1/\delta)}\right) \left( (1 - \lambda) \sqrt{\frac{1}{n_z}} + \lambda \sqrt{\frac{1}{n_g''}} \right). \end{aligned} \quad (17)$$

Then for the empirical risk, we have

$$\begin{aligned} &\frac{1 - \lambda}{n_z} \sum_{S_z} r(h_{\mathbf{W}^{(T)}, i}, y_i) + \frac{\lambda}{n_g''} \sum_{S_g''} r(h_{\mathbf{W}^{(T)}, i}, y_i) \\ &\leq \frac{1 - \lambda}{n_z} \sum_{S_z} r(y_{\text{eff}, k(x_i)}, y_i) + \frac{\lambda}{n_g''} \sum_{S_g''} r(y_{\text{eff}, k(x_i)}, y_i) + \frac{1 - \lambda}{n_z} \sum_{S_z} \|h_{\mathbf{W}^{(T)}, i} - y_{\text{eff}, i}\| + \frac{\lambda}{n_g''} \sum_{S_g''} \|h_{\mathbf{W}^{(T)}, i} - y_{\text{eff}, i}\| \\ &\leq O(\sqrt{\epsilon}) + \frac{1 - \lambda}{n_z} \sum_{S_z} r(y_{\text{eff}, k(x_i)}, y_i) + \frac{\lambda}{n_g''} \sum_{S_g''} r(y_{\text{eff}, k(x_i)}, y_i), \end{aligned} \quad (18)$$

where the first inequality holds because of the Lipschitz continuity of risk functions such that  $r(h_{\mathbf{W}^{(T)}, i}, y_i) - r(y_{\text{eff}, k(x_i)}, y_i) \leq \|h_{\mathbf{W}^{(T)}, i} - y_{\text{eff}, k(x_i)}\|$ , and the second inequality follows from the convergence of network output in Theorem 4.1 (By Theorem 4.1, we have  $\frac{(1 - \lambda)(1 - \beta)}{n_z} \sum_{S_z} \|h_{\mathbf{W}^{(T)}, i} - y_{\text{eff}, i}\|^2 \leq O(\epsilon)$  and  $\frac{\lambda}{n_g''} \sum_{S_g''} \|h_{\mathbf{W}^{(T)}, i} - y_{\text{eff}, i}\|^2 \leq O(\epsilon)$ , and so  $\frac{1 - \lambda}{n_z} \sum_{S_z} \|h_{\mathbf{W}^{(T)}, i} - y_{\text{eff}, i}\| \leq \sqrt{\frac{1 - \lambda}{1 - \beta}} O(\sqrt{\epsilon}) = O(\sqrt{\epsilon})$  and  $\frac{\lambda}{n_g''} \sum_{S_g''} \|h_{\mathbf{W}^{(T)}, i} - y_{\text{eff}, i}\| \leq \sqrt{\lambda} O(\sqrt{\epsilon}) = O(\sqrt{\epsilon})$ ).

Next, we bound the empirical risk in terms of label and knowledge imperfectness as follows:

$$\begin{aligned} &\frac{1 - \lambda}{n_z} \sum_{S_z} r(y_{\text{eff}, k(x_i)}, y_i) + \frac{\lambda}{n_g''} \sum_{S_g''} r(y_{\text{eff}, k(x_i)}, y_i) \\ &\leq \frac{1 - \lambda}{n_z} \sum_{S_z} r(h_{\mathbf{R}, \beta, i}^*, y_i) + \frac{\lambda}{n_g''} \sum_{S_g''} r(h_{\mathbf{K}, i}, y_i) + \frac{1 - \lambda}{n_z} \sum_{S_z} \|h_{\mathbf{R}, \beta, i}^* - y_{\text{eff}, k(x_i)}\| + \frac{\lambda}{n_g''} \sum_{S_g''} \|h_{\mathbf{K}, i}^* - y_{\text{eff}, k(x_i)}\| \quad (19) \\ &\leq (1 - \lambda) \widehat{Q}_{\mathbf{R}, S_z, S_g'}(\beta) + \lambda \widehat{Q}_{\mathbf{K}, S_g''} + O(\sqrt{\epsilon}) \end{aligned}$$

where the first inequality comes from the Lipschitz continuity of risk functions, and the concavity of squared root, the last inequality holds by Definitions 3 and 4, and Lemma B.8 and the assumption of  $\phi$  such that  $\phi \log^{1/2}(1/\phi) \leq O(\epsilon^2 L^{-5/2} \log^{-1/2}(m))$ . Concretely, by Lemma B.8 (a), we have  $\frac{1 - \beta}{n_z} \sum_{S_z} \|h_{\mathbf{R}, \beta, i}^* - y_{\text{eff}, k(x_i)}\|^2 \leq \widetilde{O}(L^{5/4} \phi^{1/2} \log^{1/4}(m))$  and  $\frac{1}{n_g''} \sum_{S_g''} \|h_{\mathbf{K}, i}^* - y_{\text{eff}, k(x_i)}\|^2 \leq \widetilde{O}(L^{5/4} \phi^{1/2} \log^{1/4}(m))$ , and so it holds that  $\frac{1 - \lambda}{n_z} \sum_{S_z} \|h_{\mathbf{R}, \beta, i}^* - y_{\text{eff}, k(x_i)}\| \leq \frac{1 - \lambda}{\sqrt{1 - \beta}} \widetilde{O}(L^{5/8} \phi^{1/4} \log^{1/8}(m))$  and  $\frac{\lambda}{n_g''} \sum_{S_g''} \|h_{\mathbf{K}, i}^* - y_{\text{eff}, k(x_i)}\| \leq \lambda \widetilde{O}(L^{5/8} \phi^{1/4} \log^{1/8}(m))$ . Thus we obtain the last inequality of (19) by the assumption  $\phi \log^{1/2}(1/\phi) \leq O(\epsilon^2 L^{-5/2} \log^{-1/2}(m))$ .

Substituting Eqns. (19) and (18) into Eqn. (17), we have

$$\begin{aligned} R(\mathbf{W}^{(T)}) &= (1 - \lambda)\mathbb{E}[r(h_{\mathbf{W}^{(T)}}(x), y)] + \lambda\mathbb{E}[r(h_{\mathbf{W}^{(T)}}(x), y)] \\ &\leq O(\sqrt{\epsilon}) + (1 - \lambda)\widehat{Q}_{R, S_z, S'_g}(\beta) + \lambda\widehat{Q}_{K, S''_g} + O\left(\Phi + \sqrt{\log(1/\delta)}\right) \left( (1 - \lambda)\sqrt{\frac{1}{n_z}} + \lambda\sqrt{\frac{1}{n''_g}} \right). \end{aligned}$$

□

### B.5.3. PROOF OF COROLLARY 5.2

*Proof.* First, following (17), with probability at least  $1 - O(\phi) - \delta$ ,  $\delta \in (0, 1)$ , it holds that

$$\begin{aligned} R(\mathbf{W}^{(T)}) &\leq \frac{(1 - \lambda)}{n_z} \sum_{S_z} r(h_{\mathbf{W}^{(T)}, i}, y_i) + \frac{\lambda}{n''_g} \sum_{S''_g} r(h_{\mathbf{W}^{(T)}, i}, y_i) + \lambda O(\sqrt{\epsilon}) \\ &+ O\left(\Phi + \sqrt{\log(1/\delta)}\right) \left( (1 - \lambda)\sqrt{\frac{1}{n_z}} + \lambda\sqrt{\frac{1}{n''_g}} \right) \end{aligned} \quad (20)$$

With the same reason as in Eqn. (18), we have

$$\begin{aligned} &\frac{1 - \lambda}{n_z} \sum_{S_z} r(h_{\mathbf{W}^{(T)}, i}, y_i) + \frac{\lambda}{n''_g} \sum_{S''_g} r(h_{\mathbf{W}^{(T)}, i}, y_i) \\ &\leq O(\sqrt{\epsilon}) + \frac{1 - \lambda}{n_z} \sum_{S_z} r(y_{\text{eff}, k(x_i)}, y_i) + \frac{\lambda}{n''_g} \sum_{S''_g} r(y_{\text{eff}, k(x_i)}, y_i). \end{aligned} \quad (21)$$

Then, unlike in the proof of Theorem 5.1, we need to bound the risk in (21) in terms of expected label and knowledge imperfectness. Thus, replacing  $h_{R, \beta}^*$  and  $h_K^*$  in Eqn. (19) with  $\bar{h}_{R, \beta}^*$  and  $\bar{h}_K^*$ , we have

$$\begin{aligned} &\frac{1 - \lambda}{n_z} \sum_{S_z} r(y_{\text{eff}, k(x_i)}, y_i) + \frac{\lambda}{n''_g} \sum_{S''_g} r(y_{\text{eff}, k(x_i)}, y_i) \\ &\leq (1 - \lambda)\widehat{Q}_{R, S_z, S'_g}(\beta) + \lambda\widehat{Q}_{K, S''_g} + \frac{1 - \lambda}{n_z} \sum_{S_z} \|\bar{h}_{R, \beta, i}^* - y_{\text{eff}, k(x_i)}\| + \frac{\lambda}{n''_g} \sum_{S''_g} \|\bar{h}_{K, i}^* - y_{\text{eff}, k(x_i)}\| \end{aligned} \quad (22)$$

By Lemma B.8 (b), it holds that  $\frac{1}{n''_g} \sum_{S''_g} \|\bar{h}_{K, i}^* - y_{\text{eff}, k(x_i)}\|^2 \leq \tilde{O}\left(L^{5/4}\phi^{1/2}\log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n''_g}}\right)$  and  $\frac{1 - \lambda}{n_z} \sum_{S_z} \|\bar{h}_{R, \beta, i}^* - y_{\text{eff}, k(x_i)}\|^2 \leq \tilde{O}\left(L^{5/4}\phi^{1/2}\log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n_z}}\right)$ . Thus we have  $\frac{\lambda}{n''_g} \sum_{S''_g} \|\bar{h}_{K, i}^* - y_{\text{eff}, k(x_i)}\| \leq \lambda\left(\tilde{O}\left(L^{5/8}\phi^{1/4}\log^{1/8}(m)\right) + O\left(\left(\frac{\log(1/\delta)}{n''_g}\right)^{\frac{1}{4}}\right)\right) \leq \lambda\left(O(\epsilon) + O\left(\left(\frac{\log(1/\delta)}{n''_g}\right)^{\frac{1}{4}}\right)\right)$  and  $\frac{1 - \lambda}{n_z} \sum_{S_z} \|\bar{h}_{R, \beta, i}^* - y_{\text{eff}, k(x_i)}\| \leq \frac{1 - \lambda}{\sqrt{1 - \beta}}\left(\tilde{O}\left(L^{5/8}\phi^{1/4}\log^{1/8}(m)\right) + O\left(\left(\frac{\log(1/\delta)}{n_z}\right)^{\frac{1}{4}}\right)\right) \leq (1 - \lambda)\left(O(\epsilon) + O\left(\left(\frac{\log(1/\delta)}{n_z}\right)^{\frac{1}{4}}\right)\right)$ . Therefore, continuing with (22), it holds that

$$\begin{aligned} &\frac{1 - \lambda}{n_z} \sum_{S_z} r(y_{\text{eff}, k(x_i)}, y_i) + \frac{\lambda}{n''_g} \sum_{S''_g} r(y_{\text{eff}, k(x_i)}, y_i) \\ &\leq (1 - \lambda)\widehat{Q}_{R, S_z, S'_g}(\beta) + \lambda\widehat{Q}_{K, S''_g} + O(\sqrt{\epsilon}) + O\left((1 - \lambda)\left(\frac{\log(1/\delta)}{n_z}\right)^{\frac{1}{4}} + \lambda\left(\frac{\log(1/\delta)}{n''_g}\right)^{\frac{1}{4}}\right) \\ &\leq O(\sqrt{\epsilon}) + (1 - \lambda)Q_R(\beta) + \lambda Q_K + O\left((1 - \lambda)\left(\frac{\log(1/\delta)}{n_z}\right)^{\frac{1}{4}} + \lambda\left(\frac{\log(1/\delta)}{n''_g}\right)^{\frac{1}{4}}\right), \end{aligned} \quad (23)$$

where the second inequality holds by Lemma B.8 and the last inequality holds by McDiarmid's inequality. Finally, substituting Eqns. (23) and (21) into Eqn. (20), with probability at least  $1 - O(\phi) - \delta$ ,  $\delta \in (0, 1)$ , it holds that

$$R(\mathbf{W}^{(T)}) \leq O(\sqrt{\epsilon}) + (1 - \lambda)Q_R(\beta) + \lambda Q_K + O\left(\Phi + \log^{1/4}(1/\delta)\right) \sqrt{\frac{1 - \lambda}{\sqrt{n_z}} + \frac{\lambda}{\sqrt{n_g''}}}.$$

This completes the proof.  $\square$

#### B.5.4. PROOF OF COROLLARY 5.3

*Proof. Proof of (a).* If  $Q_K \leq \sqrt{\epsilon}$  and  $\lambda$  is set as 1, it holds by Corollary 5.2 that

$$\begin{aligned} R(\mathbf{W}^{(T)}) &\leq O(\sqrt{\epsilon}) + Q_K + O\left(\Phi + \log^{1/4}(1/\delta)\right) \left(\frac{1}{n_g''}\right)^{1/4} \\ &\leq O(\sqrt{\epsilon}) + O\left(\frac{1}{n_g''}\right)^{1/4}, \end{aligned}$$

where in the last inequality we absorb the scales of the last term by  $O$  notation. Thus,  $n_g'' \leq O(\frac{1}{\epsilon^2})$  guarantees that  $R(\mathbf{W}^{(T)}) \leq \sqrt{\epsilon}$ . In the proof of Theorem 5.1, we prove that the probability that a sample belongs to the region covered by the smooth sets containing at least one labeled sample is  $P_{\mathcal{X}'} = O(n_z/N) = O(n_z\phi^b) = O(\sqrt{\epsilon})$ . Thus we have  $P_{\mathcal{X}''} = 1 - P_{\mathcal{X}'} = 1 - O(\sqrt{\epsilon})$ , and so  $n_g = \frac{n_g''}{P_{\mathcal{X}''}} = n_g''/(1 - O(\epsilon)) \sim O(1/(\epsilon^2 - \epsilon^3))$ .

*Proof of (b).* If  $Q_K > \sqrt{\epsilon}$  and  $\lambda = \frac{\sqrt{\epsilon}}{Q_K}$ , then by Corollary 5.2, we have

$$\begin{aligned} R(\mathbf{W}^{(T)}) &\leq O(\sqrt{\epsilon}) + Q_R(\beta^*) - \frac{Q_R(\beta^*)}{Q_K} \sqrt{\epsilon} + \sqrt{\epsilon} + O\left(\sqrt{\left(1 - \frac{\sqrt{\epsilon}}{Q_K}\right) \frac{1}{\sqrt{n_z}} + \frac{\sqrt{\epsilon}}{Q_K} \frac{1}{\sqrt{n_g''}}}\right) \\ &\leq O(\sqrt{\epsilon}) + O\left(\sqrt{\left(1 - \frac{\sqrt{\epsilon}}{Q_K}\right) \frac{1}{\sqrt{n_z}} + \frac{\sqrt{\epsilon}}{Q_K} \frac{1}{\sqrt{n_g''}}}\right), \end{aligned}$$

where the second inequality holds because  $\frac{\sqrt{\epsilon}}{Q_K} + \frac{\sqrt{\epsilon}}{Q_R(\beta^*)} \geq 1$  such that  $Q_R(\beta^*) - \frac{Q_R(\beta^*)}{Q_K} \sqrt{\epsilon} \leq \sqrt{\epsilon}$ . Then to guarantee  $R(\mathbf{W}^{(T)}) \leq \sqrt{\epsilon}$ , we require that  $\left(1 - \frac{\sqrt{\epsilon}}{Q_K}\right) \frac{1}{\sqrt{n_z}} \leq \epsilon$  and  $\frac{\sqrt{\epsilon}}{Q_K} \frac{1}{\sqrt{n_g''}} \leq \epsilon$ . Thus, we have  $n_z \sim O\left(\frac{1}{\epsilon} - \frac{1}{(\sqrt{\epsilon}Q_K)^2}\right)$ ,  $n_g'' \sim O\left(\frac{1}{\epsilon Q_K^2}\right)$  and  $n_g = n_g''/(1 - O(\epsilon)) \sim O(1/((\epsilon - \epsilon^2)Q_K^2))$ .

