

Classical Policy Gradient: Preserving Bellman’s Principle of Optimality

Philip S. Thomas, Scott M. Jordan, Yash Chandak, Chris Nota, and James Kostas
 University of Massachusetts Amherst, College of Information and Computer Sciences

In 1954, Richard Bellman wrote [1]:

Principle of Optimality: An optimal policy has the property that whatever the initial state and initial decisions are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decisions.

This principle of optimality has endured at the foundation of reinforcement learning research, and is central to what remains the classical definition of an optimal policy [2]. Classical reinforcement learning algorithms like Q-learning [3] embody this principle by striving to act optimally in every state that occurs, regardless of *when* the state occurs.

The start-state objective function, $\rho(\theta) = \mathbf{E}[\sum_{t=0}^{\infty} \gamma^t R_t | \theta]$, prioritizes making decisions optimally *in the initial state*, not necessarily in the states resulting from the first decisions.¹ These two goals (optimizing decisions in the initial state and optimizing decisions in subsequent states) can be conflicting when using function approximation, particularly when γ is small and the initial state distribution has limited support. So, maximizing ρ does not preserve the principle of optimality.

Let $q_{\theta}(s, a) = \mathbf{E}[\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t=s, A_t=a, \theta]$ so that

$$\nabla \rho(\theta) = \mathbf{E} \left[\sum_{t=0}^{\infty} \gamma^t q_{\theta}(S_t, A_t) \frac{\partial \ln(\pi(S_t, A_t, \theta))}{\partial \theta} \right]. \quad (1)$$

The γ^t term in (1) discounts the importance of optimal behavior in states that occur at later times. Algorithms purported to update θ following estimates of $\nabla \rho(\theta)$ typically drop this γ^t term, since including it or setting $\gamma = 1$ results in poor performance. As a result, these algorithms do not capture the essence of ρ , do not maximize ρ , and are not stochastic gradient algorithms [5].

We propose a different objective function for finite-horizon episodic Markov decision processes that better captures the principle of optimality, and provide an expression for its gradient. This new objective, which we call the *classical objective function*, has the form $f(\theta) = \sum_{s \in \mathcal{S}} d_{\theta}(s) v_{\theta}(s)$, where d_{θ} is a distribution over \mathcal{S} and $v_{\theta}(s) = \mathbf{E}[\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t=s, \theta]$. This form harks back to the classical definition of an optimal policy, particularly if d_{θ} has full support on \mathcal{S} and does not depend on θ , in which case f preserves the partial ordering on policies used in the classical definition of an optimal policy.

In model-free reinforcement learning, the agent is not free to sample states from an arbitrary distribution, which makes estimating f or its gradient challenging with such a d_{θ} . So, we trade-off similarity to the classical definition of an optimal policy with the practicality of estimating the objective function and its gradient, and define d_{θ} to be the *on-policy distribution for episodic tasks* [2, page 199], but with some

probability shifted to the terminal absorbing state: $d_{\theta}(s) = \frac{1}{h} \sum_{t=0}^{h-1} \Pr(S_t=s|\theta)$, where h is the horizon. This captures the spirit of classical algorithms like Q-learning using function approximation: updates to function approximators occur when states are encountered, and are not discounted.

In the supplementary material we show that

$$\nabla f(\theta) = \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} q_{\theta}(S_t, A_t) \sum_{i=0}^t w(i, t) \frac{\partial \ln(\pi(S_i, A_i, \theta))}{\partial \theta} \right],$$

where $w(i, t) = 1$ if $i \neq t$, $w(i, t) = (1 - \gamma^{t+1}) / (1 - \gamma)$ if $i = t$ and $\gamma < 1$, and $w(i, t) = t + 1$ if $i = t$ and $\gamma = 1$.

The techniques that make estimation of $\nabla \rho$ effective, and which have been developed over 27 years [6], do not necessarily carry over to estimating ∇f . For example, it is not clear how baselines and control variates (and thus actor-critics) should be leveraged. Developing practical algorithms for (approximately) maximizing f is an open problem—we have only had success with simple REINFORCE-like algorithms.

