

Efficient Meta Learning via Minibatch Proximal Update



Pan Zhou* Xiao-Tong Yuan[†] Huan Xu[‡] Shuicheng Yan[△] Jiashi Feng*
^{*} National University of Singapore, [†] Nanjing University of Information Science & Technology, [‡] Alibaba, [△] YITU Tech (pzhou@u.nus.edu xtyuan@nuist.edu.cn Huan.xu@alibaba-inc.com {eleyans, elefjia}@nus.edu.sg)

Problem Setup

Problems: how to learn task-level knowledge from various observed tasks such that it could facilitate new task learning, e.g. few-shot learning and new policy learning in RL?

Proposed Efficient Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx: it learns a good prior model **w** from observed tasks that is close to the optimal models of new similar tasks, facilitating new task learning.

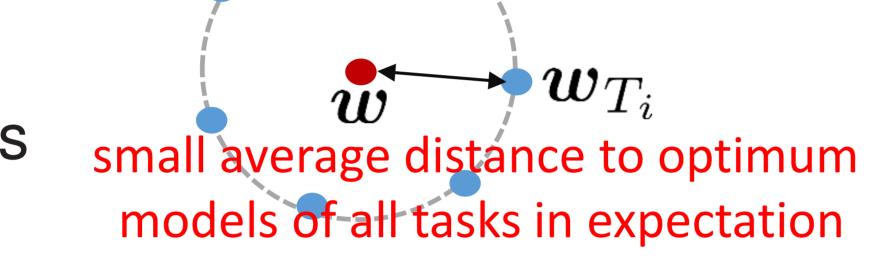
Training model: given a meta task distribution \mathcal{T} , we consider the meta-learning model:

$$\min_{\mathbf{w}} F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^{n} \min_{\mathbf{w}_{T_i}} \left\{ \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} ||\mathbf{w}_{T_i} - \mathbf{w}||_2^2 \right\}.$$
 (

where each task $T \sim \mathcal{T}$ only contains K training samples $D_{T_i} = \{(\boldsymbol{x}_i, \boldsymbol{y}_i)\}; \mathcal{L}_{D_T}(\boldsymbol{w}_T) = \frac{1}{K} \sum_{(\boldsymbol{x}, \boldsymbol{y}) \in D_T} \ell(f(\boldsymbol{w}, \boldsymbol{x}), \boldsymbol{y})$ denotes the empirical loss on D_T with predictor f and loss ℓ .

Inner level of intra-task learning: it finds the task-specific optimal model parameter \mathbf{w}_T of task T around the prior \mathbf{w} .

Outer level of inter-task learning: it uses optimal parameters w_T to tune w such that w is close to all w_T in average.



Test model: given a task $T \sim T$ with K samples $D_T = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^K$, we adapt the prior model \mathbf{w} to task T by minimizing

$$\min_{\boldsymbol{w}_T} \mathcal{L}_{D_T}(\boldsymbol{w}_T) + \frac{\lambda}{2} \|\boldsymbol{w}_T - \boldsymbol{w}\|_2^2,$$

Benefit: it requires a few data for adaptation, since prior model w is close to optimal model of T in expectation as existing and new tasks are from the same distribution and are similar.

Advantages over MAML

Higher computational efficiency. Meta-MinibatchProx only computes the gradient (see algorithm below) and is more efficient than MAML requiring Hessian computation.

More accurate solution to inner problem. MAML finds a good initialization w via

$$\mathbf{w}_{T}^{*} = \mathbf{w} - \eta \nabla \mathcal{L}_{D_{T}}(\mathbf{w}) = \underset{\mathbf{w}_{T}}{\operatorname{argmin}} \langle \nabla \mathcal{L}_{D_{T}}(\mathbf{w}), \mathbf{w}_{T} - \mathbf{w} \rangle + \frac{1}{2\eta} \|\mathbf{w}_{T} - \mathbf{w}\|_{2}^{2}.$$

