

Efficient Meta Learning via Minibatch Proximal Update

Pan Zhou

Joint work with Xiao-Tong Yuan, Huan Xu, Shuicheng Yan, Jiashi Feng

National University of Singapore
pzhou@u.nus.edu

Dec 11, 2019

Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx learns a good **prior model initialization** w from observed tasks such that
 w is close to the optimal models of new similar tasks, promoting new task learning

Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx learns a good **prior model initialization** \mathbf{w} from observed tasks such that
 \mathbf{w} is close to the optimal models of new similar tasks, promoting new task learning

- **Training model:** given a task distribution \mathcal{T} , we minimize a **bi-level** meta learning model

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

where $T_i \sim \mathcal{T}$ has K training samples $D_{T_i} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^K$

$\mathcal{L}_{D_{T_i}} = \frac{1}{K} \sum_{(\mathbf{x}, \mathbf{y}) \in D_{T_i}} \ell(f(\mathbf{w}, \mathbf{x}), \mathbf{y})$ is empirical loss with predictor f and loss ℓ .

Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx learns a good **prior model initialization** \mathbf{w} from observed tasks such that
 \mathbf{w} is close to the optimal models of new similar tasks, promoting new task learning

- **Training model:** given a task distribution \mathcal{T} , we minimize a **bi-level** meta learning model

update task-specific solution

$$\min_{\mathbf{w}} \min_{\mathbf{w}_{T_i}} \sum_{i=1}^n \boxed{\mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}_{T_i}\|_2^2},$$

where $T_i \sim \mathcal{T}$ has K training samples $D_{T_i} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^K$

$\mathcal{L}_{D_{T_i}} = \frac{1}{K} \sum_{(\mathbf{x}, \mathbf{y}) \in D_{T_i}} \ell(f(\mathbf{w}, \mathbf{x}), \mathbf{y})$ is empirical loss with predictor f and loss ℓ .

Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx learns a good **prior model initialization** \mathbf{w} from observed tasks such that
 \mathbf{w} is close to the optimal models of new similar tasks, promoting new task learning

- **Training model:** given a task distribution \mathcal{T} , we minimize a **bi-level** meta learning model
update the prior model
$$\min_{\mathbf{w}} \left[\min_{\mathbf{w}_{T_i}} \sum_{i=1}^n \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}_{T_i}\|_2^2 \right],$$

where $T_i \sim \mathcal{T}$ has K training samples $D_{T_i} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^K$

$\mathcal{L}_{D_{T_i}} = \frac{1}{K} \sum_{(\mathbf{x}, \mathbf{y}) \in D_{T_i}} \ell(f(\mathbf{w}, \mathbf{x}), \mathbf{y})$ is empirical loss with predictor f and loss ℓ .

Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx learns a good **prior model initialization** \mathbf{w} from observed tasks such that
 \mathbf{w} is close to the optimal models of new similar tasks, promoting new task learning

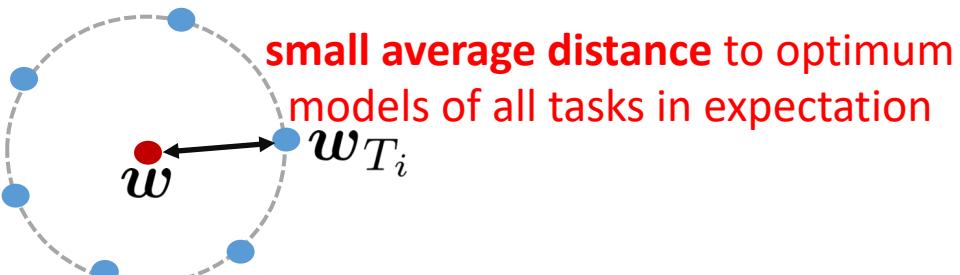
- **Training model:** given a task distribution \mathcal{T} , we minimize a **bi-level** meta learning model

update the prior model

$$\min_{\mathbf{w}} \boxed{\min_{\mathbf{w}_{T_i}} \sum_{i=1}^n \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}_{T_i}\|_2^2},$$

where $T_i \sim \mathcal{T}$ has K training samples $D_{T_i} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^K$

$\mathcal{L}_{D_{T_i}} = \frac{1}{K} \sum_{(\mathbf{x}, \mathbf{y}) \in D_{T_i}} \ell(f(\mathbf{w}, \mathbf{x}), \mathbf{y})$ is empirical loss with predictor f and loss ℓ .



Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx learns a good **prior model initialization** \mathbf{w} from observed tasks such that
 \mathbf{w} is close to the optimal models of new similar tasks, promoting new task learning

- **Test model:** given a randomly sample a task $T \sim \mathcal{T}$ consisting of K samples $D_T = \{(x_i, y_i)\}_{i=1}^K$

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

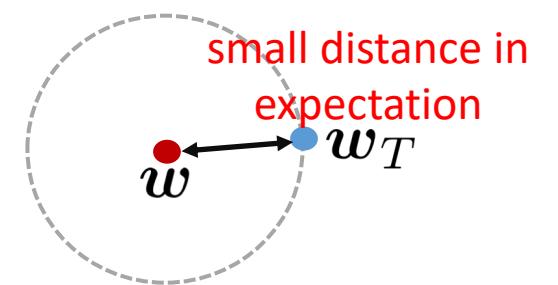
Meta Learning via Minibatch Proximal Update (Meta-MinibatchProx)

Meta-MinibatchProx learns a good **prior model initialization** \mathbf{w} from observed tasks such that
 \mathbf{w} is close to the optimal models of new similar tasks, promoting new task learning

- **Test model:** given a randomly sample a task $T \sim \mathcal{T}$ consisting of K samples $D_T = \{(x_i, y_i)\}_{i=1}^K$

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

- **Benefit:** a few data is sufficient for adaptation
 - prior model \mathbf{w} is close to optimum \mathbf{w}_T**
when training and test tasks are from a same distribution \mathcal{T} .



Optimization Algorithm

We use SGD based algorithm to solve bi-level training model :

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

Optimization Algorithm

We use SGD based algorithm to solve bi-level training model :

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

- Step1. select a mini-batch of task $\{T_i\}$.

Optimization Algorithm

We use SGD based algorithm to solve bi-level training model :

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

- Step1. select a mini-batch of task $\{T_i\}$.
- Step2. for T_i , compute an approximate minimizer:

$$\mathbf{w}_{T_i} \approx \operatorname{argmin}_{\mathbf{w}_{T_i}} \{g(\mathbf{w}_{T_i}) =: \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}_{T_i}\|_2^2\} \text{ s.t. } \|\nabla g(\mathbf{w}_{T_i})\|_2^2 \leq \epsilon_s$$

Optimization Algorithm

We use SGD based algorithm to solve bi-level training model :

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

- Step1. select a mini-batch of task $\{T_i\}$.
- Step2. for T_i , compute an approximate minimizer:
$$\mathbf{w}_{T_i} \approx \operatorname{argmin}_{\mathbf{w}_{T_i}} \{g(\mathbf{w}_{T_i}) =: \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}_{T_i}\|_2^2\} \text{ s.t. } \|\nabla g(\mathbf{w}_{T_i})\|_2^2 \leq \epsilon_s$$
- Step3. update the prior model

$$\mathbf{w} = \mathbf{w} - \eta_s \lambda \left(\mathbf{w} - \frac{1}{b_s} \sum_{i=1}^{b_s} \mathbf{w}_{T_i} \right)$$

Optimization Algorithm

We use SGD based algorithm to solve bi-level training model :

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

- Step1. select a mini-batch of task $\{T_i\}$.
- Step2. for T_i , compute an approximate minimizer:

$$\mathbf{w}_{T_i} \approx \operatorname{argmin}_{\mathbf{w}_{T_i}} \{g(\mathbf{w}_{T_i}) =: \mathcal{L}_{D_{T_i}}(\mathbf{w}_{T_i}) + \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}_{T_i}\|_2^2\} \text{ s.t. } \|\nabla g(\mathbf{w}_{T_i})\|_2^2 \leq \epsilon_s$$

- Step3. update the prior model

$$\mathbf{w} = \mathbf{w} - \eta_s \lambda \left(\mathbf{w} - \frac{1}{b_s} \sum_{i=1}^{b_s} \mathbf{w}_{T_i} \right)$$

Theorem 1 (convergence guarantees, informal).

- (1) Convex setting, i.e. convex $\phi_{D_{T_i}}(\mathbf{w})$. We prove $\mathbb{E}[\|\mathbf{w}^S - \mathbf{w}^*\|_2^2] \leq \mathcal{O}\left(\frac{1}{S}\right)$.
- (2) Nonconvex setting, i.e. smooth $\phi_{D_{T_i}}(\mathbf{w})$. We prove $\mathbb{E}_s[\|\nabla F(\mathbf{w}^s)\|_2^2] \leq \mathcal{O}\left(\frac{1}{\sqrt{S}}\right)$.

Generalization Performance Guarantee

- Ideally, for a given task T, one should train the model on the population risk
Population solution: $\mathbf{w}_P^* = \operatorname{argmin}_{\mathbf{w}_T} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim T} \ell(f(\mathbf{w}_T, \mathbf{x}), \mathbf{y}).$
- In practice, we only has K samples and adapt the prior model \mathbf{w}^* to the new task:
Empirical solution: $\mathbf{w}_T^* = \operatorname{argmin}_{\mathbf{w}_T} \mathcal{L}_{D_T}(\mathbf{w}_T) + \frac{\lambda}{2} \|\mathbf{w}^* - \mathbf{w}_T\|_2^2.$

Generalization Performance Guarantee

- Ideally, for a given task T, one should train the model on the population risk
Population solution: $\mathbf{w}_P^* = \operatorname{argmin}_{\mathbf{w}_T} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim T} \ell(f(\mathbf{w}_T, \mathbf{x}), \mathbf{y}).$
- In practice, we only have K samples and adapt the prior model \mathbf{w}^* to the new task:
Empirical solution: $\mathbf{w}_T^* = \operatorname{argmin}_{\mathbf{w}_T} \mathcal{L}_{D_T}(\mathbf{w}_T) + \frac{\lambda}{2} \|\mathbf{w}^* - \mathbf{w}_T\|_2^2.$
- **Since $\mathbf{w}_P^* \neq \mathbf{w}_T^*$, why \mathbf{w}_T^* is good for generalization in few-shot learning problem?**

Theorem 2 (generalization performance guarantee, informal).

