# 1 Proof

**Proposition 1.** *Let $\eta_t = 1/t$. Assume $\gamma < 1/\sigma_{\mathbf{x}}^2$. Both $\psi_t$ in MAML (with one inner gradient step) and $\boldsymbol{\theta}_t$ in CommonMean converge to $\bar{\mathbf{w}} = \mathbb{E}_\tau \mathbf{w}_\tau^*$.*

**Proposition 2.** *Assume that $\gamma < 1/\sigma_{\mathbf{x}}^2$. We have $\bar{\mathbf{w}} = \operatorname{argmin}_{\boldsymbol{\theta}} \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top \mathbf{w}_\tau^{(prox)} - y)^2 = \operatorname{argmin}_{\boldsymbol{\psi}} \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top \mathbf{w}_\tau^{(gd)} - y)^2.$*

We first prove Proposition 2 that the mean regressor is the unique minimizer. Then, we prove Proposition 1 by showing that MAML (with one inner gradient step) and CommonMean algorithms achieve global convergence.

## 1.1 Proof of Proposition 2

*Proof.* For each task $\tau$, let $\mathbf{v}_\tau = \mathbf{w}_\tau^* - \bar{\mathbf{w}}$, then $\{\mathbf{v}_\tau\}$ are i.i.d. random variables with zero mean. Denote $\mathbf{C}_\tau = (\lambda \mathbf{I} + \mathbf{X}_\tau^\top \mathbf{X}_\tau)^{-1}$. As $\mathbf{w}_\tau^{(prox)} = \mathbf{C}_\tau (\lambda\boldsymbol{\theta} + \mathbf{X}_\tau^\top \mathbf{y}_\tau)$ and $\mathbf{y}_\tau = \mathbf{X}_\tau \mathbf{w}_\tau^* + \boldsymbol{\xi}_\tau$, it follows that

$$\mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top \mathbf{w}_\tau^{(prox)} - y)^2$$

$$= \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\lambda\mathbf{x}^\top \mathbf{C}_\tau \boldsymbol{\theta} + \mathbf{x}^\top \mathbf{C}_\tau \mathbf{X}_\tau^\top (\mathbf{X}_\tau \mathbf{w}_\tau^* + \boldsymbol{\xi}_\tau) - \mathbf{x}^\top \mathbf{w}_\tau^* - \xi)^2$$

$$= \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\lambda\mathbf{x}^\top \mathbf{C}_\tau \boldsymbol{\theta} + \mathbf{x}^\top \mathbf{C}_\tau \mathbf{X}_\tau^\top (\mathbf{X}_\tau \bar{\mathbf{w}} + \mathbf{X}_\tau \mathbf{v}_\tau + \boldsymbol{\xi}_\tau) - \mathbf{x}^\top \bar{\mathbf{w}} - \mathbf{x}^\top \mathbf{v}_\tau - \xi)^2$$

$$= \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\lambda\mathbf{x}^\top \mathbf{C}_\tau \boldsymbol{\theta} + \mathbf{x}^\top \mathbf{C}_\tau \mathbf{X}_\tau^\top \mathbf{X}_\tau \bar{\mathbf{w}} - \mathbf{x}^\top \bar{\mathbf{w}})^2 + \text{constant} \tag{1}$$

$$= \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\lambda\mathbf{x}^\top \mathbf{C}_\tau (\boldsymbol{\theta} - \bar{\mathbf{w}}))^2 + \text{constant}$$

$$= \lambda^2 \sigma_{\mathbf{x}}^2 n_q \mathbb{E}_\tau \mathbb{E}_{S_\tau} (\boldsymbol{\theta} - \bar{\mathbf{w}})^\top \mathbf{C}_\tau^2 (\boldsymbol{\theta} - \bar{\mathbf{w}}) + \text{constant},$$

where we have used the setting that $\mathbf{x}, \xi, \mathbf{X}_\tau, \boldsymbol{\xi}_\tau$, and $\mathbf{v}_\tau$ are independent to obtain (1). Since $\mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbf{C}_\tau^2 \succeq \lambda^{-2}\mathbf{I}$, we conclude that $\boldsymbol{\theta} = \bar{\mathbf{w}}$ is the unique optima.