Notice that f is not an ideal objective since, like ρ , it does not preserve the partial ordering on policies used in the classical definition of an optimal policy, and examples exist wherein it prescribes unreasonable behavior. Still, f presents a new direction for policy gradient research, opening new questions like: **1)** are policy gradient algorithms for ρ that drop the γ^t term better viewed as algorithms for optimizing f ? **2)** How should baselines and control variates be leveraged when optimizing f ? **3)** Can practical (linear-time and generalized [7]) natural gradient algorithms be derived?² **4)** Do alternate forms for ∇f facilitate gradient estimation, e.g., writing the t -summation over $\partial \ln(\pi(S_t, A_t, \theta)) / \partial \theta$ and the inner i -summation over $q_{\theta}(S_i, A_i)$ so that the i -summation can be expressed as a new value function that measures the expected sum of state-values rather than the expected sum of rewards—a value function that might be approximated using a new TD-like algorithm, and which might allow for actor-critics for the classical objective? **5)** What are the relationships between f , ρ , and the average reward objective? For example, notice that when $\gamma = 0$, f is equivalent to ρ with $\gamma = 1$.

-
- [1] R. Bellman. The theory of dynamic programming. *Bulletin of the American Mathematical Society*, 60(6):503–515, 1954.
 - [2] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 2nd edition, 2018.
 - [3] C. Watkins. *Learning From Delayed Rewards*. PhD thesis, University of Cambridge, England, 1989.
 - [4] P. S. Thomas and B. Okal. A notation for Markov decision processes. *arXiv preprint arXiv:1512.09075v2*, 2016.
 - [5] C. Nota and P. S. Thomas. Is the policy gradient a gradient? Unpublished, 2019.
 - [6] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992.
 - [7] P. S. Thomas. GeNGA: A generalization of natural gradient ascent with positive and negative convergence results. In *ICML*, 2014.

¹We adopt notational standard MDPNv1 [4].

²Our experiments with such methods have hitherto been unsuccessful.

Supplementary Material

Here we derive the expression for $\nabla f(\theta)$ presented in the main document.

$$\nabla f(\theta) = \sum_{s \in \mathcal{S}} \frac{\partial}{\partial \theta} d_\theta(s) v_\theta(s) = \underbrace{\sum_{s \in \mathcal{S}} d_\theta(s) \frac{\partial}{\partial \theta} v_\theta(s)}_{(a)} + \underbrace{\sum_{s \in \mathcal{S}} v_\theta(s) \frac{\partial}{\partial \theta} d_\theta(s)}_{(b)}. \quad (1)$$

We will derive expressions for the two terms in (1) independently and then sum them to obtain an expression for $\nabla f(\theta)$. We begin with term (a) in (1), and start by using a property derived by Sutton et al. (2000):

$$\forall s \in \mathcal{S}, \forall t \in \mathbb{N}_{\geq 0}, \frac{\partial}{\partial \theta} v_\theta(s) = \sum_{k=0}^{h-1} \sum_{x \in \mathcal{S}} \gamma^k \Pr(S_{t+k}=x|S_t=s, \theta) \sum_{a \in \mathcal{A}} q_\theta(x, a) \frac{\partial \pi(x, a, \theta)}{\partial \theta}, \quad (2)$$

which implies that for all $t \in \mathbb{N}_{\geq 0}$,

$$\sum_{s \in \mathcal{S}} d_\theta(s) \frac{\partial}{\partial \theta} v_\theta(s) = \sum_{s \in \mathcal{S}} d_\theta(s) \sum_{k=0}^{h-1} \sum_{x \in \mathcal{S}} \gamma^k \Pr(S_{t+k}=x|S_t=s, \theta) \sum_{a \in \mathcal{A}} q_\theta(x, a) \frac{\partial \pi(x, a, \theta)}{\partial \theta} \quad (3)$$

$$= \sum_{s \in \mathcal{S}} \frac{1}{h} \sum_{t=0}^{h-1} \Pr(S_t=s|\theta) \sum_{k=0}^{h-1} \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} \gamma^k \Pr(S_{t+k}=x|S_t=s, \theta) q_\theta(x, a) \pi(x, a, \theta) \frac{\partial \ln(\pi(x, a, \theta))}{\partial \theta} \quad (4)$$

$$= \frac{1}{h} \sum_{s \in \mathcal{S}} \sum_{t=0}^{h-1} \sum_{k=0}^{h-1} \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(S_t=s|\theta) \Pr(S_{t+k}=x|S_t=s, \theta) \Pr(A_{t+k}=a|S_{t+k}=x, \theta) \gamma^k q_\theta(x, a) \frac{\partial \ln(\pi(x, a, \theta))}{\partial \theta} \quad (5)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \sum_{k=0}^{h-1} \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(S_{t+k}=x|\theta) \Pr(A_{t+k}=a|S_{t+k}=x, \theta) \gamma^k q_\theta(x, a) \frac{\partial \ln(\pi(x, a, \theta))}{\partial \theta}, \quad (6)$$