*Proof of (c).* We prove (c) by contradiction. If  $R(\mathbf{W}^{(T)}) \leq \sqrt{\epsilon}$ , we have  $(1 - \lambda)Q_R(\beta^*) \leq \sqrt{\epsilon}$  and  $\lambda Q_K \leq \sqrt{\epsilon}$ . Then  $\frac{\sqrt{\epsilon}}{Q_R(\beta^*)} + \frac{\sqrt{\epsilon}}{Q_K} \geq 1 - \lambda + \lambda = 1$ . This is contradictory to the condition  $\frac{\sqrt{\epsilon}}{Q_R(\beta^*)} + \frac{\sqrt{\epsilon}}{Q_K} \leq 1$ . Thus completes the proof.  $\square$

## C. Proofs of Lemmas in Appendix B

We now show the proofs of lemmas in Appendix B, while the proofs of lemmas newly introduced in this section are deferred to Appendix D.

### C.1. Proof of Lemma B.1

In this section, we prove the forward perturbation with respect to inputs. We first recall some important notations. For the smooth set  $k \in [N]$ , layer  $l \in [L]$ , let  $h_{l,k} = h_l(x'_k)$ ,  $h_{l,i} = h_l(x_i)$  be the activated output of  $l$ th layer, and  $f_{l,k} = \mathbf{W}_l h_{l-1}(x'_k)$ ,  $f_{l,i} = \mathbf{W}_l h_{l-1}(x_i)$  be the pre-activated output of  $l$ th layer for some weight  $\mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ . At initialization, denote  $h_{l,k}^{(0)} = h_l^{(0)}(x'_k)$ ,  $h_{l,i}^{(0)} = h_l^{(0)}(x_i)$ ,  $f_{l,k}^{(0)} = \mathbf{W}_l^{(0)} h_{l-1}^{(0)}(x'_k)$ ,  $f_{l,i}^{(0)} = \mathbf{W}_l^{(0)} h_{l-1}^{(0)}(x_i)$ , the diagonal matrices  $\mathbf{D}_{l,k}^{(0)} \in \mathbb{R}^{m \times m}$  and  $\mathbf{D}_{l,i}^{(0)} \in \mathbb{R}^{m \times m}$  with  $[\mathbf{D}_{l,k}^{(0)}]_{j,j} = \mathbb{1}([f_{l,k}^{(0)}]_j \geq 0)$  and  $[\mathbf{D}_{l,i}^{(0)}]_{j,j} = \mathbb{1}([f_{l,i}^{(0)}]_j \geq 0)$  for  $i \in \mathcal{I}_{\phi,k}$ ,  $j \in [m]$ . Then we denote for initialization  $f'_{l,i} = f_{l,i}^{(0)} - f_{l,k}^{(0)}$  and the diagonal matrix  $\mathbf{D}'_l \in \mathbb{R}^{m \times m}$  with  $[\mathbf{D}'_l]_{j,j} = [\mathbf{D}_{l,k}^{(0)}]_{j,j} - [\mathbf{D}_{l,i}^{(0)}]_{j,j}$ , omitting the notation (0) and  $i, k$ .

**Lemma C.1.** *If  $f'_{l,i}$  can be written as  $f'_{l,i} = f'_{l,i,1} + f'_{l,i,2}$  with  $\|f'_{l,i,1}\| \leq O(L^{3/2}\phi \log^{1/2}(1/\phi))$  and  $\|f'_{l,i,2}\|_\infty \leq O(L\phi^{2/3} \log^{1/2}(1/\phi)m^{-1/2})$ , then with probability at least  $1 - \exp(-\Omega(m\phi^{2/3}L))$  over the randomness of  $\mathbf{W}^{(0)}$ , we have*

$$\begin{aligned} \|\mathbf{D}'_l f'_{l,i}{}^{(0)}\|_0 &\leq \|\mathbf{D}'_l\|_0 \leq O(m\phi^{2/3}L \log^{1/2}(1/\phi)), \\ \|\mathbf{D}'_l f'_{l,i}{}^{(0)}\| &\leq O(\phi L^{3/2} \log^{1/2}(1/\phi)). \end{aligned}$$

Proof Lemma C.1 is given in Section D.1.

### Proof of Lemma B.1

*Proof.* We first prove the following three conclusions by induction under the assumptions in Lemma B.1: for  $i \in \mathcal{I}_{\phi,k}$ ,  $k \in [N]$ , with probability at least  $1 - O(\phi)$ ,

- (a)  $f'_{l,i}$  at initialization can be written as  $f'_{l,i,1} + f'_{l,i,2}$  with  $\|f'_{l,i,1}\| \leq O(L^{3/2}\phi \log^{1/2}(1/\phi))$  and  $\|f'_{l,i,2}\|_\infty \leq O(L\phi^{2/3}m^{-1/2} \log^{1/2}(1/\phi))$ .
- (b) At initialization,  $\|\mathbf{D}'_l f'_{l,i}{}^{(0)}\|_0 \leq \|\mathbf{D}'_l\|_0 \leq O(m\phi^{2/3}L \log^{1/2}(1/\phi))$ ,  $\|\mathbf{D}'_l f'_{l,i}{}^{(0)}\| \leq O(\phi L^{3/2} \log^{1/2}(1/\phi))$ .
- (c)  $\|h_{l,i}^{(0)} - h_{l,k}^{(0)}\| \leq O(L^{5/2}\phi \sqrt{\log(m) \log(1/\phi)})$  and  $\|f_{l,i}^{(0)} - f_{l,k}^{(0)}\| \leq O(L^{5/2}\phi \sqrt{\log(m) \log(1/\phi)})$ .

When  $l = 0$ , we have  $\|h_{0,i}^{(0)} - h_{0,k}^{(0)}\| = \|x_i - x_k\| \leq O(\phi)$ . Since  $[\mathbf{W}_1^{(0)}]_{j,j} \sim \mathcal{N}(0, \frac{2}{m})$ ,  $j \in [m]$ , we have  $\|f_{1,i}^{(0)} - f_{1,k}^{(0)}\| \leq O(\phi \log^{1/2}(1/\phi))$  with probability at least  $1 - O(\phi)$  over the randomness of  $\mathbf{W}_1^{(0)}$ . By Lemma C.1, the above three conclusions hold. Then we assume the conclusions (a) holds for layer  $a$ ,  $a \leq l - 1$  and prove (a)(b)(c) hold for  $l$ .

First, we re-write  $f'_{l,i}$  as

$$\begin{aligned} f'_{l,i} &= f_{l,i}^{(0)} - f_{l,k}^{(0)} = \mathbf{W}_l^{(0)} \left( \mathbf{D}_{l-1,k}^{(0)} + \mathbf{D}'_{l-1} \right) \left( f_{l-1,k}^{(0)} + f'_{l-1,i} \right) - \mathbf{W}_l^{(0)} \mathbf{D}_{l-1,k}^{(0)} f_{l-1,k}^{(0)} \\ &= \mathbf{W}_l^{(0)} \mathbf{D}'_{l-1} \left( f_{l-1,k}^{(0)} + f'_{l-1,i} \right) + \mathbf{W}_l^{(0)} \mathbf{D}_{l-1,k}^{(0)} f'_{l-1,i} \\ &= \dots \\ &= \sum_{a=2}^l \left( \prod_{b=a+1}^l \mathbf{W}_b^{(0)} \mathbf{D}_{b-1,k}^{(0)} \right) \mathbf{W}_a^{(0)} \mathbf{D}'_{a-1} \left( f_{a-1,k}^{(0)} + f'_{a-1,i} \right) + \mathbf{W}_1^{(0)} (x_i - x_k) \end{aligned}$$

By Lemma C.1, and the inductive assumption (a) for layer  $a$ ,  $a \leq l - 1$ , we have with probability at least  $1 - \exp(-\Omega(m\phi^{2/3}L))$ ,

$$\|\mathbf{D}'_a \left( f_{a,k}^{(0)} + f'_{a,i} \right)\|_0 \leq O(m\phi^{2/3}L \log^{1/2}(1/\phi)), \quad (24)$$

$$\|\mathbf{D}'_a \left( f_{a,k}^{(0)} + f'_{a,i} \right)\| \leq O(\phi L^{3/2} \log^{1/2}(1/\phi)), \quad (25)$$

so (b) holds for layer  $l$ . Then let  $q_a = \left( \prod_{b=a+1}^l \mathbf{W}_b^{(0)} \mathbf{D}_{b-1,k}^{(0)} \right) \mathbf{W}_a^{(0)} \mathbf{D}'_{a-1} \left( f_{a-1,k}^{(0)} + f'_{a-1,i} \right)$ . By Eqn.(24), (25), and Claim 8.5 ( $s = O(m\phi^{2/3}L)$ ) in (Allen-Zhu et al., 2019b), with probability at least  $1 - \exp(-\Omega(m\phi^{2/3}L \log(m)))$ , we can write  $q_a = q_{a,1} + q_{a,2}$  with

$$\|q_{a,1}\| \leq O(\phi^{4/3}L^2 \log(m) \log^{3/4}(1/\phi)) \quad \text{and} \quad \|q_{a,2}\|_\infty \leq O(\phi L^{3/2} \sqrt{\log(m)/m} \log^{1/2}(1/\phi)). \quad (26)$$

Let  $f'_{l,i,1} = \sum_{a=2}^l q_{a,1} + \mathbf{W}_1^{(0)} (x_i - x_k)$  and  $f'_{l,i,2} = \sum_{a=2}^l q_{a,2}$ . Then we have  $f'_{l,i} = f'_{l,i,1} + f'_{l,i,2}$ . Since  $\|\mathbf{W}_1^{(0)} (x_i - x_k)\| \leq O(\phi \sqrt{\log(1/\phi)})$  with probability at least  $1 - \phi$ ,  $\phi \in (0, 1)$ , by triangle inequality, we can write

$$\begin{aligned} \|f_{l,i}^{(0)} - f_{l,k}^{(0)}\| &= \|f'_{l,i}\| = \|f'_{l,i,1} + f'_{l,i,2}\| \\ &\leq O(\phi^{4/3}L^3 \log(m) \log^{3/4}(1/\phi) + \phi \sqrt{\log(1/\phi)}) + O(\phi L^{5/2} \sqrt{\log(m)} \log^{1/2}(1/\phi)) \\ &\leq O(L^{5/2}\phi \sqrt{\log(m) \log(1/\phi)}), \end{aligned}$$

where the first inequality comes from inequalities 26, and the last inequality holds by the assumption  $\phi \leq O(L^{-9/2} \log^{-3}(m) \log^{-3/4}(1/\phi))$ . Also, by the requirement of  $\phi$ , we have with  $\|f'_{l,i,1}\| \leq O(\phi^{4/3} L^3 \log(m) \log^{3/4}(1/\phi) + \phi \sqrt{\log(1/\phi)}) \leq O(L^{3/2} \phi \log^{1/2}(1/\phi))$  and  $\|f'_{l,i,2}\|_\infty \leq O(\phi L^{5/2} \sqrt{\log(m)/m} \log^{1/2}(1/\phi)) \leq O(L \phi^{2/3} m^{-1/2} \log^{1/2}(1/\phi))$ . Thus (a) holds for layer  $l$ . And by Lemma C.1, we have with probability at least  $1 - \phi$ ,

$$\begin{aligned} \|h_{l,i}^{(0)} - h_{l,k}^{(0)}\| &= \|(\mathbf{D}_{l,k}^{(0)} + \mathbf{D}'_l) (f_{l,k}^{(0)} + f_{l,i}^{(0)}) - \mathbf{D}_{l,k}^{(0)} f_{l,k}^{(0)}\| \\ &\leq \|\mathbf{D}'_l f_{l,k}^{(0)}\| + \|(\mathbf{D}_{l,k}^{(0)} + \mathbf{D}'_l) f_{l,i}^{(0)}\| \leq O(L^{5/2} \phi \sqrt{\log(m) \log(1/\phi)}), \end{aligned}$$

where the second inequality comes from Lemma C.1. Thus, conclusion (c) holds for layer  $l$ .

Finally, by Lemma 8.2 in (Allen-Zhu et al., 2019b) which gives forward perturbation regarding weights, we have with probability at least  $1 - O(\phi)$ ,

$$\begin{aligned} \|f_{l,i} - f_{l,k}\| &\leq \|f_{l,i}^{(0)} - f_{l,k}^{(0)}\| + \|f_{l,i} - f_{l,i}^{(0)}\| + \|f_{l,k} - f_{l,k}^{(0)}\| \\ &\leq O(L^{3/2} \phi \log^{1/2}(1/\phi)) + O(\tau L^{5/2} \sqrt{\log(m)}) \leq O(L^{5/2} \phi \sqrt{\log(m) \log(1/\phi)}), \end{aligned}$$

where the last probability holds by the assumption  $\tau \leq O(\phi^{3/2})$ . Similarly, we have with probability at least  $1 - O(\phi)$ ,

$$\begin{aligned} \|h_{l,i} - h_{l,k}\| &\leq \|h_{l,i}^{(0)} - h_{l,k}^{(0)}\| + \|h_{l,i} - h_{l,i}^{(0)}\| + \|h_{l,k} - h_{l,k}^{(0)}\| \\ &\leq O(L^{3/2} \phi \log^{1/2}(1/\phi)) + O(\tau L^{5/2} \sqrt{\log(m)}) \leq O(L^{5/2} \phi \sqrt{\log(m) \log(1/\phi)}). \end{aligned}$$

□

## C.2. Proofs of Lemma B.2 and Lemma B.3

### Proof of Lemma B.2

*Proof.* By the mean value theorem,  $r$  can be represented as

$$r(h') = r(h) + \nabla r(h)^\top (h' - h) + \frac{1}{2} (h' - h)^\top \nabla^2 r(z) (h' - h), \quad (27)$$

where  $z$  lies in the line segment between  $h'$  and  $h$ .

Since the maximum eigenvalue of the Hessian matrix of  $r$  is bounded by 1, for any output of the neural network  $h$  and  $h'$ , we have

$$r(h') \leq r(h) + \nabla r(h)^\top (h' - h) + \frac{1}{2} \|h' - h\|_2^2 \quad (28)$$

Let  $h' = h - \nabla r(h)$ . We have

$$r_{\min} \leq r(h') \leq r(h) - \frac{1}{2} \|\nabla r(h)\|_2^2. \quad (29)$$

Thus, we get the first inequality of the lemma  $\|\nabla r(h)\|_2^2 \leq 2(r(h) - r_{\min})$ .

By strong convexity, for any  $h$  and  $h'$  in the domain of risk function  $r$ , we have

$$\begin{aligned} r(h') &\geq r(h) + \nabla r(h)^\top (h' - h) + \frac{\rho}{2} \|h' - h\|^2 \\ &\geq r(h) - \frac{1}{2\rho} \|\nabla r(h)\|^2, \end{aligned} \quad (30)$$

where the first inequality comes from strong convexity and the second inequality holds by choosing  $h' = -\frac{\nabla r(h)}{\rho}$  that minimizes the right hand side. Then letting  $h'$  in the left hand side equals to  $h^*$  such that  $r(h^*) = r_{\min}$ , we get the second inequality of the lemma  $\|\nabla r(h)\|_2^2 \geq 2\rho(r(h) - r_{\min})$ .