Suppose each loss $\phi_{D_{T_i}}(\mathbf{w})$ is convex and is smooth. Let $D_T = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^K \sim T$. Then we have

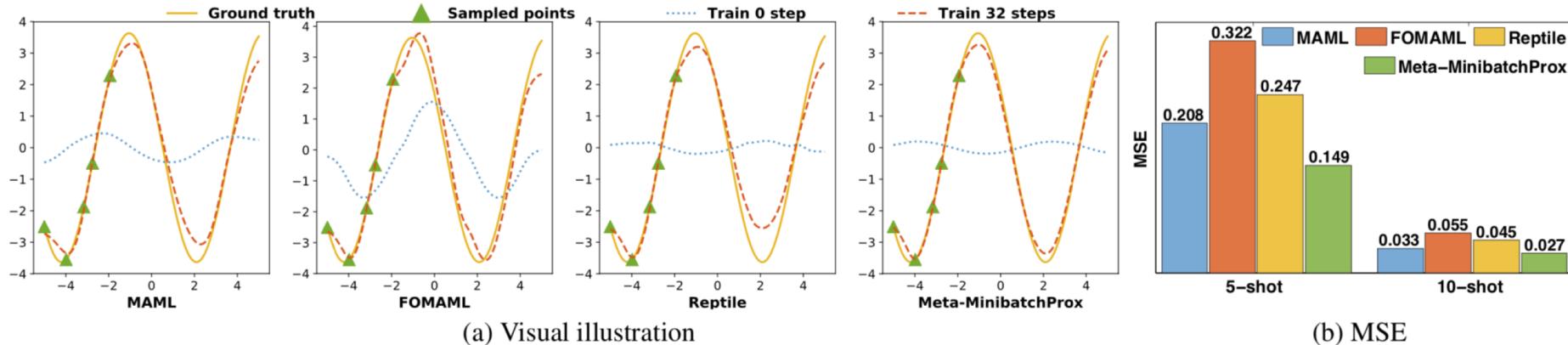
$$\mathbb{E}_{T \sim \mathcal{T}} \mathbb{E}_{D_T \sim T} (\mathcal{L}(\mathbf{w}_T^*) - \mathcal{L}(\mathbf{w}_P^*)) \leq \frac{c}{\sqrt{K}} \mathbb{E}[\|\mathbf{w}^* - \mathbf{w}_P^*\|_2^2].$$

Remark: strong generalization performance, as our training model guarantee

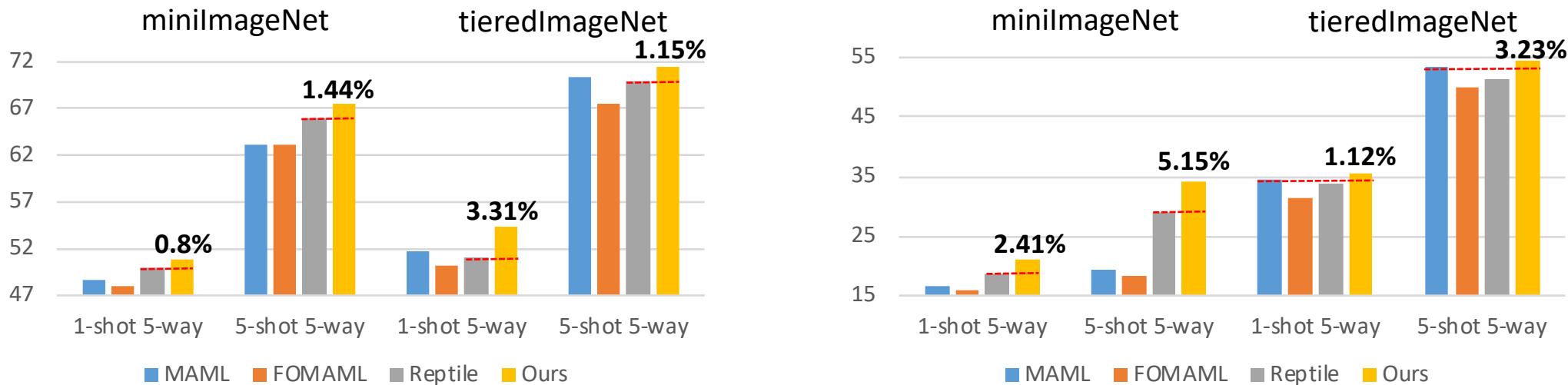
prior \mathbf{w}^* is close to the optimum model \mathbf{w}_P^* .

Experimental results

Few-shot regression : smaller mean square error (MSE) between prediction and ground truth



Few-shot classification: higher classification accuracy



POSTER # 26

05:00 -- 07:00 PM @ East Exhibition Hall B + C

Thanks!