

A. DICE derivations

In (Foerster et al., 2018), the surrogate function is given by:

$$J^{\text{DICE}} = \sum_{t=0}^{H-1} \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'})}{\perp(\pi_{\theta}(a_{t'}|s_{t'}))} \right) r(s_t, a_t), \quad (1)$$

where \perp denotes the ‘stop gradient’ operation. In expectation, the Monte Carlo estimate of the Hessian $\mathbb{E}_{\tau \sim P_{\tau}(\tau|\theta)}[\nabla_{\theta}^2 J^{\text{DICE}}]$ is equivalent to the Hessian of the inner objective:

$$\begin{aligned} \nabla_{\theta} J^{\text{DICE}} &= \sum_{t=0}^{H-1} \nabla_{\theta} \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'})}{\perp(\pi_{\theta}(a_{t'}|s_{t'}))} \right) r(s_t, a_t) \\ &= \sum_{t=0}^{H-1} \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'})}{\perp(\pi_{\theta}(a_{t'}|s_{t'}))} \right) \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right) r(s_t, a_t) \end{aligned} \quad (2)$$

Once again we take the derivative in order to obtain the Hessian of J^{DICE} :

$$\begin{aligned} \nabla_{\theta}^2 J^{\text{DICE}} &= \sum_{t=0}^{H-1} \nabla_{\theta} \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'})}{\perp(\pi_{\theta}(a_{t'}|s_{t'}))} \right) \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right) \end{aligned} \quad (3)$$

$$r(s_t, a_t) + \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'})}{\perp(\pi_{\theta}(a_{t'}|s_{t'}))} \right) \quad (4)$$

$$\nabla_{\theta} \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right) r(s_t, a_t) \quad (5)$$

$$\rightarrow \sum_{t=0}^{H-1} \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right) \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right)^{\top} \quad (6)$$

$$r(s_t, a_t) + \left(\sum_{t'=0}^t \nabla_{\theta}^2 \log \pi_{\theta}(a_{t'}|s_{t'}) \right) r(s_t, a_t) \quad (7)$$

It is easy to show that

$$\mathbb{E}_{\tau \sim P_{\tau}(\tau|\theta)}[\nabla_{\theta}^2 J^{\text{DICE}}] \quad (8)$$

$$= \mathbb{E}_{\tau \sim P_{\tau}(\tau|\theta)} \left[\sum_{t=0}^{H-1} \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right) \right] \quad (9)$$

$$\left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right)^{\top} r(s_t, a_t) \quad (10)$$

$$+ \left(\sum_{t'=0}^t \nabla_{\theta}^2 \log \pi_{\theta}(a_{t'}|s_{t'}) \right) r(s_t, a_t) \quad (11)$$

$$= \nabla_{\theta}^2 J_{\text{inner}} \quad (12)$$

B. LVC derivations

LVC (Rothfuss et al., 2019) is a method that trade-off bias for lower variance estimation of gradient.

$$\nabla_{\theta}^2 J^{\text{LVC}} \quad (13)$$

$$= \nabla_{\theta} \sum_{t=0}^{H-1} \frac{\pi_{\theta}(a_t|s_t)}{\perp(\pi_{\theta}(a_t|s_t))} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \left(\sum_{t'=t}^{H-1} r(s_{t'}, a_{t'}) \right) \quad (14)$$

$$= \sum_{t=0}^{H-1} \frac{\pi_{\theta}(a_t|s_t)}{\perp(\pi_{\theta}(a_t|s_t))} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)^{\top} \quad (15)$$

$$\left(\sum_{t'=t}^{H-1} r(s_{t'}, a_{t'}) \right) + \frac{\pi_{\theta}(a_t|s_t)}{\perp(\pi_{\theta}(a_t|s_t))} \nabla_{\theta}^2 \log \pi_{\theta}(a_t|s_t) \quad (16)$$