Also, letting  $h'$  in Eqn. (30) be  $h^*$ , we have

$$r_{\min} = r(h^*) \geq r(h) - \|\nabla r(h)\| \|h^* - h\| + \frac{\rho}{2} \|h^* - h\|^2 \quad (31)$$

By the fact that  $r_{\min} \leq r(h)$ , we have

$$\|h^* - h\| \leq \frac{2}{\rho} \|\nabla r(h)\|. \quad (32)$$

We thus get the third inequality.  $\square$

### Proof of Lemma B.3

*Proof.* Denote the risk of the  $k$ th cell  $k \in [N]$  with respect to the input  $x'_k \in \mathcal{X}_\phi$  for hypothesis  $h \in \mathcal{H}$  as

$$\bar{r}_{\mathbf{I},k}^\circ = \sum_{i \in \mathcal{I}_{\phi,k}} r_{\mathbf{I},i}(h(x'_k)).$$

Recall that  $M_k = \sum_{i \in \mathcal{I}_{\phi,k}} (\mu_i + \lambda_i)$ . By 1-Lipschitz continuity of risk functions and their gradients, we have with probability at least  $1 - O(\phi)$ ,

$$|\bar{r}_{\mathbf{I},k} - \bar{r}_{\mathbf{I},k}^\circ| \leq \sum_{i \in \mathcal{I}_{\phi,k}} (\mu_i + \lambda_i) \|h(x_i) - h(x'_k)\| \leq \tilde{O}(M_k L^{5/2} \phi \log^{1/2}(m)),$$

$$\text{and } \|\nabla_h \bar{r}_{\mathbf{I},k} - \nabla_h \bar{r}_{\mathbf{I},k}^\circ\| \leq \sum_{i \in \mathcal{I}_{\phi,k}} (\mu_i + \lambda_i) \|h(x_i) - h(x'_k)\| \leq \tilde{O}(M_k L^{5/2} \phi \log^{1/2}(m)).$$

Since the eigenvalues of  $\nabla_h^2 \bar{r}_{\mathbf{I},k}^\circ$  is no less than  $M_k \rho$  and  $r_{\text{eff},k}^*$  is the minimum value of  $\sum_{i \in \mathcal{I}_{\phi,k}} r_{\mathbf{I},i}(h)$ , by Lemma B.2, we have

$$\|\nabla_h \bar{r}_{\mathbf{I},k}^\circ\|^2 \geq 2M_k \rho (\bar{r}_{\mathbf{I},k}^\circ - r_{\text{eff},k}^*) \geq 2M_k \rho (\bar{r}_{\mathbf{I},k} - r_{\text{eff},k}^*) - \tilde{O}(M_k^2 \rho L^{5/2} \phi \log^{1/2}(m)).$$

Therefore, we have

$$\begin{aligned} \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i(h(x_i)) \right\|^2 &= \|\nabla_h \bar{r}_{\mathbf{I},k}\|^2 \geq \|\nabla_h \bar{r}_{\mathbf{I},k}^\circ\|^2 - \tilde{O}(M_k^2 L^{5/2} \phi \log^{1/2}(m)) \\ &\geq 2M_k \rho (\bar{r}_{\mathbf{I},k} - r_{\text{eff},k}^*) - \tilde{O}(M_k^2 L^{5/2} \phi \log^{1/2}(m)) \end{aligned}$$

Also, since the eigenvalues of  $\nabla_h^2 \bar{r}_{\mathbf{I},k}^\circ$  is no larger than  $M_k$ , we have

$$\|\nabla_h \bar{r}_{\mathbf{I},k}^\circ\|^2 \leq 2M_k (\bar{r}_{\mathbf{I},k}^\circ - r_{\text{eff},k}) \leq 2M_k (\bar{r}_{\mathbf{I},k} - r_{\text{eff},k}) + \tilde{O}(M_k^2 L^{5/2} \phi \log^{1/2}(m)).$$

Therefore, it holds that

$$\begin{aligned} \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i(h(x_i)) \right\|^2 &= \|\nabla_h \bar{r}_{\mathbf{I},k}\|^2 \leq \|\nabla_h \bar{r}_{\mathbf{I},k}^\circ\|^2 + \tilde{O}(M_k^2 \rho L^{5/2} \phi \log^{1/2}(m)) \\ &\leq 2M_k (\bar{r}_{\mathbf{I},k} - r_{\text{eff},k}) + \tilde{O}(M_k^2 L^{5/2} \phi \log^{1/2}(m)). \end{aligned}$$

Applying Lemma B.2 for  $\bar{r}_{\mathbf{I},k}^\circ$ , we have

$$\|h(x'_k) - y_{\text{eff},k}\|^2 \leq \frac{4}{M_k^2 \rho^2} \|\nabla_h \bar{r}_{\mathbf{I},k}^\circ\|^2 \leq \frac{1}{M_k \rho^2} O\left(\bar{r}_{\mathbf{I},k} - r_{\text{eff},k} + \tilde{O}(M_k L^{5/2} \phi \log^{1/2}(m))\right).$$

By applying Lemma B.1 to  $h(x_i)$ , we have

$$\begin{aligned} \|h(x_i) - y_{\text{eff},k}\|^2 &\leq 2\|h(x_i) - h(x'_k)\|^2 + 2\|h(x'_k) - y_{\text{eff},k}\|^2 \\ &\leq \frac{1}{M_k \rho^2} O\left(\bar{r}_{\mathbf{I},k} - r_{\text{eff},k} + \tilde{O}(M_k L^{5/2} \phi \log^{1/2}(m))\right) + \tilde{O}(L^5 \phi^2 \log(m)). \end{aligned}$$

Taking weighted summation in the cell  $\mathcal{I}_{\phi,k}$ , since  $\sum_{i \in \mathcal{I}_{\phi,k}} (\mu_i + \lambda_i) = M_k$ , we have

$$\sum_{i \in \mathcal{I}_{\phi,k}} (\mu_i + \lambda_i) \|h(x_i) - y_{\text{eff},k}\|^2 \leq \frac{1}{\rho^2} O\left(\bar{r}_{\mathbf{I},k} - r_{\text{eff},k} + \tilde{O}(M_k L^{5/2} \phi \log^{1/2}(m))\right).$$

$\square$

### C.3. Proof of Lemma B.4

**Lemma C.2.** Suppose that  $m \geq \Omega(N^2 d^2 \phi^{-1})$ . For any  $u_i : \|u_i\| \leq \mu_i + \lambda_i, i \in [n']$  and  $v_j \sim \mathcal{N}(0, (1/d)\mathbf{I})$ ,  $w_j \sim \mathcal{N}(0, (2/m)\mathbf{I})$ , with Assumption 1 satisfied, with probability at least  $1 - O(\phi)$ , we have

$$\sum_{j=1}^m \left\| \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi,k}} \langle u_i, v_j \rangle \sigma' \left( \langle w_j, h_{L-1}^{(0)}(x_i) \rangle \right) h_{L-1}^{(0)}(x_i) \right\|^2 \geq \Omega \left( \frac{\alpha m \phi}{Nd} \right) \left( \sum_{k=1}^N \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\|^2 - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right).$$

Lemma C.2 is proved in Section D.2

**Lemma C.3** (Lemma 8.7, Lemma 8.2c in (Allen-Zhu et al., 2019b)). For any  $\mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ , with probability at least  $1 - \exp(-O(m\tau^{2/3}L))$ ,

$$\left\| \mathbf{V} \mathbf{D}_{i,L} - \mathbf{V} \mathbf{D}_{i,L}^{(0)} \right\|_2 \leq O \left( \tau^{1/3} L^2 \sqrt{m \log(m)/d} \right),$$

and  $\forall l \in [L]$ ,

$$\left\| h_{i,l} - h_{i,l}^{(0)} \right\| \leq O \left( \tau L^{5/2} \sqrt{\log(m)} \right).$$

### Proof of Lemma B.4

*Proof.* Denote  $u_i = u_i(h_{\mathbf{W}}(x_i))$ . The gradient of the empirical informed risk can be expressed as

$$\nabla_{W_l} \hat{R}_l(\mathbf{W}) = \sum_{i=1}^{n'} (u_i \mathbf{V} \mathbf{D}_{L,i} \mathbf{W}_{L,i} \cdots \mathbf{W}_{l+1,i} \mathbf{D}_l)^\top h_{l-1,i}. \quad (33)$$

Let  $\mathbf{G} = \nabla_{W_L} \hat{R}_l(\mathbf{W}^{(0)}) = \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi,k}} \left( u_i \mathbf{V} \mathbf{D}_{L,i}^{(0)} \right)^\top h_{L-1,i}^{(0),\top}$ . By Lemma C.2, with probability at least  $1 - O(\phi)$ , we have

$$\begin{aligned} \|\mathbf{G}\|_F^2 &\geq \Omega \left( \frac{\alpha m \phi}{Nd} \right) \left( \sum_{k=1}^N \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\|^2 - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right) \\ &\geq \Omega \left( \frac{\alpha m \phi \rho}{dN} \right) \left( \sum_{k=1}^N M_k (\bar{r}_{1,k} - r_{\text{eff},k}) - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right) \\ &\geq \Omega \left( \frac{\alpha m \phi \rho \bar{\lambda}}{dN^2} \right) \left( \hat{R}_l - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right), \end{aligned} \quad (34)$$

where  $\hat{R}_{\text{eff}} = \sum_{k=1}^N r_{\text{eff},k}$  and the second inequality comes from Lemma B.3 and the last inequality holds because  $\sum_{k=1}^N M_k (\bar{r}_{1,k} - r_{\text{eff},k}) \geq \bar{M} \sum_{k=1}^N (\bar{r}_{1,k} - r_{\text{eff},k})$  with  $\bar{M} = \min_k M_k$ ,  $M_k = \sum_{i \in \mathcal{I}_{\phi,k}} (\mu_i + \lambda_i)$ , and  $N\bar{M} = \bar{\lambda}$ .

Here, we need to discuss more about  $\bar{\lambda}$  which is different for different objectives. Denote  $\bar{p}_z = \min_k |S_z \cap S_{\mathcal{I}_{\phi,k}}|$ ,  $\bar{p}_g = \min_k |S_g \cap S_{\mathcal{I}_{\phi,k}}|$  and  $\bar{p}_{g/z} = \min_k |(S_g \setminus S_z) \cap S_{\mathcal{I}_{\phi,k}}|$ . When  $\lambda \neq 1$  or  $\lambda \neq 0$ ,  $\bar{\lambda} = N \min \left\{ \frac{(1-\lambda)\bar{p}_z}{n_z} + \frac{\lambda\bar{p}_g}{n_g}, \frac{\lambda\bar{p}_g}{n_g} \right\} \geq \Omega(\lambda)$  for objective (3),  $\bar{\lambda} = N \min \left\{ \frac{(1-\lambda)(1-\beta)\bar{p}_z}{n_z} + \frac{(1-\lambda)\beta\bar{p}_g}{n_g}, \frac{\lambda\bar{p}_g}{n_g} \right\} \geq \Omega(\min(1-\lambda, \lambda))$  for objective (5)<sup>1</sup>. Beside, the cases when  $\lambda = 0$  or  $\lambda = 1$  mean the corresponding datasets are empty (e.g. when  $\lambda = 1$  in (5),  $S_z = \emptyset$  and  $S'_g = \emptyset$ ), so we have  $\bar{\lambda} = 1$ . In conclusion, we have  $\bar{\lambda} = \Omega(\min(1-\lambda, \lambda)\mathbf{1}(\lambda \in (0, 1)) + \mathbf{1}(\lambda \in \{0, 1\}))$  for two objectives.

<sup>1</sup>Here,  $S_{\mathcal{I}_{\phi,k}}$  is the set of samples with their indices in  $\mathcal{I}_{\phi,k}$ . Thus, there exists a constant  $C$  such that  $n_z \leq CN\bar{p}_z$ ,  $n_g \leq CN\bar{p}_g$ ,  $(n_g - n_z) \leq CN\bar{p}_{g/z}$  where  $C$  relies on the input distribution.

Next we bound the difference of  $\|\mathbf{G}\|$  and  $\|\nabla_{W_L} \hat{R}_I(\mathbf{W})\|$  with  $\mathbf{W} \in \mathcal{B}(\mathbf{W}^{(0)}, \tau)$ . By definition, we have

$$\begin{aligned} \left\| \mathbf{G} - \nabla_{W_L} \hat{R}_I(\mathbf{W}) \right\|_F &= \left\| \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi,k}} \left( u_i \mathbf{V} \mathbf{D}_{L,i}^{(0)} \right)^\top h_{L-1,i}^{(0),\top} - \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi,k}} \left( u_i \mathbf{V} \mathbf{D}_{L,i} \right)^\top h_{L-1,i}^\top \right\|_F \\ &\leq \left\| \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi,k}} \left( u_i \mathbf{V} \mathbf{D}_{L,i}^{(0)} - u_i \mathbf{V} \mathbf{D}_{L,i} \right)^\top h_{L-1,i}^{(0),\top} \right\|_F + \left\| \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi,k}} \left( u_i \mathbf{V} \mathbf{D}_{L,i} \right)^\top \left( h_{L-1,i}^{(0)} - h_{L-1,i} \right)^\top \right\|_F \end{aligned} \quad (35)$$

(a) (b)

For the term (a) in the above inequality, denoting  $h_{L-1,k} = h_{L-1}(x'_k)$  and letting  $(a)_k$  be the  $k$ th item in the summation, we have

$$\begin{aligned} (a)_k &\leq \left\| \sum_{i \in \mathcal{I}_{\phi,k}} \left( \mathbf{D}_{L,k}^{(0)} - \mathbf{D}_{L,k} \right) \mathbf{V}^\top u_i^\top h_{L-1,k}^{(0),\top} \right\|_F \\ &\quad + \left\| \sum_{i \in \mathcal{I}_{\phi,k}} \left( \mathbf{D}_{L,i}^{(0)} - \mathbf{D}_{L,i} \right) \mathbf{V}^\top u_i^\top h_{L-1,i}^{(0),\top} - \sum_{i \in \mathcal{I}_{\phi,k}} \left( \mathbf{D}_{L,k}^{(0)} - \mathbf{D}_{L,k} \right) \mathbf{V}^\top u_i^\top h_{L-1,k}^{(0),\top} \right\|_F \\ &\leq O\left(\tau^{1/3} L^2 \sqrt{m \log(m)/d}\right) \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| + \sum_{i \in \mathcal{I}_{\phi,k}} \left\| \mathbf{V} \left( \mathbf{D}_{L,i}^{(0)} - \mathbf{D}_{L,i} \right) \right\|_2 \|u_i\| \left\| h_{L-1,i}^{(0)} \right\| \\ &\quad + \sum_{i \in \mathcal{I}_{\phi,k}} \left\| \mathbf{V} \left( \mathbf{D}_{L,k}^{(0)} - \mathbf{D}_{L,k} \right) \right\|_2 \|u_i\| \left\| h_{L-1,k}^{(0)} \right\| \\ &\leq O\left(\tau^{1/3} L^2 \sqrt{m \log(m)/d}\right) \left( \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| + M_k \right), \end{aligned}$$

where the second inequality comes from Lemma C.3 and Cauchy-Schwartz inequality, and the last inequality comes from Lemma C.3 and  $\sum_{i \in \mathcal{I}_{\phi,k}} \|u_i\| \leq M_k$  and Lemma B.6 such that  $\left\| h_{L-1,i}^{(0)} \right\| \leq \tilde{O}(1)$  with probability at least  $1 - O(\phi)$ .

For the term (b), it holds that

$$\begin{aligned} (b)_k &\leq \left\| \sum_{i \in \mathcal{I}_{\phi,k}} \left( \mathbf{V} \mathbf{D}_{L,k} \right)^\top u_i^\top \left( h_{L-1,k}^{(0)} - h_{L-1,k} \right)^\top \right\|_F \\ &\quad + \left\| \sum_{i \in \mathcal{I}_{\phi,k}} \left( \mathbf{V} \mathbf{D}_{L,i} \right)^\top u_i^\top \left( h_{L-1,i}^{(0)} - h_{L-1,i} \right)^\top - \sum_{i \in \mathcal{I}_{\phi,k}} \left( \mathbf{V} \mathbf{D}_{L,k} \right)^\top u_i^\top \left( h_{L-1,k}^{(0)} - h_{L-1,k} \right)^\top \right\|_F \\ &\leq O\left(\tau L^{5/2} \sqrt{m \log(m)/d}\right) \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| + \sum_{i \in \mathcal{I}_{\phi,k}} \left\| \mathbf{V} \mathbf{D}_{L,i} \right\|_2 \|u_i\| \left\| h_{L-1,i}^{(0)} - h_{L-1,i} \right\| \\ &\quad + \sum_{i \in \mathcal{I}_{\phi,k}} \left\| \mathbf{V} \mathbf{D}_{L,k} \right\|_2 \|u_i\| \left\| h_{L-1,k}^{(0)} - h_{L-1,k} \right\| \\ &\leq O\left(\tau L^{5/2} \sqrt{m \log(m)/d}\right) \left( \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| + M_k \right), \end{aligned}$$

where the second inequality comes from Lemma C.3 and Cauchy-Schwartz inequality, and the last inequality comes from Lemma C.3 and  $\sum_{i \in \mathcal{I}_{\phi,k}} \|u_i\| \leq M_k$  and Lemma B.6 such that  $\left\| h_{L-1,i}^{(0)} \right\| \leq \tilde{O}(1)$  with probability at least  $1 - O(\phi)$ .

