

Learning an Explicit Hyper-parameter Prediction Function Conditioned on Tasks

Jun Shu

XJTUSHUJUN@GMAIL.COM

School of Mathematics and Statistics and Ministry of Education Key Lab of Intelligent Networks and Network Security, Xi'an Jiaotong University, Xi'an, Shaan'xi Province, P. R. China
Pazhou Lab (Huangpu), Guangzhou, Guangdong Province, P. R. China

Deyu Meng*

DYMENG@MAIL.XJTU.EDU.CN

School of Mathematics and Statistics and Ministry of Education Key Lab of Intelligent Networks and Network Security, Xi'an Jiaotong University, Xi'an, Shaan'xi Province, P. R. China
Pazhou Lab (Huangpu), Guangzhou, Guangdong Province, P. R. China
Macau Institute of Systems Engineering, Macau University of Science and Technology
Taipa, Macau, P. R. China.

Zongben Xu

ZBXU@MAIL.XJTU.EDU.CN

School of Mathematics and Statistics and Ministry of Education Key Lab of Intelligent Networks and Network Security, Xi'an Jiaotong University, Xi'an, Shaan'xi Province, P. R. China
Pazhou Lab (Huangpu), Guangzhou, Guangdong Province, P. R. China

Editor: Samuel Kaski

Abstract

Meta learning has attracted much attention recently in machine learning community. Contrary to conventional machine learning aiming to learn inherent prediction rules to predict labels for new query data, meta learning aims to learn the learning methodology for machine learning from observed tasks, so as to generalize to new query tasks by leveraging the meta-learned learning methodology. In this study, we achieve such learning methodology by learning an explicit hyper-parameter prediction function shared by all training tasks, and we call this learning process as *Simulating Learning Methodology* (SLeM). Specifically, this function is represented as a parameterized function called meta-learner, mapping from a training/test task to its suitable hyper-parameter setting, extracted from a pre-specified function set called meta learning machine. Such setting guarantees that the meta-learned learning methodology is able to flexibly fit diverse query tasks, instead of only obtaining fixed hyper-parameters by many current meta learning methods, with less adaptability to query task's variations. Such understanding of meta learning also makes it easily succeed from traditional learning theory for analyzing its generalization bounds with general losses/tasks/models. The theory naturally leads to some feasible controlling strategies for ameliorating the quality of the extracted meta-learner, verified to be able to finely ameliorate its generalization capability in some typical meta learning applications, including few-shot regression, few-shot classification and domain generalization. The source code of our method is released at <https://github.com/xjtushujun/SLeM-Theory>.

Keywords: Meta learning, simulating learning methodology, statistical learning theory, few-shot learning, domain generalization, structural risk minimization, meta-regularization

*. Corresponding author

1 Introduction

The core goal of machine learning is to learn the inherent generalization rule underlying data from some empirical observations, so as to make label predictions for new query samples. Such a generalization rule is generally modeled as a parameterized function (i.e., learner), and extracted from a pre-specified function set (i.e., learning machine) (Jordan and Mitchell, 2015). In the recent decade, deep learning approaches, equipped with highly parameterized learning machine (deep neural network architectures), have made great successes in a variety of fields (He et al., 2016; Silver et al., 2016; Devlin et al., 2019), strongly substantiating the validity of this learning framework.

However, nowadays gradually more deficiencies have been emerging for this conventional machine learning framework. On the one hand, its success largely relies on vast quantities of pre-collected annotated data, and simultaneously huge computation resources. However, most applications in real world have intrinsically rare or expensive data, or limited computation resources. This inclines to largely degenerate the capability of conventional machine learning, especially deep learning, expected by general users. On the other hand, current deep learning methods are always designed with complicated architectures and possess huge amounts of hyper-parameters, making them easily trapped into the overfitting issues, and possibly perform poorly on the test domain.

The above limitations can be mainly attributed to highly complicated hyper-parametric configurations involved in almost every single component of the learning process in machine learning (Jordan and Mitchell, 2015), e.g., data screening, model constructing, loss function presetting and algorithm designing, etc (see Section 2.1 and Table 1), for handling a practical learning task. The conventional assumption is that these hyperparameters are pre-specified when executing machine learning algorithm. However, the specification of hyperparameters can drastically affect performance measures like accuracy or data efficiency. Therefore, seeking a hyper-parameter configuration where the machine learning algorithm will produce a model that generalizes well to new data attracts much recent research attention.

Conventional machine learning literatures have raised many elegant approaches for such hyper-parameter selection issue, like Akaike Information Criterion (AIC) (Akaike, 1974), Bayesian Information Criterion (BIC) (Schwarz et al., 1978), Minimum Description Length (MDL) (Barron and Cover, 1991), and validation set based approaches (Stone, 1974). Their availability, however, is mostly restricted to the problems with relatively small-scale hyper-parametric structures. When facing complicated and massive hyper-parametric configurations, especially those related to a deep neural network, these concise manners are generally incapable of taking effect. Instead, in most cases the task still highly relies on manual attempts by human experts to the problem. By deeply understanding the problem and accumulating heuristic experience of hyper-parameter tuning, such human-designed manner does be able to make effect to specific tasks. However, when the task is with dynamic variations, the learning method is always required to be re-designed from scratch based on the new understandings of humans to the varying task. Especially, a learning method with carefully modulated hyper-parameters might be excessively good to the investigated learning tasks, while hardly to be readily generalized to a new query task with insightful correlations but also evident variations, like data modalities, network architectures, loss formulations and utilized algorithms, with the trained tasks. Then such laborious process

Table 1: Taxonomy of some typical recent literatures on meta learning based on their specified hyper-parameters to learn.

Taxonomy	Hyper-parameters to learn in the learning process
Data collection	data simulator (Ruiz et al., 2018), dataset distillation (Wang et al., 2018), instance weights (Shu et al., 2019, 2020a, 2023b,a), label corrector (Wu et al., 2021; Zheng et al., 2021), exploration policy (Xu et al., 2018; Garcia and Thomas, 2019), noise generator (Madaan et al., 2020) data annotator (Konyushkova et al., 2017), data augmentation policy (Cubuk et al., 2019)
Model construct	neural architecture (Zoph and Le, 2017; Liu et al., 2019), activation function (Ramachandran et al., 2017), feature modulation function (Ryu et al., 2020; Wang et al., 2020c), neural modules (Alet et al., 2019) dropout (Lee et al., 2020), neural processes (Garnelo et al., 2018; Requeima et al., 2019), attention (Kanika et al., 2021), batch normalization (Bronskill et al., 2020; Yingjun et al., 2021)
Loss function preset	loss predictor (Houthoofd et al., 2018; Huang et al., 2019; Gonzalez and Miikkulainen, 2020) metric (Sung et al., 2018; Lee and Choi, 2018), auxiliary loss (Li et al., 2019b; Veeriah et al., 2019), critic predictor (Sung et al., 2017), regularization (Balaji et al., 2018; Denevi et al., 2020), robust loss (Shu et al., 2020b,a; Ding et al., 2023; Rui et al., 2022)
Algorithm design	gradient generator (Andrychowicz et al., 2016; Wichrowska et al., 2017; Ravi and Larochelle, 2017) initialization (Finn et al., 2017; Fakoor et al., 2020; Song et al., 2020; Wang et al., 2020d), preconditioning matrix (Flennerhag et al., 2020), curvature (Park and Oliva, 2019), learning rate schedule (Li et al., 2017b; Shu et al., 2022), minimax optimizer (Shen et al., 2021)

for hyper-parameter tuning has to be started again, which then inclines to result in most core issues encountered by current machine learning. It has been gradually more widely recognized nowadays that appropriately specifying these hyper-parameters contained in the entire learning process of machine learning algorithm has become significantly more difficult, time-consuming and laborious than directly running a well designed learning process itself for a machine learning algorithm (Hospedales et al., 2020).

1.1 Background of Meta-learning Research

Meta learning, or learning to learn, seeks to improve conventional machine learning by nesting two search problems (Hospedales et al., 2020; Franceschi et al., 2018): at the inner level (or with-task level) we seek a good task-specific model (as in standard machine learning), while at the outer level (or meta level) we seek a good hyper-parameters configuration rather than assuming it is pre-specified and fixed, that ensures the produced model at the inner level to generalize well. The realization strategy is to learn the common hyper-parameter specification principle among a set of training tasks, instead of training samples as conventional machine learning. The aim is to get the hyper-parameter setting rule shared by training tasks, which is thus expected to improve future task learning performance. Such a learning manner can lead to a variety of benefits such as improved data and compute efficiency of machine learning, and being better aligned with human learning, that can learn new concepts quickly from few examples (Hospedales et al., 2020). In the recent years, many meta learning studies have been raised, constructed on specific problems, e.g., AutoML (Yao et al., 2018) and algorithm selection (Vanschoren, 2018), or particular applications, e.g., few-shot learning (Finn et al., 2017; Wang et al., 2020e; Shu et al., 2018), neural architecture search (NAS) (Elsken et al., 2019), or hyper-parameter optimization (Franceschi et al., 2018). Typical researches along this line are listed in Table 1.

Recent studies on meta learning, however, still have two major issues. Firstly, many of current meta learning works put emphasis on learning fixed hyper-parameters shared from the training tasks, and then directly applying them to adapt to new query tasks. Typical

works include MAML (Finn et al., 2017) and its variants (Li et al., 2017b; Nichol et al., 2018; Antoniou et al., 2019; Finn et al., 2019), AutoML (Yao et al., 2018), hyper-parameter optimization (Franceschi et al., 2018), and also the meta learning framework summarized in the recent excellent survey work in (Hospedales et al., 2020). Relying on the shared hyper-parameters is challenging for complex task distributions (e.g., those with distribution shift), since different tasks may require to set substantially different hyper-parameters. This makes it always infeasible to find a set of common hyper-parameters performable for all tasks, which may fail to adapt to heterogenous environments of tasks.

Secondly, although some solid theoretical justifications have been presented along this research line, most existing meta learning theories are developed under conventional machine learning framework and put emphasis on evaluating the generalization capability of the traditional learning model (i.e. the learner). And the most theoretical results are seldom well-aligned with the current meta learning practice with innovatory support/query episodic training mode (Vinyals et al., 2016; Finn et al., 2017)¹. More importantly, the current theoretical results have less mentioned theory-inspired controlling strategies for meta-learner to improve the capability of meta learning, and thus are less functioned to feedback meta learning models for helping improve their practical generalization performance.

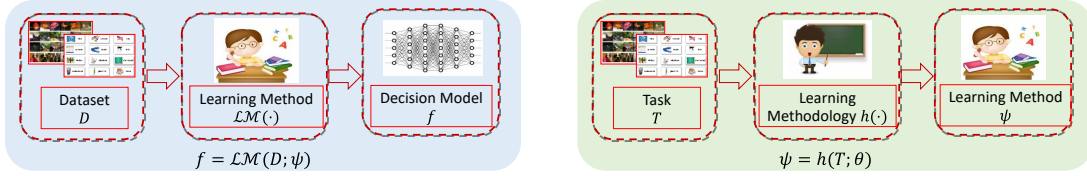
1.2 Our Contributions

To alleviate the aforementioned deficiencies of current meta learning, we can revisit the hyper-parameter setting for machine learning. In fact, if we take the entire learning process as an implicit function² mapping from an input training dataset to a decision model (i.e., a learner extracted from the learning machine by regular machine learning), as shown in Fig.1(a), these hyper-parameters involved in the entire learning process then constitute the parameters of this function. Note that such hyper-parametric configurations are involved in almost every single component of the learning process in machine learning (Jordan and Mitchell, 2015), e.g., data screening, model constructing, loss function presetting and algorithm designing, etc (see Section 2.1 and Table 1), containing broader essence compared with simple hyper-parameters needed to be tuned in conventional machine learning, like regularization strength or learning rate. From this perspective, we can interpret such hyper-parameter configuration as the “learning method” of the learning task at hand, since their proper settings essentially determine the ultimate capability, especially the generalization, of the extracted learner by this function through implementing the corresponding learning process equipped with the hyper-parameter configuration.

Compared with many current meta learning methods aiming to learn hyper-parameter configurations themselves, we propose to learn an explicit hyper-parameter setting function for predicting proper hyper-parameter configurations when adapted to new query tasks, as shown in Fig. 1(b). Specifically, equipped with such an explicit function, a novel machine

1. More details are presented in Section 3.

2. In many conventional literatures on learning theory, this implicit function is often called a learning algorithm (Maurer and Jaakkola, 2005; Chen et al., 2020). Yet it intrinsically represents the whole learning process from input data to output learner. In current machine learning, this process should contain more general learning component settings besides the pre-specified learning algorithm itself, like the designing of network architecture and loss function. To avoid possible clutters of readers, in this paper we call it as the function of learning method, or $\mathcal{LM}(\cdot)$ briefly.



(a) A learning method producing a decision model (b) A learning methodology producing a learning method from a query task.

Figure 1: (a) depicts the executive process when we deal with a machine learning problem. The learning method can be seen as an implicit function mapping from an input training dataset to a decision model. This function represents the entire learning process equipped with proper hyper-parameter configurations, guaranteeing the machine learning system to produce a decision function by executing it. (b) is the principle of how to determine the learning method. There exists a common learning methodology among various learning tasks, which can be seen as a function that maps an input task to the corresponding learning method. We propose a SLeM framework to learn the learning methodology function. Once we achieve the function, it can allocate a proper learning method for a specific query task. The machine learning system can then automatically produce the decision model for the task by running learning pipeline (a), equipped with the learning method produced by the learning methodology function.

learning system can automatically produce proper learning method for different tasks without need of much extra human intervention. More strictly speaking, what we want is a mapping from the learning task space to the hyper-parameter space covering the whole learning process. From this perspective, learning an automatic hyper-parameter setting function conditioned on tasks is expected to get the “learning methodology” shared among various learning tasks, and directly be used to allocate learning method for new tasks. In this way, such a “learning methodology function” is hopeful to be employed to finely adapt varying query tasks from heterogeneous environment with less computation/data costs, as well as fewer human interventions (Biggs, 1985; Schrier, 1984; Schmidhuber, 1987; Thrun and Pratt, 2012). We call this learning process as *Simulating Learning Methodology* (SLeM).

Formally, we can study the statistical guarantee of our novel SLeM framework to understand the generalization capability of such hyper-parameter prediction function. Especially, just succeeded from the conventional machine learning framework, this function is mathematically modelled as a parameterized function, called meta-learner, mapping from a learning task to its proper hyper-parameter configuration, and extracted from a pre-specified function set, called meta learning machine. On the one hand, such meta-learner setting facilitates a better flexibility of the learned hyper-parameters adaptable to diverse query tasks than fixed hyper-parameters. Specifically, the extracted meta-learner can be readily used to produce proper hyper-parameters conditioned on new query tasks. On the other hand, such SLeM framework with explicit hyper-parameter prediction function conditioned on learning tasks can be seen as a substantial but homologous extension from the conventional machine learning framework imposed on the pre-specified learning machine, and the corresponding statistical learning theory can be subsequently derived. Especially, similar to the structural risk minimization (SRM) principle in the conventional statistical learning theory (Vapnik, 1999; Shawe-Taylor et al., 1998), some beneficial theoretical results can then be achieved for revealing and controlling the intrinsic generalization capability of the preset meta learning machine. In summary, the main contributions of this work are as follows:

(1) We propose a new meta learning method aiming at learning an explicit hyper-parameter prediction function conditioned on learning tasks. This formulation provides a general framework to understand meta learning, and reveals that the essential character of meta learning is to construct a meta-learner for simulating the learning methodology, and transferably use it to help fulfill new query tasks.

(2) We introduce a problem-agnostic definition of meta-learner to realize learning the hyper-parameter prediction function conditioned on learning tasks. The parameterized meta-learner can be easily integrated into the traditional machine learning framework to provide a fresh understanding and extension of the original machine learning framework. We further provide generalization bounds for the new SLeM framework with general losses, tasks, and models.

(3) We provide general-purpose bounds for SLeM with decoupling the complexity of learning the task-specific learner from that of learning the task-transferrable meta-learner. The meta-learner is extracted from a pre-specified function set (i.e., meta learning machine), and the theoretical results facilitate some feasible controlling strategies on the meta-learner, capable of being easily embedded into current off-the-shelf meta learning programs and yet helping generally improve its generalization capability.

(4) We highlight the utility of our SLeM framework for obtaining the learning guarantees of some typical meta learning applications, including few-shot regression, few-shot classification and domain generalization. The theory-induced meta-regularization control effects of the meta-learner are empirically verified to be effective for consistently improving its generalization capability on new query tasks.

The rest of the paper is organized as follows. The proposed SLeM meta learning framework is formulated in Section 2, and in Section 3, we derive the generalization bounds of this meta learning manner during the meta-training and meta-test stages, respectively. Related works are introduced in Section 4. We further instantiate our general theoretical framework on few-shot regression in Section 5, few-shot classification in Section 6, and domain generalization in Section 7. Then Section 8 discusses the online version of our method and its theoretical rationality. The conclusion and discussion are finally made.

2 Exploring a Task-Transferable Meta-learner for Meta-learning

For clarity, in the following we focus on supervised learning. We firstly recall the conventional machine learning framework, and then discuss the deficiencies of existing meta learning methods. Finally, we present the proposed SLeM meta-learning framework and methods for addressing these deficiencies, which can be seen as a substantial but homologous extension from the conventional machine learning framework.

Let's first introduce some necessary notations. \mathbf{X}^T denotes the transpose of a matrix \mathbf{X} . Let $\mathcal{D} = \mathcal{X} \times \mathcal{Y}$ be the data space, where $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} \subset \mathbb{R}$ (regression) or $\mathcal{Y} = \{0, 1, \dots, K-1\}$ (multi-class classification) are the input and output spaces, respectively. We use the bracketed notation $[k] = \{1, 2, \dots, k\}$ as shorthand for index sets. The norm $\|\cdot\|$ appearing on a vector or a matrix refers to its ℓ_2 norm or Frobenius norm.

2.1 Machine Learning: Learning a Learner from Learning Machine

Machine learning (Jordan and Mitchell, 2015) aims to extract a learner (decision model) from a pre-specified learning machine based on a set of training observations. It generally includes the following ingredients.

Training dataset. It is usually composed of a finite paired dataset $D = \{(x_i, y_i), i \in [n]\}$ drawn i.i.d. from a task (probability distribution) μ , which simulates the input and output of the learner to be estimated.

Learner and learning machine. The learner corresponds to a mapping $f(x; \omega) : \mathcal{X} \rightarrow \mathcal{Y}$, which is requested to produce a prediction rule from $x \in \mathcal{X}$ to its label $y \in \mathcal{Y}$, and $\omega \in \Omega$ represents the parameters of f to be estimated. f is chosen from a preset learning machine (i.e., hypothesis space) \mathcal{F} , constituting a set of candidate learners.

Performance measure. The performance is measured by a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, which reflects the extent of how well the learner fits training data. The expected risk $R_\mu(f)$ and empirical risk $\hat{R}_D(f)$ are respectively defined as:

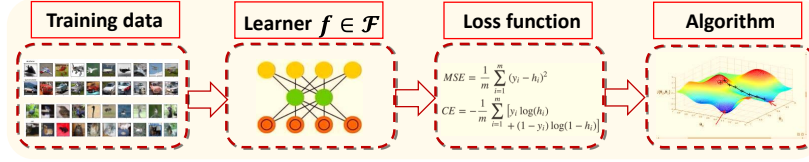
$$R_\mu(f) = \mathbb{E}_{(x,y) \sim \mu} \ell(f(x), y), \quad \hat{R}_D(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i). \quad (1)$$

Optimization algorithm. The learning algorithm A is employed from optimization toolkits to extract a learner f from \mathcal{F} guided by the performance measure. E.g., the commonly used algorithm for deep learning methods is SGD or Adam (Goodfellow et al., 2016).

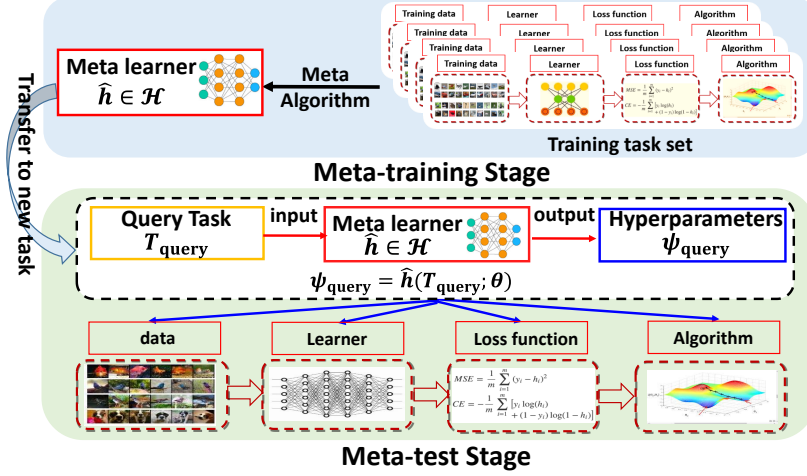
The overall learning process of machine learning is shown in Fig.2(a). Actually, before executing machine learning for a given learning task in practice, it needs to firstly determine complicated hyper-parametric configurations involved in all components of the learning process. Then the whole machine learning process can be executed to produce the decision model. If we take the entire learning process as an implicit function mapping from an input training dataset to a decision model, all involved hyper-parameters then constitute the parameters of this function. We can represent this function as $\mathcal{LM}(D; \psi) : \mathcal{D} \rightarrow \mathcal{F}$. We denote ψ as all hyper-parameters involved in the learning method³. For example, we can write $\psi = (\psi_D, \psi_f, \psi_\ell, \psi_A)$ to record hyper-parameters in the learning configurations with respect to the data D , the learner f , the loss function ℓ and the optimization algorithm A , respectively. The ultimate capability, especially the generalization, of the decision model then essentially depends on the proper hyper-parameter settings ψ .

After the learning process outputs the decision model, it can produce the label predictions for new query samples. Current great successes of deep learning have strongly substantiated the validity of this learning framework for summarizing the underlying label prediction rules from data. Assume that the optimal learner underlying the task μ is $f^* = \arg \min_{f \in \mathcal{F}} R_\mu(f)$, and the estimated learner through minimizing the empirical risk from the training data D is $\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}_D(f)$. Then we generally use $R_\mu(\hat{f}) - R_\mu(f^*)$ to measure the ‘‘closeness’’ between the estimated \hat{f} and the optimal f^* , whose upper bound has been proved as follows (Mohri et al., 2018):

3. In practice, the hyper-parameters to be learned by meta learning generally only contain a certain subset of hyper-parameters among all involved ones in the learning process. Some typical hyper-parameters meta learning aims to learn are listed in Table 1.



(a) Machine learning framework.



(b) Task-Transferable Meta-learner for Meta Learning.

Figure 2: (a) A general machine learning framework. (b) The meta-training and meta-test stages of the proposed meta learning framework.

Theorem 1 Let $\mathcal{D} = \mathcal{X} \times \mathcal{Y}$ be the data space, hypothesis (learning machine) \mathcal{F} be the function class of the mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$, and $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, B]$ be the loss function. Assume that the loss $\ell(\cdot, y)$ is L -Lipschitz for any $y \in \mathcal{Y}$. Then for any $\delta > 0$, with probability at least $1 - \delta$, we have

$$R_\mu(\hat{f}) - R_\mu(f^*) \leq 6L\mathcal{G}_n(\mathcal{F}) + 2B\sqrt{\frac{\ln(1/\delta)}{2n}}, \quad (2)$$

where $\mathcal{G}_n(\mathcal{F}) := \mathbb{E}_{D \sim \mu^n}[\hat{\mathcal{G}}_D(\mathcal{F})]$ is the Gaussian complexity, and $\hat{\mathcal{G}}_D(\mathcal{F})$ is the empirical Gaussian complexity of \mathcal{F} w.r.t D defined as

$$\hat{\mathcal{G}}_D(\mathcal{F}) := \mathbb{E}_{\mathbf{g}} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n g_i f(x_i) \right], \quad g_i \sim \mathcal{N}(0, 1)$$

where $\mathcal{N}(0, 1)$ is the Gaussian distribution with zero mean and unit variance.

Since $\mathcal{G}_n(\mathcal{F}) \sim \sqrt{C(\mathcal{F})/n}$ (Mohri et al., 2018), where $C(\cdot)$ measures the intrinsic complexity of the function class (e.g., VC dimension), Theorem 1 delivers the information that the distance between \hat{f} and f^* will decrease as n becomes large for a given learning machine \mathcal{F} . In other word, given enough data, model \hat{f} obtained by minimizing $\hat{R}_D(f)$ is arbitrarily close to optimal model f^* .

2.2 Deficiencies of Existing Meta Learning Methods

Despite their success, the effect of contemporary machine learning models, especially deep learning models, largely relies on vast quantities of pre-collected annotated data, and

simultaneously huge computation resources. This inclines to significantly degenerate the capability of machine learning for real-world applications with intrinsically rare or expensive data, or limited computation resources. Besides, deep learning methods with complicated architectures and huge amounts of hyper-parameters tend to easily suffer from the overfitting issue, and possibly perform poorly on the test domain.

Recently, meta learning methods are proposed to improve conventional machine learning by distilling the common hyper-parameter specification rule among a set of training tasks, and then using this rule to improve future task learning performance. It can be described as a nested optimization scheme: at the with-task level, an inner algorithm performs task-specific machine learning algorithm with the current hyper-parameter configuration ψ , i.e., $\mathcal{LM}(D; \psi)$, where D is the task-specific training dataset. And at the meta-level, a meta-algorithm updates the aforementioned hyper-parameter configuration ψ by leveraging the experience accumulated from the tasks observed so far. Such a learning manner can lead to a variety of benefits such as improved data and compute efficiency of machine learning, and it is better aligned with human learning, that can learn new concepts quickly from few examples (Hospedales et al., 2020). To achieve this goal, recently a large quantity of meta learning methods have been raised, that have helped machine learning improve the data efficiency (Shu et al., 2018), algorithm automation (Andrychowicz et al., 2016; Liu et al., 2019; Cubuk et al., 2019; Shu et al., 2023b, 2020a), and generalization (Finn et al., 2017; Tripuraneni et al., 2020). Typical researches along this line are listed in Table 1.

Despite their success, most meta learning methods may fail to adapt to heterogenous environments of tasks, in which the complexity of the tasks’ distribution cannot be captured by a single fixed hyper-parameter configuration shared from the training tasks. In fact, heterogenous tasks may require to substantially different hyper-parameters. To address this limitation, Denevi et al. (2020, 2022); Wang et al. (2020d) propose conditional meta learning aiming at learning a conditioning function that maps a task’s dataset into a hyper-parameter vector that is appropriate for the task at hand. While they are mainly delicately designed for specific problems (e.g., task conditional model initialization, biased regularization and linear representation, etc), and their conditioning functions are simply set as linear functions, which may limit the applicability to other meta-learning problems.

From the theoretical perspective, most existing meta learning theories along above research line are developed under conventional machine learning framework and put emphasis on evaluating the generalization capability of the traditional learning model (i.e. the learner). And the most theoretical results are seldom well-aligned with the current meta learning practice with innovatory support/query episodic training mode (Vinyals et al., 2016; Finn et al., 2017). Besides, we notice that the regularization techniques have been effective in improving capability of machine learning, which often impose controlling strategies on the task-specific model (learner) developed from statistical learning theory via structural risk minimization principle. However, in the meta learning framework, most previous theoretical works have seldom emphasized theory-inspired meta-regularization strategies for meta-learner to improve the capability of meta learning, and thus are less functioned to feedback meta learning models for helping improve their practical generalization performance.

2.3 Proposed Method: Learning a Meta-learner from Meta-learning Machine

To address the deficiencies of existing meta learning methods, as shown in Fig.1(b), we propose to learn an explicit hyper-parameter prediction function for predicting proper hyper-parameter configurations when adapted to new query tasks. We can mathematically model the function as $h(T; \theta) : \mathcal{T} \rightarrow \Psi$, called the meta-learner, which maps from the learning task space \mathcal{T} to the hyper-parameter space Ψ . Then the extracted meta-learner can be readily used to produce proper hyper-parameter configurations conditioned on new query tasks. Therefore, the goal of our novel SLeM meta learning method is to determine the parameter $\theta \in \Theta$ contained in the meta-learner h , instead of the hyper-parameter configuration ψ itself like many previous meta learning methods, e.g., weights initialization (Finn et al., 2017; Nichol et al., 2018; Rusu et al., 2019; Antoniou et al., 2019; Finn et al., 2019), learning rate (Li et al., 2017b; Baydin et al., 2019), and other hyper-parameters (Franceschi et al., 2018). These works have been comprehensively introduced in the recent excellent survey paper (Hospedales et al., 2020). Relying on the shared and unique hyper-parameters is challenging for adapting complex heterogenous task distributions, especially those with distribution shift. Comparatively, our hyper-parameter prediction function facilitates a better flexibility of the learned hyper-parameters adaptable to diverse query tasks, so as to effectively address the issue of single fixed hyper-parameters. In fact, when we set the meta-learner as a constant function, these hyper-parameter-fixed approaches can be categorized as special cases of this more general SLeM meta learning framework.

Different from the learner $f(x; \omega)$ in machine learning, attempting to extract the label prediction rule among data, the meta-learner $h(T; \theta)$ is requested to extract the hyper-parameter prediction rule for learning process shared by training tasks, which is then expected to finely generalize to new query tasks to specify their hyper-parameters for a learning method. Nevertheless, such meta learning framework with explicit hyper-parameter prediction function conditioned on learning tasks can be seen as a substantial but homologous extension from the conventional machine learning framework imposed on the pre-specified learning machine. To achieve such a meta-learner in practice, we usually assume access to the following ingredients corresponding to the conventional machine learning in Section 2.1:

Training task set. It usually assumes access to a set of source tasks $\Gamma = \{\mu_t = \{\mu_t^s, \mu_t^q\}, t \in [T]\}$, to learn the meta-learner mapping, where μ_t denotes the probability distribution for the t -th task and i.i.d. sampled from an environment distribution η (Baxter, 2000). Here μ_t^s and μ_t^q denote the distributions sampling support/training set D_t^{tr} and query/validation dataset D_t^{val} for the task, respectively. This follows the support/query meta learning setting proposed by (Vinyals et al., 2016). Specifically, for the t -th task, it is composed of $D_t = (D_t^{tr}, D_t^{val})$, where $D_t^{tr} = \{(x_{ti}^{(s)}, y_{ti}^{(s)})\}_{i=1}^{m_t} \sim (\mu_t^s)^{m_t}$, $D_t^{val} = \{(x_{tj}^{(q)}, y_{tj}^{(q)})\}_{j=1}^{n_t} \sim (\mu_t^q)^{n_t}$, and m_t and n_t are the sizes of training set and validation set for the t -th task, respectively. We denote $\mathbf{D} = \{D_t\}_{t=1}^T$ as the entire training task dataset.

