

$$2a) i) \theta'_i = \theta - 2 \nabla_{\theta} L_{T_i}(f_{\theta})$$

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum L_{T_i}(f_{\theta'_i})$$

$$= \theta - \beta \sum_{T_i \sim p(T)} \nabla_{\theta} L_{T_i}(f_{\theta'_i})$$

$$= \theta - \beta \sum_{T_i \sim p(T)} (\nabla_{\theta} \theta'_i) \nabla_{\theta'_i} L_{T_i}(f_{\theta'_i})$$

$$= \theta - \beta \sum_{T_i \sim p(T)} \underbrace{(\mathbf{I} - 2 \nabla_{\theta}^2 L_{T_i}(f_{\theta}))}_{\text{Hessian matrix calculation is required}} \nabla_{\theta'_i} L_{T_i}(f_{\theta'_i})$$

Hessian matrix calculation
is required

ii) This is because the adaptation of MAML requires hessian matrix calculation, as shown in 2a(i).

As we know, hessian matrix is a technique for calculating the second derivative of an n -dimensional function.

$H(x) \in \mathbb{R}^{n \times n}$ (the second order derivative and a symmetric matrix).

Therefore, this makes MAML computationally expensive.

iii) In 1st order approximation, we regard

$$\mathbf{I} - 2 \nabla_{\theta}^2 L_{T_i}(f_{\theta})$$

as identity matrix \mathbf{I} .

Update rule of MAML with 1st order approximation:

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{T \in \mathcal{P}(T)} L_{T_i} (f_{\theta} - \alpha \delta)$$

where $\delta \leftarrow \nabla L_{T_i} (f_{\theta})$

c) I Selected Meta-SGD: learning to learn quickly for Few-Shot Learning

i) previous work, such as meta-LSTM, are difficult to train. Moreover, each parameter of the learner is updated independently in each step and thus limits its potential.

ii) Meta-SGD has a much higher capacity by learning to learn not just the learner initialization, but also the learner update direction and learning rate in a single meta-learning process.

iii) The objective function of meta-SGD is as follows:

$$\min_{\theta, \alpha} E_{T \in \mathcal{P}(T)} [L_{\text{Test}(T)}(\theta')] = E_{T \in \mathcal{P}(T)} [L_{\text{Test}(T)}(\theta - \alpha \nabla L_{\text{Train}(T)}(\theta))]$$

where θ and α are (meta-) parameters of the meta-learner to be learned. To be more precise, θ represents the state of a learner that can be used to initialize the learner for any new task, and α is a vector of the same size as θ that decides both the update direction and learning rate. The adaptation term $\alpha \cdot \nabla L_T(\theta)$ is a vector whose direction represents the update direction and whose length represents the learning rate.

Therefore, given a few examples $T = \{(x_i, y_i)\}$ for a few shot learning problem, meta-SGD first initializes the learner with θ and then adapts it to θ' in just one step, in a new direction $\alpha \cdot \nabla L_T(\theta)$ different from the gradient $\nabla L_T(\theta)$ and using a learning rate implicitly implemented in $\alpha \cdot \nabla L_T(\theta)$.

(iv) Because in meta-SGD, the initialization, update direction and learning rate are all learned via meta-learning.