

REINFORCEMENT LEARNING

Sample Solution 6



1 Line Walker

Algorithm 1 Monte Carlo Policy Gradient

```

1: procedure REINFORCE
2:   initialize parameters  $\theta$  of policy  $\pi$  arbitrarily
3:   for each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$  do
4:     for  $t = 1$  to  $T-1$  do
5:        $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) G_t$ 
6:   return  $\theta$ 

```

Solve this task using pen and paper. Imagine an agent that walks along a line. The states are described by the position of the agent and the goal it has to reach, i.e. $s_t = \begin{pmatrix} x_t \\ g_t \end{pmatrix}$, where x_t is the position and g_t is the goal. The agent is always starting in state $s_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$. The action space is continuous, so the agent follows a Gaussian policy, $a \sim \mathcal{N}(\mu(s), \sigma^2)$, with a parameterized mean $\mu(s) = s^T \theta$, where θ are the parameters of the policy, and a fixed variance of 0.1, i.e. $\sigma^2 = 0.1$. The agent gets in state $s_{t+1} = s_t + a_t$ after applying action a_t in state s_t .

If the agent reaches $[0.95, 1.05]$, it gets a reward of 1 and the episode ends. If the agent is in a state with $x_i > 1.05$, it gets a reward of -1 and the episode ends. The same holds, if the agent is in a state with $x_i < 0$. The agent gets 0 reward otherwise.

Assume the weights are initialized with $\theta_0 = \begin{pmatrix} -0.4 \\ 0.6 \end{pmatrix}$. Given the two trajectories:

$$\text{traj}_1 = \left(\begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}, \begin{pmatrix} 1.1 \\ 1 \end{pmatrix} \right)$$

and

$$\text{traj}_2 = \left(\begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}, \begin{pmatrix} 1.0 \\ 1 \end{pmatrix} \right),$$

provide all updates of θ following the REINFORCE algorithm with a learning rate of $\alpha = 0.1$ and no baseline.

Solution $\pi_\theta(a|s) \sim \mathcal{N}(s^T \theta, \sigma^2)$, i.e. $\pi_\theta(a|s) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2} \frac{(s^T \theta - a)^2}{\sigma^2})$.

The log yields

$$\log \pi(a|s) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (s^T \theta - a)^2$$

and the gradient

$$\nabla_{\theta} \log \pi(a|s) = -\frac{1}{2\sigma^2}(s^T \theta - a)2s = \frac{(a - s^T \theta)s}{\sigma^2}.$$

tra_j₁:

- $\theta_0 = \begin{pmatrix} -0.4 \\ 0.6 \end{pmatrix}$, $s_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$, $a_0 = 0.5$
 - $\nabla_{\theta_0} \log \pi(a_{00}|s_{00}) = \frac{(0.5 - (0,1)\begin{pmatrix} -0.4 \\ 0.6 \end{pmatrix})\begin{pmatrix} 0 \\ 1 \end{pmatrix}}{0.1} = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$
 - return is -1, learning rate is 0.1
 - $\theta_1 = \theta_0 + 0.1\begin{pmatrix} 0 \\ -1 \end{pmatrix}(-1) = \begin{pmatrix} -0.4 \\ 0.7 \end{pmatrix}$
- $\theta_1 \rightarrow \theta_2$ analogously
 - $\nabla_{\theta_1} \log \pi(a_{01}|s_{01}) = \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$
 - return is -1, learning rate is 0.1
 - $\theta_2 = \begin{pmatrix} -0.45 \\ 0.6 \end{pmatrix}$

tra_j₂:

- $\theta_2 = \begin{pmatrix} -0.45 \\ 0.6 \end{pmatrix}$, $s_{10} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$, $a_0 = 0.5$
 - $\nabla_{\theta_0} \log \pi(a_{10}|s_{10}) = \frac{(0.5 - (0,1)\begin{pmatrix} -0.45 \\ 0.6 \end{pmatrix})\begin{pmatrix} 0 \\ 1 \end{pmatrix}}{0.1} = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$
 - return is 1, learning rate is 0.1
 - $\theta_3 = \theta_2 + 0.1\begin{pmatrix} 0 \\ -1 \end{pmatrix}1 = \begin{pmatrix} -0.4 \\ 0.5 \end{pmatrix}$
- $\theta_3 \rightarrow \theta_4$ analogously
 - $\nabla_{\theta_3} \log \pi(a_1|s_1) = \begin{pmatrix} 1.125 \\ 2.25 \end{pmatrix}$
 - return is 1, learning rate is 0.1
 - $\theta_4 = \begin{pmatrix} -0.3375 \\ 0.725 \end{pmatrix}$