Meta-learner and meta learning machine. Formally, the meta-learner $h(T; \theta) : \mathcal{T} \rightarrow \Psi$ is required to extract the rule from $T \in \mathcal{T}$ to its hyper-parameter setting $\psi \in \Psi$ for the learning method, and $\theta \in \Theta$ is the parameter contained in h to be estimated. h is selected from a pre-designed parameterized function set \mathcal{H} , called meta learning machine, constituting a function class of candidate meta-learners h .

Table 2: Illustration of the parametrization of tasks and the hyper-parameter prediction function from some typical examples of existing meta learning methods.

Methods	Specified hyper-parameters to learn	Parametrization of tasks	Parametrization of hyper-parameter prediction function
Conditional meta learning Denevi et al. (2020)	Bias vector	Task side information	Linear functions
MetaOptNet Lee et al. (2019)	Feature representation	Sample original input	Deep neural networks
R2D2 Bertinetto et al. (2019)	Feature representation	Sample original input	Deep neural networks
MMAML Vuorio et al. (2019)	Task modulation vector	Task embedding	MLP
LSTM optimizer Andrychowicz et al. (2016)	Dynamic update rules of optimizer	Gradient information	LSTM
MW-Net Shu et al. (2019)	Sample weights	Sample loss information	MLP
Feature-Critic Li et al. (2019b)	Loss function	Feature information	MLP
LEO Rusu et al. (2019)	Model adapted parameters	Training dataset	Encoding & decoding network

Generally, the input of a meta-learner h is a representation conveyed by the task $T \in \mathcal{T}$. It can be some static information for task, e.g., directly obtained from the provided dataset D for the task, or dynamic knowledge extracted from the learning process for handling task T , e.g., gradient information (Andrychowicz et al., 2016), sample loss information (Shu et al., 2019), feature information (Li et al., 2019b), etc. Table 2 shows some typical examples of existing meta learning methods to illustrate the practical parametrization of task representations and the hyper-parameter prediction functions. It is seen that conditional meta learning (Denevi et al., 2020, 2022) could be a special case under our general meta-learner formulation, where the meta-learner is chose from linear functions class. In this paper, we mainly focus on the mathematical formulation and the statistical learning theory to understand the proposed meta learning framework, and thus we do not especially emphasize a specific practical parametrized hyper-parameter setting function. For simplicity and convenience, we formulate the hyper-parameter prediction function as $h(D; \theta)$ in our study, where we assume the representation of task is directly obtained from the provided dataset D for the task at hand, and the hyper-parameter prediction function h is selected from a pre-designed parameterized function set \mathcal{H} . For simplicity, we often briefly denote it as h .

Performance measure. The goal of SLeM is to learn a common hyper-parameter prediction function suitable for a family of learning tasks. For intuitive understanding, we rewrite the empirical risk $\hat{R}_D(f)$ as $\mathbf{L}(f, D) : \mathcal{F} \times \mathcal{D} \rightarrow \mathbb{R}^+$ to characterize the performance of the decision model on the task. Our aim is then to seek the best hyper-parameter prediction function h , which is employed to set the with-task level learning process on D^{tr} , so that the obtained learner can achieve optimal generalization performance on D^{val} . Specifically, the risk so described can be formally written as:

$$R_\eta(h) = \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D^{val} \sim (\mu^q)^n} \mathbb{E}_{D^{tr} \sim (\mu^s)^m} \mathbf{L}(\mathcal{LM}(D^{tr}; h), D^{val}), \quad (3)$$

where $\mu = (\mu^s, \mu^q)$ is the task distribution sampled from an environment distribution η , and D^{tr} and D^{val} are training and validation sets i.i.d. generated from μ^s and μ^q , respectively.

In practice, we have access to a set of training task set \mathbf{D} , and the meta-learner is often learned by (approximately) minimizing the following expected empirical risk on the meta learning task:

$$\hat{R}_D(h) = \frac{1}{T} \sum_{t=1}^T \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val}). \quad (4)$$

Optimization algorithm. The goal then is to achieve a meta-learner h from \mathcal{H} by minimizing the empirical task average risk $\hat{R}_D(h)$. There are multiple effective optimization

Table 3: Comparison of the main notations used in conventional machine learning and the proposed meta learning frameworks.

Compared concepts	Machine Learning	Meta Learning
Source distribution	Task distribution μ	Environment η
Training instance	$D = \{(x_i, y_i)\}_{i=1}^n \sim \mu^n$	$D = \{D_t = (D_t^{tr}, D_t^{val})\}_{t=1}^T, D_t^{tr} \sim (\mu_t^s)^{m_t}, D_t^{val} \sim (\mu_t^q)^{n_t}, \mu_t = (\mu_t^s, \mu_t^q) \sim \eta$
Test instance	$(x, y) \sim \mu$	$D_\mu^{tr} \sim (\mu^s)^{m_\mu}, (x, y) \in \mu^q, \mu = (\mu^s, \mu^q) \sim \eta$
Learning objective	Learner $f : \mathcal{X} \rightarrow \mathcal{Y}$	Meta-learner $h : \mathcal{T} \rightarrow \Psi$
Output of the (meta)-learner	$y = f(\mathbf{x}; \omega)$	$\psi = h(T; \theta)$
Performance measure	$R_\mu(f) = \mathbb{E}_{(x,y) \sim \mu} \ell(f(x), y)$	$R_\eta(h) = \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D^{val} \sim (\mu^q)^n} \mathbb{E}_{D^{tr} \sim (\mu^s)^m} \mathcal{L}(\mathcal{LM}(D^{tr}; h), D^{val})$
Goal of the generalization	Predicting label y for query sample x	Predicting learning method ψ for query task T

algorithms designed for solving the problem, e.g., the bilevel optimization techniques (Finn et al., 2017; Shu et al., 2019; Hospedales et al., 2020).

Contrary to conventional machine learning whose performance can be directly computed on the given annotated data, meta learning algorithm is relatively hard to directly compute the performance on given tasks. Generally, it needs a two-stage procedure to evaluate the performance of a meta-learner (Finn et al., 2017; Shu et al., 2019). Firstly, a meta-learner predicts the configuration of the hyper-parameters after observing the training set D_t^{tr} in each task of the task set \mathbf{D} . Then the machine learning system equipped with such hyper-parameter configurations can be executed to automatically produce the proper decision model for D_t^{tr} . Then in the second stage, the task loss can be computed as the average over all errors calculated on the validation set $D_t^{val}, t \in [T]$ with respect to the corresponding decision model obtained from D_t^{tr} , reflecting the hyper-parameter prediction capability of such meta-learner. This process actually is exactly expressed as the empirical task average risk $\hat{R}_{\mathbf{D}}(h)$ as formulated above. As shown in Fig.2(b), after achieving the meta-learner h , we can further transfer it to new query tasks, which can help set the learning process (i.e., parameter of learning method) for new query tasks.

The above SLeM meta learning framework is evidently succeeded from but also with evident difference with conventional machine learning framework in Section 2.1. To better clarify this point, we list the main notations of both frameworks to easily compare their main differences as introduced in Table 3. It should be noted that in the previous researches, meta learning is interpreted to improve performance by learning ‘how to learn’ (Thrun and Pratt, 2012). The main study of this work is expected to solidify such ‘learning to learn’ manner as learning an explicit hyper-parameter prediction function (i.e., a meta-learner) conditioned on learning tasks. From this view, a meta-learner helps formulate the learning method to tell the machine learning system how to automatically learn the decision model. Thus the insight of this formulation is to simulate humans to master the tricks for the learning methodology (thus we call it simulating learning methodology (SLeM)), and most importantly, further transfer it to help fulfill new tasks in the future. Therefore, achieving the meta-learner is hopeful to be employed to finely adapt varying query tasks from heterogeneous environment with less computation/data costs, as well as fewer human interventions (Biggs, 1985; Schrier, 1984; Schmidhuber, 1987; Thrun and Pratt, 2012).

Algorithm 1 Meta-Training: Learning the Methodology

Input: Training task set $\Gamma = \{D_t, t \in [T]\}$.

Do:

(1) Run algorithm to obtain $\hat{\mathbf{f}}^{(h)} = (\hat{f}_1^{(h)}, \dots, \hat{f}_T^{(h)})$, where $\hat{f}_t^{(h)}$ is obtained by calculating $\mathcal{LM}(D_t^{tr}; h(D_t^{tr}))$.

(2) Put $\hat{\mathbf{f}}^{(h)}$ into Eq.(4), and then minimize the error $\hat{R}_{\mathcal{D}}(\hat{\mathbf{f}}^{(h)})$ to obtain \hat{h} .

Return: The meta-learner \hat{h} .

3 Learning Theory

We use two phases to characterize the learning paradigm of the above introduced meta learning framework. The first is the meta-training phase, aiming to learn the hyper-parameter prediction function, and the second is the meta-test phase, purposing to generalize this function for setting learning methods on new query tasks. We firstly present some preliminaries, and then present learning theory results for meta-training and meta-test phases, respectively.

3.1 Preliminaries

We use the Gaussian complexity to measure the complexity of a function class. For a generic vector-valued function class \mathcal{E} , containing a set of functions $e : \mathbb{R}^d \rightarrow \mathbb{R}^r$. Given a data set $\mathbf{U} = \{\mathbf{u}_i \in \mathbb{R}^d, i \in [N]\}$ i.i.d. drawn from a task μ , the empirical Gaussian complexity is then defined as (Bartlett and Mendelson, 2002):

$$\hat{\mathcal{G}}_{\mathbf{U}}(\mathcal{E}) = \mathbb{E}_{\mathbf{g}} \left[\sup_{e \in \mathcal{E}} \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^r g_{ik} e_k(\mathbf{u}_i) \right], \quad g_{ik} \sim \mathcal{N}(0, 1),$$

where $\mathbf{g} = \{g_{ik}\}_{i \in [N], k \in [r]}$, and $e_k(\cdot)$ is the k -th coordinate of the vector-valued function $e(\cdot)$. The corresponding population Gaussian complexity is defined as $\mathcal{G}_N(\mathcal{E}) = \mathbb{E}_{\mathbf{U} \sim \mu^N} [\hat{\mathcal{G}}_{\mathbf{U}}(\mathcal{E})]$, where the expectation is taken over the distribution of \mathbf{U} . In comparison, for empirical Rademacher complexity $\hat{\mathfrak{R}}_{\mathbf{U}}(\mathcal{E})$, g_{ik} are replaced by uniform $\{-1, 1\}$ -distributed variables, and the corresponding population Rademacher complexity is $\mathfrak{R}_N(\mathcal{E}) = \mathbb{E}_{\mathbf{U} \sim \mu^N} [\hat{\mathfrak{R}}_{\mathbf{U}}(\mathcal{E})]$.

To prove the main learning theory results, we require the following assumptions, all being usually satisfied.

Assumption 1 (Bounded Inputs) $\mathcal{X} \subset \mathcal{B}(0, R)$, for $R > 0$, where $\mathcal{B}(0, R) = \{x \in \mathbb{R}^d : \|x\| \leq R\}$.

Assumption 2 (Bounded and Lipschitz Loss Function) The loss function $\ell(\cdot, \cdot)$ is B -bounded, and $\ell(\cdot, y)$ be L -Lipschitz for any $y \in \mathcal{Y}$.

3.2 Meta-Training Stage: Learning the Hyper-parameter Prediction Function

In the meta-training stage, the training task dataset $\mathcal{D} = \{D_t, t \in [T]\}$ is available for learning, where $D_t = (D_t^{tr}, D_t^{val})$, $D_t^{tr} \sim (\mu_t^s)^{m_t}$, $D_t^{val} \sim (\mu_t^q)^{n_t}$, $\mu_t = (\mu_t^s, \mu_t^q) \sim \eta$. We aim to learn the hyper-parameter prediction function from this task dataset. The overall learning

process is listed in Algorithm 1. The quantity of interest in meta-training is the expectation of task average generalization error as follows:

$$R_{train}(h) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^{n_t}} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val}). \quad (5)$$

Denote the optimal meta-learner obtained by minimizing the theoretical risk $R_\eta(h)$ and empirical risk \hat{R}_D in meta learning as:

$$h^* = \arg \min_{h \in \mathcal{H}} R_\eta(h); \quad \hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_D(h), \quad (6)$$

where $R_\eta(h)$ and \hat{R}_D are defined in Eqs. (3) and (4), respectively. Then we can naturally use $R_{train}(\hat{h}) - R_{train}(h^*)$ to measure the “closeness” between the estimated \hat{h} and the true underlying h^* , which captures the extent of how much the two meta-learners \hat{h} and h^* differ in aggregating over the T training tasks.

We can then present a theoretical upper bound estimation for this “closeness” measure. Before showing the theorem, we first introduce an important notation, $d_{\mathcal{F}}(\mu_t^s, \mu_t^q)$, necessary for the theoretical result. $d_{\mathcal{F}}(\mu_t^s, \mu_t^q)$ denotes the discrepancy divergence (Ben-David et al., 2010) between support and query data with respect to their sampled probability distributions μ_t^s and μ_t^q imposed on the hypothesis class \mathcal{F} :

$$d_{\mathcal{F}}(\mu_t^s, \mu_t^q) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^{n_t}} \mathbf{L}(f, D_t^{val}) - \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(f, D_t^{tr}) \right|.$$

Theorem 2 *If Assumptions 1 and 2 hold, for any $\delta > 0$, with probability at least $1 - \delta$, we have*

$$\begin{aligned} R_{train}(\hat{h}) - R_{train}(h^*) &\leq 768L \log(4 \sum_{t=1}^T n_t) \cdot L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) + \frac{6L}{T} \sum_{t=1}^T \hat{\mathcal{G}}_{D_t^{tr}}(\mathcal{F}) \\ &+ \frac{4}{T} \sum_{t=1}^T d_{\mathcal{F}}(\mu_t^s, \mu_t^q) + \frac{6B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t} \sqrt{\frac{\log \frac{2}{\delta}}{2}}} + \frac{6B}{T} \sum_{t=1}^T \sqrt{\frac{\log \frac{2}{\delta}}{m_t}} + \frac{12L \text{Dis}(D^{val})}{(\sum_{t=1}^T n_t)^2} \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}}, \end{aligned}$$

where $\text{Dis}(D^{val}) = \sup_{h, h'} \rho_{D^{val}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')})$, $\rho_{D^{val}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) = \frac{1}{\sum_{t=1}^T n_t} \sum_{t=1}^T \sum_{j=1}^{n_t} (f_t^{(h)}(x_{tj}^{(q)}) - f_t^{(h')}(x_{tj}^{(q)}))^2$, $D^{val} = \{D_t^{val}\}_{t=1}^T$, $\mathbf{f}^{(h)} = (f_1^{(h)}, f_2^{(h)}, \dots, f_T^{(h)})$, $f_t^{(h)} = \mathcal{LM}(D_t^{tr}; h)$, and $L(\mathcal{F})$ is the Lipschitz constant of $\mathbf{f}^{(h)}$ with respect to h , $\beta_t = \frac{n_t T}{\sum_{t=1}^T n_t}$.

The above bound for excess task average risk can be more concisely formulated by dominating its three main terms as follows:

$$R_{train}(\hat{h}) - R_{train}(h^*) \lesssim \tilde{\mathcal{O}} \left(\sqrt{\frac{C(\mathcal{H})}{\sum_{t=1}^T n_t}} + \frac{1}{T} \sum_{t=1}^T \sqrt{\frac{C(\mathcal{F})}{m_t}} + \frac{1}{T} \sum_{t=1}^T d_{\mathcal{F}}(\mu_t^s, \mu_t^q) \right).$$

They can be interpreted as: the complexity of learning the meta-learner h (the first term), the complexity of learning the task-specific learners \mathbf{f} (the second term), and the distribution

shift between support and query sets (the third term). As aforementioned, it generally holds that $\hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) \sim \sqrt{C(\mathcal{H})/\sum_{t=1}^T n_t}$, and $\hat{\mathcal{G}}_{D_t^{tr}}(\mathcal{F}) \sim \sqrt{C(\mathcal{F})/m_t}$, where $C(\cdot)$ measures the intrinsic complexity of the function class (e.g., VC dimension). Thus the first term on the right hand side above is of the order $\tilde{\mathcal{O}}(1/\sqrt{\sum_{t=1}^T n_t})$, and the second term above is of the order $1/T \sum_{t=1}^T \mathcal{O}(1/\sqrt{m_t}) \leq \mathcal{O}(1/\sqrt{m})$, $m = \min\{m_1, \dots, m_T\}$, where $\tilde{\mathcal{O}}$ denotes an expression that hides polylogarithmic factors in all problem parameters. Note that the leading term capturing the complexity of learning the meta-learner h decays in terms of the number of query samples ($\sum_{t=1}^T n_t$) among all training tasks. This implies that even with insufficient number of training tasks, improving the amount of query samples in the tasks can also be helpful to the final performance of the extracted meta-learner. Comparatively, such bounds deduced in some previous works, e.g., (Maurer et al., 2016), have mostly not involved this amount, and thus could not reflect such common sense fact that increasing the number of query samples in training tasks should be helpful for improving generalization of meta learning on new tasks. Actually, in many current meta learning applications, good performance can be obtained on the basis of not very large training task set (even only with one training task), e.g., hyper-parameter learning (Franceschi et al., 2018), neural architecture search (NAS) (Elsken et al., 2019), etc. Therefore, our result could provide a more rational and comprehensive theoretical explanation for these methods.

Another point is necessary to be illustrated. Many of existing meta learning theories employ traditional empirical error to develop the error bounds (Maurer et al., 2016; Tripuraneni et al., 2020), whose training strategy is not episodic (i.e., support/query training strategy) (Vinyals et al., 2016). This will be somehow inapplicable for modern meta learning algorithm, e.g., gradient-based meta-algorithms (Finn et al., 2017; Franceschi et al., 2018). Besides, some recent theoretical investigations (Denevi et al., 2019b; Balcan et al., 2019) study model-agnostic meta-algorithm (Finn et al., 2017), and explore the convergence guarantees for gradient-based meta learning. Specifically, (Denevi et al., 2018, 2019a; Khodak et al., 2019) pay attention to specific meta-algorithms and propose the corresponding generalization guarantees. However, they study the case that the loss function is convex or the mapping is linear, which might be relatively hard to make analysis on deep neural network. Also, the training strategy studied is not episodic, which tends not to be easily used to train practical popular meta-algorithms. The most related work (Yin et al., 2020; Chen et al., 2020) considered the support/query episodic training strategy. Albeit beneficial to illustrate some generalization insight of meta learning, their theory has not been used to conduct some feasible meta-regularization terms that can help improve the meta-learner. Different from the aforementioned works, we developed the theoretical results followed by the support/query training strategy. Besides, the theoretical guarantees contain the complexities of learning meta-learner h and learner \mathbf{f} , which can easily induce control effects of the meta-learner for improving generalization capability of meta-algorithms.

Particularly, it is seen that there exists a term $d_{\mathcal{F}}(\mu_t^s, \mu_t^q)$ in the error bound, which describes the distribution shift between training/support set and validation/query set among tasks. For applications without such domain shift, this term is zero and can be omitted. However, there are many meta learning applications with such support/query shift. A typical example is domain generalization (as demonstrated in Section 7), aiming to learn how to set a learning method implemented on the source domain data (i.e., support data), but

able to yield a learner which could perform well on different target domain data (i.e., query data). To this purpose, various training tasks need to be collected, each containing the simulated training-to-testing domain shift by splitting source domains into virtual training and testing (i.e., validation) domains. And then meta-learning can be executed to improve model training, e.g., MetaReg Balaji et al. (2018), Feature-Critic Li et al. (2019b). We believe that our theory with support/query discrepancy can then be used to well explain the theoretical insight of meta learning in such cases.

3.3 Meta-Test Stage: Generalizing to New Query Tasks

In the meta-test stage, we aim to transfer the extracted meta-learner \hat{h} to help set the learning method on new query tasks. As shown in Algorithm 2, in the meta-test stage, we have the training set $D_\mu^{tr} \sim (\mu^s)^{m_\mu}$, and for evaluating the performance we also assume to have a validation set $D_\mu^{val} \sim (\mu^q)^{n_\mu}$. The two sets are sampled from the query task $\mu = (\mu^s, \mu^q)$ drawn from the environment distribution η . Apply \hat{h} obtained in the meta-training stage to set the learning process on D_μ^{tr} to obtain the decision model, and the test performance on D_μ^{val} reflects the generalization performance of the meta-learner \hat{h} on the new task μ . In a nutshell, the test error $R_{test}(h)$ can be calculated as⁴

$$R_{test}(h) = \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_\mu^{val} \sim (\mu^q)^{n_\mu}} \mathbf{L}(\mathcal{LM}(D_\mu^{tr}; h), D_\mu^{val}). \quad (7)$$

We further employ the following excess transfer risk to measure the “closeness” between the estimated \hat{h} and the true underlying h^* , as defined in Eq. (6), in terms of helping fulfill new query task μ :

$$R_{test}(\hat{h}) - R_{test}(h^*).$$

To build a bridge between the meta-training and meta-test processes, similar as recent works (Du et al., 2021; Tripuraneni et al., 2020), we make the following assumption:

Assumption 3 (Task diversity) *Given the meta learning machine \mathcal{H} , and \hat{h} and h^* are defined in Eq. (6), it holds that*

$$R_\eta(\hat{h}) - R_\eta(h^*) \leq \alpha \left(R_{train}(\hat{h}) - R_{train}(h^*) \right) + \beta, \quad (8)$$

where $R_\eta(h)$ and $R_{train}(h)$ are defined in Eqs. (3) and (5), respectively.

In fact, $R_\eta(\hat{h}) - R_\eta(h^*)$ presents to measure the closeness between the estimated meta-learner \hat{h} and the underlying optimal meta-learner h^* in terms of an arbitrary new query task, and $R_{train}(\hat{h}) - R_{train}(h^*)$ defines this measure by virtue of the information extracted from training tasks. However, we do not have direct access to the underlying information of the meta-learner. One can only indirectly extract partial information from observed training

4. Here R_{test} denotes the meta-test error with respect to the meta-learner, and the missing expectation w.r.t. D_{tr} means that we use the empirical estimation of the task-specific learner to compute the meta-test error on D_μ^{val} . Comparatively, R_η denotes that we use the ideal expected estimation of the task-specific learner to compute the generalization performance.

Algorithm 2 Meta-Test: Generalization to New Query Tasks

Input: Query task $\mu = (\mu^s, \mu^q)$ drawn from environment η ; meta-learned meta-learner \hat{h} .
Do:

(1) Draw training set D_μ^{tr} from μ^s .

(2) For given \hat{h} , calculate $\mathcal{LM}(D_\mu^{tr}; \hat{h})$ to obtain \hat{f}_μ .

Return: The task-specific learner \hat{f}_μ .

tasks, and thus we assume Eq.(8) holds, where $\alpha \in \mathbb{R}$ reflects the task similarity between source meta-training tasks and target meta-test task, and $\beta \in \mathbb{R}$ denotes the small additive error. It is anticipated that transferring the learning methodology will not be possible when the new query tasks are entirely different from those training ones (i.e., α is large).

In other words, α could essentially encode the ratio of these two quantities, i.e., how well the meta-training tasks can cover the space captured by the \hat{h} needed to predict on new meta-test tasks. So it can be regarded as the task diversity measure of meta-training tasks. If all training tasks were quite similar, then it could only be expected that the meta-training stage can learn about a narrow slice of the learning methodology, which makes such task-transferring aim difficult. Generally, when the task diversity is large, then α is small; otherwise, large α implies relatively more similar training tasks.

When we instantiate our SLeM theoretical framework into typical meta learning applications, we could provide some specific value of α, β . Specifically, for the few-shot regression model (see Section 5), we can set $\alpha = \frac{M}{\sigma_{d_L}(\mathbf{K})}$ and $\beta = 0$ in Proposition 4, where $\sigma_{d_L}(\cdot)$ denotes the d_L -th singular value of a matrix at a decreasing order. Also, for the few-shot classification model (see Section 6), we can set $\alpha = \frac{M}{\sum_{k=1}^K \sigma_{d_L}((\mathbf{K})_k)}$ and $\beta = 0$ in Proposition 7. These examples can verify the reasonability and realizability of Assumption 3.

We can then present the theoretical result for the meta-test phase of meta learning.

Theorem 3 *If Assumptions 1 - 3 hold, for any $\delta > 0$, with probability at least $1 - \delta$, we have*

$$\begin{aligned} R_{test}(\hat{h}) - R_{test}(h^*) &\leq \alpha \left(R_{train}(\hat{h}) - R_{train}(h^*) \right) + \beta \\ &\quad + 6L\hat{\mathcal{G}}_{D_\mu^{tr}}(\mathcal{F}) + 2\mathbb{E}_{\mu \sim \eta} d_{\mathcal{F}}(\mu^s, \mu^q) + 6B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}. \end{aligned} \quad (9)$$

Theorem 3 provides an excess transfer risk bound of the meta-learned \hat{h} for query task μ , which can be more concisely expressed with three dominant terms as follows:

$$R_{test}(\hat{h}) - R_{test}(h^*) \lesssim \alpha \left(R_{train}(\hat{h}) - R_{train}(h^*) \right) + \mathcal{O} \left(\sqrt{\frac{C(\mathcal{F})}{m_\mu}} \right) + \mathbb{E}_{\mu \sim \eta} d_{\mathcal{F}}(\mu^s, \mu^q).$$

They can be interpreted as: one upper bounded by the task average excess risk in the meta-training stage (the first term), the complexity of learning the task-specific learner \hat{f}_μ for a new query task (the second term), as well as the distribution shift between training set and test set of query task (the third term). It is seen that the excess transfer risk stems

from two main components: one arises from the bias of using an imperfect meta-learner estimate to transfer knowledge across training tasks, which accounts for the meta-training risk (i.e., the excess task average risk in Theorem 2); the other arises from the difficulty of learning task-specific learner of query task with the estimated meta-learner (i.e., the last two terms of the upper bound).

In a word, the leading-order terms of the transfer risk for learning meta-learner h scales as $\tilde{\mathcal{O}}\left(\sqrt{C(\mathcal{H})/(\sum_{t=1}^T n_t)} + \sqrt{C(\mathcal{F})/m} + \sqrt{C(\mathcal{F})/m_\mu}\right)$, $m = \min\{m_1, \dots, m_T\}$. A naive algorithm which learns the new task in isolation, ignoring the training tasks, has an excess risk scaling $\mathcal{O}\left(\sqrt{(C(\mathcal{H}) + C(\mathcal{F}))/m_\mu}\right)$. This theoretically explained the fact that when $\sum_{t=1}^T n_t$ are sufficiently large compared with m_μ (e.g., the setting of few-shot learning, m_μ is always relatively small), the performance of meta learning should be evidently better than the baseline of learning in isolation.

3.4 Remarks

In the following we present two necessary remarks on our method.

The ERM estimation for meta-learner. We use empirical risk minimization (ERM) principle to derive generalization bounds for our framework with general losses, tasks, and models. In fact, we just assume that the ERM within and outer tasks estimators can be exactly computed, and our theoretical results are shown for the global minimizer of empirical risks, just following previous works, e.g., (Maurer et al., 2016; Tripuraneni et al., 2020; Xu and Tewari, 2021). And the qualitative predictions still hold true for gradient-based optimization algorithms as verified by our simulations on deep neural networks. When the loss function is not convex, we can compute the ERM within and outer task estimators by the aid of advances of bi-level optimization algorithms, e.g., (Ji et al., 2021; Liu et al., 2021a,b). Specifically, Ji et al. (2021) studied the bi-level optimization where the upper-level objective function is nonconvex, and Liu et al. (2021a,b) studied the bi-level optimization where inner-level objective function is nonconvex. Our theoretical results are hopeful to benefit from these latest bi-level optimization algorithms, and the bounds could be further improved with relatively loose conditions, which will be investigated in our future research.

Theory-induced meta-regularization. Since we introduce an explicit parameterized meta-learner to extract the hyper-parameter prediction function, it is easy to control and improve the learning of the meta-learner with proper meta-regularization techniques, just like regularization techniques being effective in improving capability of machine learning, which often impose controlling strategies on the task-specific model (learner) developed from statistical learning theory via structural risk minimization principle (Vapnik, 2013), e.g., the large margin regularizer for learner. To bridge the gap between meta-learning theory and its practical use, we will highlight the utility of our meta learning framework for obtaining the learning guarantees of some typical meta learning applications, including few-shot regression (Section 5), few-shot classification (Section 6) and domain generalization (Section 7). We can empirically substantiate that the theory-induced meta-regularization strategies are capable of helping consistently improve generalization capability of meta-learner on new query tasks.

4 Related Works

While there are a large number of meta learning literatures recently emerging in the field, most of them focused on practical feasible techniques against certain problems. In this section we mainly review the related studies considering the intrinsic understanding of the fundamental meta learning concepts as well as its basic learning theories. More related papers can be referred to in the recently proposed comprehensive surveys (Vanschoren, 2018; Hospedales et al., 2020).

Meta learning understandings and taxonomies. Meta learning has a long history in psychology (Ward, 1937), and was described by (Maudsley, 1980) as “the process by which learners become aware of and increasingly in control of habits of perception, inquiry, learning, and growth that they have internalized”. Then Schmidhuber (1987); Bengio et al. (1990) introduced it into computer science to train a meta-learner that learns how to update the parameters of the learner. Afterwards, this approach has been applied to learning to optimize deep networks (Andrychowicz et al., 2016). Besides, Vilalta and Drissi (2002) used meta learning to improve the learning bias dynamically through experience by the continuous accumulation of meta-knowledge. Lemke et al. (2015) further employed meta learning to extract meta-knowledge from different domains or problems. A recent survey paper of (Vanschoren, 2018) mainly attributed the capability of meta learning to its leveraging prior learning experience, and interpreted that it can learn new tasks more quickly. These literatures, however, have not presented a general mathematical formulation for meta learning, which could yet be useful to clearly help distinguish meta learning from previous related learning manners, like transfer learning and multi-task learning. Very recently, Hospedales et al. (2020) proposed a comprehensive survey paper, introducing a taxonomy of meta learning along three independent axes, i.e., meta-representation, meta-optimizer, and meta-objective. However, as aforementioned, this paper summarized the task of meta learning as learning fixed hyper-parameters instead of a hyper-parameter-prediction-function, making the learned methodology with insufficient flexibility adapt to new query tasks.

