
Fine-Grained Analysis of Stability and Generalization for Modern Meta Learning Algorithms

Jiechao Guan^{1,3}

Yong Liu^{2,3}

Zhiwu Lu^{2,3,*}

¹ School of Information, Renmin University of China, Beijing, China

² Gaoling School of Artificial Intelligence, Renmin University of China, Beijing, China

³ Beijing Key Laboratory of Big Data Management and Analysis Methods

{2014200990, liuyonggsai, luzhiwu}@ruc.edu.cn

Abstract

The support/query episodic training strategy has been widely applied in modern meta learning algorithms. Supposing the n training episodes and the test episodes are sampled independently from the same environment, previous work has derived a generalization bound of $O(1/\sqrt{n})$ for smooth non-convex functions via algorithmic stability analysis. In this paper, we provide fine-grained analysis of stability and generalization for modern meta learning algorithms by considering more general situations. Firstly, we develop matching lower and upper stability bounds for meta learning algorithms with two types of loss functions: (1) nonsmooth convex functions with α -Hölder continuous subgradients ($\alpha \in [0, 1)$); (2) smooth (including convex and non-convex) functions. Our tight stability bounds show that, in the nonsmooth convex case, meta learning algorithms can be inherently less stable than in the smooth convex case. For the smooth non-convex functions, our stability bound is sharper than the existing one, especially in the setting where the number of iterations is larger than the number n of training episodes. Secondly, we derive improved generalization bounds for meta learning algorithms that hold with high probability. Specifically, we first demonstrate that, under the independent episode environment assumption, the generalization bound of $O(1/\sqrt{n})$ via algorithmic stability analysis is near optimal. To attain faster convergence rate, we show how to yield a deformed generalization bound of $O(\ln n/n)$ with the curvature condition of loss functions. Finally, we obtain a generalization bound for meta learning with dependent episodes whose dependency relation is characterized by a graph. Experiments on regression problems are conducted to verify our theoretical results.

1 Introduction

The last decade has witnessed the success of deep learning techniques in machine learning community [30, 26, 12]. However, the need of large amount of annotated data hinders their application in real-life scenarios. To alleviate this issue, meta learning [4], which employs knowledge from past tasks to facilitate adaptation to the new task, has emerged as a promising direction to reduce annotation cost.

Traditional meta learning algorithms directly minimize the empirical error over all samples in the training tasks [3, 36, 37, 41, 38]. To improve the generalization ability of meta learning algorithms, recent works propose the support/query (S/Q) episodic training strategy [48, 21, 44]. Specifically, in modern meta learning algorithms, each episode/task is split into two non-overlapped parts: support set and query set. The support set is used to learn a hypothesis, and the query set is used to measure the performance of the learned hypothesis on that episode. Therefore, the S/Q episodic strategy regards

*Corresponding Author

each task as a training instance and updates the meta learning model by implementing episode-level stochastic gradient descent (SGD). Supposing the n training episodes and the test episodes are sampled independently from the same environment, previous work [9] has derived a high-probability generalization bound of $O(1/\sqrt{n})$ for modern meta learning. Such bound is obtained via algorithmic stability analysis [7] for smooth non-convex loss functions. However, it is still unknown whether such generalization bound of $O(1/\sqrt{n})$ is optimal, and whether we can obtain sharper bounds for modern meta learning. Further, there is still lack of comprehensive comparisons between the bounds obtained via S/Q episodic training and the bounds obtained via traditional empirical risk minimizing (ERM).

In this work, we will address the above problems via algorithmic stability analysis. Algorithmic stability, roughly speaking, bounds the change in the model output by the algorithm when a single data in the dataset is replaced. Our goal is to provide fine-grained analysis of stability and generalization for modern meta learning algorithms by considering more general situations. Firstly, we develop matching lower and upper stability bounds for meta learning algorithms with two types of loss functions: (1) nonsmooth convex functions with α -Hölder continuous subgradients where $0 \leq \alpha < 1$; (2) smooth (including convex and non-convex) functions. Our tight stability bounds demonstrate that, in the nonsmooth convex case, modern meta learning algorithms can be less stable than in the smooth convex case. In particular, the lower stability bound for nonsmooth convex functions is vacuous even if we train modern meta learning algorithms with a relatively small constant step size in SGD. In the smooth non-convex case, our derived bound is sharper than the existing one [9], especially in the setting where the number of SGD iterations is larger than the number n of training episodes. Secondly, we provide high-probability generalization bounds for modern meta learning algorithms with the aforementioned two types of loss functions. Specifically, we first demonstrate that, under the independent episode environment assumption, the bound of $O(1/\sqrt{n})$ is near optimal and is independent of the sample size m per episode. We thus show that, in terms of the sharpness of the generalization bounds, the S/Q episodic training strategy is superior to the traditional ERM strategy for meta learning (see Remark 6). To obtain faster convergence rate, we next show how to yield a deformed generalization bound of $O(\ln n/n)$ with additional curvature assumption (i.e., Polyak-Łojasiewicz condition [49]) of the loss function. Finally, we use the graph approximation technique [50] to obtain a bound for meta learning with dependent episodes whose dependency relation can be characterized by a graph. To the best of our knowledge, this is the first bound that captures how the dependency between episodes can affect the generalization behavior of meta learning algorithms.

Overall, our contributions are four-fold: (1) We provide matching lower and upper stability bounds for modern meta learning algorithms with general loss functions. The stability bound for nonsmooth convex functions implies that modern meta learning algorithms are not stable enough; and the stability bound in the smooth non-convex case is sharper than the existing one. (2) We develop a near-optimal high-probability bound of $O(1/\sqrt{n})$ on the transfer error in meta learning. Such bound is also used to reveal the advantage of the S/Q episodic strategy for meta learning over the traditional ERM strategy. (3) We derive a deformed generalization bound of $O(\ln n/n)$ with additional curvature condition of loss functions. (4) We obtain the first bound for meta learning with dependent episodes. Experiments on regression problems are conducted to validate the convergence of our generalization bounds.

2 Related Work

Algorithmic Stability Theory. Algorithmic stability analysis is an important tool to provide theoretical guarantee for the learnability of machine learning models. [43] has shown that there are non-trivial problems where traditional uniform convergence analysis (i.e., empirical process theory [47]) fails to hold, but stability can be identified as the sufficient and necessary condition for learnability. There are two main groups in this direction: (1) The first group develop different notations of stability and connect their relation to the generalization of specific algorithms. Among them, uniform stability is the most widely used notation and has been utilized to analyze the stability and generalization of regularized ERM algorithms [7]. Hypothesis stability is a weaker notation and has been used to show the stability of k -Nearest Neighbor model [13]. Both of the above algorithmic stability notions have been extended to the randomized setting to demonstrate the stability of Bagging algorithm [16]. In recent years, different notations have been employed to analyze the stability and in-expectation generalization bounds of stochastic gradient descent method, which include uniform stability [25], on-average stability [31], uniform argument stability [35, 2], on-average model stability [33] and locally elastic stability [11]. (2) The second group aims to derive tight high-probability generalization bounds for uniformly stable algorithm in single-task learning. The first high-probability bound has been derived by [7], and has been improved in [19]. Recently, nearly

optimal generalization bounds of $O(1/\sqrt{n})$ have been established in [20, 8], where n is the size of training dataset. Further, with additional Bernstein condition, [29] derives a generalization bound of $O(1/n)$. In this work, we aim to provide tight stability bounds and improved high-probability generalization bounds for episodic meta learning algorithms. The key step to achieve our goal is to reveal the equivalence of notations between single-task learning and episodic meta learning, hence we can extend the demonstration techniques from [8, 29, 50] to the episodic meta learning setting.

Generalization Bounds for Meta Learning. Supposing the n training tasks and the novel tasks are sampled independently from the same environment, [4] derives the first generalization bound on the *transfer error* over the novel task for meta learning. Under the independent task environment assumption, we can categorize existing transfer error bounds into three main groups: **(1)** transfer error bounds of hypothesis space. Such bounds are always achieved via covering number analysis [4] or VC theory [6], and hence are always dimension-dependent. The latest upper bound in this group is of $O(1/\sqrt{nm} + 1/\sqrt{m})$ in [23, Theorem 5], where m is the sample size per task. **(2)** transfer error bounds of the hyper-distribution of prior. Such bounds are obtained via PAC-Bayes analysis [41, 42, 14]. The tightest bound in this group is of order $O(1/\sqrt{n} + 1/m)$ in [18, Theorem 3]. **(3)** transfer error bounds of the algorithm. Such bounds are obtained via algorithmic stability analysis [36, 1]. The tightest bound in this group is of $O(1/\sqrt{n})$ in [9, Theorem 4] for episodic meta learning algorithms. Detailed comparisons between different transfer error bounds can be found in Table A.2 of Appendix A. There also exist other works without the task environment assumption. Instead, they choose to bound the excess risk on the novel task by proposing task-similarity measurement [15, 46], or using the total variation distance as the diversity measurement between novel task and training tasks [17]. In this work, we take the task environment assumption and follow the work of [9]. Our first improvement is to demonstrate that the bound of $O(1/\sqrt{n})$ is near optimal. Besides, we show how to obtain a deformed generalization bound of $O(\ln n/n)$ with additional curvature assumption of the loss function. Further, we derive a bound with dependent training episodes, revealing how dependency relation between episodes can affect the generalization of meta learning algorithms.

3 Problem Formulation

In supervised learning, a sample space $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ is a product space of an input space \mathcal{X} and an output space \mathcal{Y} . $\mathcal{H} = \{h_w : w \in \mathcal{W}\}$ is the hypothesis space where the hypothesis $h_w \in \mathcal{H}$ is parameterized by parameter w in the parameter space \mathcal{W} . A measurable function $f : \mathcal{H} \times \mathcal{Z} \rightarrow [0, M]$ ($M > 0$) is defined as a nonnegative and bounded loss function, whose loss of a hypothesis h_w over a sample z is denoted by $f(h_w, z)$ or $f(w, z)$. Let $\mathcal{M}_1(A)$ denote the set of probability measures over the set A .

Loss Functions. Throughout we assume that the parameter space $\mathcal{W} \subset \mathbb{R}^d$. Thus, we use unambiguously $\|\cdot\| = \|\cdot\|_2$ as the Euclidean norm. Let $\text{Proj}_{\mathcal{W}}$ be the Euclidean projection onto \mathcal{W} , which is nonexpansive $\|\text{Proj}_{\mathcal{W}}(u) - \text{Proj}_{\mathcal{W}}(v)\| \leq \|u - v\|$. For any fixed $z \in \mathcal{Z}$, a function $f(\cdot, z) : \mathcal{W} \rightarrow \mathbb{R}$ is convex if for all $u, v \in \mathcal{W}$, $f(u, z) \geq f(v, z) + \langle g, u - v \rangle$, where $g \in \partial f(v, z)$, and $\partial f(v, z)$ denotes the set of subgradients of $f(\cdot, z)$ at v . Let $\partial^0 f(v, z)$ denote the subgradient with the least norm. If $f(\cdot, z)$ is differentiable, $\partial f(\cdot, z)$ denotes the gradient of $f(\cdot, z)$, i.e., $\partial f(\cdot, z) = \{\nabla f(\cdot, z)\}$. For any $z \in \mathcal{Z}$, $f(\cdot, z)$ is σ -Lipschitz if $\forall u, v \in \mathcal{W}$, $|f(u, z) - f(v, z)| \leq \sigma \|u - v\|$. For any $z \in \mathcal{Z}$, $f(\cdot, z)$ is G -smooth if $\forall u, v \in \mathcal{W}$, $\|\partial f(u, z) - \partial f(v, z)\| \leq G \|u - v\|$. We also give the definition of function with (α, G) -Hölder continuous subgradient as follows. We may refer to such functions as (α, G) -Hölder smooth or α -Hölder smooth function for simplicity when the context is clear.

Definition 1 Let $G > 0$, $\alpha \in [0, 1]$. For any $z \in \mathcal{Z}$, a function $f(\cdot, z)$ is called (α, G) -Hölder smooth if its subgradient is (α, G) -Hölder continuous, i.e., $\partial f(\cdot, z)$ satisfies the following conditions:

$$\forall u, v \in \mathcal{W}, \quad \|\partial f(u, z) - \partial f(v, z)\| \leq G \|u - v\|^\alpha. \quad (1)$$

If (1) holds with $\alpha = 1$, then $f(\cdot, z)$ is a G -smooth function; if (1) holds with $\alpha = 0$, this implies the subgradient boundedness of $f(\cdot, z)$. The examples of loss functions in machine learning satisfying (1) include the q -norm hinge-loss $f(w, z) = (\max(0, 1 - y\langle w, x \rangle))^q$ for classification and the q -th power absolute distance loss $f(w, z) = |y - \langle w, x \rangle|^q$ for regression, whose subgradients are both $(q - 1, C)$ -Hölder continuous for some $C > 0$ if $q \in [1, 2]$ (see [10]). For (α, G) -Hölder smooth function, define $c_\alpha = (1 + 1/\alpha)^{\frac{\alpha}{1+\alpha}} G^{\frac{1}{1+\alpha}}$ if $\alpha \in (0, 1]$; and $c_\alpha = \sup_z \|\partial f(0, z)\| + G$, if $\alpha = 0$.

Single-Task Learning. The training dataset $S = \{z_j = (x_j, y_j)\}_{j=1}^m$ is given by m independent draws from an unknown distribution D on \mathcal{Z} (i.e., $D \in \mathcal{M}_1(\mathcal{Z})$). An algorithm A takes S as input

and outputs a hypothesis $A(S)$ in \mathcal{H} . The set of such algorithms depends only on \mathcal{H} and \mathcal{Z} and will be denoted by $\mathcal{A}(\mathcal{H}, \mathcal{Z})$. In single-task learning, a hypothesis $A(S)$ is always obtained by minimizing the empirical error on S : $\hat{L}(A(S), S) \triangleq \frac{1}{m} \sum_{j=1}^m f(A(S), z_j)$. The performance of the returned hypothesis $A(S)$ is measured by the expected/generalization error with respect to (w.r.t.) the data distribution D : $L(A(S), D) \triangleq \mathbb{E}_{z \sim D} f(A(S), z)$. The goal of learning theory is thus to give a (lower or upper) bound on the expected error based on the empirical error on the training dataset S .

Meta Learning. Following existing theoretical works for meta learning [4, 36, 41, 9], we assume that the distributions $\{D_i\}_{i=1}^n$ associated with different training tasks are drawn from the same task environment τ , which is a probability distribution over the set of all data distributions on \mathcal{Z} (i.e., $\tau \in \mathcal{M}_1(\mathcal{M}_1(\mathcal{Z}))$). During the training process, a meta-sample $\mathbf{S} = \{S_i = S_i^{tr} \cup S_i^{ts}\}_{i=1}^n$ is available, where $S_i^{tr} \stackrel{\text{i.i.d.}}{\sim} D_i^K$ of size K is the training set, and $S_i^{ts} \stackrel{\text{i.i.d.}}{\sim} D_i^q$ of size q is the test set of the i -th training task. In this work, we assume that $K + q = m$ for notation convenience. The training set and the test set are also called support set and query set [9], respectively. Let $\mathcal{A}(\mathcal{A}(\mathcal{H}, \mathcal{Z}), \mathcal{Z}^m)$ be the set of meta learning algorithms. For any $\mathbf{A} \in \mathcal{A}(\mathcal{A}(\mathcal{H}, \mathcal{Z}), \mathcal{Z}^m)$, it takes the meta-sample $\mathbf{S} = \{S_i\}_{i=1}^n$ as input and outputs an algorithm (inner-task algorithm) $\mathbf{A}(\mathbf{S}) : \cup_{m=1}^{\infty} \mathcal{Z}^m \rightarrow \mathcal{H}$. The performance of the learned inner-task algorithm is measured by the expectation of the generalization error w.r.t. the task environment τ , which is defined as the *transfer error* by [36, 9] as follows:

$$er(\mathbf{A}(\mathbf{S}), \tau) \triangleq \mathbb{E}_{D \sim \tau} \mathbb{E}_{S^{tr} \sim D^K} \mathbb{E}_{z \sim D} f(\mathbf{A}(\mathbf{S})(S^{tr}), z). \quad (2)$$

Actually, the environment τ can define an induced distribution $\mathbf{D}_\tau \in \mathcal{M}_1(\mathcal{Z}^m)$, by setting $\mathbf{D}_\tau(F) = \mathbb{E}_{D \sim \tau} D^m(F)$ for any measurable set $F \subseteq \mathcal{Z}^m$. Define the estimator $\mathbf{l}(\mathbf{A}(\mathbf{S}), S) \triangleq \hat{L}(\mathbf{A}(\mathbf{S})(S^{tr}), S^{ts})$, where $S = S^{tr} \cup S^{ts}$, $S \stackrel{\text{i.i.d.}}{\sim} D^m$. Then we can rewrite the transfer error as a simple form: $er(\mathbf{A}(\mathbf{S}), \tau) = \mathbb{E}_{S \sim \mathbf{D}_\tau} \mathbf{l}(\mathbf{A}(\mathbf{S}), S)$. This means that, the training error $\mathbf{l}(\mathbf{A}(\mathbf{S}), S)$ is the unbiased version of the transfer error $er(\mathbf{A}(\mathbf{S}), \tau) = \mathbb{E}_{S \sim \mathbf{D}_\tau} \mathbf{l}(\mathbf{A}(\mathbf{S}), S)$. This is similar to the fact that, in single-task learning, the empirical error $f(A(S), z)$ is the unbiased version of the generalization error $L(A(S), D) = \mathbb{E}_{z \sim D} f(A(S), z)$. Therefore, a transfer error bound is formally equivalent to a single-task generalization error bound under the identifications $\mathcal{Z} \leftrightarrow \mathcal{Z}^m$, $f \leftrightarrow \mathbf{l}$, $A \leftrightarrow \mathbf{A}$. The relation of the notations between single-task learning and meta learning is listed in Table B.1 in Appendix B. In practice, it is difficult to minimize $er(\mathbf{A}(\mathbf{S}), \tau)$ directly as we have no information of the environment distribution τ . Instead, we choose to minimize the following empirical risk based on the S/Q episodic training strategy. The goal of meta learning theory is thus to give a bound on the transfer error, based on the *empirical multi-task error* on the meta-sample \mathbf{S} :

$$\hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) \triangleq \frac{1}{n} \sum_{i=1}^n \hat{L}(\mathbf{A}(\mathbf{S})(S_i^{tr}), S_i^{ts}) = \frac{1}{n} \sum_{i=1}^n \mathbf{l}(\mathbf{A}(\mathbf{S}), S_i). \quad (3)$$

Uniform Stability of Meta Learning Algorithms. We say two meta-samples $\mathbf{S} = \{S_i\}_{i=1}^n$ and $\mathbf{S}' = \{S'_i\}_{i=1}^n$ are neighboring, denoted by $\mathbf{S} \simeq \mathbf{S}'$, if they only differ on a single entry, i.e., there exists $i \in [n]$ s.t. $\forall j \neq i, S_j = S'_j$; and $S_i \neq S'_i$. We also define $\mathbf{S}^i = \{S_1, \dots, S'_i, \dots, S_n\}$ as the neighboring meta sample of \mathbf{S} that differs only on the i -th entry. We next define the uniform stability of meta algorithms with episodic training strategy, which is formulated explicitly in [9, Definition 3].