Therefore, we can bound Eqn. (35) as

$$\begin{aligned}
 & \left\| \mathbf{G} - \nabla_{W_L} \hat{R}_I(\mathbf{W}) \right\|_F \leq O\left(\tau^{1/3} L^{5/2} \sqrt{m \log(m)/d}\right) \left( \sum_{k=1}^N \left\| \sum_{i \in \mathcal{I}_{\phi, k}} u_i \right\| + 1 \right) \\
 & \leq O\left(N^{1/2} \tau^{1/3} L^{5/2} \sqrt{m \log(m)/d}\right) \left( \sqrt{\left(\hat{R}_I - \hat{R}_{\text{eff}} + \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right)} + 1/\sqrt{N} \right) \\
 & \leq O\left(N^{1/2} \tau^{1/3} L^{5/2} \sqrt{m \log(m)/d}\right) \left( \sqrt{\left(\hat{R}_I - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right)} \right),
 \end{aligned} \tag{36}$$

where the second inequality holds by Lemma B.3 and the last inequality holds because  $\phi$  is small enough such that  $\hat{R}_I(\mathbf{W}) - \hat{R}_{\text{eff}} + \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \leq 2(\hat{R}_I(\mathbf{W}) - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)))$ .

Combining Eqn. (36) with Eqn. (34), we have

$$\begin{aligned}
 & \left\| \nabla_{W_L} \hat{R}_I(\mathbf{W}) \right\|_F \geq \|\mathbf{G}\|_F - \left\| \mathbf{G} - \nabla_{W_L} \hat{R}_I(\mathbf{W}) \right\|_F \\
 & \geq \Omega\left(\sqrt{\frac{\alpha m \phi \rho \bar{\lambda}}{d N^2}} - O\left(\tau^{1/3} N^{1/2} L^{5/2} \sqrt{m \log(m)/d}\right)\right) \left(\hat{R}_I - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right) \\
 & \geq \Omega\left(\sqrt{\frac{\alpha m \phi \rho \bar{\lambda}}{d N^2}} \left(\hat{R}_I - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right)\right),
 \end{aligned}$$

where the last inequality holds by the choice of  $m \geq \Omega\left(N^{11} L^{15} d \phi^{-4} \rho^{-4} \bar{\lambda}^{-4} \alpha^{-4} \log^3(m)\right)$  and the weight update range in the proof of Theorem 4.1 such that  $\tau^{1/3} = O\left(N^{-3/2} \phi^{1/2} \rho^{1/2} \bar{\lambda}^{1/2} \alpha^{1/2} L^{-5/2} \log^{-1/2}(m)\right)$ .

$$\left\| \nabla_{\mathbf{W}} \hat{R}_I(\mathbf{W}) \right\|_F^2 \geq \left\| \nabla_{W_L} \hat{R}_I(\mathbf{W}) \right\|_F^2 \geq \Omega\left(\frac{\alpha m \phi \rho \bar{\lambda}}{d N^2}\right) \left(\hat{R}_I - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2} \phi \log^{1/2}(m))\right).$$

□

#### C.4. Proof of Lemma B.5

##### Proof of Lemma B.5

*Proof.* Since the maximum eigenvalue of the second order derivation of the informed risk function  $r_{I,i}(h_{\mathbf{W},i})$  with respect to  $h$  is less than  $\mu_i + \lambda_i$ , we have

$$r_{I,i}(h_{\mathbf{W}',i}) - r_{I,i}(h_{\mathbf{W},i}) \leq u_i(h_{\mathbf{W},i})^\top (h_{\mathbf{W}',i} - h_{\mathbf{W},i}) + O\left(\frac{\mu_i + \lambda_i}{2} \|h_{\mathbf{W}',i} - h_{\mathbf{W},i}\|^2\right). \tag{37}$$

Then denote  $\widehat{\mathbf{W}} = \mathbf{W}' - \mathbf{W}$ . We have

$$\begin{aligned}
 & \sum_{i \in \mathcal{I}_{\phi, k}} r_{I,i}(h_{\mathbf{W}',i}) - r_{I,i}(h_{\mathbf{W},i}) - \left\langle \nabla_{\mathbf{W}} r_{I,i}(h_{\mathbf{W},i}), \widehat{\mathbf{W}} \right\rangle \\
 & \leq \sum_{i \in \mathcal{I}_{\phi, k}} u_i(h_{\mathbf{W},i})^\top \left( h_{\mathbf{W}',i} - h_{\mathbf{W},i} - \left\langle \nabla_{\mathbf{W}} h_{\mathbf{W},i}, \widehat{\mathbf{W}} \right\rangle \right) + O\left(\sum_{i \in \mathcal{I}_{\phi, k}} \frac{\mu_i + \lambda_i}{2} \|h_{\mathbf{W}',i} - h_{\mathbf{W},i}\|^2\right) \\
 & \leq \left( \sum_{i \in \mathcal{I}_{\phi, k}} u_i(h_{\mathbf{W},i}) \right)^\top \left( h_{\mathbf{W}',k} - h_{\mathbf{W},k} - \left\langle \nabla_{\mathbf{W}} h_{\mathbf{W},k}, \widehat{\mathbf{W}} \right\rangle \right) + O\left(\sum_{i \in \mathcal{I}_{\phi, k}} \frac{\mu_i + \lambda_i}{2} \|h_{\mathbf{W}',i} - h_{\mathbf{W},i}\|^2\right) \\
 & \quad + \sum_{i \in \mathcal{I}_{\phi, k}} \|u_i(h_{\mathbf{W},i})\| \left\| \left[ \left( h_{\mathbf{W}',i} - h_{\mathbf{W},i} - \left\langle \nabla_{\mathbf{W}} h_{\mathbf{W},i}, \widehat{\mathbf{W}} \right\rangle \right) - \left( h_{\mathbf{W}',k} - h_{\mathbf{W},k} - \left\langle \nabla_{\mathbf{W}} h_{\mathbf{W},k}, \widehat{\mathbf{W}} \right\rangle \right) \right] \right\|,
 \end{aligned} \tag{38}$$

where Cauchy-Schwartz inequality is used in the last inequality. By Theorem 4 in (Allen-Zhu et al., 2019b), we have with probability at least  $1 - \exp(-\Omega(m\tau^{2/3}L))$ ,

$$\begin{aligned} & \left\| h_{\mathbf{W}',i} - h_{\mathbf{W},i} - \left\langle \nabla_{\mathbf{W}} h_{\mathbf{W},i}, \widehat{\mathbf{W}} \right\rangle \right\| \\ & \leq O\left(\tau^{1/3}L^{5/2}\sqrt{m\log(m)}d^{-1/2}\right) \left\| \widehat{\mathbf{W}} \right\| + O\left(L^2\sqrt{m/d}\left\| \widehat{\mathbf{W}} \right\|^2\right) \end{aligned}$$

By Claim 11.2 in (Allen-Zhu et al., 2019b), we have

$$\|h_{\mathbf{W}',i} - h_{\mathbf{W},i}\| \leq O(L\sqrt{m/d}) \left\| \widehat{\mathbf{W}} \right\|.$$

Thus since  $\|u_i(h_{\mathbf{W},i})\| \leq O(\mu_i + \lambda_i)$ , we have

$$\begin{aligned} & \sum_{i \in \mathcal{I}_{\phi,k}} r_{1,i}(h_{\mathbf{W}',i}) - r_{1,i}(h_{\mathbf{W},i}) - \left\langle \nabla_{\mathbf{W}} r_{1,i}(h_{\mathbf{W},i}), \widehat{\mathbf{W}} \right\rangle \\ & \leq O\left(\sqrt{M_k(\bar{r}_{1,k} - r_{\text{eff},k} + \tilde{O}(L^{5/2}\phi\log^{1/2}(m)))} + M_k\right) O\left(\tau^{1/3}L^{5/2}\sqrt{m\log(m)}d^{-1/2}\right) \left\| \widehat{\mathbf{W}} \right\| \\ & \quad + O(M_kL^2m/d) \left\| \widehat{\mathbf{W}} \right\|^2. \end{aligned} \quad (39)$$

where the inequality comes from Lemma B.3. Taking summation over  $i \in [N]$ , we have

$$\begin{aligned} \hat{R}_1(\mathbf{W}') - \hat{R}_1(\mathbf{W}) & \leq \left\langle \nabla_{\mathbf{W}} \hat{R}_1(\mathbf{W}), \mathbf{W}' - \mathbf{W} \right\rangle + O(L^2m/d) \left\| \widehat{\mathbf{W}} \right\|^2 \\ & \quad + \left( \sqrt{\left(\hat{R}_1(\mathbf{W}) - \hat{R}_{\text{eff}} + \tilde{O}(L^{5/2}\phi\log^{1/2}(m))\right)} + 1/\sqrt{N} \right) O\left(N^{1/2}\tau^{1/3}L^{5/2}\sqrt{m\log(m)}d^{-1/2}\right) \left\| \widehat{\mathbf{W}} \right\| \\ & \leq \left\langle \nabla_{\mathbf{W}} \hat{R}_1(\mathbf{W}), \mathbf{W}' - \mathbf{W} \right\rangle + O(L^2m/d) \left\| \widehat{\mathbf{W}} \right\|^2 \\ & \quad + \left( \sqrt{\left(\hat{R}_1(\mathbf{W}) - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2}\phi\log^{1/2}(m))\right)} \right) O\left(N^{1/2}\tau^{1/3}L^{5/2}\sqrt{m\log(m)}d^{-1/2}\right) \left\| \widehat{\mathbf{W}} \right\|, \end{aligned}$$

where the second inequality comes from the choice of  $\phi$  such that  $\hat{R}_1(\mathbf{W}) - \hat{R}_{\text{eff}} + \tilde{O}(L^{5/2}\phi\log^{1/2}(m)) \leq 2(\hat{R}_1(\mathbf{W}) - \hat{R}_{\text{eff}} - \tilde{O}(L^{5/2}\phi\log^{1/2}(m)))$  and  $1/\sqrt{N} \leq \sqrt{\phi} \leq \sqrt{\left(\hat{R}_1 - \hat{R}_{\text{eff}} + \tilde{O}(L^{5/2}\phi\log^{1/2}(m))\right)}$ .  $\square$

### C.5. Proof of Lemma B.6

*Proof.* By Lemma 7.1 in (Allen-Zhu et al., 2019b), with probability at least  $1 - O(NL)\exp(-\Omega(m/L))$ , we have  $\forall k \in [N], \left\| h_{k,L}^{(0)} \right\| \leq 2$ . Thus by Lemma B.1, we have with probability at least  $1 - O(\phi)$ ,

$$\forall k \in [N], \forall i \in \mathcal{I}_{\phi,k}, \left\| h_{i,L}^{(0)} \right\| \leq 2 + \tilde{O}(L^{5/2}\phi\log^{1/2}(m)).$$

Then since each entry of  $\mathbf{V}$  satisfies  $\mathcal{N}(0, \frac{1}{d})$  and  $O(NL)\exp(-\Omega(m/L)) \leq O(\phi)$ , we have with probability at least  $1 - O(\phi)$ ,

$$\|h_{\mathbf{W}^{(0)},i}\| = \left\| \mathbf{V}h_{i,L}^{(0)} \right\| \leq 2 \left\| h_{i,L}^{(0)} \right\| \sqrt{\log(1/\phi)} = O\left(\sqrt{\log(1/\phi)}\right).$$

Let  $r_{1,i}(y_{\text{eff},k(x_i)}) = \mu_i r(y_{\text{eff},k(x_i)}, y_i) + \lambda_i r_K(y_{\text{eff},k(x_i)}, g(x_i))$ . Thus with probability at least  $1 - O(\phi)$ , by 1-Lipschitz continuity of risk functions, we have

$$\begin{aligned} r_{1,i}^{(0)} - r_{1,i}(y_{\text{eff},k(x_i)}) & \leq (\mu_i + \lambda_i) \left\| h_{\mathbf{W}^{(0)},i} - y_{\text{eff},k(x_i)} \right\| \\ & \leq (\mu_i + \lambda_i) (\left\| h_{\mathbf{W}^{(0)},i} \right\| + \left\| y_{\text{eff},k(x_i)} \right\|) \\ & \leq O\left((\mu_i + \lambda_i)\log^{1/2}(1/\phi)\right). \end{aligned}$$

Taking summation for  $i \in [n']$ , we have

$$\hat{R}_I(\mathbf{W}^{(0)}) - \hat{R}_{\text{eff}} \leq O\left(\log^{1/2}(1/\phi)\right).$$

□

### C.6. Proof of Lemma B.8

*Proof. Proof of (a):* Denote  $\mathcal{U}'_g = \mathcal{U}_\phi(S_z) = \{k \in [N] \mid \exists x \in S_z, x \in \mathcal{C}_{\phi,k}\}$  as the index collection of smooth sets that contain at least one labeled sample, and  $\mathcal{U}''_g = [N] \setminus \mathcal{U}'_g$  as the index collection of smooth sets that only contain knowledge-supervised samples. Denote  $h_{K,i}^* = h_K^*(x_i)$  for notation simplicity, and recall that  $x'_k \in \mathcal{X}_\phi$  is the representative input of the smooth set  $k$ , so we have

$$\begin{aligned} & \frac{1}{n''_g} \sum_{S'_g} r_K(h_{K,i}^*, g_i) = \frac{1}{n''_g} \sum_{k \in \mathcal{U}'_g} \sum_{\mathcal{I}_{\phi,k}} r_K(h_{K,i}^*, g_i) \\ & \geq \frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} \sum_{\mathcal{I}_{\phi,k}} r_K(h_K^*(x'_k), g_i) - \tilde{O}\left(L^{5/2}\phi \log^{1/2}(m)\right) \\ & \geq \frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} \sum_{\mathcal{I}_{\phi,k}} r_K(y_{\text{eff},k}, g_i) + \frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} \left\langle \nabla h \sum_{\mathcal{I}_{\phi,k}} r_K(y_{\text{eff},k}, g_i), h_K^*(x'_k) - y_{\text{eff},k} \right\rangle \\ & \quad + \frac{\rho}{2n''_g} \sum_{k \in \mathcal{U}''_g} |\mathcal{I}_{\phi,k}| \|h_K^*(x'_k) - y_{\text{eff},k}\|^2 - \tilde{O}\left(L^{5/2}\phi \log^{1/2}(m)\right), \end{aligned} \tag{40}$$

where the first inequality holds by Lemma B.1 and Lipschitz continuity of the risk function, and the second inequality holds by the strongly convexity of  $\sum_{\mathcal{I}_{\phi,k}} r_K(h, g_i)$  with respect to  $h$ . By subtracting  $\frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} \sum_{\mathcal{I}_{\phi,k}} r_K(y_{\text{eff},k}, g_i)$  from both sides of (40), we have

$$\begin{aligned} & \frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} |\mathcal{I}_{\phi,k}| \|h_K^*(x'_k) - y_{\text{eff},k}\|^2 \\ & \leq \frac{2}{\rho n''_g} \sum_{k \in \mathcal{U}''_g} \left\| \nabla h \sum_{\mathcal{I}_{\phi,k}} r_K(y_{\text{eff},k}, g_i) \right\| \|h_K^*(x'_k) - y_{\text{eff},k}\| + \tilde{O}\left(L^{5/2}\phi \log^{1/2}(m)\right) \\ & \leq O\left(\frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} \left\| \nabla h \sum_{\mathcal{I}_{\phi,k}} r_K(y_{\text{eff},k}, g_i) \right\|\right) + \tilde{O}\left(L^{5/2}\phi \log^{1/2}(m)\right) \\ & \leq \sum_{k \in \mathcal{U}''_g} \left[ \sqrt{\frac{2|\mathcal{I}_{\phi,k}|}{n''_g} (r_K(y_{\text{eff},k}, g_i) - r_K(y_{\text{eff},k}, g_i))} + \frac{|\mathcal{I}_{\phi,k}|}{n''_g} \tilde{O}\left(L^{5/4}\phi^{1/2} \log^{1/4}(m)\right) \right] + \tilde{O}\left(L^{5/2}\phi \log^{1/2}(m)\right) \\ & = \tilde{O}\left(L^{5/4}\phi^{1/2} \log^{1/4}(m)\right) + \tilde{O}\left(L^{5/2}\phi \log^{1/2}(m)\right) \leq \tilde{O}\left(L^{5/4}\phi^{1/2} \log^{1/4}(m)\right), \end{aligned} \tag{41}$$

where the first inequality holds since  $\{h_{K,i}^*, i \in S''_g\}$  minimizes  $\frac{1}{n''_g} \sum_{S''_g} r_K(h(x_i), g_i)$ , the second inequality holds since  $\|h_K^*(x'_k) - y_{\text{eff},k}\| \leq \|h_K^*(x'_k)\| + \|y_{\text{eff},k}\| \leq \tilde{O}(1)$  by Lemma B.6 and Lemma 8.2(c) in (Allen-Zhu et al., 2019b), and the third inequality holds by applying Lemma B.3 for  $\sum_{\mathcal{I}_{\phi,k}} r_K(y_{\text{eff},k}, g_i)$  with  $r_{\text{eff},k} = \sum_{\mathcal{I}_{\phi,k}} r_K(y_{\text{eff},k}, g_i)$ . Therefore, by Lemma B.1, we have

$$\frac{1}{n''_g} \sum_{S''_g} \|h_{K,i}^* - y_{\text{eff},k(x_i)}\|^2 \leq \frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} |\mathcal{I}_{\phi,k}| \|h_K^*(x'_k) - y_{\text{eff},k}\|^2 + \tilde{O}\left(L^{5/2}\phi \log^{1/2}(m)\right) \leq \tilde{O}\left(L^{5/4}\phi^{1/2} \log^{1/4}(m)\right). \tag{42}$$

