

---

**Algorithm 1** On-policy policy gradient with Monte Carlo estimator

---

- 1: Initialize  $\theta_0$
  - 2: **for** *iteration*  $k \in [0, \dots, K]$  **do**
  - 3:   sample trajectories  $\{\tau_i\}$  by running  $\pi_{\theta_k}(\mathbf{a}_t|\mathbf{s}_t)$     $\triangleright$  each  $\tau_i$  consists of  
     $\mathbf{s}_{i,0}, \mathbf{a}_{i,0}, \dots, \mathbf{s}_{i,H}, \mathbf{a}_{i,H}$
  - 4:   compute  $\mathcal{R}_{i,t} = \sum_{t'=t}^H \gamma^{t'-t} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$
  - 5:   fit  $b(\mathbf{s}_t)$  to  $\{\mathcal{R}_i, t\}$   $\triangleright$  use constant  $b_t = \frac{1}{N} \sum_i \mathcal{R}_i, t$ , or fit  $b(\mathbf{s}_t)$  to  $\{\mathcal{R}_i, t\}$
  - 6:   compute  $\hat{A}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) = \mathcal{R}_{i,t} - b(\mathbf{s}_t)$
  - 7:   estimate  $\nabla_{\theta_k} J(\pi_{\theta_k}) \approx \sum_{i,t} \nabla_{\theta_k} \log \pi_{\theta_k}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t}) \hat{A}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$
  - 8:   update parameters:  $\theta_{k+1} \leftarrow \theta_k + \alpha \nabla_{\theta_k} J(\pi_{\theta_k})$
  - 9: **end for**
-