For MAML with one gradient step $\mathbf{w}_\tau^{(gd)} = \boldsymbol{\psi} - \gamma\mathbf{X}_\tau^\top (\mathbf{X}_\tau \boldsymbol{\psi} - \mathbf{y}_\tau)$, it follows that

$$\mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top \mathbf{w}_\tau^{(gd)} - y)^2$$

$$= \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top (\mathbf{I} - \gamma\mathbf{X}_\tau^\top \mathbf{X}_\tau)\boldsymbol{\psi} + \gamma\mathbf{x}^\top \mathbf{X}_\tau^\top \mathbf{y}_\tau - y)^2$$

$$= \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top (\mathbf{I} - \gamma\mathbf{X}_\tau^\top \mathbf{X}_\tau)\boldsymbol{\psi} + \gamma\mathbf{x}^\top \mathbf{X}_\tau^\top (\mathbf{X}_\tau \bar{\mathbf{w}} + \mathbf{X}_\tau \mathbf{v}_\tau + \boldsymbol{\xi}_\tau) - \mathbf{x}^\top \bar{\mathbf{w}} - \mathbf{x}^\top \mathbf{v}_\tau - \xi)^2$$

$$= \mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top (\mathbf{I} - \gamma\mathbf{X}_\tau^\top \mathbf{X}_\tau)(\boldsymbol{\psi} - \bar{\mathbf{w}}))^2 + \text{constant}$$

$$= n_q \sigma_{\mathbf{x}}^2 \mathbb{E}_\tau \mathbb{E}_{S_\tau} \|(\mathbf{I} - \gamma\mathbf{X}_\tau^\top \mathbf{X}_\tau)(\boldsymbol{\psi} - \bar{\mathbf{w}})\|^2 + \text{constant}.$$

As $\gamma < 1/\sigma_{\mathbf{x}}^2$, we conclude that $\boldsymbol{\psi} = \bar{\mathbf{w}}$ is the unique optima. $\qquad\square$

## 1.2 Proof of Proposition 1

*Proof.* (i) Notice that $\mathbf{w}_\tau^{(prox)}$ is affine in $\boldsymbol{\theta}$, thus, $\mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top \mathbf{w}_\tau^{(prox)} - y)^2$ is convex in $\boldsymbol{\theta}$. The CommonMean algorithm is using stochastic gradient descent to minimize the population risk, and the global convergence of $\boldsymbol{\theta}_t$ follows from the stochastic convex optimization [1].

(ii) Similarly, $\mathbf{w}_\tau^{(gd)}$ is affine in $\boldsymbol{\psi}$, thus, the loss $\mathbb{E}_\tau \mathbb{E}_{S_\tau} \mathbb{E}_{Q_\tau} \sum_{(\mathbf{x},y)\in Q_\tau} (\mathbf{x}^\top \mathbf{w}_\tau^{(gd)} - y)^2$ is convex in $\boldsymbol{\psi}$. Using stochastic gradient descent, $\boldsymbol{\psi}_t$ achieves global convergence [1]. By Proposition 2, $\bar{\mathbf{w}}$ is the unique optima, and we finish the proof. $\qquad\square$

## 1.3 Proof of Proposition 4

The task index $\tau'$ will be omitted for simplifying notations in Proposition 4.

**Proposition 4.** $\mathbb{E}_{\boldsymbol{\xi}} \|\mathbf{w}^{(prox)} - \mathbf{w}^*\|^2 = \|\tilde{\mathbf{b}}\|^2 + \sum_{j=1}^{n_s} \left( \frac{\lambda \tilde{\mathrm{a}}_j}{\lambda + \nu_j^2} \right)^2 + \sum_{j=1}^{n_s} \left( \frac{\sigma_{\xi}}{(\lambda/\nu_j) + \nu_j} \right)^2$, where the expectation is over the label noise vector $\boldsymbol{\xi}$.

*Proof.* The ridge regression has a closed-form solution $\mathbf{w}^{(prox)} = \left( \lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X} \right)^{-1} \left( \lambda \boldsymbol{\theta} + \mathbf{X}^\top \mathbf{y} \right)$. Using the SVD decomposition of $\mathbf{X} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$ and $\mathbf{y} = \mathbf{X} \mathbf{w}^* + \boldsymbol{\xi}$, we obtain

$$\mathbf{w}^{(prox)} = \left( \mathbf{I} + \lambda^{-1} \mathbf{V} \boldsymbol{\Sigma}^2 \mathbf{V}^\top \right)^{-1} \left( \mathbf{V} \mathbf{a}_0 + \mathbf{V}^\perp \mathbf{b}_0 + \lambda^{-1} \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^\top \mathbf{y} \right)$$