Definition 2 (*Uniform stability of modern meta learning algorithms*) A meta algorithm \mathbf{A} has uniform stability w.r.t. the loss function \hat{L} if the following holds for any meta-sample \mathbf{S} and for any $i \in [n]$, any $D \sim \tau$, $S^{tr} \sim D^K$, $S^{ts} \sim D^q$: $|\hat{L}(\mathbf{A}(\mathbf{S})(S^{tr}), S^{ts}) - \hat{L}(\mathbf{A}(\mathbf{S}^i)(S^{tr}), S^{ts})| \leq \gamma$.

Since $\mathbf{l}(\mathbf{A}(\mathbf{S}), S) = \hat{L}(\mathbf{A}(\mathbf{S})(S^{tr}), S^{ts})$, we can also define the uniform stability of \mathbf{A} as: $\forall S \sim \mathbf{D}_\tau, \forall i \in [n], |\mathbf{l}(\mathbf{A}(\mathbf{S}), S) - \mathbf{l}(\mathbf{A}(\mathbf{S}^i), S)| \leq \gamma$. Such definition is analogous to the uniform stability of an inner-task algorithm A in single-task learning (see Definition D.1 in Appendix D) under the identifications: $\mathbf{l} \leftrightarrow f$, $\mathbf{A} \leftrightarrow A$, $\mathbf{S} \leftrightarrow S$. Thus, we can directly apply the existing uniform stability based generalization bound from single-task learning [7, Theorem 12] to obtain the uniform stability based transfer bound for episodic meta learning [9, Theorem 2], without lengthy and somewhat duplicate proof in [9]. We list such fundamental uniform stability based transfer error bound in Theorem 1 for later comparison. To derive sharper transfer error bounds, our **key step** is to utilize the equivalent relation between the notations of single-task learning and episodic meta learning, thus extending fast-rate bounds in single-task learning [8, 29, 50] to the episodic meta learning setting.

Theorem 1 Suppose the S/Q episodic meta learning algorithm \mathbf{A} has uniform stability γ w.r.t. the estimator $\mathbf{l}(\cdot, S)$ bounded by M . Then, for any task distribution $\tau \in \mathcal{M}_1(\mathcal{M}_1(\mathcal{Z}))$, any $\delta \in (0, 1)$, the following inequality holds with probability at least $1 - \delta$ over the draw of meta sample \mathbf{S} :

$$er(\mathbf{A}(\mathbf{S}), \tau) \leq \hat{er}(\mathbf{A}(\mathbf{S}), S) + \gamma + (2n\gamma + M)\sqrt{\frac{\ln(1/\delta)}{2n}}.$$

4 Uniform Argument Stability Bounds of Meta Learning Algorithms

For a modern meta learning algorithm with deep neural networks [21, 44], we always employ stochastic gradient descent (SGD) method to minimize the empirical error $\frac{1}{n} \sum_{i=1}^n \hat{L}(\mathbf{A}(\mathbf{S})(S_i^{tr}), S_i^{ts})$ to learn a good feature embedding. Formally, we define the sampling-with-replacement gradient update rule at $(t + 1)$ -th step as: $w_{t+1} = \text{Proj}_{\mathcal{W}}[w_t - \eta_t \partial_{w_t} \hat{L}(\mathbf{A}(\mathbf{S})(S_{i_t}^{tr}), S_{i_t}^{ts})]$, where i_t is independently and identically drawn (i.i.d.) from the uniform distribution $Unif([n])$. Therefore, although $\hat{L}(\mathbf{A}(\mathbf{S})(S_{i_t}^{tr}), S_{i_t}^{ts})$ is the loss only calculated over the query samples $S_{i_t}^{ts}$, it is still related to the support samples $S_{i_t}^{tr}$, and the updated parameter w_{t+1} is also related to $S_{i_t}^{tr}$. Therefore, we define an equivalent empirical loss $\hat{R}(\mathbf{A}(\mathbf{S})(S), S) \triangleq \hat{L}(\mathbf{A}(\mathbf{S})(S^{tr}), S^{ts})$ to indicate that: the empirical loss over the episode $S = S^{tr} \cup S^{ts}$ is related to the whole episode sample S , and so is the output hypothesis $\mathbf{A}(\mathbf{S})(S)$. Therefore, for the empirical error $\frac{1}{n} \sum_{i=1}^n \hat{R}(\mathbf{A}(\mathbf{S})(S_i), S_i)$, the episode-level SGD update rule is: $w_{t+1} = \text{Proj}_{\mathcal{W}}[w_t - \eta_t \partial_{w_t} \hat{R}(w_t, S_{i_t})]$. The pseudo code as well as several examples of modern meta learning algorithms can be found in Subsection 4.1. In this section, we provide lower and upper stability bounds for meta learning with sampling-with-replacement SGD method. First, we give the definition of uniform argument stability of episodic meta learning algorithms.

Definition 3 (Uniform argument stability of meta learning algorithms). Given a meta learning algorithm \mathbf{A} , any neighboring meta samples \mathbf{S}, \mathbf{S}' , and any training episode $S \in \mathcal{Z}^m$, we define the uniform argument stability random variable of \mathbf{A} as $\delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) = \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}')(S)\|$.

\mathbf{A} is defined as a uniform argument β -stable meta learning algorithm if for some $\beta > 0$, we have $\sup_{\mathbf{S} \approx \mathbf{S}', S} \delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) \leq \beta$ or $\sup_{\mathbf{S} \approx \mathbf{S}', S} \mathbb{E}_{\mathbf{A}} \delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) \leq \beta$, where $\mathbb{E}_{\mathbf{A}}$ denote the expectation w.r.t. the internal randomness of \mathbf{A} . For a meta learning algorithm with SGD method, the internal randomness of \mathbf{A} comes from the randomness of sampling at each iteration. Note that if $\hat{R}(\cdot, S)$ is a Lipschitz function for any $S \in \mathcal{Z}^m$, the uniform argument stability of \mathbf{A} implies the uniform stability of \mathbf{A} in Definition 2. In this work, we investigate the stability of modern meta learning algorithms with sampling-with-replacement SGD training strategy. Therefore, we will derive lower and upper bounds on $\mathbb{E}_{\mathbf{A}} \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}')(S)\|$ across different settings in Subsections 4.2-4.3.

4.1 Pseudo Code of Modern Meta Learning Algorithms

Algorithm 1 Support/Query Episodic Training based Meta Learning Algorithm

- 1: **Input:** training dataset $\mathbf{S} = \{S_i\}_{i=1}^n$ with $S_i = \{S_i^{tr}, S_i^{ts}\}$, # of iterations T , learning rates η_t ($t \in [T]$).
 - 2: **Initialize:** the parameters of deep neural networks w_1 .
 - 3: **for** $t = 1$ to T **do**
 - 4: Uniformly sample one of n training episodes with replacement. Let i_t be the episode index.
 - 5: $w_{t+1} = \text{Proj}_{\mathcal{W}}(w_t - \eta_t \partial \hat{R}(w_t, S_{i_t}))$ // episode-level SGD update
 - 6: **end for**
 - 7: **return** w_{T+1}
-

We provide several specific meta learning algorithms to illustrate the calculation of loss $\hat{R}(w_t, S_{i_t})$ on the episode S_{i_t} at the i_t -th iteration, where w_t is always the parameters of the feature extractor that is shared across different episodes. For the metric-learning based algorithms [44, 48] in classification, h_{w_t} is regarded as the feature extractor to output the feature vector $h_{w_t}(x)$ with data x as input. Then

$$\hat{R}(w_t, S_{i_t}) = \frac{1}{q} \sum_{(x,y) \in S_{i_t}^{ts}} -\log \frac{\exp\{-d(h_{w_t}(x), c_y)\}}{\sum_k \exp\{-d(h_{w_t}(x), c_k)\}},$$

where $c_k = \frac{1}{N\sigma m} \sum_{(x,y) \in S_{i_t}^{tr}, y=k} h_{w_t}(x)$ is the prototype (i.e., averaged vector) of the sample features in the support set $S_{i_t}^{tr}$ with the same class label k ; $d(\cdot, \cdot)$ is the distance between two feature vectors, e.g. the Euclidean distance in ProtoNet [44], and the Cosine distance in MatchingNet [48]. For the classifier-learning based meta learning algorithm MetaOptNet [32],

$$\hat{R}(w_t, S_{i_t}) = \frac{1}{q} \sum_{(x,y) \in S_{i_t}^{ts}} -\log \frac{\exp\{\lambda \langle h_{w_t}(x), \phi_y \rangle\}}{\sum_k \exp\{\lambda \langle h_{w_t}(x), \phi_k \rangle\}},$$

where $\{\phi_k\}_{k=1}^K$ are the parameters of the classifier returned by supervised learning algorithms (e.g. SVM) on the support set $S_{i_t}^{tr}$, $\langle \cdot, \cdot \rangle$ represents the inner product. For the gradient-learning based meta algorithm MAML [21], let α_t be the learning rate of the inner-task algorithm at the t -th iteration, then

$$\hat{R}(w_t, S_{i_t}) = \frac{1}{q} \sum_{z \in S_{i_t}^{ts}} f(w_t - \frac{\alpha_t}{K} \sum_{z' \in S_{i_t}^{tr}} \partial f(w_t, z'), z).$$

4.2 Stability Bounds for Nonsmooth Functions with α -Hölder Continuous Subgradients

In this subsection, we provide lower and upper stability bounds for episodic meta learning algorithm whose loss function is nonsmooth convex and has α -Hölder continuous subgradient with $0 \leq \alpha < 1$.

Theorem 2 \forall fixed $S \in \mathcal{Z}^m$, let $\hat{R}(\cdot, S)$ be a convex and (α, G) -Hölder smooth function, where $\alpha \in [0, 1)$. Let \mathbf{A} be a meta learning algorithm with sampling-with-replacement SGD. Denote by w_j and w'_j the outputs after j ($j \in [T]$) steps of SGD on \mathbf{S} and \mathbf{S}^i , respectively. Define $R_{\mathbf{S}}(w) = n^{-1} \sum_{i=1}^n \hat{R}(w, S_i)$, $\forall w \in \mathcal{W}$. Then $\forall S \in \mathcal{Z}^m$, $\mathbb{E}_{\mathbf{A}} \delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S)$ is upper bounded by

$$\sqrt{2}c_{\alpha} \left[\sum_{j=1}^T \eta_j^2 \mathbb{E} [R_{\mathbf{S}}^{\frac{2\alpha}{1+\alpha}}(w_j) + R_{\mathbf{S}^i}^{\frac{2\alpha}{1+\alpha}}(w'_j)] \right]^{\frac{1}{2}} + \frac{2c_{\alpha}}{n} \sum_{j=1}^T \eta_j [\hat{R}^{\frac{\alpha}{1+\alpha}}(w_j, S_i) + \hat{R}^{\frac{\alpha}{1+\alpha}}(w_j, S'_i)]. \quad (4)$$

In addition, if $\hat{R}(\cdot, S)$ is bounded by M and the step size $\eta_j = \eta \forall j \in [T]$, we can obtain the lower and upper bounds of the uniform argument stability of \mathbf{A} : $c_{\alpha} M^{\frac{\alpha}{1+\alpha}} (\min\{1, \frac{T}{n}\} \eta \sqrt{T} + \frac{\eta T}{n}) \leq \sup_{\mathbf{S}, \mathbf{S}', S} \mathbb{E}_{\mathbf{A}} \delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) \leq 4c_{\alpha} M^{\frac{\alpha}{1+\alpha}} (\min\{1, \frac{T}{n}\} \eta \sqrt{T} + \frac{\eta T}{n})$.

Remark 1 Our upper stability bound in Eq. (4) depends on the empirical risk during the optimization process. Formally, Eq. (4) shows that, the stability of modern meta algorithm increases if we find good parameters w_j with small empirical risk $R_{\mathbf{S}}(w_j)$ at the j -th optimization step. This illustrates a key insight that optimization is beneficial to improve the generalization of algorithms. Besides, our stability upper bound also implies the importance of a good embedding [45] (which may have a good initialization and low empirical risk during the first several optimization steps) to generalization.

Remark 2 We additionally suppose the function to be bounded by M such that the stability bounds can be used to analyze the generalization bounds in the next section where the loss function is always assumed to be bounded. Note that when $\alpha = 0$, $\hat{R}(\cdot, S)$ is a nonsmooth c_{α} -Lipschitz convex function, and our lower and upper stability bounds recover the results in [2]. For bounded convex α -Hölder smooth functions, the lower stability bound implies that modern meta learning algorithms are not stable enough even if we train them with a relatively small constant step size in each SGD iteration.

Remark 3 The above result shows that, for bounded convex α -Hölder smooth function ($\alpha \in [0, 1)$), the uniform argument stability parameter $\beta = O(c_{\alpha} M^{\frac{\alpha}{1+\alpha}} (\min\{1, T/n\} \eta \sqrt{T} + \eta T/n))$. Another work [33] also focuses on convex α -Hölder smooth function. Using the technique from [33], we obtain the upper stability bounds for bounded convex α -Hölder smooth function under the same conditions: $\beta \leq O(c_{\alpha} M^{\frac{\alpha}{1+\alpha}} (\eta^{\frac{1}{1-\alpha}} T + \eta T/n))$ or $\beta \leq O(c_{\alpha} M^{\frac{\alpha}{1+\alpha}} (\eta^{\frac{1}{1-\alpha}} \sqrt{T} + \eta \sqrt{T/n}))$ (see Theorems C.1-C.2 in Appendix C.2.2), both of which are larger than our tight stability bound in Theorem 2 under the setting $T \leq n$. When $T > n$, the upper bound $\beta \leq O(c_{\alpha} M^{\frac{\alpha}{1+\alpha}} (\eta^{\frac{1}{1-\alpha}} \sqrt{T} + \eta \sqrt{T/n}))$ in Theorem C.2 is slightly sharper than the upper bound $O(c_{\alpha} M^{\frac{\alpha}{1+\alpha}} (\eta \sqrt{T} + \eta T/n))$ in Theorem 2.

4.3 Stability Bounds for Smooth Functions

In this subsection, we provide lower and upper uniform argument stability bounds for modern meta learning algorithms with smooth functions. First, we consider smooth convex functions.

Theorem 3 \forall fixed $S \in \mathcal{Z}^m$, let $\hat{R}(\cdot, S)$ be a G -smooth convex function. Let \mathbf{A} be a meta learning algorithm with sampling-with-replacement SGD. Denote by w_j and w'_j the outputs after j ($j \in [T]$) steps of SGD on neighboring meta samples \mathbf{S} and \mathbf{S}^i , respectively. Then $\forall S \in \mathcal{Z}^m$, $\eta_j \leq 2/G$,

$$\mathbb{E}_{\mathbf{A}} \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}^i)(S)\| \leq \frac{\sqrt{2G}}{n} \sum_{j=1}^T \eta_j \mathbb{E}_{\mathbf{A}} [\sqrt{\hat{R}(w_j, S_i)} + \sqrt{\hat{R}(w'_j, S'_i)}].$$

In addition, if $\hat{R}(\cdot, S)$ is bounded by M , we can obtain the lower and upper bounds of the uniform argument stability of \mathbf{A} : $\frac{1}{n} \sum_{j=1}^T \eta_j \leq \sup_{\mathbf{S}, \mathbf{S}', S} \mathbb{E}_{\mathbf{A}} \delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) \leq \frac{2\sqrt{2MG}}{n} \sum_{j=1}^T \eta_j$.

If we set all $\eta_j = \eta$, then for bounded convex functions, the tight stability bound of $O(\frac{\eta T}{n})$ under the smooth case is sharper than the stability bound of $O(\min\{1, \frac{T}{n}\}\eta\sqrt{T} + \frac{\eta T}{n})$ in Theorem 2 under the nonsmooth case. This indicates that in the smooth case, meta learning algorithms are more stable than in the nonsmooth case. Finally, we give stability bounds for smooth non-convex functions.

