

Env

type 1

\prod policy θ_1^+

possible type:

type 4

type 5 type 6

type 2

\prod policy θ_2

type 3

\prod policy θ_3

initialize $\theta_1, \theta_2, \theta_3$

done indicate how many round we want to update.

while not done:

for each building L in Env:

for each task:

↳ create new env by switching each building for each loop.

Ex.

type 4.

\prod policy θ_1

type 2

\prod policy θ_2

type 3

\prod policy θ_3

↳ sample for 1 year and add data into replay buffer.

↳ update $\theta_{i, task_i}$ using below equation

$$\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\tau_i}(f_{\theta})$$

↳ with updated policy sample for 1 more year and update replay buffer
end loop

MAML

Update $\theta_i \leftarrow \theta - \beta \nabla_{\theta} \sum_{\tau_i \sim p(\tau)} \mathcal{L}_{\tau_i}(f_{\theta_i})$ using each \mathcal{D}'_i and \mathcal{L}_{τ_i} in Equation 4

$\theta_i^+ \leftarrow \theta_i$

The updated θ_i will be assign to the building from original environment.

Ex.

If θ_1^* has been updated then it will be assign to building 1.

Also, for updating θ_2 using MAML algorithm, building will use θ_1^* to make action.

Thus, our update of policy is similar to gibbs sampling.

Question.

↳ add buffer.

↳ do we always start new or keep them.

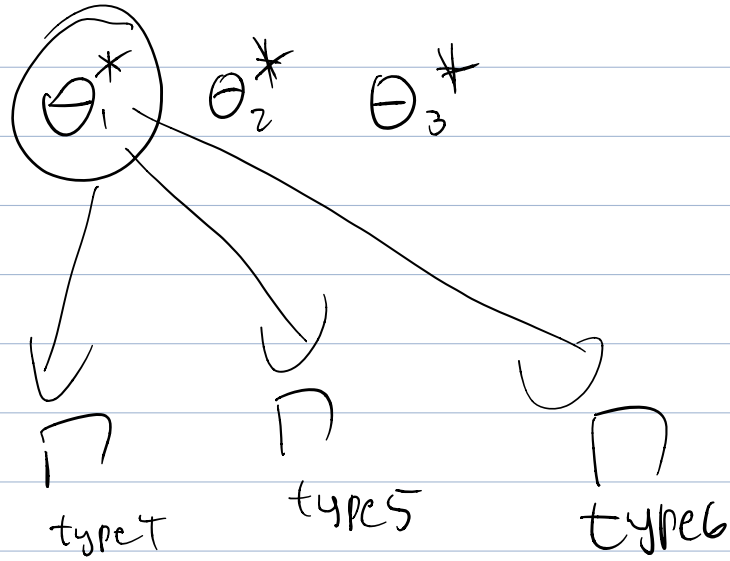
(within MAML and outside of MAML)

add on zeros

padding to image

add zero to take account

$\begin{bmatrix} 3 \\ 4 \end{bmatrix} \rightarrow$



reintorce ~~po~~ A2C \leftarrow

vanilla policy gradient.