

- Objective function

$$J(\theta) = \sum_s d(s) \cdot V_{\pi}(s) = \sum_s d(s) \cdot \sum_a \pi_{\theta}(a|s) \cdot q_{\pi}$$

$$\rightarrow d(s) = \sum_{t=0}^{\infty} P[s_t = s | s_0, \pi_{\theta}]$$

- Policy Gradient Theorem

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \sum_s d(s) \cdot V_{\pi}(s) = \nabla_{\theta} \sum_s d(s) \cdot \sum_a \pi_{\theta}(a|s) \cdot q_{\pi_{\theta}}$$

$$\propto \sum_s d(s) \cdot \sum_a q_{\pi_{\theta}} \nabla_{\theta} \pi_{\theta}(a|s) \quad (\because d(s), q_{\pi_{\theta}} \rightarrow \theta \text{에 무관})$$

- Stationary state distribution

$$d^{\pi}(s) := \lim_{k \rightarrow \infty} P[s_k = s]$$

→ π 를 따라 k step 움직일 때 $s_k = s$ 일 확률...

- Visitation probability

$$n(s) := \sum_{k=0}^{\infty} P^{\pi}(s_0 \rightarrow s, k)$$

→ π 를 따라 갔을 때 state s 를 방문할 확률

- Visitation probability → stationary state distribution

$$d^{\pi}(s) = \frac{n(s)}{\sum_s n(s)}$$

[Proof]

$$\begin{aligned} \textcircled{1} \quad \nabla_{\theta} V_{\pi_{\theta}} &= \nabla_{\theta} \sum_a \pi_{\theta}(a|s) q_{\pi_{\theta}} \\ &= \sum_a (\nabla_{\theta} \pi_{\theta}(a|s) q_{\pi_{\theta}} + \pi_{\theta}(a|s) \cdot \nabla_{\theta} q_{\pi_{\theta}}) \quad \rightarrow q_{\pi} = R_s^A + \gamma \cdot \sum_{s'} P_{ss'}^A \cdot V_{\pi}(s') \Rightarrow \gamma=1 \text{로... 실제 구현에서는 사용!} \\ &= \sum_a (\nabla_{\theta} \pi_{\theta}(a|s) q_{\pi_{\theta}} + \pi_{\theta}(a|s) \cdot \nabla_{\theta} [R_s^A + \sum_{s'} P_{ss'}^A \cdot V_{\pi}(s')]) \\ &= \sum_a (\nabla_{\theta} \pi_{\theta}(a|s) q_{\pi_{\theta}} + \pi_{\theta}(a|s) \cdot \nabla_{\theta} [\sum_{s'} P_{ss'}^A \cdot V_{\pi}(s')]) \quad \leftarrow R_s^A \text{는 } \theta \text{에 무관} \\ &= \sum_a (\nabla_{\theta} \pi_{\theta}(a|s) q_{\pi_{\theta}} + \pi_{\theta}(a|s) \cdot \sum_{s'} P_{ss'}^A \cdot \nabla_{\theta} V_{\pi}(s')) \end{aligned}$$

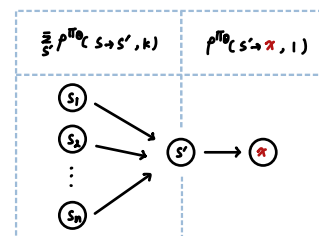
$$\dots \quad \nabla_{\theta} V_{\pi_{\theta}} = \sum_a (\nabla_{\theta} \pi_{\theta}(a|s) q_{\pi_{\theta}} + \pi_{\theta}(a|s) \cdot \sum_{s'} P_{ss'}^A \cdot \nabla_{\theta} V_{\pi}(s'))$$

→ $\nabla_{\theta} V_{\pi_{\theta}}$ 가 Recursive하게 풀어질 수 있다.

② π_{θ} 에 따라 state s 이며 k 번 움직여 state x 에 도착하는 전이확률 : $P^{\pi}(s \rightarrow x, k)$

$$\begin{array}{ccccccc} s & \xrightarrow{\quad} & s' & \xrightarrow{\quad} & s'' & \xrightarrow{\quad} & \dots \\ \downarrow & & \downarrow & & \downarrow & & \\ a \sim \pi_{\theta}(\cdot|s) & & a \sim \pi_{\theta}(\cdot|s') & & a \sim \pi_{\theta}(\cdot|s'') & & \end{array}$$