Meta-MinibatchProx finds the task-specific optimal model by minimizing

$$\mathbf{w}_{T}^{*} = \min_{\mathbf{w}_{T}} \mathcal{L}_{D_{T}}(\mathbf{w}_{T}) + \frac{\lambda}{2} \|\mathbf{w}_{T} - \mathbf{w}\|_{2}^{2} = \min_{\mathbf{w}_{T}} \langle \nabla \mathcal{L}_{D_{T}}(\mathbf{w}), \mathbf{w}_{T} - \mathbf{w} \rangle + \frac{\lambda}{2} \|\mathbf{w}_{T} - \mathbf{w}\|_{2}^{2} + \frac{1}{2} \langle \nabla^{2} \mathcal{L}_{D_{T}}(\mathbf{w})(\mathbf{w}_{T} - \mathbf{w}), (\mathbf{w}_{T} - \mathbf{w}) \rangle + \frac{1}{6} \langle \nabla^{3} \mathcal{L}_{D_{T}}(\mathbf{w}), (\mathbf{w}_{T} - \mathbf{w})^{\otimes^{3}} \rangle + \cdots$$

So Meta-MinibatchProx uses higher-order information of $\mathcal{L}_{D_{\tau}}$ beyond gradient to search.

More flexible regularization. Meta-MinibatchProx can handle different possible structures in model parameter space by using proper $\ell_{p,q}$ -regularizer $\|\mathbf{w}_{T_i} - \mathbf{w}\|_p^q$, while MAML cannot due to its fixed update rule. E.g., Meta-MinibatchProx could use $\ell_{2,1}$ -norm to handle outlier-tasks.

Stochastic Gradient Meta-Optimization

$$\min_{\boldsymbol{w}} \frac{1}{n} \sum_{i=1}^{n} \phi_{D_{T}}(\boldsymbol{w}) \quad \text{where} \quad \phi_{D_{T}}(\boldsymbol{w}) = \min_{\boldsymbol{w}_{T}} \mathcal{L}_{D_{T}}(\boldsymbol{w}_{T}) + \frac{\lambda}{2} \|\boldsymbol{w}_{T} - \boldsymbol{w}\|_{2}^{2}. \tag{3}$$

The difficulty in applying SGD is how to compute the gradient of meta loss $\phi_{D_T}(\mathbf{w})$ w.r.t. \mathbf{w} .

Stochastic gradient computation of $\phi_{D_T}(\mathbf{w})$. Assume \mathcal{L}_{D_T} is differentiable and $\mathbf{w}_T^* = \operatorname{argmin}_{\mathbf{w}_T} \mathcal{L}_{D_T}(\mathbf{w}_T) + \frac{\lambda}{2} ||\mathbf{w}_T - \mathbf{w}^*||_2^2$. Then the gradient of the meta-loss $\phi_{D_T}(\mathbf{w})$ is given by $\nabla \phi_{D_T}(\mathbf{w}) = \lambda(\mathbf{w} - \mathbf{w}_T^*)$.

Stochastic Gradient Meta-Optimization

In practice, we only approximately solve the inner problem: $\mathbf{w}_T^* \approx \operatorname{argmin}_{\mathbf{w}_T} \mathcal{L}_{D_T}(\mathbf{w}_T)$.

Algorithm 1: SGD for Meta-MinibatchProx

Input: Initial point \mathbf{w}_0 , learning rate $\{\eta_s\}$.

for s = 0 to S - 1 do

Randomly select a mini-batch of tasks $\{T_i\}$ of size b_s from the observed n tasks. **for** $T_i \in \{T_i\}$

Compute an ϵ_s -approximate stable minimizer $\mathbf{w}_{T_i}^s$ to the within-meta-task problem $\min_{\mathbf{w}_{T_i}} g(\mathbf{w}_{T_i}) := \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} \|\mathbf{w}_{T_i} - \mathbf{w}^s\|_2^2$ such that $\|\nabla g(\mathbf{w}_{T_i}^s)\|_2^2 \le \epsilon_s$. **end**

Update the meta parameter $\mathbf{w}^{s+1} = \mathbf{w}^s - \eta_s \lambda (\mathbf{w}^s - \frac{1}{b_s} \sum_{i=1}^{b_s} \mathbf{w}_{T_i}^s)$.

end

Output: the parameter initialization \mathbf{w}^{S} of model f.