Conditional meta-learning. Several recent works (Wang et al., 2020d; Denevi et al., 2020, 2022; Rusu et al., 2019; Wang et al., 2019) addressed the issue of learning fixed hyper-parameters by conditioning hyper-parameters on target tasks. Wang et al. (2020d) used a task-specific objective functions by weighting meta-training data on target tasks based on structured prediction, to achieve task adaptive model initialization. TASML (Wang et al., 2020d) is a non-parametric approach and requires access to training tasks at test time for the task-specific objectives on target tasks. Considering the limitation for non-convex formulation of TASML, Denevi et al. (2020) proposed another conditional meta-learning framework for biased regularization and fine tuning, which learns a conditioning function mapping task’s side information into a task’s specific bias vector. Denevi et al. (2022) further proposed to learn a conditioning function mapping task’s side information into a linear representation for better performance in more scenarios. To sum up, they are against specific problems (task conditional model initialization, biased regularization and linear representation, respectively), which can be actually seen as special hyper-parameter configurations cases under the investigated general hyper-parameter setting framework (as can be seen in Table 1). Besides, our meta-learner is a general parameterized structure, e.g., we consider deep neural networks as the meta-learner in our experiments. Thus their

meta-learner could be a special case under our general meta-learner formulation, e.g., the meta-learner is chose from linear functions class.

Our SLeM framework can be easily integrated into the traditional machine learning framework to provide a fresh understanding and extension of the original machine learning framework. This facilitates us capable of readily establishing learning theory upon the traditional statistical learning theory, and provide generalization bounds for the new framework with general losses, tasks, and models. Comparatively, (Wang et al., 2020d; Denevi et al., 2020, 2022) have not specifically emphasized the relationship with conventional machine learning regime, and they mainly pay attention to the advantageous performance of conditional meta-learning approach compared with the standard unconditional counterpart. Specifically, similar to the structural risk minimization (SRM) principle in the conventional statistical learning theory, we can further analyze general theory-inducted control strategies for meta-learner to ameliorate its generalization capability according to the derived generalization bounds. While (Wang et al., 2020d; Denevi et al., 2020, 2022) does not develop corresponding theory-inducted controlling strategies based on their formulation and theory.

Theory of gradient-based meta-learning. GBML stems from the Model-Agnostic Meta-Learning (MAML) algorithm (Finn et al., 2017) and has been widely used in practice. Finn et al. (2019) showed that MAML meta-initialization is learnable via follow the meta leader method under strong-convexity and smoothness assumptions. Balcan et al. (2019); Denevi et al. (2019a) provided finite-sample meta-test-time performance guarantees in the convex setting for SGD-based Reptile algorithm or biased regularization via online learning. ARUBA (Khodak et al., 2019) further improved upon the bound of (Balcan et al., 2019; Denevi et al., 2019a), and obtained better statistical guarantees. And Denevi et al. (2019b) explored Online-Within-Online meta-learning approach via online mirror descent algorithm, and derived a cumulative error bound for their method. Fallah et al. (2020) developed a convergence analysis for one-step MAML for a general nonconvex objective, and Wang et al. (2020b); Ji et al. (2022) studied the convergence of multi-step MAML cases. Wang et al. (2021); Kao et al. (2022); Collins et al. (2022) provided theoretical analysis for two well-known GBML methods, MAML and ANIL. Recently, Fallah et al. (2021) studied the generalization properties of MAML using the connections between algorithmic stability and generalization bounds of algorithms. Huang et al. (2022) analyzed the generalization properties of MAML trained with SGD in the overparameterized regime under a mixed linear regression model. Yao et al. (2021) proposed task augmentation strategies to address meta-overfitting problem. Chen et al. (2022) further showed that overfitting is more likely to happen in MAML and its variants than in ERM, especially when the data heterogeneity across tasks is relatively high.

For all these GBML based methods, the bounds presented in the related literatures are restricted to MAML algorithms and its variants, i.e., unconditional meta-learning, which, however, may be likely to fail to adapt to heterogenous environments of tasks. Comparatively, in this study, we learn an explicit hyper-parameter prediction function conditioned on training tasks, aiming at more flexibly dealing with dynamic and changing environments of tasks. We then provide general-purpose bounds with decoupling the complexity of learning the task-specific learner from that of learning the task-transferrable hyper-parameter prediction function. Thus such new bounds should have a more comprehensive adaptability on real

complicated environments. This can also be evidently verified by the superiority of our method compared with MAML-based ones in our experiments.

Statistical learning to learn (LTL). The seminal work of (Baxter, 2000) introduced a framework to study the statistical benefits of meta learning, subsequently used to show excess risk bounds for ERM-like methods using techniques like covering numbers (Baxter, 2000), algorithmic stability (Maurer and Jaakkola, 2005; Chen et al., 2020) and Gaussian complexities (Maurer et al., 2016). Afterwards, Du et al. (2021); Tripuraneni et al. (2020, 2021) proposed a general framework for obtaining the generalization bound, and made analyses on the benefit of representation learning for reducing the sample complexity of the target task. Arora et al. (2020) proposed a representation learning approach for imitation learning via bilevel optimization, and demonstrated the improved sample complexity brought by representation learning. Zhou et al. (2019) statistically demonstrated the importance of prior hypothesis in reducing the excess risk via a regularization approach. Chua et al. (2021) further provided risk bounds on predictors found by fine-tuning via gradient descent, using an initial representation from a MAML-like algorithm. Bai et al. (2021) studied the implicit regularization effect of the practical design choice of train-validation splitting popular in meta learning. Sun et al. (2021) showed that the optimal representation for representation-based meta learning is overparameterized and provides sample complexity for the method of moment estimator. Bernacchia (2021) discovered that the optimal step size of overparameterized MAML during training is negative.

The PAC-Bayes framework has been extended to understand this LTL approach, facilitating a PAC-Bayes meta-population risk bound (Pentina and Lampert, 2014, 2015; Amit and Meir, 2018; Ding et al., 2021; Farid and Majumdar, 2021; Rothfuss et al., 2021a,b; Nguyen et al., 2022; Rezazadeh, 2022). These works mostly focus on the case where the meta learning model is underparameterized; that is, the total number of meta training data from all tasks is larger than the dimension of the model parameter. Comparatively, our meta learning model is problem-agnostic, which could be overparameterized or underparameterized. Besides, information theoretical bounds have been recently proposed in (Chen et al., 2021; Jose and Simeone, 2021; Jose et al., 2021; Hellström and Durisi, 2022), which bound the generalization error in terms of mutual information between the input training data and the output of the meta-learning algorithms. Particularly, Hellström and Durisi (2022) applied their bounds to a representation learning setting, and yields a bound coincides with the one reported in (Tripuraneni et al., 2020).

Compared with these bound estimation theories, the generalization bounds we derive have several differences. Firstly, previous works assume the traditional learning model to train the meta-learner, which is somehow deviated from the current meta learning practice with support/query episodic training mode such as MAML (Finn et al., 2017) and ProtoNet (Snell et al., 2017). Comparatively, we derive the bounds based on the commonly used support/query meta-training manner in meta learning practice. Besides, although previous works achieve solid theoretical justifications, they seldom emphasized theory-inspired meta-regularization terms to conduct practically feasible controlling strategies for further ameliorating the meta-learner performance, especially its generalization capability across diverse testing query tasks. In our study, however, our estimated bounds could be readily employed to induce such expected generalization-controlling strategies for general meta-learning algorithms. Specifically, our theory can be used to deduce several meta-

regularization strategies on soundly rectifying the learning track of the meta-learner to improve its generalization ability. The effectiveness of such strategy has been comprehensively substantiated with typical meta learning applications on few-shot regression, few-shot classification and domain generalization, as demonstrated in Sections 5, 6 and 7, respectively.

Hyper-parameter tuning/optimization. The hyper-parameter tuning/optimization techniques (Feurer and Hutter, 2019) often search some relatively small-scale hyper-parameters like regularization strength, learning rate, etc. While when faced with complicated and massive hyper-parametric configurations, especially those related to a deep neural network, conventional hyper-parameter tuning approaches are generally incapable of taking effect. Yet our framework aims to study highly complicated hyper-parametric configurations constituted in general components of the learning process, e.g., data screening, model constructing, loss function presetting and algorithm designing, etc. Please see Table 1 in the main manuscript, which lists some typical examples on this motivation. The hyper-parameter setting issues considered in this study thus significantly exceed those involved in conventional hyper-parameter tuning literatures.

Besides, hyper-parameter tuning techniques mainly focus on extracting hyper-parameters themselves for the investigated tasks, and often employ validation set based approaches to search proper hyper-parameters, e.g., random search (Bergstra and Bengio, 2012), Bayesian hyper-parameter optimization (Snoek et al., 2012), gradient-based hyper-parameter optimization (Franceschi et al., 2018), etc. While the framework emphasized in this study employs an explicit hyper-parameter prediction function (meta-learner) for predicting proper hyper-parameters, which can more flexibly and adaptively fit query tasks. Besides, it often adopts bi-level optimization tools to calculate hyper-parameters under a sound optimization framework, which is always more efficient than purely searching strategies.

Furthermore, the hyper-parameter tuning techniques usually focus on finding optimal hyper-parametric configurations for tasks at hand. Despite their success, the searched hyper-parameters should be mainly appropriate for the investigated learning tasks, but can be hardly finely generalized to wider range of learning problems with heterogenous environments, which may require to set substantially different hyper-parameters. Comparatively, the explored hyper-parameter setting function in the suggested meta-learning framework of this study aims to extract the hyper-parameter setting rule shared by training tasks sampled from a task distribution, which is expected to facilitate a better flexibility of the learned hyper-parameters setting function performing well on new tasks sampled from this task distribution. Thus this hyper-parameter setting function should potentially own better task generalization ability, which can be readily used to predict proper hyper-parameters conditioned on new query tasks. Especially, we have provided the theoretical evidence to support such task transfer generalization capability of our hyper-parameter prediction function, to make this essential characteristics more solid and convincing.

Multi-task/Transfer Learning. Here we want to shortly clarify the main difference of the meta learning framework as aforementioned from conventional multi-task/transfer learning approaches Weiss et al. (2016); Zhuang et al. (2020); Pan and Yang (2010); Zhang and Yang (2018); Ruder (2017). Although the latter also makes use of the intrinsic relationship among multiple tasks for improving the generalization performance for the learning results, they still fall into the scope of conventional machine learning, which considers the problem at the data/learner level. Specifically, most conventional multi-task/transfer learning works aim to

improve the task learning performance by virtue of ameliorating the parameters of multiple learners imposed on different tasks by leveraging their underlying relevance. Meta-learning, however, considers the problem at the task/meta-learner level, and aims to learn the common hyper-parameter setting function imposed on the parameters of one single meta-learner functioned on different tasks. This way inclines to better improve its flexibility, available range and potential capability in practice than conventional learning manners. In particular, a feasible meta learning method allows the learners on different training/test tasks with very different forms, like variant input data modalities (Cubuk et al., 2019), model architectures (Zoph and Le, 2017), feature dimensions (Ryu et al., 2020), etc, and only requires them to possess certain shared common knowledge in learning method setting, which yet should be hardly handled by conventional multi-task/transfer learning manners. Such consideration in the learning-methodology perspective inclines to make this learning manner be with a more widely potential usage than many previous machine learning fashions.

5 Application I: Few-shot Regression

In this section, we instantiate the proposed meta learning framework for the few-shot regression model.

5.1 Basic Setting

Here we consider the few-shot regression setting where $\mathcal{X} = \mathbb{R}^d, \mathcal{Y} = \mathbb{R}$. Let $\|x\| \leq R, |y| \leq B, \forall x \in \mathcal{X}, y \in \mathcal{Y}$. The output of the meta-learner is the representation of each sample. The loss functions ℓ is chosen as the square loss $\ell(\mathbf{w}^\top h(x), y) = (y - \mathbf{w}^\top h(x))^2$. The task-specific learning machine \mathcal{F} are chosen as linear regression function maps, and the meta learning machine \mathcal{H} is parameterized as a depth- L vector-valued deep neural network (DNN) to extract the common representation for various regression tasks. Concretely, $h(x)$ can be written as:

$$h(x) = \phi_L(\mathbf{W}_L(\phi_{L-1}(\mathbf{W}_{L-1} \cdots \phi_1(\mathbf{W}_1 x)))), \quad (10)$$

where $\mathbf{W}_k, k \in [L]$ is the parameter matrix, and $\phi_k, k \in [L]$ is the activation function. Here we assume that the activation functions in all layers are element-wise 1-Lipschitz and zero-centered as assumed in (Golowich et al., 2018) which is typically satisfied, like the ReLU. Formally, \mathcal{F} and \mathcal{H} are written as:

$$\begin{aligned} \mathcal{F} &= \{f|f(z) = \mathbf{w}^\top z, \mathbf{w}, z \in \mathbb{R}^{d_L}, \|\mathbf{w}\| \leq M\}, \\ \mathcal{H} &= \{h|h(x) \in \mathbb{R}^{d_L} \text{ as defined in (10)}, x \in \mathcal{X}\}. \end{aligned} \quad (11)$$

Following the setting in (Finn et al., 2017), we assume that there are T tasks $\mathbf{D} = \{D_t\}_{t=1}^T$ available for learning, and $D_t = (D_t^{tr}, D_t^{val})$, where $D_t^{tr} = \{z_{ti}^{(s)} = (x_{ti}^{(s)}, y_{ti}^{(s)})\}_{i=1}^m$, $D_t^{val} = \{z_{tj}^{(q)} = (x_{tj}^{(q)}, y_{tj}^{(q)})\}_{j=1}^n$. Here, the number of training/validation set samples for each task is

identical. The few-shot regression model is then written as:

$$\begin{aligned} \mathbf{W}^* &= \arg \min_{\mathbf{W}} \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n \ell(\mathbf{w}_t^{*\top} h(x_{tj}^{(q)}), y_{tj}^{(q)}) \\ \text{s.t., } \mathbf{w}_t^* &= \arg \min_{\mathbf{w}_t} \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{w}_t^T h(x_{ti}^{(s)}), y_{ti}^{(s)}), t \in [T], \end{aligned} \quad (12)$$

where $\mathbf{W} = \{\mathbf{W}_k, k \in [L]\}$ represents the collection of parameter matrices of h and \mathbf{w}, h is chosen from \mathcal{F}, \mathcal{H} as defined in Eq. (11). For notation convenience, we denote $\mathbf{P} = (\mathbf{w}_1, \dots, \mathbf{w}_T)^\top \in \mathbb{R}^{T \times d_L}$, and $\mathbf{P}^* = (\mathbf{w}_1^*, \dots, \mathbf{w}_T^*)^\top$ as its theoretical optimal solution.

5.2 Theoretical Analysis

In this part, we will instantiate Theorem 3 for few-shot regression model as defined in Eq. (12). The $Dis(D^{val})$ can be computed as

$$\begin{aligned} Dis(D^{val}) &= \sup_{h, h'} \rho_{D^{val}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) \leq 4 \sup_{h, x \in \mathcal{X}} \|\mathbf{w}^\top h(x)\| \\ &\leq \sup_{h, x} 4M \|h(x)\| \leq 4MD \cdot \prod_{k=1}^L \|\mathbf{W}_k\|, \end{aligned}$$

where $\rho_{D^{val}}(\cdot, \cdot)$ is defined in Theorem 2, and $\|h(x)\|$ is bounded by (Let $\mathbf{r}_k(\cdot)$ denote the vector-valued output of the k -layer for $k \in [L]$):

$$\begin{aligned} \|h(x)\| &= \|\mathbf{r}_L(x)\| = \|\phi_L(\mathbf{W}_L \mathbf{r}_{L-1}(x))\| \\ &\leq \|\mathbf{W}_L \mathbf{r}_{L-1}(x)\| \leq \|\mathbf{W}_L\| \|\mathbf{r}_{L-1}(x)\| \leq D \prod_{k=1}^L \|\mathbf{W}_k\|. \end{aligned} \quad (13)$$

Now we can verify that Assumptions 1 - 3 holds. Observe that $\nabla_{\hat{y}} \ell(\hat{y}, y) = 2(\hat{y} - y) \leq 2(B + M\|h(\mathbf{x})\|)$, and thus the loss function is Lipschitz continuous with $L = 2(B + M\|h(x)\|)$, and $\|h(x)\|$ is bounded by Eq.(13). Besides, since $|\ell(\hat{y}, y)| \leq B^2 + 2BM\|h(x)\| + M^2\|h(x)\|^2$, the loss function is bounded. The following result verifies that Assumption 3 holds,

Proposition 4 *Consider the few-shot regression model defined in Eq.(12), and the loss function $\ell(\cdot, \cdot)$ is chosen as the squared loss. The feature representation and the linear regression function are designed as in Eq.(11). Then Assumption 3 holds with $\alpha = \frac{M}{\sigma_{d_L}(\mathbf{K})}$ and $\beta = 0$, where $\mathbf{K} = \mathbf{P}^{*\top} \mathbf{P}^*/T$, $\sigma_{d_L}(\mathbf{K})$ denote the d_L -th singular value of matrix \mathbf{K} at a decreasing order.*

It can be seen that α reflects the diversity of training task set. The larger the similarity of task-specific learners is (i.e., the smaller $\sigma_{d_L}(\mathbf{K})$ is), the large α value is. Namely, the larger the similarity of task-specific learners is, the harder the meta-learner could be transferably used to new query tasks.

Now we can compute the leading-order terms in Eq.(9) for the parameterized meta-learner and learners as defined in Eq.(11). We firstly show the Rademacher complexity of \mathcal{H} in the following theorem.

Theorem 5 (Theorem 1 in (Golowich et al., 2018)) *Let \mathcal{H} be the class of real-valued DNN as defined in Eq.(10) while requiring $d_L = 1$ over $\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\| \leq R\}$, where each parameter matrix \mathbf{W}_i , $i \in [L]$ has Frobenius norm at most B_i , and the activation function ϕ_i , $i \in [L]$ is 1-Lipschitz, with $\phi_i(0) = 0$, and applied element-wise. Then we have:*

$$\hat{\mathfrak{R}}_N(\mathcal{H}) \leq \frac{R \left(\sqrt{2 \log(2)L} + 1 \right) \prod_{i=1}^L B_i}{\sqrt{N}}.$$

(1) For the Gaussian complexity of meta-learner h :

$$\begin{aligned} \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) &= \mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n \sum_{k=1}^{d_L} g_{tjk} h_k(x_{tj}^{(q)}) \\ &\leq \sum_{k=1}^{d_L} \hat{\mathcal{G}}_{D^{val}}(h_k) \leq 2\sqrt{\log(nT)} \sum_{k=1}^{d_L} \hat{\mathfrak{R}}_{D^{val}}(h_k) \\ &\leq 2d_L \sqrt{\log(nT)} \cdot \frac{R \left(\sqrt{2 \log(2)L} + 1 \right) \prod_{i=1}^L B_i}{\sqrt{nT}}, \end{aligned} \quad (14)$$

where the second inequality holds since $\hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{H}) \leq 2\sqrt{\log(N)}\hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{H})$ for any function class \mathcal{H} when \mathbf{X} has N data points (Ledoux and Talagrand, 2013).

(2) For the Gaussian complexity of the task-specific learner:

$$\begin{aligned} \hat{\mathcal{G}}_{D_t^{val}}(\mathcal{F}) &= \mathbb{E} \sup_{\mathbf{w} \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m g_{ti} \mathbf{w}^\top h(x_{ti}^{(s)}) \\ &\leq \frac{\|\mathbf{w}\|}{m} \mathbb{E} \left\| \sum_{i=1}^m g_{ti} h(x_{ti}^{(s)}) \right\| \leq \frac{\|\mathbf{w}\|}{m} \sqrt{\sum_{i=1}^m \|h(x_{ti}^{(s)})\|^2} \\ &\leq \frac{\|\mathbf{w}\|}{\sqrt{m}} \cdot \max_{x_{ti}^{(s)} \in D_t^{tr}} \|h(x_{ti}^{(s)})\|. \end{aligned} \quad (15)$$

Thus, the transfer error defined in Eq.(9) now can be written as

$$\begin{aligned} &R_{test}(\hat{f}_\mu, \hat{h}) - R_{test}(f_\mu^*, h^*) \\ &\leq \frac{M}{\sigma_{d_L}(\mathbf{K})} \left(768L \log(4nT) L(\mathcal{F}) 2d_L \sqrt{\log(nT)} \frac{R \left(\sqrt{2 \log(2)L} + 1 \right) \prod_{i=1}^L B_i}{\sqrt{nT}} \right. \\ &\quad \left. + \frac{6L\|\mathbf{w}\|}{\sqrt{mT}} \sum_{t=1}^T \max_{x_{ti}^{(s)} \in D_t^{val}} \|h(x_{ti}^{(s)})\| + 6B\sqrt{\frac{\log \frac{2}{\delta}}{2nT}} + 6B\sqrt{\frac{\log \frac{2}{\delta}}{m}} + \frac{48L \sup_{h,x} M \|h(x)\|}{n^2 T^2} \right) \\ &\quad + \frac{6L\|\mathbf{w}\|}{\sqrt{m_\mu}} \cdot \max_{x_i^{(s)} \in D_\mu^{val}} \|h(x_i^{(s)})\| + 6B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}, \end{aligned} \quad (16)$$

where $\sigma_1(\mathbf{X}), \dots, \sigma_r(\mathbf{X})$ denote the singular values of a rank r matrix \mathbf{X} at a decreasing order. Note that we assume that there exists no distribution shift between the training and validation sets, and thus the term $d_{\mathcal{F}}(D_t^{(tr)}, D_t^{(val)})$ is zero and can be omitted.

5.3 Theory-inspired Meta-Regularization

It can be seen that the complexity of the meta-learner is normal without extra term needed to be restricted, while the complexity of task-specific meta-learner has additional term $\max_{x_{ti}^{(s)} \in D_t^{val}} \|h(x_{ti}^{(s)})\|$, which is dependent on the output range of the meta-learner. The transfer error bound in Eq.(16) can then naturally inspire the following three controlling strategies for improving the generalization capability of the extracted meta-learner h .

(i) Control the output range of the meta-learner h . Conventional models usually set all activation functions of meta-learner easily as ReLU. If we revise the last activation function as Tanh (i.e., $\phi_L = \frac{e^z - e^{-z}}{e^z + e^{-z}}$), then we have

$$\|h(\mathbf{x})\| = \|\mathbf{r}_L(\mathbf{x})\| = \|\phi_L(\mathbf{W}_L \mathbf{r}_{L-1}(\mathbf{x}))\| \leq \sqrt{d_L}.$$

Generally, $\sqrt{d_L}$ is smaller than $D \prod_{k=1}^L \|\mathbf{W}_k\|_F$, and the complexity can thus be expected to substantially decrease, and the generalization of the calculated meta-learner is hopeful to be improved (more analysis and discussion could be found in Appendix H.1).

(ii) Minimize $\|\mathbf{w}\|$ of the learners. The terms $\frac{3L\|\mathbf{w}\|}{\sqrt{mT}} \sum_{t=1}^T \max_{x_{ti}^{(s)} \in D_t^{val}} \|h(x_{ti}^{(s)})\|$ and $\frac{3L\|\mathbf{w}\|}{\sqrt{m\mu}}$ imply that minimizing the $\|\mathbf{w}\|$ also tends to decrease the transfer error in Eq.(16).

(iii) Maximize the $\sigma_{d_L}(\mathbf{K})$. The term $\frac{M}{\sigma_{d_L}(\mathbf{K})}$ accounts for the gravity of the meta-training error influencing the final transfer error in Eq.(16), and thus maximizing the $\sigma_{d_L}(\mathbf{K})$ also inclines to reduce the transfer error.

It should be emphasized that the above training strategy (i) corresponds to a meta-regularized control manner imposed on the meta-learner, while the training strategies (ii) and (iii) put controls on the task-specific learners. And they could lead to different training behaviors as can be observed in the next subsection.

The training strategy (i) is easy to be implemented by directly replacing the last activation function of learners from conventional ReLU to Tanh function. The training strategy (ii) can be achieved by adding a $\|\mathbf{w}\|$ regularizer into the meta-training objective in Eq. (12) as:

$$\begin{aligned} \mathbf{W}^* &= \arg \min_{\mathbf{W}} \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n \ell(\mathbf{w}_t^{\top} h(\mathbf{x}_{tj}^{(q)}), y_{tj}^{(q)}) \\ s.t., \mathbf{w}_t^* &= \arg \min_{\mathbf{w}_t} \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{w}_t^{\top} h(\mathbf{x}_{ti}^{(s)}), y_{ti}^{(s)}) + \lambda_1 \|\mathbf{w}_t\|^2, t \in [T]. \end{aligned}$$

In the meta-test stage, it solves the following objective for a new query task given \hat{h} :

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{M_0} \sum_{i=1}^{M_0} \ell(\mathbf{w}^{\top} \hat{h}(\mathbf{x}_i^{(s)}), y_i^{(s)}) + \lambda_2 \|\mathbf{w}\|^2.$$

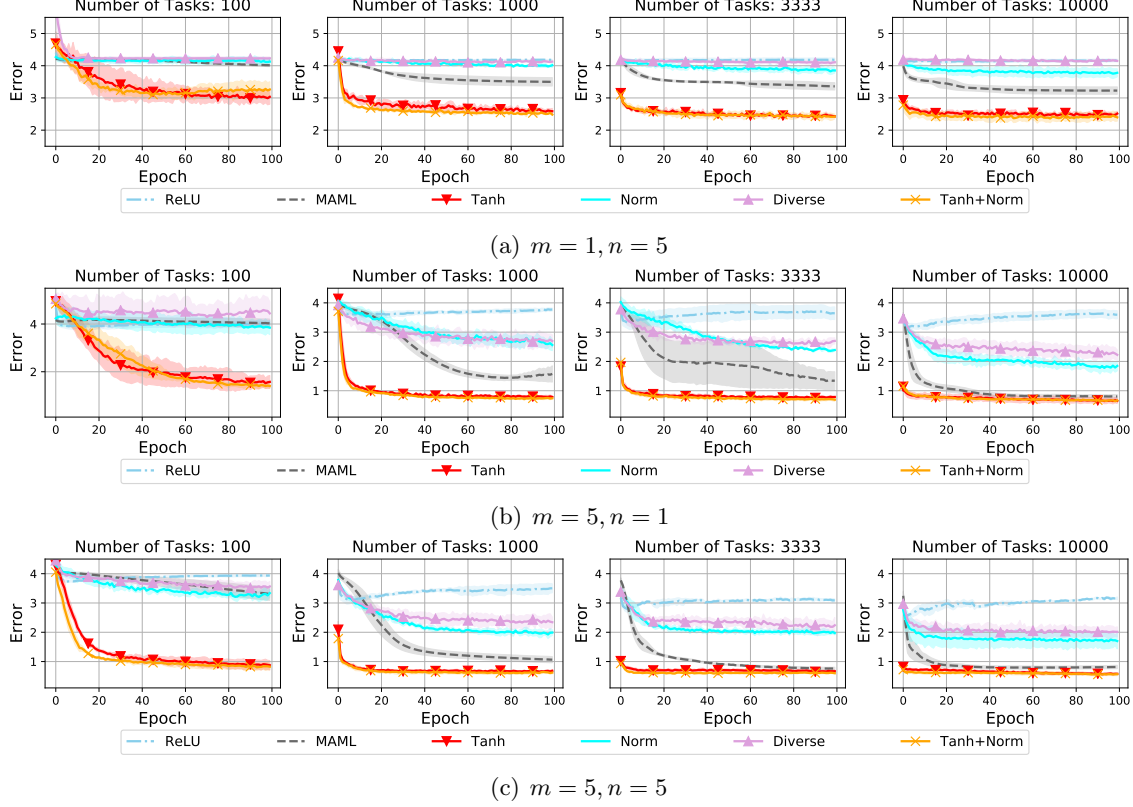


Figure 3: Transfer error changing curves of different training strategies with $T = 100, 1000, 3333, 10000$ tasks and different m, n values. The variance of each curve over 3 independent runs has also depicted along the curve to show the stability of each method.

The training strategy (iii) can be realized by solving the following meta-training objective:

$$\begin{aligned} \mathbf{W}^* &= \arg \min_{\mathbf{W}} \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n \ell(\mathbf{w}_t^{*\top} h(\mathbf{x}_{tj}^{(q)}), y_{tj}^{(q)}) \\ s.t., \mathbf{P}^* &= \arg \min_{\mathbf{P}} \frac{1}{mT} \sum_{t=1}^T \sum_{i=1}^m \ell(\mathbf{w}_t^{\top} h(\mathbf{x}_{ti}^{(s)}), y_{ti}^{(s)}) - \lambda_3 \sigma_{d_L}(\mathbf{P}^{\top} \mathbf{P} / T). \end{aligned}$$

And the $\lambda_1, \lambda_2, \lambda_3$ are the hyper-paramters of the above regularization problems.

5.4 Numerical Experiments

In this section, we test the effectiveness of the theory-induced training strategies on few-shot regression benchmark. We follow the experimental setting of MAML (Finn et al., 2017). Each task involves regression from the input to the output of a sine wave, where the amplitude and phase of the sinusoid are varied among tasks. Thus the problem aims to approximate a family of sine functions $f(x) = \alpha \sin(\beta x)$. The task distribution η is the joint distribution $p(\alpha, \beta)$ of the amplitude parameter α and the phase parameter β . We set $p(\alpha) = U[0.1, 5]$ and $p(\beta) = U[0, \pi]$. All the meta-training and meta-test tasks are randomly generated from the task distribution $p(\alpha, \beta)$. The function class \mathcal{H} is set as an MLP with two hidden

layers of size 40 with ReLU activation function, and \mathcal{F} is a linear layer with bias False. Both the input and the output layers have dimensionality 1. The loss is the mean-squared error between the prediction $\mathbf{w}^\top h(x)$ and true value. The transfer error is averaged over 600 meta-test tasks with varying shots and queries. It uses episodic training strategy for meta-training, i.e., the meta-algorithm observes a set of training tasks sequentially and applies stochastic gradient descent with one task per batch.

