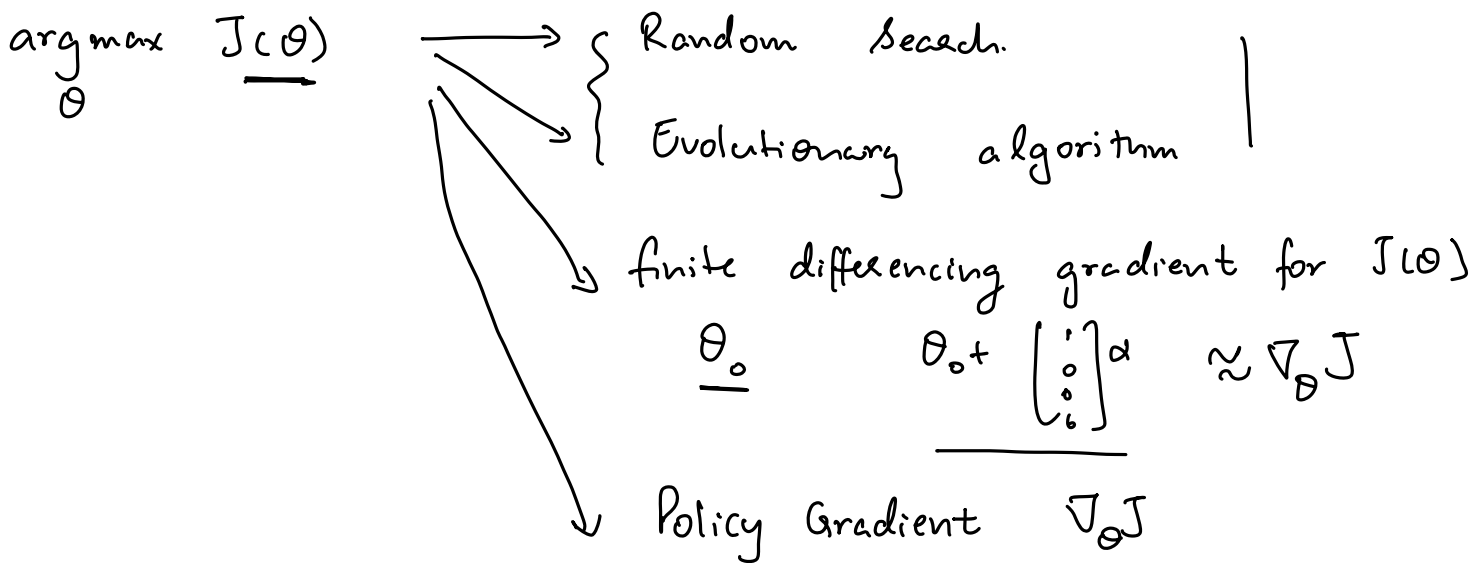


# Direct Policy Optimization

- let's assume that we have a policy  $\pi$  parameterized by  $\theta$ .  $\pi(a|s) = \theta^T s$ .
- we will assume  $\gamma=1$  & terminated episodes.
- $\underline{J(\theta)} = v_{\pi_\theta}(s_0)$  MDP starts in state  $s_0$ .
- Consider  $v_\pi(s)$  this is a function of  $\theta$ .



we want to compute  $\nabla_{\theta} v_{\pi}(s)$

$$\begin{aligned}
 \nabla_{\theta} v_{\pi}(s) &= \nabla_{\theta} \left[ \sum_a \pi(a|s) q_{\pi}(s,a) \right] \\
 &= \sum_a \left[ \nabla_{\theta} (\pi(a|s) q_{\pi}(s,a)) \right] \\
 &= \sum_a \left[ \underbrace{q_{\pi}(s,a)}_{\pi(a|s)} \nabla_{\theta} \pi(a|s) + \pi(a|s) \nabla_{\theta} \underbrace{q_{\pi}(s,a)}_{\sum_{s',a'} p(s'|s,a) [r + v_{\pi}(s')]} \right] \\
 &= \pi(a|s) \nabla_{\theta} \left[ \sum_{s',a'} \underbrace{p(s',a|s,a)}_{\sum_{s'} p(s'|s,a)} [r + v_{\pi}(s')] \right] \\
 &= \pi(a|s) \nabla_{\theta} \left[ \sum_{s'} p(s'|s,a) v_{\pi}(s') \right] \\
 &= \pi(a|s) \sum_{s'} p(s'|s,a) \nabla_{\theta} v_{\pi}(s')
 \end{aligned}$$