Convergence guarantees. Assume $\mathcal{L}_{D_T}(\mathbf{w}_T)$ is differentiable and for each task, its optimum $\mathbf{w}_T^* = \operatorname{argmin}_{\mathbf{w}_T} \mathcal{L}_{D_T}(\mathbf{w}_T) + \frac{\lambda}{2} ||\mathbf{w}_T - \mathbf{w}||_2^2$ satisfies $\mathbb{E}[||\mathbf{w}_T^* - \mathbf{w}||_2^2] \leq \sigma^2$. Let $\mathbf{w}_T^* = \operatorname{argmin}_{\mathbf{w}} F(\mathbf{w})$ and $\mathbf{w}_{T_i}^{*S} = \operatorname{argmin}_{\mathbf{w}_{T_i}} \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} ||\mathbf{w}_{T_i} - \mathbf{w}^S||_2^2$.

(1) Convex setting. Assume $\mathcal{L}_{D_T}(\mathbf{w}_T)$ is convex. Then by setting $\eta_s = \frac{2}{s\lambda}$, $\epsilon_s = \frac{c}{S}$, $\alpha = \frac{8S\lambda^2\sigma^2}{S-1} + c(1 + \frac{8}{S-1})$ with a constant c, we have

$$\mathbb{E}[\|\boldsymbol{w}^{S} - \boldsymbol{w}^{*}\|_{2}^{2}] \leq \frac{\alpha}{\lambda^{2}S} \quad \text{and} \quad \mathbb{E}\left[\left\|\frac{1}{n}\sum_{i=1}^{n}\boldsymbol{w}_{T_{i}}^{*S} - \boldsymbol{w}^{S}\right\|_{2}^{2}\right] \leq \frac{L^{2}\alpha}{(\lambda + L)^{2}S}.$$
(2) Non-convex setting. Assume $\mathcal{L}_{D_{T}}(\boldsymbol{w}_{T})$ is L -smooth. Then by setting $\lambda > L$,

(2) Non-convex setting. Assume $\mathcal{L}_{D_T}(\mathbf{w}_T)$ is L-smooth. Then by setting $\lambda > \eta_s = \sqrt{\frac{\Delta}{\gamma S}}$, $\epsilon_s = \frac{c}{\sqrt{S}}$ with $\gamma = \frac{\lambda^3 L}{(\lambda + L)} \left(\sigma^2 + \frac{c}{\sqrt{S}(\lambda - L)^2}\right)$ and $\Delta = F(\mathbf{w}^0) - F(\mathbf{w}^*)$, we have

 $\min_{\mathcal{S}} \mathbb{E}[\|\nabla F(\mathbf{w}^{\mathcal{S}})\|_{2}^{2}] = \lambda^{2} \min_{\mathcal{S}} \mathbb{E}\Big[\Big\|\frac{1}{n}\sum_{i=1}^{n} \mathbf{w}_{T_{i}}^{*S} - \mathbf{w}^{\mathcal{S}}\Big\|_{2}^{2}\Big] \leq \frac{1}{\sqrt{S}}\Big[4\sqrt{\Delta\gamma} + \frac{2c\lambda^{2}}{(\lambda - L)^{2}}\Big].$

Statistical Justification: Benefit of Prior Model in Meta Learning

Assume we have learned an optimal prior hypothesis $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} F(\mathbf{w})$. Here we view \mathbf{w}^* as a deterministic hypothesis.

For any $T \sim \mathcal{T}$ and $D_T = \{(\boldsymbol{x}_i, \boldsymbol{y}_i)\}_{i=1}^K \sim T$, define the population solution as $\boldsymbol{w}_{T,E}^* \in \operatorname{argmin} \big\{ \mathcal{L}(\boldsymbol{w}_T) := \mathbb{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim T} \big[\ell(f(\boldsymbol{w}_T,\boldsymbol{x}),\boldsymbol{y}) \big] \big\}$

and empirical solution in Meta-MinibatchProx as

$$\mathbf{w}_T^* = \underset{\mathbf{w}_T}{\operatorname{argmin}} \mathcal{L}_{D_T}(\mathbf{w}_T) + \frac{\lambda}{2} \|\mathbf{w}_T - \mathbf{w}^*\|_2^2.$$

Benefit of prior w^* **for new task learning**. Suppose $\ell(f(w, x), y)$ is G-Lipschitz continuous and convex w.r.t. w. Then we have

 $\mathbb{E}_{T \sim T} \mathbb{E}_{D_{T} \sim T} \left[\mathcal{L}(\mathbf{w}_{T}^{*}) - \mathcal{L}(\mathbf{w}_{T,E}^{*}) \right] \leq \frac{\sqrt{2}G}{\sqrt{K}} \sqrt{\mathbb{E}_{T \sim T} \left[\|\mathbf{w}^{*} - \mathbf{w}_{T,E}^{*}\|^{2} \right]}.$

Remark. Training model guarantees \mathbf{w}^* to be close to \mathbf{w}_{TF}^* in expectation.