Theorem 4 \forall fixed $S \in \mathcal{Z}^m$, let $\hat{R}(\cdot, S)$ be a σ -Lipschitz and G -smooth function. Let \mathbf{A} be a meta learning algorithm. Denote by w_j and w'_j the outputs after j ($j \in [T]$) steps of SGD on \mathbf{S} and \mathbf{S}^i , respectively. Define the learning rate $\eta_j = \frac{a}{j^G}$, $\forall j \in [T]$ with $a > 0$. Then $\forall S \in \mathcal{Z}^m$, the lower and upper stability bounds of \mathbf{A} satisfy: $\frac{T^a}{6n^{1+a}} \leq \sup_{\mathbf{S}, \mathbf{S}', S} \mathbb{E}_{\mathbf{A}} \delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) \leq \frac{11 \ln(n)\sigma T^a}{n^{1+a}}$.

Under the same step size setting, existing upper uniform argument stability bound in [9, Theorem3] or in [35, Proposition 4] for non-convex, smooth and Lipschitz function is of $O(T^{\frac{a}{1+a}}/n)$. Our bound of order $O(T^a/n^{1+a})$ is improved over the existing bound when $T^{\frac{a}{1+a}} \leq n$. Besides, our stability bound can be non-vacuous for multi-pass SGD setting (i.e., when $T = kn$, $k \in \mathbb{N}$) where the number T of SGD iterations is larger than n , as long as $k \leq n^{1/a}$.

5 High Probability Transfer Error Bounds for Meta Learning

In this section, we establish high probability bounds for transfer error $er(\mathbf{A}(\mathbf{S}), \tau)$. Specifically, we still consider two kinds of loss function: (1) convex and (α, G) -Hölder smooth function ($\alpha \in [0, 1]$); (2) non-convex, σ -Lipschitz and G -smooth function. We always assume that the loss function $\hat{R}(\cdot, \cdot)$ is bounded by M . Define $\sigma_{\alpha} = c_{\alpha} M^{\frac{\alpha}{1+\alpha}}$ if $\hat{R}(w, S)$ is a convex and (α, G) -Hölder smooth function; $\sigma_{\alpha} = \sigma$ if $\hat{R}(w, S)$ is a σ -Lipschitz and G -smooth function. We just exhibit the generalization bounds of randomized algorithm \mathbf{A} by supposing $\mathbb{P}_{\mathbf{A}}[\delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) > \beta] \leq \delta_0$. We provide an example to illustrate how to calculate δ_0 in Example D.1 in Appendix D. The generalization bound for deterministic meta algorithm (e.g. with gradient descent) can be stated by setting $\delta_0 = 0$.

5.1 Near Optimal Transfer Error Bound for Meta Learning with Independent Episodes

We denote by $a \lesssim b$ the existence of some universal constant $c > 0$ such that $a \leq cb$. Then we obtain the following near optimal bound of $O(1/\sqrt{n})$ under the independent task environment assumption.

Theorem 5 Let $\mathbf{A} \in \mathcal{A}(\mathcal{A}(\mathcal{H}, \mathcal{Z}), \mathcal{Z}^m)$ be a uniform argument β -stable meta algorithm, i.e., $\sup_{\mathbf{S} \simeq \mathbf{S}', S} \mathbb{E}_{\mathbf{A}} \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}')(S)\| \leq \beta$. For any $S \in \mathcal{Z}^m$, let $\hat{R}(\cdot, S)$ be $[0, M]$ -valued, and satisfy one of the two following conditions: (1) $\hat{R}(\cdot, S)$ is convex and (α, G) -Hölder smooth ($\alpha \in [0, 1]$); (2) $\hat{R}(\cdot, S)$ is σ -Lipschitz and G -smooth. Suppose $\mathbb{P}_{\mathbf{A}}[\delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) > \beta] \leq \delta_0$. Then for any independent task environment $\tau \in \mathcal{M}_1(\mathcal{M}_1(\mathcal{Z}))$, any $\delta \in (0, 1)$, the following holds with probability at least $1 - \delta - \delta_0$ over the draw of \mathbf{S} and the internal randomness of \mathbf{A} :

$$\sigma_{\alpha} \beta \ln \frac{1}{\delta} + \frac{M}{\sqrt{n}} \sqrt{\ln(1/\delta)} \lesssim er(\mathbf{A}(\mathbf{S}), \tau) - \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) \lesssim \sigma_{\alpha} \beta \ln \frac{n}{\delta} + \frac{M}{\sqrt{n}} \sqrt{\ln(1/\delta)}.$$

Remark 4 Our transfer error bound in Theorem 5 has three advantages over the bound in Theorem 1 from [9]: (1) Theorem 1 gives a high-probability upper bound of $O(\sqrt{n}\gamma + M/\sqrt{n})$ for transfer error, where γ is the uniform stability parameter and always scales as $O(1/n)$; in contrast, our upper bound of $O(\beta \ln n + M/\sqrt{n})$ is improved by replacing the \sqrt{n} factor before the stability parameter with $\ln n$. (2) In [9], the uniform stability $\gamma = O(T^{\frac{\alpha}{\alpha+1}}/n)$, whereas our uniform argument stability $\beta = O(T^a/n^{1+a})$ is tighter when $T^{\frac{\alpha}{\alpha+1}} \leq n$, i.e., when the uniform stability bound $\gamma = O(T^{\frac{\alpha}{\alpha+1}}/n)$ is non-vacuous. (3) Our high-probability transfer error bound of order $O(1/\sqrt{n})$ is near optimal.

Remark 5 We uncover two limitations of stability-based meta learning theory: (1) Recall the lower stability bound for meta learning algorithms with convex α -Hölder smooth function ($\alpha \in [0, 1)$) in Theorem 2, we find that the lower transfer error bound in Theorem 5 is $er(\mathbf{A}(\mathbf{S}), \tau) - \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) \gtrsim \sigma_\alpha \ln(1/\delta) c_\alpha M^{\frac{\alpha}{1-\alpha}} (\eta\sqrt{T} + \eta T/n)$ when $T \geq n$. This indicates that the lower transfer error bound is greater than a constant and will not converge to zero with the increase of n . Thus, the stability-based transfer error bound is vacuous and cannot provide asymptotic guarantees for convex Hölder smooth functions. (2) The stability-based transfer error bound of $O(1/\sqrt{n})$ in Theorem 5 is near optimal. Such result is assistant with the observation in [38, Section 2] that under the (i.i.d.) task environment assumption, the term $O(1/\sqrt{n})$ in the generalization bound is unavoidable. Thus, to obtain sharper generalization bounds for meta learning (e.g. the bound of $O(1/\sqrt{mn})$ or even $O(1/mn)$), we need to consider other stability notions (e.g. [17]), or suppose stronger task relatedness in the environment (e.g. [6, 23]), or even drop the task environment assumption (e.g. [15, 46]).

Remark 6 Under the independent task environment assumption, we compare our bound of $O(1/\sqrt{n})$ via S/Q episodic training strategy with other transfer error bounds that are obtained via traditional ERM strategy over all samples in training tasks. In detail, the bound from [36, Theorems 2 and 6] via algorithmic stability analysis is of $O(1/m + 1/\sqrt{n})$; the bounds from [41, Theorem 1] and [42, Theorem 2] via PAC-Bayes analysis are of $O(1/\sqrt{n} + 1/\sqrt{m})$; the bound from [23, Theorem 5] via covering number analysis is of $O(1/\sqrt{nm} + 1/\sqrt{m})$. All of these bounds via ERM strategy involve a term $O(1/\sqrt{m})$, and such term can be large when m is relatively small (e.g. $m = 5$ or $m = 10$ in the few-shot learning setting). Thus, in terms of the tightness of transfer error bounds, the S/Q episodic training strategy is superior to the ERM strategy for meta learning, when $m \ll n$. Such result was also pointed out by [9] and is more rigorously demonstrated in this work. Detailed comparisons between different transfer error bounds for meta learning are shown in Table A.2 in Appendix A.

5.2 Fast Transfer Error Bound of $O(\ln n/n)$ for Meta Learning with Independent Episodes

To obtain faster convergence rate, we need to take additional assumption. The generalized Bernstein condition is one of the most widely used condition to attain fast convergence rate of generalization bound in single-task learning [35, 29]. Next, we extend the generalized Bernstein condition to the meta learning setting, where we study the optimal algorithm A^* instead of the optimal hypothesis.

Definition 4 (Generalized Bernstein Condition for Meta Learning) Assume that $A^*(\mathcal{H}, \mathcal{Z}) = \text{Argmin}_{A \in \mathcal{A}(\mathcal{H}, \mathcal{Z})} er(A, \tau)$ is a set of risk minimizers in a closed set. We say that an algorithm A together with the environment τ and the empirical estimator $\mathbf{1}$ satisfies the generalized Bernstein condition if for some $B > 0$, $\forall A \in \mathcal{A}(\mathcal{H}, \mathcal{Z})$, there is a $A^* \in A^*(\mathcal{H}, \mathcal{Z})$, such that

$$\mathbb{E}_{S \sim \mathcal{D}_\tau} (\mathbf{1}(A, S) - \mathbf{1}(A^*, S))^2 \leq B(er(A, \tau) - er(A^*, \tau)). \quad (5)$$

[29] has shown that in single-task learning, a strongly-convex and Lipschitz function satisfies the generalized Bernstein condition. In this work, we relax the strong-convexity condition by considering the following Polyak-Łojasiewicz condition, one of the weakest curvature conditions of functions.

Definition 5 (Polyak-Łojasiewicz [49]) Any function $f : \mathcal{W} \rightarrow \mathbb{R}$ satisfies the Polyak-Łojasiewicz (PL) condition on \mathcal{W} with parameter $\mu > 0$ if for all $w \in \mathcal{W}$, $f(w) - f(w^*) \leq \frac{1}{2\mu} \|\partial^0 f(w)\|_2^2$, where w^* denotes the Euclidean projection of w onto the set of global minimizer of f in \mathcal{W} .

A key insight into the PL condition is that it is the sufficient and necessary condition to guarantee the linear convergence of gradient descent methods for smooth convex optimization problem [40]. Such PL condition can also be satisfied by many non-convex neural network models, such as the two-layer neural networks with ReLU activation functions [34] and the deep linear residual networks [24]. We will show that if the functions in Theorem 5 additionally satisfy the PL condition, then the loss functions in meta learning also satisfy the generalized Bernstein condition in Definition 4. Thus, we can derive a "deformed" transfer error bound of $O(\ln n/n)$ for modern meta learning algorithms.

Theorem 6 Under the same conditions of Theorem 5, for any fixed $S \in \mathcal{Z}^m$, let $\hat{R}(\cdot, S)$ additionally satisfy Polyak-Łojasiewicz condition in Definition 5. Suppose $\mathbb{P}_{\mathbf{A}}[\delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) > \beta] \leq \delta_0$. Then, there exist $c > 0$, such that for any environment $\tau \in \mathcal{M}_1(\mathcal{M}_1(\mathcal{Z}))$, and any $\delta \in (0, 1)$, the following holds with probability at least $1 - \delta - \delta_0$ over the draw of \mathbf{S} and the internal randomness of \mathbf{A} :

$$er(\mathbf{A}(\mathbf{S}), \tau) \leq (1 + \eta)\hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) + c(1 + 1/\eta) \left(\sigma_\alpha \beta \ln n + \frac{M}{n} \right) \ln \frac{1}{\delta}.$$

Recall in Section 4, our stability parameter is always $\beta = O(1/n)$. Hence, when the empirical error in the RHS of the above bound is close to zero, the transfer error bound always scales as $O(\ln n/n)$.

5.3 Transfer Error Bound for Meta Learning with Dependent Episodes

In this subsection, we investigate the generalization bound for meta learning algorithms with dependent episodes whose dependency relation can be characterized by a graph. The approach undertaken to establish our results is based on the forest approximation of the dependency graph [27]. Formally, a dependency graph is an undirected graph $\Gamma = (V, E)$ of a random vector $\mathbf{S} = (S_1, \dots, S_n)$ if the following two conditions are satisfied: (1) $V(\Gamma) = [n]$; (2) if $I, J \subset [n]$ are non-adjacent in Γ , then $\{S_i\}_{i \in I}$ and $\{S_j\}_{j \in J}$ are independent. We next give the concept of forest approximation.

Definition 6 (Forest Approximation [50]) *Given a graph Γ , a forest F , and a mapping $\phi : V(\Gamma) \rightarrow V(F)$, if $\phi(u) = \phi(v)$ or $\langle \phi(u), \phi(v) \rangle \in E(F)$ for any $\langle u, v \rangle \in E(\Gamma)$, we say that (ϕ, F) is a forest approximation of Γ . Let $\Phi(\Gamma)$ denote the set of forest approximations of Γ .*

Intuitively, a forest approximation transform a graph into a forest by merging vertices and removing self-loops. We then give the definition of forest complexity, which measures how a dependency graph looks like a forest, and hence measures the strength of dependency among random variables in \mathbf{S} .

Definition 7 (Forest Complexity [50]) *Given a graph Γ and any forest approximation $(\phi, F) \in \Phi(\Gamma)$ with F consisting of trees $\{T_i\}_{i \in [k]}$. Define $\lambda_{(\phi, F)} = \sum_{\langle u, v \rangle \in E(F)} (|\phi^{-1}(u)| + |\phi^{-1}(v)|)^2 + \sum_{i=1}^k \min_{u \in V(T_i)} |\phi^{-1}(u)|^2$. We call $\Lambda(\Gamma) = \min_{(\phi, F) \in \Phi(\Gamma)} \lambda_{(\phi, F)}$ the forest complexity of the graph $\Gamma = (V, E)$. Here, $\phi^{-1}(u)$ is the set of pre-images of the element u .*

For sample \mathbf{S} whose components are independent, we choose the identity map and its dependency graph as the forest approximation. Hence $\Lambda(\Gamma) = \sum_{i=1}^n 1^2 = n$. For sample \mathbf{S} whose dependency graph Γ is a tree, the identity map between Γ and itself is a forest approximation of Γ . Then $\Lambda(\Gamma) \leq |E(\Gamma)|(1+1)^2 + 1 = 4n - 3 = O(n)$. More examples of forest approximation can be found in [50, Section 3.3]. We next give a forest-complexity based transfer error bound for meta learning.

Theorem 7 *Under the same conditions of Theorem 5, except that \mathbf{S} is a meta sample of size n with dependency graph Γ . Let the maximum degree of the graph Γ is Δ . Suppose $\mathbb{P}_{\mathbf{A}}[\delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) > \beta] \leq \delta_0$. Then, for any environment $\tau \in \mathcal{M}_1(\mathcal{M}_1(\mathcal{Z}))$, any $\delta \in (0, 1)$, the following holds with probability at least $1 - \delta - \delta_0$ over the draw of \mathbf{S} and the internal randomness of \mathbf{A} :*

$$er(\mathbf{A}(\mathbf{S}), \tau) \leq \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) + \sigma_{\alpha} \beta (\Delta + 1) + \left(2\sigma_{\alpha} \beta + \frac{M}{n}\right) \sqrt{\frac{\Lambda(\Gamma) \ln 1/\delta}{2}},$$

When \mathbf{S} is an independent sample, the forest complexity $\Lambda(\Gamma) = n$, the maximum degree $\Delta = 0$, and the above forest-complexity based generalization bound degenerates to the bound in Theorem 1 for meta learning with independent episodes. When \mathbf{S} is a dependent sample, $\Lambda(\Gamma)$ will be greater than n . Both the complexity $\Lambda(\Gamma)$ and the maximum degree Δ will increase with more dependency relation between samples in \mathbf{S} (i.e., with more adjacent edges in its dependency graph $\Gamma = (V, E)$). In the next section, we conduct experiments on regression problems to show the convergence performance of the generalization bound for meta learning with dependent episodes. The corresponding forest-complexity based generalization bound for such problem is provided in Example D.2 in Appendix D.3.

6 Experiments

To verify our theoretical analysis, we conduct experiments on few-shot regression problems to show the convergence performance of our generalization bounds with independent and dependent episodes.

Experimental Settings. We follow the experimental setting in [21, 9]. The problem aims to approximate the distribution of parameters of function $f(x) = \alpha \sin(\beta x)$. The task environment τ is the joint distribution $p(\alpha, \beta)$ of the parameters α and β . We set $p(\alpha) = U[-5, 5]$, $p(\beta) = U[0, \pi]$. For independent setting, we construct training episodes by sampling pairs (α, β) from the task distribution $\tau = p(\alpha, \beta)$; for dependent setting, we construct the first half training episodes by sampling pairs (α, β) from $\tau = p(\alpha, \beta)$ independently, and construct the rest half training episodes by setting $(-\alpha, \pi - \beta)$ with (α, β) from the first half training episodes. Each training episode contains

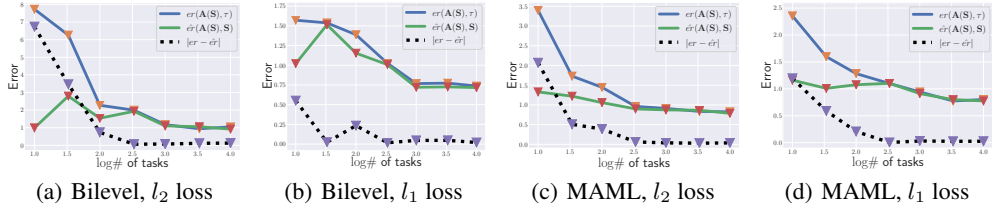


Figure 1: Convergence analysis of generalization gaps for independent tasks.

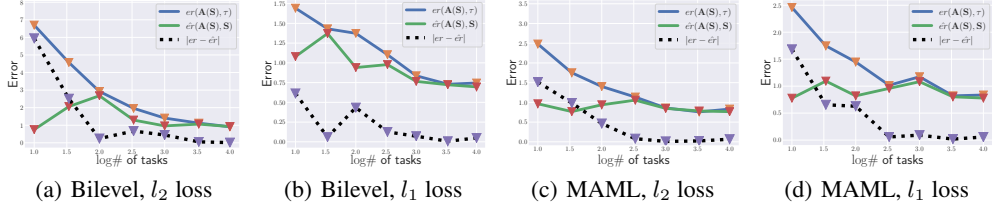


Figure 2: Convergence analysis of generalization gaps for dependent tasks.