We implement the MAML (Finn et al., 2017) as the baseline method in Eq.(12) solved by Bilevel Programming Franceschi et al. (2018). We denote the later by “ReLU”, since the last activation function of the meta-learner is ReLU. The ‘ReLU’ and ‘MAML’ are implemented based on the code located at <https://github.com/jiaxinchen666/meta-theory> released by the paper (Chen et al., 2020). And ‘Tanh’, ‘Norm’, ‘Diverse’ and ‘Tanh+Norm’ denote the training strategies (i), (ii), (iii), and both (i), (ii) applied to the baseline method, respectively. We set $\lambda_1 = \lambda_2 = 1$ for training strategy (ii) and $\lambda_3 = 10$ for training strategy (iii). The task-specific optimizer is set as Adam optimizer with learning rate 0.01, and the meta-optimizer is set as Adam optimizer with learning rate 0.001. The other hyper-parameters of task-specific optimizer and meta-optimizer are the default settings of Adam optimizer.

Fig.3 shows the transfer error of different training strategies with varied numbers of training tasks, support samples and query samples. It can be seen that the proposed training strategies can help consistently improve the performance of the baseline method on different meta-training tasks or support/query samples. As shown, training strategy (i) achieves an evidently larger improvement compared with (ii), (iii). Furthermore, we combine training strategies (i) and (ii), hoping to achieve further improvement compared with only one training strategy (i). However, there exists an unsubstantial improvement when using both training strategies (i) and (ii). This implies that a proper meta-regularization technique for meta-learner can bring more essential improvements compared with the regularization techniques for task-specific learners. Specifically, when m is small ($m = 1$), the training strategies (ii), (iii) bring little improvement compared with ‘ReLU’, while training strategy (i) produces consistent improvement with different training tasks and m, n values.

Note that the training strategies (ii) and (iii) put controls on the task-specific learners, which are important to improve the performance for the given training tasks. While the training strategy (i) adds a meta-regularized control imposed on the meta-learner, which is verified to be more important to improve the performance for transferring to the new query tasks. Such phenomenon is rational and has been observed comprehensively in our experiments. We thus will only study the control strategies for the meta-learner in the latter sections.

6 Application II: Few-shot Classification

In this section, we instantiate our meta learning framework for the few-shot classification problem.

6.1 Basic Setting

For the few-shot classification issue, usually one considers the K -way N -shot setting, in which each task contains NK examples from K classes with N examples for each class. Thus the task-specific predictor machine \mathcal{F} is often a K -class linear classifier. Let $\mathcal{X} =$

$\mathbb{R}^d, \mathcal{Y} = \{0, 1, \dots, K-1\}$ and $\|x\| \leq R, \forall x \in \mathcal{X}$. The output of the meta-learner is the feature representation of each sample, and the meta learning machine \mathcal{H} is parameterized as a depth- L vector-valued convolutional neural network (CNN) to extract the common feature representation for various classification tasks. Concretely, $h(x)$ can be written as:

$$h(x) = \phi_L(\mathbf{W}_L(\phi_{L-1}(\mathbf{W}_{L-1} \cdots \phi_1(\mathbf{W}_1 x))), \quad (17)$$

where $\mathbf{W}_k, k \in [L]$ is the parameter matrix, and $\phi_k, k \in [L]$ is the 1-Lipschitz activation function. Formally, we additionally require both the norm of the weight matrix of each layer and the distance between the weights and the starting point weights are bounded, i.e.,

$$\|\mathbf{W}_j\|_F \leq B_j, \|\mathbf{W}_j^0\|_F \leq B_j, \|\mathbf{W}_j - \mathbf{W}_j^0\|_F \leq D_j, \quad j \in [L], \quad (18)$$

where \mathbf{W}_j^0 is the initialized parameter matrix of the model h . The distance to initialization has been observed to have substantial influence to generalization in deep learning (Bartlett et al., 2017; Neyshabur et al., 2019; Nagarajan and Kolter, 2019), and we introduce it here to develop the control strategy of the meta-learner. Specifically, \mathcal{F} and \mathcal{H} are chosen as:

$$\begin{aligned} \mathcal{F} &= \{f|f(z) = \mathbf{A}^\top z, \mathbf{A} \in \mathbb{R}^{d_L \times K}, \|\mathbf{A}\| \leq M\}, \\ \mathcal{H} &= \{h|h(x) \in \mathbb{R}^{d_L} \text{ as defined in (17)}, x \in \mathcal{X}\}. \end{aligned} \quad (19)$$

The conditional distribution of classification is

$$P(y = k|f(h(x))) = \frac{\exp(a_k)}{\sum_j \exp(a_j)} := \text{Softmax}(f(h(x))), a_k = (f(h(x)))_k, k \in [K]. \quad (20)$$

The loss functions ℓ is chosen as the cross-entropy loss $\ell(\mathbf{A}^\top h(x), y) = -\sum_{k=1}^K y_k \log(a_k)$, where y is the one-hot encoding of the ground-truth label, and a_k is defined as in Eq. (20).

Following the setting in (Finn et al., 2017; Lee et al., 2019), we assume that there are T tasks $\Gamma = \{D_t\}_{t=1}^T$ available for learning, and $D_t = (D_t^{tr}, D_t^{val})$, where $D_t^{tr} = \{z_{ti}^{(s)}\}_{i=1}^m, D_t^{val} = \{z_{tj}^{(q)}\}_{j=1}^n$. We set identical sample numbers for training/validation data sets for each task. The few-shot classification model is then written as

$$\begin{aligned} \mathbf{W}^* &= \arg \min_{\mathbf{W}} \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n \ell(\mathbf{A}_t^{*\top} h(x_{tj}^{(q)}), y_{tj}^{(q)}) \\ s.t., \mathbf{A}_t^* &= \arg \min_{\mathbf{A}_t} \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{A}_t^\top h(x_{ti}^{(s)}), y_{ti}^{(s)}), t \in [T], \end{aligned} \quad (21)$$

where $\mathbf{W} = \{\mathbf{W}_k, k \in [L]\}$ is the collection of the parameter matrices of h , and \mathbf{A}, h is chosen from \mathcal{F}, \mathcal{H} defined in Eq. (19). We denote $\mathbf{Q} = (\mathbf{A}_1, \dots, \mathbf{A}_T)^\top \in \mathbb{R}^{T \times d_L \times K}$, and $\mathbf{Q}^* = (\mathbf{A}_1^*, \dots, \mathbf{A}_T^*)^\top \in \mathbb{R}^{T \times d_L \times K}$ as its theoretical optimal solution.

6.2 Theoretical Analysis

Here, we will instantiate the general Theorem 3 for few-shot classification model defined in Eq. (21). The $Dis(D^{val})$ can be computed as

$$\begin{aligned} Dis(D^{val}) &= \sup_{h, h'} \rho_{2, D^{val}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) \leq 4 \sup_{h, x \in \mathcal{X}} \|\mathbf{A}^\top h(x)\|_2 \\ &\leq \sup_{h, x} 4M \|h(x)\| \leq 4MD \cdot \prod_{k=1}^L \|\mathbf{W}_k\|_F. \end{aligned}$$

Now we can verify that Assumptions 1 - 3 hold. The following proposition implies the Lipschitz continuity for the cross-entropy function.

Proposition 6 *The loss function $\ell(f(h(x)), y)$ is 1-Lipschitz with respect to $f(h(x))$, where ℓ is cross-entropy loss.*

According to proposition 6, we have

$$|\ell(f(h(x)), y)| \leq \|f(h(x))\| = \|\mathbf{A}^\top h(x)\| \leq M \|h(x)\| \leq MD \prod_{k=1}^L \|\mathbf{W}_k\|_F,$$

and thus the loss function is bounded. The following result then verifies Assumption 3:

Proposition 7 *Consider the few-shot classification model defined in Eq.(21), and the loss function $\ell(\cdot, \cdot)$ chosen as the cross-entropy loss. The meta learning machine \mathcal{H} and task-specific predictor machine are designed as in Eq.(19). Then the α, β in Assumption 3 are $\frac{M}{\sum_{k=1}^K \sigma_{d_L}((\mathbf{K})_k)}$ and 0, where $(\mathbf{K})_k = (\mathbf{Q}^*)^\top_k (\mathbf{Q}^*)_k / T$, and $(\cdot)_k$ denote the k -th element, $\sigma_{d_L}(\mathbf{K}_k)$ denote the d_L -th singular value of matrix \mathbf{K}_k at a decreasing order*

It should be indicated that the value of α reflects the similarity of task-specific learners for given training task set. And minimizing α attempts to increase the diversity of learned task-specific learners, which can cover the space as much as possible captured by the h needed to be predicted on new tasks, and thus help improve the generalization ability of the extracted meta-learner for new query tasks.

Now we can compute the leading-order terms in Eq.(9) for the parameterized meta-learner and learners defined in Eq.(19). We firstly show the Rademacher complexity of \mathcal{H} in the following theorem.

Theorem 8 (Theorem 2 in (Gouk et al., 2021)) *Let \mathcal{H} be the class of real-valued DNN as defined in Eq.(17) and (18) over $\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\| \leq R\}$, where each parameter matrix \mathbf{W}_i , $i \in [L]$ has Frobenius norm at most B_i , and the activation function $\phi_i, i \in [L]$ is 1-Lipschitz, with $\phi_i(0) = 0$, and applied element-wise. Then we have:*

$$\hat{\mathfrak{R}}_N(\mathcal{H}) \leq \frac{2\sqrt{2}d_L R \sum_{j=1}^L \frac{D_j}{2B_j \prod_{i=1}^j \sqrt{c_i}} \prod_{j=1}^L 2B_j \sqrt{c_j}}{\sqrt{N}},$$

where c_i is the number of columns in \mathbf{W}_i .

(1) The Gaussian complexity of meta-learner h can be computed as

$$\begin{aligned}
 \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) &= \mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n \sum_{k=1}^{d_L} g_{tjk} h_k(x_{tj}^{(q)}) \\
 &\leq \sum_{k=1}^{d_L} \hat{\mathcal{G}}_{D^{val}}(h_k) \leq 2\sqrt{\log(nT)} \sum_{k=1}^{d_L} \hat{\mathfrak{R}}_{D^{val}}(h_k) \\
 &\leq 2\sqrt{\log(nT)} \cdot \frac{2\sqrt{2}d_L R \sum_{j=1}^L \frac{D_j}{2B_j \prod_{i=1}^j \sqrt{c_i}} \prod_{j=1}^L 2B_j \sqrt{c_j}}{\sqrt{nT}}.
 \end{aligned}$$

(2) The Gaussian complexity of the task-specific learner can be computed as:

$$\begin{aligned}
 \hat{\mathcal{G}}_{D_t^{val}}(\mathcal{F}) &= \mathbb{E} \sup_{\mathbf{A}_t \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K g_{tik} ((\mathbf{A}_t)_k^\top h(x_{ti}^{(s)})) \\
 &= \frac{1}{m} \mathbb{E} \sup_{\mathbf{A}_t \in \mathcal{F}} \sum_{k=1}^K ((\mathbf{A}_t)_k^\top \sum_{i=1}^m g_{tik} h(x_{ti}^{(s)})) \\
 &\leq \frac{M}{m} \mathbb{E} \sum_{k=1}^K \left\| \sum_{i=1}^m g_{tik} h(x_{ti}^{(s)}) \right\| \\
 &= \frac{MK}{\sqrt{m}} \sqrt{\frac{1}{m} \sum_{i=1}^m \|h(x_{ti}^{(s)})\|_2^2} \\
 &\leq \frac{MK}{\sqrt{m}} \cdot \max_{x_{ti}^{(s)} \in D_t^{tr}} \|h(x_{ti}^{(s)})\|.
 \end{aligned}$$

Different from few-shot regression, there exists an additional multiplicative factor with term K , which is the number of classes. This implies the complexity of the task-specific learner increases as the number of classes increases, complying with general experience in common practice.

Thus, the transfer error defined in Eq.(9) now can be written as

$$\begin{aligned}
 &R_{test}(\hat{f}_\mu, \hat{h}) - R_{test}(f_\mu^*, h^*) \\
 &\leq \frac{M}{\sum_{k=1}^K \sigma_{d_L}((\mathbf{K})_k)} \left(768L \log(4nT) L(\mathcal{F}) \cdot \frac{2\sqrt{2}d_L R \sum_{j=1}^L \frac{D_j}{2B_j \prod_{i=1}^j \sqrt{c_i}} \prod_{j=1}^L 2B_j \sqrt{c_j}}{\sqrt{nT}} \right. \\
 &\quad \left. + \frac{6LMK}{\sqrt{mT}} \sum_{t=1}^T \max_{x_{ti}^{(s)} \in D_t^{val}} \|h(x_{ti}^{(s)})\| + 6B \sqrt{\frac{\log \frac{2}{\delta}}{2nT}} + 6B \sqrt{\frac{\log \frac{2}{\delta}}{m}} + \frac{48L \sup_{h,x} M \|h(x)\|}{n^2 T^2} \right) \\
 &\quad + \frac{6LMK}{\sqrt{m_\mu}} \cdot \max_{x_i^{(s)} \in D_\mu^{val}} \|h(x_i^{(s)})\| + 6B \sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}.
 \end{aligned}$$

Note that we assume there exists no distribution shift between the training and validation sets, and thus the term $d_{\mathcal{F}}(D_t^{(tr)}, D_t^{(val)})$ is zero and omitted.

6.3 Theory-Inspired Meta-Regularization

As aforementioned, we develop the following two controlling strategies for meta learner h to improve its methodology transferable generalization capability among tasks.

(i) Control the output range of the meta-learner h . All activation functions of CNN defined in Eq.(17) are usually assumed to be ReLU. We take the same control strategy as in Section 5.3, and revise the last activation function as Tanh (i.e., $\phi_L = \frac{e^z - e^{-z}}{e^z + e^{-z}}$).

(ii) Minimize the distance between the weights from the starting point weights in Eq.(18). This control strategy can decrease the complexity of meta-learner h .

The training strategy (ii) can be achieved by adding the penalty terms corresponding to each layer into the meta-training objective in Eq. (12),

$$\begin{aligned} \mathbf{W}^* &= \arg \min_{\mathbf{W}} \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n \ell(\mathbf{A}_t^{\top} h(\mathbf{x}_{tj}^{(q)}), y_{tj}^{(q)}) + \lambda \sum_{j=1}^L \|\mathbf{W}_j - \mathbf{W}_j^0\|_F^2 \\ s.t., \mathbf{A}_t^* &= \arg \min_{\mathbf{A}_t} \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{A}_t^{\top} h(\mathbf{x}_{ti}^{(s)}), y_{ti}^{(s)}), t \in [T], \end{aligned}$$

where λ is the hyper-paramters of the meta-regularization. Note that this regularization scheme has been used to help improve the performance of transfer learning proposed by (Li et al., 2018b, 2019a). We firstly employ this regularization scheme for few-shot classification problems. A detailed comparison with current meta learning methods (Zhou et al., 2019; Rajeswaran et al., 2019) could be found in Appendix H.2.

6.4 Numerical Experiments

In this section, we test the effectiveness of the theory-induced training strategies on few-shot classification benchmark.

6.4.1 IMPLEMENTATION DETAILS

We parameterize the meta-learner h as a ResNet-12 network and a standard 4-layer convolutional network followed by (Snell et al., 2017; Lee et al., 2019). Also, we use DropBlock regularization (Ghiasi et al., 2018) for the ResNet-12 network. The compared methods include ProtoNet, MetaOptNet-RR and MetaOptNet-SVM, whose task-specific learners are nearest-neighbor classifier (Snell et al., 2017), ridge regression classifier (Bertinetto et al., 2019) and SVM classifier (Lee et al., 2019), respectively. We denote ‘‘Tanh’’, ‘‘ L^2 -SP’’ and ‘‘Tanh + L^2 -SP’’ as training strategies (i), (ii) and combining training strategy (i) and (ii).

We follow the setting in (Lee et al., 2019) to conduct the few-shot classification experiments. To optimize the meta-learners, we use SGD with Nesterov momentum of 0.9 and weight decay of 0.0005. Each mini-batch consists of 8 episodes. The model was meta-trained for 60 epochs, with each epoch consisting of 1000 episodes. The learning rate was initially set to 0.1, and then changed to 0.006, 0.0012, and 0.00024 at epochs 20, 40 and 50, respectively. During the meta-training stage, we adopt horizontal flip, random crop, and color (brightness, contrast, and saturation) jitter data augmentation techniques. We use 5-way classification in both meta-training and meta-test stages. Each class contains 6 test (query) samples during meta-training and 15 test samples during meta-testing. Our meta-trained model

was chosen based on 5-way 5-shot test accuracy on the meta-validation set. Meanwhile, during meta-training, we set training shot to 15 for miniImageNet with ResNet-12; 5 for miniImageNet with 4-layer CNN; 10 for tieredImageNet; 5 for CIFAR-FS; and 15 for FC100. We set the hyper-parameter λ as 0.1, and keep the default hyper-parameter setting for the compared baselines in the original papers. Our implementation is based on the code provided on <https://github.com/kjunelee/MetaOptNet>.

6.4.2 EXPERIMENTS ON CIFAR DERIVATIVES

The **CIFAR-FS** dataset (Bertinetto et al., 2019) is a recently proposed few-shot image classification benchmark, consisting of all 100 classes from CIFAR-100. The classes are randomly split into 64, 16 and 20 for meta-training, meta-validation, and meta-testing, respectively. Each class contains 600 images of size 32×32 .

The **FC100** dataset is another dataset derived from CIFAR-100 (Oreshkin et al., 2018), containing 100 classes which are grouped into 20 superclasses. These classes are partitioned into 60 classes from 12 superclasses for meta-training, 20 classes from 4 superclasses for meta-validation, and 20 classes from 4 superclasses for meta-testing. The goal is to minimize semantic overlap between classes. All images are also of size 32×32 .

Results. Table 4 summarizes the results on the 5-way CIFAR-FS and FC100 classification tasks with different shots. It can be seen that our theory-induced training strategy does help improve generalization error from SOTA results in most cases. Specifically, we shows the results where we vary the meta-learner for two different embedding architectures. In both cases that the dimension of the feature representation is low (1600, a standard 4-layer convolutional network), and the dimension is much higher (16000, ResNet12), our proposed meta-regularization schemes yield better few-shot accuracy than baselines. This implies that our proposed meta-regularization scheme is model-agnostic, in the sense that it can be directly applied to meta-regularize the learning for different meta-learners. Meanwhile, even on the harder FC100 dataset, our strategies can still help increase the accuracy for new query few-shot tasks, which highlights the advantage of our problem-agnostic meta-regularization schemes.

6.4.3 EXPERIMENTS ON IMAGENET DERIVATIVES

The **miniImageNet** dataset (Vinyals et al., 2016) is a standard benchmark for few-shot image classification benchmark, consisting of 100 randomly chosen classes from ILSVRC-2012 (Russakovsky et al., 2015). The meta-training, meta-validation, and meta-testing sets contain 64, 16 and 20 classes randomly split from 100 classes, respectively. Each class contains 600 images of size 84×84 . We use the commonly-used split proposed by (Ravi and Larochelle, 2017).

The **tieredImageNet** benchmark (Ren et al., 2019) is a larger subset of ILSVRC-2012 (Russakovsky et al., 2015), composed of 608 classes grouped into 34 high-level categories. These categories are then split into 3 disjoint sets: 20 categories for meta-training, 6 for meta-validation, and 8 for meta-test. This corresponds to 351, 97 and 160 classes for meta-training, meta-validation, and meta-testing, respectively. This dataset aims to minimize the semantic similarity between the splits similar to FC100. All images are of size 84×84 .

Table 4: Average few-shot classification accuracies (%) with 95% confidence intervals on CIFAR-FS and FC-100 meta-test splits.

model	CIFAR-FS 5-way		FC100 5-way	
	1-shot	5-shot	1-shot	5-shot
4-layer conv (feature dimension=1600)				
ProtoNet	59.89 \pm 0.70	80.14 \pm 0.53	36.13 \pm 0.53	50.69 \pm 0.55
ProtoNet(Tanh)	63.55 \pm 0.75 4.66\uparrow	80.30 \pm 0.52 0.16\uparrow	36.64 \pm 0.55 0.51\uparrow	50.53 \pm 0.56 0.16\downarrow
ProtoNet(L^2 -SP)	60.71 \pm 0.71 0.82\uparrow	80.01 \pm 0.53 0.13\downarrow	36.52 \pm 0.54 0.39\uparrow	51.10 \pm 0.56 0.01\uparrow
ProtoNet(Tanh + L^2 -SP)	63.58 \pm 0.75 4.66\uparrow	80.20 \pm 0.52 0.06\uparrow	36.74 \pm 0.55 0.61\uparrow	50.77 \pm 0.54 0.08\uparrow
MetaOptNet-RR	63.02 \pm 0.71	79.60 \pm 0.52	35.75 \pm 0.52	51.20 \pm 0.55
MetaOptNet-RR(Tanh)	63.45 \pm 0.70 0.43\uparrow	80.11 \pm 0.53 0.51\uparrow	38.61 \pm 0.56 2.86\uparrow	52.29 \pm 0.56 1.09\uparrow
MetaOptNet-RR(L^2 -SP)	63.66 \pm 0.72 0.64\uparrow	79.49 \pm 0.52 0.11\downarrow	36.55 \pm 0.52 0.80\uparrow	51.27 \pm 0.53 0.07\uparrow
MetaOptNet-RR(Tanh + L^2 -SP)	63.81 \pm 0.71 0.69\uparrow	80.44 \pm 0.50 0.84\uparrow	37.28 \pm 0.53 1.53\uparrow	52.31 \pm 0.52 1.11\uparrow
MetaOptNet-SVM	61.77 \pm 0.73	79.08 \pm 0.52	35.82 \pm 0.54	50.60 \pm 0.54
MetaOptNet-SVM(Tanh)	63.18 \pm 0.71 1.41\uparrow	79.87 \pm 0.53 0.49\uparrow	37.10 \pm 0.58 1.28\uparrow	51.62 \pm 0.57 1.02\uparrow
MetaOptNet-SVM(L^2 -SP)	62.70 \pm 0.71 0.93\uparrow	79.95 \pm 0.53 0.87\uparrow	37.13 \pm 0.68 1.31\uparrow	51.42 \pm 0.53 0.82\uparrow
MetaOptNet-SVM(Tanh + L^2 -SP)	63.24 \pm 0.72 1.47\uparrow	80.20 \pm 0.51 1.12\uparrow	37.29 \pm 0.58 1.47\uparrow	51.71 \pm 0.55 1.11\uparrow
ResNet-12 (feature dimension=16000)				
ProtoNet	68.00 \pm 0.74	83.50 \pm 0.51	38.43 \pm 0.58	52.56 \pm 0.56
ProtoNet(Tanh)	69.71 \pm 0.76 1.71\uparrow	83.62 \pm 0.51 0.12\uparrow	40.12 \pm 0.58 1.69\uparrow	55.02 \pm 0.54 2.06\uparrow
ProtoNet(L^2 -SP)	69.01 \pm 0.73 1.01\uparrow	83.57 \pm 0.51 0.07\uparrow	39.16 \pm 0.57 0.73\uparrow	53.60 \pm 0.56 1.04\uparrow
ProtoNet(Tanh + L^2 -SP)	71.22 \pm 0.74 3.22\uparrow	83.89 \pm 0.51 0.39\uparrow	40.70 \pm 0.58 2.43\uparrow	55.33 \pm 0.54 2.77\uparrow
MetaOptNet-RR	68.58 \pm 0.73	84.75 \pm 0.50	38.98 \pm 0.56	54.46 \pm 0.55
MetaOptNet-RR(Tanh)	71.64 \pm 0.72 3.06\uparrow	84.81 \pm 0.49 0.06\uparrow	40.13 \pm 0.58 1.15\uparrow	54.48 \pm 0.54 0.02\uparrow
MetaOptNet-RR(L^2 -SP)	71.40 \pm 0.70 2.82\uparrow	84.91 \pm 0.50 0.16\uparrow	39.79 \pm 0.57 0.81\uparrow	54.80 \pm 0.55 0.34\uparrow
MetaOptNet-RR(Tanh + L^2 -SP)	72.15 \pm 0.72 3.57\uparrow	85.03 \pm 0.49 0.28\uparrow	40.86 \pm 0.57 1.88\uparrow	55.58 \pm 0.54 1.12\uparrow
MetaOptNet-SVM	68.81 \pm 0.74	83.80 \pm 0.51	40.24 \pm 0.58	54.71 \pm 0.56
MetaOptNet-SVM(Tanh)	71.12 \pm 0.71 2.31\uparrow	85.19 \pm 0.49 1.39\uparrow	39.99 \pm 0.57 0.25\downarrow	53.92 \pm 0.56 0.79\downarrow
MetaOptNet-SVM(L^2 -SP)	70.84 \pm 0.72 2.03\uparrow	84.22 \pm 0.48 0.42\uparrow	40.41 \pm 0.57 0.17\uparrow	55.53 \pm 0.55 0.82\uparrow
MetaOptNet-SVM(Tanh + L^2 -SP)	71.52 \pm 0.73 2.71\uparrow	84.77 \pm 0.48 0.97\uparrow	40.83 \pm 0.58 0.59\uparrow	55.75 \pm 0.56 1.04\uparrow

Results. Table 5 summarizes the results on the 5-way classification tasks with different shots on miniImageNet and tieredImageNet benchmarks. Compared with CIFAR derivatives dataset, this benchmark contains more classes and more natural images. It can be observed that our proposed meta-regularization schemes can also help improve generalization error in most cases from SOTA baselines method without adding these schemes with different embedding architectures of meta-learners or different semantic similarity between the splits.

It should be emphasized that we have just directly replaced the last activation function (ReLU) of feature extractor (meta-learner) as Tanh, or add the meta-regularization scheme to the meta loss, and keep the original configurations without implementing any further fine-tuning executions and introducing extra computation cost. This way allows us to clearly see the effect of proposed meta-regularization schemes on the compared baselines. The empirical results verify that our theory-induced meta-regularization schemes are rational and arguably simple to improve the baseline performance.

6.4.4 CROSS-DOMAIN FEW-SHOT CLASSIFICATION

To further justify the effectiveness of our proposed meta-regularization schemes in producing robust features for unseen tasks, we further perform experiments under cross-domain few-shot classification settings, where the meta-test tasks are substantially different from meta-train tasks. We report the results in Table 6, using the same experiment settings, firstly introduced in (Chen et al., 2020), where miniImageNet is used as meta-train set and CUB dataset (Wah et al., 2011) as meta-test set.

Table 6 shows the cross-domain performance of the baselines and our proposed regularization schemes, which exhibits the similar tendency to few-shot classification results from Table 5. As is known, the CUB is a fine-grained classification dataset with more intra-class variations, and has a large domain shift to the miniImageNet dataset. There-

Table 5: Average few-shot classification accuracies (%) with 95% confidence intervals on miniImageNet and tieredImageNet meta-test splits.

model	miniImageNet 5-way		tieredImageNet 5-way	
	1-shot	5-shot	1-shot	5-shot
4-layer conv (feature dimension=1600)				
ProtoNet	47.69 ± 0.64	70.51 ± 0.52	50.10 ± 0.70	71.85 ± 0.57
ProtoNet(Tanh)	53.13 ± 0.64 5.44↑	71.22 ± 0.52 0.71↑	54.07 ± 0.69 3.93↑	71.05 ± 0.57 0.60↓
ProtoNet(L^2 -SP)	48.15 ± 0.63 0.46↑	70.67 ± 0.52 0.16↑	50.84 ± 0.71 0.74↑	72.50 ± 0.50 0.64↑
ProtoNet(Tanh + L^2 -SP)	54.43 ± 0.71 6.74↑	70.88 ± 0.52 0.37↑	54.06 ± 0.69 3.93↑	71.25 ± 0.57 0.40↓
MetaOptNet-RR	52.66 ± 0.63	69.80 ± 0.53	54.04 ± 0.68	72.09 ± 0.56
MetaOptNet-RR(Tanh)	52.70 ± 0.62 0.04↑	70.72 ± 0.52 0.98↑	54.23 ± 0.67 0.19↑	72.35 ± 0.55 0.26↑
MetaOptNet-RR(L^2 -SP)	52.38 ± 0.62 0.28↓	68.95 ± 0.52 0.85↓	55.66 ± 0.68 1.62↑	72.15 ± 0.57 0.06↑
MetaOptNet-RR(Tanh + L^2 -SP)	52.47 ± 0.62 0.19↓	70.17 ± 0.51 0.37↑	54.45 ± 0.69 0.41↑	72.40 ± 0.55 0.31↑
MetaOptNet-SVM	52.50 ± 0.62	69.66 ± 0.52	54.38 ± 0.70	71.48 ± 0.56
MetaOptNet-SVM(Tanh)	52.75 ± 0.62 0.25↑	70.79 ± 0.50 1.13↑	54.66 ± 0.69 0.28↑	71.57 ± 0.57 0.09↑
MetaOptNet-SVM(L^2 -SP)	52.20 ± 0.60 0.30↓	69.04 ± 0.52 0.62↓	55.52 ± 0.71 1.14↑	71.49 ± 0.57 0.01↑
MetaOptNet-SVM(Tanh + L^2 -SP)	52.54 ± 0.61 0.04↑	70.27 ± 0.52 0.61↑	54.66 ± 0.69 0.28↑	71.62 ± 0.50 0.14↑
ResNet-12 (feature dimension=16000)				
ProtoNet	56.13 ± 0.67	74.57 ± 0.53	63.02 ± 0.72	80.11 ± 0.56
ProtoNet(Tanh)	57.55 ± 0.67 1.42↑	74.62 ± 0.67 0.05↑	65.11 ± 0.73 2.09↑	79.69 ± 1.39 0.42↓
ProtoNet(L^2 -SP)	56.51 ± 0.65 0.38↑	75.15 ± 0.53 0.58↑	63.48 ± 0.72 0.46↑	80.65 ± 0.54 0.54↑
ProtoNet(Tanh + L^2 -SP)	58.10 ± 0.66 1.97↑	74.87 ± 0.51 0.30↑	65.68 ± 0.73 2.66↑	80.59 ± 0.51 0.48↑
MetaOptNet-RR	60.27 ± 0.67	76.19 ± 0.50	66.35 ± 0.71	81.25 ± 0.53
MetaOptNet-RR(Tanh)	60.44 ± 0.64 0.17↑	76.66 ± 0.47 0.47↑	66.97 ± 0.71 0.62↑	81.42 ± 0.52 0.17↑
MetaOptNet-RR(L^2 -SP)	60.05 ± 0.65 0.22↓	77.02 ± 0.48 0.83↑	66.77 ± 0.71 0.42↑	81.48 ± 0.53 0.23↑
MetaOptNet-RR(Tanh + L^2 -SP)	60.41 ± 0.65 0.14↑	76.98 ± 0.48 0.79↑	66.85 ± 0.71 0.50↑	81.49 ± 0.52 0.24↑
MetaOptNet-SVM	60.58 ± 0.66	75.99 ± 0.51	66.40 ± 0.72	81.18 ± 0.54
MetaOptNet-SVM(Tanh)	60.59 ± 0.65 0.01↑	77.05 ± 0.47 1.06↑	66.43 ± 0.70 0.03↑	81.30 ± 0.52 0.12↑
MetaOptNet-SVM(L^2 -SP)	61.17 ± 0.64 0.59↑	77.40 ± 0.49 1.41↑	66.45 ± 0.71 0.05↑	81.52 ± 0.53 0.34↑
MetaOptNet-SVM(Tanh + L^2 -SP)	60.99 ± 0.65 0.41↑	77.43 ± 0.48 1.44↑	66.49 ± 0.09 0.09↑	81.43 ± 0.52 0.25↑

fore, the extracted learning methodology (feature representation) from the meta-training dataset is mostly irrelevant to the meta-test datasets, which makes the learning tasks from new different domains relatively difficult, as suggested in (Baik et al., 2020). Even so, our proposed meta-regularization schemes can also evidently improves the performance of baselines. This further validates the effectiveness and validness of the developed theory and the theory-induced meta-regularization schemes to learn a more robust feature extractor with better adaptability to new tasks.