since $\Pr(A_{t+k}=a|S_{t+k}=x, \theta) = \Pr(A_{t+k}=a|S_{t+k}=x, S_t=s, \theta)$ and by the law of total probability. Continuing, starting with the fact that $\Pr(S_{t+k}=x|\theta) \Pr(A_{t+k}=a|S_{t+k}=x, \theta) = \Pr(S_{t+k}=x, A_{t+k}=a|\theta)$, we have that:

$$\sum_{s \in \mathcal{S}} d_\theta(s) \frac{\partial}{\partial \theta} v_\theta(s) = \frac{1}{h} \sum_{t=0}^{h-1} \sum_{k=0}^{h-1} \sum_{x \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(S_{t+k}=x, A_{t+k}=a|\theta) \gamma^k q_\theta(x, a) \frac{\partial \ln(\pi(x, a, \theta))}{\partial \theta} \quad (7)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \sum_{k=0}^{h-1} \mathbf{E} \left[\gamma^k q_\theta(S_{t+k}, A_{t+k}) \frac{\partial \ln(\pi(S_{t+k}, A_{t+k}, \theta))}{\partial \theta} \middle| \theta \right] \quad (8)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \sum_{i=t}^{h-1} \mathbf{E} \left[\gamma^{i-t} q_\theta(S_i, A_i) \frac{\partial \ln(\pi(S_i, A_i, \theta))}{\partial \theta} \middle| \theta \right], \quad (9)$$

by substitution of the variable $i = t + k$. Since S_i is the terminal absorbing state for $i \geq h$, we have that $q_\theta(S_i, a_i) = 0$ for $i \geq h$, and thus the sum over i can stop at $h - 1$ rather than $h - 1 + t$. Continuing, starting with this change and then reordering the summations over t and i , we have:

$$\sum_{s \in \mathcal{S}} d_\theta(s) \frac{\partial}{\partial \theta} v_\theta(s) = \frac{1}{h} \sum_{t=0}^{h-1} \sum_{i=t}^{h-1} \mathbf{E} \left[\gamma^{i-t} q_\theta(S_i, A_i) \frac{\partial \ln(\pi(S_i, A_i, \theta))}{\partial \theta} \middle| \theta \right] \quad (10)$$

$$= \mathbf{E} \left[\frac{1}{h} \sum_{i=0}^{h-1} \sum_{t=0}^i \gamma^{i-t} q_\theta(S_i, A_i) \frac{\partial \ln(\pi(S_i, A_i, \theta))}{\partial \theta} \middle| \theta \right] \quad (11)$$

$$= \mathbf{E} \left[\frac{1}{h} \sum_{i=0}^{h-1} w(i, i) q_\theta(S_i, A_i) \frac{\partial \ln(\pi(S_i, A_i, \theta))}{\partial \theta} \middle| \theta \right], \quad (12)$$

since $\sum_{t=0}^i \gamma^{i-t}$ is equal to $w(i, i)$, which is $\frac{1-\gamma^{i+1}}{1-\gamma}$ if $\gamma < 1$ and $i + 1$ otherwise. Replacing the symbol i with the symbol t we have:

$$\sum_{s \in \mathcal{S}} d_\theta(s) \frac{\partial}{\partial \theta} v_\theta(s) = \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} w(t, t) q_\theta(S_t, A_t) \frac{\partial \ln(\pi(S_t, A_t, \theta))}{\partial \theta} \middle| \theta \right]. \quad (13)$$

Notice that (13) closely resembles the policy gradient for the start-state setting with the γ^t term removed. Also notice that the left side of (13) captures how changes to the policy parameters change the value of states, but does not capture how changes to the policy parameters change the state distribution. This suggests that removing the γ^t term from the policy gradient theorem for the start-state objective function results in an update that does not properly account for how changes to the policy change the state distribution.

We now simplify term **(b)** in (1). Let T_t be a trajectory (a random variable) that includes the states and actions (not the rewards) up until (and including) time t . That is $T_t = (S_0, A_0, S_1, A_1, \dots, S_t, A_t)$. Let \mathcal{T}_t be the set of all possible values for T_t . Using this notation, we simplify term **(b)** in (1):

$$\sum_{s \in \mathcal{S}} v_\theta(s) \frac{\partial}{\partial \theta} d_\theta(s) = \sum_{s \in \mathcal{S}} v_\theta(s) \frac{\partial}{\partial \theta} \frac{1}{h} \sum_{t=0}^{h-1} \Pr(S_t=s|\theta) \quad (14)$$

$$= \frac{1}{h} \sum_{s \in \mathcal{S}} v_\theta(s) \frac{\partial}{\partial \theta} \sum_{t=0}^{h-1} \sum_{\tau_{t-1} \in \mathcal{T}_{t-1}} \Pr(T_{t-1}=\tau_{t-1}, S_t=s|\theta) \quad (15)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \sum_{s_t \in \mathcal{S}} \sum_{\tau_{t-1} \in \mathcal{T}_{t-1}} v_\theta(s_t) \frac{\partial}{\partial \theta} \Pr(T_{t-1}=\tau_{t-1}, S_t=s_t|\theta). \quad (16)$$

