

Bilevel optimization in MAML

meta-learner $\rightarrow \theta$ (meta-parameters)fast-learner \rightarrow copies θ & adapts on each task

$$\theta' = \theta - \alpha \frac{\partial L_{\text{train}}}{\partial \theta}$$

 \rightarrow meta-steps - how many times algo is run $\rightarrow 2000$ in our case

↳ outer loop

 \rightarrow sample tasks from distribution $p(i)$ \rightarrow compute val loss before adaptation

each task has

 \rightarrow K inner loop steps↳ produces adapted θ' training data - D_i^{train} validation data - D_i^{val}

For each method :

inner loop (Task Adaptation)

$$\theta_i^{(1)} = \theta - \alpha \nabla_{\theta} L_{\text{train}}(\theta)$$

↳ do this for each task

 K steps

$$\theta^{(0)} = \theta$$

$$\theta^{(1)} = \theta^{(0)} - \alpha \nabla_{\theta} L_{\text{train}}(\theta^{(0)})$$

$$\theta^{(2)} = \theta^{(1)} - \alpha \nabla_{\theta} L_{\text{train}}(\theta^{(1)})$$

⋮

$$\theta^{(K)} = \theta^{(K-1)} - \alpha \nabla_{\theta} L_{\text{train}}(\theta^{(K-1)})$$

} parameters update

start with θ , predict outputcompute loss wrt this loop's parameter θ

backward pass (to find gradient)

gradient descent update.

final $\theta^{(K)}$ \rightarrow task adapted parameters

$$\hat{y}_{\text{val}} = f_{\theta^{(K)}}(x_{\text{val}})$$

val loss after adaptation:

 $L_{\text{val}}(\theta^{(K)}) \rightarrow$ how good initialisation θ
 was at letting us adapt to task

Now in explicit maml:

for meta-gradient step:

$$\nabla_{\theta} L_{\text{val}}(\theta^{(k)})$$

since $\theta_i^{(k)}$ depends on θ through k inner steps:

$$\theta_i^{(k)} = \theta - \alpha \nabla_{\theta} L_{\text{train}}(\theta) - \alpha \nabla_{\theta^{(1)}} L_{\text{train}}(\theta^{(1)}) - \dots$$

 $\Rightarrow \frac{\partial L_{\text{val}}}{\partial \theta}$ chain rule through all k steps

This introduces Hessian terms:

$$\nabla_{\theta} L_{\text{val}}(\theta^{(k)}) = \frac{\partial L_{\text{val}}}{\partial \theta^{(k)}} \cdot \frac{\partial \theta^{(k)}}{\partial \theta}$$

} Repeat this process for each task and average the auto-gradients:

$$g_{\text{meta}} = \frac{1}{B} \sum_i \nabla_{\theta} L_{\text{val}}(\theta_i^{(k)})$$

Then update meta parameters:

$$\theta \leftarrow \theta - \beta g_{\text{meta}}$$

↳ meta learning rate

One meta iteration

our code runs for 2000 meta-steps

pre-update loss - loss before inner loop gradient step } train

post update loss - loss after inner loop K steps } valval loss for each outer loop. \hookrightarrow how well model fits task after adaptation

averaged over all tasks

FOMAML : full maml meta-gradient:

$$\nabla_{\theta} L_{\text{val}}(\theta') = \nabla_{\theta} L_{\text{val}}(\theta) \cdot (I - \alpha \nabla_{\theta} L_{\text{train}}(\theta))$$

FOMAML : drops the Hessians:

$$\nabla_{\theta} L_{\text{val}}(\theta') \approx \nabla_{\theta} L_{\text{val}}(\theta')$$

Implicit MAML: solves a minimization problem

we treat $\theta' \rightarrow$ adapted parameters as solution toan optimization problem $\theta' \rightarrow$

so gradient = 0

$$\theta_i^* = \underset{\theta}{\operatorname{argmin}} \left(L_{\text{train}}(\theta) + \frac{\lambda}{2} \| \theta - \theta_i \|_2^2 \right)$$

so now:

$$\nabla_{\theta} L_{\text{train}}(\theta_i^*) + \lambda (\theta_i^* - \theta) = 0$$

differentiating wrt θ :

$$\nabla_{\theta} L_{\text{train}}(\theta_i^*) \cdot \frac{\partial \theta_i^*}{\partial \theta} + \lambda \left(\frac{\partial \theta_i^*}{\partial \theta} - I \right) = 0$$

This gives a linear system:

$$\left(\nabla_{\theta} L_{\text{train}}(\theta_i^*) + \lambda I \right) \cdot \frac{\partial \theta_i^*}{\partial \theta} = \lambda I$$

$$\frac{\partial \theta_i^*}{\partial \theta} = (H_{\text{train}} + \lambda I)^{-1} \lambda I$$

plugging this jacobian into our gradient:

$$\nabla_{\theta} L_{\text{val}}(\theta_i^*) = \lambda (H + \lambda I)^{-1} \nabla_{\theta} L_{\text{val}}(\theta_i^*)$$

computed by solving and Hessian vector product

Reptile:

1. Starts with θ 2. K inner loops \rightarrow get adapted parameter ϕ

3. outer update

$$\theta \leftarrow \theta + \epsilon(\phi - \theta)$$

Reptile is very close to FOMAML

Newmann:

$$\nabla_{\theta} L_{\text{val}}(\phi) = \nabla_{\phi} L_{\text{val}}(\phi) \frac{\partial \phi}{\partial \theta}$$

$$\frac{\partial \phi}{\partial \theta} = (I - \alpha H)^{-1}$$

↳ Hessian of training loss

$$(I - \alpha H)^{-1} \approx I + \alpha H + \alpha^2 H^2 \dots$$

$$\Rightarrow (I - \alpha H)^{-1} \approx I + \alpha H + \alpha^2 H^2 \dots$$

$$\nabla_{\theta} L_{\text{val}}(\phi) = g_{\text{val}} + \alpha H g_{\text{val}} + \alpha^2 H^2 g_{\text{val}} + \dots$$

↳ uses Hessian vector product

Penalty method: bilevel approximation where meta objective is:

$$L_{\text{meta}}(\theta) = L_{\text{val}}(\theta) + \lambda \| \nabla_{\theta} L_{\text{train}}(\theta) \|_2^2$$

if inner loop converged perfectly, then:

$$\nabla_{\theta} L_{\text{train}}(\phi^*) = 0$$

uses Hessian vector product