$$= \left( \mathbf{I} + \lambda^{-1} \mathbf{V} \boldsymbol{\Sigma}^2 \mathbf{V}^\top \right)^{-1} \left( \mathbf{V} \mathbf{a}_0 + \mathbf{V}^\perp \mathbf{b}_0 + \lambda^{-1} \mathbf{V} \boldsymbol{\Sigma}^2 \mathbf{a}^* + \lambda^{-1} \mathbf{V} \boldsymbol{\Sigma} \mathbf{U} \boldsymbol{\xi} \right) \quad (2)$$

$$= \mathbf{V}^\perp \mathbf{b}_0 + \mathbf{V} (\mathbf{I} + \lambda^{-1} \boldsymbol{\Sigma}^2)^{-1} \left( \mathbf{a}_0 + \lambda^{-1} \boldsymbol{\Sigma}^2 \mathbf{a}^* \right) + \mathbf{V} \left( \lambda \boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma} \right)^{-1} \mathbf{U}^\top \boldsymbol{\xi}, \quad (3)$$

where we have used $\mathbf{U}^\top \mathbf{y} = \mathbf{U}^\top (\mathbf{X} \mathbf{w}^* + \boldsymbol{\xi}) = \mathbf{U}^\top \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top (\mathbf{V} \mathbf{a}^* + \mathbf{V}^\perp \mathbf{b}^*) + \mathbf{U}^\top \boldsymbol{\xi} = \boldsymbol{\Sigma} \mathbf{a}^* + \mathbf{U}^\top \boldsymbol{\xi}$ in (2) and the Woodbury identity in (3). Then the estimation error is

$$\mathbf{w}^{(prox)} - \mathbf{w}^* = \mathbf{V}^\perp (\mathbf{b}_0 - \mathbf{b}^*) + \mathbf{V} (\mathbf{I} + \lambda^{-1} \boldsymbol{\Sigma}^2)^{-1} (\mathbf{a}_0 - \mathbf{a}^*) + \mathbf{V} \left( \lambda \boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma} \right)^{-1} \mathbf{U}^\top \boldsymbol{\xi}.$$

Taking the square $\ell_2$-norm and then expectation over $\boldsymbol{\xi}$ on both sides, we have

$$\mathbb{E}_{\boldsymbol{\xi}} \|\mathbf{w}^{(prox)} - \mathbf{w}^*\|^2$$

$$= \|\mathbf{V}^\perp (\mathbf{b}_0 - \mathbf{b}^*)\|^2 + \|\mathbf{V} (\mathbf{I} + \lambda^{-1} \boldsymbol{\Sigma}^2)^{-1} (\mathbf{a}_0 - \mathbf{a}^*)\|^2 + \mathbb{E}_{\boldsymbol{\xi}} \|\mathbf{V} \left( \lambda \boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma} \right)^{-1} \mathbf{U}^\top \boldsymbol{\xi}\|^2 \quad (4)$$

$$= \|\mathbf{b}_0 - \mathbf{b}^*\|^2 + \|(\mathbf{I} + \lambda^{-1} \boldsymbol{\Sigma}^2)^{-1} (\mathbf{a}_0 - \mathbf{a}^*)\|^2 + \mathbb{E}_{\boldsymbol{\xi}} \| \left( \lambda \boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma} \right)^{-1} \mathbf{U}^\top \boldsymbol{\xi}\|^2$$

$$= \|\tilde{\mathbf{b}}\|^2 + \sum_{j=1}^{n_s} \left( \frac{\lambda \tilde{\mathrm{a}}_j}{\lambda + \nu_j^2} \right)^2 + \sum_{j=1}^{n_s} \left( \frac{\nu_j \sigma_{\xi}}{\lambda + \nu_j^2} \right)^2,$$

where (4) follows from the fact that $\mathbf{V}^\perp$ is $\mathbf{V}$'s orthogonal complement and $\boldsymbol{\xi}$ is independent with $\mathbf{X}$ (also the $\boldsymbol{\Sigma}$, $\mathbf{U}$ and $\mathbf{V}$). $\qquad \square$

## 1.4 Proof of Theorem 1

**Lemma 1.** $\mathcal{L}_{meta}(\boldsymbol{\theta}, \boldsymbol{\phi})$ is Lipschitz-smooth w.r.t. $(\boldsymbol{\theta}, \boldsymbol{\phi})$ with a Lipschitz constant $\beta_{meta}$.

Lipschitz-smoothness is a basic assumption to establish convergence of gradient descent algorithms in stochastic non-convex optimization [4, 8] and meta-learning in non-convex settings [2, 11].