First-order optimality. Suppose $\ell(f(\boldsymbol{w}, \boldsymbol{x}), \boldsymbol{y})$ is *G*-Lipschitz continuous and *L*-smooth w.r.t. \boldsymbol{w} . Then for $\lambda > L$, it holds that

$$\mathbb{E}_{T \sim \mathcal{T}} \Big[\Big\| \mathbb{E}_{D_T \sim T} \left[\nabla \mathcal{L}(\mathbf{w}_T^*) \right] \Big\|^2 \Big] \leq \frac{32G^2L^2}{(\lambda - L)^2K^2} + \frac{8G^2}{(\lambda - L)\beta K} + \frac{2}{\beta} \mathbb{E}_{T \sim \mathcal{T}} \left[\mathcal{L}(\mathbf{w}^*) - \mathcal{L}(\mathbf{w}_T^*) \right],$$
 where $\beta = \frac{1}{\lambda} \Big[1 - \frac{L}{2\lambda} \Big]$ is a constant.

Remark. w* helps obtain small gradient.

Comparison with Existing Works

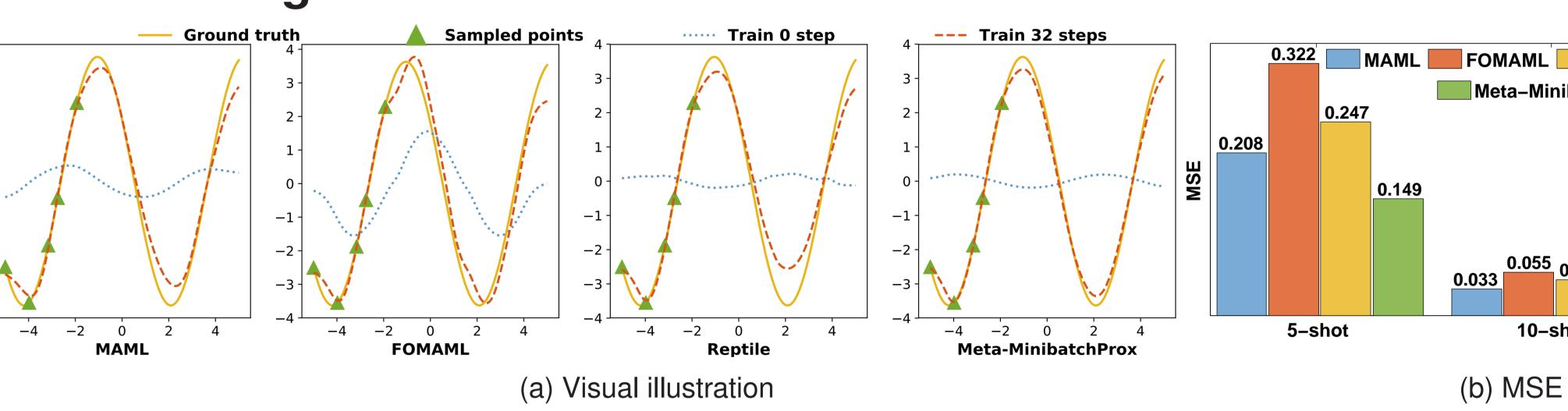
Recently, Denevi et al. (2018, 2019) and Khodak et al. (2019) also consider similar prior hypothesis biased regularized empirical risk minimization.

But Denevi et al. focus on convex linear models in contrast ours which is developed for both convex and non-convex learning problems.

Different from online convex meta-learning framework developed in (Khodak et al. 2019), we use a simple yet scalable paradigm within SGD for stochastic meta-optimization which is particularly friendly for computational and statistical complexity analysis in both convex and non-convex settings.

Experiments (code is available at https://panzhous.github.io/)

Few-shot regression.



Few-shot classification.