7 Application III: Domain Generalization

In this section, we instantiate the proposed meta learning framework for the domain generalization problem.

7.1 Basic Setting and Theoretical Analysis

Here we consider the domain generalization (DG) setting. We present the generally used meta learning setting (Li et al., 2018a, 2019b) for heterogeneous DG. Usually, assuming that we have T domains (datasets) $\Gamma = \{D_1, D_2, \dots, D_T\}$ for meta-training, each domain (task) has both meta-training and meta-validation data with $D_t = (D_t^{tr}, D_t^{val})$, where $D_t^{tr} = \{z_{ti}^{(s)}\}_{i=1}^{m_t}$, $D_t^{val} = \{z_{tj}^{(q)}\}_{j=1}^{n_t}$. Note that the label space is not shared between training/support and validation/query domains, but it is straightforwardly applicable to conventional (homogeneous) DG when assuming that the label space is shared between different domains.

There exist several differences between DG and few-shot classification problem. (1) The former often assumes that the meta-training and meta-validation sets of each domain (task) share the same label space, but have distribution shift. However, the latter often

assumes that they are with the same distribution. (2) The latter pays more attention to solving new query tasks with limited data, while the former not only deals with limited data in the meta-test domain, but also attempts to eliminate the distribution shift between meta-training and meta-test domain.

We take the same setting for \mathcal{F} and \mathcal{H} in few-shot classification, and borrow the theoretical analysis therein. Thus the transfer error defined in Eq.(9) can be written as

$$\begin{aligned}
& R_{test}(\hat{f}_\mu, \hat{h}) - R_{test}(f_\mu^*, h^*) \\
& \leq \frac{M}{\sum_{k=1}^K \sigma_{d_L}((\mathbf{K})_k)} \left(768L \log(4nT) L(\mathcal{F}) \cdot \frac{2\sqrt{2}d_L R \sum_{j=1}^L \frac{D_j}{2B_j \prod_{i=1}^j \sqrt{c_i}} \prod_{j=1}^L 2B_j \sqrt{c_j}}{\sqrt{nT}} \right. \\
& \quad + \frac{6LMK}{\sqrt{mT}} \sum_{t=1}^T \max_{x_{ti}^{(s)} \in D_t^{val}} \|h(x_{ti}^{(s)})\| + \frac{4}{T} \sum_{t=1}^T d_{\mathcal{F}}(D_t^{(tr)}, D_t^{(val)}) + 6B \sqrt{\frac{\log \frac{2}{\delta}}{2nT}} + 6B \sqrt{\frac{\log \frac{2}{\delta}}{m}} \\
& \quad \left. + \frac{48L \sup_{h,x} M \|h(x)\|}{n^2 T^2} \right) + \frac{6LMK}{\sqrt{m_\mu}} \cdot \max_{x_i^{(s)} \in D_\mu^{val}} \|h(x_i^{(s)})\| + \mathbb{E}_{\mu \sim \eta} d_{\mathcal{F}}(D_\mu^{(tr)}, D_\mu^{(val)}) + 6B \sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}.
\end{aligned}$$

Note that there exists an additional distribution shift term $d_{\mathcal{F}}(D_t^{(tr)}, D_t^{(val)})$ in the error bound, characterizing the distribution shift between the meta-training and meta-validation domains, as well as the term $d_{\mathcal{F}}(D_\mu^{(tr)}, D_\mu^{(val)})$ characterizing the distribution shift between the training and test sets of meta-test set. It is easy to see that proposed meta-regularization schemes in the few-shot classification can also be applied to DG, i.e., (i) Control the output range of the meta-learner h ; (ii) Minimize the distance of the weights from the starting point weights as Eq.(18).

Considering the difference from few-shot classification, DG should require proper training strategy to eliminate the distribution shift. Recently, Li et al. (2019b) proposed a Feature-Critic Networks technique to produce an auxiliary loss function to guide learning on the meta-training set to produce more robust and effective feature extractor on the meta-validation set. From this view, Feature-Critic Networks attempts to simulate the distribution shift between the meta-training and meta-validation sets, i.e., $d_{\mathcal{F}}(D_t^{(tr)}, D_t^{(val)})$. Due to its rationality and efficacy, in the following part, we employ Feature-Critic Networks as the baseline method, and further verify the effectiveness of proposed meta-regularization schemes for DG problem.

7.2 Numerical Experiments

In this section, we verify the effectiveness of the theory-induced training strategies on homogeneous and heterogeneous DG benchmark, respectively.

7.2.1 HOMOGENEOUS DG EXPERIMENTS

Dataset. The **PACS** dataset is a recent object recognition benchmark for domain generalisation (Li et al., 2017a). PACS contains 9991 images of size 224×224 from four different domains - Photo (P), Art painting (A), Cartoon (C) and Sketch (S). It has 7 categories across these domains: dog, elephant, giraffe, guitar, house, horse and person. It can be

Table 6: Average few-shot classification accuracies (%) with 95% confidence intervals on 5-way cross-domain classification.

model	miniImageNet \rightarrow CUB	
	1-shot	5-shot
4-layer conv (feature dimension=1600)		
ProtoNet	36.20 \pm 0.54	55.44 \pm 0.52
ProtoNet(Tanh)	39.33 \pm 0.55 3.13\uparrow	56.30 \pm 0.52 0.86\uparrow
ProtoNet(L^2 -SP)	36.40 \pm 1.80 0.20\uparrow	58.40 \pm 0.52 2.96\uparrow
ProtoNet(Tanh + L^2 -SP)	39.23 \pm 0.54 3.03\uparrow	57.46 \pm 0.53 2.02\uparrow
MetaOptNet-RR	41.27 \pm 0.55	58.15 \pm 0.52
MetaOptNet-RR(Tanh)	41.34 \pm 0.56 0.07\uparrow	59.82 \pm 0.51 1.67\uparrow
MetaOptNet-RR(L^2 -SP)	41.23 \pm 0.56 0.04\downarrow	59.11 \pm 0.51 0.96\uparrow
MetaOptNet-RR(Tanh + L^2 -SP)	42.13 \pm 0.55 0.86\uparrow	60.27 \pm 0.52 2.12\uparrow
MetaOptNet-SVM	41.47 \pm 0.56	58.21 \pm 0.53
MetaOptNet-SVM(Tanh)	41.89 \pm 0.56 0.42\uparrow	61.26 \pm 0.52 3.05\uparrow
MetaOptNet-SVM(L^2 -SP)	41.00 \pm 0.58 0.47\downarrow	58.81 \pm 0.51 0.60\uparrow
MetaOptNet-SVM(Tanh + L^2 -SP)	41.49 \pm 0.54 0.02\uparrow	59.44 \pm 0.50 1.23\uparrow
ResNet-12 (feature dimension=16000)		
ProtoNet	41.86 \pm 0.59	58.41 \pm 0.56
ProtoNet(Tanh)	42.23 \pm 0.63 0.37\uparrow	61.61 \pm 0.56 3.20\uparrow
ProtoNet(L^2 -SP)	43.60 \pm 0.60 1.74\uparrow	62.59 \pm 0.55 4.18\uparrow
ProtoNet(Tanh + L^2 -SP)	43.53 \pm 0.63 1.67\uparrow	60.03 \pm 2.56 1.62\uparrow
MetaOptNet-RR	44.43 \pm 0.59	64.12 \pm 0.51
MetaOptNet-RR(Tanh)	45.14 \pm 0.58 0.71\uparrow	64.23 \pm 0.52 0.11\uparrow
MetaOptNet-RR(L^2 -SP)	45.00 \pm 0.60 0.57\uparrow	64.44 \pm 0.53 0.32\uparrow
MetaOptNet-RR(Tanh + L^2 -SP)	45.24 \pm 0.60 0.81\uparrow	64.95 \pm 0.50 0.83\uparrow
MetaOptNet-SVM	44.60 \pm 0.59	64.53 \pm 0.53
MetaOptNet-SVM(Tanh)	45.09 \pm 0.57 0.49\uparrow	65.31 \pm 0.52 0.78\uparrow
MetaOptNet-SVM(L^2 -SP)	44.68 \pm 0.59 0.08\uparrow	64.85 \pm 0.53 0.32\uparrow
MetaOptNet-SVM(Tanh + L^2 -SP)	44.69 \pm 0.46 0.09\uparrow	64.90 \pm 0.52 0.37\uparrow

downloaded at <http://sketchx.eecs.qmul.ac.uk/>. We follow the standard protocol and perform leave-one-domain-out evaluation.

Implementation details. We parameterize the shared feature extractor (meta-learner) h as a pre-trained AlexNet (Krizhevsky et al., 2012) network followed by (Li et al., 2019b), and the task-specific learners as linear classifier. The compared methods include Reptile (Nichol et al., 2018), CrossGrad (Shankar et al., 2018), MetaReg (Balaji et al., 2018) and FC (Li et al., 2019b). We follow the experimental setting in (Li et al., 2019b). Specifically, the feature extractor and task-specific learners are trained with M-SGD optimizer (batch size/per meta-train domain=32, batch size/per meta-test domain=16, lr=0.0005, weight decay=0.00005, momentum=0.9) for 45K iterations. The Feature-Critic (set embedding variant) is adopted as MLP type, and is trained with the M-SGD optimizer (lr=0.001, weight decay=0.00005, momentum=0.9). We denote “Tanh”, “ L^2 -SP” and “Tanh + L^2 -SP” as the training strategies (i), (ii) and combining training strategy (i) and (ii). Our implementation is based on the code provided on https://github.com/liyiyang/Feature_Critic/.

Results. The comparison with state-of-the-art methods on PACS dataset is shown in Table 7. The results of Reptile, CrossGrad and MetaReg are directly taken from (Li et al., 2019b). We re-implement the result of AGG (Li et al., 2019b) and FC. Note that compared with those reported in (Li et al., 2019b), the depicted results of these two methods have a

Table 7: Cross-domain recognition accuracy (%) on PACS using train split in (Li et al., 2017a) for training.

Target	Reptile	CrossGrad	MetaReg	AGG	FC	FC+Tanh	FC+ L^2 -SP	FC+Tanh+ L^2 -SP
A	63.4	61.0	63.5	62.2	63.0	64.2	63.4	63.5
C	67.5	67.2	69.5	66.2	67.2	70.6	68.2	68.3
P	88.7	87.6	87.4	87.0	87.9	87.5	88.0	87.6
S	55.9	55.9	59.1	54.9	59.5	60.9	62.1	61.4
Ave.	68.9	67.9	69.9	67.6	69.4	70.8	70.4	70.2

slight drop on A, C, P domain, but a slight rise on S domain. These minor differences are due to our used higher Pytorch version. From the table, it can be easily observed that our proposed meta-regularization schemes consistently outperform other competing methods. Especially, as compared with the SOTA method FC, aiming to learn representations that generalise to new domains through training a supervised loss and explicitly simulates domain shift, our proposed meta-regularization scheme puts emphasis on decreasing the complexity of the feature extractor (meta-learner). Therefore, it can be combinationally used with FC to compensate its performance.

7.2.2 HETEROGENEOUS DG EXPERIMENTS

Datasets. The **Visual Decathlon** (VD) dataset consists of ten well-known datasets from multiple visual domains (Rebuffi et al., 2017). **FGVC-Aircraft Benchmark** (Maji et al., 2013) contains 10,000 images of aircraft, with 100 images for each of 100 different aircraft model variants. **CIFAR100** (Krizhevsky et al., 2009) contains colour images for 100 object categories. **Daimler Mono Pedestrian Classification Benchmark (DPed)** (Munder and Gavrila, 2006) consists of 50,000 grayscale pedestrian and non-pedestrian images. **Describable Texture Dataset (DTD)** (Cimpoi et al., 2014) is a texture database, consisting of 5640 images, organized according to a list of 47 categories such as bubbly, cracked, marbled. **The German Traffic Sign Recognition (GTSR) Benchmark** (Stallkamp et al., 2012) contains cropped images for 43 common traffic sign categories in different image resolutions. **Flowers102** (Nilsback and Zisserman, 2008) is a fine-grained classification task which contains 102 flower categories from the UK. **ILSVRC12 (ImageNet)** (Russakovsky et al., 2015) is the largest dataset in our benchmark, containing 1000 categories and 1.2 million images. **Omniglot** (Lake et al., 2015) consists of 1623 different handwritten characters from 50 different alphabets. **The Street View House Numbers (SVHN)** (Netzer et al., 2011) is a real-world digit recognition dataset with around 70,000 images. **UCF101** (Soomro et al., 2012) is an action recognition dataset of realistic human action videos, collected from YouTube. It contains 13,320 videos for 101 action categories. Each video has been summarized into an image based on a ranking principle using the Dynamic Image encoding of (Bilen et al., 2016). The images have been pre-processed to the size of 72×72 . Following the benchmark for DG in (Li et al., 2019b), we take the six larger datasets (CIFAR-100, DPed, GTSR, Omniglot, SVHN and ImageNet) as source domains and hold out the four smaller datasets (FGVC-Aircraft, DTD, Flowers102 and UCF101) as target domains.

Implementation details. We parameterize the shared feature extractor (meta-learner) h as a pre-trained ResNet-18 (He et al., 2016) network followed by (Li et al., 2019b), and the task-specific learners as the SVM classifier. The competing methods include Reptile

Table 8: Cross-domain recognition accuracy (%) of four held out target datasets on VD using train split in (Li et al., 2019b) for training.

Target	Reptile	CrossGrad	MR	MR-FL	AGG	FC	FC+Tanh	FC+ L^2 -SP	FC+Tanh+ L^2 -SP
FGVC-Aircraft	19.62	19.92	20.91	18.18	18.84	19.02	19.35	19.80	19.60
DTD	37.39	36.54	32.34	35.69	37.52	39.01	39.20	39.04	40.37
Flowers102	58.26	57.84	35.49	53.04	57.64	58.13	58.62	58.43	60.78
UCF101	49.85	45.80	47.34	48.10	47.88	51.12	51.38	51.48	52.12
Ave.	41.28	40.03	34.02	38.75	40.47	41.82	42.14	42.19	43.22

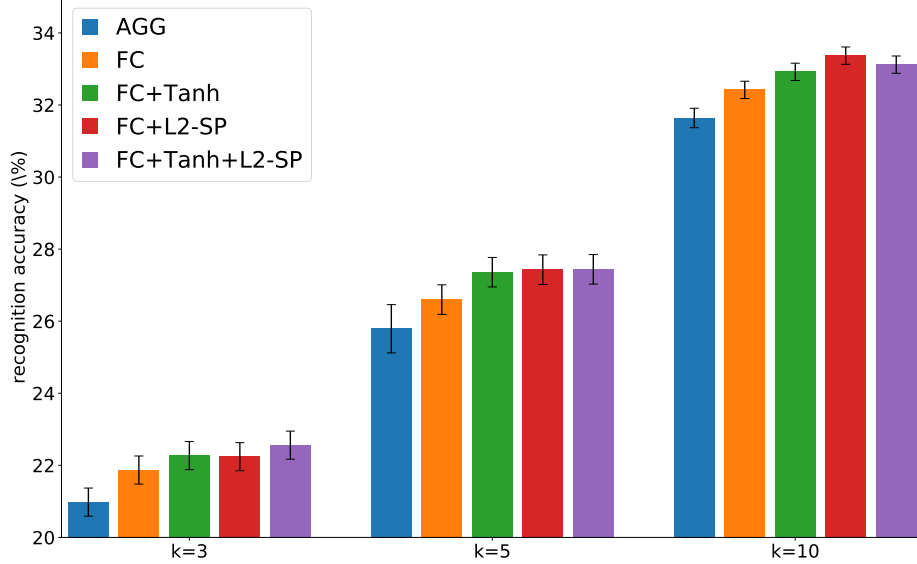


Figure 4: Recognition accuracies (%) of all competing methods averaged over 5 test runs on VD K-shot learning

(Nichol et al., 2018), CrossGrad (Shankar et al., 2018), MetaReg (Balaji et al., 2018; Li et al., 2019b) and FC (Li et al., 2019b). We also follow the experimental setting in (Li et al., 2019b), and train all components end-to-end using the AMSGrad (batch-size/per meta-train domain=64, batch-size/per meta-test domain=32, lr=0.0005, weight decay=0.0001) for 30k iterations where the learning rate decayed in 5K, 12K, 15K, 20K iterations by a factor 5, 10, 50, 100, respectively. Similar to MetaReg (Balaji et al., 2018), after the parameters are trained via meta learning, we fine-tune the network on all source datasets for the final 10k iterations. The Feature-Critic (set embedding variant) is adopted as the MLP type, and is trained with AMSGrad optimizer (lr=0.0001, weight decay=0.00001).

Results. Table 8 shows the classification accuracy on four hold-out target domain. The results of Reptile, CrossGrad, MR, MR-FL are directly taken from (Li et al., 2019b). We re-implemented the result of AGG (Li et al., 2019b) and FC. Different from homogeneous DG, heterogeneous DG assumes that the label spaces among different domains are different, and thus the task-specific learners can not be shared. Thus learning a robust off-the-shelf feature extractor is very important for final performance of heterogeneous DG. From the table, it can be seen that our proposed meta-regularization schemes produce more robust and effective feature extractor for unseen domains compared with other competing methods. Especially, although our methods are easily rebuilt from FC, they can attain evident improvements.

Algorithm 3 Online Methodology-Learning Algorithm: FTML**Input:** T number of tasks $\{D_t = (D_t^{tr}, D_t^{val})\}$.**Initialization:** $h_0 \in \mathcal{H}$ **For** $t = 0$ **to** $T - 1$: (1) The meta-learner incrementally receives a task dataset D_t ; (2) Run the inner online algorithm with meta-learner h_t on D_t^{tr} , returning the predictor model $f_t^{ht} = f_t - \alpha \frac{\partial \mathcal{L}^{task}(f, h_t, D_t^{tr})}{\partial f}$; (3) Incrementally incur the errors $\ell_t(h_t) = L(f_t^{ht}, D_t^{val})$;

(4) Update the meta-learner via

$$h_{t+1} = \arg \min_h \sum_{k=1}^t \ell_k(h). \quad (22)$$

End For

Furthermore, we repeat the evaluation assuming that few-shot (3-shot, 5-shot, 10-shot) samples of training splits are available for SVM training in the meta-test stage. Fig. 4 reports the target domain test accuracies under these settings. It can be easily observed that our proposed meta-regularization schemes provide a consistent improvement over the baseline, and produce a superior off-the-shelf feature representation. It can thus be substantiated that the proposed meta-regularization strategy is hopeful to generally improve the training quality of current meta learning methods to produce more robust and effective feature extractor.

8 Online Methodology-Learning Strategy

Here we further consider a sequential setting where an agent is faced with tasks one after another. Specifically, let's denote T be the number of tasks and for each task $t \in \{1, \dots, T\}$, and $D_t = (D_t^{tr}, D_t^{val})$ be the corresponding task sequence. Our goal is to find an estimator of $h \in \mathcal{H}$ that improves incrementally as the number of observed tasks T increases. Algorithm 3 shows the online methodology-learning process, in which we take inspiration from the form of the “follow the leader algorithm” (FTL) presented in the following Reference Shalev-Shwartz et al. (2012). Correspondingly, we call our advanced algorithm in the meta-level as the “follow the meta-leader algorithm” (FTML) as follows.

The FTML algorithm template can be interpreted as the agent playing the best meta-learner in hindsight if the learning process was to stop at the round t . We will show that under standard assumptions on the losses of analyzing bi-level optimization, the online learning-the-methodology algorithm has the corresponding regret guarantees. Note that we may not have full access to $\ell_t(h_t)$ when it is the population risk, and we only have a finite sized dataset. We can draw upon stochastic optimization strategies to solve the optimization problem.

We make the following assumptions about each loss function in the learning problem for all t . Let $z = (f, h)$ denote all parameters. Assumption 4 (1)-(3) are largely standard in online learning (Shalev-Shwartz et al., 2012; Cesa-Bianchi and Lugosi, 2006), and Assumption 4 (4)-(5) are largely used in bi-level optimization analysis, e.g., (Ji et al., 2021).

Let's make the following assumptions about each loss function in the learning problem for all t . Let $z = (f, h)$ denote all parameters. We first list the following necessary assumptions. Note that in the following, assumptions (1)-(3) are largely standard in online learning (see References (Shalev-Shwartz et al., 2012; Cesa-Bianchi and Lugosi, 2006)), and assumptions (4)-(5) are always used in bi-level optimization analysis, e.g., (Ji et al., 2021).

- Assumption 4** (1) The meta-loss \mathbf{L} has gradients bounded by G , i.e., $\|\nabla \mathbf{L}(f)\| \leq G$.
 (2) The meta-loss \mathbf{L} is β -smooth, i.e., $\|\nabla \mathbf{L}(f) - \nabla \mathbf{L}(f')\| \leq \beta\|f - f'\|$.
 (3) The meta-loss \mathbf{L} is μ -strongly convex, i.e., $\|\nabla \mathbf{L}(f) - \nabla \mathbf{L}(f')\| \geq \mu\|f - f'\|$.
 (4) Suppose the task-specific loss \mathbf{L}^{task} is ρ -strongly convex, i.e., $\left\| \frac{\partial \mathbf{L}^{task}(z)}{\partial f} - \frac{\partial \mathbf{L}^{task}(z')}{\partial f} \right\| \geq \rho\|z - z'\|, \forall z, z'$.
 (5) Suppose the task-specific loss \mathbf{L}^{task} is τ -smooth, i.e., $\left\| \frac{\partial \mathbf{L}^{task}(z)}{\partial f} - \frac{\partial \mathbf{L}^{task}(z')}{\partial f} \right\| \leq \tau\|z - z'\|, \forall z, z'$.

Based on the aforementioned assumptions, we can then deduce the following theorem:

Theorem 9 Suppose that for all t , meta loss and task-specific loss satisfy Assumption 4. Let h_T be the output of Algorithm 3, then FTML enjoys the following regret guarantee

$$\mathcal{R}_T = \sum_{t=1}^T \ell_t(h_t) - \min_h \sum_{t=1}^T \ell_t(h) \leq \frac{4G^2(1 + \log(T))}{\tau}. \quad (23)$$

The above theorem provides the online methodology-learning bound, which tells us that there exists a large family of online meta-learning algorithms that can enjoy sub-linear regret under some mild conditions. Note that our theorem enjoys the same regret guarantees (up to constant factors) as FTL does in the comparable setting (with strongly convex losses) (Shalev-Shwartz et al., 2012).

9 Conclusion

This study has introduced a SLeM framework for understanding and formulating meta learning, in which a meta-learner is extracted to get the hyper-parameter prediction function for machine learning over training task set. The meta-learner is represented as an explicit function mapping from task information to the hyper-parameters involved in the learning process, facilitating the meta-learned meta-learner able to be readily transferred to new tasks to adaptively set their hyper-parametric configurations. Besides, the corresponding learning theory is developed for this SLeM meta learning framework. Very similar to the SRM principle used to improve generalization capability of the extracted learner in conventional machine learning, this theory can help conduct some useful meta-regularization strategies for ameliorating the generalization capability of the extracted meta-learner as a diverse-task-transferable learning methodology. We have further substantiated the beneficial effects brought by these meta-regularization strategies in typical meta learning applications, including few-shot regression, few-shot classification, and domain generalization. Especially, the new meta-regularization schemes can be easily embedded into the current meta learning programs by directly replacing the form of the output activation function in the learning

model, or revising the parameter updating step for the learner or meta-learner from the unregularized solution to the regularized one. These meta-regularized strategies are thus hopeful to be easily and extensively used in more comprehensively meta learning tasks.

In future research, we'll further ameliorate the presented learning theory of SLeM, and especially try to deduce much tighter upper bound for evaluating the task transfer error to reveal more intrinsic generalization insight of this SLeM meta learning framework. Besides, we'll investigate more recent bi-level optimization techniques specifically considering the non-convexity of the inner and outer losses of our meta-learning framework, and ameliorate our theoretical results to improve the learning bounds with relatively loose conditions. Furthermore, we'll make endeavor to explore more helpful meta-regularization schemes useful for more comprehensive and diverse meta learning problems, e.g., (Shu et al., 2023c). Specifically, we'll continue to explore how to effectively design the structure and rectify the parameter learning of meta-learner by enforcing meta-regularization on the model through certain meta-regularization tricks, so as to improve its generalization capability among variant tasks. We believe this will be potentially beneficial to the meta learning field, just like the significance of the SRM principle in the conventional machine learning field.

Acknowledgments and Disclosure of Funding

We thank the anonymous reviewers for insightful comments and suggestions. This research was supported in part by the National Key Research and Development Program of China under Grant 2020YFA0713900, in part by the National Natural Science Foundation of China (NSFC) Project under Grants 61721002, and in part by the Macao Science and Technology Development Fund under Grant 061/2020/A2.

Appendix A. Theoretical Tools

In this part, we will show some important auxiliary theoretical results and assumptions preparing for the proofs of the main theorems.

A.1 Some Properties of Gaussian Complexity

Here, we present some properties of Gaussian complexity, and the proofs can refer to the references therein.

Proposition 10 (Ledoux-Talagrand contraction principle (Ledoux and Talagrand, 2013)) *Let \mathcal{X} be any set, \mathcal{F} be a class of functions: $f : \mathcal{X} \rightarrow \mathbb{R}^d$. Then for N data points, $\mathbf{X} = (x_1, \dots, x_N)^\top$, and a fixed, centered L -Lipschitz function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$, we have*

$$\hat{\mathcal{G}}_{\mathbf{X}}(\phi(\mathcal{F})) \leq L\hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F}); \quad \hat{\mathfrak{R}}_{\mathbf{X}}(\phi(\mathcal{F})) \leq L\hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}),$$

where $\hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F})$ and $\hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F})$ denote the empirical Gaussian complexity and empirical Rademacher complexity, respectively.

Proposition 11 (The relationship between Gaussian complexity and Rademacher complexity (Wainwright, 2019)) *The empirical Rademacher complexity can be lower/upper bounded by empirical Gaussian complexity as follows:*

$$\hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F}) \leq 2\sqrt{\log(N)} \cdot \hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}), \quad \hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}) \leq \sqrt{\frac{\pi}{2}} \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F}) \leq \frac{3}{2} \hat{\mathcal{G}}_{\mathbf{X}}(\mathcal{F}),$$

where $\mathbf{X} = (x_1, \dots, x_N)^\top$.

A.2 Task-averaged Estimation Error

For the single-task setting, the uniform bounds on the estimation error based on the Rademacher complexity can be found in (Mohri et al., 2018). Here we show the estimation error for the multi-task setting.

Theorem 12 *Let \mathcal{F} be a family of function mappings $f_t : \mathcal{X} \rightarrow [0, B]$, and let $\mu_1, \mu_2, \dots, \mu_T$ be probability measures on \mathcal{X} with $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_T) \sim \prod_{t=1}^T (\mu_t)^{n_t}$, where $\mathbf{X}_t = (x_{t1}, \dots, x_{t,n_t})$ for each $t \in [T]$. Let $\mathbf{f} = (f_1, \dots, f_T)$, we define $\hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}] = \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} f_t(x_{tj})$, and $\mathbb{E}[\mathbf{f}] = \mathbb{E}_{\mathbf{X}}[\hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}]]$. Then, for any $\delta > 0$, with probability at least $1 - \delta$ we have*

$$\mathbb{E}[\mathbf{f}] \leq \hat{\mathbb{E}}_S[\mathbf{f}] + 2\hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}^{\otimes T}) + 3\frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t} \log \frac{2}{\delta}},$$

where $\hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}^{\otimes T}) = \mathbf{E}_{\sigma}[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} \sigma_{tj} f_t(x_{tj})]$, and $\mathcal{F}^{\otimes T} = \underbrace{\mathcal{F} \times \mathcal{F} \cdots \mathcal{F}}_T$.

