

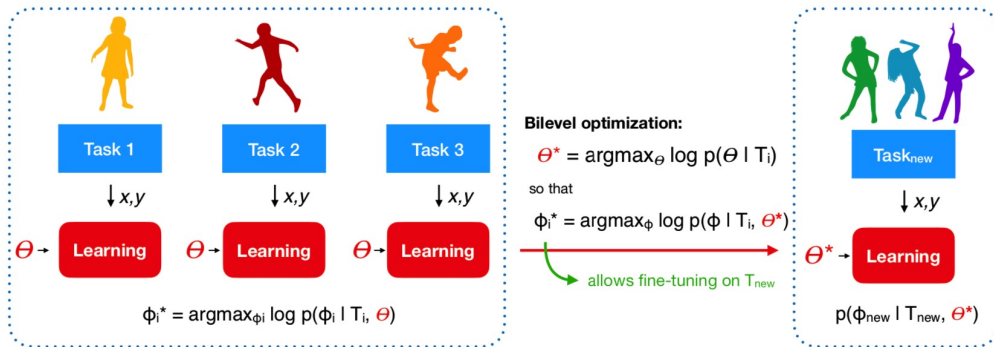
AutoML: Meta-Learning

Gradient-based meta-learning

Bernd Bischl Frank Hutter Lars Kotthoff
Marius Lindauer Joaquin Vanschoren

Gradient-based meta-learning

parameterize some aspect of the learner that we want to learn as
meta-parameters θ meta-learn θ across tasks



$\Theta(Prior)$, could encode an initialization ϕ , the hyperparameters λ , the optimizer,...

Learned θ^ should learn T_{new} from small amount of data, yet generalize to a large number of tasks*

Model agnostic meta-learning (MAML)

Meta-training

- Current initialization θ , model f_θ
- On i tasks, perform k SGD steps to find ϕ_i^* , then evaluate

derivative of test-set loss

$$\nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- Update task-specific parameters: $\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$

- Update θ to minimize sum of per-task losses, repeat α, β : learning rates

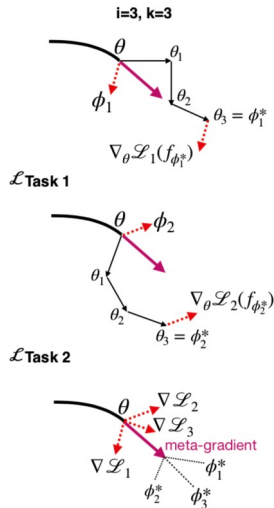
$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f_{\phi_i})$$

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f(\theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})))$$

meta-gradient: second-order gradient + backpropagate
compute how changes in θ affect the gradient at new θ

Meta-testing

- Training data of new task D_{train}
- θ^* : pre-trained parameters
- Finetune: $\phi = \theta^* - \alpha \nabla_{\theta} \mathcal{L}(f_\theta)$



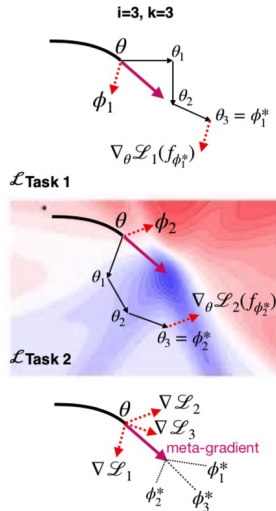
Model agnostic meta-learning (MAML)

Meta-training

- Current initialization θ , model f_θ
- On i tasks, perform k SGD steps to find ϕ_i^* , then evaluate
 $\phi_i = \theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
derivative of test-set loss $\nabla_\theta \mathcal{L}_i(f_{\phi_i^*})$
- Update task-specific parameters:
- Update θ to minimize sum of per-task losses, repeat
 $\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f_{\phi_i^*})$
 $\theta \leftarrow \theta - \beta \nabla_\theta \sum_i \mathcal{L}_i(f(\theta - \alpha \nabla_\theta \mathcal{L}_i(f_{\phi_i^*})))$
 α, β : learning rates
meta-gradient: second-order gradient + backpropagate
compute how changes in θ affect the gradient at new θ

Meta-testing

- Training data of new task D_{train}
- θ^* : pre-trained parameters
- Finetune: $\phi = \theta^* - \alpha \nabla_\theta \mathcal{L}(f_\theta)$



Model agnostic meta-learning (MAML)

Meta-training

- Current initialization θ , model f_θ
- On i tasks, perform k SGD steps to find ϕ_i^* , then evaluate

derivative of test-set loss

$$\nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- Update task-specific parameters:
$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$
- Update θ to minimize sum of per-task losses, repeat