*Proof of Lemma 1.* As $\mathcal{L}_{meta}(\boldsymbol{\theta}, \boldsymbol{\phi}) \equiv \sum_{\tau \in \mathcal{T}} \sum_{(\mathbf{x}, y) \in Q_\tau} \ell(\hat{y}, y)$, it suffices to show that $\ell(\hat{y}, y)$ is Lipschitz-smooth in $(\boldsymbol{\theta}, \boldsymbol{\phi})$.

Using the chain rule, we have

$$\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \ell(\hat{y}, y) = \nabla_1 \ell(\hat{y}, y) \nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \hat{y}, \quad (5)$$

$$\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \hat{y} = \nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} f_{\boldsymbol{\theta}}(\mathbf{z}) + (\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \mathcal{K}(\mathbf{Z}_\tau, \mathbf{z}))^\top \boldsymbol{\alpha}_\tau + (\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \boldsymbol{\alpha}_\tau)^\top \mathcal{K}(\mathbf{Z}_\tau, \mathbf{z}). \quad (6)$$

The Lipschitz properties of direct derivatives $\nabla_1 \ell(\hat{y}, y)$, $\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} f_{\boldsymbol{\theta}}(\mathbf{z})$, $\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \mathcal{K}(\mathbf{Z}_\tau, \mathbf{z})$, and $\mathcal{K}(\mathbf{Z}_\tau \mathbf{z})$ follow from the Assumption 1. It remains to claim $\boldsymbol{\alpha}_\tau$ and $\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \boldsymbol{\alpha}_\tau$ are Lipschitz. Let $\mathbf{p} = [f_{\boldsymbol{\theta}}(\mathbf{z}_1); \ldots; f_{\boldsymbol{\theta}}(\mathbf{z}_{n_s}); \mathcal{K}(\mathbf{Z}_\tau, \mathbf{z}_1); \ldots; \mathcal{K}(\mathbf{Z}_\tau, \mathbf{z}_{n_s})] \in \mathbb{R}^{n_s + n_s^2}$ be the input of the dual problem.

(i) Claim: $\boldsymbol{\alpha}_\tau$ is Lipschitz w.r.t. $(\boldsymbol{\theta}, \boldsymbol{\phi})$ and $\boldsymbol{\alpha}_\tau(\mathbf{p})$ is Lipschitz-smooth w.r.t. $\mathbf{p}$. To show $\boldsymbol{\alpha}_\tau$ is Lipschitz w.r.t. $(\boldsymbol{\theta}, \boldsymbol{\phi})$, it suffices to show that $\|\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \boldsymbol{\alpha}_\tau\|$ is bounded. By the chain rule, $\nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \boldsymbol{\alpha}_\tau = \nabla_{\mathbf{p}} \boldsymbol{\alpha}_\tau \nabla_{(\boldsymbol{\theta}, \boldsymbol{\phi})} \mathbf{p}$. Denote the dual objective by $g(\mathbf{p}, \boldsymbol{\alpha})$. By the implicit function theorem [9], $\nabla_{\mathbf{p}} \boldsymbol{\alpha}_\tau = - \left( \nabla_{\boldsymbol{\alpha}}^2 g(\mathbf{p}, \boldsymbol{\alpha}_\tau) \right)^{-1} \frac{\partial^2}{\partial \mathbf{p} \partial \boldsymbol{\alpha}} g(\mathbf{p}, \boldsymbol{\alpha}_\tau)$, where $\nabla_{\boldsymbol{\alpha}}^2 g(\mathbf{p}, \boldsymbol{\alpha}_\tau) = \sum_{(\mathbf{x}_i, y_i) \in S_\tau} \nabla_1^2 \ell(f_\tau(\mathbf{z}_i), y_i) \mathcal{K}(\mathbf{Z}_\tau, \mathbf{z}_i) \mathcal{K}(\mathbf{Z}_\tau, \mathbf{z}_i)^\top + \mathcal{K}(\mathbf{Z}_\tau, \mathbf{Z}_\tau)$, $\frac{\partial^2}{\partial \mathbf{p} \partial \boldsymbol{\alpha}} g(\mathbf{p}, \boldsymbol{\alpha}_\tau) = [\mathcal{K}(\mathbf{Z}_\tau, \mathbf{Z}_\tau) \mathbf{D} \mid (\mathcal{K}(\mathbf{Z}_\tau, \mathbf{Z}_\tau) \mathbf{D}) \otimes \boldsymbol{\alpha}_\tau^\top + \mathbf{v}^\top \otimes \mathbf{I} + \mathbf{I} \otimes \boldsymbol{\alpha}_\tau^\top]$, $\mathbf{D} =$