Proof We define the function $\Phi : \mathbf{X} \rightarrow \mathbb{R}$ as

$$\Phi(\mathbf{X}) = \sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} (\mathbb{E}[\mathbf{f}] - \hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}]).$$

Let \mathbf{X} and \mathbf{X}' be two samples differing by exactly one point, say x_{tj} in \mathbf{X} and x'_{tj} in \mathbf{X}' . Then, since the difference of suprema does not exceed the supremum of the difference, we have

$$\begin{aligned} |\Phi(\mathbf{X}') - \Phi(\mathbf{X})| &\leq \left| \sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} (\hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}] - \hat{\mathbb{E}}_{\mathbf{X}'}[\mathbf{f}]) \right| \\ &= \frac{1}{Tn_t} \left| \sup_{f_t \in \mathcal{F}} f_t(\mathbf{x}'_{tj}) - f_t(\mathbf{x}_{tj}) \right| \leq \frac{B}{Tn_t}. \end{aligned}$$

Then, based on the McDiarmid's inequality (Mohri et al., 2018), we have:

$$\begin{aligned} \mathbb{P}(\Phi(\mathbf{X}) - \mathbb{E}_{\mathbf{X}}[\Phi(\mathbf{X})] \geq \epsilon) &\leq \exp \left(\frac{-2\epsilon^2}{\sum_{t=1}^T \sum_{j=1}^{n_t} (\frac{B}{Tn_t})^2} \right) \\ &= \exp \left(\frac{-2T^2\epsilon^2}{B^2 \sum_{t=1}^T \frac{1}{n_t}} \right). \end{aligned}$$

For $\delta > 0$, setting the right-hand side above to be $\delta/2$, with probability at least $1 - \delta/2$, the following holds:

$$\Phi(\mathbf{X}) \leq \mathbb{E}_{\mathbf{X}}[\Phi(\mathbf{X})] + \frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t}} \sqrt{\frac{\log \frac{2}{\delta}}{2}}.$$

We bound the expectation of the right-hand side as follows:

$$\mathbb{E}_{\mathbf{X}}[\Phi(\mathbf{X})] = \mathbb{E}_{\mathbf{X}} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \mathbb{E}[\mathbf{f}] - \hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}] \right] \quad (24)$$

$$= \mathbb{E}_{\mathbf{X}} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \mathbb{E}_{\mathbf{X}'}[\hat{\mathbb{E}}_{\mathbf{X}'}[\mathbf{f}] - \hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}]] \right] \quad (25)$$

$$\leq \mathbb{E}_{\mathbf{X}, \mathbf{X}'} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \hat{\mathbb{E}}_{\mathbf{X}'}[\mathbf{f}] - \hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}] \right] \quad (26)$$

$$= \mathbb{E}_{\mathbf{X}, \mathbf{X}'} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} (f_t(\mathbf{x}'_{tj}) - f_t(\mathbf{x}_{tj})) \right] \quad (27)$$

$$= \mathbb{E}_{\sigma, \mathbf{X}, \mathbf{X}'} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} \sigma_{tj} (f_t(\mathbf{x}'_{tj}) - f_t(\mathbf{x}_{tj})) \right] \quad (28)$$

$$\leq \mathbb{E}_{\sigma, \mathbf{X}'} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} \sigma_{tj} f_t(\mathbf{x}'_{tj}) \right] + \mathbb{E}_{\sigma, \mathbf{X}} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} -\sigma_{tj} f_t(\mathbf{x}_{tj}) \right] \quad (29)$$

$$= 2\mathbb{E}_{\sigma, \mathbf{X}} \left[\sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}} \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} \sigma_{tj} f_t(\mathbf{x}_{tj}) \right] = 2\hat{\mathfrak{N}}_{\mathbf{X}}(\mathcal{F}^{\otimes T}). \quad (30)$$

We introduce a ghost sample (denoted by S'), drawn from the same distribution as our original sample (denoted by S), and thus $\mathbb{E}[\mathbf{f}] = \mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[\mathbf{f}]]$ (Eq.(25)). The inequality (26) holds due to the sub-additivity of the supremum function. In Eq. (28), we introduce Rademacher variables σ_{tj} , which are uniformly distributed independent random variables taking values in $\{-1, 1\}$. Eq. (29) holds by the sub-additivity of the supremum function, and Eq. (30) stems from the definition of Rademacher complexity and the fact that the variables σ_{tj} and $-\sigma_{tj}$ are distributed in the same way.

By using again the McDiarmid's inequality, with probability $1 - \delta/2$ the following holds

$$\mathfrak{R}_{\mathbf{X}}(\mathcal{F}^{\otimes T}) \leq \hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}^{\otimes T}) + \frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t}} \sqrt{\frac{\log \frac{2}{\delta}}{2}}.$$

Finally, we use the union bound to combine above inequalities, which yields that with probability at least $1 - \delta$ it holds:

$$\mathbb{E}[\mathbf{f}] \leq \hat{\mathbb{E}}_S[\mathbf{f}] + 2\hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}^{\otimes T}) + 3\frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t}} \sqrt{\frac{\log \frac{2}{\delta}}{2}}. \quad (31)$$

■

Remark 13 When $n_t = n$ for $t = 1, 2, \dots, T$, the conclusion becomes

$$\mathbb{E}[\mathbf{f}] \leq \hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}] + 2\hat{\mathfrak{R}}_{\mathbf{X}}(\mathcal{F}^{\otimes T}) + 3B \sqrt{\frac{\log \frac{2}{\delta}}{2nT}},$$

where $\hat{\mathbb{E}}_{\mathbf{X}}[\mathbf{f}] = \frac{1}{nT} \sum_{t=1}^T \sum_{j=1}^n f_t(\mathbf{x}_{ti})$. This recovers the result in Theorem 9 of (Maurer et al., 2016).

Remark 14 When $T = 1$, it recovers the result of the single task case (Mohri et al., 2018).

A.3 A Chain Rule for Gaussian Complexity $\mathcal{F}^{\otimes T}(\mathcal{H})$

Suppose we have the training task dataset $\mathbf{D} = \{D_t, t \in [T]\}$, with $D_t = (D_t^{tr}, D_t^{val})$, where $D_t^{tr} = \{(x_{ti}^{(s)}, y_{ti}^{(s)})\}_{i=1}^{m_t}$, $D_t^{val} = \{(x_{tj}^{(q)}, y_{tj}^{(q)})\}_{j=1}^{n_t}$. In this section, we aim to decouple the complexity of learning the class $\mathcal{F}^{\otimes T}(\mathcal{H})$ based on the chain rule tools (Maurer, 2016; Wainwright, 2019; Tripuraneni et al., 2020). We let $f_t^{(h)} = \mathcal{LM}(D_t^{tr}; h)$ in the following for simplicity.

The empirical Gaussian complexity of function class $\mathcal{F}^{\otimes T}(\mathcal{H})$ can be written as

$$\hat{\mathcal{G}}_{D^{val}}(\mathcal{F}_{\mathcal{H}}^{\otimes T}) = \mathbb{E} \sup_{f \in \mathcal{F}^{\otimes T}, h \in \mathcal{H}} \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{j=1}^{n_t} g_{tj} f_t^{(h)}(x_{tj}^{(q)}); \quad g_{tj} \sim \mathcal{N}(0, 1), \quad (32)$$

where $D^{val} = \{D_t^{val}\}_{t=1}^T$. Generally, the Ledoux-Talagrand contraction principle in Proposition 10 can be applied to decoupling the complexity. However, we learn the meta-learner

h and the task-specific learners $f_t, t \in [T]$ simultaneously after observing the data, where learners are not fixed.

To present the result, we firstly list some important theorems for the proof.

Theorem 15 (Dudley’s entropy integral bound (Wainwright, 2019)) *Let $\{\mathbf{X}_\vartheta, \vartheta \in \mathbb{T}\}$ be a zero-mean sub-Gaussian process with respect to induced pseudometric $\rho_{2,\mathbf{X}}$, and $\text{Dis}(\mathbf{X}) = \sup_{\vartheta, \vartheta' \in \mathbb{T}} \rho_{2,\mathbf{X}}(\vartheta, \vartheta')$. Then for any $\delta \in [0, E]$, we have*

$$E \left[\sup_{\vartheta, \vartheta' \in \mathbb{T}} (\mathbf{X}_\vartheta - \mathbf{X}_{\vartheta'}) \right] \leq 2E \left[\sup_{\gamma, \gamma' \in \mathbb{T}, \rho_{2,\mathbf{X}}(\gamma, \gamma') \leq \delta} (\mathbf{X}_\gamma - \mathbf{X}_{\gamma'}) \right] + 32 \int_{\delta/4}^E \sqrt{\log N_{2,\mathbf{X}}(\mu, \rho_{2,\mathbf{X}}, \mathbb{T})} d\mu,$$

where $N_{2,\mathbf{X}}(\mu, \rho_{2,\mathbf{X}}, \mathbb{T})$ denotes the μ -covering number of \mathbb{T} in the metric $\rho_{2,\mathbf{X}}$.

Theorem 16 (Sudakov minoration (Wainwright, 2019)) *Let $\{\mathbf{X}_\vartheta, \vartheta \in \mathbb{T}\}$ be a zero-mean sub-Gaussian process defined on the non-empty set \mathbb{T} . Then*

$$\mathbb{E} \left[\sup_{\vartheta \in \mathbb{T}} \mathbf{X}_\vartheta \right] \geq \sup_{\delta > 0} \frac{\delta}{2} \sqrt{\log M_{2,\mathbf{X}}(\delta, \rho_{2,\mathbf{X}}, \mathbb{T})},$$

where $M_{2,\mathbf{X}}(\delta, \rho_{2,\mathbf{X}}, \mathbb{T})$ is the δ -packing number of \mathbb{T} in the metric $\rho_{2,\mathbf{X}}$.

Now, the following gives the decomposition theorem for Gaussian complexity $\hat{\mathcal{G}}_{D^{\text{val}}}(\mathcal{F}_{\mathcal{H}}^{\otimes T})$.

Theorem 17 *Let the function class \mathcal{F} consist of functions that are ℓ_2 -Lipschitz with constant $L(\mathcal{F})$. Define $\text{Dis}(D^{\text{val}}) = \sup_{h, h'} \rho_{2,D^{\text{val}}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')})$, $\mathbf{f}^{(h)} = (f_1^{(h)}, f_2^{(h)}, \dots, f_T^{(h)})$, and then the empirical Gaussian complexity of function class $\mathcal{F}_{\mathcal{H}}^{\otimes T}$ satisfies*

$$\hat{\mathcal{G}}_{D^{\text{val}}}(\mathcal{F}_{\mathcal{H}}^{\otimes T}) \leq \frac{2\text{Dis}(D^{\text{val}})}{(\sum_{t=1}^T n_t)^2} \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}} + 128 \log(4 \sum_{t=1}^T n_t) \cdot L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{\text{val}}}(\mathcal{H}),$$

where $\rho_{2,D^{\text{val}}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) = \frac{1}{\sum_{t=1}^T n_t} \sum_{i=1}^T \sum_{j=1}^{n_t} (f_t^{(h)}(\mathbf{z}_{tj}^{(q)}) - f_t^{(h')}(\mathbf{z}_{tj}^{(q)}))^2$, $\mathbf{f}^{(h)} = (f_1^{(h)}, f_2^{(h)}, \dots, f_T^{(h)})$, $f_t^{(h)} = \mathcal{LM}(D_t^{\text{tr}}; h)$, and $L(\mathcal{F})$ is the Lipschitz constant of $f^{(h)}$ with respect to h .

Proof For ease of notation we define $N = nT = \sum_{t=1}^T n_t$, $n_t = \beta_t n$ in the following. We define the mean-zero stochastic process $Z_{\mathbf{f}^{(h)}} = \frac{1}{\sqrt{\sum_{t=1}^T n_t}} \sum_{i=1}^T \sum_{j=1}^{n_t} g_{tj} / \beta_t \cdot f_t^{(h)}(x_{tj}^{(q)})$ for a sequence of data points $x_{tj}^{(q)}$. Note that the process $Z_{\mathbf{f}^{(h)}}$ has sub-Gaussian increments, and then $Z_{\mathbf{f}^{(h)}} - Z_{\mathbf{f}^{(h'')}}$ is a sub-Gaussian random variable with parameter $\rho_{2,D^{\text{val}}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) = \frac{1}{\sum_{t=1}^T n_t} \sum_{i=1}^T \sum_{j=1}^{n_t} (f_t^{(h)}(x_{tj}^{(q)}) - f_t^{(h')}(x_{tj}^{(q)}))^2$. Since $Z_{\mathbf{f}^{(h)}}$ is a mean-zero stochastic process, we have

$$\mathbb{E} \left[\sup_{\mathbf{f}^{(h)} \in \mathcal{F}_{\mathcal{H}}^{\otimes T}} Z_{\mathbf{f}^{(h)}} \right] = \mathbb{E} \left[\sup_{\mathbf{f}^{(h)} \in \mathcal{F}_{\mathcal{H}}^{\otimes T}} Z_{\mathbf{f}^{(h)}} - Z_{\mathbf{f}^{(h')}} \right] \leq \mathbb{E} \left[\sup_{\mathbf{f}^{(h)}, \mathbf{f}^{(h')} \in \mathcal{F}_{\mathcal{H}}^{\otimes T}} Z_{\mathbf{f}^{(h)}} - Z_{\mathbf{f}^{(h')}} \right].$$

According to the Dudley entropy integral bound (Theorem 15), we have

$$\begin{aligned} & \mathbb{E} \left[\sup_{\mathbf{f}^{(h)}, \mathbf{f}^{(h')} \in \mathcal{F}_{\mathcal{H}}^{\otimes T}} Z_{\mathbf{f}^{(h)}} - Z_{\mathbf{f}^{(h')}} \right] \\ & \leq 2\mathbb{E} \left[\sup_{\substack{\mathbf{f}^{(h)}, \mathbf{f}^{(h')} \in \mathcal{F}_{\mathcal{H}}^{\otimes T} \\ \rho_{2, \mathbf{Z}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) \leq \delta}} Z_{\mathbf{f}^{(h)}} - Z_{\mathbf{f}^{(h')}} \right] + 32 \int_{\delta/4}^{Dis(D^{val})} \sqrt{\log N_{D^{val}}(\mu; \rho_{2, D^{val}}, \mathcal{F}_{\mathcal{H}}^{\otimes T})} d\mu. \end{aligned}$$

The first term in the right hand can be computed as:

$$\begin{aligned} & \mathbb{E} \sup_{\substack{\mathbf{f}^{(h)}, \mathbf{f}^{(h')} \in \mathcal{F}_{\mathcal{H}}^{\otimes T} \\ \rho_{2, D^{val}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) \leq \delta}} Z_{\mathbf{f}^{(h)}} - Z_{\mathbf{f}^{(h')}} \\ & = \mathbb{E} \sup_{\substack{\mathbf{f}^{(h)}, \mathbf{f}^{(h')} \in \mathcal{F}_{\mathcal{H}}^{\otimes T} \\ \rho_{2, D^{val}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) \leq \delta}} \frac{1}{\sqrt{\sum_{t=1}^T n_t}} \sum_{i=1}^T \sum_{j=1}^{n_t} g_{tj} / \beta_t \cdot (f_t^{(h)}(x_{tj}^{(q)}) - f_t^{(h')}(x_{tj}^{(q)})) \\ & \leq \mathbb{E} \sup_{\mathbf{v}: \|\mathbf{v}\|_2 \leq \delta} \sqrt{\sum_{t=1}^T \sum_{j=1}^{n_t} (g_{tj} / \beta_t)^2 \cdot \|\mathbf{v}\|} \\ & \leq \sqrt{\sum_{t=1}^T \sum_{j=1}^{n_t} \frac{1}{\beta_t^2}} \cdot \delta \leq \delta \sqrt{n \sum_{t=1}^T \frac{1}{\beta_t}}. \end{aligned}$$

Let $C_{\mathcal{H}}$ be a covering of the function space \mathcal{H} in the empirical ℓ_2 -norm at scale ϵ_1 with respect to D^{val} . We will claim that $C_{\mathcal{F}_{\mathcal{H}}^{\otimes T}} = \cup_{h \in \mathcal{H}} C_{\mathcal{F}_h^{\otimes T}}$ is an $\epsilon_1 \cdot L(\mathcal{F})$ -covering for the function space $\mathcal{F}_{\mathcal{H}}^{\otimes T}$ in the empirical ℓ_2 -norm. To see this, for any $h \in \mathcal{H}, \mathbf{f} \in \mathcal{F}^{\otimes T}$, we let $h' \in C_{\mathcal{H}}$ be ϵ_1 -close to h in \mathcal{H} . By construction we have $\mathbf{f}^{(h')} \in C_{\mathcal{F}_{\mathcal{H}}^{\otimes T}}$. The following inequality establishes the claim,

$$\begin{aligned} \rho_{2, D^{val}}(\mathbf{f}^{(h)}, \mathbf{f}^{(h')}) &= \frac{1}{\sum_{t=1}^T n_t} \sum_{t=1}^T \sum_{j=1}^{n_t} (f_t^{(h)}(x_{tj}^{(q)}) - f_t^{(h')}(x_{tj}^{(q)}))^2 \\ &\leq L(\mathcal{F}) \cdot \frac{1}{\sum_{t=1}^T n_t} \sum_{t=1}^T \sum_{j=1}^{n_t} (h(D_t^{val}) - h'(D_t^{val}))^2 \\ &= L(\mathcal{F}) \cdot \rho_{2, D^{val}}(h, h') \leq \epsilon_1 \cdot L(\mathcal{F}), \end{aligned}$$

where the first inequality holds since the Lipschitz property of the learning method \mathcal{LM} with respect to h . Now, the cardinaity of $C_{\mathcal{F}_{\mathcal{H}}^{\otimes T}}$ can be bounded as

$$|C_{\mathcal{F}_{\mathcal{H}}^{\otimes T}}| = \sum_{h \in C_{\mathcal{H}}} |C_{\mathcal{F}_h^{\otimes T}}| \leq |C_{\mathcal{H}}| \cdot \max_{h \in C_{\mathcal{H}}} |C_{\mathcal{F}_h^{\otimes T}}|.$$

Thus, it follows that

$$\log N_{2,D^{val}}(\epsilon_1 \cdot L(\mathcal{F}), \rho_{2,D^{val}}, C_{\mathcal{F}_h}^{\otimes T}) \leq \log N_{2,D^{val}}(\epsilon_1, \rho_{2,D^{val}}, \mathcal{H}).$$

We define $\epsilon_1 = \frac{\epsilon}{L(\mathcal{F})}$ to show that

$$\int_{\delta/4}^{Dis(D^{val})} \sqrt{\log N_{2,D^{val}}(\epsilon, \rho_{2,D^{val}}, C_{\mathcal{F}_h}^{\otimes T})} d\epsilon \leq \int_{\delta/4}^{Dis(D^{val})} \sqrt{\log N_{2,D^{val}}(\epsilon/(L(\mathcal{F})), \rho_{2,D^{val}}, \mathcal{H})} d\epsilon.$$

Observing that the covering number can be bounded by packing number $M(\epsilon, \rho_{2,D^{val}}, \mathcal{H})$, i.e., $N(\epsilon, \rho_{2,D^{val}}, \mathcal{H}) \leq M(\epsilon, \rho_{2,D^{val}}, \mathcal{H})$, we can then employ Sudakov minoration Theorem (Theorem 16) to upper bound the covering number by Gaussian complexity. For the covering number of \mathcal{H} , $N_{2,D^{val}}(\epsilon/(L(\mathcal{F})), \rho_{2,D^{val}}, \mathcal{H})$, we can apply the theorem with mean-zero Gaussian process $\frac{1}{\sqrt{\sum_{t=1}^T n_t}} \sum_{t=1}^T \sum_{j=1}^{n_t} g_{tj} \cdot h(D_t^{val})$, with $g_{tj} \sim \mathcal{N}(0, 1)$,

$$\log N_{2,D^{val}}(\epsilon/(L(\mathcal{F})), \rho_{2,D^{val}}, \mathcal{H}) d\epsilon \leq 4 \left(\frac{\sqrt{\sum_{t=1}^T n_t} \hat{\mathcal{G}}_{D^{val}}(\mathcal{H})}{\epsilon/(L(\mathcal{F}))} \right)^2 = 4 \left(\frac{\sqrt{\sum_{t=1}^T n_t} L(\mathcal{F}) \hat{\mathcal{G}}_{D^{val}}(\mathcal{H})}{\epsilon} \right)^2.$$

Finally, we can show that

$$\begin{aligned} \hat{\mathcal{G}}_{D^{val}}(\mathcal{F}_h^{\otimes T}) &= \mathbb{E} \left[\frac{1}{\sqrt{nT}} \sup_{\mathbf{f} \in \mathcal{F}^{\otimes T}, h \in \mathcal{H}} Z_{\mathbf{f}(h)} \right] \\ &\leq \frac{1}{\sqrt{nT}} \left(2\delta \sqrt{n \sum_{t=1}^T \frac{1}{\beta_t}} + 64L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) \cdot \sqrt{\sum_{t=1}^T n_t \int_{\delta/4}^{Dis(D^{val})} \frac{1}{\mu} d\mu} \right) \\ &\leq 2\delta \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}} + 64 \log\left(\frac{4Dis(D^{val})}{\delta}\right) \left(L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) \sqrt{\frac{\sum_{t=1}^T n_t}{nT}} \right) \\ &= 2\delta \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}} + 64 \log\left(\frac{4Dis(D^{val})}{\delta}\right) \cdot L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}). \end{aligned}$$

Let $\delta = \frac{Dis(D^{val})}{(nT)^2}$, we have

$$\begin{aligned} \hat{\mathcal{G}}_{D^{val}}(\mathcal{F}_h^{\otimes T}) &\leq \frac{2Dis(D^{val})}{(nT)^2} \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}} + 128 \log(4nT) \cdot L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) \\ &= \frac{2Dis(D^{val})}{(\sum_{t=1}^T n_t)^2} \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}} + 128 \log(4 \sum_{t=1}^T n_t) \cdot L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}). \end{aligned}$$

■

Appendix B. Proof of Theorem 2 : Excess Task Average Risk in Meta-Training Stage

Proof Recall that

$$R_\Gamma(h) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^{n_t}} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val}),$$

$$\hat{R}_D(h) = \frac{1}{T} \sum_{t=1}^T \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val}),$$

and we denote $\tilde{R}_\Gamma(h), \bar{R}_\Gamma(h)$ as

$$\tilde{R}_\Gamma(h) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^{n_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val})$$

$$\bar{R}_\Gamma(h) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val}).$$

For the task average excess risk $R_{train}(\hat{\mathbf{f}}, \hat{h}) - R_{train}(\mathbf{f}^*, h^*)$, we have the following decomposition

$$\begin{aligned} & R_{train}(\hat{h}) - R_{train}(h^*) \\ &= R_\Gamma(\hat{h}) - R_\Gamma(h^*) \\ &= \underbrace{R_\Gamma(\hat{h}) - \bar{R}_\Gamma(\hat{h})}_a + \underbrace{\bar{R}_\Gamma(\hat{h}) - \hat{R}_\Gamma(\hat{h})}_b + \underbrace{\hat{R}_\Gamma(\hat{h}) - \hat{R}_\Gamma(h^*)}_c + \underbrace{\hat{R}_\Gamma(h^*) - \bar{R}_\Gamma(h^*)}_d + \underbrace{\bar{R}_\Gamma(h^*) - R_\Gamma(h^*)}_e. \end{aligned}$$

For the terms (a) and (e), we have

$$\begin{aligned} & (a) + (e) \\ & \leq \sup_{h \in \mathcal{H}} 2 \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^{n_t}} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val}) - \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h), D_t^{val}) \right| \\ & \leq 4L\hat{\mathfrak{R}}_{D^{val}}(\mathcal{F}_{\mathcal{H}}^{\otimes T}) + 3\frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t}} \sqrt{\frac{\log \frac{2}{\delta}}{2}} \leq 6L\hat{\mathcal{G}}_{D^{val}}(\mathcal{F}_{\mathcal{H}}^{\otimes T}) + 6\frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t}} \sqrt{\frac{\log \frac{2}{\delta}}{2}} \\ & \leq 768L \log(4 \sum_{t=1}^T n_t) \cdot L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{D^{val}}(\mathcal{H}) + \frac{12LDis(D^{val})}{(\sum_{t=1}^T n_t)^2} \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}} + 6\frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t}} \sqrt{\frac{\log \frac{2}{\delta}}{2}}, \end{aligned}$$

where $D^{val} = \{D_t^{val}\}_{t=1}^T$, and $|D_t^{val}| = n_t$, for $t \in [T]$. The second inequality we use the task-averaged estimation error in Theorem 12 and the Ledoux-Talagrand contraction principle in Proposition 10. Proposition 11 is employed for the third inequality. And the last inequality holds for Theorem 17.

For the term (b), we have

$$\begin{aligned}
& \bar{R}_\Gamma(\hat{h}) - \hat{R}_\Gamma(\hat{h}) \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{val}) - \frac{1}{T} \sum_{t=1}^T \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{val}) \\
&= \underbrace{\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{val}) - \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{tr})}_{b1} \\
&+ \underbrace{\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{tr}) - \frac{1}{T} \sum_{t=1}^T \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{tr})}_{b2} \\
&+ \underbrace{\frac{1}{T} \sum_{t=1}^T \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{tr}) - \frac{1}{T} \sum_{t=1}^T \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{val})}_{b3}.
\end{aligned}$$

For the terms (b1) and (b3), we have

$$(b1) + (b3) \leq \frac{2}{T} \sum_{t=1}^T d_{\mathcal{F}}(\mu_t^s, \mu_t^q),$$

where $d_{\mathcal{F}}(\mu_t^s, \mu_t^q)$ denotes the discrepancy divergence (Ben-David et al., 2010) between samples from the probability distributions μ_t^s and μ_t^q with respect to the hypothesis class \mathcal{F} ,

$$d_{\mathcal{F}}(\mu_t^s, \mu_t^q) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^{n_t}} \mathbf{L}(f, D_t^{val}) - \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(f, D_t^{tr}) \right|.$$

For the term (b2), we denote $\hat{f}_t^{(h^*)} = \mathcal{LM}(D_t^{tr}; h^*)$, and then we have

$$\begin{aligned}
b2 &= \frac{1}{T} \sum_{t=1}^T \mathbf{L}(\hat{f}_t^{(h^*)}, D_t^{tr}) - \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^{m_t}} \mathbf{L}(\hat{f}_t^{(h^*)}, D_t^{tr}) \\
&\leq \frac{1}{T} \sum_{t=1}^T \left(2L\hat{\mathfrak{R}}_{D_t^{tr}}(\mathcal{F}) + 3B\sqrt{\frac{\log \frac{2}{\delta}}{m_t}} \right) \leq \frac{3L}{T} \sum_{t=1}^T \hat{\mathcal{G}}_{D_t^{tr}}(\mathcal{F}) + \frac{3B}{T} \sum_{t=1}^T \sqrt{\frac{\log \frac{2}{\delta}}{m_t}}.
\end{aligned}$$

Thus we have

$$b \leq \frac{2}{T} \sum_{t=1}^T d_{\mathcal{F}}(\mu_t^s, \mu_t^q) + \frac{3L}{T} \sum_{t=1}^T \hat{\mathcal{G}}_{D_t^{tr}}(\mathcal{F}) + \frac{3B}{T} \sum_{t=1}^T \sqrt{\frac{\log \frac{2}{\delta}}{m_t}}. \quad (33)$$

Similar process can also be applied to the term (d), and (d) can be bounded by

$$d \leq \frac{2}{T} \sum_{t=1}^T d_{\mathcal{F}}(\mu_t^s, \mu_t^q) + \frac{3L}{T} \sum_{t=1}^T \hat{\mathcal{G}}_{D_t^{tr}}(\mathcal{F}) + \frac{3B}{T} \sum_{t=1}^T \sqrt{\frac{\log \frac{2}{\delta}}{m_t}}. \quad (34)$$

For the term (c), according to $\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_D(h)$, we have $c \leq 0$.
Combining the above results from the terms (a) to (e), we have

$$\begin{aligned} R_{train}(\hat{h}) - R_{train}(h^*) &\leq 768L \log(4 \sum_{t=1}^T n_t) \cdot L(\mathcal{F}) \cdot \hat{\mathcal{G}}_{\Gamma(q)}(\mathcal{H}) + \frac{6L}{T} \sum_{t=1}^T \hat{\mathcal{G}}_{D_t^{tr}}(\mathcal{F}) \\ &+ \frac{4}{T} \sum_{t=1}^T d_{\mathcal{F}}(D_t^{(tr)}, D_t^{(val)}) + 6 \frac{B}{T} \sqrt{\sum_{t=1}^T \frac{1}{n_t} \sqrt{\log \frac{2}{\delta}}} + \frac{6B}{T} \sum_{t=1}^T \sqrt{\frac{\log \frac{2}{\delta}}{m_t}} + \frac{12L \text{Dis}(D^{val})}{(\sum_{t=1}^T n_t)^2} \sqrt{\sum_{t=1}^T \frac{1}{\beta_t T}}. \end{aligned}$$

■

Appendix C. Proof of Theorem 3: Excess Transfer Error in Meta-Test Stage

Proof Recall that

$$\begin{aligned} &R_{test}(\hat{h}) - R_{test}(h^*) \\ &= \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{val}) - \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; h^*), D_{\mu}^{val}) \\ &= \underbrace{\mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{val}) - \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{val})}_{a} \\ &+ \underbrace{\mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{val}) - \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; h^*), D_{\mu}^{val})}_{b} \\ &+ \underbrace{\mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; h^*), D_{\mu}^{val}) - \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; h^*), D_{\mu}^{val})}_{c}. \end{aligned}$$

We now bound the first term (a) by

$$\begin{aligned} (a) &= \underbrace{\mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{val}) - \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{tr'} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'})}_{a1} \\ &+ \underbrace{\mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{tr'} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'}) - \mathbb{E}_{\mu \sim \eta} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'})}_{a2} \\ &+ \underbrace{\mathbb{E}_{\mu \sim \eta} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'}) - \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'})}_{a3} \\ &+ \underbrace{\mathbb{E}_{\eta} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'}) - \mathbb{E}_{\eta} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbb{E}_{D_{\mu}^{tr'} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'})}_{a4} \\ &+ \underbrace{\mathbb{E}_{\eta} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbb{E}_{D_{\mu}^{tr'} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{tr'}) - \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbf{L}(\mathcal{LM}(D_{\mu}^{tr}; \hat{h}), D_{\mu}^{val})}_{a5}, \end{aligned}$$

where $D_{\mu}^{tr'}$ is equivalent to D_{μ}^{tr} .

For the terms (a1) and (a5), we have

$$(a1) + (a5) \leq 2\mathbb{E}_{\mu \sim \eta} d_{\mathcal{F}}(\mu^s, \mu^q),$$

where $d_{\mathcal{F}}(\mu^s, \mu^q)$ denotes the discrepancy divergence (Ben-David et al., 2010) between samples from the probability distributions μ^s and μ^q with respect to the hypothesis class \mathcal{F} ,

$$d_{\mathcal{F}}(\mu^s, \mu^q) = \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{D_{\mu}^{val} \sim (\mu^q)^{n_{\mu}}} \mathbf{L}(f, D_{\mu}^{val}) - \mathbb{E}_{D_{\mu}^{tr} \sim (\mu^s)^{m_{\mu}}} \mathbf{L}(f, D_{\mu}^{tr}) \right|.$$

For the term (a2), we have

$$(a2) \leq 2L\hat{\mathfrak{R}}_{D_\mu^{tr}}(\mathcal{F}) + 3B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}} \leq 3L\hat{\mathcal{G}}_{D_\mu^{tr}}(\mathcal{F}) + 3B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}.$$

Similarly, For the term (a4), we have

$$(a4) \leq 3L\hat{\mathcal{G}}_{D_\mu^{tr}}(\mathcal{F}) + 3B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}.$$

For the term (a3), suppose that the outer loss can be upper bounded by the inner loss, and then we have $a3 < 0$ according to the definition of \mathcal{LM} function, i.e., $\mathcal{LM}(D_\mu^{tr}; \hat{h}(D_\mu^{tr}))$ minimizes the $\mathbf{L}(f, D_\mu^{tr})$ equipped with hyper-parameters $\hat{h}(D_\mu^{tr})$.