5 support samples and 1 query samples. In both settings, the 600 test episodes are constructed by sampling (α, β) from τ independently, each containing 5 support samples and 15 query samples. We implement meta learning algorithms MAML [21] and Bilevel Programming [22] with square loss (l_2) and absolute loss (l_1). The neural network has two hidden layers of size 40 with ReLU activation functions. The generalization gap is the absolute distance between the training error and test error.

Experimental Results. From Figures 1-2, we can observe that: **(1)** The generalization gap in both independent and weakly dependent episode settings can converge to 0 with the increase of the training episodes, demonstrating the asymptotic behaviour of our transfer error bounds in Theorems 5 and 7. **(2)** The generalization gap with independent episodes can converge to zero more quickly than the gap with dependent episodes. The test error with independent episodes also always converge to the lower level than the one with dependent episodes. The better convergence performance with independent episodes truly demonstrate how the dependency between episodes can affect the generalization of meta learning algorithms. **(3)** With non-convex neural network models, both square loss and nonsmooth absolute loss can lead to similar convergence performance of generalization bounds.

7 Conclusion and Future Work

In this work, we provide fine-grained analysis of stability and generalization for modern meta learning algorithms. From the perspective of stability, our tight stability bounds implies that in the nonsmooth convex case the meta learning algorithm is less stable than in the smooth convex case. The stability bounds in the smooth non-convex case enjoys an order of $O(1/n)$ even for the multi-pass SGD setting. From the perspective of generalization, we demonstrate that the high-probability transfer error bound of $O(1/\sqrt{n})$ is optimal. Based on this bound, we uncover the limitations of algorithmic stability analysis for meta learning, and reveal the advantage of episodic training strategy for meta learning over tradition ERM training strategy. Further, by extending the generalized Bernstein condition to the meta learning setting, we obtain a deformed generalization bound of $O(\ln n/n)$ with additional Polyak-Łojasiewicz condition. Finally, we derive a generalization bound for meta learning with dependent episodes. Experiments are also provided to show the convergence performance of generalization error with independent and dependent episodes. In the future, we will explore new stability notions to see whether we can develop sharper generalization bounds for meta learning.

Acknowledgments and Disclosure of Funding

Jiechao would like to thank Dr. Mingxue Quan from School of Mathematics in Renmin University of China for helpful discussions. We thank all reviewers for their constructive suggestions to improve the quality of this paper. This work was supported in part by National Natural Science Foundation of China (61976220 and 61832017), Beijing Outstanding Young Scientist Program (BJJWZYJH012019100020098), and Large-Scale Pre-Training Program 468 of Beijing Academy of Artificial Intelligence (BAAI). Prof. Zhiwu Lu is the corresponding author of this paper.

References

- [1] M. Al-Shedivat, L. Li, E. Xing, and A. Talwalkar. On data efficiency of meta-learning. In *International Conference on Artificial Intelligence and Statistics (AISTAT)*, pages 1369–1377, 2021.
- [2] R. Bassily, V. Feldman, C. Guzmán, and K. Talwar. Stability of stochastic gradient descent on nonsmooth convex losses. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [3] J. Baxter. Learning internal representations. In *Conference on Learning Theory (COLT)*, pages 311–320, 1995.
- [4] J. Baxter. A model of inductive bias learning. *Journal of Artificial Intelligence Research*, 12:149–198, 2000.
- [5] A. Beck. *First-Order Methods in Optimization*. SIAM-Society for Industrial and Applied Mathematics, 2017.
- [6] S. Ben-David and R. Schuller. Exploiting task relatedness for multiple task learning. In *Conference on Learning Theory (COLT)*, pages 567–580, 2003.
- [7] O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research (JMLR)*, 2:499–526, 2002.
- [8] O. Bousquet, Y. Klochkov, and N. Zhivotovskiy. Sharper bounds for uniformly stable algorithms. In *Conference on Learning Theory (COLT)*, pages 610–626, 2020.
- [9] J. Chen, X. Wu, Y. Li, Q. LI, L. Zhan, and F. Chung. A closer look at the training strategy for modern meta-learning. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 396–406, 2020.
- [10] A. Christmann and I. Steinwart. *Support Vector Machines*. Springer, 2008.
- [11] Z. Deng, H. He, and W. J. Su. Toward better generalization bounds with locally elastic stability. In *International Conference on Machine Learning (ICML)*, pages 2590–2600, 2021.
- [12] J. Devlin, M. Chang, K. Lee, and K. Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 4171–4186, 2019.
- [13] L. Devroye and T. J. Wagner. Distribution-free inequalities for the deleted and holdout error estimates. *IEEE Transactions on Information Theory*, 25(2):202–207, 1979.
- [14] N. Ding, X. Chen, T. Levinboim, S. Goodman, and R. Soricut. Bridging the gap between practice and PAC-Bayes theory in few-shot meta-learning. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 29506–29516, 2021.
- [15] S. S. Du, W. Hu, S. M. Kakade, J. D. Lee, and Q. Lei. Few-shot learning via learning the representation, provably. In *International Conference on Learning Representations (ICLR)*, 2021.
- [16] A. Elisseeff, T. Evgeniou, and M. Pontil. Stability of randomized learning algorithms. *Journal of Machine Learning Research (JMLR)*, 6:55–79, 2005.
- [17] A. Fallah, A. Mokhtari, and A. E. Ozdaglar. Generalization of model-agnostic meta-learning algorithms: Recurring and unseen tasks. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [18] A. Farid and A. Majumdar. Generalization bounds for meta-learning via PAC-Bayes and uniform stability. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 2173–2186, 2021.

- [19] V. Feldman and J. Vondrák. Generalization bounds for uniformly stable algorithms. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 9770–9780, 2018.
- [20] V. Feldman and J. Vondrák. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Conference on Learning Theory (COLT)*, pages 1270–1279, 2019.
- [21] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, pages 1126–1135, 2017.
- [22] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning (ICML)*, pages 1563–1572, 2018.
- [23] J. Guan and Z. Lu. Task relatedness-based generalization bounds for meta learning. In *International Conference on Learning Representations (ICLR)*, 2022.
- [24] M. Hardt and T. Ma. Identity matters in deep learning. In *International Conference on Learning Representations (ICLR)*, 2017.
- [25] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning (ICML)*, pages 1225–1234, 2016.
- [26] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [27] S. Janson. Large deviations for sums of partly dependent random variables. *Random Structures & Algorithms*, 24(3):234–248, 2004.
- [28] H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *European Conference on Machine Learning (ECML)*, pages 795–811, 2016.
- [29] Y. Klochkov and N. Zhivotovskiy. Stability and deviation optimal risk bounds with convergence rate $o(1/n)$. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [30] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 1106–1114, 2012.
- [31] I. Kuzborskij and C. H. Lampert. Data-dependent stability of stochastic gradient descent. In *International Conference on Machine Learning (ICML)*, pages 2820–2829, 2018.
- [32] K. Lee, S. Maji, A. Ravichandran, and S. Soatto. Meta-learning with differentiable convex optimization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10657–10665, 2019.
- [33] Y. Lei and Y. Ying. Fine-grained analysis of stability and generalization for stochastic gradient descent. In *International Conference on Machine Learning (ICML)*, pages 5809–5819, 2020.
- [34] Y. Li and Y. Yuan. Convergence analysis of two-layer neural networks with relu activation. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 597–607, 2017.
- [35] T. Liu, G. Lugosi, G. Neu, and D. Tao. Algorithmic stability and hypothesis complexity. In *International Conference on Machine Learning (ICML)*, pages 2159–2167, 2017.
- [36] A. Maurer. Algorithmic stability and meta-learning. *Journal of Machine Learning Research (JMLR)*, 6:967–994, 2005.
- [37] A. Maurer. Transfer bounds for linear feature learning. *Machine Learning*, 75(3):327–350, 2009.
- [38] A. Maurer, M. Pontil, and B. Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research (JMLR)*, 17:81:1–81:32, 2016.

- [39] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. Adaptive computation and machine learning. MIT Press, 2012.
- [40] I. Necoara, Y. E. Nesterov, and F. Glineur. Linear convergence of first order methods for non-strongly convex optimization. *Mathematical Programming*, 175(1-2):69–107, 2019.
- [41] A. Pentina and C. H. Lampert. A PAC-Bayesian bound for lifelong learning. In *International Conference of Machine Learning (ICML)*, pages 991–999, 2014.
- [42] J. Rothfuss, V. Fortuin, M. Josifoski, and A. Krause. PACOH: Bayes-optimal meta-learning with PAC-Guarantees. In *International Conference on Machine Learning (ICML)*, pages 9116–9126, 2021.
- [43] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Learnability, stability and uniform convergence. *Journal of Machine Learning Research (JMLR)*, 11:2635–2670, 2010.
- [44] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 4077–4087, 2017.
- [45] Y. Tian, Y. Wang, D. Krishnan, J. B. Tenenbaum, and P. Isola. Rethinking few-shot image classification: A good embedding is all you need? In *European Conference on Computer Vision (ECCV)*, pages 266–282, 2020.
- [46] N. Tripuraneni, M. I. Jordan, and C. Jin. On the theory of transfer learning: The importance of task diversity. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 7852–7862, 2020.
- [47] A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Berlin: Springer, 1996.
- [48] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra. Matching networks for one shot learning. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 3630–3638, 2016.
- [49] H. Zhang. New analysis of linear convergence of gradient-type methods via unifying error bound conditions. *Mathematical Programming*, 180(1):371–416, 2020.
- [50] R. R. Zhang, X. Liu, Y. Wang, and L. Wang. Mcdiarmid-type inequalities for graph-dependent variables and stability bounds. In *Conference on Neural Information Processing Systems (NeurIPS)*, pages 10889–10899, 2019.
- [51] Y. Zhang, W. Zhang, S. Bald, V. Pingali, C. Chen, and M. Goswami. Stability of SGD: tightness analysis and improved bounds. *arXiv preprint arXiv:2102.05274*, 2021.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Remark 5.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Sections 4–5.
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See Appendix.
3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See Supplementary Material.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Section 6.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] See Figures 1-2.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Supplementary Material.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 6.
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

APPENDIX

A Different Transfer Error Bounds for Meta Learning

Table A.1: Three types of transfer error and the corresponding empirical multi-task error for meta learning. The empirical multi-task error is calculated over all samples in the training tasks.

Object	Transfer Error	Empirical Multi-Task Error
Hypothesis Space \mathcal{H}	$er(\mathcal{H}, \tau) = \mathbb{E}_{D \sim \tau} \inf_{h \in \mathcal{H}} L(h, D)$	$\hat{er}(\mathcal{H}, \mathbf{S}) = \frac{1}{n} \sum_{j=1}^n \inf_{h \in \mathcal{H}} \hat{L}(h, S_j)$
Hyper-Posterior \mathcal{Q}	$er(\mathcal{Q}, \tau) = \mathbb{E}_{P \sim \mathcal{Q}} \mathbb{E}_{D \sim \tau} \mathbb{E}_{S \sim D^m} \mathbb{E}_{h \sim Q(S, P)} L(h, D)$	$\hat{er}(\mathcal{Q}, \mathbf{S}) = \frac{1}{n} \sum_{j=1}^n \mathbb{E}_{P \sim \mathcal{Q}} \mathbb{E}_{h \sim Q(S_j, P)} \hat{L}(h, S_j)$
Algorithm A	$er(A, \tau) = \mathbb{E}_{D \sim \tau} \mathbb{E}_{S \sim D^m} L(A(S), D)$	$\hat{er}(A, \mathbf{S}) = \frac{1}{n} \sum_{j=1}^n \hat{L}(A(S_j), S_j)$

We describe in detail existing three main types of transfer error for meta learning in Table A.1. They include: (1) transfer error of hypothesis space \mathcal{H} in model capacity theory for meta learning [4, 23]; (2) transfer error of hyper-posterior \mathcal{Q} in PAC-Bayesian analysis for meta learning [41, 18]; (3) transfer error of an algorithm A in algorithmic stability analysis for meta learning [36]. The corresponding empirical multi-task errors are all calculated over all samples in the training tasks. Denote by $L(h, D) = \mathbb{E}_{z \sim D} f(h, z)$ the expected error over the underlying distribution D , $\hat{L}(h, S) = \frac{1}{m} \sum_{i=1}^m f(h, z_i)$ the empirical error over the training sample S . Let $\{S_j\}_{j=1}^n$ denote the training samples, where each $S_j \sim D_j^m$ contains m samples, D_j is the distribution sampled from the environment τ . In PAC-Bayesian meta learning theory, $Q(S, P)$ is the posterior distribution over the hypothesis space \mathcal{H} , i.e., $Q(S, P) \in \mathcal{M}_1(\mathcal{H})$. $Q(S, P)$ is the output of the PAC-Bayesian algorithm that takes sample S and prior distribution $P (\in \mathcal{M}_1(\mathcal{H}))$ as input. In algorithmic stability analysis for meta learning, an algorithm A takes sample S as input and then returns a hypothesis $A(S) \in \mathcal{H}$.

We also provide different transfer error bounds in Table A.2 to show the improvements of our bounds over existing ones. Under the task environment assumption, traditional transfer error bounds are obtained by minimizing the empirical risk over all samples in the training tasks; modern transfer error bounds are derived for support/query episodic training based meta learning algorithms. The comparisons between different transfer error bounds demonstrate that: the episodic training strategy can always lead to a generalization bound unrelated to the sample size m per task, thus outperforming the traditional ERM training strategy under the few-shot learning regime where m is small.

Table A.2: Different bounds on the generalization gap in meta learning, where generalization gap = transfer error - empirical multi-task error. The T-ERM training strategy represents traditional empirical risk minimization strategy over all samples in training tasks; the S/Q training strategy represents the support/query episodic training strategy for modern meta learning. All error bounds hold for $[0, M]$ -valued bounded functions, under the independent task environment assumption, with n training tasks and m samples per task. For the bound in [23, Theorem 5], c_1, c_2 represent norm-based complexity of neural networks. For the PAC-Bayes bounds in [41, 18], $\text{KL}(Q||P)$ represents the KL-divergence between distributions Q and P . For the stability-based bounds in [36, 9], $\gamma_n = O(\frac{1}{n})$ represents the uniform stability of a meta learning algorithm; $\gamma_m = O(\frac{1}{m})$ represents the uniform stability of an inner-task algorithm. In our bounds, $\beta_n = O(\frac{1}{n})$ represents the uniform argument stability of a meta learning algorithm, and the generalization bound of $O(\frac{1}{n})$ is obtained with additional Polyak-Łojasiewicz condition.

Existing Works	Object	Training Strategy	Generalization Gap	Bounds on Generalization Gap
[23, Theorem 5]	\mathcal{H}	T-ERM	$er(\mathcal{H}, \tau) - \hat{er}(\mathcal{H}, \mathbf{S})$	$O(\frac{c_1}{\sqrt{nm}} + \frac{c_2}{\sqrt{m}})$
[41, Theorem 1]	\mathcal{Q}	T-ERM	$er(\mathcal{Q}, \tau) - \hat{er}(\mathcal{Q}, \mathbf{S})$	$O(\frac{\text{KL}(\mathcal{Q} P)}{\sqrt{n}} + \frac{\mathbb{E}_{P \sim \mathcal{Q}} \text{KL}(Q(S_i, P) P)}{\sqrt{m}})$
[18, Theorem 3]	\mathcal{Q}	T-ERM	$er(\mathcal{Q}, \tau) - \hat{er}(\mathcal{Q}, \mathbf{S})$	$O(\sqrt{\frac{\text{KL}(\mathcal{Q} P)}{n}} + \gamma_m)$
[36, Theorem 6]	$\mathbf{A}(\mathbf{S})$	T-ERM	$er(\mathbf{A}(\mathbf{S}), \tau) - \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S})$	$O(\gamma_n \sqrt{n} + \frac{M}{\sqrt{n}} + \gamma_m)$
[9, Theorem 1]	$\mathbf{A}(\mathbf{S})$	S/Q	$er(\mathbf{A}(\mathbf{S}), \tau) - \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S})$	$O(\gamma_n \sqrt{n} + \frac{M}{\sqrt{n}})$
Our Theorem 5	$\mathbf{A}(\mathbf{S})$	S/Q	$er(\mathbf{A}(\mathbf{S}), \tau) - \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S})$	$O(\beta_n \ln n + \frac{M}{\sqrt{n}})$
Our Theorem 6	$\mathbf{A}(\mathbf{S})$	S/Q	$er(\mathbf{A}(\mathbf{S}), \tau) - \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S})$	$O(\beta_n \ln n + \frac{M}{n})$

Remark A.1 *It is still hard to directly compare our bounds with the latest stability-based generalization bound for meta learning [17], due to the following three reasons: (1) We focus on different bounding objective. Our work aims to bound the transfer error over the novel task (under the task environment assumption), whereas [17] aims to bound the (expected) excess risk over the novel task (without the task environment assumption, cf. its Corollary 2). (2) The generalization bounds hold with different forms. The bounds in our Theorem 5-7 all hold with high probability, but the generalization bounds in [17] (i.e. the bound on the gap between the expected multi-task error and empirical multi-task error in its Theorem 1, as well as the bound on the excess risk on the novel task) hold in expectation (w.r.t. all training samples). (3) We take different assumption of the loss function. In Assumption 1 of [17], they assume the loss function satisfy 4 conditions: strong convexity, Lipschitzness, smoothness and Hessian Lipschitzness. But our work only takes one or two conditions to derive stability of meta algorithm. Consider the above reasons, we believe it is not suitable to directly compare our in-probability generalization bound with the in-expectation bound of [17].*

Remark A.2 *We could potentially derive a sharper generalization bound for meta learning w.r.t. m in two aspects: (1) Under the task environment assumption: actually we can extend the algorithmic stability notions in single-task learning (e.g. uniform stability [7], uniform argument stability [35, 2], on-average stability [31], on-average model stability [33]) to the episodic meta learning setting by defining an algorithmic stability in a way that the whole dataset corresponding to a task changes (see our Definition 2 and Definition 3). However, no matter which stability notion we use, our Remark 5 tells us that the stability-based transfer error bound will not be tighter than $O(1/\sqrt{n})$. Therefore, to derive sharper transfer error bound (e.g. of $O(1/\sqrt{nm})$) for meta learning under the task environment assumption, we should leverage the tools of other theories (e.g. model-capacity theory in [23]), instead of the tool of algorithmic stability analysis. (2) Without the task environment assumption: without such assumption, we cannot define the transfer error of a meta algorithm on the novel task, so we should focus on the excess risk bound on the novel task. In this case, we can define a more elaborate algorithmic stability notion in a way that the part (not the whole) of dataset in a task change, and may derive a sharper bound that is related to the number m of samples per task (like the expected multi-task error bound of $O(\frac{1}{nm})$ in [17]).*

B Comparisons between Single-Task Learning and Modern Meta Learning

In this section, we provide the comparisons of notations between single-task learning and episodic meta learning in Table B.1. From this table, we can find the equivalence relation between these two learning paradigms, and hence can directly apply existing generalization bounds [8, 29, 50] from

single-task learning to the episodic meta learning setting. We next give formal definitions of the generalization error bound for single-task learning and the transfer error bound for meta learning.