$\mathrm{diag}([\nabla_1^2\ell(f_\tau(\mathbf{z}_1),y_1);\ldots;\nabla_1^2\ell(f_\tau(\mathbf{z}_{n_s}),y_{n_s})])$, $\mathbf{v} = [\nabla_1\ell(f_\tau(\mathbf{z}_1),y_1);\ldots;\nabla_1\ell(f_\tau(\mathbf{z}_{n_s}),y_{n_s})]$, where $\otimes$ is the Kronecker product. It follows from the Assumption 1 that both $\nabla_\alpha^2 g(\mathbf{p},\alpha_\tau)$ and $\frac{\partial^2}{\partial\mathbf{p}\partial\alpha}g(\mathbf{p},\alpha_\tau)$ are Lipschitz w.r.t. $\mathbf{p}$. Hence, we conclude that $\nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p})$ is Lipschitz, $\alpha_\tau(\mathbf{p})$ is Lipschitz-smooth w.r.t. $\mathbf{p}$, and $\|\nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p})\|$ is bounded. Again, the boundedness of $\nabla_{(\theta,\phi)}\mathbf{p}$ follows from the Lipschitz-smoothness of $\mathbf{p}$ w.r.t. $(\theta,\phi)$. We conclude that $\alpha_\tau$ is Lipschitz w.r.t. $(\theta,\phi)$.

(ii) Claim: $\nabla_{(\theta,\phi)}\alpha_\tau$ is Lipschitz w.r.t. $(\theta,\phi)$. Given $(\theta,\phi)$ and $(\theta',\phi')$, we show $\|\nabla_{(\theta,\phi)}\alpha_\tau(\theta,\phi) - \nabla_{(\theta,\phi)}\alpha_\tau(\theta',\phi')\| \leq \beta\|(\theta,\phi) - (\theta',\phi')\|$ for some $\beta > 0$. For notation simplicity, let $\varphi = (\theta,\phi)$ and $\varphi' = (\theta',\phi')$, then we have

$$\|\nabla_\varphi\alpha_\tau(\varphi) - \nabla_\varphi\alpha_\tau(\varphi')\|$$
$$=\|\nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi))\nabla_\varphi\mathbf{p}(\varphi) - \nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi'))\nabla_\varphi\mathbf{p}(\varphi')\|$$
$$=\|\nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi))\nabla_\varphi\mathbf{p}(\varphi) - \nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi'))\nabla_\varphi\mathbf{p}(\varphi') \pm \nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi))\nabla_\varphi\mathbf{p}(\varphi')\|$$
$$\leq\|\nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi))\|\|\nabla_\varphi\mathbf{p}(\varphi) - \nabla_\varphi\mathbf{p}(\varphi')\| + \|\nabla_\varphi\mathbf{p}(\varphi')\|\|\nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi)) - \nabla_{\mathbf{p}}\alpha_\tau(\mathbf{p}(\varphi'))\|.$$

As $\mathbf{p}(\varphi)$ and $\alpha_\tau(\mathbf{p})$ are Lipschitz-smooth, there exists $\beta > 0$ such that

$$\|\nabla_\varphi\alpha_\tau(\varphi) - \nabla_\varphi\alpha_\tau(\varphi')\| \leq \beta\|\varphi - \varphi'\| + \beta\|\mathbf{p}(\varphi) - \mathbf{p}(\varphi')\|$$
$$\leq \beta\|\varphi - \varphi'\| + \beta\|\varphi - \varphi'\|$$
$$= 2\beta\|\varphi - \varphi'\|.$$

We conclude that $\nabla_\varphi\alpha_\tau$ is Lipschitz.

By (i) and (ii), $\ell$ is Lipschitz-smooth w.r.t. the meta-parameters $\varphi$. Therefore, $\mathcal{L}_{\mathrm{meta}}(\varphi)$ is Lipschitz-smooth w.r.t. $\varphi$ with a Lipschitz constant $\beta_{\mathrm{meta}} > 0$. $\qquad\square$

**Theorem 1.** *Let the step size be* $\eta_t = \min(1/\sqrt{T}, 1/\beta_{meta})$. *Algorithm 3 satisfies* $\min_{1\leq t\leq T}\mathbb{E}\|\nabla_{(\theta_t,\phi_t)}\mathcal{L}_{meta}(\theta_t,\phi_t)\|^2 = \mathcal{O}\left(\sigma_{\mathbf{g}}^2/\sqrt{T}\right)$, *where the expectation is taken over the random training samples.*

The proof is similar to non-convex stochastic programming [4].