Continuing, with $\tau_{t-1} = (s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1})$, using p to denote the state-transition function, and writing $p(s_{-1}, a_{-1}, s_0)$ to denote $\Pr(S_0=s_0)$, we have:

$$\sum_{s \in \mathcal{S}} v_\theta(s) \frac{\partial}{\partial \theta} d_\theta(s) = \frac{1}{h} \sum_{t=0}^{h-1} \sum_{s_t \in \mathcal{S}} \sum_{\tau_{t-1} \in \mathcal{T}_{t-1}} v_\theta(s_t) \frac{\partial}{\partial \theta} \left[\left(\prod_{i=0}^{t-1} p(s_{i-1}, a_{i-1}, s_i) \pi(s_i, a_i, \theta) \right) p(s_{t-1}, a_{t-1}, s_t) \right] \quad (17)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \sum_{s_t \in \mathcal{S}} \sum_{\tau_{t-1} \in \mathcal{T}_{t-1}} v_\theta(s_t) \left(\prod_{i=0}^t p(s_{i-1}, a_{i-1}, s_i) \right) \frac{\partial}{\partial \theta} \prod_{i=0}^{t-1} \pi(s_i, a_i, \theta) \quad (18)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \sum_{s_t \in \mathcal{S}} \sum_{\tau_{t-1} \in \mathcal{T}_{t-1}} \left(\prod_{i=0}^t p(s_{i-1}, a_{i-1}, s_i) \right) \left(\prod_{i=0}^{t-1} \pi(s_i, a_i, \theta) \right) v_\theta(s_t) \sum_{i=0}^{t-1} \frac{\partial}{\partial \theta} \ln(\pi(s_i, a_i, \theta)) \quad (19)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \sum_{s_t \in \mathcal{S}} \sum_{\tau_{t-1} \in \mathcal{T}_{t-1}} \Pr(T_{t-1}=\tau_{t-1}, S_t=s_t|\theta) v_\theta(s_t) \sum_{i=0}^{t-1} \frac{\partial \ln(\pi(s_i, a_i, \theta))}{\partial \theta} \quad (20)$$

$$= \frac{1}{h} \sum_{t=0}^{h-1} \mathbf{E} \left[v_\theta(S_t) \sum_{i=0}^{t-1} \frac{\partial \ln(\pi(s_i, a_i, \theta))}{\partial \theta} \middle| \theta \right] \quad (21)$$

$$= \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} v_\theta(S_t) \sum_{i=0}^{t-1} \frac{\partial \ln(\pi(s_i, a_i, \theta))}{\partial \theta} \middle| \theta \right] \quad (22)$$

$$= \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} q_\theta(S_t, A_t) \sum_{i=0}^{t-1} \frac{\partial \ln(\pi(s_i, a_i, \theta))}{\partial \theta} \middle| \theta \right]. \quad (23)$$

Summing (13) and (23) and using the fact that $w(i, t) = 1$ if $i \neq t$, we obtain an expression for the sum of terms **(a)** and **(b)** in (1):

$$\nabla f(\theta) = \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} w(t, t) q_\theta(S_t, A_t) \frac{\partial \ln(\pi(S_t, A_t, \theta))}{\partial \theta} \middle| \theta \right] + \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} q_\theta(S_t, A_t) \sum_{i=0}^{t-1} \frac{\partial \ln(\pi(s_i, a_i, \theta))}{\partial \theta} \middle| \theta \right] \quad (24)$$

$$= \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} q_\theta(S_t, A_t) w(t, t) \frac{\partial \ln(\pi(S_t, A_t, \theta))}{\partial \theta} + \frac{1}{h} \sum_{t=0}^{h-1} q_\theta(S_t, A_t) \sum_{i=0}^{t-1} w(i, t) \frac{\partial \ln(\pi(s_i, a_i, \theta))}{\partial \theta} \middle| \theta \right] \quad (25)$$

$$= \mathbf{E} \left[\frac{1}{h} \sum_{t=0}^{h-1} q_\theta(S_t, A_t) \sum_{i=0}^t w(i, t) \frac{\partial \ln(\pi(s_i, a_i, \theta))}{\partial \theta} \middle| \theta \right]. \quad (26)$$

References

- R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems 12*, pages 1057–1063, 2000.