Therefore, combining the above analysis, the term (a) can be upper bounded as

$$(a) \leq 2\mathbb{E}_{\mu \sim \eta} d_{\mathcal{F}}(\mu^s, \mu^q) + 6L\hat{\mathcal{G}}_{D_\mu^{tr}}(\mathcal{F}) + 6B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}.$$

For the term (b), according to the task diversity definition (Assumption 3), we have

$$(b) = R_\eta(\hat{h}) - R_\eta(h^*) \leq \alpha \left(R_{train}(\hat{h}) - R_{train}(h^*) \right) + \beta.$$

For the term (c), according to definition of h^* , we have $(c) \leq 0$.

Combining above results from term (a) to (c), we have

$$\begin{aligned} & R_{test}(\hat{h}) - R_{test}(h^*) \\ & \leq \alpha \left(R_{train}(\hat{\mathbf{f}}, \hat{h}) - R_{train}(\mathbf{f}^*, h^*) \right) + \beta + 6L\hat{\mathcal{G}}_{D_\mu^{tr}}(\mathcal{F}) + 2\mathbb{E}_{\mu \sim \eta} d_{\mathcal{F}}(\mu^s, \mu^q) + 6B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}. \end{aligned}$$

■

Appendix D. Proof of the Proposition 4

Proof We denote

$$\tilde{\mathbf{w}} = \arg \min_{\mathbf{w}} \mathbb{E}_{(x^{(s)}, y^{(s)}) \in \mu^s} \ell(\mathbf{w}^T \hat{h}(x^{(s)}), y^{(s)})$$

and

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathbb{E}_{(x^{(s)}, y^{(s)}) \in \mu^s} \ell(\mathbf{w}^T h^*(x^{(s)}), y^{(s)}).$$

Observe that

$$\begin{aligned}
 & R_\eta(\hat{h}) - R_\eta(h^*) \\
 &= \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_\mu^{val} \sim (\mu^q)^n} \mathbb{E}_{D_\mu^{tr} \sim (\mu^s)^m} \mathbf{L}(\mathcal{LM}(D_\mu^{tr}; \hat{h}), D_\mu^{val}) - \\
 &\quad \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{D_\mu^{val} \sim (\mu^q)^n} \mathbb{E}_{D_\mu^{tr} \sim (\mu^s)^m} \mathbf{L}(\mathcal{LM}(D_\mu^{tr}; h^*), D_\mu^{val}) \\
 &= \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} \left\{ \left| \tilde{\mathbf{w}}^\top \hat{h}(x) - \mathbf{w}^{*\top} h^*(x) \right|^2 \right\}, \\
 &= \sup_{\mathbf{w}_0} \inf_{\mathbf{w}} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} \left\{ \left| \mathbf{w}^\top \hat{h}(x) - \mathbf{w}_0^\top h^*(x) \right|^2 \right\} \\
 &= \sup_{\mathbf{w}_0} \inf_{\tilde{\mathbf{w}}} \left\{ [\tilde{\mathbf{w}}^\top \quad -\mathbf{w}_0^\top] \Lambda \begin{bmatrix} \tilde{\mathbf{w}} \\ -\mathbf{w}_0 \end{bmatrix} \right\},
 \end{aligned}$$

where Λ is defined as

$$\Lambda(\hat{h}, h^*) = \begin{bmatrix} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [\hat{h}(x) \hat{h}(x)^\top] & \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [\hat{h}(x) h^*(x)^\top] \\ \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [h^*(x) \hat{h}(x)^\top] & \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [h^*(x) h^*(x)^\top] \end{bmatrix} = \begin{bmatrix} \mathbf{G}_{\hat{h}\hat{h}} & \mathbf{G}_{\hat{h}h^*} \\ \mathbf{G}_{h^*\hat{h}} & \mathbf{G}_{h^*h^*} \end{bmatrix}.$$

According to the partial minimization of a convex quadratic form (Boyd et al., 2004), we have

$$\inf_{\tilde{\mathbf{w}}} \left\{ [\tilde{\mathbf{w}}^\top \quad -\mathbf{w}_0^\top] \Lambda \begin{bmatrix} \tilde{\mathbf{w}} \\ -\mathbf{w}_0 \end{bmatrix} \right\} = \mathbf{w}_0^\top \Lambda_S(\hat{h}, h^*) \mathbf{w}_0,$$

where $\Lambda_S(\hat{h}, h^*) = \mathbf{G}_{h^*h^*} - \mathbf{G}_{h^*\hat{h}}(\mathbf{G}_{\hat{h}\hat{h}})^\dagger \mathbf{G}_{\hat{h}h^*}$ is the generalized Schur complement of the representation of h^* with respect to \hat{h} . Furthermore, according to the variational characterization of the singular value, we have

$$\sup_{\mathbf{w}_0: \|\mathbf{w}_0\| \leq M} \mathbf{w}_0^\top \Lambda_S \mathbf{w}_0 = M \sigma_1(\Lambda_S(\hat{h}, h^*)),$$

where σ_1 denotes the maximal singular value.

Moreover, we denote

$$\tilde{\mathbf{w}}_t = \arg \min_{\mathbf{w}_t} \mathbb{E}_{(x_t^{(s)}, y_t^{(s)}) \in \mu_t^s} \ell(\mathbf{w}_t^T \hat{h}(x_t^{(s)}), y_t^{(s)}), t \in [T]$$

and

$$\mathbf{w}_t^* = \arg \min_{\mathbf{w}_t} \mathbb{E}_{(x_t^{(s)}, y_t^{(s)}) \in \mu^s} \ell(\mathbf{w}_t^T h^*(x_t^{(s)}), y_t^{(s)}), t \in [T].$$

Thus we have the following derivation:

$$\begin{aligned}
& R_{train}(\hat{h}) - R_{train}(h^*) \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^n} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^m} \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{val}) - \\
&\quad \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^n} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^m} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h^*), D_t^{val}) \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\mu_t \sim \eta} \mathbb{E}_{x \sim \mu_t^q} \left\{ \left| \tilde{\mathbf{w}}_t^\top \hat{h}(x) - \mathbf{w}_t^{*T} h^*(x) \right|^2 \right\} \\
&= \frac{1}{T} \sum_{t=1}^T \inf_{\mathbf{w}_t} \mathbb{E}_{\mu_t \sim \eta} \mathbb{E}_{x \sim \mu_t^q} \left\{ \left| \mathbf{w}_t^\top \hat{h}(x) - \mathbf{w}_t^{*T} h^*(x) \right|^2 \right\} \\
&= \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t^{*T} \Lambda_S(\hat{h}, h^*) \mathbf{w}_t^* = \text{tr}(\Lambda_S(\hat{h}, h^*) \mathbf{P}^\top \mathbf{P} / T) = \text{tr}(\Lambda_S(\hat{h}, h^*) \mathbf{K}),
\end{aligned}$$

where $\mathbf{K} = \mathbf{P}^\top \mathbf{P} / T$. Since $\Lambda_S \succeq 0$, and $\mathbf{K} \succeq 0$, through the Von-Neumann trace inequality, we have that

$$\begin{aligned}
\text{tr}(\Lambda_S(\hat{h}, h^*) \mathbf{K}) &\geq \sum_{i=1}^{d_L} \sigma_i(\Lambda_S(\hat{h}, h^*)) \sigma_{d_L-i+1}(\mathbf{K}) \geq \sum_{i=1}^{d_L} \sigma_i(\Lambda_S(\hat{h}, h^*)) \sigma_{d_L}(\mathbf{K}) \\
&= \text{tr}(\Lambda_S(\hat{h}, h^*)) \sigma_{d_L}(\mathbf{K}) \geq \sigma_1(\Lambda_S(\hat{h}, h^*)) \sigma_{d_L}(\mathbf{K}).
\end{aligned}$$

Now, we can deduce that

$$R_\eta(\hat{h}) - R_\eta(h^*) \leq \frac{M}{\sigma_{d_L}(\mathbf{K})} \left\{ R_{train}(\hat{h}) - R_{train}(h^*) \right\}.$$

Based on the above derivation, we can have the conclusion. ■

Appendix E. Proof of the Proposition 6

Proof Denote $z_i(x) = f(h(x))_i$, the loss function can be written as

$$\ell(f(h(x)), y) = - \sum_i y_i \log \left(\frac{e^{z_i}}{\sum_k e^{z_k}} \right). \quad (35)$$

Let $a_i = \frac{e^{z_i}}{\sum_k e^{z_k}}$,

$$\frac{\partial \ell}{\partial z_i} = \sum_j \left(\frac{\partial \ell}{\partial a_j} \frac{\partial a_j}{\partial z_i} \right), \quad (36)$$

where $\frac{\partial \ell}{\partial a_j} = -\frac{y_j}{a_j}$, for $\frac{\partial a_j}{\partial z_i}$, if $i = j$, $\frac{\partial a_j}{\partial z_i} = \frac{\partial(\frac{e^{z_i}}{\sum_k e^{z_k}})}{\partial z_i} = \frac{e^{z_i} \sum_k e^{z_k} - (e^{z_i})^2}{(\sum_k e^{z_k})^2} = a_i(1 - a_i)$; otherwise, $i \neq j$, $\frac{\partial a_j}{\partial z_i} = \frac{\partial(\frac{e^{z_j}}{\sum_k e^{z_k}})}{\partial z_i} = \frac{-e^{z_i} e^{z_j}}{(\sum_k e^{z_k})^2} = -a_i a_j$. Therefore,

$$\begin{aligned} \frac{\partial \ell}{\partial z_i} &= -\sum_{j \neq i} \left(\frac{y_j}{a_j} (-a_i a_j) \right) - \frac{y_i}{a_i} (1 - a_i) \\ &= a_i \sum_j y_j - y_i = a_i - y_i, \end{aligned}$$

since $a_i \in [0, 1]$, $y_i \in \{0, 1\}$, we can obtain $\frac{\partial \ell}{\partial z_i} \in [-1, 1]$. Therefore, we can demonstrate that the loss function is 1-Lipschitz with respect to $f(h(X))$. \blacksquare

Appendix F. Proof of the Proposition 7

We first present some essential theoretical results preparing for proving the Proposition 7.

Lemma 18 *Consider the following generalized linear model*

$$p(\mathbf{y}|\boldsymbol{\eta}) = h(\mathbf{y}) \exp(\boldsymbol{\eta}^T t(\mathbf{y}) - a(\boldsymbol{\eta})), \quad (37)$$

Then we have

$$(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \frac{a''(\mathbf{c}_1)}{2} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \leq KL(p(\mathbf{y}|\boldsymbol{\eta})|p(\mathbf{y}|\hat{\boldsymbol{\eta}})) \leq (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \frac{a''(\mathbf{c}_2)}{2} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}), \quad (38)$$

where $\mathbf{c}_1 = \inf_{\mathbf{c} \in [\hat{\boldsymbol{\eta}}, \boldsymbol{\eta}]} a''(\mathbf{c})$, $\mathbf{c}_2 = \sup_{\mathbf{c} \in [\hat{\boldsymbol{\eta}}, \boldsymbol{\eta}]} a''(\mathbf{c})$.

Proof Observe that

$$\begin{aligned} KL(p(\mathbf{y}|\boldsymbol{\eta})|p(\mathbf{y}|\hat{\boldsymbol{\eta}})) &= \int (p(\mathbf{y}|\boldsymbol{\eta}) \log \frac{p(\mathbf{y}|\boldsymbol{\eta})}{p(\mathbf{y}|\hat{\boldsymbol{\eta}})}) d\mathbf{y} \\ &= \int p(\mathbf{y}|\boldsymbol{\eta})^T [t(\mathbf{y})(\boldsymbol{\eta} - \hat{\boldsymbol{\eta}}) + a(\hat{\boldsymbol{\eta}}) - a(\boldsymbol{\eta})] d\mathbf{y} \\ &= a'(\boldsymbol{\eta})^T (\boldsymbol{\eta} - \hat{\boldsymbol{\eta}}) + a(\hat{\boldsymbol{\eta}}) - a(\boldsymbol{\eta}), \end{aligned}$$

where $a'(\boldsymbol{\eta}) = \int t(\mathbf{y})^T p(\mathbf{y}|\boldsymbol{\eta}) d\mathbf{y}$. Based on the Taylor's theorem we have that

$$a(\hat{\boldsymbol{\eta}}) = a(\boldsymbol{\eta}) + a'(\boldsymbol{\eta})^T (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) + (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \frac{a''(\mathbf{c})}{2} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}),$$

where $\mathbf{c} \in [\hat{\boldsymbol{\eta}}, \boldsymbol{\eta}]$. Therefore, we can obtain

$$KL(p(\mathbf{y}|\boldsymbol{\eta})|p(\mathbf{y}|\hat{\boldsymbol{\eta}})) = (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \frac{a''(\mathbf{c})}{2} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}).$$

We can then obtain

$$(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \frac{a''(\mathbf{c}_1)}{2} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \leq KL(p(\mathbf{y}|\boldsymbol{\eta})|p(\mathbf{y}|\hat{\boldsymbol{\eta}})) \leq (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \frac{a''(\mathbf{c}_2)}{2} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}),$$

where $\mathbf{c}_1 = \inf_{\mathbf{c} \in [\hat{\boldsymbol{\eta}}, \boldsymbol{\eta}]} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T a''(\mathbf{c}) (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})$, $\mathbf{c}_2 = \sup_{\mathbf{c} \in [\hat{\boldsymbol{\eta}}, \boldsymbol{\eta}]} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T a''(\mathbf{c}) (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})$. \blacksquare

Remark 19 If the data generating model satisfies the conditional likelihood defined in Eq. (20), for the cross-entropy loss we have

$$\begin{aligned} & \mathbb{E}_{(x,y) \in D_t^{val}} [\ell(f(h(x)), y) - \ell(f(h(x)), y)] \\ &= \mathbb{E}_x [KL(\text{Multi}(\text{Softmax}(f(h(x)))) | KL(\text{Multi}(\text{Softmax}(f(h(x)))))], \end{aligned}$$

where Multi denote the multinomial distribution.

When consider the multinomial distribution, the generalized linear model in Eq. (37) satisfies that $h(\mathbf{y}) = 1, t(\mathbf{y}) = \mathbf{y}, a(\boldsymbol{\eta}) = \log(\sum_{k=1}^K \exp(\eta_k)), \boldsymbol{\eta} = (\eta_1, \dots, \eta_K)^\top$. Then, we have

$$\begin{aligned} a'(\boldsymbol{\eta}) &= \left(\frac{\exp(\eta_1)}{\sum_{k=1}^K \exp(\eta_k)}, \dots, \frac{\exp(\eta_K)}{\sum_{k=1}^K \exp(\eta_k)} \right)^\top, \\ a''(\boldsymbol{\eta}) &= \begin{pmatrix} \frac{\frac{\exp(\eta_1)}{S} - \frac{\exp(\eta_1)^2}{S^2}}{S^2} & -\frac{\exp(\eta_1)\exp(\eta_2)}{S^2} & \dots & -\frac{\exp(\eta_1)\exp(\eta_{K-1})}{S^2} \\ -\frac{\exp(\eta_2)\exp(\eta_1)}{S^2} & \frac{\exp(\eta_2)}{S} - \frac{\exp(\eta_2)^2}{S^2} & \dots & -\frac{\exp(\eta_2)\exp(\eta_K)}{S^2} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{\exp(\eta_K)\exp(\eta_1)}{S^2} & -\frac{\exp(\eta_K)\exp(\eta_2)}{S^2} & \dots & \frac{\exp(\eta_K)}{S} - \frac{\exp(\eta_K)^2}{S^2} \end{pmatrix}, \end{aligned}$$

where $S = [\sum_{k=1}^K \exp(\eta_k)]$. We denote $p_i = \exp(\eta_i)/S$, Then we have

$$\begin{aligned} & \mathbb{E}_{(x,y) \in D_t^{val}} [\ell(f(h(x)), y) - \ell(f(h(x)), y)] = (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^\top \frac{a''(\mathbf{c})}{2} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \\ &= \sum_{k=1}^K [p_k - p_k^2] \eta_k^2 - \sum_{k=1}^K \sum_{j=1, j \neq k}^K p_k p_j \eta_k \eta_j \\ &\geq \min_{k,j \in [K], k \neq j} \{p_k - p_k^2, p_k p_j\} \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}\|_2^2 := C \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}\|_2^2. \end{aligned}$$

Proof of Proposition 7

Proof We denote

$$\tilde{\mathbf{A}} = \arg \min_{\mathbf{A}} \mathbb{E}_{(x^{(s)}, y^{(s)}) \sim \mu^s} \ell(\mathbf{A}^\top \hat{h}(x^{(s)}), y^{(s)})$$

and

$$\mathbf{A}^* = \arg \min_{\mathbf{A}} \mathbb{E}_{(x^{(s)}, y^{(s)}) \sim \mu^s} \ell(\mathbf{A}^\top h^*(x^{(s)}), y^{(s)}).$$

Observe that

$$\begin{aligned}
 & R_\eta(\hat{h}) - R_\eta(h^*) \\
 &= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^n} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^m} \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{val}) - \\
 &\quad \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^n} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^m} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h^*), D_t^{val}) \\
 &= \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{(x,y) \sim \mu^q} \left\{ \ell(\tilde{\mathbf{A}}^\top \hat{h}(x), y) - \ell(\mathbf{A}^{*\top} h^*(x), y) \right\} \\
 &\leq \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} \|\tilde{\mathbf{A}}^\top \hat{h}(x) - \mathbf{A}^{*\top} h^*(x)\|_2^2 \\
 &= \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} \left\{ \sum_{k=1}^K \left| (\tilde{\mathbf{A}})_k^\top \hat{h}(x) - (\mathbf{A})_k^{*\top} h^*(x) \right|^2 \right\} \\
 &= \sup_{\mathbf{A}'} \inf_{\mathbf{A}} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} \left\{ \sum_{k=1}^K \left| (\mathbf{A})_k^\top \hat{h}(x) - (\mathbf{A})_k'^\top h^*(x) \right|^2 \right\} \\
 &= \sum_{k=1}^K \left\{ \sup_{(\mathbf{A})_k'} \inf_{(\mathbf{A})_k} [(\mathbf{A})_k^\top - (\mathbf{A})_k'^\top] \Lambda \begin{bmatrix} (\mathbf{A})_k \\ -(\mathbf{A})_k' \end{bmatrix} \right\},
 \end{aligned}$$

where Λ is defined as

$$\Lambda(\hat{h}, h^*) = \begin{bmatrix} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [\hat{h}(x) \hat{h}(x)^\top] & \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [\hat{h}(x) h^*(x)^\top] \\ \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [h^*(x) \hat{h}(x)^\top] & \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{x \sim \mu^q} [h^*(x) h^*(x)^\top] \end{bmatrix} = \begin{bmatrix} \mathbf{G}_{\hat{h}\hat{h}} & \mathbf{G}_{\hat{h}h^*} \\ \mathbf{G}_{h^*\hat{h}} & \mathbf{G}_{h^*h^*} \end{bmatrix}.$$

The first inequality holds since the 1-Lipschitz continuity of the cross-entropy loss. According to the partial minimization of a convex quadratic form (Boyd et al., 2004), we have

$$\inf_{(\mathbf{A})_k} \left\{ [(\mathbf{A})_k^\top - (\mathbf{A})_k'^\top] \Lambda \begin{bmatrix} (\mathbf{A})_k \\ -(\mathbf{A})_k' \end{bmatrix} \right\} = (\mathbf{A})_k'^\top \Lambda_S(\hat{h}, h^*) (\mathbf{A})_k',$$

where $\Lambda_S(\hat{h}, h^*) = \mathbf{G}_{h^*h^*} - \mathbf{G}_{h^*\hat{h}}(\mathbf{G}_{\hat{h}\hat{h}})^\dagger \mathbf{G}_{\hat{h}h^*}$ is the generalized Schur complement of the representation of h^* with respect to \hat{h} . Furthermore, according to the variational characterization of the singular value, we have

$$\sum_{k=1}^K \sup_{(\mathbf{A})_k' : \|(\mathbf{A})_k'\| \leq M_k} (\mathbf{A})_k'^\top \Lambda_S(\hat{h}, h^*) (\mathbf{A})_k' = \sum_{k=1}^K M_k \sigma_1(\Lambda_S(\hat{h}, h^*)) \leq M \sigma_1(\Lambda_S(\hat{h}, h^*)),$$

where σ_1 denotes the maximal singular value.

Moreover, we denote

$$\tilde{\mathbf{A}}_t = \arg \min_{\mathbf{A}_t} \mathbb{E}_{(x_t^{(s)}, y_t^{(s)}) \sim \mu_t^s} \ell(\mathbf{A}_t^\top \hat{h}(x_t^{(s)}), y_t^{(s)}), t \in [T]$$

and

$$\mathbf{A}_t^* = \arg \min_{\mathbf{A}_t} \mathbb{E}_{(x_t^{(s)}, y_t^{(s)}) \sim \mu_t^s} \ell(\mathbf{A}_t^\top h^*(x_t^{(s)}), y_t^{(s)}), t \in [T].$$

Thus we have the following derivation

$$\begin{aligned}
& R_{train}(\hat{h}) - R_{train}(h^*) \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^n} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^m} \mathbf{L}(\mathcal{LM}(D_t^{tr}; \hat{h}), D_t^{val}) - \\
&\quad \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{D_t^{val} \sim (\mu_t^q)^n} \mathbb{E}_{D_t^{tr} \sim (\mu_t^s)^m} \mathbf{L}(\mathcal{LM}(D_t^{tr}; h^*), D_t^{val}) \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{(x,y) \sim \mu_t^q} \left\{ \ell(\tilde{\mathbf{A}}_t^\top \hat{h}(x), y) - \ell(\mathbf{A}_t^{*\top} h^*(x), y) \right\} \\
&\geq \frac{C}{T} \sum_{t=1}^T \mathbb{E}_{\mu_t \sim \eta} \mathbb{E}_{x \sim \mu_t^q} \left\{ \|\tilde{\mathbf{A}}_t^\top \hat{h}(x) - \mathbf{A}_t^{*\top} h^*(x)\|_2^2 \right\} \\
&= \frac{C}{T} \sum_{t=1}^T \inf_{\mathbf{A}_t} \mathbb{E}_{\mu_t \sim \eta} \mathbb{E}_{x \sim \mu_t^q} \left\{ \|\tilde{\mathbf{A}}_t^\top \hat{h}(x) - \mathbf{A}_t^{*\top} h^*(x)\|_2^2 \right\} \\
&= \frac{C}{T} \sum_{t=1}^T \inf_{\mathbf{A}_t} \mathbb{E}_{\mu_t \sim \eta} \mathbb{E}_{x \sim \mu_t^q} \left\{ \sum_{k=1}^K |(\tilde{\mathbf{A}}_t)_k^\top \hat{h}(x) - (\mathbf{A}_t)_k^{*\top} h^*(x)|^2 \right\} \\
&= \frac{C}{T} \sum_{t=1}^T \sum_{k=1}^K (\mathbf{A}_t)_k^{*\top} \Lambda_S(\hat{h}, h^*) (\mathbf{A}_t)_k^* = \sum_{k=1}^K \text{tr}(\Lambda_S(\hat{h}, h^*) (\mathbf{Q})_k^\top (\mathbf{Q})_k / T) = \sum_{k=1}^K \text{tr}(\Lambda_S(\hat{h}, h^*) (\mathbf{K})_k),
\end{aligned}$$

where $(\mathbf{K})_k = (\mathbf{Q})_k^\top (\mathbf{Q})_k / T$. Since $\Lambda_S \succeq 0$, and $(\mathbf{K})_k \succeq 0$, through the Von-Neumann trace inequality, we have that

$$\begin{aligned}
\sum_{k=1}^K \text{tr}(\Lambda_S(\hat{h}, h^*) (\mathbf{K})_k) &\geq \sum_{k=1}^K \sum_{i=1}^{d_L} \sigma_i(\Lambda_S(\hat{h}, h^*)) \sigma_{d_L-i+1}((\mathbf{K})_k) \geq \sum_{k=1}^K \sum_{i=1}^{d_L} \sigma_i(\Lambda_S(\hat{h}, h^*)) \sigma_{d_L}((\mathbf{K})_k) \\
&= \sum_{k=1}^K \text{tr}(\Lambda_S(\hat{h}, h^*)) \sigma_{d_L}((\mathbf{K})_k) \geq \sigma_1(\Lambda_S(\hat{h}, h^*)) \sum_{k=1}^K \sigma_{d_L}((\mathbf{K})_k).
\end{aligned}$$

Now, we can deduce that

$$R_\eta(\hat{h}) - R_\eta(h^*) \leq \frac{M}{\sum_{k=1}^K \sigma_{d_L}((\mathbf{K})_k)} \left\{ R_{train}(\hat{h}) - R_{train}(h^*) \right\}.$$

Based on the above derivation, the conclusion can then obtained. ■

Appendix G. Proof of Online Methodology-Learning algorithm

G.1 Proof of Theorem 9

We firstly present Lemma 20 and then provide the proof of Theorem 9.

Lemma 20 *If we set the step size $\frac{2\tau}{\mu\rho L'(\mathcal{F})} \leq \alpha \leq \min\{\frac{1}{G}, \frac{1}{\beta L(\mathcal{F})}\}$, then $\ell_t(h)$ is convex, where $L'(\mathcal{F})\|h - h'\| \leq \|f_t^{(h)} - f_t^{(h')}\| \leq L(\mathcal{F})\|h - h'\|$. Furthermore, it is also 2τ -smooth and τ -strongly convex.*

Proof Using the chain rule and our definitions, we have

$$\begin{aligned} \nabla \ell_t(h) - \nabla \ell_t(h') &= \frac{\partial f_t^{(h)}}{\partial h} \nabla \mathbf{L}(f_t^{(h)}, D_t^{val}) - \frac{\partial f_t^{(h')}}{\partial h} \nabla \mathbf{L}(f_t^{(h')}, D_t^{val}) \\ &= \left(\frac{\partial f_t^{(h)}}{\partial h} - \frac{\partial f_t^{(h')}}{\partial h} \right) \nabla \mathbf{L}(f_t^{(h)}, D_t^{val}) + \frac{\partial f_t^{(h')}}{\partial h} \left(\nabla \mathbf{L}(f_t^{(h)}, D_t^{val}) - \nabla \mathbf{L}(f_t^{(h')}, D_t^{val}) \right) \end{aligned}$$

Taking the norm on both sides, and noticing that $f_t^{h_t} = f_t - \alpha \frac{\partial \mathcal{L}^{task}(f, h_t, D_t^{tr})}{\partial f}$, for the specified α , we have:

$$\begin{aligned} &\|\nabla \ell_t(h) - \nabla \ell_t(h')\| \\ &\leq \left\| \left(\frac{\partial f_t^{(h)}}{\partial h} - \frac{\partial f_t^{(h')}}{\partial h} \right) \nabla \mathbf{L}(f_t^{(h)}, D_t^{val}) \right\| + \left\| \frac{\partial f_t^{(h')}}{\partial h} \left(\nabla \mathbf{L}(f_t^{(h)}, D_t^{val}) - \nabla \mathbf{L}(f_t^{(h')}, D_t^{val}) \right) \right\| \\ &\leq \alpha \tau G \|h - h'\| + \alpha \tau \beta L(\mathcal{F}) \|h - h'\| \\ &\leq \tau \|h - h'\| + \tau \|h - h'\| = 2\tau \|h - h'\| \end{aligned}$$

Similarly, we obtain the following lower bound

$$\begin{aligned} &\|\nabla \ell_t(h) - \nabla \ell_t(h')\| \\ &\geq \left\| \frac{\partial f_t^{(h')}}{\partial h} \left(\nabla \mathbf{L}(f_t^{(h)}, D_t^{val}) - \nabla \mathbf{L}(f_t^{(h')}, D_t^{val}) \right) \right\| - \left\| \left(\frac{\partial f_t^{(h)}}{\partial h} - \frac{\partial f_t^{(h')}}{\partial h} \right) \nabla \mathbf{L}(f_t^{(h)}, D_t^{val}) \right\| \\ &\geq \alpha \mu L'(\mathcal{F}) \rho \|h - h'\| - \alpha \tau G \|h - h'\| \\ &\geq 2\tau \|h - h'\| - \tau \|h - h'\| = \tau \|h - h'\|, \end{aligned}$$

which completes the proof. ■

Then we present the proof of Theorem 9.

Proof The FTML algorithm is identical to FTL on the sequence of loss functions $a_t, t \in [T]$, which has a $\frac{4A^2 \log(1+T)}{B}$ regret bound, where a_t is A -Lipschitz, and B -strongly convex (see (Cesa-Bianchi and Lugosi, 2006) Theorem 3.1). Using $A = G, B = \tau$ completes the proof. ■

Appendix H. Comparing Our Meta-Regularization with Related Works

H.1 Comparing Tanh with Neural Architecture Search

Recall the bounds that we instantiate our general-purpose bounds in Theorem 3 for few-shot regression in Section 5,

$$\begin{aligned}
& R_{test}(\hat{f}_\mu, \hat{h}) - R_{test}(f_\mu^*, h^*) \\
& \leq \frac{M}{\sigma_{d_L}(\mathbf{K})} \left(768L \log(4nT) L(\mathcal{F}) 2d_L \sqrt{\log(nT)} \frac{R \left(\sqrt{2\log(2)L} + 1 \right) \prod_{i=1}^L B_i}{\sqrt{nT}} \right. \\
& \quad \left. + \frac{6L\|\mathbf{w}\|}{\sqrt{mT}} \sum_{t=1}^T \max_{x_{ti}^{(s)} \in D_t^{val}} \|h(x_{ti}^{(s)})\| + 6B\sqrt{\frac{\log \frac{2}{\delta}}{2nT}} + 6B\sqrt{\frac{\log \frac{2}{\delta}}{m}} + \frac{48L \sup_{h,x} M \|h(x)\|}{n^2 T^2} \right) \\
& \quad + \frac{6L\|\mathbf{w}\|}{\sqrt{m_\mu}} \cdot \max_{x_i^{(s)} \in D_\mu^{val}} \|h(x_i^{(s)})\| + 6B\sqrt{\frac{\log \frac{2}{\delta}}{m_\mu}}.
\end{aligned} \tag{39}$$

It can be seen that the transfer error defined in Eq.(39) contains an important term $\|h(x)\|$, i.e., the output range of the meta-learner $h(x)$ in our few-shot regression setting defined as follows:

$$h(x) = \phi_L(\mathbf{W}_L(\phi_{L-1}(\mathbf{W}_{L-1} \cdots \phi_1(\mathbf{W}_1 x)))). \tag{40}$$

Thus $\|h(x)\|$ can be bounded by:

$$\begin{aligned}
\|h(x)\| &= \|\mathbf{r}_L(x)\| = \|\phi_L(\mathbf{W}_L \mathbf{r}_{L-1}(x))\| \\
&\leq \|\mathbf{W}_L \mathbf{r}_{L-1}(x)\| \leq \|\mathbf{W}_L\| \|\mathbf{r}_{L-1}(x)\| \leq D \prod_{k=1}^L \|\mathbf{W}_k\|,
\end{aligned} \tag{41}$$

where $\mathbf{r}_k(\cdot)$ denotes the vector-valued output of the k -layer for $k \in [L]$. The above bound assumes that the activation function $\phi_k, k \in [L]$ is ReLU. *If we revise the last activation function as Tanh (i.e., $\phi_L = \frac{e^z - e^{-z}}{e^z + e^{-z}}$), then we have*

$$\|h(\mathbf{x})\| = \|\mathbf{r}_L(\mathbf{x})\| = \|\phi_L(\mathbf{W}_L \mathbf{r}_{L-1}(\mathbf{x}))\| \leq \sqrt{d_L}.$$

Generally, $\sqrt{d_L}$ is much smaller than $D \prod_{k=1}^L \|\mathbf{W}_k\|_F$, and the complexity can thus be expected to substantially decrease. This means that the generalization capability of the calculated meta-learner is hopeful to be improved by replacing the last activation function from ReLU with Tanh.