Similarly, denote  $r_{R,\beta}(h(x_i)) = \frac{1-\beta}{n_z} r(h(x_i), y_i) \mathbb{1}(x_i \in S_z) + \frac{\beta}{n_g} r_K(h(x_i), g(x_i)) \mathbb{1}(x_i \in S'_g)$ . We have

$$\begin{aligned}
 & \sum_{S_z \cup S'_g} r_{R,\beta}(h_{R,\beta}^*(x_i)) \geq \sum_{k \in \mathcal{U}'_g} \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(h_{R,\beta}^*(x'_k)) - \tilde{O}\left(L^{5/2} \phi \log^{1/2}(m)\right) \\
 & \geq \sum_{k \in \mathcal{U}'_g} \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(y_{\text{eff},k}) - \sum_{k \in \mathcal{U}'_g} \left\| \nabla_h \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(y_{\text{eff},k}) \right\| \left\| h_{R,\beta}^*(x'_k) - y_{\text{eff},k} \right\| \\
 & \quad + \frac{\rho}{2} \sum_{k \in \mathcal{U}'_g} M_k \left\| h_{R,\beta}^*(x'_k) - y_{\text{eff},k} \right\|^2 - \tilde{O}\left(L^{5/2} \phi \log^{1/2}(m)\right),
 \end{aligned} \tag{43}$$

where  $M_k = \sum_{i \in \mathcal{I}_{\phi,k}} \left[ \frac{1-\beta}{|S_z|} \mathbb{1}(x_i \in S_z) + \frac{\beta}{|S'_g|} \mathbb{1}(x_i \in S'_g) \right]$ , the first inequality holds by Lemma B.1 and Lipschitz continuity of the risk function, and the second inequality holds by the strongly convexity of  $r_{R,\beta}(h_{R,\beta}^*)$  with respect to  $h_{R,\beta}^*$ . Then, subtracting  $\sum_{k \in \mathcal{U}'_g} \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(y_{\text{eff},k})$  from both sides of (43), similarly as Eqn. (41), it holds that

$$\begin{aligned}
 & \sum_{k \in \mathcal{U}'_g} M_k \left\| h_{R,\beta}^*(x'_k) - y_{\text{eff},k} \right\|^2 \\
 & \leq \frac{2}{\rho} \sum_{k \in \mathcal{U}'_g} \left\| \nabla_h \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(y_{\text{eff},k}) \right\| \left\| h_{R,\beta}^*(x'_k) - y_{\text{eff},k} \right\| + \tilde{O}\left(L^{5/2} \phi \log^{1/4}(m)\right) \\
 & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right),
 \end{aligned} \tag{44}$$

where the first inequality holds because  $\sum_{S_z \cup S'_g} r_{R,\beta}(h_{R,\beta}^*(x_i)) - \sum_{k \in \mathcal{U}'_g} \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(y_{\text{eff},k}) \leq 0$ , and the second inequality holds since  $\left\| h_{R,\beta}^*(x'_k) - y_{\text{eff},k} \right\| \leq \tilde{O}(1)$  by Lemma B.6 and Lemma 8.2(c) in (Allen-Zhu et al., 2019b) and then applying Lemma B.3. Therefore, by Lemma B.1, we have

$$\begin{aligned}
 & \frac{1-\beta}{n_z} \sum_{S_z} \left\| h_{R,\beta}^*(x_i) - y_{\text{eff},k(x_i)} \right\|^2 + \frac{\beta}{n'_g} \sum_{S'_g} \left\| h_{R,\beta}^*(x_i) - y_{\text{eff},k(x_i)} \right\|^2 \leq O\left(\sum_{k \in \mathcal{U}'_g} M_k \left\| h_{R,\beta}^*(x'_k) - y_{\text{eff},k} \right\|^2\right) \\
 & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right).
 \end{aligned}$$

**Proof of (b):** Replacing  $h_{R,\beta}^*$  in Eqn. (40) with  $\bar{h}_{\mathbb{K}}^*$  and applying the second and third inequality in Eqn. (41), we have

$$\begin{aligned}
 & \frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} |\mathcal{I}_{\phi,k}| \left\| \bar{h}_{\mathbb{K}}^*(x'_k) - y_{\text{eff},k} \right\|^2 \\
 & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + \frac{1}{n''_g} \sum_{S''_g} r_{\mathbb{K}}(\bar{h}_{\mathbb{K}}^*(x_i), g_i) - \frac{1}{n''_g} \sum_{k \in \mathcal{U}''_g} \sum_{\mathcal{I}_{\phi,k}} r_{\mathbb{K}}(y_{\text{eff},k}, g_i) \\
 & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + \mathbb{E}[r_{\mathbb{K}}(\bar{h}_{\mathbb{K}}^*(x), g(x))] - \mathbb{E}[r_{\mathbb{K}}(y_{\text{eff},k(x)}, g(x))] + O\left(\sqrt{\frac{\log(1/\delta)}{n''_g}}\right) \\
 & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n''_g}}\right),
 \end{aligned}$$

where the second inequality follows from McDiarmid's inequality and the last inequality is because  $\bar{h}_{\mathbb{K}}^*(x)$  minimizes  $\mathbb{E}[r_{\mathbb{K}}(h, g(x))]$ . Therefore, by Lemma B.1 and with the same reason as (42), we have

$$\frac{1}{n''_g} \sum_{S''_g} \left\| \bar{h}_{\mathbb{K},i}^* - y_{\text{eff},k(x_i)} \right\|^2 \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n''_g}}\right).$$

Similarly, replacing  $h_{R,\beta}$  in Eqns. (43) with  $h_{R,\beta}^*$  and with the same reason as the second inequality in (44), we have

$$\begin{aligned} & \sum_{k \in \mathcal{U}'_g} M_k \|\bar{h}_{R,\beta}^*(x'_k) - y_{\text{eff},k}\|^2 \\ & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + \sum_{S_z \cup S'_g} r_{R,\beta}(\bar{h}_{R,\beta}^*(x_i)) - \sum_{k \in \mathcal{U}'_g} \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(y_{\text{eff},k}). \end{aligned} \quad (45)$$

Continuing with (45) and by McDiarmid's inequality, we have

$$\begin{aligned} & \sum_{k \in \mathcal{U}'_g} M_k \|\bar{h}_{R,\beta}^*(x'_k) - y_{\text{eff},k}\|^2 \\ & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + \mathbb{E} \left[ \sum_{S_z \cup S'_g} r_{R,\beta}(\bar{h}_{R,\beta}^*(x_i)) \right] - \mathbb{E} \left[ \sum_{k \in \mathcal{U}'_g} \sum_{\mathcal{I}_{\phi,k}} r_{R,\beta}(y_{\text{eff},k}) \right] + O\left(\sqrt{\frac{\log(1/\delta)}{n_z}}\right) \\ & \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n_z}}\right), \end{aligned}$$

where the second inequality is because  $\bar{h}_{R,\beta}^*(x)$  minimizes  $\mathbb{E} \left[ \sum_{S_z \cup S'_g} r_{R,\beta}(\bar{h}_{R,\beta}^*(x_i)) \right]$ . And thus by Lemma B.1, we have

$$\frac{1-\beta}{n_z} \sum_{S_z} \|\bar{h}_{R,\beta}^*(x_i) - y_{\text{eff},k(x_i)}\|^2 + \frac{\beta}{n'_g} \sum_{S'_g} \|\bar{h}_{R,\beta}^*(x_i) - y_{\text{eff},k(x_i)}\|^2 \leq \tilde{O}\left(L^{5/4} \phi^{1/2} \log^{1/4}(m)\right) + O\left(\sqrt{\frac{\log(1/\delta)}{n_z}}\right).$$

□

## D. Proof of Lemmas in Appendix C

### D.1. Proof of Lemma C.1

*Proof.* We simply use  $\mathbf{D}'$  to denote  $\mathbf{D}'_j$ . If for some  $j \in [m]$ ,  $[\mathbf{D}']_{j,j} \neq 0$ , then it holds that

$$|[f'_{l,i}]_j| = |[f'_{l,i,1}]_j + [f'_{l,i,2}]_j| > |[f_{l,k}^{(0)}]_j|. \quad (46)$$

Let  $\xi : \xi \leq \frac{1}{2\sqrt{m}}$  and  $\xi \leq 2\|f'_{l,i,2}\|_\infty$  be a parameter to be chosen later. We then discuss the zero norm of  $\mathbf{D}'$  in the following two cases.

First, we consider the case that  $|[f_{l,k}^{(0)}]_j| \leq \xi$ . In this case, (46) is easy to be satisfied. Denote  $S_1 = \{j \in [m] \mid |[f_{l,k}^{(0)}]_j| \leq \xi\}$ . Since  $[f_{l,k}^{(0)}]_j \sim \mathcal{N}(0, \frac{2}{m})$ , we have  $\mathbb{P}\{|[f_{l,k}^{(0)}]_j| \leq \xi\} \leq O(\xi\sqrt{m})$ . Since  $|S_1| = \sum_{j=1}^m \mathbf{1}(\{|[f_{l,k}^{(0)}]_j| \leq \xi\})$ , we have  $\mathbb{E}[\exp(|S_1|)] \leq \exp(\xi m^{3/2}(e-1))$ . Thus, by Chernoff bound,  $\mathbb{P}(|S_1| \geq 2\xi m^{3/2}) \leq \frac{\mathbb{E}[\exp(|S_1|)]}{\exp(\xi m^{3/2})} \leq \exp(\xi m^{3/2}(e-3))$ . Hence, with probability at least  $1 - \exp(-\Omega(m^{3/2}\xi))$ , we have

$$|S_1| \leq O(\xi m^{3/2}).$$

Then, for  $j \in S_1$  such that  $[\mathbf{D}']_{j,j} \neq 0$ , we have  $|[\mathbf{D}' f'_{l,i}]_j| \leq |[f_{l,k}^{(0)}]_j| + |[f'_{l,i,1}]_j| + |[f'_{l,i,2}]_j| \leq |[f'_{l,i,1}]_j| + 3\xi/2$ . Further, we have

$$\sum_{j \in S_1} [\mathbf{D}' f'_{l,i}]_j^2 \leq O(\|f'_{l,i,1}\|^2 + \xi^2 |S_1|) \leq O(\|f'_{l,i,1}\|^2 + \xi^2 |S_1|) \leq O(\|f'_{l,i,1}\|^2 + \xi^3 m^{3/2}).$$

Second, we consider the case that  $|[f_{l,k}^{(0)}]_j| > \xi$ . Denote  $S_2 = \{j \in [m] \mid |[f_{l,k}^{(0)}]_j| > \xi, [\mathbf{D}']_{j,j} \neq 0\}$ . Then, (46) requires that

$$|[f'_{l,i,1}]_j| = |[f'_{l,i}]_j - [f'_{l,i,2}]_j| \geq |[f'_{l,i}]_j| - |[f'_{l,i,2}]_j| \geq |[f_{l,k}^{(0)}]_j| - |[f'_{l,i,2}]_j| \geq \xi - \|f'_{l,i,2}\|_\infty \geq \xi/2.$$



Thus we have

$$|S_2| \leq \frac{4\|f'_{l,i,1}\|^2}{\xi^2}.$$

Then since for  $j \in S_2$  such that  $[D']_{j,j} \neq 0$ , the signs of  $[f_{l,k}^{(0)}]_j + [f'_{l,i,1}]_j + [f'_{l,i,2}]_j$  and  $[f_{l,k}^{(0)}]_j$  are opposite, we have

$$|[D'f_{l,i}^{(0)}]_j| = |[f_{l,k}^{(0)}]_j + [f'_{l,i,1}]_j + [f'_{l,i,2}]_j| \leq |[f'_{l,i,1}]_j + [f'_{l,i,2}]_j| \leq |[f'_{l,i,1}]_j| + \xi/2 \leq 2|[f'_{l,i,1}]_j|.$$

Therefore, it holds that

$$\sum_{j \in S_2} [D'f_{l,i}^{(0)}]_j^2 \leq 4 \sum_{j \in S_2} |[f'_{l,i,1}]_j|^2 \leq 4\|f'_{l,i,1}\|^2.$$

Combining the two cases, we have

$$\|D\|_0 \leq |S_1| + |S_2| \leq O\left(\xi m^{3/2} + \frac{4\|f'_{l,i,1}\|^2}{\xi^2}\right),$$

$$\|D'f_{l,i}^{(0)}\|^2 \leq O(\|f'_{l,i,1}\|^2 + \xi^3 m^{3/2}).$$

Choosing  $\xi = \max\left\{2\|f'_{l,i,2}\|_\infty, \Theta\left(\frac{\|f'_{l,i,1}\|^{2/3}}{m^{1/2}}\right)\right\}$ , and recalling  $\|f'_{l,i,1}\| \leq O(L^{3/2}\phi \log^{1/2}(1/\phi))$  and  $\|f'_{l,i,2}\|_\infty \leq O(L\phi^{2/3}/m^{1/2} \log^{1/2}(1/\phi))$ , we get  $\|D\|_0 \leq O(mL\phi^{2/3} \log^{1/2}(1/\phi))$ . Choosing  $\xi = 2\|f'_{l,i,2}\|_\infty$ , we get  $\|D'f_{l,i}^{(0)}\| \leq O(\phi L^{3/2} \log^{1/2}(1/\phi))$ .  $\square$

## D.2. Proof of Lemma C.2

**Lemma D.1** (Lemma B.1 in (Zou et al., 2020)). *Assume  $m > \Omega(L \log(NL))$ . For any  $x'_i, x'_j \in \mathcal{X}_\phi$ ,  $i, j \in [N], l \in [L]$ , with probability at least  $1 - \exp(-O(m/L))$  over the randomness of  $\mathbf{W}^{(0)}$ , it holds that  $1/2 \leq \|h_l(x'_i)\| \leq 2$  and  $\|h_l(x'_i)/\|h_l(x'_i)\| - h_l(x'_j)/\|h_l(x'_j)\|\| \geq \phi/2$ , where  $h_l(x'_i)$  is the output of the  $l$ -th layer at initialization.*

Denote  $b_i = h_{\mathbf{W}^{(0)}, L-1}(x_i)$  and  $\bar{b}_i = b_i/\|b_i\|$  for  $x_i \in \mathcal{X}$ , and  $b'_i = h_{\mathbf{W}^{(0)}, L-1}(x'_i)$  and  $\bar{b}'_i = b'_i/\|b'_i\|$  for  $x'_i \in \mathcal{X}_\phi$ . By Lemma D.1, we have  $\forall i \notin \mathcal{I}_{\phi,k}, \|\bar{b}_i - \bar{b}'_k\| \geq \phi/4$ . Moreover, by Lemma B.1, we have  $\forall i \in \mathcal{I}_{\phi,k}, \|b_i - b'_j\| \leq \tilde{O}(L^{5/2}\phi \log^{1/2}(m))$ .

Then we construct several sets for the vector  $w \in \mathcal{R}^m$  subject to  $\mathcal{N}(0, (2/m)\mathbf{I})$ . Given  $\bar{b}'_k$ , we construct an orthogonal matrix  $Q_k = [\bar{b}'_k, Q'_k] \in \mathcal{R}^{m \times m}$  and let  $q_k = Q_k^\top w \sim \mathcal{N}(0, (2/m)\mathbf{I})$ . In this way, the vector  $w$  is decomposed as two orthogonal vector:  $w = Q_k q_k = q_k^{(1)} \bar{b}'_k + Q'_k q'_k$  where  $q_k^{(1)}$  is the first element of  $q_k$ . Letting  $\gamma = \sqrt{2\pi\phi}/(32N\sqrt{m})$ , we construct the set

$$\mathcal{W}_k = \left\{w \in \mathcal{R}^d \mid |q_k^{(1)}| \leq \gamma, |\langle Q'_k q'_k, \bar{b}'_j \rangle| \geq 2\gamma, \forall j \neq k\right\}, \quad (47)$$

where  $[q_k^{(1)}, q'_k] = q_k$ .