$$\left(\sum_{t'=t}^{H-1} r(s_{t'}, a_{t'}) \right) \quad (17)$$

$$\rightarrow \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)^{\top} \quad (18)$$

$$\left(\sum_{t'=t}^{H-1} r(s_{t'}, a_{t'}) \right) + \nabla_{\theta}^2 \log \pi_{\theta}(a_t|s_t) \left(\sum_{t'=t}^{H-1} r(s_{t'}, a_{t'}) \right) \quad (19)$$

$$= \sum_{t=0}^{H-1} \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)^{\top} \right) r(s_t, a_t) \quad (20)$$

$$+ \left(\sum_{t'=0}^t \nabla_{\theta}^2 \log \pi_{\theta}(a_{t'}|s_{t'}) \right) r(s_t, a_t) \quad (21)$$

It can be prove that $\mathbb{E}_{\tau \sim P_{\tau}(\tau|\theta)} [J^{\text{LVC}}]$ is a biased approximation of the $\nabla_{\theta}^2 J_{\text{inner}}(\theta)$ but with lower variance, see (Rothfuss et al., 2019) for more details.

$$\mathbb{E}_{\tau \sim P_{\tau}(\tau|\theta)} [J^{\text{LVC}}] \quad (22)$$

$$= \mathbb{E}_{\tau \sim P_{\tau}(\tau|\theta)} \left[\sum_{t=0}^{H-1} \left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}) \right) \right] \quad (23)$$

$$\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)^{\top} r(s_t, a_t) \quad (24)$$

$$\left. \right] + \quad (25)$$

$$\mathbb{E}_{\tau \sim P_{\tau}(\tau|\theta)} \left[\sum_{t=0}^{H-1} \left(\sum_{t'=0}^t \nabla_{\theta}^2 \log \pi_{\theta}(a_{t'}|s_{t'}) \right) r(s_t, a_t) \right] \quad (26)$$

$$\approx \nabla_{\theta}^2 J_{\text{inner}} \quad (27)$$

C. Proof of Theorem 1.

Firstly, we derive first order gradient of surrogate objective function.

$$\nabla_{\theta} J_i^{\text{TMAML}} \quad (28)$$

$$= \nabla_{\theta} \sum_{t=0}^{H-1} \left[1 - \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \right) \right] \left(1 - \frac{\pi_{\theta}(a_t|s_t, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_t|s_t, \mathcal{T}_i))} \right) b(s_t, \mathcal{T}_i) \quad (29)$$

$$= \sum_{t=0}^{H-1} b(s_t, \mathcal{T}_i) \left[- \frac{\pi_{\theta}(a_t|s_t, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_t|s_t, \mathcal{T}_i))} \right] \quad (30)$$

$$\nabla_{\theta} \log \pi_{\theta}(a_t|s_t, \mathcal{T}_i) \left[1 - \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \right) \right] \quad (31)$$

$$- \left(1 - \frac{\pi_{\theta}(a_t|s_t, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_t|s_t, \mathcal{T}_i))} \right) \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \right) \quad (32)$$

$$\left(\sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i) \right) \quad (33)$$

Once again we take the derivative of Equation (33) in order to obtain the Hessian of J_i^{TMAML} ,

$$\nabla_{\theta}^2 J_i^{\text{TMAML}} \quad (34)$$

$$= \sum_{t=0}^{H-1} b(s_t, \mathcal{T}_i) \left[- \left(1 - \prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \right) \right] \quad (35)$$

$$\nabla_{\theta}^2 \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} - \left(1 - \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \right) \quad (36)$$

$$\nabla_{\theta}^2 \prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \quad (37)$$

$$+ 2 \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \nabla_{\theta} \log \pi(a_t|s_t, \mathcal{T}_i) \quad (38)$$

$$\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \quad (39)$$

$$\rightarrow 2 \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi(a_t|s_t, \mathcal{T}_i) \quad (40)$$

$$\left[\sum_{t'=t}^{H-1} \nabla_{\theta} \log \pi(a_{t'}|s_{t'}, \mathcal{T}_i) b(s_{t'}, \mathcal{T}_i) \right] \quad (41)$$