*Proof of Theorem 1.* Let $\varphi = (\theta,\phi)$. Let $\zeta_t = \nabla_{\varphi_t}\mathcal{L}_{\mathrm{meta}}(\varphi_t) - \frac{1}{b}\sum_{\tau\in\mathcal{B}_t}\mathbf{g}_\tau$, where $\frac{1}{b}\sum_{\tau\in\mathcal{B}_t}\mathbf{g}_\tau$ is an unbiased estimation of $\nabla_{\varphi_t}\mathcal{L}_{\mathrm{meta}}(\varphi_t)$, Using the Taylor expansion, we have

$$\mathcal{L}_{\mathrm{meta}}(\varphi_{t+1})$$
$$\leq \mathcal{L}_{\mathrm{meta}}(\varphi_t) + \nabla_{\varphi_t}\mathcal{L}_{\mathrm{meta}}(\varphi_t)^\top(\varphi_{t+1} - \varphi_t) + \frac{1}{2}\beta_{\mathrm{meta}}\|\varphi_{t+1} - \varphi_t\|^2$$
$$\leq \mathcal{L}_{\mathrm{meta}}(\varphi_t) - \eta_t(1 - \frac{\beta_{\mathrm{meta}}\eta_t}{2})\|\nabla_{\varphi_t}\mathcal{L}_{\mathrm{meta}}(\varphi_t)\|^2 + \eta_t\nabla_{\varphi_t}^\top\mathcal{L}_{\mathrm{meta}}(\varphi_t)\zeta_t + \frac{1}{2}\beta_{\mathrm{meta}}\eta_t^2\sigma_{\mathbf{g}}^2.$$

Taking conditional expectation over $\zeta_{t-1}$ on both sides and then take the expectation over the random training samples, we have

$$\mathbb{E}\mathcal{L}_{\mathrm{meta}}(\varphi_{t+1}) \leq \mathbb{E}\mathcal{L}_{\mathrm{meta}}(\varphi_t) - \frac{\eta_t}{2}\mathbb{E}\|\nabla_{\varphi_t}\mathcal{L}_{\mathrm{meta}}(\varphi_t)\|^2 + \frac{1}{2}\beta_{\mathrm{meta}}\eta_t^2\sigma_{\mathbf{g}}^2, \qquad (7)$$

where we have used $1 - \frac{\beta_{\mathrm{meta}}\eta_t}{2} \geq \frac{1}{2}$. Rearranging the above inequality and summing over $t$, we have

$$\sum_{t=1}^T \frac{\eta_t}{2}\mathbb{E}\|\nabla_{\varphi_t}\mathcal{L}_{\mathrm{meta}}(\varphi_t)\|^2 \leq \mathbb{E}\mathcal{L}_{\mathrm{meta}}(\varphi_1) + \beta_{\mathrm{meta}}\sigma_{\mathbf{g}}^2\sum_{t=1}^T \eta_t^2. \qquad (8)$$

Since $\eta_t = \min(1/\sqrt{T}, 1/2\beta_{\mathrm{meta}})$, we have $\sum_{t=1}^T \eta_t^2 \leq 1$. Diving both sides by $1/\sqrt{T}$, we conclude that $\min_{1\leq t\leq T}\mathbb{E}\|\nabla_{\varphi_t}\mathcal{L}_{\mathrm{meta}}(\varphi_t)\|^2 = \mathcal{O}\left(\sigma_{\mathbf{g}}^2/\sqrt{T}\right)$. $\qquad\square$

## 1.5 Proof of Theorem 2

**Theorem 2.** *Assume that* $\mathcal{M}(\theta,\phi)$ *is uniform conditioning. (i) Let* $\eta_t = \min(1/\sqrt{T}, 1/2\beta_{meta})$. *Algorithm 3 satisfies* $\min_{1\leq t\leq T}\mathbb{E}\mathcal{L}_{meta}(\theta_t,\phi_t) - \min_{(\theta,\phi)}\mathcal{L}_{meta}(\theta,\phi) = \mathcal{O}\left(\sigma_{\mathbf{g}}^2/\sqrt{T}\right)$, *where the expectation is taken over the random training samples. (ii) Let* $\eta_t = \eta < \min(1/2\beta_{meta}, 4|\mathcal{T}|/\rho\mu)$ *and* $\mathcal{B}_t = \mathcal{T}$. *Algorithm 3 satisfies* $\mathcal{L}_{meta}(\theta_t,\phi_t) - \min_{(\theta,\phi)}\mathcal{L}_{meta}(\theta,\phi) = \mathcal{O}((1 - \eta\rho\mu/4|\mathcal{T}|)^t)$.