$$= \sum_a \left[ \underbrace{q_{\pi}(s,a)}_{\uparrow} \nabla_{\theta} \pi(a|s) + \pi(a|s) \sum_{s'} \underbrace{p(s'|s,a)}_{\downarrow} \nabla_{\theta} \underbrace{v_{\pi}(s')}_{\uparrow} \right]$$

keep expanding

$$\sum_{a'} q_{\pi}(s',a') \nabla_{\theta} \pi(a'|s') + \dots$$

$$= \sum_a q_{\pi}(s,a) \nabla_{\theta} \pi(a|s).$$

$$+ \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) \sum_{a'} q_{\pi}(s',a') \nabla_{\theta} \pi(a'|s') -$$

$$+ \sum_a \sum_{s'} \sum_{a'} \sum_{s''} \sum_{a''} q_{\pi}(s'',a'') \nabla_{\theta} \pi(a''|s'')$$

$$+ \sum_{\pi} \sum_{i=0}^{\infty} Pr(s \rightarrow \pi, i, \pi) \sum_a q_{\pi}(s,a) \nabla_{\theta} \pi(a|s)$$

$$\underline{J(\theta)} = v_{\pi_{\theta}}(s_0)$$

$$\underline{\nabla J(\theta)} = \nabla v_{\pi_{\theta}}(s_0)$$

$$= \sum_s \sum_{k=0}^{\infty} Pr(s_0 \rightarrow s, k, \pi) \sum_a q_{\pi}(s,a) \nabla \pi(a|s).$$

$$= \sum_s \underbrace{\eta_{\pi}(s)}_{\uparrow \uparrow} \sum_a \underline{q_{\pi}(s,a)} \nabla \underline{\pi(a|s)} \quad \text{--- (1)}$$

how often do you end up in state  $s$  when following policy  $\pi$  & starting from  $s_0$

$$= E_{\pi} \left[ \sum_a q_{\pi}(s, a) \nabla \pi(a|s) \right]$$

$$= E_{\pi} \left[ \sum_a q_{\pi}(s, a) \frac{\nabla_{\theta} \pi(a|s)}{\pi(a|s)} \pi(a|s) \right]$$

$$= E_{\pi} \left[ \sum_a q_{\pi}(s, a) \nabla \log \pi(a|s) \pi(a|s) \right]$$

$$\nabla_{\theta} J(\theta) = E_{\pi} \left[ q_{\pi}(s, A) \nabla \log \pi(A|s) \right]$$

$\pi(a|s)$  function that outputs probability of executing action  $a$  when in state  $s$ .

REINFORCE [Williams 92]

Loop forever

generate an episode  $S_0 A_0 \dots S_T$

for each step  $t$ .

$$G_t \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$$\theta \leftarrow \theta + \alpha G_t \nabla_{\theta} \ln \pi(A_t | S_t, \theta)$$

(b) REINFORCE w/ Baseline.

$$E_{\pi} \left[ q_{\pi}(s, a) \nabla \log \pi_{\theta}(a|s) \right]$$

$$= E_{\pi} \left[ (q_{\pi}(s, a) - b(s)) \nabla \log \pi_{\theta}(a|s) \right]$$

$$E_{\pi} (b(s) \nabla \log \pi_{\theta}(a|s)) = 0$$

$V(s)$

$$\begin{aligned}
& \sum_a b(s) \pi(a|s) \nabla \log \pi_\theta(a|s) \\
&= \sum_a b(s) \nabla \pi_\theta(a|s) \\
&= \nabla \sum_a b(s) \pi_\theta(a|s) \\
&= \nabla \cdot \underline{b(s)} = 0
\end{aligned}$$

$$\begin{aligned}
& \underline{Q_\pi(s,a)} - \underline{V_\pi(s)} \\
&= \text{Advantage function} \\
& \quad \downarrow \\
& \quad A(s,a)
\end{aligned}$$

$$= E_{\pi_\theta} [ A(s,a) \nabla_\theta \log \pi(a|s) ]$$

$$\begin{aligned}
\nabla_\theta J(\theta) &= E_\pi ( \underline{V_t} \nabla \log \pi_\theta(a|s) ) \quad \text{REINFORCE} \\
&= E_{\pi_\theta} ( \underline{Q(s,a)} \nabla \log \pi_\theta ) \quad \text{Actor critic} \\
&= E_{\pi_\theta} ( \underline{A(s,a)} \nabla \log \pi_\theta ) \quad \text{Advantage actor critic} \\
& \quad \downarrow \\
& \quad \text{A2C}
\end{aligned}$$