
MAML Problem formulation

- Each task is a MDP with horizon H .
- Each task \mathcal{T}_i contains initial state distribution $q_i(x_1)$ and transition distribution $q_i(x_{t+1} \mid x_t, a_t)$
- Model being learned $f_\theta : X_t \rightarrow A_t$, where X_t is the set of states at time t and A_t is the set of actions at time t . Such that $t \in \{1, \dots, H\}$
- The loss for task \mathcal{T}_i and model f_ϕ is given by:

$$L_{\mathcal{T}_i}(f_\phi) = -\mathbb{E}_{x_t, a_t \sim f_\phi, q_{\mathcal{T}_i}} \left[\sum_{t=1}^H R_i(x_t, a_t) \right]$$

- Uses Policy gradient methods to estimate gradient both for model and meta-optimization, because the dynamics are unknown. Also have the option of using trust region optim.

VPG Pseudo Code (Spinning Up)

<https://spinningup.openai.com/en/latest/algorithms/vpg.html>

Algorithm 1 Vanilla Policy Gradient Algorithm

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 7: Compute policy update, either using standard gradient ascent,

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k,$$

or via another gradient ascent algorithm like Adam.

- 8: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 9: **end for**
-

TRPO Pseudo Code (Spinning Up)

<https://spinningup.openai.com/en/latest/algorithms/trpo.html>

Algorithm 1 Trust Region Policy Optimization

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: Hyperparameters: KL-divergence limit δ , backtracking coefficient α , maximum number of backtracking steps K
- 3: **for** $k = 0, 1, 2, \dots$ **do**
- 4: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 5: Compute rewards-to-go \hat{R}_t .
- 6: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 7: Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 8: Use the conjugate gradient algorithm to compute

$$\hat{x}_k \approx \hat{H}_k^{-1} \hat{g}_k,$$

where \hat{H}_k is the Hessian of the sample average KL-divergence.

- 9: Update the policy by backtracking line search with

$$\theta_{k+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{\hat{x}_k^T \hat{H}_k \hat{x}_k}} \hat{x}_k,$$

where $j \in \{0, 1, 2, \dots, K\}$ is the smallest value which improves the sample loss and satisfies the sample KL-divergence constraint.

- 10: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 11: **end for**
-

Pseudo Code

```
1 def task_rollout(task):
2
3     cumulative_disc_rew = 0
4
5     for i in range(run_length):
6
7         observation = get_observation()
8         action = self.model(observation)
9
10        disc_reward, log_prob = step(action) # Discounted reward, log prob from action distrib
11        cumulative_disc_rew += disc_reward * log_prob ### UNSURE ABOUT THIS #####
12
13        if run is done: # if in terminal state, or run_length reached
14            return cumulative_disc_rew
15
16
17 def sample_trajectories(task, K):
18
19     rewards_for_task = torch.zeros(K) # Cumulative sum of rewards * logprob for each K task
20     for i in range(K):
21         cumulative_rew = task_rollout(task)
22         rewards_for_task[i] = cumulative_rew
23
24     return rewards_for_task
25
26
27 def train_maml():
28     Tasks = [t_1, t_2, t_3, ..., t_n] # List of tasks
29     K: int # Number of rollouts per task
30     kl_div = [] # KL divergence for each model distrib. after update
31     updated_model_reward = []
32
33     while not done:
34         for task in Tasks:
35
36             self.model = copy.deepcopy(self.meta) # copy meta model for EACH sample
37             exp_rewards_for_tasks = sample_trajectories(task, K)
38
39             ### Update self.model with VPG
40             loss = - exp_rewards_for_tasks.mean() # (Eq 4 in paper, and line 6 in VPG)
41             loss.backward()
42             self.optim.step()
43
44             exp_rews_updated_model = sample_trajectories(task, K)
45             updated_model_reward.append(exp_rews_updated_model)
46
47             # KL Div. of updated and old policy
48             kl = torch.kl_divergence(self.model._distrib, self.meta._distrib)
49             kl_div.append(kl)
50
51             ### Update self.meta with TRPO
52             meta_grad = 0
53             for exp_rew in updated_model_reward:
54                 meta_grad += - exp_rew.mean()
55             meta_grad.backward() # This is g_k in TRPO, line 10 in MAML code
56             # meta_grad will be same size as network self.model
57             ### UNSURE HERE ^^^^^ #####
58
59             # KL Divergence can be computed for some non closed form sample, we can say model is
60             # actually some unknown distrib.
61             kl_average = kl_div.mean(axis=1) # axis=1, i.e. (10x5) -> (10x1)
62             kl_hessian = compute_hessian(kl_average) # size: ((obs_dim x obs_dim) x action_dim)???
63             x = inverse(kl_hessian) * meta_grad # Line 8 in TRPO pseudo code, exchange w/ conj. grad?
64             ### UNSURE HERE ^^^^^ #####
65
66             update_val = backtrack_line_search(kl_hessian, x) # Line 9 in TRPO pseudo code
67             self.meta = update_model(update_val) # Update model parameters
68 self.meta = torch_network()
69 self.model = None
```

```
69 train()
70
71
72 ### Questions ---
73 """
74
75
76
77 Paper says: "In order to avoid computing third derivatives, we use finite differences to compute
78 the Hessian-vectorproducts for TRPO.
79 Where is this 3rd derivative? If I'm not mistaken it should be in the backtracking line search?
80 """
```