*Proof of Theorem 2.* Let $\boldsymbol{\varphi} = (\boldsymbol{\theta}, \boldsymbol{\phi})$. By the chain rule, we have

$$\nabla_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}) = \frac{1}{|\mathcal{T}|} \sum_{\tau \in \mathcal{T}} \sum_{(\mathbf{x},y) \in Q_\tau} \nabla_1 \ell(\hat{y}, y) \nabla_{\boldsymbol{\varphi}} \hat{y} \tag{9}$$

$$= \frac{1}{|\mathcal{T}|} \mathcal{G}(\boldsymbol{\varphi})^\top \nabla_{\boldsymbol{\varphi}} \mathcal{M}(\boldsymbol{\varphi}), \tag{10}$$

where $\mathcal{G}(\boldsymbol{\varphi}) \equiv [\cdots \quad \nabla_1 \ell(\hat{y}, y) \quad \cdots] \in \mathbb{R}^{n_q |\mathcal{T}|}$ stacks all gradients of the losses on query examples as a vector. Hence, we establish the Polyak-Lojasiewicz (PL) inequality [7] as follows

$$\|\nabla_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi})\|^2 = \frac{1}{|\mathcal{T}|^2} \left\| \mathcal{G}(\boldsymbol{\varphi})^\top \nabla_{\boldsymbol{\varphi}} \mathcal{M}(\boldsymbol{\varphi}) \right\|^2$$

$$= \frac{1}{|\mathcal{T}|^2} \mathcal{G}(\boldsymbol{\varphi})^\top \nabla_{\boldsymbol{\varphi}} \mathcal{M}(\boldsymbol{\varphi}) \nabla_{\boldsymbol{\varphi}}^\top \mathcal{M}(\boldsymbol{\varphi}) \mathcal{G}(\boldsymbol{\varphi})$$

$$\geq \frac{\mu}{|\mathcal{T}|^2} \|\mathcal{G}(\boldsymbol{\varphi})\|^2 \qquad \text{(uniform conditioning)}$$

$$= \frac{\mu}{|\mathcal{T}|^2} \sum_{\tau \in \mathcal{T}} \sum_{(\mathbf{x},y) \in Q_\tau} (\nabla_1 \ell(\hat{y}, y))^2$$

$$\geq \frac{\mu \rho}{2|\mathcal{T}|^2} \sum_{\tau \in \mathcal{T}} \sum_{(\mathbf{x},y) \in Q_\tau} \left( \ell(\hat{y}, y) - \min_{y'} \ell(y', y) \right) \qquad \text{(strongly convex)}$$

$$\geq \frac{\mu \rho}{2|\mathcal{T}|} \left( \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}) - \min_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}) \right).$$

The PL inequality is commonly used in proving the global convergence of nonconvex optimization [5, 6]. Then, $\min_{1 \leq t \leq T} \mathbb{E} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t) - \min_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}) = \mathcal{O}\left( \sigma_{\mathbf{g}}^2 / \sqrt{T} \right)$ follows directly from Theorem 1.

For full gradient descent, the gradient noise $\boldsymbol{\zeta}_t = \nabla_{\boldsymbol{\varphi}_t} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t) - \frac{1}{b} \sum_{\tau \in \mathcal{B}_t} \mathbf{g}_\tau = \mathbf{0}$, thus, the noisy gradient will be the true gradient. By the Taylor expansion, it follows that

$$\mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_{t+1}) - \min_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi})$$

$$\leq \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t) + \nabla_{\boldsymbol{\varphi}_t}^\top \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t)(\boldsymbol{\varphi}_{t+1} - \boldsymbol{\varphi}_t) + \frac{\beta_{\text{meta}}}{2} \|\boldsymbol{\varphi}_{t+1} - \boldsymbol{\varphi}_t\|^2 - \min_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi})$$

$$= \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t) - \eta \|\nabla_{\boldsymbol{\varphi}_t} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t)\|^2 + \frac{\eta^2 \beta_{\text{meta}}}{2} \|\nabla_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t)\|^2 - \min_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi})$$