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f_{\phi_i})$$

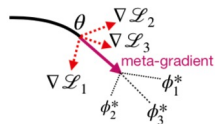
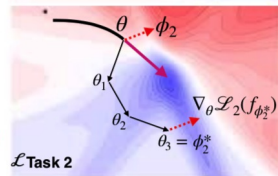
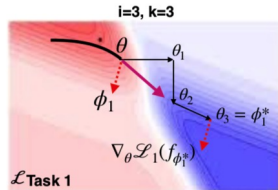
α, β : learning rates

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f(\theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})))$$

meta-gradient: second-order gradient + backpropagate
compute how changes in θ affect the gradient at new θ

Meta-testing

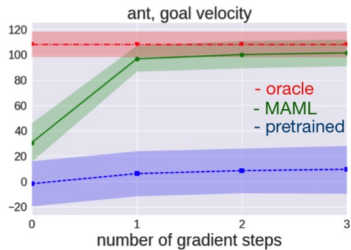
- Training data of new task D_{train}
- θ^* : pre-trained parameters
- Finetune: $\phi = \theta^* - \alpha \nabla_{\theta} \mathcal{L}(f_\theta)$



Model agnostic meta-learning (MAML)

Example of reinforcement learning:

- Goal: reach certain velocity in certain direction



Other gradient-based techniques

- Changing update rule yield different variants:

- MAML ^{1,6}

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_i(f_{\phi_i})$$

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- MetaSGD ²

$$\phi_i = \theta - \alpha \text{diag}(w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

w: weight per parameter

$$\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*}, w)$$

- Tnet ³

$$\phi_i = \theta - \alpha B(\theta, w) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$$

- Meta curvature ⁴ $\phi_i = \theta - \alpha P(\theta, \phi_i) \nabla_{\theta} \mathcal{L}_i(f_{\phi_i^*})$

- Online MAML (Follow the Meta Leader) ⁷

- WarpGrad ⁵

- Minimizes regret

- Robust, but computation costs grow over time

**All use second-order gradients,
Meta-learn a transformation of the
gradient for better adaptation**

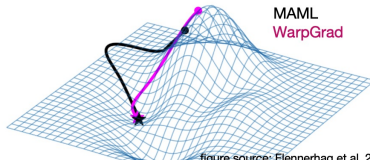


figure source: Flennerhag et al. 2019

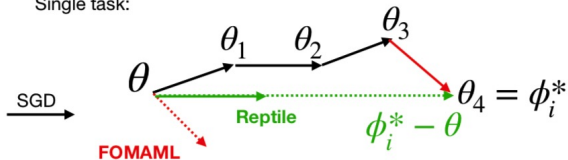
Scalability

- Backpropagating derivatives of ϕ_i wrt θ is compute + memory intensive (for large models)
- First order approximations of MAML:
 - First order MAML¹ uses only the last inner gradient update:

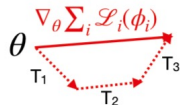
$$\theta \leftarrow \theta - \beta \sum_i \mathcal{L}_i(f_{\phi_i^*})$$
 - Reptile²: iterate over tasks, update θ in direction of ϕ_i^*

$$\theta \leftarrow \theta - \beta (\phi_i^* - \theta)$$

Single task:



FOMAML meta-gradient:



Generalizability

¹ [Finn et al. 2018](#)

² [Raghu et al. 2020](#)

³ [Tian et al. 2020](#)

⁴ [Stadie et al. 2019](#)

- MAML is more resilient to overfitting than many other meta-learning techniques ¹
- Effectiveness seems mainly due to feature reuse (finding θ) ²
 - Fine-tuning *only the last layer* equally good
- On few-shot learning benchmarks, a good embedding outperforms most meta-learning ³
 - Learn representation on entire meta-dataset (merged into single task)
 - Train a linear classifier on embedded few-shot D_{train} , predict D_{test}



- For meta-RL, also learn how to *explore* (how to sample new environments)
- E-MAML: add exploration to meta-objective (allows longer term goals) ⁴