Table B.1: The relation between the notations of single-task learning and support/query (S/Q) episodic training based meta learning. In meta learning, the empirical estimator $\mathbf{I}(\mathbf{A}(\mathbf{S}), S) = \hat{R}(\mathbf{A}(\mathbf{S})(S), S)$. Such Table is an adaptive version from [36, Table 1]. We list it here to show the equivalence relation between traditional single-task learning and modern episodic meta learning.

	Single-Task Learning	S/Q Training based Meta Learning
Sample	$z \in \mathcal{Z}$	$S = (z_1, \dots, z_m) \in \mathcal{Z}^m$
Training Set	$S = (z_1, \dots, z_m) \in \mathcal{Z}^m$	$\mathbf{S} = (S_1, \dots, S_n) \in (\mathcal{Z}^m)^n$
Hypothesis	$h \in \mathcal{H}$	$A \in \mathcal{A}(\mathcal{H}, \mathcal{Z})$
Algorithm	$A \in \mathcal{A}(\mathcal{H}, \mathcal{Z})$	$\mathbf{A} \in \mathcal{A}(\mathcal{A}(\mathcal{H}, \mathcal{Z}), \mathcal{Z}^m)$
Learning Task	$D \in \mathcal{M}_1(\mathcal{Z})$	$\mathbf{D} \in \mathcal{M}_1(\mathcal{Z}^m)$, typically $\mathbf{D} = \mathbf{D}_\tau$ is induced by the environment $\tau \in \mathcal{M}_1(\mathcal{M}_1(\mathcal{Z}))$.
Loss Estimator	$f : \mathcal{H} \times \mathcal{Z} \rightarrow [0, M]$	$\mathbf{I} : \mathcal{A}(\mathcal{H}, \mathcal{Z}) \times \mathcal{Z}^m \rightarrow [0, M]$
Empirical Error	$\hat{L}(A(S), S) = \frac{1}{m} \sum_{i=1}^m f(A(S), z_i)$	$\hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) = \frac{1}{n} \sum_{i=1}^n \mathbf{I}(\mathbf{A}(\mathbf{S}), S_i)$
Expected Error	$L(A(S), D) = \mathbb{E}_{z \sim D} f(A(S), z)$	$er(\mathbf{A}(\mathbf{S}), \tau) = \mathbb{E}_{S \sim \mathbf{D}_\tau} \mathbf{I}(\mathbf{A}(\mathbf{S}), S)$
Probability Bound	$D^m \{S : L(A(S), D) \geq B(\delta, S)\} \leq \delta$	$\mathbf{D}^n \{\mathbf{S} : R(\mathbf{A}(\mathbf{S}), \tau) \geq \Pi(\delta, \mathbf{S})\} \leq \delta$

Generalization Error and Transfer Error Bounds. A function $B : (0, 1) \times \cup_{m=1}^{\infty} \mathcal{Z}^m \rightarrow [0, M]$ is a generalization error bound for an algorithm $A \in \mathcal{A}(\mathcal{H}, \mathcal{Z})$ w.r.t. the loss function f if and only if

$$\forall D \in \mathcal{M}_1(\mathcal{Z}), \forall \delta \in (0, 1), D^m \{S : L(A(S), D) \leq B(\delta, S)\} \geq 1 - \delta. \quad (6)$$

A function $\Pi : (0, 1) \times \cup_{m=1}^{\infty} (\mathcal{Z}^m)^n \rightarrow [0, M]$ is a transfer error bound for a meta learning algorithm $\mathbf{A} \in \mathcal{A}(\mathcal{A}(\mathcal{H}), \mathcal{Z}, \mathcal{Z}^m)$ w.r.t. the estimator $\mathbf{I} : \mathcal{A}(\mathcal{H}, \mathcal{Z}) \times \mathcal{Z}^m \rightarrow [0, M]$ if and only if

$$\forall \mathbf{D} \in \mathcal{M}_1(\mathcal{Z}^m), \forall \delta \in (0, 1), \mathbf{D}^n \{\mathbf{S} : er(\mathbf{A}(\mathbf{S}), \tau) \leq \Pi(\delta, \mathbf{S})\} \geq 1 - \delta. \quad (7)$$

A transfer error bound is formally equivalent to an ordinary generalization error bound under the identifications $\mathcal{Z} \leftrightarrow \mathcal{Z}^m, h \leftrightarrow A, f \leftrightarrow \mathbf{I}, D \leftrightarrow \mathbf{D}, L \leftrightarrow er, B(\delta, S) \leftrightarrow \Pi(\delta, \mathbf{S})$.

C Proof of Stability Bounds for Meta Learning Algorithms

In this section, we will provide stability bounds for meta learning algorithms. According to the equivalence relation of the notations between single task learning and episodic meta learning, we will use the loss function $f(w, z)$ over the sample z and the loss function $\hat{R}(w, S)$ over the episode S interchangeably to compute the uniform argument stability bounds of (episode-level) SGD.

C.1 Proof of Sharp Stability Bounds for Convex Losses

In this section, we provide upper stability bounds for convex losses. The upper bounds are particularly tight for convex and nonsmooth losses (i.e., for convex and Lipschitz function, as well as for convex and α -Hölder smooth function ($\alpha \in (0, 1)$)). The proof technique is originated in [2] for obtaining upper stability bounds for convex and Lipschitz function, and in this work we generalize it to the convex setting for deriving upper stability bounds for convex and α -Hölder smooth functions ($\alpha \in [0, 1)$). Note that when $\alpha = 0$, the α -Hölder smooth property implies the Lipschitz property of the loss function. Therefore, our results can be used to analyze more general loss functions.

C.1.1 Proof of Sharp Stability Bounds for Convex Losses

The proof for the following lemma is almost the same as that for the Lemma 3.1 in [2], except that we do not assume the boundedness of the (sub)gradient of the loss function. We still include the detailed proof for the completeness of this paper. Such upper stability bound is particularly sharp

when $T > n$ for Lipschitz continuous function; when $T \leq n$, we will derive a sharper stability bound for Lipschitz continuous function in the next subsection.

Lemma C.1 *Let $(v_t)_{t \in [T]}$ and $(w_t)_{t \in [T]}$ with $v_1 = w_1$, be online (sub)gradient descent trajectories for convex objective $(f_t)_{i \in [T-1]}$ and $(f'_t)_{t \in [T-1]}$ respectively; i.e.,*

$$v_{t+1} = \text{Proj}_{\mathcal{W}}[v_t - \eta_t \partial f_t(v_t)], \quad w_{t+1} = \text{Proj}_{\mathcal{W}}[w_t - \eta_t \partial f'_t(w_t)],$$

for all $t \in [T-1]$. Suppose for every $t \in [T-1]$, let $\Delta_t = \|\partial f_t(v_t) - \partial f'_t(w_t)\|$, $\theta_t = \|\partial f_t(v_t) - \partial f'_t(v_t)\|$, $\delta_t = \|v_t - w_t\|$. Then if $t_0 = \inf\{t \geq 1 : f_t \neq f'_t\}$, we have

$$\|v_T - w_T\| \leq \sqrt{\sum_{j=t_0}^{T-1} \eta_j^2 \Delta_j^2} + 2 \sum_{j=t_0+1}^{T-1} \eta_j \theta_j.$$

Proof. By the definition of t_0 , we have $\delta_1 = \dots = \delta_{t_0}$. For $t = t_0 + 1$, we have $\delta_{t_0+1} = \|\eta_{t_0}(\partial f_{t_0}(v_{t_0}) - \partial f'_{t_0}(w_{t_0}))\|$. And

$$\begin{aligned} \delta_{t+1}^2 &= \|\text{Proj}_{\mathcal{W}}[v_t - \eta_t \partial f_t(v_t)] - \text{Proj}_{\mathcal{W}}[w_t - \eta_t \partial f'_t(w_t)]\|^2 \\ &\leq \|v_t - \eta_t \partial f_t(v_t) - w_t - \eta_t \partial f'_t(w_t)\|^2 \\ &= \delta_t^2 + \eta_t^2 \|\partial f_t(v_t) - \partial f'_t(w_t)\|^2 - 2\eta_t \langle v_t - w_t, \partial f_t(v_t) - \partial f'_t(w_t) \rangle \\ &= \delta_t^2 + \eta_t^2 \|\partial f_t(v_t) - \partial f'_t(w_t)\|^2 - 2\eta_t \langle v_t - w_t, \partial f_t(v_t) - \partial f'_t(v_t) \rangle - 2\eta_t \langle v_t - w_t, \partial f'_t(v_t) - \partial f'_t(w_t) \rangle \\ &\leq \delta_t^2 + \eta_t^2 \|\partial f_t(v_t) - \partial f'_t(w_t)\|^2 + 2\eta_t \|v_t - w_t\| \cdot \|\partial f_t(v_t) - \partial f'_t(v_t)\| \\ &= \delta_t^2 + \eta_t^2 \Delta_t^2 + 2\eta_t \delta_t \theta_t, \end{aligned} \tag{8}$$

where the last inequality holds due to the monotonicity of the subgradient of convex objectives. Unraveling the recursion and noticing $\delta_{t_0+1} = \eta_{t_0} \Delta_{t_0}$ gives

$$\delta_t^2 \leq \sum_{j=t_0}^{t-1} \eta_j^2 \Delta_j^2 + 2 \sum_{j=t_0+1}^{t-1} \eta_j \delta_j \theta_j. \tag{9}$$

Next we prove the following inequality by induction, which directly gives our results.

$$\delta_t \leq \sqrt{\sum_{j=t_0}^{t-1} \eta_j^2 \Delta_j^2} + 2 \sum_{j=t_0+1}^{t-1} \eta_j \theta_j.$$

The above claim clearly holds for $t = t_0$. Suppose the claim holds for any $t \in [T-1]$. For the $(t+1)$ -th step, we consider two cases. First when $\delta_{t+1} \leq \max_{j \in [t]} \delta_j$, by induction hypothesis,

$$\delta_{t+1} \leq \max_{j \in [t]} \delta_j \leq \sqrt{\sum_{j=t_0}^{t-1} \eta_j^2 \Delta_j^2} + 2 \sum_{j=t_0+1}^{t-1} \eta_j \theta_j \leq \sqrt{\sum_{j=t_0}^t \eta_j^2 \Delta_j^2} + 2 \sum_{j=t_0+1}^t \eta_j \theta_j.$$

In the other case when $\delta_{t+1} > \max_{s \in [t]} \delta_s$, we use the result in Eq. (9), and have

$$\delta_{t+1}^2 \leq \sum_{j=t_0}^t \eta_j^2 \Delta_j^2 + 2 \sum_{j=t_0+1}^t \eta_j \delta_j \theta_j \leq \sum_{j=t_0}^t \eta_j^2 \Delta_j^2 + 2\delta_{t+1} \left(\sum_{j=t_0+1}^t \eta_j \theta_j \right).$$

Rearrange the terms in the above inequality, we have

$$\left(\delta_{t+1} - \sum_{j=t_0+1}^t \eta_j \theta_j \right)^2 \leq \sum_{j=t_0}^t \eta_j^2 \Delta_j^2 + \left(\sum_{j=t_0+1}^t \eta_j \theta_j \right)^2.$$

Take square root of both sides and use the subadditivity of function $t \rightarrow \sqrt{t}$, we finish the proof. \square

C.1.2 Proof of Sharp Stability Bounds for Convex Losses with $T \leq n$

Lemma C.2 Suppose $T \leq n$. Let $S = (z_i)_{i \in [n]}$, $S^i = (z'_i)_{i \in [n]}$ be neighboring datasets that only differ on the i -th entry. Let $(i_t)_{t \geq 0} \stackrel{i.i.d.}{\sim} \text{Unif}([n])$. Let B_t denote the event that $i_s = i$ for some $s \leq t$. Let $(v_t)_{t \in [T]}$ and $(w_t)_{t \in [T]}$ with $v_1 = w_1$, be online (sub)gradient descent trajectories for convex objective, respectively; i.e.,

$$v_{t+1} = \text{Proj}_{\mathcal{W}}[v_t - \eta_t \partial f(v_t, z_{i_t})], \quad w_{t+1} = \text{Proj}_{\mathcal{W}}[w_t - \eta_t \partial f(w_t, z'_{i_t})],$$

for all $t \in [T-1]$. For every $t \in [T-1]$, let $\Delta_t = \|\partial f(v_t, z_{i_t}) - \partial f(w_t, z'_{i_t})\|$, $\theta_t = \|\partial f(v_t, z_{i_t}) - \partial f(v_t, z'_i)\|$, $\bar{\theta}_t = \|\partial f(v_t, z_i) - \partial f(v_t, z'_i)\|$, $\delta_t = \|v_t - w_t\|$. Then we have

$$\mathbb{E}[\delta_T] \leq \frac{T-1}{n} \left[\sqrt{\sum_{s=1}^{T-1} \mathbb{E}[\eta_s^2 \Delta_s^2 | B_s]} + \frac{2}{n} \sum_{s=1}^{T-1} \eta_s \sqrt{\mathbb{E}[\bar{\theta}_s^2 | B_{s-1}]} \right].$$

Proof. Note that B_t is the event that the index i is sampled at least once in the first t iterations. Denote by \bar{B}_t the complement of B_t . Then for any $t \in [T]$, we have

$$\mathbb{P}[B_t] = 1 - \mathbb{P}[\bar{B}_t] = 1 - \left(1 - \frac{1}{n}\right)^t \leq 1 - \left(1 - \frac{t}{n}\right) \leq \min\left(1, \frac{t}{n}\right).$$

Then we have

$$\begin{aligned} \mathbb{E}[\delta_T] &= \mathbb{E}[\delta_T | B_{T-1}] \cdot \mathbb{P}[B_{T-1}] + \mathbb{E}[\delta_T | \bar{B}_{T-1}] \cdot \mathbb{P}[\bar{B}_{T-1}] \\ &= \mathbb{E}[\delta_T | B_{T-1}] \cdot \mathbb{P}[B_{T-1}] \leq \min\left(1, \frac{T-1}{n}\right) \mathbb{E}[\delta_T | B_{T-1}]. \end{aligned} \quad (10)$$

For the rest of the proof we bound $\mathbb{E}[\delta_T | B_{T-1}]$. To this end, we need to derive recurrence for $\mathbb{E}[\delta_{t+1} | B_t]$. For convenience we consider bounding $\mathbb{E}[\delta_{t+1}^2 | B_t]$. As shown in Eq. (8), we have

$$\begin{aligned} \delta_{t+1}^2 &\leq \|v_t - w_t - \eta_t [\partial f(v_t, z_{i_t}) - \partial f(w_t, z'_{i_t})]\|^2 \\ &\leq \delta_t^2 + \eta_t^2 \|\partial f(v_t, z_{i_t}) - \partial f(w_t, z'_{i_t})\|^2 + 2\eta_t \|v_t - w_t\| \cdot \|\partial f(v_t, z_{i_t}) - \partial f(v_t, z'_{i_t})\| \\ &= \delta_t^2 + \eta_t^2 \Delta_t^2 + 2\eta_t \delta_t \theta_t. \end{aligned}$$

Note that $B_t = \left(\{i_t = i\} \cap \bar{B}_{t-1}\right) \cup \left(B_{t-1}\right)$. Then by the law of total expectation, we have

$$\begin{aligned} &\mathbb{E}[\delta_{t+1}^2 | B_t] \\ &= \mathbb{E}[\delta_{t+1}^2 | i_t = i, \bar{B}_{t-1}] \cdot \mathbb{P}[i_t = i, \bar{B}_{t-1} | B_t] + \mathbb{E}[\delta_{t+1}^2 | B_{t-1}] \cdot \mathbb{P}[B_{t-1} | B_t] \\ &\leq \mathbb{E}[\eta_t^2 \Delta_t^2 | i_t = i, \bar{B}_{t-1}] \cdot \mathbb{P}[i_t = i, \bar{B}_{t-1} | B_t] + \mathbb{E}[\delta_t^2 + \eta_t^2 \Delta_t^2 + 2\eta_t \delta_t \theta_t | B_{t-1}] \cdot \mathbb{P}[B_{t-1} | B_t] \\ &\leq \mathbb{E}[\eta_t^2 \Delta_t^2 | B_t] + \mathbb{E}[\delta_t^2 | B_{t-1}] + 2\eta_t \mathbb{E}[\delta_t \theta_t | B_{t-1}]. \end{aligned} \quad (11)$$