$$\leq \left( 1 - \frac{\eta \mu \rho}{4|\mathcal{T}|} \right) \left( \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}_t) - \min_{\boldsymbol{\varphi}} \mathcal{L}_{\text{meta}}(\boldsymbol{\varphi}) \right),$$

and we obtain the exponential convergence. $\square$

# 2 Additional Experiments

## 2.1 Compared with MAML using a wide network on *Sine*

As the network width is critical to MAML, we perform few-shot regression experiments on *Sine* using the setting in [10]. We compare MetaProx with MAML that uses a larger (denoted by LargeMAML) and wider (denoted by VeryWideMAML) network. As can be seen from Table 1, MetaProx achieves the best performance.

## 2.2 MetaProx with RBF kernel on *Sine*

In this section, we evaluate the performance of MetaProx with the radial basis function (RBF) kernel on *Sine*. The RBF kernel is $\mathcal{K}(\mathbf{x}_1, \mathbf{x}_2) = \exp\left( -\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2} \right)$, where $\sigma > 0$. Table 2 reports the results when $\sigma$ varies from $\{0.01, 0.05, 0.1, 0.5, 1.0, 5.0\}$. As can be seen, a simple linear kernel is better.

4

Table 1: Average MSE (with 95% confidence intervals) of few-shot regression on *Sine* using the settings in [10]. Results of baselines are from [10].

| method | 5-shot | 10-shot |
|---|---|---|
| OriginalMAML [3] | $0.390 \pm 0.156$ | $0.114 \pm 0.010$ |
| LargeMAML | $0.208 \pm 0.009$ | $0.061 \pm 0.004$ |
| VeryWideMAML | $0.205 \pm 0.013$ | $0.059 \pm 0.010$ |
| MetaFun [10] | $0.040 \pm 0.008$ | $0.017 \pm 0.005$ |
| MetaProx (proposed) | $\mathbf{0.010 \pm 0.001}$ | $\mathbf{0.002 \pm 0.001}$ |

Table 2: Average MSE (with 95% confidence intervals) of MetaProx with different base kernels on *Sine* (noise-free).

| kernel | 2-shot | 5-shot |
|---|---|---|
| RBF (0.01) | $2.92 \pm 0.19$ | $2.78 \pm 0.18$ |
| RBF (0.05) | $2.72 \pm 0.18$ | $2.36 \pm 0.17$ |
| RBF (0.1) | $2.50 \pm 0.17$ | $2.25 \pm 0.14$ |
| RBF (0.5) | $2.38 \pm 0.16$ | $1.71 \pm 0.13$ |
| RBF (1.0) | $2.36 \pm 0.16$ | $1.68 \pm 0.12$ |
| RBF (5.0) | $2.38 \pm 0.15$ | $1.72 \pm 0.13$ |
| linear | $\mathbf{0.11 \pm 0.01}$ | $\mathbf{0.01 \pm 0.00}$ |

# References

[1] J. C. Duchi. Introductory lectures on stochastic optimization. *The mathematics of data*, 25:99, 2018.

[2] A. Fallah, A. Mokhtari, and A. Ozdaglar. On the convergence theory of gradient-based model-agnostic meta-learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pages 1082–1092, 2020.

[3] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pages 1126–1135, 2017.

[4] S. Ghadimi and G. Lan. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.

[5] H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer, 2016.

[6] C. Liu, L. Zhu, and M. Belkin. Toward a theory of optimization for over-parameterized systems of non-linear equations: the lessons of deep learning. Preprint arXiv:2003.00307, 2020.

[7] B. T. Polyak. Gradient methods for minimizing functionals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 3(4):643–653, 1963.

[8] S. J. Reddi, A. Hefny, S. Sra, B. Poczos, and A. Smola. Stochastic variance reduction for nonconvex optimization. In *International Conference on Machine Learning*, pages 314–323, 2016.

[9] W. Rudin. *Principles of Mathematical Analysis*. McGraw-Hill, New York, 1976.

[10] J. Xu, J.-F. Ton, H. Kim, A. Kosiorek, and Y. W. Teh. MetaFun: Meta-learning with iterative functional updates. In *International Conference on Machine Learning*, pages 10617–10627, 2020.

[11] P. Zhou, X. Yuan, H. Xu, S. Yan, and J. Feng. Efficient meta learning via minibatch proximal update. In *Neural Information Processing Systems*, pages 1534–1544, 2019.