Table 1: Few-shot classification accuracy (%) on minilmageNet.

method	1-shot 5-way	5-shot 5-way	1-shot 20-way	5-shot 20-way
Matching Net	43.56 ± 0.84	55.31 ± 0.73	17.31 ± 0.22	22.69 ± 0.20
Meta-LSTM	43.33 ± 0.77	60.60 ± 0.71	16.70 ± 0.23	26.06 ± 0.25
MAML	46.21 ± 1.76	61.12 ± 1.01	16.01 ± 0.52	18.34 ± 0.33
FOMAML	45.53 ± 1.58	61.02 ± 1.12	15.21 ± 0.54	17.67 ± 0.47
Reptile	47.07 ± 0.26	62.74 ± 0.37	18.27 ± 0.16	28.71 ± 0.19
Meta-MinibatchProx	$\textbf{48.51} \pm \textbf{0.92}$	$\textbf{64.15} \pm \textbf{0.92}$	$\textbf{20.50} \pm \textbf{0.35}$	$\textbf{33.61} \pm \textbf{0.41}$
MAML + Transduction	48.70 ± 1.84	63.11 ± 0.92	16.49 ± 0.58	19.29 ± 0.29
FOMAML + Transduction	48.07 ± 1.75	63.15 ± 0.91	15.80 ± 0.61	18.15 ± 0.43
Reptile + Transduction	49.97 ± 0.32	65.99 ± 0.58	18.76 ± 0.17	29.15 ± 0.22
Meta-MinibatchProx + Transduction	$\textbf{50.77} \pm \textbf{0.90}$	$\textbf{67.43} \pm \textbf{0.89}$	$\textbf{21.17} \pm \textbf{0.38}$	$\textbf{34.30} \pm \textbf{0.41}$

Table 2: Few-shot classification accuracy (%) on tieredImageNet.

method	1-shot 5-way	5-shot 5-way	1-shot 10-way	5-shot 10-way
Matching Net	34.95 ± 0.89	43.95 ± 0.85	22.46 ± 0.34	31.19 ± 0.30
Meta-LSTM	33.71 ± 0.76	46.56 ± 0.79	22.09 ± 0.43	35.65 ± 0.39
MAML	49.60 ± 1.83	66.58 ± 1.78	33.18 ± 1.23	49.05 ± 1.32
FOMAML	48.01 ± 1.74	64.07 ± 1.72	30.31 ± 1.12	46.54 ± 1.24
Reptile	49.12 ± 0.43	65.99 ± 0.42	31.79 ± 0.28	47.82 ± 0.30
Meta-MinibatchProx	$\textbf{50.14} \pm \textbf{0.92}$	$\textbf{68.30} \pm \textbf{0.91}$	$\textbf{33.68} \pm \textbf{0.64}$	$\textbf{51.84} \pm \textbf{0.65}$
MAML + Transduction	51.67 ± 1.81	70.30 ± 1.75	34.44 ± 1.19	53.32 ± 1.33
FOMAML + Transduction	50.12 ± 1.82	67.43 ± 1.80	31.53 ± 1.08	49.99 ± 1.36
Reptile + Transduction	51.06 ± 0.45	69.94 ± 0.42	33.79 ± 0.29	51.27 ± 0.31
Meta-MinibatchProx + Transduction	$\textbf{54.37} \pm \textbf{0.93}$	$\textbf{71.45} \pm \textbf{0.94}$	$\textbf{35.56} \pm \textbf{0.60}$	$\textbf{54.50} \pm \textbf{0.71}$

Outlier-Corrupted Tasks.

Assume there are a few outlier-tasks $\mathcal{O} = \{T_o\}$ whose optima $\{\mathbf{w}_o\}$ are quite far from the optima $\{\mathbf{w}_s\}$ of inlier/normal-tasks $\mathcal{S} = \{T_s\}$.

Here we add 5% outlier images with zero pixels into each training class in minilmageNet. If a sampled task T contains these outlier images, it is an outlier-task.

Meta-MinibatchProx uses the robust $\ell_{2,1}$ -norm $\frac{1}{n}\sum_{i=1}^{n}\|\mathbf{w}_{T_i}-\mathbf{w}\|_2$ that tolerates relatively large distances between \mathbf{w} and $\{\mathbf{w}_o\}$.