In expectation, $\mathbb{E}_{\tau \sim P_{\mathcal{T}}(\tau|\theta)} [\nabla_{\theta}^2 J_i^{\text{TMAML}}]$ is zero, to show this, we use the form of $\nabla_{\theta}^2 J_i^{\text{TMAML}}$ just derived, for each t ,

we have

$$\mathbb{E}_{\tau \sim P_{\mathcal{T}}(\tau|\theta)} \left[\nabla_{\theta} \log \pi(a_t|s_t) \left[\sum_{t'=t+1}^{H-1} b(s_{t'}, \mathcal{T}_i) \nabla_{\theta} \log \pi(a_{t'}|s_{t'}, \mathcal{T}_i) \right] \right] = \mathbb{E}_{\tau \sim P_{\mathcal{T}}(\tau|\theta)} \left[\nabla_{\theta} \log \pi(a_t|s_t, \mathcal{T}_i) \right]$$

$$\left[\sum_{t'=t+1}^{H-1} b(s_{t'}, \mathcal{T}_i) \mathbb{E}_{a_{t'}, s_{t'}|a_t, s_t, \mathcal{T}_i} [\nabla_{\theta} \log \pi(a_{t'}|s_{t'}, \mathcal{T}_i)] \right] = 0 \text{ (score function expectation equals 0) ,}$$

Base on this, sum over all time step t will get the expectation $\nabla_{\theta}^2 J_i^{\text{TMAML}}$ equals 0,

$$\mathbb{E}_{\tau \sim P_{\mathcal{T}}(\tau|\theta)} [\nabla_{\theta}^2 J_i^{\text{TMAML}}] = 0$$

D. Proof of Theorem 2.

Here we show how J_i^{TMAML} can introduce control variate into the Hessian estimation. The intuition is getting a sum of reward in the Hessian formula instead of the sum of gradient which is impossible to introduce control variates. In order to do this, we first decompose the expectation of J_i^{DICE} into the desired form.

$$\nabla_{\theta}^2 (J_i^{\text{DICE}} + J_i^{\text{TMAML}}) \quad (42)$$

$$= \nabla_{\theta}^2 J_i^{\text{TMAML}} + \quad (43)$$

$$\sum_{t=0}^{H-1} \nabla_{\theta} \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \right) r(s_t, a_t, \mathcal{T}_i) \quad (44)$$

$$= \nabla_{\theta}^2 J_i^{\text{TMAML}} + \quad (45)$$

$$\sum_{t=0}^{H-1} \left(\prod_{t'=0}^t \frac{\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i)}{\perp(\pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i))} \right) r(s_t, a_t, \mathcal{T}_i) \quad (46)$$

$$\left[\left(\sum_{t'=t}^{H-1} \nabla_{\theta} \log \pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i) \right)^2 \right] \quad (47)$$

$$+ \sum_{t'=t}^{H-1} \nabla_{\theta}^2 \log \pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i) \quad (48)$$

$$\quad (49)$$

Then we substitute $\nabla_{\theta}^2 J_i^{\text{TMAML}}$ with Equation (34), we can derive,

$$\nabla_{\theta}^2 (J_i^{\text{DICE}} + J_i^{\text{TMAML}})$$

$$\rightarrow \sum_{t=0}^{H-1} \left[\nabla_{\theta}^2 \log \pi_{\theta}(a_t|s_t, \mathcal{T}_i) \left(\sum_{t'=0}^t r(a_{t'}, s_{t'}, \mathcal{T}_i) \right) \right]$$

$$+ 2 \sum_{t=0}^{H-1} \left[\nabla \log \pi_{\theta}(a_t|s_t, \mathcal{T}_i)^{\top} \left(\sum_{t'=t}^{H-1} \nabla \log \pi_{\theta}(a_{t'}|s_{t'}, \mathcal{T}_i) \right) \right]$$