**Lemma D.2** (Lemma C.1 in (Zou & Gu, 2019)). *For any  $\mathcal{W}_j$  and  $\mathcal{W}_k$ ,  $j \neq k$ , we have  $\mathcal{W}_j \cap \mathcal{W}_k = \emptyset$  and  $\mathbb{P}(w \in \mathcal{W}_k) \geq \frac{\phi}{N32\sqrt{2e}}$ .*

**Lemma D.3.** *Let  $f(w_j) = \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi,k}} a_i \sigma'(\langle w_j, b_i \rangle) b_i$  where  $w_j, j \in [m]$  is drawn from  $\mathcal{N}(0, (2/m)\mathbf{I})$ ,  $|a_i| \leq O((\mu_i + \lambda_i)/\sqrt{d})$ . If for each smooth set  $k$ , there exists a subset  $\mathcal{G}_{k,\alpha} \in [m]$  with size  $\alpha m, \alpha \in (0, 1)$  such that  $\forall i \in \mathcal{I}_{\phi,k}, \forall j \in \mathcal{G}_{k,\alpha}, \sigma'(\langle w_j, b_i \rangle) = \sigma'(\langle w_j, b'_k \rangle)$  and  $\forall j \notin \mathcal{G}_{k,\alpha}, |\langle w_j, b_i \rangle| \geq \frac{3\sqrt{2\pi}\phi}{16N\sqrt{m}}$ , we have for any  $j \in \mathcal{G}_{k,\alpha}$ ,  $\mathbb{P}\left(\|f(w_j)\| \geq |A_k|/4 - M_k/\sqrt{d} \tilde{O}(L^{5/2}\phi \log^{1/2}(m)) \mid w_j \in \mathcal{W}_k\right) \geq 1/2$  where  $A_k = \sum_{i \in \mathcal{I}_{\phi,k}} a_i, k \in [N]$ .*

*Proof.* For  $j \in \mathcal{G}_{k,\alpha}$ , let  $q_k = Q_k^\top w_j \sim \mathcal{N}(0, (2/m)\mathbf{I})$ . Then we have  $w_j = Q_k q_k = q_k^{(1)} \bar{b}'_k + Q'_k q'_k$ . We decompose

$f(w_j)$  as

$$\begin{aligned}
 f(w_j) &= \sum_{i \in \mathcal{I}_{\phi, k}} a_i \sigma'(\langle w, b_i \rangle) b_i + \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle w, b_i \rangle) b_i \\
 &= \sum_{i \in \mathcal{I}_{\phi, k}} a_i \sigma'(\langle w, b'_k \rangle) b_i + \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle w, b_i \rangle) b_i \\
 &= \sum_{i \in \mathcal{I}_{\phi, k}} a_i \sigma'(q_k^{(1)}) b_i + \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle w, b_i \rangle) b_i,
 \end{aligned} \tag{48}$$

where the second equality holds by the assumption  $\forall j \in \mathcal{G}_{k, \alpha}$ ,  $\sigma'(\langle w_j, b_i \rangle) = \sigma'(\langle w_j, b'_k \rangle)$ .

Then for the second term of (48), if  $j \in \mathcal{G}_{k', \alpha}$ , we have for  $i \in \mathcal{I}_{\phi, k'}$ ,  $\sigma'(\langle w_j, b_i \rangle) = \sigma'(\langle w_j, b'_{k'} \rangle)$  and thus

$$\begin{aligned}
 &\sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle w, b_i \rangle) b_i = \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle w, b'_{k'} \rangle) b_i \\
 &= \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(q_k^{(1)} \langle \bar{b}'_k, b'_{k'} \rangle + \langle Q'_k q'_k, b'_{k'} \rangle) b_i \\
 &= \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle Q'_k q'_k, b'_{k'} \rangle) b_i
 \end{aligned}$$

where the last equality holds by the condition  $w_j \in \mathcal{W}_k$  such that for  $k' \neq k$ ,  $|\langle Q'_k q'_k, b'_{k'} \rangle| \geq 2\gamma \|b'_{k'}\| \geq |q_k^{(1)}| \|b'_{k'}\| \geq |q_k^{(1)}| \langle \bar{b}'_k, b'_{k'} \rangle > |q_k^{(1)}| \langle \bar{b}'_k, b'_{k'} \rangle$  and thus the sign is determined by  $\langle Q'_k q'_k, b'_{k'} \rangle$ . Therefore, if  $j \in \mathcal{G}_{k', \alpha}$ , we can write (48) as

$$f(w_j) = \sum_{i \in \mathcal{I}_{\phi, k}} a_i \sigma'(q_k^{(1)}) b_i + \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle Q'_k q'_k, b'_{k'} \rangle) b_i. \tag{49}$$

In the other case with  $j \notin \mathcal{G}_{k', \alpha}$ , by assumption  $\forall i \in \mathcal{I}_{\phi, k'}$ ,  $|\langle w_j, b_i \rangle| \geq \frac{3\sqrt{2\pi}\phi}{16N\sqrt{m}} = 6\gamma$ , we have with probability at least  $1 - \exp(-O(m/L))$ ,  $|\langle w_j, \bar{b}_i \rangle| = |\langle w_j, b_i \rangle| \frac{1}{\|\bar{b}_i\|} \geq 3\gamma$  by Lemma D.1. Then  $\forall i \in \mathcal{I}_{\phi, k'}$ , we have

$$\begin{aligned}
 &|\langle Q'_k q'_k, \bar{b}_i \rangle| = |\langle w_j, \bar{b}_i \rangle - \langle q_k^{(1)} \bar{b}'_k, \bar{b}_i \rangle| \\
 &\geq |\langle w_j, \bar{b}_i \rangle| - |\langle q_k^{(1)} \bar{b}'_k, \bar{b}_i \rangle| \geq |\langle w_j, \bar{b}_i \rangle| - |q_k^{(1)}| \\
 &\geq |\langle w_j, \bar{b}_i \rangle| - \gamma \geq 2\gamma \\
 &\geq 2\gamma \|b'_{k'}\| \geq |q_k^{(1)}| \|b'_{k'}\| \geq |q_k^{(1)}| \langle \bar{b}'_k, b'_{k'} \rangle > |q_k^{(1)}| \langle \bar{b}'_k, b'_{k'} \rangle,
 \end{aligned}$$

where the first inequality comes from triangle inequality, the second inequality holds by  $|\langle \bar{b}'_k, \bar{b}_i \rangle| \leq 1$ , and the last inequality holds by the condition  $w_j \in \mathcal{W}_k$ . Therefore if  $j \notin \mathcal{G}_{k', \alpha}$ , we can write the second term in (48) as

$$\begin{aligned}
 &\sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle w, b_i \rangle) b_i \\
 &= \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(q_k^{(1)} \langle \bar{b}'_k, b_i \rangle + \langle Q'_k q'_k, b_i \rangle) b_i \\
 &= \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle Q'_k q'_k, b_i \rangle) b_i
 \end{aligned}$$

Therefore, if  $j \notin \mathcal{G}_{k', \alpha}$ , we can write (48) as

$$f(w_j) = \sum_{i \in \mathcal{I}_{\phi, k}} a_i \sigma'(q_k^{(1)}) b_i + \sum_{k' \neq k} \sum_{i \in \mathcal{I}_{\phi, k'}} a_i \sigma'(\langle Q'_k q'_k, b_i \rangle) b_i. \tag{50}$$

Note that (49) and (50) are different only in terms of whether  $b'_k$  or  $b_i, i \in \mathcal{I}_{\phi, k'}$  determines the second term, but for both of them, the second term does not rely on  $q_k^{(1)}$ . We thus proceed as follows.

Since  $q_k^{(1)} > 0$  and  $q_k^{(1)} < 0$  occurs with equal probability conditioned on the event  $w \in \mathcal{W}_k$ , we have

$$\mathbb{P} \left[ \|f(w_j)\|_2 \geq \inf_{q_1 > 0, q_2 < 0} \max \{ \|f(q_1 \bar{b}'_k + Q'_k q'_k)\|, \|f(q_2 \bar{b}'_k + Q'_k q'_k)\| \} \mid w \in \mathcal{W}_k \right] \geq 1/2.$$

Thus, with probability at least  $1/2$  conditioned on the event  $w \in \mathcal{W}_k$ , we have

$$\begin{aligned} \|f(w_j)\| &\geq \inf_{q_1 > 0, q_2 < 0} \max \{ \|f(q_1 \bar{b}'_k + Q'_k q'_k)\|, \|f(q_2 \bar{b}'_k + Q'_k q'_k)\| \} \\ &\geq \inf_{q_1 > 0, q_2 < 0} \|f(q_1 \bar{b}'_k + Q'_k q'_k) - f(q_2 \bar{b}'_k + Q'_k q'_k)\| / 2 \\ &= \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b_i \right\| \end{aligned}$$

Since  $|a_i| \leq O((\mu_i + \lambda_i)/\sqrt{d})$  and  $\|b_i - b'_k\| \leq \tilde{O}(L^{5/2} \phi \log^{1/2}(m))$  for  $i \in \mathcal{I}_{\phi, k}$ , we have,  $\left| \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b_i \right\| - \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b'_k \right\| \right| \leq \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i (b_i - b'_k) \right\| \leq \sum_{i \in \mathcal{I}_{\phi, k}} |a_i| \|b_i - b'_k\| \leq M_k \sqrt{d} \tilde{O}(L^{5/2} \phi \log^{1/2}(m))$ . Thus,

$$\begin{aligned} \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b_i \right\| &\geq \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b'_k \right\| - \left| \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b_i \right\| - \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b'_k \right\| \right| \\ &\geq \left\| \sum_{i \in \mathcal{I}_{\phi, k}} a_i b'_k \right\| - M_k \sqrt{d} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \\ &\geq |A_k|/4 - M_k / \sqrt{d} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)), \end{aligned}$$

where the last inequality follows from Lemma D.1. The proof is completed.  $\square$

**Lemma D.4** (Bernstein inequality). *Let  $X_1, \dots, X_n$  be independent zero-mean random variables. If  $|X_i| \leq 1$  almost surely for all  $i$ , then  $\forall t > 0$ ,*

$$\mathbb{P} \left( \sum_{i=1}^n X_i \geq t \right) \leq \exp \left( -\frac{1}{2} t^2 / \left( \sum_{i=1}^n \mathbb{E}[X_i^2] + \frac{1}{3} t \right) \right).$$

### Proof of Lemma C.2

*Proof.* Denote  $b_i = h_{L-1}^{(0)}(x_i)$ , so  $[h_L^{(0)}(x_i)]_j = \langle w_j, b_i \rangle$ . For any fixed  $[u_1, \dots, u_{n'}]$ , denote  $a_i(v_j) = \langle u_i, v_j \rangle, i \in [n'], j \in [m]$  and  $A_k(v_j) = \sum_{i \in \mathcal{I}_{\phi, k}} a_i(v_j)$ . Let  $f(v_j, w_j) = \sum_{k=1}^N \sum_{i \in \mathcal{I}_{\phi, k}} a_i(v_j) \sigma'(\langle w_j, b_i \rangle) b_i$ . Define the event for  $k \in [N]$

$$\mathcal{E}_k = \left\{ j \in \mathcal{G}_{k, \alpha} : w_j \in \mathcal{W}_k, \|f(v_j, w_j)\| \geq \left\| \sum_{i \in \mathcal{I}_{\phi, k}} u_i \right\| / (4\sqrt{d}) - M_k d^{-1/2} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right\}.$$

Since  $v_j \sim \mathcal{N}(0, (1/d)\mathbf{I})$ , we have  $A_k(v_j) = \langle \sum_{i \in \mathcal{I}_{\phi, k}} u_i, v_j \rangle \sim \mathcal{N}(0, \|\sum_{i \in \mathcal{I}_{\phi, k}} u_i\|^2/d)$ . Thus, we have

$$\mathbb{P} \left( \left\| \sum_{i \in \mathcal{I}_{\phi, k}} u_i \right\| / \sqrt{d} \leq |A_k(v_j)| \leq 2 \left\| \sum_{i \in \mathcal{I}_{\phi, k}} u_i \right\| / \sqrt{d} \right) \geq 1/4.$$

Note that when  $|A_k(v_j)| \leq 2 \left\| \sum_{i \in \mathcal{I}_{\phi, k}} u_i \right\| / \sqrt{d}$ , we have  $\forall i \in \mathcal{I}_{\phi, k}, |a_i(v_j)| \leq |A_k(v_j)| / |\mathcal{I}_{\phi, k}| - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \leq 2 \left\| \sum_{i \in \mathcal{I}_{\phi, k}} u_i \right\| / \sqrt{d} / |\mathcal{I}_{\phi, k}| - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \leq 2(\mu_i + \lambda_i) / \sqrt{d} - 3\tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \leq 3(\mu_i + \lambda_i) \sqrt{d}$  when  $\phi$  is small enough, so the condition about  $|a_i(v_j)|$  in Lemma D.3 is met.

Since Assumption 1 is satisfied, we have for each smooth set  $k$ , there exists a subset  $\mathcal{G}_{k,\alpha} \in [m]$  with size  $\alpha m$ ,  $\alpha \in (0, 1)$  such that  $\forall i \in \mathcal{I}_{\phi,k}, \forall j \in \mathcal{G}_{k,\alpha}, \sigma'(\langle w_j, b_i \rangle) = \sigma'(\langle w_j, b'_k \rangle)$  and  $\forall j \notin \mathcal{G}_{k,\alpha}, |\langle w_j, b_i \rangle| \geq \frac{3\sqrt{2}\pi\phi}{16N\sqrt{m}}$ , so the assumption in Lemma D.3 is satisfied. Then by Lemma D.2, Lemma D.3 and the fact that  $w_j$  and  $v_j$  are independent, we have for  $j \in \mathcal{G}_{k,\alpha}$

$$\begin{aligned} \mathbb{P}(j \in \mathcal{E}_k) &= \mathbb{P} \left\{ \|f(v_j, w_j)\| \geq |A_k(v_j)|/4 - M_k d^{-1/2} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \mid w_j \in \mathcal{W}_k \right\} \\ &\cdot \mathbb{P} \{w_j \in \mathcal{W}_k\} \mathbb{P} \left( \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| / \sqrt{d} \leq |A_k(v_j)| \leq 2 \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| / \sqrt{d} \right) \geq \frac{\phi}{N 256 \sqrt{2} e} = p_\phi. \end{aligned}$$

and  $\mathcal{E}_{k_1} \cap \mathcal{E}_{k_2} = \emptyset$  for any  $k_1 \neq k_2$ .