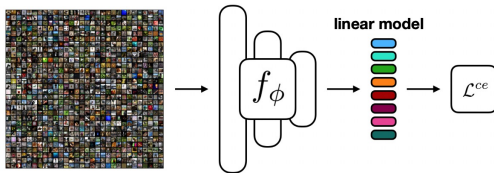
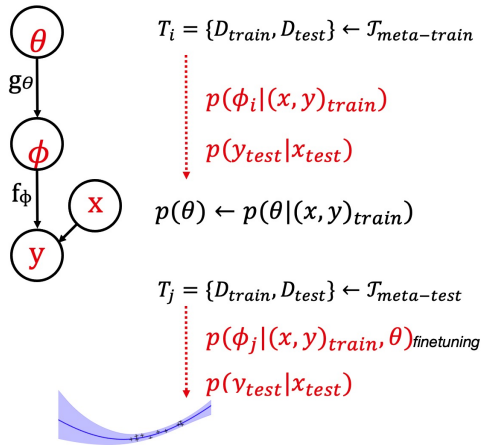
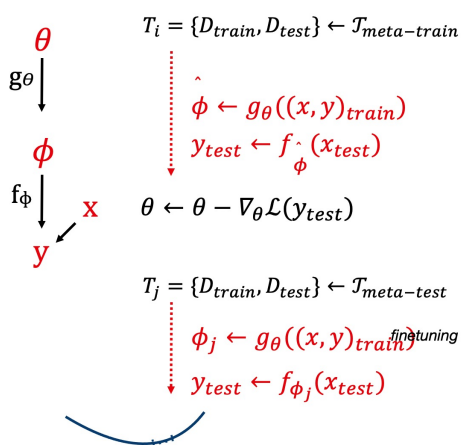


figure source: [Tian et al. 2020](#)

Bayesian meta-learning

Can meta-learning reason about **uncertainty** in the task distribution?



Fully Bayesian meta-learning

- Use approximation methods to represent uncertainty over ϕ_i
 - Sampling technique + variational inference: PLATIPUS ¹, BMAML ², ABML ³
 - Laplace approximation: LLAMA ⁴
 - Variational approximation of posterior:
 - Neural Statistician ⁵, Neural Processes ⁶
- What if our tasks are not IID?
 - Impose additional structure, e.g. with task-specific variable z ⁷

¹ [Finn et al. 2019](#)

² [Kim et al. 2018](#)

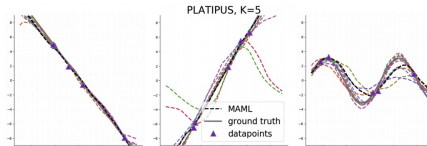
³ [Ravi et al. 2019](#)

⁴ [Grant et al. 2018](#)

⁵ [Edwards et al. 2017](#)

⁶ [Garnelo et al. 2018](#)

⁷ [Jerfel et al. 2019](#)



mini-ImageNet with filters for non-homogeneous tasks:



(a) plain



(b) blur

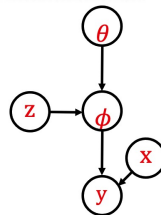


(c) night



(d) pencil

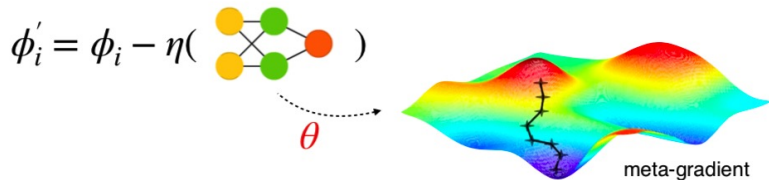
Figure source: [Jerfel et al. 2019](#)



Meta-learning optimizers

Our brains probably don't do backprop, instead:

- weights: networks that continuously modify the weights of another network
- Gradient based: parameterize the update rule using a neural network
 - ▶ Learn meta-parameters across tasks, by gradient descent



- Represent update rule as an LSTM, hierarchical RNN

Meta-learning optimizers

- Meta-learned (RNN) optimizers ‘rediscover’ momentum, gradient clipping, learning rate schedules, learning rate adaptation,...
- RL-based optimizers: represent updates as a policy, learn using guided policy search
- Combined with MAML:
 - ▶ learn per-parameter learning rates
 - ▶ learn precondition matrix (to ‘warp’ the loss surface)
- Black-box optimizers: meta-learned with an RNN, or with user-defined priors
- Speed up backpropagation by meta-learning sparsity and weight sharing

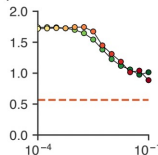
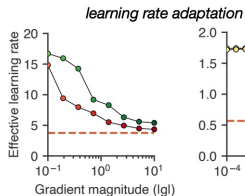
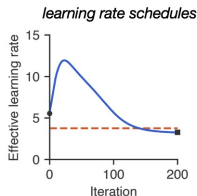


Figure source:
Maheswaranathan et al. 2020