Recall that the structural risk minimization (SRM) principle in conventional machine learning is often used to improve generalization capability of the extracted learner through inducing some rational meta-regularization strategies. Similarly, changing the last activation function from ReLU as conventional to Tanh can be thought as a meta-regularization strategy for ameliorating the generalization capability of the extracted meta-learner. This

meta-regularization strategy is completely inspired from our derived learning theory, and this should be the first time to use such meta-regularization strategy in meta-learning, to the best of our knowledge.

Besides, we want to clarify that our proposed method should be evidently different from the approaches that apply neural architecture search to meta-learning, e.g., (Lian et al., 2020; Wang et al., 2020a). Specifically, T-NAS (Lian et al., 2020) learns a meta-architecture that is capable of adapting to a new task quickly through a few gradient steps, and M-NAS (Wang et al., 2020a) proposes meta neural architecture search to learn a controller for architecture generation and then fast adapt to new tasks. Compared with current meta-learning methods using existing backbone networks, both of them attempt to learn a new architecture using neural architecture search techniques for the specific tasks. Comparatively, our method proposes to revise the activation function inspired from statistical learning theory perspective, rather than treat activation function as hyper-parameters to learn or tune from training data. This means that we explore a new orthogonal research direction to understand and improve meta-learning. Though T-NAS (Lian et al., 2020), M-NAS (Wang et al., 2020a) can obtain nice performance improvement, our method is designed upon more comprehensively solid theoretical guarantee and with more efficient computation implementation. In this sense, we believe that it is still rational to say our meta-learning framework and its statistical learning theory, as well as its deduced meta-regularization strategies, should be meaningful to the field.

H.2 Comparing L^2 -SP with Meta-MinibatchProx and iMAML

Meta-MinibatchProx proposed in (Zhou et al., 2019), as well as iMAML proposed in (Rajeswaran et al., 2019), mainly focuses on improving the computation efficiency of the well-known MAML algorithm. The original MAML seeks to find the task-specific model by fine-tuning meta-parameters θ with stochastic gradient descent, i.e.,

$$\mathcal{Alg}^*(\theta, \mathcal{D}_i^{tr}) = \theta - \alpha \nabla \mathbf{L}(\theta, \mathcal{D}_i^{tr}), \quad (42)$$

where α is the learning rate. To have sufficient learning in the task-specific model learning while also avoiding over-fitting, Meta-MinibatchProx and iMAML consider the following proximal regularization for the task-specific model:

$$\mathcal{Alg}^*(\theta, \mathcal{D}_i^{tr}) = \arg \min_{\phi \in \Phi} \mathbf{L}(\phi, \mathcal{D}_i^{tr}) + \frac{\lambda}{2} \|\phi - \theta\|^2, \quad (43)$$

where θ and ϕ are the meta-parameters and model-parameters, respectively. The overall bi-level meta-learning objective can be written as

$$\theta^* = \arg \min_{\theta \in \Theta} F(\theta), \text{ where } F(\theta) = \frac{1}{T} \sum_{i=1}^T \mathbf{L}_i(\mathcal{Alg}_i^*(\theta), \mathcal{D}_i^{test}), \quad (44)$$

$$\mathcal{Alg}_i^*(\theta) = \arg \min_{\phi \in \Phi} G(\phi, \theta), \text{ where } G(\phi, \theta) = \mathbf{L}_i(\phi, \mathcal{D}_i^{tr}) + \frac{\lambda}{2} \|\phi - \theta\|^2, \quad (45)$$

where $\mathcal{Alg}_i^*(\theta) = \mathcal{Alg}_i^*(\theta, \mathcal{D}_i^{tr})$.

Comparatively, our L^2 -SP meta-regularizer is with evident differences with the above regularizers mainly from the following aspects:

- *Formulation.* The proximal regularization in Eq.(43) encourages the model-parameters to remain close to θ , which regulates the *inner-level model capacity*. While our L^2 -SP meta-regularizer minimizes the distance between the weights from the starting point weights of meta-model, which regulates the *outer-level meta-model (meta-learner) capacity*. For example, if we use $L^2 - SP$, the outer-level objective in Eq.(44,45) can be written as

$$\theta^* = \arg \min_{\theta \in \Theta} \tilde{F}(\theta), \text{ where } \tilde{F}(\theta) = \frac{1}{T} \sum_{i=1}^T \mathbf{L}_i(\text{Alg}_i^*(\theta), \mathcal{D}_i^{\text{test}}) + \lambda \|\theta - \theta_0\|^2, \quad (46)$$

$$\text{Alg}_i^*(\theta) = \arg \min_{\phi \in \Phi} G(\phi, \theta), \text{ where } G(\phi, \theta) = \mathbf{L}_i(\phi, \mathcal{D}_i^{\text{tr}}) \quad (47)$$

where θ_0 is the starting point weights of θ .

- *Theory.* The proximal regularization explores memory and computation efficient optimization steps via the efficient proximal learning from the *optimization theory perspective*, and bring better results on few-shot learning benchmarks. Comparatively, our method induces the L^2 -SP meta-regularization from the *statistical learning theory perspective*. Specifically, the upper bound of Rademacher complexity for meta-model in Theorem 5 of our manuscript is given by

$$\hat{\mathfrak{R}}_N(\mathcal{H}) \leq \frac{2\sqrt{2}d_L R \sum_{j=1}^L \frac{D_j}{2B_j \prod_{i=1}^j \sqrt{c_i}} \prod_{j=1}^L 2B_j \sqrt{c_j}}{\sqrt{N}}, \quad (48)$$

where D_j controls the distance between the weights and the starting point weights of j -th layer, i.e., $\|\mathbf{W}_j - \mathbf{W}_j^0\|_F \leq D_j, j \in [L]$. Therefore, to reduce the meta generalization error of meta-learner, we can minimize $\|\mathbf{W}_j - \mathbf{W}_j^0\|_F$ to reduce $D_j, j \in [L]$, so as to improve the transferable generalization capability of meta-learner.

- *Algorithm.* Compared with the original MAML, Meta-MinibatchProx and iMAML need to re-design inner-level algorithm for new formulation. Comparatively, our method can be easily integrated into existing meta-learning methods, and still use their out-level algorithms to update meta-learner, and do not revise the inner-level algorithm.

References

- Hirofugu Akaike. A new look at the statistical model identification. *IEEE transactions on automatic control*, 19(6):716–723, 1974.
- Ferran Alet, Erica Weng, Tomás Lozano-Pérez, and Leslie Kaelbling. Neural relational inference with fast modular meta-learning. In *NeurIPS*, 2019.
- Ron Amit and Ron Meir. Meta-learning by adjusting priors based on extended pac-bayes theory. In *ICML*, pages 205–214, 2018.
- Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. In *NeurIPS*, 2016.
- Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. In *ICLR*, 2019.
- Sanjeev Arora, Simon Du, Sham Kakade, Yuping Luo, and Nikunj Saunshi. Provable representation learning for imitation learning via bi-level optimization. In *ICML*, pages 367–376, 2020.
- Yu Bai, Minshuo Chen, Pan Zhou, Tuo Zhao, Jason Lee, Sham Kakade, Huan Wang, and Caiming Xiong. How important is the train-validation split in meta-learning? In *ICML*, pages 543–553, 2021.
- Sungyong Baik, Seokil Hong, and Kyoung Mu Lee. Learning to forget for meta-learning. In *CVPR*, 2020.
- Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. Metareg: Towards domain generalization using meta-regularization. In *NeurIPS*, 2018.
- Maria-Florina Balcan, Mikhail Khodak, and Ameet Talwalkar. Provable guarantees for gradient-based meta-learning. In *ICML*, pages 424–433, 2019.
- Andrew R Barron and Thomas M Cover. Minimum complexity density estimation. *IEEE transactions on information theory*, 37(4):1034–1054, 1991.
- Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Peter L Bartlett, Dylan J Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for neural networks. volume 2017, 2017.
- Jonathan Baxter. A model of inductive bias learning. *Journal of artificial intelligence research*, 12:149–198, 2000.
- Atilim Gunes Baydin, Robert Cornish, David Martinez Rubio, Mark Schmidt, and Frank Wood. Online learning rate adaptation with hypergradient descent. In *ICML*, 2019.

- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- Yoshua Bengio, Samy Bengio, and Jocelyn Cloutier. *Learning a synaptic learning rule*. 1990.
- James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of machine learning research*, 13(2), 2012.
- Alberto Bernacchia. Meta-learning with negative learning rates. In *ICLR*, 2021.
- Luca Bertinetto, Joao F Henriques, Philip Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. In *ICLR*, 2019.
- John B Biggs. The role of metalearning in study processes. *British journal of educational psychology*, 55(3):185–212, 1985.
- Hakan Bilen, Basura Fernando, Efstratios Gavves, Andrea Vedaldi, and Stephen Gould. Dynamic image networks for action recognition. In *CVPR*, 2016.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- John Bronskill, Jonathan Gordon, James Requeima, Sebastian Nowozin, and Richard Turner. Tasknorm: Rethinking batch normalization for meta-learning. In *ICML*, pages 1153–1164, 2020.
- Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- Jiaxin Chen, Xiao-Ming Wu, Yanke Li, Qimai Li, Li-Ming Zhan, and Fu-lai Chung. A closer look at the training strategy for modern meta-learning. In *NeurIPS*, volume 33, 2020.
- Lisha Chen, Songtao Lu, and Tianyi Chen. Understanding benign overfitting in gradient-based meta learning. In *NeurIPS*, 2022.
- Qi Chen, Changjian Shui, and Mario Marchand. Generalization bounds for meta-learning: An information-theoretic analysis. In *NeurIPS*, volume 34, pages 25878–25890, 2021.
- Kurtland Chua, Qi Lei, and Jason D Lee. How fine-tuning allows for effective meta-learning. In *NeurIPS*, volume 34, pages 8871–8884, 2021.
- Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *CVPP*, 2014.
- Liam Collins, Aryan Mokhtari, Sewoong Oh, and Sanjay Shakkottai. Maml and anil provably learn representations. In *ICML*, pages 4238–4310, 2022.
- Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. In *CVPR*, 2019.

- Giulia Denevi, Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Learning to learn around a common mean. In *NeurIPS*, volume 31, 2018.
- Giulia Denevi, Carlo Ciliberto, Riccardo Grazi, and Massimiliano Pontil. Learning-to-learn stochastic gradient descent with biased regularization. In *ICML*, pages 1566–1575, 2019a.
- Giulia Denevi, Dimitris Stamos, Carlo Ciliberto, and Massimiliano Pontil. Online-within-online meta-learning. In *NeurIPS*, 2019b.
- Giulia Denevi, Massimiliano Pontil, and Carlo Ciliberto. The advantage of conditional meta-learning for biased regularization and fine tuning. In *NeurIPS*, 2020.
- Giulia Denevi, Carlo Ciliberto, et al. Conditional meta-learning of linear representations. In *NeurIPS*, 2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *ACL*, 2019.
- Kehui Ding, Jun Shu, Deyu Meng, and Zongben Xu. Improve noise tolerance of robust loss via noise-awareness. *arXiv preprint arXiv:2301.07306*, 2023.
- Nan Ding, Xi Chen, Tomer Levinboim, Sebastian Goodman, and Radu Soricut. Bridging the gap between practice and pac-bayes theory in few-shot meta-learning. In *NeurIPS*, volume 34, pages 29506–29516, 2021.
- Simon S Du, Wei Hu, Sham M Kakade, Jason D Lee, and Qi Lei. Few-shot learning via learning the representation, provably. In *ICLR*, 2021.
- Thomas Elsken, Jan Hendrik Metzen, Frank Hutter, et al. Neural architecture search: A survey. *JMLR*, 20(55):1–21, 2019.
- Rasool Fakoor, Pratik Chaudhari, Stefano Soatto, and Alexander J Smola. Meta-q-learning. In *ICLR*, 2020.
- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. On the convergence theory of gradient-based model-agnostic meta-learning algorithms. In *AISTATS*, pages 1082–1092, 2020.
- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Generalization of model-agnostic meta-learning algorithms: Recurring and unseen tasks. In *NeurIPS*, volume 34, pages 5469–5480, 2021.
- Alec Farid and Anirudha Majumdar. Generalization bounds for meta-learning via pac-bayes and uniform stability. In *NeurIPS*, volume 34, pages 2173–2186, 2021.
- Matthias Feurer and Frank Hutter. Hyperparameter optimization. *Automated machine learning: Methods, systems, challenges*, pages 3–33, 2019.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, pages 1126–1135, 2017.

- Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *ICML*, pages 1920–1930, 2019.
- Sebastian Flennerhag, Andrei A Rusu, Razvan Pascanu, Francesco Visin, Hujun Yin, and Raia Hadsell. Meta-learning with warped gradient descent. In *ICLR*, 2020.
- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, pages 1568–1577, 2018.
- Francisco Garcia and Philip S Thomas. A meta-mdp approach to exploration for lifelong reinforcement learning. In *NeurIPS*, volume 32, pages 5691–5700, 2019.
- Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *ICML*, pages 1704–1713. PMLR, 2018.
- Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. In *NeurIPS*, 2018.
- Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, pages 297–299, 2018.
- Santiago Gonzalez and Risto Miikkulainen. Improved training speed, accuracy, and data utilization through loss function optimization. In *2020 IEEE Congress on Evolutionary Computation (CEC)*, pages 1–8, 2020.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- Henry Gouk, Hospedales Subhransu, and Massimiliano Pontil. Distance-based regularisation of deep networks for fine-tuning. In *ICLR*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Fredrik Hellström and Giuseppe Durisi. Evaluated cmi bounds for meta learning: Tightness and expressiveness. In *NeurIPS*, 2022.
- Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*, 2020.
- Rein Houthoofd, Richard Y Chen, Phillip Isola, Bradley C Stadie, Filip Wolski, Jonathan Ho, and Pieter Abbeel. Evolved policy gradients. In *NeurIPS*, 2018.
- Chen Huang, Shuangfei Zhai, Walter Talbott, Miguel Bautista Martin, Shih-Yu Sun, Carlos Guestrin, and Josh Susskind. Addressing the loss-metric mismatch with adaptive loss alignment. In *ICML*, pages 2891–2900, 2019.
- Yu Huang, Yingbin Liang, and Longbo Huang. Provable generalization of overparameterized meta-learning trained with sgd. In *NeurIPS*, 2022.

- Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced design. In *ICML*, pages 4882–4892, 2021.
- Kaiyi Ji, Junjie Yang, and Yingbin Liang. Theoretical convergence of multi-step model-agnostic meta-learning. *Journal of machine learning research*, 2022.
- Michael I Jordan and Tom M Mitchell. Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245):255–260, 2015.
- Sharu Theresa Jose and Osvaldo Simeone. Information-theoretic generalization bounds for meta-learning and applications. *Entropy*, 23(1):126, 2021.
- Sharu Theresa Jose, Osvaldo Simeone, and Giuseppe Durisi. Transfer meta-learning: Information-theoretic bounds and information meta-risk minimization. *IEEE Transactions on Information Theory*, 68(1):474–501, 2021.
- Madan Kanika, Ke Nan Rosemary, Goyal Anirudh, Schölkopf Bernhard, and Bengio Yoshua. Meta attention networks: Meta-learning attention to modulate information between recurrent independent mechanisms. In *ICLR*, 2021.
- Chia-Hsiang Kao, Wei-Chen Chiu, and Pin-Yu Chen. Maml is a noisy contrastive learner in classification. In *ICLR*, 2022.
- Mikhail Khodak, Maria-Florina F Balcan, and Ameet S Talwalkar. Adaptive gradient-based meta-learning methods. In *NeurIPS*, volume 32, pages 5917–5928, 2019.
- Ksenia Konyushkova, Raphael Sznitman, and Pascal Fua. Learning active learning from data. In *NeurIPS*, 2017.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012.
- Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.
- Hae Beom Lee, Taewook Nam, Eunho Yang, and Sung Ju Hwang. Meta dropout: Learning to perturb latent features for generalization. In *ICLR*, 2020.
- Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *CVPR*, 2019.
- Yoonho Lee and Seungjin Choi. Gradient-based meta-learning with learned layerwise metric and subspace. In *ICML*, pages 2927–2936, 2018.

- Christiane Lemke, Marcin Budka, and Bogdan Gabrys. Metalearning: a survey of trends and technologies. *Artificial intelligence review*, 44(1):117–130, 2015.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *ICCV*, pages 5542–5550, 2017a.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy Hospedales. Learning to generalize: Meta-learning for domain generalization. In *AAAI*, 2018a.
- Xingjian Li, Haoyi Xiong, Hanchao Wang, Yuxuan Rao, Liping Liu, and Jun Huan. Delta: Deep learning transfer using feature map with attention for convolutional networks. In *ICLR*, 2019a.
- Xuhong Li, Grandvalet Yves, and Davoine Franck. Explicit inductive bias for transfer learning with convolutional networks. In *ICML*, 2018b.
- Yiying Li, Yongxin Yang, Wei Zhou, and Timothy Hospedales. Feature-critic networks for heterogeneous domain generalization. In *ICML*, 2019b.
- Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few-shot learning. *arXiv preprint arXiv:1707.09835*, 2017b.
- Dongze Lian, Yin Zheng, Yintao Xu, Yanxiong Lu, Leyu Lin, Peilin Zhao, Junzhou Huang, and Shenghua Gao. Towards fast adaptation of neural architectures with meta learning. In *ICLR*, 2020.
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *ICLR*, 2019.
- Risheng Liu, Xuan Liu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A value-function-based interior-point method for non-convex bi-level optimization. In *ICML*, pages 6882–6892, 2021a.
- Risheng Liu, Yaohua Liu, Shangzhi Zeng, and Jin Zhang. Towards gradient-based bilevel optimization with non-convex followers and beyond. In *NeurIPS*, volume 34, pages 8662–8675, 2021b.
- Divyam Madaan, Jinwoo Shin, and Sung Ju Hwang. Learning to generate noise for robustness against multiple perturbations. *arXiv preprint arXiv:2006.12135*, 2020.
- Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- Donald B Maudsley. A theory of meta-learning and principles of facilitation: An organismic perspective. 1980.
- Andreas Maurer. A chain rule for the expected suprema of gaussian processes. *Theoretical Computer Science*, 650:109–122, 2016.
- Andreas Maurer and Tommi Jaakkola. Algorithmic stability and meta-learning. *Journal of Machine Learning Research*, 6(6), 2005.

- Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81):1–32, 2016.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- Stefan Munder and Dariu M Gavrilă. An experimental study on pedestrian classification. *IEEE transactions on pattern analysis and machine intelligence*, 28(11):1863–1868, 2006.
- Vaishnavh Nagarajan and J Zico Kolter. Generalization in deep networks: The role of distance from initialization. *arXiv preprint arXiv:1901.01672*, 2019.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of over-parametrization in generalization of neural networks. In *ICLR*, 2019.
- Cuong Nguyen, Thanh-Toan Do, and Gustavo Carneiro. Pac-bayes meta-learning with implicit task-specific posteriors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):841–851, 2022.
- Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729, 2008.
- Boris N Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. In *NeurIPS*, 2018.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.
- Eunbyung Park and Junier B Oliva. Meta-curvature. In *NeurIPS*, 2019.
- Anastasia Pentina and Christoph Lampert. A pac-bayesian bound for lifelong learning. In *ICML*, pages 991–999, 2014.
- Anastasia Pentina and Christoph H Lampert. Lifelong learning with non-iid tasks. In *NeurIPS*, volume 28, pages 1540–1548, 2015.
- Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. In *NeurIPS*, volume 32, 2019.
- Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017.

- Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017.
- Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with residual adapters. In *NeurIPS*, pages 506–516, 2017.
- Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classification. In *ICLR*, 2019.
- James Requeima, Jonathan Gordon, John Bronskill, Sebastian Nowozin, and Richard E Turner. Fast and flexible multi-task classification using conditional neural adaptive processes. In *NeurIPS*, 2019.
- Arezou Rezazadeh. A unified view on pac-bayes bounds for meta-learning. In *ICML*, pages 18576–18595, 2022.
- Jonas Rothfuss, Vincent Fortuin, Martin Josifoski, and Andreas Krause. Pacoh: Bayes-optimal meta-learning with pac-guarantees. In *ICML*, pages 9116–9126, 2021a.
- Jonas Rothfuss, Dominique Heyn, Andreas Krause, et al. Meta-learning reliable priors in the function space. In *NeurIPS*, volume 34, pages 280–293, 2021b.
- Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.
- Xiangyu Rui, Xiangyong Cao, Jun Shu, Qian Zhao, and Deyu Meng. A hyper-weight network for hyperspectral image denoising. *arXiv preprint arXiv:2301.06081*, 2022.
- Nataniel Ruiz, Samuel Schuler, and Manmohan Chandraker. Learning to simulate. In *ICLR*, 2018.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *ICLR*, 2019.
- Jeong Un Ryu, JaeWoong Shin, Hae Beom Lee, and Sung Ju Hwang. Metaperturb: Transferable regularizer for heterogeneous tasks and architectures. In *NeurIPS*, volume 33, 2020.
- Jürgen Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- Allan M Schrier. Learning how to learn: The significance and current status of learning set formation. *Primates*, 25:95–102, 1984.

- Gideon Schwarz et al. Estimating the dimension of a model. *Annals of statistics*, 6(2): 461–464, 1978.
- Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- Shiv Shankar, Vihari Piratla, Soumen Chakrabarti, Siddhartha Chaudhuri, Preethi Jyothi, and Sunita Sarawagi. Generalizing across domains via cross-gradient training. In *ICLR*, 2018.
- John Shawe-Taylor, Peter L Bartlett, Robert C Williamson, and Martin Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE transactions on Information Theory*, 44(5):1926–1940, 1998.
- Jiayi Shen, Xiaohan Chen, Howard Heaton, Tianlong Chen, Jialin Liu, Wotao Yin, and Zhangyang Wang. Learning a minimax optimizer: A pilot study. In *ICLR*, 2021.
- J Shu, D Meng, and Z Xu. Meta self-paced learning. *Scientia Sinica Informationis*, 50(6): 781–793, 2020a.
- Jun Shu, Zongben Xu, and Deyu Meng. Small sample learning in big data era. *arXiv preprint arXiv:1808.04572*, 2018.
- Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng. Meta-weight-net: Learning an explicit mapping for sample weighting. In *NeurIPS*, 2019.
- Jun Shu, Qian Zhao, Keyu Chen, Zongben Xu, and Deyu Meng. Learning adaptive loss for robust learning with noisy labels. *arXiv preprint arXiv:2002.06482*, 2020b.
- Jun Shu, Yanwen Zhu, Qian Zhao, Deyu Meng, and Zongben Xu. Meta-lr-schedule-net: Learned lr schedules that scale and generalize. *TPAMI*, 2022.
- Jun Shu, Xiang Yuan, and Deyu Meng. Cmw-net: an adaptive robust algorithm for sample selection and label correction. *National Science Review*, 10(6):nwad084, 2023a.
- Jun Shu, Xiang Yuan, Deyu Meng, and Zongben Xu. Cmw-net: Learning a class-aware sample weighting mapping for robust deep learning. *TPAMI*, 2023b.
- Jun Shu, Xiang Yuan, Deyu Meng, and Zongben Xu. Dac-mr: Data augmentation consistency based meta-regularization for meta-learning. *arXiv preprint arXiv:2305.07892*, 2023c.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, pages 4080–4090, 2017.
- Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. In *NeurIPS*, volume 25, 2012.

- Xingyou Song, Wenbo Gao, Yuxiang Yang, Krzysztof Choromanski, Aldo Pacchiano, and Yunhao Tang. Es-maml: Simple hessian-free meta learning. In *ICLR*, 2020.
- Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- Johannes Stalldkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks*, 32:323–332, 2012.
- Mervyn Stone. Cross-validatory choice and assessment of statistical predictions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 36(2):111–133, 1974.
- Yue Sun, Adhyayan Narang, Ibrahim Gulluk, Samet Oymak, and Maryam Fazel. Towards sample-efficient overparameterized meta-learning. In *NeurIPS*, volume 34, pages 28156–28168, 2021.
- Flood Sung, Li Zhang, Tao Xiang, Timothy Hospedales, and Yongxin Yang. Learning to learn: Meta-critic networks for sample efficient learning. *arXiv preprint arXiv:1706.09529*, 2017.
- Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *CVPR*, 2018.
- Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 2012.
- Nilesh Tripuraneni, Michael Jordan, and Chi Jin. On the theory of transfer learning: The importance of task diversity. In *NeurIPS*, volume 33, 2020.
- Nilesh Tripuraneni, Chi Jin, and Michael Jordan. Provable meta-learning of linear representations. In *ICML*, pages 10434–10443, 2021.
- Joaquin Vanschoren. Meta-learning: A survey. *arXiv preprint arXiv:1810.03548*, 2018.
- Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 2013.
- Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5):988–999, 1999.
- Vivek Veeriah, Matteo Hessel, Zhongwen Xu, Janarthanan Rajendran, Richard L Lewis, Junhyuk Oh, Hado van Hasselt, David Silver, and Satinder Singh. Discovery of useful questions as auxiliary tasks. In *NeurIPS*, 2019.
- Ricardo Vilalta and Youssef Drissi. A perspective view and survey of meta-learning. *Artificial intelligence review*, 18(2):77–95, 2002.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *NeurIPS*, pages 3637–3645, 2016.

- Risto Vuorio, Shao-Hua Sun, Hexiang Hu, and Joseph J Lim. Multimodal model-agnostic meta-learning via task-aware modulation. In *NeurIPS*, 2019.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. *Technical report*, 2011.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Haoxiang Wang, Han Zhao, and Bo Li. Bridging multi-task learning and meta-learning: Towards efficient training and effective adaptation. In *ICML*, pages 10991–11002, 2021.
- Jiaxing Wang, Jiaxiang Wu, Haoli Bai, and Jian Cheng. M-nas: Meta neural architecture search. In *AAAI*, volume 34, pages 6186–6193, 2020a.
- Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. On the global optimality of model-agnostic meta-learning. In *ICML*, pages 9837–9846, 2020b.
- Renzhen Wang, Kaiqin Hu, Yanwen Zhu, Jun Shu, Qian Zhao, and Deyu Meng. Meta feature modulator for long-tailed recognition. *arXiv preprint arXiv:2008.03428*, 2020c.
- Ruohan Wang, Yiannis Demiris, and Carlo Ciliberto. Structured prediction for conditional meta-learning. In *NeurIPS*, 2020d.
- Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018.
- Xin Wang, Fisher Yu, Ruth Wang, Trevor Darrell, and Joseph E Gonzalez. Tafe-net: Task-aware feature embeddings for low shot learning. In *CVPR*, 2019.
- Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM Computing Surveys (CSUR)*, 53(3):1–34, 2020e.
- Lewis B Ward. Reminiscence and rote learning. *Psychological Monographs*, 49(4):i, 1937.
- Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. A survey of transfer learning. *Journal of Big data*, 3(1):1–40, 2016.
- Olga Wichrowska, Niru Maheswaranathan, Matthew W Hoffman, Sergio Gomez Colmenarejo, Misha Denil, Nando Freitas, and Jascha Sohl-Dickstein. Learned optimizers that scale and generalize. In *ICML*, 2017.
- Yichen Wu, Jun Shu, Qi Xie, Qian Zhao, and Deyu Meng. Learning to purify noisy labels via meta soft label corrector. In *AAAI*, 2021.
- Tianbing Xu, Qiang Liu, Liang Zhao, and Jian Peng. Learning to explore with meta-policy gradient. In *ICML*, 2018.
- Ziping Xu and Ambuj Tewari. Representation learning beyond linear prediction functions. In *NeurIPS*, volume 34, pages 4792–4804, 2021.

- Huaxiu Yao, Long-Kai Huang, Linjun Zhang, Ying Wei, Li Tian, James Zou, Junzhou Huang, et al. Improving generalization in meta-learning via task augmentation. In *ICML*, pages 11887–11897, 2021.
- Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Yu-Feng Li, Wei-Wei Tu, Qiang Yang, and Yang Yu. Taking human out of learning applications: A survey on automated machine learning. *arXiv preprint arXiv:1810.13306*, 2018.
- Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, and Chelsea Finn. Meta-learning without memorization. In *ICLR*, 2020.
- Du Yingjun, Zhen Xiantong, Shao Ling, and Snoek Cees G. M. Metanorm: Learning to normalize few-shot batches across domains. In *ICLR*, 2021.
- Yu Zhang and Qiang Yang. An overview of multi-task learning. *National Science Review*, 5(1):30–43, 2018.
- Guoqing Zheng, Ahmed Hassan Awadallah, and Susan Dumais. Meta label correction for noisy label learning. In *AAAI*, 2021.
- Pan Zhou, Xiaotong Yuan, Huan Xu, Shuicheng Yan, and Jiashi Feng. Efficient meta learning via minibatch proximal update. In *NeurIPS*, volume 32, 2019.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020.
- Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. In *ICLR*, 2017.