For smooth set  $k$ , denote Bernoulli random variables  $\mathbb{1}(j \notin \mathcal{E}_k)$  for  $j \in \mathcal{G}_{k,\alpha}$ . Then we have  $\mathbb{E}[\mathbb{1}(j \notin \mathcal{E}_k)] = 1 - p_\phi$  and  $\text{var}[\mathbb{1}(j \notin \mathcal{E}_k)] = p_\phi(1 - p_\phi)$ . By Bernstein inequality in Lemma D.4 for random variables  $\mathbb{1}(j \notin \mathcal{E}_k) - (1 - p_\phi), j \in [m]$ , it holds that

$$\mathbb{P} \left( \sum_{j \in \mathcal{G}_{k,\alpha}} \mathbb{1}(j \notin \mathcal{E}_k) - \alpha m (1 - p_\phi) \geq \frac{\alpha m p_\phi}{2} \right) \leq \exp \left( -\frac{\frac{1}{4} \alpha m p_\phi}{1 - p_\phi + \frac{1}{6}} \right) \leq \exp \left( -\frac{3}{14} \alpha m p_\phi \right).$$

Thus, by union bounds, with probability at least  $1 - O(N) \exp(-O(\alpha m \phi / N))$ , we have for any  $k \in [N], \sum_{j \in \mathcal{G}_{k,\alpha}} \mathbb{1}(j \notin \mathcal{E}_k) \leq \alpha m - \alpha m p_\phi / 2$  and

$$\left| \mathcal{G}_{k,\alpha} \cap \mathcal{E}_k \right| = \sum_{j \in \mathcal{G}_{k,\alpha}} \mathbb{1}(j \in \mathcal{E}_k) = \alpha m - \sum_{j \in \mathcal{G}_{k,\alpha}} \mathbb{1}(j \notin \mathcal{E}_k) \geq \alpha m p_\phi / 2. \quad (51)$$

Therefore, with probability at least  $1 - O(\phi)$ , it holds that

$$\begin{aligned} &\sum_{j=1}^m \|f(v_j, w_j)\|^2 \geq \sum_{j=1}^m \|f(v_j, w_j)\|^2 \sum_{k=1}^N \mathbb{1}(j \in \mathcal{E}_k) = \sum_{k=1}^N \sum_{j=1}^m \|f(v_j, w_j)\|^2 \mathbb{1}(j \in \mathcal{E}_k) \\ &\geq \sum_{k=1}^N \sum_{j \in \mathcal{G}_{k,\alpha}} \|f(v_j, w_j)\|^2 \mathbb{1}(j \in \mathcal{E}_k) = \sum_{k=1}^N \sum_{j \in \mathcal{G}_{k,\alpha}, j \in \mathcal{E}_k} \|f(v_j, w_j)\|^2 \\ &\geq \sum_{k=1}^N \sum_{j \in \mathcal{G}_{k,\alpha}, j \in \mathcal{E}_k} \left( \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| / (4\sqrt{d}) - M_k d^{-1/2} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right)^2 \\ &\geq \sum_{k=1}^N \sum_{j \in \mathcal{G}_{k,\alpha}, j \in \mathcal{E}_k} \left( \frac{1}{16d} \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\|^2 - M_k^2 d^{-1} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right) \quad (52) \\ &= \sum_{k=1}^N \left( \frac{1}{16d} \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\|^2 - M_k^2 d^{-1} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right) \left| \mathcal{G}_{k,\alpha} \cap \mathcal{E}_k \right| \\ &\geq \frac{\alpha m p_\phi}{2} \sum_{k=1}^N \left( \frac{1}{16d} \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\|^2 - M_k^2 d^{-1} \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right) \\ &\geq \Omega \left( \frac{\alpha m \phi}{Nd} \right) \left( \sum_{k=1}^N \left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\|^2 - \tilde{O}(L^{5/2} \phi \log^{1/2}(m)) \right), \end{aligned}$$

where the first inequality comes from the fact that  $\mathcal{E}_{k_1} \cap \mathcal{E}_{k_2} = \emptyset$  such that  $\sum_{k=1}^N \mathbb{1}(j \in \mathcal{E}_k) \leq 1$  and the second inequality comes from the fact that  $\mathcal{G}_{k,\alpha} \in [m]$ , and the third inequality holds by the definition of event  $\mathcal{E}_k$ , and the fourth inequality comes from the fact that  $(a - b)^2 \geq a^2 - 2ab$  and  $\left\| \sum_{i \in \mathcal{I}_{\phi,k}} u_i \right\| / (4\sqrt{d}) \leq M_k d^{-1/2} / 4$ , and the fifth inequality comes from (51) and the last inequality holds by the fact that  $\sum_{k=1}^N M_k^2 \leq (\sum_{k=1}^N M_k)^2 = 1$ .  $\square$

## E. Preliminaries on Informed Machine Learning

Informed machine learning is rapidly emerging as a broad paradigm that incorporates domain knowledge, either directly or indirectly, to augment the purely data-driven approach and better accomplish a machine learning task. We provide a summary of how domain knowledge is integrated with machine learning (von Rueden et al., 2021).

- *Training Dataset.* A straightforward approach to utilizing domain knowledge is to generate (sometimes synthetic) data and enlarge the otherwise limited training dataset. For example, based on the simple knowledge of image invariance, cropping (Gao et al., 2018a), scaling (Zhang et al., 2018), flipping (Benaim & Wolf, 2018) and many other image pre-processing methods have been used to augment the training data for image classification tasks. As another example, in reinforcement learning (e.g., robot control and autonomous driving) where initial pre-training is crucial to avoid arbitrarily bad decisions in the real world, simulated environments can be built based on domain knowledge, providing simulations or demonstrations to generate training data (Gao et al., 2018b; Hester et al., 2018). Additionally, generative models constructed based on specific knowledge have been shown useful for increasing training data to improve model performance and robustness (Gao et al., 2018a; Goodfellow et al., 2016).
- *Hypothesis Set.* The goal of a machine learning task is to search for an optimal hypothesis that correctly expresses the relationships between input and output. To reduce the training complexity, the target hypothesis set (decided by, e.g., different neural architectures) should contain the optimal hypothesis and preferably be small enough. Thus, domain knowledge can be employed for hypothesis set selection. For example, (Chen et al., 2021a) makes use of the prior knowledge from the existing neural architectures to design new architectures (and hence, new hypothesis sets) for DNNs. As implicit domain knowledge, long short-term memory recurrent neural networks are commonly used for time series prediction (Goodfellow et al., 2016). Also, the structure of a knowledge graph helps to determine the hypothesis set of graph learning (Marino et al., 2016; Battaglia et al., 2018), while (Towell & Shavlik, 1994) maps the domain knowledge represented in propositional logic into neural networks.
- *Model Training.* Domain knowledge can be integrated, either implicitly or explicitly, with the model training procedure in various ways. First, domain knowledge can assist with the initialization of training. For example, (Ramsey & Grefenstette, 1993) provides a case-based method to initialize genetic algorithms (i.e., generating the initial population based on different cases), while (Husken & Goerick, 2000; Kurata et al., 2016; Humbird et al., 2018) initialize neural network training with various domain knowledge such as label co-occurrence and decision trees. Second, domain knowledge can be used to better tune the hyper-parameters (Bardnet et al., 2013; Van Rijn & Hutter, 2018; Maher & Sakr, 2019; Bamler et al., 2020). In (Bardnet et al., 2013), implicit knowledge from previous training is incorporated to improve hyper-parameter tuning, and (Van Rijn & Hutter, 2018) extracts knowledge from multiple datasets to determine the most important hyper-parameters. In addition, a more explicit way to integrate domain knowledge is to directly modify the training objective function (i.e., risk function) based on rigorous characterization of the model output (von Rueden et al., 2021). For example, in (Muralidhar et al., 2018), the knowledge of constraints is incorporated into neural networks expressing the knowledge based loss by the ReLU function. For another example, when learning to optimally schedule transmissions for rate maximization in multi-user wireless networks, the communication channel capacity can be added as domain knowledge to the standard label-based loss to guide scheduling decisions; in physics, the analytical expression of a partial differential equation can be utilized as domain knowledge on top of labeled data to better learn the solution to the equation given different inputs; more examples are shown in Section F. Such integration of explicit and rigorous domain knowledge can significantly benefit machine learning tasks (e.g., fewer labels needed than otherwise). Thus, it is crucial and being actively studied in informed machine learning (von Rueden et al., 2021; Willard et al., 2020), which is also the focus of our work. Note that using domain knowledge to generate pseudo labeled data to augment the training dataset is a special case of integrating domain knowledge into the training risk function (i.e., the knowledge-based risk is the same as the data-based risk, except that its labels are generated based on domain knowledge).
- *Final Hypothesis.* Domain knowledge can also be used for consistency check on the final learnt hypothesis or model (von Rueden et al., 2021). For example, (Karpatne et al., 2017) employs physics domain knowledge to construct the final model, (Pfrommer et al., 2018) builds simulators to validate results of learned model, and (Fang et al., 2017) leverages semantic consistency is used to refine the predicted probabilities.

## F. Application Examples

We now present a few application examples to explain domain knowledge-informed DNNs.

### F.1. Learning for resource management in communications networks

Optimizing resource management is crucial to improve the system performance in communications networks (Chiang et al., 2008; Goldsmith, 2005; Zappone et al., 2019). Well-known examples include power allocation (Hong & Luo, 2014; Chiang et al., 2008; Liang et al., 2019), link scheduling (Gore & Karandikar, 2010; Cui et al., 2019), antenna or beam selection (Sanayei & Nosratinia, 2004; Klautau et al., 2018), among others. While many of the problems were studied using theoretical model-based approaches in the past, machine learning has been increasingly employed, in view of the rapidly growing complexity of communications technologies that theoretical models are often incapable of capturing accurately (Zappone et al., 2019). Let us take power allocation in multi-user wireless interference networks as an example. The recent work (Sun et al., 2018) uses a pure data-driven approach for power allocation to maximize the sum rate: a labeled dataset containing channel state information (CSI) and the corresponding power allocation decisions is collected in advance, and a neural network is trained to learn the optimal power allocation. On the other hand, Shannon-based transmission rate has been extensively used as an analytical objective function to optimize power allocation, and (Liang et al., 2019) exploits this domain knowledge to train an ensemble of neural networks that directly learn the optimal power allocation for Shannon rate maximization.

The data-driven approach (Sun et al., 2018) can maximize the practically achievable rate (if labels are collected from real systems), but is significantly constrained by the limited amount of training samples. Meanwhile, the knowledge-based approach (Liang et al., 2019) can utilize a large number of input samples (at the expense of higher training complexity), but the resulting power allocation decisions may not maximize the sum rate in real systems. The reason is that the Shannon formula for interference channels, albeit commonly used for analysis, only represents an approximation of the achievable rate which is subject to finite channel code lengths and modulation schemes (Goldsmith, 2005). In other words, even an oracle DNN that minimizes this knowledge-based loss may not maximize the achievable rate in practice.

To reap the benefit of both labeled data and domain knowledge, informed machine learning can be adopted, resulting in a new informed loss as follows:

$$\min_{h \in \mathcal{H}} (1 - \gamma) \left\{ \frac{1}{n} \sum_{(x,y) \in S} [h(x) - y]^2 \right\} + \gamma \cdot \left\{ -\frac{1}{\tilde{n}} \sum_{x \in \tilde{S}_X} \text{Shannon\_rate}[h(x)] + \text{constant} \right\}, \quad (53)$$

where  $x$  is the input (e.g., channel state information),  $h(x)$  is the learned power allocation given  $x$ , the two loss terms represent label-based loss and knowledge-based loss, and  $n$  and  $\tilde{n}$  are the numbers of labeled data samples and (possibly unlabeled) knowledge samples, respectively. The detailed Shannon formula for wireless networks can be found in (Liang et al., 2019; Goldsmith, 2005).

### F.2. Image classification based on semantic knowledge

Typical image classifiers rely on labeled training data, but labels can be difficult and expensive to collect in practice (Goodfellow et al., 2016). As a result, few-shot learning (Wang et al., 2020; Sung et al., 2018; Garcia & Bruna, 2018) that only needs a small number of labeled samples has been proposed. Informed machine learning under our consideration can be viewed as few-shot learning. Concretely, semantic knowledge formulated as the first-order logic clauses/sentences (Xu et al., 2018; Diligenti et al., 2017b) can be incorporated to improve learning performance given limited labeled samples. An example logic clause is “if it is an animal and has wings, then it is a bird”. By a logic clause  $K$ , a knowledge-based loss can be defined as  $F_K(h(x), g(x))$  for an (possibly unlabeled) input image  $x$  and a certain logic clause  $g(x)$  that the output class  $h(x)$  needs to satisfy. Then, combining the standard label-based loss with knowledge-based loss, the model performance can be improved by minimizing the informed loss Eqn. (3) given limited labeled samples.

### F.3. Learning to solve PDEs in scientific and engineering fields

Partial differential equations (PDEs) are classic problems in many scientific and engineering fields, such as physics and mechanical engineering, but are notoriously difficult to solve in most practical settings (Institute, 2020; Baker et al., 2019). In recent years, physics knowledge-informed machine learning has been suggested as a promising approach to augment or even

replace classic PDE solution approaches (Deng et al., 2020; Willard et al., 2020; Khoo et al., 2021; Beck et al., 2019; Raissi et al., 2019; Lu et al., 2021). For example, (Raissi et al., 2019; Lu et al., 2021) proposes a physics-informed neural network (PINN) to solve PDEs by minimizing the PDE residual and penalties of boundary/initial conditions, which correspond to the knowledge-based loss  $F_K(h, g(h))$  in our framework. Additionally, we can combine the knowledge-based loss with labeled-based loss, achieving faster convergence and better performances in practice (especially when the PDE-based knowledge does not perfectly represent the real physical world). Take magnetic field strength estimation for magnetic materials as an example. If a few measured magnetic field strengths are provided as labels combined with the knowledge of Maxwell equations, the model trained by minimizing the informed loss can perform better in the real world. The measured labels can partly correct the imperfectness of physics knowledge, while the knowledge can improve the generalization in the presence of limited labels.

#### E.4. Knowledge distillation and transfer

Knowledge distillation (Hinton et al., 2014; Furlanello et al., 2018; Phuong & Lampert, 2019; Allen-Zhu & Li, 2020) is an important technique to transfer prior knowledge from a pre-trained neural network (a.k.a. teacher network) to another network (a.k.a. student network), with the same or different architectures. Typically, given an (possibly unlabeled) input, knowledge distillation is performed by matching the output of the student network with the output of the teacher network. In addition, labeled samples can also be included to introduce a label-based loss. Thus, by formulating  $g(X)$  as the output of the teacher network, knowledge distillation can be viewed as a particular instance of informed machine learning, where the knowledge comes from a teacher network and is usually assumed to be perfect.

## G. Numerical Results

We consider two specific applications — learning a multi-dimensional Bohachevsky function and learning to manage wireless spectrum.

### G.1. Settings of Learning with Constraint Knowledge in Section 7

We consider an informed DNN with domain knowledge in the form of constraints to learn a Bohachevsky function. The learning task is to learn a relationship  $y(x)$ . The learner is provided with a dataset with labeled samples  $S_z = \{(x_i, z_i), i \in [n_z]\}$ , having possibly noisy labels

$$z_i = y(x_i) + n_i, n_i \sim \mathcal{N}(0, \sigma_z^2),$$

and an unlabeled dataset  $S_g = \{(x_i), i \in [n_g]\}$ . Additionally, the learner is informed with the constraint knowledge, which includes an upper bound  $g_{\text{ub}}(x)$  and a lower bound  $g_{\text{lb}}(x)$  on the true label corresponding to input  $x$ , i.e.  $g_{\text{lb}}(x) \leq y(x) \leq g_{\text{ub}}(x)$ . A neural network  $h_{\mathbf{W}}(x)$  is used to learn the relationship  $y(x)$ , and the metric of interest is the mean square error (MSE) of the network output  $h_{\mathbf{W}}(x)$  with respect to the true label  $y(x)$  on a test dataset  $S_t$ , which is expressed as

$$\hat{R}_{S_t}(h_{\mathbf{W}}) = \frac{1}{2|S_t|} \sum_{(x_i, y_i) \in S_t} \text{mse}(h_{\mathbf{W}}(x_i), y_i),$$

where  $\text{mse}(h_{\mathbf{W}}(x_i), y_i) = (h_{\mathbf{W}}(x_i) - y_i)^2$  with  $y_i$  as the true test label with respect to  $x_i$ . Assume that the relationship to be learned is governed by a multi-dimensional Bohachevsky function

$$y(x) = x\mathbf{A}\mathbf{A}^\top x^\top - c \cos(a^\top x) + c,$$

where  $\mathbf{A}$  is a  $b \times b$  matrix,  $a$  is a  $b$ -dimensional vector and  $c$  is a constant. The learner has no access to the values of these parameters or the exact form of the relationship, but is empowered with the constraint knowledge in the form of an upper bound model

$$g_{\text{ub}}(x) = x\mathbf{A}\mathbf{A}^\top x^\top + ub$$

with  $ub \geq 2c$ , and a lower bound model

$$g_{\text{lb}}(x) = x\mathbf{A}\mathbf{A}^\top x^\top + lb.$$

with  $lb \leq 0$ . While it is not strongly convex and hence deviates from the assumptions in our theoretical analysis, we use ReLU as the knowledge-based risk function, i.e., the knowledge-based risk is written as

$$r_K(h_{\mathbf{W}}(x)) = \text{relu}(h_{\mathbf{W}}(x) - g_{\text{ub}}(x)) + \text{relu}(g_{\text{lb}}(x) - h_{\mathbf{W}}(x)).$$

And the label supervised risk given a sample pair  $(x, z)$  is  $r(h_{\mathbf{W}}(x)) = \text{mse}(h_{\mathbf{W}}(x), z)$ .

To show the performance under different levels of imperfectness, we consider labels with different noise variances and different knowledge-informed constraints. For training, the labeled dataset  $S_z$  contains  $n_z \in \{200, 400\}$  labeled samples with label noise variance  $\sigma_z^2 \in \{0, 0.1\}$ , and the unlabeled dataset  $S_g$  for the knowledge risk contains  $n_g = 1000$  input samples. The parameters for knowledge-informed constraint models include  $lb = 0$  and  $ub \in \{0.6, 0.8\}$ . Naturally, the higher variance  $\sigma_z^2$ , the worse label quality; and the greater  $ub$ , the worse knowledge quality. The test dataset  $S_t$  contains 1000 samples with labels calculated as  $y_i = y(x_i), x_i \in S_t$ .