For simplicity, denote by $i_1^t = (i_1, \dots, i_t)$ the indices vector in the first t iterations. For the term $\mathbb{E}[\delta_t \theta_t | B_{t-1}]$ in the right-hand-side of the above inequality, using the independence between i_t and δ_t , as well as the independence between i_t and B_{t-1} , we have

$$\begin{aligned} \mathbb{E}[\delta_t \theta_t | B_{t-1}] &= \mathbb{E}_{i_1^t}[\delta_t \theta_t | B_{t-1}] \\ &= \mathbb{E}_{i_1^{t-1}}[\delta_t \mathbb{E}_{i_t}[\|\partial f(v_t, z_{i_t}) - \partial f(v_t, z'_{i_t})\|] | B_{t-1}] \\ &= \frac{1}{n} \mathbb{E}_{i_1^{t-1}}[\delta_t \|\partial f(v_t, z_i) - \partial f(v_t, z'_i)\| | B_{t-1}] \\ &= \frac{1}{n} \mathbb{E}_{i_1^{t-1}}[\delta_t \bar{\theta}_t | B_{t-1}]. \end{aligned}$$

Plug the above result into Eq. (11), and unravel the recursion, we have

$$\begin{aligned} \mathbb{E}[\delta_{t+1}^2 | B_t] &\leq \mathbb{E}[\eta_t^2 \Delta_t^2 | B_t] + \mathbb{E}[\delta_t^2 | B_{t-1}] + \frac{2\eta_t}{n} \mathbb{E}[\delta_t \bar{\theta}_t | B_{t-1}] \\ &\leq \sum_{s=1}^t \mathbb{E}[\eta_s^2 \Delta_s^2 | B_s] + \frac{2}{n} \sum_{s=1}^t \eta_s \mathbb{E}[\delta_s \bar{\theta}_s | B_{s-1}] \\ &\leq \sum_{s=1}^t \mathbb{E}[\eta_s^2 \Delta_s^2 | B_s] + \frac{2}{n} \sum_{s=1}^t \eta_s \left(\mathbb{E}[\delta_s^2 | B_{s-1}]\right)^{\frac{1}{2}} \left(\mathbb{E}[\bar{\theta}_s^2 | B_{s-1}]\right)^{\frac{1}{2}}, \end{aligned} \quad (12)$$

where the last inequality holds due to the Hölder inequality. We next proceed to prove the following claim by induction.

$$\sqrt{\mathbb{E}[\delta_{t+1}^2|B_t]} \leq \sqrt{\sum_{s=1}^t \eta_s^2 \mathbb{E}[\Delta_s^2|B_s]} + \frac{2}{n} \sum_{s=1}^t \eta_s \left(\mathbb{E}[\bar{\theta}_s^2|B_{s-1}] \right)^{\frac{1}{2}}. \quad (13)$$

For the base case $t = 0$, the claim clearly holds. Assume the claim holds for any $t \in [T - 1]$. For the $(t + 1)$ -th step, we consider two separate cases as follow.

(I): $\sqrt{\mathbb{E}[\delta_{t+1}^2|B_t]} \leq \max_{s \in [t]} \sqrt{\mathbb{E}[\delta_s^2|B_{s-1}]}$.

By the induction hypothesis, we have

$$\begin{aligned} \sqrt{\mathbb{E}[\delta_{t+1}^2|B_t]} &\leq \max_{s \in [t]} \sqrt{\mathbb{E}[\delta_s^2|B_{s-1}]} \\ &\leq \sqrt{\sum_{s=1}^{t-1} \eta_s^2 \mathbb{E}[\Delta_s^2|B_s]} + \frac{2}{n} \sum_{s=1}^{t-1} \eta_s \left(\mathbb{E}[\bar{\theta}_s^2|B_{s-1}] \right)^{\frac{1}{2}} \\ &\leq \sqrt{\sum_{s=1}^t \eta_s^2 \mathbb{E}[\Delta_s^2|B_s]} + \frac{2}{n} \sum_{s=1}^t \eta_s \left(\mathbb{E}[\bar{\theta}_s^2|B_{s-1}] \right)^{\frac{1}{2}}. \end{aligned}$$

(II): $\sqrt{\mathbb{E}[\delta_{t+1}^2|B_t]} > \max_{s \in [t]} \sqrt{\mathbb{E}[\delta_s^2|B_{s-1}]}$.

For this case, we use the result in Eq. (12).

$$\begin{aligned} \mathbb{E}[\delta_{t+1}^2|B_t] &\leq \sum_{s=1}^t \mathbb{E}[\eta_s^2 \Delta_s^2|B_s] + \frac{2}{n} \sum_{s=1}^t \eta_s \left(\mathbb{E}[\delta_s^2|B_{s-1}] \right)^{\frac{1}{2}} \left(\mathbb{E}[\bar{\theta}_s^2|B_{s-1}] \right)^{\frac{1}{2}} \\ &\leq \sum_{s=1}^t \mathbb{E}[\eta_s^2 \Delta_s^2|B_s] + \frac{2}{n} \left(\mathbb{E}[\delta_{t+1}^2|B_t] \right)^{\frac{1}{2}} \sum_{s=1}^t \eta_s \left(\mathbb{E}[\bar{\theta}_s^2|B_{s-1}] \right)^{\frac{1}{2}}. \end{aligned}$$

Rearrange the terms in the above inequality, we have

$$\left(\sqrt{\mathbb{E}[\delta_{t+1}^2|B_t]} - \frac{1}{n} \sum_{s=1}^t \eta_s \left(\mathbb{E}[\bar{\theta}_s^2|B_{s-1}] \right)^{\frac{1}{2}} \right)^2 \leq \sum_{s=1}^t \mathbb{E}[\eta_s^2 \Delta_s^2|B_s] + \left(\frac{1}{n} \sum_{s=1}^t \eta_s \left(\mathbb{E}[\bar{\theta}_s^2|B_{s-1}] \right)^{\frac{1}{2}} \right)^2.$$

Taking square root of both sides and using the subadditivity of function $t \rightarrow \sqrt{t}$ finish the proof for the claim in Eq. (13). From the Jensen inequality we have $\mathbb{E}[\delta_{t+1}|B_t] \leq \sqrt{\mathbb{E}[\delta_{t+1}^2|B_t]}$. Plugging the above results into Eq. (10) completes the whole proof. \square

Remark C.1 (Compare the upper bounds in Lemmas C.1-C.2 for convex and Lipschitz loss)
Recall Lemma C.1 and Lemma C.2, we have

$$T > n : \quad \mathbb{E}[\delta_T] \leq \sqrt{\sum_{s=1}^{T-1} \eta_s^2 \mathbb{E}[\Delta_s^2]} + 2 \sum_{s=1}^{T-1} \eta_s \mathbb{E}[\theta_s] = \sqrt{\sum_{s=1}^{T-1} \eta_s^2 \mathbb{E}[\Delta_s^2]} + \frac{2}{n} \sum_{s=1}^{T-1} \eta_s \mathbb{E}[\bar{\theta}_s],$$

$$T \leq n : \quad \mathbb{E}[\delta_T] \leq \frac{T-1}{n} \left[\sqrt{\sum_{s=1}^{T-1} \mathbb{E}[\eta_s^2 \Delta_s^2|B_s]} + \frac{2}{n} \sum_{s=1}^{T-1} \eta_s \sqrt{\mathbb{E}[\bar{\theta}_s^2|B_{s-1}]} \right].$$

where $\Delta_t = \|\partial f(v_t, z_{i_t}) - \partial f(w_t, z'_{i_t})\|$, $\theta_t = \|\partial f(v_t, z_{i_t}) - \partial f(v_t, z'_{i_t})\|$, $\bar{\theta}_t = \|\partial f(v_t, z_i) - \partial f(v_t, z'_i)\|$, $\delta_t = \|v_t - w_t\|$. If we further assume that the convex loss function $f(w, z)$ is σ -Lipschitz w.r.t. w , then we have $\|\partial f(w, z)\| \leq \sigma$, for any $z \in \mathcal{Z}$. Then $\Delta_t \leq 2\sigma$, $\theta_t \leq 2\sigma$. Plug these results into the above two inequalities, we have when $T > n$: $\mathbb{E}[\delta_T] \leq 2\sigma \sqrt{\sum_{s=1}^{T-1} \eta_s^2} + \frac{4\sigma}{n} \sum_{s=1}^{T-1} \eta_s$; when $T \leq n$: $\mathbb{E}[\delta_T] \leq \frac{T-1}{n} \left[2\sigma \sqrt{\sum_{s=1}^{T-1} \eta_s^2} + \frac{4\sigma}{n} \sum_{s=1}^{T-1} \eta_s \right]$. The above inequalities indicates that $\mathbb{E}[\delta_T] \leq 4\sigma \left(\min \left\{ 1, \frac{T}{n} \right\} \sqrt{\sum_{s=1}^T \eta_s^2} + \frac{1}{n} \sum_{s=1}^T \eta_s \right)$, which recovers the result in [2, Theorem 3.3]. And such upper bound is truly tight when compared with the lower stability bound in Lemma C.5. \square

Remark C.2 (Compare the upper bounds in Lemmas C.1-C.2 for convex and α -Hölder smooth loss ($\alpha \in (0, 1)$))

When $T \leq n$, we can obtain a sharper stability bound for convex and σ -Lipschitz loss function in Lemma C.2 than that in Lemma C.1. Concretely speaking, the sharper in Lemma C.2 is obtained by using the law of total expectation $\mathbb{E}[\delta_T] = \mathbb{E}[\delta_T|B_{T-1}] \cdot \mathbb{P}[B_{T-1}]$ to obtain the less-than-one factor $\frac{T-1}{n}$, (i.e., $\mathbb{P}[B_{T-1}] \leq \frac{T-1}{n}$), and using the subgradient boundedness to bound the conditional expectation $\mathbb{E}[\Delta_s^2|B_s] \leq 4\sigma^2$. However, for convex and α -Hölder smooth function, we have no more the subgradient boundedness assumption, and hence could not bound the conditional expectation $\mathbb{E}[\Delta_s^2|B_s]$ directly. By using the upper bound in Lemma C.2, and the total expectation formula $\mathbb{E}[\Delta_s^2|B_s]\mathbb{P}(B_s) = \mathbb{E}[\Delta_s^2]$, we have

$$\mathbb{E}[\delta_T] \leq \mathbb{P}(B_{T-1}) \left[\sqrt{\sum_{s=1}^{T-1} \eta_s^2 \frac{\mathbb{E}[\Delta_s^2]}{\mathbb{P}(B_s)}} + \frac{2}{n} \sum_{s=1}^{T-1} \eta_s \sqrt{\mathbb{E}[\bar{\theta}_s^2|B_{s-1}]} \right].$$

Since the term $\frac{\mathbb{P}(B_{T-1})}{\sqrt{\mathbb{P}(B_s)}}$ may be greater than 1, the bound for α -Hölder smooth function in Lemma C.2 is not necessarily sharper than that in Lemma C.1 when $T \leq n$. Actually, the "total-expectation-expansion" demonstration strategy in Lemma C.2 can lead to the same upper bound in Lemma C.1 for convex and α -Hölder smooth function ($\alpha \in (0, 1)$). Recall Eq. (11) we have

$$\begin{aligned} \mathbb{E}[\delta_{t+1}^2] &= \mathbb{E}[\delta_{t+1}^2|B_t]\mathbb{P}(B_t) \\ &= \left(\mathbb{E}[\delta_{t+1}^2|i_t = i, \overline{B_{t-1}}] \cdot \mathbb{P}[i_t = i, \overline{B_{t-1}}|B_t] + \mathbb{E}[\delta_{t+1}^2|B_{t-1}] \cdot \mathbb{P}[B_{t-1}|B_t] \right) \mathbb{P}(B_t) \\ &\leq \left(\mathbb{E}[\eta_t^2 \Delta_t^2|i_t = i, \overline{B_{t-1}}] \cdot \mathbb{P}[i_t = i, \overline{B_{t-1}}|B_t] + \mathbb{E}[\delta_t^2 + \eta_t^2 \Delta_t^2 + 2\eta_t \delta_t \theta_t|B_{t-1}] \cdot \mathbb{P}[B_{t-1}|B_t] \right) \mathbb{P}(B_t) \\ &= \mathbb{E}[\eta_t^2 \Delta_t^2] + \mathbb{E}[\delta_t^2] + 2\eta_t \mathbb{E}[\delta_t \theta_t] \\ &\leq \sum_{s=1}^t \mathbb{E}[\eta_s^2 \Delta_s^2] + \frac{2}{n} \sum_{s=1}^t \mathbb{E}[\delta_s \bar{\theta}_s], \end{aligned}$$

which is the same as the upper stability bound in Lemma C.1. \square

C.2 Proof of Stability Bounds for Convex and Hölder Smooth Losses

Lemma C.3 [5, Theorem 3.61] Let $f : \mathcal{W} \rightarrow (-\infty, +\infty]$ be a proper and convex function. Denote by $\partial f(w)$ the set of subgradients of f at $w \in \mathcal{W}$. Then if $\|g\| \leq \sigma$ for any $g \in \partial f(w)$, we have $|f(u) - f(v)| \leq \sigma \|u - v\|$ for any $u, v \in \mathcal{W}$.

Lemma C.4 [33, Lemma A.1] Define

$$c_\alpha = \begin{cases} (1 + 1/\alpha)^{\frac{\alpha}{1+\alpha}} G^{\frac{1}{1+\alpha}}, & \text{if } \alpha \in (0, 1] \\ \sup_z \|\partial f(0, z)\| + G, & \text{if } \alpha = 0. \end{cases} \quad (14)$$

Assume for all $z \in \mathcal{Z}$, the map $w \mapsto f(w, z)$ is nonnegative, and $w \mapsto \partial f(w, z)$ is (α, G) -Hölder continuous with $\alpha \in [0, 1]$. Then for c_α defined as in (14) we have

$$\|\partial f(w, z)\| \leq c_\alpha f^{\frac{\alpha}{1+\alpha}}(w, z), \quad \forall w \in \mathcal{W}, z \in \mathcal{Z}.$$

Lemma C.5 [2, Theorem 4.2] Let $\mathcal{W} \subseteq \mathbb{R}^d$ be a space uniformly bounded by θ , $f : \mathcal{W} \rightarrow \mathbb{R}$ be a σ -Lipschitz function (hence $f(w) \in [-\theta\sigma, \theta\sigma]$). For the sampling-with-replacement stochastic (sub)gradient descent algorithm with constant step size $\eta > 0$, there exist neighboring datasets $S \simeq S'$ of size n such that the uniform argument stability satisfies $\mathbb{E}\delta_A(S, S') \geq \sigma \left(\min\{1, \frac{T}{n}\} \eta \sqrt{T} + \frac{\eta T}{n} \right)$.

C.2.1 Proof of Sharp Stability Bounds for Convex and Hölder Smooth Losses

Proposition C.1 Let $S \simeq S^i$ be neighboring datasets, with $S = (z_1, \dots, z_i, \dots, z_n)$, $S^i = (z'_1, \dots, z'_i, \dots, z'_n)$ and $z_i \neq z'_i$. Let the empirical loss $R_S(w) = \frac{1}{n} \sum_{i=1}^n f(w, z_i)$, where $f : \mathcal{W} \times \mathcal{Z} \rightarrow [0, +\infty)$ is a (α, G) -Hölder smooth function ($\alpha \in [0, 1)$) w.r.t. w for any fixed $z \in \mathcal{Z}$

. Let the random sequence of indices used by the stochastic subgradient descent (SSD) algorithm $(i_t)_{t \geq 0} \stackrel{i.i.d.}{\sim} \text{Unif}([n])$. Under the conditions of Lemma C.1, let $f_t(\cdot) = f(\cdot, z_{i_t})$, $f'_t(\cdot) = f'(\cdot, z'_{i_t})$. Then the uniform argument stability of SSD can be bounded as follow:

$$\mathbb{E}_{A_{SSD}} \delta(S, S^i) \leq \left[2c_\alpha^2 \sum_{s=t_0}^T \eta_s^2 \mathbb{E} [R_S^{\frac{2\alpha}{1+\alpha}}(v_s) + R_{S^i}^{\frac{2\alpha}{1+\alpha}}(w_s)] \right]^{\frac{1}{2}} + \frac{2c_\alpha}{n} \sum_{s=t_0+1}^T \eta_s \mathbb{E} [f^{\frac{\alpha}{1+\alpha}}(v_s, z_i) + f^{\frac{\alpha}{1+\alpha}}(v_s, z'_i)]$$

Proof. Recalling Lemma C.1, and using the Jensen's inequality of square root function, we have

$$\mathbb{E}_A \|v_T - w_T\| \leq \sqrt{\sum_{s=t_0}^{T-1} \eta_s^2 \mathbb{E}_A \Delta_s^2} + 2 \sum_{s=t_0+1}^{T-1} \eta_s \mathbb{E}_A \theta_s,$$

where $\Delta_s = \|\partial f(v_s, z_{i_s}) - \partial f(w_s, z'_{i_s})\|$, $\theta_s = \|\partial f(v_s, z_{i_s}) - \partial f(v_s, z'_{i_s})\|$. We next bound $\mathbb{E} \Delta_s^2$ and $\mathbb{E} \theta_s$ respectively. For $\mathbb{E} \Delta_s^2$, using inequality $(a+b)^2 \leq 2(a^2 + b^2)$ and the self-bounding property of Hölder smooth function in Lemma C.4, we have $\Delta_s^2 = \|\partial f(v_s, z_{i_s}) - \partial f(w_s, z'_{i_s})\|^2 \leq 2(\|\partial f(v_s, z_{i_s})\|^2 + \|\partial f(w_s, z'_{i_s})\|^2) \leq 2c_\alpha^2 (f^{\frac{2\alpha}{1+\alpha}}(v_s, z_{i_s}) + f^{\frac{2\alpha}{1+\alpha}}(w_s, z'_{i_s}))$. Then we have

$$\begin{aligned} \mathbb{E} \Delta_s^2 &\leq 2c_\alpha^2 \mathbb{E} [f^{\frac{2\alpha}{1+\alpha}}(v_s, z_{i_s}) + f^{\frac{2\alpha}{1+\alpha}}(w_s, z'_{i_s})], \\ &\leq 2c_\alpha^2 \left[(\mathbb{E} f(v_s, z_{i_s}))^{\frac{2\alpha}{1+\alpha}} + (\mathbb{E} f(w_s, z'_{i_s}))^{\frac{2\alpha}{1+\alpha}} \right] \\ &= 2c_\alpha^2 \left[(\mathbb{E} R_S(v_s))^{\frac{2\alpha}{1+\alpha}} + (\mathbb{E} R_{S^i}(w_s))^{\frac{2\alpha}{1+\alpha}} \right], \end{aligned}$$

where the second inequality uses Jensen's inequality of the concave function $t \rightarrow t^{\frac{2\alpha}{1+\alpha}}$ ($\alpha \in [0, 1]$), the last equality holds since the random variable $v_s(w_s)$ is independent of i_s .

