

Modular Meta-learning: Variants of MAML

Reference - *Modular Meta-Learning with Shrinkage*

In general $w = \{\theta_1 \dots \theta_M\}$ e.g., different layers of a network. Variants of MAML learn a prior for w that is adapted for each task; but do all layers need to adapt? E.g. if only one layer is adapted, perhaps it could be trained for many more steps per task without risk of over-fitting. This paper learns to *differently* adapt each layer: assuming each θ_m is normally distributed as $\mathcal{N}(\phi_m, \sigma_m^2)$. **Layers with small or zero σ_m^2 will not adapt.** To learn ϕ, σ^2 we take Bayesian view:

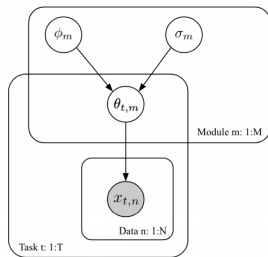
$$p(w^{1:T}, \mathcal{D} | \phi, \sigma^2) = \prod_{t=1}^T \prod_{m=1}^M \mathcal{N}(\theta_m^t | \phi_m, \sigma_m^2) \prod_{t=1}^T p(\mathcal{D}_t | t)$$

using the MAML approach to update ϕ, σ^2 :

the inner loop computes:

$$\hat{\theta}^t(\phi, \sigma^2) \equiv \operatorname{argmin}_{w^t} [-\log p(\mathcal{D}_t^{\text{Train}} | w^t) - \log p(w^t | \phi, \sigma^2)]$$

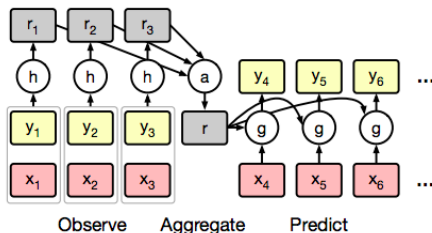
& the outer loop minimizes $\frac{1}{T} \sum_{t=1}^T -\log p(\mathcal{D}_t^{\text{Test}} | \hat{w}^t)$.



Model-based Meta-learning (cont) 1

Reference - [Conditional Neural Processes](#)

Training examples are passed through an MLP to generate representations; class-specific representations are aggregated and passed to a second classification MLP concatenated with query examples (from both test and train). Entire network is trained end-to-end on multiple tasks. Recent paper and seems to beat many baselines.



Model-based Meta-learning (cont) 2

Reference - *Meta-Learning with Memory-Augmented Neural Networks*

Datasets are presented as sequences, $\{(x_t, y_{t-1})\}$ as in the RNN-based approach. A 'memory-augmented LSTM/FF network' is the meta-learner.

Network updates rows of a memory matrix M with keys $k_t = \phi(x_t)^T$:
 $M_t = M_{t-1} + w^w k_t$; w^w = write weights
Retrieved memory $r_t = w^r M_t$ used to predict \hat{y}_t using a feedforward layer.
(read weights $w^r = \text{Softmax}(M_t \cdot k_t)$)
usage weights: $w_t^u = \gamma u_{t-1}^u + w_t^r + w_t^w$
used to compute w^w as follows:

'least used' $w^{lu} = 1$ for t smallest elements of w^u , 0 otherwise
 $w_t^w = \delta w_{t-1}^r + (1 - \delta) w_{t-1}^{lu}$; prior to writing the least used row of M is zeroed.
Read/write to 'memory' - we can view this as a 'neural Turing machine'.