$$\left(\sum_{k=t'}^{H-1} r(a_k, s_k) - b(s_t, \mathcal{T}_i) \right) \Bigg],$$

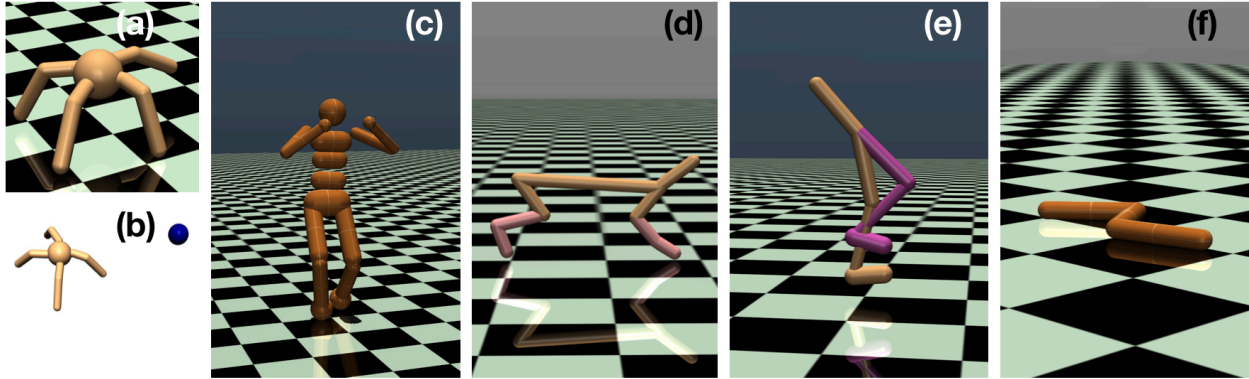


Figure 1. Meta-RL environments based on Mujoco (Todorov et al., 2012) and OpenAI Gym (Brockman et al., 2016): (a) Ant chasing random goal location, the tasks are generated by sampling the target positions from the uniform distribution on $[-3, 3]^2$; (b) An example of Ant chasing goal task, where green ball denotes a goal, the goal is randomly sampled at the beginning of each episode; (c) Humanoid; (d) HalfCheetah; (e) Walker; (f) Swimmer; We build different tasks including random velocity, random direction, and random goal location on each environment, see experiments in main paper for more details. In order to get high reward in these environments, agents must be able to do learning to learning or meta-learning from past experience.

Hyperparameters	
Policy Hidden Layer Sizes	64^2 (Humanoid: 128^2)
Num Adapt Steps	1 (Ant, rand goal: 2)
Inner Step Size α	0.01
Outer LR β	0.001
Meta CV Inner Step Size α'	0.01
Meta CV Outer LR β'	0.01
Meta-Tasks Per Iteration	40
Num Traj Per Meta-Task	20
Grad Steps Per Iteration	5
Outer Clip Ratio ϵ	0.3
KL Penalty Coefficient η	0.0005
Traj Length	100

Table 1. Hyperparameter settings used in each algorithm

ceedings of Machine Learning Research, pp. 1529–1538, Stockholm, Sweden, 10–15 Jul 2018. PMLR.

Rothfuss, J., Lee, D., Clavera, I., Asfour, T., and Abbeel, P. ProMP: Proximal meta-policy search. In *International Conference on Learning Representations*, 2019.

Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pp. 5026–5033. IEEE, 2012.

E. Hyperparameters and Environments Setup

In every experiment, each algorithm is tested with 5 random seeds. An example of tested environment is shown in Figure 1. Table 1 contains the hyperparameter settings used for the different algorithms.

References

- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Foerster, J., Farquhar, G., Al-Shedivat, M., Rocktäschel, T., Xing, E., and Whiteson, S. DiCE: The infinitely differentiable Monte Carlo estimator. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Pro-*