For $\mathbb{E} \theta_s$, notice that with probability $1 - \frac{1}{n}$, $z_{i_s} = z'_{i_s}$, and then $\|\partial f(v_s, z_{i_s}) - \partial f(v_s, z'_{i_s})\| = 0$; with probability $\frac{1}{n}$, $z_{i_s} \neq z'_{i_s}$ (i.e., $z_{i_s} = z_i, z'_{i_s} = z'_i$), and then $\|\partial f(v_s, z_{i_s}) - \partial f(v_s, z'_{i_s})\| \leq \|\partial f(v_s, z_{i_s})\| + \|\partial f(v_s, z'_{i_s})\| \leq c_\alpha (f^{\frac{\alpha}{1+\alpha}}(v_s, z_i) + f^{\frac{\alpha}{1+\alpha}}(v_s, z'_i))$. Therefore we have

$$\mathbb{E} \theta_s \leq \frac{c_\alpha}{n} \mathbb{E} (f^{\frac{\alpha}{1+\alpha}}(v_s, z_i) + f^{\frac{\alpha}{1+\alpha}}(v_s, z'_i)).$$

Combining the above analysis completes the whole proof. \square

Proof of Theorem 2 in the main paper. For the first part of the result, the proof can be found in Proposition C.1. For the second part of the result when the loss function is uniformly bounded by M , we first give the upper stability. Recalling Proposition C.1 and letting $t_0 = 1$, we have,

$$\mathbb{E} \delta_T \leq \sqrt{4c_\alpha^2 \sum_{s=1}^T \eta_s^2 M^{\frac{2\alpha}{1+\alpha}} + \frac{4M^{\frac{\alpha}{1+\alpha}} c_\alpha}{n} \sum_{s=1}^T \eta_s} \leq 4c_\alpha M^{\frac{\alpha}{1+\alpha}} \left(\sqrt{\sum_{s=1}^T \eta_s^2} + \frac{1}{n} \sum_{s=1}^T \eta_s \right). \quad (15)$$

Actually, when the loss function $\hat{R}(w, S)$ is bounded by M and convex (α, G) -Hölder smooth, we can use the self-bounding property of the (α, G) -Hölder smooth function in Lemma C.4, and derive the subgradient bounded norm $\|\partial \hat{R}(w, S)\| \leq c_\alpha M^{\frac{\alpha}{1+\alpha}}$. Recalling Lemma C.3, we know that the bounded convex (α, G) -Hölder smooth function also satisfies $c_\alpha M^{\frac{\alpha}{1+\alpha}}$ -Lipschitz property. Therefore, $\hat{R}(w, S)$ is a convex Lipschitz but nonsmooth function (i.e., α -Hölder smooth ($\alpha \in [0, 1]$) is nonsmooth), and we can use the results in Lemma C.5 and Remark C.1 to derive lower and upper stability bounds for convex (α, G) -Hölder smooth function $\hat{R}(w, S)$: $c_\alpha M^{\frac{\alpha}{1+\alpha}} (\min\{1, \frac{T}{n}\} \eta \sqrt{T} + \frac{\eta T}{n}) \leq \sup_{S \simeq S', S^{tr}} \mathbb{E}_A \delta_A(S, S'; S^{tr}) \leq 4c_\alpha M^{\frac{\alpha}{1+\alpha}} (\min\{1, \frac{T}{n}\} \sqrt{\sum_{j=1}^T \eta_j^2} + \frac{1}{n} \sum_{j=1}^T \eta_j)$. Note that such upper bound also match the upper bound obtained in Eq. (15) by directly setting $\hat{R}(w, S) \leq M$ in Proposition C.1 without using the Lipschitz property of the bounded convex Hölder smooth function, somewhat implying the tightness of the upper bound in Eq. (15). \square

C.2.2 An Alternative Proof of Upper Stability Bounds for Convex and Hölder Smooth Losses

In this subsection, we use the technique from [33] to derive upper stability bounds for convex and α -Hölder smooth function. We will derive two upper bounds, where the former is likely to be sharper when $T \leq n$, and the second one will be a little sharper when $T > n$.

Lemma C.6 [33, Lemma D.3] Assume for all $z \in \mathcal{Z}$, the map $w \rightarrow f(w, z)$ is convex, and $w \rightarrow \partial f(w, z)$ is (α, G) -Hölder continuous with $\alpha \in [0, 1)$. Define $G_\alpha \triangleq \sqrt{\frac{1-\alpha}{1+\alpha}} (2^{-\alpha} G)^{\frac{1}{1-\alpha}}$. Then for all $u, v \in \mathcal{W}$ and $\eta > 0$ there holds

$$\|u - \eta \partial f(u, z) - v + \eta \partial f(v, z)\|^2 \leq \|u - v\|^2 + \eta^{\frac{2}{1-\alpha}} G_\alpha^2.$$

Theorem C.1 \forall fixed $S \in \mathcal{Z}^m$, let $\hat{R}(\cdot, S)$ be a convex and (α, G) -Holder smooth function, where $\alpha \in [0, 1)$. Let \mathbf{A} be a meta learning algorithm, $\mathbf{S} \simeq \mathbf{S}^i$ be any neighboring meta samples that differ only on the i -th entry. Denote by w_j and w'_j the outputs after j ($j \in [T]$) steps of SGD on \mathbf{S} and \mathbf{S}^i , respectively. Then $\forall S^{tr} \in \mathcal{Z}^K$,

$$\mathbb{E}_{\mathbf{A}} \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}^i)(S)\| \leq G_\alpha \sum_{j=1}^T \eta_j^{\frac{1}{1-\alpha}} + \frac{c_\alpha}{n} \sum_{j=1}^T \eta_j \mathbb{E}_{\mathbf{A}} [\hat{R}^{\frac{\alpha}{1+\alpha}}(w_j, S_i) + \hat{R}^{\frac{\alpha}{1+\alpha}}(w'_j, S'_i)].$$

If we set $\eta_t = \eta$, for any $t \in [T]$, and assume the loss function \hat{R} is bounded by M , we have

$$\mathbb{E}_{\mathbf{A}} \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}^i)(S)\| \leq O\left(G_\alpha \eta^{\frac{1}{1-\alpha}} T + \frac{c_\alpha M^{\frac{\alpha}{1+\alpha}} \eta T}{n}\right).$$

Proof. With probability $(1 - \frac{1}{n})$, $i_t \neq i$, and

$$\begin{aligned} \|w_{t+1} - w'_{t+1}\| &\leq \|w_t - \eta_t \partial \hat{R}(w_t, S_{i_t}) - w'_t + \eta_t \partial \hat{R}(w'_t, S_{i_t})\| \\ &\leq \|w_t - w'_t\| + G_\alpha \eta_t^{\frac{1}{1-\alpha}}, \end{aligned}$$

where the second inequality holds due to Lemma C.6.

With probability $\frac{1}{n}$, $i_t = i$, and

$$\begin{aligned} \|w_{t+1} - w'_{t+1}\| &\leq \|w_t - \eta_t \partial \hat{R}(w_t, S_i) - w'_t + \eta_t \partial \hat{R}(w'_t, S'_i)\| \\ &\leq \|w_t - w'_t\| + \eta_t \|\partial \hat{R}(w_t, S_i)\| + \eta_t \|\partial \hat{R}(w'_t, S'_i)\| \\ &\leq \|w_t - w'_t\| + \eta_t c_\alpha [\hat{R}^{\frac{\alpha}{1+\alpha}}(w_t, S_i) + \hat{R}^{\frac{\alpha}{1+\alpha}}(w'_t, S'_i)], \end{aligned}$$

where the last inequality holds due to the self-bounding property of (α, G) -Hölder smooth function in Lemma C.4. Combining the above results and the iteration rules, we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{A}} \|w_{t+1} - w'_{t+1}\| \\ &\leq \mathbb{E}_{\mathbf{A}} \|w_t - w'_t\| + \frac{\eta_t c_\alpha}{n} \mathbb{E}_{\mathbf{A}} [\hat{R}^{\frac{\alpha}{1+\alpha}}(w_t, S_i) + \hat{R}^{\frac{\alpha}{1+\alpha}}(w'_t, S'_i)] + (1 - \frac{1}{n}) G_\alpha \eta_t^{\frac{1}{1-\alpha}} \\ &\leq (1 - \frac{1}{n}) G_\alpha \sum_{j=1}^t \eta_j^{\frac{1}{1-\alpha}} + \frac{c_\alpha}{n} \sum_{j=1}^t \eta_j \mathbb{E}_{\mathbf{A}} [\hat{R}^{\frac{\alpha}{1+\alpha}}(w_j, S_i) + \hat{R}^{\frac{\alpha}{1+\alpha}}(w'_j, S'_i)]. \quad \square \end{aligned}$$

Theorem C.2 Under the same conditions of Theorem C.1, if we set the step size $\eta_t = \eta$, for any $t \in [T]$, and assume the loss function \hat{R} is bounded by M , we have

$$\mathbb{E}_{\mathbf{A}} \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}^i)(S)\| \leq O\left(G_\alpha \eta^{\frac{1}{1-\alpha}} \sqrt{T} + \frac{c_\alpha M^{\frac{\alpha}{1+\alpha}} \eta \sqrt{T}}{\sqrt{n}}\right)$$

Proof. With probability $(1 - \frac{1}{n})$, $i_t \neq i$, and hence $\|w_{t+1} - w'_{t+1}\|^2 \leq \|w_t - w'_t\|^2 + G_\alpha \eta_t^{\frac{2}{1-\alpha}}$, where the inequality holds due to Lemma C.6.

With probability $\frac{1}{n}$, $i_t = i$, and

$$\begin{aligned} \|w_{t+1} - w'_{t+1}\|^2 &\leq \|w_t - \eta_t \partial \hat{R}(w_t, S_i) - w'_t + \eta_t \partial \hat{R}(w'_t, S'_i)\|^2 \\ &\leq (1+p) \|w_t - w'_t\|^2 + 2(1+p^{-1}) \eta_t^2 (\|\partial \hat{R}(w_t, S_i)\|^2 + \|\partial \hat{R}(w'_t, S'_i)\|^2), \end{aligned}$$

where the second inequality holds due to $(a+b)^2 \leq (1+p)a^2 + (1+p^{-1})b^2$, for any $p > 0$. Then,

$$\begin{aligned} &\mathbb{E} \|w_{t+1} - w'_{t+1}\|^2 \\ &\leq (1 + \frac{p}{n}) \left(\mathbb{E} \|w_t - w'_t\|^2 + G_\alpha^2 \eta_t^{\frac{2}{1-\alpha}} \right) + \frac{2(1+p^{-1})c_\alpha^2 \eta_t^2}{n} \mathbb{E} [\hat{R}^{\frac{2\alpha}{1+\alpha}}(w_t, S_i) + \hat{R}^{\frac{2\alpha}{1+\alpha}}(w'_t, S'_i)]. \end{aligned}$$

Multiply both sides with $(1 + \frac{p}{n})^{-(t+1)}$, and denote by $\xi_t = \mathbb{E}[\hat{R}^{\frac{2\alpha}{1+\alpha}}(w_t, S_t) + \hat{R}^{\frac{2\alpha}{1+\alpha}}(w'_t, S'_t)]$ for simplicity, we have

$$\begin{aligned} & \mathbb{E}(1 + \frac{p}{n})^{-(t+1)} \|w_{t+1} - w'_{t+1}\|^2 \\ & \leq (1 + \frac{p}{n})^{-t} \left(\mathbb{E} \|w_t - w'_t\|^2 + G_\alpha^2 \eta_t^{\frac{2}{1-\alpha}} \right) + \frac{2(1+p^{-1})c_\alpha^2 \eta_t^2 (1 + \frac{p}{n})^{-(t+1)}}{n} \xi_t \\ & \leq \sum_{s=1}^t (1 + \frac{p}{n})^{-s} G_\alpha^2 \eta_s^{\frac{2}{1-\alpha}} + \sum_{s=1}^t \frac{2(1+p^{-1})c_\alpha^2 \eta_s^2 (1 + \frac{p}{n})^{-(s+1)}}{n} \xi_s. \end{aligned}$$

Therefore, we have

$$\mathbb{E} \|w_{t+1} - w'_{t+1}\|^2 \leq \sum_{s=1}^t (1 + \frac{p}{n})^{t+1-s} G_\alpha^2 \eta_s^{\frac{2}{1-\alpha}} + \sum_{s=1}^t \frac{2(1+p^{-1})c_\alpha^2 \eta_s^2 (1 + \frac{p}{n})^{t-s}}{n} \xi_s.$$

Take the Jensen inequality of the function $t \rightarrow \sqrt{t}$, we have

$$\begin{aligned} \mathbb{E} \|w_{t+1} - w'_{t+1}\| & \leq \sqrt{\mathbb{E} \|w_{t+1} - w'_{t+1}\|^2} \\ & \leq G_\alpha \sqrt{\sum_{s=1}^t (1 + \frac{p}{n})^{t+1-s} \eta_s^{\frac{2}{1-\alpha}}} + \frac{c_\alpha}{\sqrt{n}} \sqrt{\sum_{s=1}^t 2(1+p^{-1})\eta_s^2 (1 + \frac{p}{n})^{t-s} \xi_s}. \end{aligned}$$

Plug the step size $\eta_s = \eta$, and $\hat{R} \leq M$ into the above inequality, we have

$$\begin{aligned} \mathbb{E} \|w_{t+1} - w'_{t+1}\| & \leq G_\alpha \eta^{\frac{1}{1-\alpha}} \sqrt{\sum_{s=1}^t (1 + \frac{p}{n})^{t+1-s}} + \frac{2\sqrt{(1+p^{-1})}\eta\sigma_\alpha}{\sqrt{n}} \sqrt{\sum_{s=1}^t (1 + \frac{p}{n})^{t-s}} \\ & \leq G_\alpha \eta^{\frac{1}{1-\alpha}} \sqrt{\frac{n+p}{p} (1 + \frac{p}{n})^t} + \frac{2\sqrt{(1+p^{-1})}\eta\sigma_\alpha}{\sqrt{n}} \sqrt{\frac{n}{p} (1 + \frac{p}{n})^t}. \end{aligned}$$

Setting $p = \frac{n}{t-1}$ to minimize $\frac{n}{p} (1 + \frac{p}{n})^t$, we obtain

$$\mathbb{E} \|w_{t+1} - w'_{t+1}\| \leq O\left(G_\alpha \eta^{\frac{1}{1-\alpha}} \sqrt{t} + \frac{\eta c_\alpha M^{\frac{\alpha}{1+\alpha}} \sqrt{t}}{\sqrt{n}}\right). \quad \square$$

C.3 Proof of Stability Bounds for Convex and Smooth Losses

Lemma C.7 [25, Lemma 3.6] Assume for all $z \in \mathcal{Z}$, the map $w \mapsto f(w, z)$ is G -smooth. Then for all $u, v \in \mathcal{W}$,

$$(1) \|u - \eta \partial f(u, z) - v + \eta \partial f(v, z)\| \leq (1 + \eta G) \|u - v\|.$$

(2) If in addition the map $w \mapsto f(w, z)$ is convex, then for any $\eta \leq 2/G$ we have:

$$\|u - \eta \partial f(u, z) - v + \eta \partial f(v, z)\| \leq \|u - v\|.$$

Lemma C.8 [51, Theorem 1][25, Theorem 3.8] Assume for all $z \in \mathcal{Z}$, $f(\cdot, z)$ is convex, G -smooth, σ -Lipschitz. Let w_t, w'_t be the outputs of SGD on neighboring datasets S, S' respectively, with the step size $\eta_t \leq 2/G$. Then the uniform stability parameter β of SGD satisfies

$$\frac{\sigma}{n} \sum_{t=1}^T \eta_t \leq \beta \leq \frac{2\sigma^2}{n} \sum_{t=1}^T \eta_t$$

Proof of Theorem 3 in the main paper. We analyze the situation at the $(t+1)$ -th iteration of SGD. According to the sampling-with-replacement strategy, with probability $(1 - \frac{1}{n})$ we have $i_t \neq i$, and

$$\|w_{t+1} - w'_{t+1}\| \leq \|w_t - \eta_t \partial \hat{R}(w_t, S_{i_t}) - w'_t + \eta_t \partial \hat{R}(w'_t, S_{i_t})\| \leq \|w_t - w'_t\|,$$

where the first inequality holds due to the nonexpansiveness of the projection operator $\text{Proj}_{\mathcal{W}}(\cdot)$, and the last inequality holds due to Lemma C.7.