For training, we use a neural network with two hidden layers, each having 2048 neurons and ReLU activations. Note that for the large network width needed for analysis to gain insights is not necessary in practice. The network is initialized based on Algorithm A. The training procedure is performed by Adam optimizer for 3000 steps with batch size 100. The learning rate is set as  $10^{-6}$  for the first 2000 steps,  $5 \times 10^{-5}$  for the following 500 steps, and  $10^{-5}$  for the remaining 500 steps. We run the network training with 10 random seeds. We run the simulations on a HPC cluster with GPUs of type P100.

## G.2. Learning for Resource Management in Wireless Networks

We apply an informed DNN to the problem of learning for resource management in wireless networks — wireless link scheduling in interference channels. We first describe problem setup, then present our method by informed DNN, and finally show the experiment results.

### G.2.1. PROBLEM SETUP

Link scheduling is a classic and important problem in wireless interference channels, with the objective of maximizing the sum throughput of wireless links. Consider a time-slotted wireless network consisting of a transmitter-receiver set  $\mathcal{U} = \{1, 2, \dots, N\}$  with  $N$  links (i.e., transmitter-receiver pairs) subject to cross-link interference. At the beginning of each time slot, the scheduler needs to decide a subset of links  $\mathcal{U}_S \subseteq \mathcal{U}$  to transmit depending on the channel state information (CSI).

We assume Rayleigh fading channels with interference across different links. If a link  $u \in \mathcal{U}$  is scheduled, the channel gain is  $g_{u,u}$  subject to Rayleigh fading. For notational convenience, we omit the time slot index. Multiple links can be scheduled at the same time slot, creating interference to each other. For example, if link  $u$  and link  $v$  are scheduled simultaneously, the interference channel gain from the transmitter  $u$  to receiver  $v$  is  $g_{u,v}$ , and the interference channel gain from the transmitter  $v$  to receiver  $u$  is  $g_{v,u}$ . Thus, the received signal at receiver  $u$  can be expressed as  $g_{u,u}s_u + \sum_{v \in \mathcal{U}_S/u} g_{v,u}s_v + \text{noise}_u$ , where  $\text{noise}_u \sim \mathcal{N}(0, \sigma_n^2)$  is an additive white Gaussian noise and the transmit signals  $s_u$  and  $s_v$  are normalized with unit power. Considering a centralized setting as in (Liang et al., 2019), the scheduler has access to the direct transmit channel gains as well as interference channel gains at the beginning of each time slot, which are contained in a  $N \times N$  dimensional CSI vector  $x = [g_{1,1}, \dots, g_{1,N}, g_{2,1}, \dots, g_{N-1,N}, g_{N,1}, \dots, g_{N,N}]$ .

The scheduling decision can be represented by a  $N$  dimensional scheduling vector  $y$ . Specifically, if the link  $u$  is scheduled, then the  $u$ -th entry of  $y$  is one, and zero otherwise. By the Shannon rate formula in the communications theory (Goldsmith, 2005), the achievable rate for link  $u$  can be expressed as

$$C_{\text{Shannon}}^u(x, y, \mu) = \log \left( 1 + \frac{\mu y(u) \|g_{u,u}\|^2}{\sigma_n^2 + \sum_{v \in \mathcal{U}/u} y(v) \|g_{v,u}\|^2} \right), \quad (54)$$

where  $\mu(0, 1]$  is a parameter subject to real communication systems, with  $C_{\text{Shannon}}^u(x, y, 1)$  representing the standard Shannon rate (i.e., when  $\mu = 1$ ). The sum rate is  $C_{\text{Shannon}}(x, y) = \sum_{u \in \mathcal{U}} C_{\text{Shannon}}^u(x, y)$ .

In practice, given the CSI vector  $x$  and the corresponding decision vector  $y$ , the real sum rate is denoted as  $C_{\text{real}}(x, y) = \sum_{u \in \mathcal{U}} C_{\text{real}}^u(x, y)$ . The real rate is difficult to express analytically in view of the complex factors in real environments including various schemes of modulation, finite channel coding and quality of service (QoS) guarantee. In fact, except for a few special cases, the *exact* channel capacity for general interference channels (even for two links) is still an open problem. Thus, while the Shannon rate is useful and has been utilized to design various systems, it only represents an approximation of the practically achievable rate.



Next, we formulate the link scheduling problem as

$$\max_y \sum_{u \in \mathcal{U}} C_{\text{real}}^u(x, y), \quad \text{s.t. } y(u) \in \{0, 1\}, u \in \mathcal{U}. \quad (55)$$

The scheduling objective is the real sum rate in a practical environment. The challenge of this problem is that the real rate in terms of the CSI  $x$  and scheduling decision  $y$  is too complex to express precisely, let alone the longstanding challenges of deriving the exact interference channel capacity (Goldsmith, 2005).

### G.2.2. INFORMED DNN FOR WIRELESS LINK SCHEDULING

DNNs have strong representation power to learn the optimal scheduling decisions given CSI input (Sun et al., 2018), but they typically require a large number of labeled samples  $(x, y)$  for training. On the other hand, domain knowledge (i.e., Shannon rate formula) is also useful, but it may not capture the real achievable rate in practice (Liang et al., 2019). Thus, informed DNN, which exploits domain knowledge to complement labeled samples, has the potential to reap the benefits of both approaches.

Concretely, we use a DNN to represent the relationship between the scheduling decision  $y$  and CSI  $x$ . Given  $N$  links, the input dimension is  $N \times N$ , which is the dimension of vectorized CSI  $x$  and the scheduling decision  $y$  is a  $N$ -dimensional binary vector. The training is based on a labeled dataset  $S_y = \{(x_i, y_i), i = 1, 2, \dots, n_y\}$  collected from real systems or field studies, where  $y_i$  is the true label (i.e., optimal scheduling decision) given  $x_i$ , along with the domain knowledge of Shannon rate. Also, we use  $Y_{\text{comb}} \in \{0, 1\}^{I_{\text{max}} \times N}$ ,  $I_{\text{max}} = 2^N - 1$  to represent all the possible decision combinations. Denote  $I(y)$  as the index of a scheduling decision  $y$  in  $Y_{\text{comb}}$ , i.e.  $y = [Y_{\text{comb}}]_{I(y)}$ . The output dimension of the DNN is  $I_{\text{max}} = 2^N - 1$  with each entry representing an index for a scheduling decision.

The label-based risk is the cross-entropy loss between the output of the DNN and one-hot encoding labels, which is expressed as

$$\hat{R}_{S_y}(\mathbf{W}) = \frac{1}{n_y} \sum_{i=1}^{n_y} \text{cross\_entropy}(\text{softmax}(h_{\mathbf{W}}(x_i)), \text{one\_hot}(I(y_i))), \quad (56)$$

where  $\text{one\_hot}(I(y_i))$  is the one-hot encoding of the index of  $y_i$ . Given an CSI input  $x$  and setting  $\mu = \mu_K$  based on domain experience, we can compute the sum rate of all possible scheduling decisions by the Shannon equation in Eqn. (54) as  $C_{\text{Shannon}}(x, [Y_{\text{comb}}]_j, \mu_K)$ ,  $j \in [I_{\text{max}}]$  and get the vector of sum rate as  $\mathbf{c}(x) = [C_{\text{Shannon}}(x, [Y_{\text{comb}}]_1, \mu_K), \dots, C_{\text{Shannon}}(x, [Y_{\text{comb}}]_{I_{\text{max}}}, \mu_K)]$ . Taking the softmax operation on  $T\mathbf{c}(x)$  with  $T$  as a scaling hyper-parameter, we get  $\text{softmax}(T\mathbf{c}(x))$ , which is essentially soft encoding of scheduling decisions based on the Shannon rate knowledge. Therefore, given an input dataset  $S_g = \{x_i, i = 1, 2, \dots, n_g\}$ , the knowledge-based risk is designed as

$$\hat{R}_K(\mathbf{W}) = -\frac{1}{n_g} \sum_{i=1}^{n_g} \text{cross\_entropy}(\text{softmax}(h_{\mathbf{W}}(x_i)), \text{softmax}(T\mathbf{c}(x))). \quad (57)$$

Thus, the DNN can be trained to minimize the informed risk combining both label-based and knowledge-based risks:  $\hat{R}_I(\mathbf{W}) = (1 - \lambda)\hat{R}_{S_y}(\mathbf{W}) + \lambda\hat{R}_K(\mathbf{W})$ . That is, the informed DNN uses hard labels for direct supervision, while exploiting domain knowledge in the form of soft labels for indirect supervision on unlabeled inputs. After training the network, the scheduling decision for CSI  $x$  is calculated as  $y_{\mathbf{W}}(x) = [Y_{\text{comb}}]_{I_{\mathbf{W}}(x)}$  with  $I_{\mathbf{W}}(x) = \arg \max_{j \in [I_{\text{max}}]} [h_{\mathbf{W}}(x)]_j$

### G.2.3. RESULTS

Now, we show the simulation results for the wireless link scheduling problem based on our informed DNN. We first give the simulation settings and then show the results of classification accuracy as well as the sum rate.

**Simulation Settings.** For illustration, we consider a simulation scenario with  $N = 4$  wireless links for scheduling, which is a reasonable setting for many practical ad hoc networks (Goldsmith, 2005). Given the CSI, the scheduler needs to choose one out of 15 scheduling combinations. To evaluate the performance of our informed DNN when the domain knowledge of Shannon rate is not perfect, we construct a synthetic dataset as the ground truth. The direct link channel gain of a wireless link  $g_{u,u}$  is subject to Rayleigh distribution, with an expected power gain of 100 dB. The cross-link interference channel gain is also subject to Rayleigh distribution with an expected power gain of 10 dB. The labels in the labeled training dataset

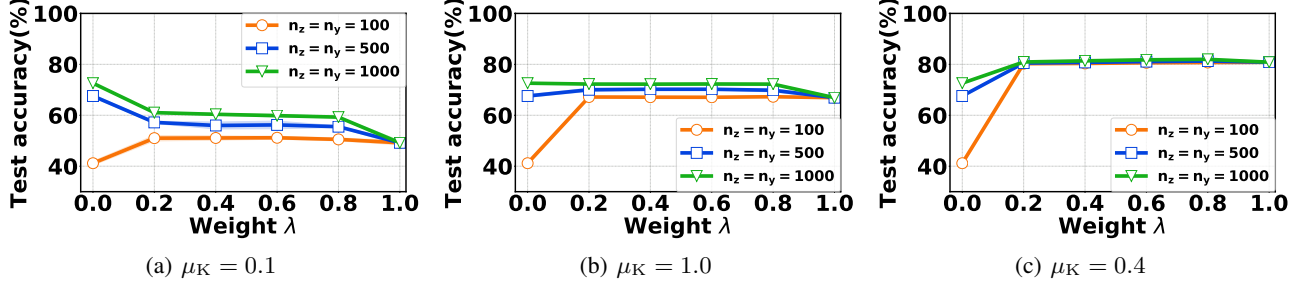


Figure 2. Test accuracy under different knowledge qualities and numbers of labels.

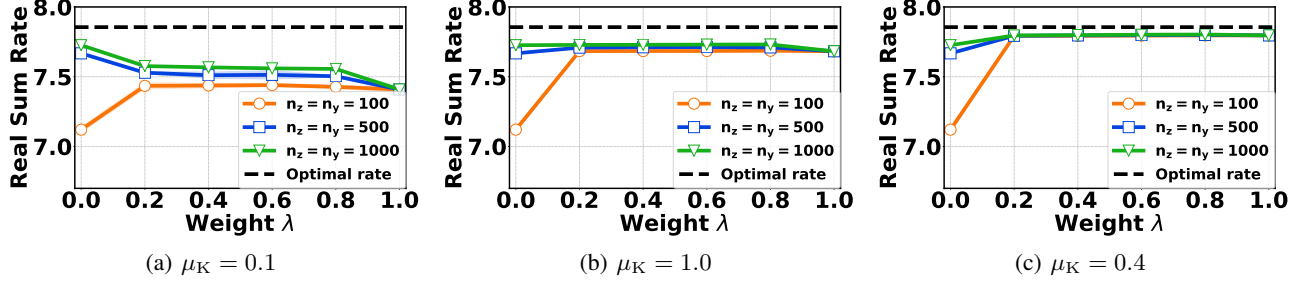


Figure 3. Sum rate under different knowledge qualities and numbers of labels.

and test dataset are generated by a pseudo-real rate expression to reflect some practical constraints:

$$C_{\text{pseudo-real}} = C_{\text{Shannon}}^u(x, y, \mu_R), \quad (58)$$

which differs from the standard Shannon formula by using a factor  $\mu_R \in (0, 1)$  to account for achievable rate degradation. Note that the pseudo-real rate is only defined to generate synthetic real rate different from the standard Shannon rate for evaluation purposes. In practice, the achievable rate is even more complex. In the simulations, we set  $\mu_R = 0.5$  to generate the training and testing labels as ground truth, while the value of  $\mu_R = 0.5$  is not available to the learner.

Based on the pseudo-real rate expression, we find the optimal labels (i.e., optimal scheduling decision  $y$ ) via exhaustive search, while labels are actually be collected by field measurement in a practical environment. We have  $n_g = 2000$  unlabeled CSI input samples in the training dataset  $S_g$  for knowledge-based supervision, and  $n_t = 10000$  samples in the test dataset. The test accuracy is defined as the percentage of DNN outputs that are identical to the optimal scheduling decision label, i.e. for samples in the test dataset,  $\text{acc} = \sum_{i=1}^{n_t} \mathbb{1}(I_{\mathbf{W}}(x_i) = I(y_i)) / n_t$ . We compare the results when the labeled training dataset has 100, 500 and 1000 samples, respectively. Also, we compare the results obtained by setting different parameters  $\mu_K \in \{1.0, 0.4, 0.1\}$  in the knowledge-based Shannon rate in Eqn. (54). The parameter  $\mu_K \in \{1.0, 0.4, 0.1\}$  results in a test accuracy of  $\{71.4\%, 91.2\%, 52.8\%\}$ , which is the maximum test accuracy obtained by directly solving the scheduling problem based on Eqn. (54) and can be used to informally indicate the knowledge quality. Thus,  $\mu_K = 0.4$  represents the best knowledge quality, whereas  $\mu_K = 0.1$  is the worst.

Now we list the settings for training. The neural network has three hidden layers with 512, 1024 and 512 neurons, respectively, followed by ReLU activations. The network is initialized based on Algorithm 1. The training is performed by the Adam optimizer with learning rate  $10^{-5}$  for 2000 steps on a HPC cluster with GPU type P100. We use 5 random seeds for each setting to evaluate the performance error.

**Results.** The results, including the test accuracy and the test sum rate under different knowledge quality, numbers of labels and weights  $\lambda$ , are shown in Fig. 2 and Fig. 3. The test sum rate is the (pseudo) real sum rate defined in Eqn. (58) with  $\mu_R = 0.5$ . We can find that the sum rate expectedly increases if the test accuracy increases. From Fig. 2(a) and Fig. 3(a), we see that if the domain knowledge quality is only 52.8% (i.e.,  $\mu_K = 0.1$ ), it has bad effects on learning when labels are enough. Nevertheless, it still benefits the performance when there are only 100 labels and, if we place a less weight on the knowledge-based risk, the accuracy and sum rate is higher.

If the knowledge quality is 71.4% (i.e.,  $\mu_K = 1.0$ ), as shown in Fig. 2(b) and Fig. 3(b), the domain knowledge has significant

benefits when there are only 100 labeled samples. When there are 500 labeled samples, the domain knowledge and labels complement each other and get a better performance than pure label-based and knowledge-based learning. When the number of labeled samples is even higher and reaches 1000, the integration of domain knowledge cannot benefit the learning further. In Fig. 2(c) and Fig. 3(c), when the domain knowledge quality further improves, we can see that the domain knowledge can still bring benefits even in the presence of 1000 labeled samples.

From these results, we see that labels and domain knowledge can complement each other. The domain knowledge plays an important role when labels are relatively scarce, while labels, even only a few, help improve the learning performance when domain knowledge has a low quality. Additionally, it is important to achieve a balance between label-based supervision and knowledge-based supervision. In general, we place more weight on the knowledge-based risk if knowledge quality is good enough and the number of labels is small, and vice versa.