



POLITECNICO
MILANO 1863

META LEARNING THE STEP SIZE IN POLICY GRADIENT METHODS

L. SABBIONI, F. CORDA AND M. RESTELLI

{luca.sabbioni, marcello.restelli}@polimi.it, francesco.corda@mail.polimi.it

MOTIVATION

- We want to apply **Hyperparameter Optimization** for Policy Gradient methods in Reinforcement Learning.
- A dynamic learning rate may allow a **fast and stable** learning process.
- If the process depends on **external variables**, hyperparameters can be adapted to the context.

CONTRIBUTIONS

1. **Meta-MDP**: Definition of an abstract class of hyperparameter decision processes called **Meta-MDP**, a framework for modeling meta learning;
2. **Bounds for Lipschitz Contextual-MDPs**: Derivation of **general guarantees** on the return with respect to the parameterization of the context under smoothness assumptions;
3. **FQI on Meta-MDP**: **Fitted Q-Iteration** algorithm [Ernst et al., 2005] on the meta-MDP, using Natural Gradient Ascent as update rule, used to choose an adaptive step size throughout the learning process.

PROPERTIES

Assumptions: Lipschitz MDPs

$$\mathcal{K}(P_\omega(\cdot|s, a), P_\omega(\cdot|\bar{s}, \bar{a})) \leq L_P d_{\mathcal{S} \times \mathcal{A}}((s, a), (\bar{s}, \bar{a})),$$

$$|r_\omega(s, a) - r_\omega(\bar{s}, \bar{a})| \leq L_r d_{\mathcal{S} \times \mathcal{A}}((s, a), (\bar{s}, \bar{a})).$$

Assumptions: Lipschitz Policy

$$\mathcal{K}(\pi(\cdot|s), \pi(\cdot|\bar{s})) \leq L_\pi d_{\mathcal{S}}(s, \bar{s}),$$

- These assumptions allow the Value Function to be Lipschitz continuous with constant L_V [Pirotta et al., 2015]:

Assumptions: Context Lipschitz Continuity

$$\mathcal{K}(P_\omega(\cdot|s, a), P_{\hat{\omega}}(\cdot|s, a)) \leq L_{\omega_P} d_\Omega(\omega, \hat{\omega})$$

$$|r_\omega(s, a) - r_{\hat{\omega}}(s, a)| \leq L_{\omega_r} d_\Omega(\omega, \hat{\omega});$$

Results: Lipschitz Value functions/return

$$|Q_\omega(s, a) - Q_{\hat{\omega}}(s, a)| \leq L_{\omega_Q} d_\Omega(\omega, \hat{\omega}); \quad |j_\omega - j_{\hat{\omega}}| \leq L_{\omega_j} d_\Omega(\omega, \hat{\omega})$$

$$L_{\omega_Q} = \frac{L_{\omega_r} + \gamma L_{\omega_P} L_V}{1 - \gamma}$$

REFERENCES

- Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6(Apr):503–556, 2005.
- Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv preprint arXiv:1502.02259*, 2015.
- Matteo Pirotta, Marcello Restelli, and Luca Bascetta. Policy gradient in lipschitz markov decision processes. *Machine Learning*, 100(2):255–283, 2015.

META MDP

CMDP

Contextual Markov Decision Process [Hallak et al., 2015]:

$$\mathcal{M}(\omega) = \langle \mathcal{S}, \mathcal{A}, P_\omega, R_\omega, \gamma, \mu \rangle$$

- Transition and reward processes depend on context ω
- Expected return $j_\omega(\theta) = \mathbb{E}_{P_\omega, \pi_\theta} [\sum \gamma^t R_\omega(s_t, a_t)]$

Policy Gradient Update

$$\nabla_\theta j_\omega(\theta) = \mathbb{E} \left[\nabla_\theta \log \pi_\theta(a|s) Q_\omega^\pi(s, a) \right],$$

- Stochastic Gradient Ascent (SGA) with step size h

$$\theta_{t+1} = \theta_t + h \nabla j_\omega(\theta_t)$$

Meta-MDP

Meta Observation:	Task and policy info	$x_t = \langle \theta_t, \omega, \nabla j_\omega(\theta_t) \rangle$
Meta Action:	Hyperparameter choice	h
Meta reward:	One step Learning	$l = j_\omega(\theta_{t+1}) - j_\omega(\theta_t)$
Meta discount factor:		$\tilde{\gamma}$

General framework:

- Update rule f
- Hyperparameter h

$$\theta_{t+1} = f(\theta_t, h)$$

$$l = j_\omega(f(\theta_t, h)) - j_\omega(\theta_t)$$

DATASET GENERATION

$\mathcal{F} = \{ \}$

For $k = 1, \dots, K$

Sample context ω

Sample initial policy θ_0

Estimate $j_\omega(\theta_0), \nabla j_\omega(\theta_0)$ in \mathcal{M}_ω

For $t = 0, \dots, T$

Sample hyperparameter h_t ;

Update $\theta_{t+1} = \theta_t + h_t \nabla j_\omega(\theta_t)$

Estimate $j_\omega(\theta_{t+1}), \nabla j_\omega(\theta_{t+1})$ in \mathcal{M}_ω

Set x_t, x_{t+1} and $l = j_\omega(\theta_{t+1}) - j_\omega(\theta_t)$

Append (x_t, h_t, x_{t+1}, l) to \mathcal{F}

end for

end for

DOUBLE CLIPPED FITTED Q ITERATION

Set regression algorithm f (Extra-trees);

Initialize $Q_1^0, Q_2^0, N = 0$;

While stopping conditions not reached do:

Build training set $TS = \{(i, o)\}$ from \mathcal{F} where:

$$\underline{Q}(h) = \min(Q_1^N(x_{t+1}, h), Q_2^N(x_{t+1}, h)),$$

$$\overline{Q}(h) = \max(Q_1^N(x_{t+1}, h), Q_2^N(x_{t+1}, h)),$$

$$i = (x_t, h_t); \quad o = l + \tilde{\gamma} \max[\lambda \underline{Q}(h) + (1 - \lambda) \overline{Q}(h)];$$

Split TS in TS_1, TS_2 ;

Perform regressions with f on $TS_{1,2}$ to learn Q_1^{N+1}, Q_2^{N+1} .

end while

EMPIRICAL RESULTS