With probability $\frac{1}{n}$, $i_t = i$, and

$$\begin{aligned} \|w_{t+1} - w'_{t+1}\| &\leq \|w_t - \eta_t \partial \hat{R}(w_t, S_{i_t}) - w'_t + \eta_t \partial \hat{R}(w'_t, S_{i_t})\| \\ &\leq \|w_t - w'_t\| + \|\eta_t \partial \hat{R}(w_t, S_{i_t}) - \eta_t \partial \hat{R}(w'_t, S_{i_t})\| \\ &\leq \|w_t - w'_t\| + \eta_t \|\partial \hat{R}(w_t, S_{i_t})\| + \eta_t \|\partial \hat{R}(w'_t, S'_{i_t})\| \\ &\leq \|w_t - w'_t\| + \sqrt{2G} \eta_t (\sqrt{\hat{R}(w_t, S_{i_t})} + \sqrt{\hat{R}(w'_t, S'_{i_t})}), \end{aligned}$$

where the last inequality holds due to the self-bounding property of smooth function in Lemma C.4. Combining the above analysis and the iteration rules, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{A}} \|w_{t+1} - w'_{t+1}\| &\leq \mathbb{E}_{\mathbf{A}} \|w_t - w'_t\| + \frac{\sqrt{2G} \eta_t}{n} \mathbb{E}_{\mathbf{A}} (\sqrt{\hat{R}(w_t, S_{i_t})} + \sqrt{\hat{R}(w'_t, S'_{i_t})}) \\ &\leq \frac{\sqrt{2G}}{n} \sum_{j=1}^t \eta_j \mathbb{E}_{\mathbf{A}} (\sqrt{\hat{R}(w_j, S_{i_j})} + \sqrt{\hat{R}(w'_j, S'_{i_j})}), \end{aligned}$$

which finishes the proof of the first part.

For the second part of the proof, if we additionally assume the loss function $\hat{R}(w, S)$ is uniformly bounded by M , then from Lemma C.4, we have $\alpha = 1$ and $\forall w \in \mathcal{W}, \|\partial \hat{R}(w, S)\| \leq c_\alpha \hat{R}^{\frac{1-\alpha}{1+\alpha}}(w, S) \leq \sqrt{2GM}$. Recalling Lemma C.3 and the convexity of the loss function, $\hat{R}(w, S)$ is a $\sqrt{2GM}$ -Lipschitz function w.r.t. w . Combining the $\sqrt{2GM}$ -Lipschitzness, G -smoothness of $\hat{R}(\cdot, S)$, and Lemma C.8, we obtain lower and upper bounds of uniform argument stability of \mathbf{A} . \square

C.4 Proof of Stability Bounds in Non-Convex Case

Lemma C.9 [51, Theorems 4-6] Assume $f(\cdot, z)$ is G -smooth and σ -Lipschitz for all $z \in \mathcal{Z}$. Running $T \geq n$ iterations of SGD with step size $\eta_t = \frac{\alpha}{Gt}$, then the uniform stability γ of SGD satisfies

$$\frac{\sigma T^\alpha}{6n^{1+\alpha}} \leq \gamma \leq 11 \ln(n) \sigma^2 \frac{T^\alpha}{n^{1+\alpha}}.$$

Proof of Theorem 4 in the main paper. Note that the loss function $\hat{R}(w, S)$ is σ -Lipschitz and G -smooth. Then using the techniques from Lemma C.9, we can obtain the lower and upper bound of uniform argument stability of \mathbf{A} . \square

D Proof of High Probability Transfer Error Bounds for Meta Learning

In this section, we just present the transfer error bounds for deterministic meta algorithm \mathbf{A} with the following form: $\forall \delta \in (0, 1), \mathbb{P}_{\mathbf{S}} \{\mathbf{S} : \text{er}(\mathbf{A}(\mathbf{S}), \tau) \leq \Pi(\delta, \mathbf{S})\} \geq 1 - \delta$. If \mathbf{A} is a randomized meta learning algorithm, for any neighboring meta samples \mathbf{S}, \mathbf{S}' , we suppose the *uniform argument stability* random variable satisfies $\mathbb{P}_{\mathbf{A}}[\delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) > \beta] \leq \delta_0$. Then, we can obtain the following probability bound for randomized meta learning algorithm \mathbf{A} with union bound technique:

$$\mathbb{P}_{\mathbf{A}, \mathbf{S}} \{\text{er}(\mathbf{A}(\mathbf{S}), \tau) \geq \Pi(\delta, \mathbf{S})\} \leq \delta + \delta_0.$$

We next give an example to illustrate how to obtain the probability bound $\mathbb{P}_{\mathbf{A}}[\delta_{\mathbf{A}}(\mathbf{S}, \mathbf{S}'; S) > \beta] \leq \delta_0$, where $\delta_0 = \exp\{-\frac{n}{2}\}$, n is the number of training episodes.

Example D.1 Under the same conditions of Lemma C.1, we further assume the function $f(\cdot, z)$ is bounded by M and has (α, G) -Hölder continuous subgradient. We run stochastic subgradient descent with constant step size η on neighboring datasets S and S^i . Then with probability at least $1 - \exp\{-\frac{n}{2}\}$, we have $\|v_T - w_T\| \leq O(c_\alpha M^{\frac{\alpha}{1+\alpha}} \eta (\sqrt{T-1} + \frac{T-1}{n}))$.

Proof. From Section C.2, we know $f(\cdot, z)$ has a subgradient bounded by $c_\alpha M^{\frac{\alpha}{1+\alpha}}$. Then from Lemma C.1 we have $\|v_T - w_T\| \leq 2c_\alpha M^{\frac{\alpha}{1+\alpha}} \eta \sqrt{T-1} + 4c_\alpha M^{\frac{\alpha}{1+\alpha}} \eta \sum_{j=1}^{T-1} \mathbf{r}_j$. Define random variable $\mathbf{r}_j = \mathbf{1}_{\{i_j=i\}}$, we have the expectation $\mathbb{E} \mathbf{r}_j = \frac{1}{n}$, and the variance $\mathbb{D}(\mathbf{r}_j) = \frac{1}{n}(1 - \frac{1}{n})$. Then according to the Bernstein inequality for rv with Gaussian behavior [39, Section D.4], we have

$$\mathbb{P} \left[\sum_{j=1}^{T-1} \mathbf{r}_j - \frac{T-1}{n} \geq \sqrt{T-1} \right] \leq \exp \left\{ -\frac{(T-1) \left(\frac{1}{\sqrt{T-1}} \right)^2}{2 \frac{1}{n} \left(1 - \frac{1}{n} \right)} \right\} \leq \exp \left\{ -\frac{n}{2} \right\}.$$

Therefore, with probability at least $1 - \exp\{-\frac{n}{2}\}$, we have $\sum_{j=1}^{T-1} \mathbf{r}_j \leq \frac{T-1}{n} + \sqrt{T-1}$, and

$$\|v_T - w_T\| \leq 2c_\alpha M^{\frac{\alpha}{1+\alpha}} \eta \sqrt{T-1} + 4c_\alpha M^{\frac{\alpha}{1+\alpha}} \eta \left(\frac{T-1}{n} + \sqrt{T-1} \right). \quad \square$$

D.1 Proof of Bound with Near Optimal Rate

Definition D.1 (Uniformly Stable Algorithm [7]) An algorithm is called γ -uniformly stable if for any S, S' that differ only at most one entry, for any $z \in \mathcal{Z}$, $|f(A(S), z) - f(A(S'), z)| \leq \gamma$.

Lemma D.1 [8, Corollary 8 and Proposition 9] Let A be a γ -uniformly stable algorithm. Let $f(\cdot, \cdot)$ be a loss function uniformly bounded by M . Then $\forall \delta \in (0, 1)$, with probability at least $1 - \delta$ over the sample S ,

$$\gamma \ln \frac{1}{\delta} + \frac{M}{\sqrt{n}} \sqrt{\ln \frac{1}{\delta}} \lesssim L(A(S), D) - \hat{L}(A(S), S) \lesssim \gamma \ln n \ln \frac{1}{\delta} + \frac{M}{\sqrt{n}} \sqrt{\ln \frac{1}{\delta}}.$$

Proof of Theorem 5 in the main paper. The result clearly holds for σ -Lipschitz and G -smooth function. Next, just consider the convex and (α, G) -Hölder smooth function. By utilizing the self-bounding property of (α, G) -Hölder smooth function in Lemma C.4, we have

$$\|\partial \hat{R}(w, S)\|_2 \leq c_\alpha \hat{R}^{\frac{\alpha}{1+\alpha}}(w, S) \leq c_\alpha M^{\frac{\alpha}{1+\alpha}}.$$

Therefore, $\hat{R}(w, S)$ has a bounded subgradient. Then recalling the convexity of the loss function and Lemma C.3, we know $\hat{R}(w, S)$ is a $c_\alpha M^{\frac{\alpha}{1+\alpha}}$ -Lipschitz function w.r.t. w . Thus, for any neighboring meta samples \mathbf{S}, \mathbf{S}' , $\forall S^{tr} \in \mathcal{Z}^K$,

$$\hat{R}(\mathbf{A}(\mathbf{S})(S), S) - \hat{R}(\mathbf{A}(\mathbf{S}')(S), S) \leq c_\alpha M^{\frac{\alpha}{1+\alpha}} \|\mathbf{A}(\mathbf{S})(S) - \mathbf{A}(\mathbf{S}')(S)\| \leq c_\alpha M^{\frac{\alpha}{1+\alpha}} \beta.$$

Denote by $\sigma_\alpha = c_\alpha M^{\frac{\alpha}{1+\alpha}}$ for simplicity. Thus, the meta learning algorithm \mathbf{A} is also uniformly $\sigma_\alpha \beta$ -stable w.r.t. the loss function \hat{L} as defined in Definition 2. Recalling the relationship between generalization error bound for single-task learning and transfer error bound for meta learning in Table B.1, and utilizing the near optimal bound in Lemma D.1, we obtain the following near-optimal high probability bound for transfer error:

$$\sigma_\alpha \beta \ln \frac{1}{\delta} + \frac{M}{\sqrt{n}} \sqrt{\ln \frac{1}{\delta}} \lesssim er(\mathbf{A}(\mathbf{S}), \tau) - \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) \lesssim \sigma_\alpha \beta \ln \frac{n}{\delta} + \frac{M}{\sqrt{n}} \sqrt{\ln \frac{1}{\delta}}. \quad \square$$

D.2 Proof of Bound with Fast Rate $O(\ln n/n)$

Definition D.2 (Generalized Bernstein Condition for Single-Task Learning) Assume that $\mathcal{W}^* = \text{Argmin}_{w \in \mathcal{W}} L(w, D)$ is a set of risk minimizers in a closed set \mathcal{W} . We say that \mathcal{W} together with the measure D and the loss function f satisfy the generalized Bernstein assumption if for some $B > 0$ for any $w \in \mathcal{W}$, there is a $w^* \in \mathcal{W}^*$ such that

$$\mathbb{E}_{z \sim D} (f(w, z) - f(w^*, z))^2 \leq B(L(w, D) - L(w^*, D)).$$

Definition D.3 (Quadratic Growth Condition [40]) Any function $f : \mathcal{W} \rightarrow \mathbb{R}$ satisfies the quadratic growth (QG) condition on \mathcal{W} with parameter $\mu > 0$ if for all $w \in \mathcal{W}$, $f(w) - f(w^*) \geq \frac{\mu}{2} \|w - w^*\|_2^2$, where w^* denotes the Euclidean projection of w onto the set of global minimizer of f in \mathcal{W} .

Lemma D.2 [29, Theorem 1.2] There is a constant $c > 0$ such that the following holds. Let A be a uniformly γ -stable algorithm, and assume that the loss function $f(\cdot, \cdot)$ is bounded by M and satisfies the generalized Bernstein condition in Def D.2. Fix any $\eta > 0$. Then, with probability at least $1 - \delta$,

$$L(A(S), D) \leq (1 + \eta) \hat{L}(A(S), S) + c \left(1 + \frac{1}{\eta}\right) \left(\gamma \ln n + \frac{M}{n}\right) \ln \frac{1}{\delta}.$$

Recall that $\sigma_\alpha = c_\alpha M^{\frac{\alpha}{1+\alpha}}$ if $\hat{R}(w, S)$ is a convex and (α, G) -Hölder smooth function; $\sigma_\alpha = \sigma$ if $\hat{R}(w, S)$ is a σ -Lipschitz and G -smooth function. Then we have the following proposition.

Proposition D.1 *Suppose that for any fixed $S \in \mathcal{Z}^m$, the loss function $\hat{R}(w, S)$ is bounded by M , has a quadratic growth on w with parameter μ (see Definition D.3), and satisfies one the two following conditions: (1) (α, G) -Hölder smooth; (2) L -Lipschitz and G -smooth. Then we can derive the generalized Bernstein condition for meta learning in Definition 4 with parameter $B = \frac{2\sigma_\alpha^2}{\mu}$.*

Proof. Recall the proof for Theorem 5, we know that for any fixed $S \in \mathcal{Z}^m$, $\hat{R}(w, S)$ is a σ_α -Lipschitz function w.r.t. w . Then we have

$$\begin{aligned} \mathbb{E}_{S \sim \mathcal{D}_\tau} \left(\mathbf{I}(A, S) - \mathbf{I}(A^*, S) \right)^2 &= \mathbb{E}_{S \sim \mathcal{D}_\tau} \left(\hat{R}(A(S), S) - \hat{R}(A^*(S), S) \right)^2 \\ &\leq \mathbb{E}_{S \sim \mathcal{D}_\tau} \sigma_\alpha^2 \|A(S) - A^*(S)\|_2^2 \\ &\leq \mathbb{E}_{S \sim \mathcal{D}_\tau} \frac{2}{\mu} \sigma_\alpha^2 \left(\hat{R}(A(S), S) - \hat{R}(A^*(S), S) \right) \\ &= \frac{2\sigma_\alpha^2}{\mu} \left[er(A, \tau) - er(A^*, \tau) \right], \end{aligned}$$

where the first inequality holds due to the Lipschitz property of function \hat{R} , the second inequality holds due to the Quadratic Growth of function \hat{R} w.r.t. its first argument. \square

Proof of Theorem 6 in the main paper. We first show that the loss function $\hat{R}(w, S)$ also satisfies the quadratic growth (QG) condition in Definition D.3. (1) when $\hat{R}(w, S)$ is a convex α -Hölder smooth function: from [49, Corollary 2] we know that the PL condition is equivalent to the QG condition. Therefore, $\hat{R}(w, S)$ satisfies the QG condition with parameter μ ; (2) when $\hat{R}(w, S)$ is a non-convex smooth function: the smooth function $\hat{R}(w, S)$ satisfies the PL condition for all $w \in \mathcal{W}$, $\hat{R}(w, S) - \hat{R}(w^*, S) \leq \frac{1}{2\mu} \|\nabla \hat{R}(w, S)\|_2^2$. From [28, Appendix A], we also have $\hat{R}(w, S)$ satisfies the QG condition in Definition D.3 with parameter μ . From Proposition D.1, we know $\hat{R}(\cdot, S)$ satisfies generalized Bernstein condition with parameter $\frac{2\sigma_\alpha^2}{\mu}$. From Theorem 5, we know the meta learning algorithm \mathbf{A} is uniformly stable with parameter $\sigma_\alpha\beta$. Therefore, recalling Lemma D.2, as well as the ‘equivalence’ relationship between generalization error bound and transfer error bound listed in Table B.1, we complete the whole proof. \square

D.3 Proof of Bound with Dependent Learning Tasks

Lemma D.3 [50, Theorem 4.4] *Given a sample S of size n with dependency graph Γ , assume that the learning algorithm A is γ -uniformly stable. Suppose the maximum degree of Γ is Δ , and the loss function f is bounded by M . For any $\delta \in (0, 1)$, with probability at least $1 - \delta$, it holds that*

$$L(A(S), \mathcal{D}) \leq \hat{L}(A(S), S) + \gamma(\Delta + 1) + \left(2\gamma + \frac{M}{n}\right) \sqrt{\frac{\Lambda(\Gamma) \ln(1/\delta)}{2}}.$$

Proof of Theorem 7 in the main paper. According to the proof for Theorem 5, the meta learning algorithm \mathbf{A} is uniformly $(\sigma_\alpha\beta)$ -stable w.r.t. the loss function \hat{R} . Combining the above result with Lemma D.3, we have with probability at least $1 - \delta$ over the draw of meta sample \mathbf{S} ,

$$er(\mathbf{A}(\mathbf{S}), \tau) \leq \hat{er}(\mathbf{A}(\mathbf{S}), \mathbf{S}) + \sigma_\alpha\beta(\Delta + 1) + \left(2\sigma_\alpha\beta + \frac{M}{n}\right) \sqrt{\frac{\Lambda(\Gamma) \ln 1/\delta}{2}},$$

which completes the whole proof. \square

Example D.2 *In the experiment section, we run meta learning algorithms to approximate the distribution $p(\alpha, \beta)$ of parameters α and β , with $p(\alpha) = U[-5, 5]$, $p(\beta) = U[0, \pi]$. To construct dependent episodes, we first independently sample n pairs of parameters (α, β) from $p(\alpha, \beta)$ to form the first n training episodes, and set $(-\alpha, \pi - \beta)$ with pair (α, β) from the first n training episodes to form the rest n training episodes. In this problem, the dependency graph $\Gamma = (V, E)$ of \mathbf{S} satisfies: $V = [2n]$, $|E| = n$ with each edge in E connecting only two vertices. Then we can connect all edges in E in series to construct a tree F of depth $(2n - 1)$ as one forest approximation of Γ . Then the maximum degree $\Delta = 1$, the forest complexity of Γ satisfies $\Lambda(\Gamma) \leq |F|(1 + 1)^2 + 1 = 8n - 3 = O(n)$. Thus, the forest-complexity based bounds for meta algorithms with dependent episodes is of $O(\beta\sqrt{n} + M/\sqrt{n})$.*