- $k=0$: $P^{\pi_{\theta}}(s \rightarrow s, 0) = 1$
- $k=1$: $P^{\pi_{\theta}}(s \rightarrow s', 1) = \sum_a \pi_{\theta}(a|s) \cdot P_{ss'}^a$
- \vdots
- $k=k+1$: $P^{\pi_{\theta}}(s \rightarrow x, k+1) = \sum_{s'} P^{\pi_{\theta}}(s \rightarrow s', k) \times P^{\pi_{\theta}}(s' \rightarrow x, 1) \rightarrow$



$$\begin{aligned}
\textcircled{3} \quad V_{\theta} V_{\pi_{\theta}}(s) &= \sum_a \left(V_{\theta} \pi_{\theta}(a|s) \cdot q_{\pi_{\theta}} + \pi_{\theta}(a|s) \sum_{s'} p_{ss'}^a \cdot V_{\theta} V_{\pi_{\theta}}(s') \right) \\
&\quad \xrightarrow{\text{ } \phi(s) \text{로 간략화!} \dots} \\
&= \phi(s) + \sum_{s'} \sum_a \pi_{\theta}(a|s) p_{ss'}^a \cdot V_{\theta} V_{\pi_{\theta}}(s') \\
&= \phi(s) + \sum_{s'} \rho(s \rightarrow s', k=1) \cdot V_{\theta} V_{\pi_{\theta}}(s') \quad \left[\because \sum_a \pi_{\theta}(a|s) p_{ss'}^a = \rho(s \rightarrow s', k=1) \right] \\
&= \phi(s) + \sum_{s'} \rho(s \rightarrow s', k=1) \cdot \left[\phi(s') + \sum_{s''} \rho(s' \rightarrow s'', k=1) \cdot V_{\theta} V_{\pi_{\theta}}(s'') \right] \\
&= \phi(s) + \sum_{s'} \rho(s \rightarrow s', k=1) \cdot \phi(s') + \sum_{s'} \rho(s \rightarrow s', k=1) \cdot \sum_{s''} \rho(s' \rightarrow s'', k=1) \cdot V_{\theta} V_{\pi_{\theta}}(s'') \\
&= \phi(s) + \sum_{s'} \rho(s \rightarrow s', k=1) \cdot \phi(s') + \sum_{s''} \rho(s \rightarrow s'', k=2) \cdot V_{\theta} V_{\pi_{\theta}}(s'') \quad \left[\because s \rightarrow s' \rightarrow s'' \right] \\
&\quad \vdots \\
&= \sum_{k=0}^{\infty} \sum_{s''} \rho^{k+1}(s \rightarrow s'', k) \phi(s'')
\end{aligned}$$

✓ Sutton book p325

$J(\theta)$ ———→ ① episodic env : $J(\theta) = V_{\pi_{\theta}}(s_0)$ + start Value
 ———→ ② Continuous " : $J_{\text{env}}(\theta) = \sum_s d^{\pi}(s) \cdot V_{\pi_{\theta}}(s)$

④ $V_{\theta} J(\theta) = V_{\theta} V_{\pi_{\theta}}(s_0)$ ← random state!!

$$\begin{aligned}
&= \sum_s \sum_{k=0}^{\infty} \rho^k(s_0 \rightarrow s, k) \phi(s) \\
&= \sum_s n(s) \phi(s) \\
&= \sum_s n(s) \cdot \sum_a \left(\frac{n(s)}{\sum_s n(s)} \right) \cdot \phi(s) \\
&= \sum_s n(s) \cdot \sum_a d^{\pi}(s) \cdot \phi(s) \quad \left(\sum_s n(s) : \text{상수} \right) \\
&\propto \sum_s d^{\pi}(s) \cdot \phi(s) \\
&= \sum_s d^{\pi}(s) \cdot \sum_a V_{\theta} \pi_{\theta}(a|s) \cdot q_{\pi} \quad \rightarrow \therefore V_{\theta} J(\theta) \propto \sum_s d^{\pi}(s) \cdot \sum_a V_{\theta} \pi_{\theta}(a|s) \cdot q_{\pi_{\theta}}
\end{aligned}$$

$\hookrightarrow V_{\theta} J_{\text{env}}(\theta) = \sum_s V_{\theta} d^{\pi}(s) \cdot V_{\pi_{\theta}}(s) + \sum_s d^{\pi}(s) \cdot V_{\theta} V_{\pi_{\theta}}(s)$
 $= \sum_s d^{\pi}(s) \cdot V_{\theta} V_{\pi_{\theta}}(s)$
 $= \sum_s d^{\pi}(s) \cdot V_{\theta} \left(\sum_a \pi_{\theta}(a) \cdot q_{\pi} \right)$
 $= \sum_s d^{\pi}(s) \cdot \sum_a V_{\theta} \pi_{\theta}(a) \cdot q_{\pi}$
① ② 모두 동일 결과..

$$\begin{aligned}
\textcircled{5} \quad V_{\theta} J(\theta) &\propto \sum_s d^{\pi}(s) \cdot \sum_a V_{\theta} \pi_{\theta}(a|s) \cdot q_{\pi} \\
&= \sum_s d^{\pi}(s) \cdot \sum_a \frac{V_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} \pi_{\theta}(a|s) \cdot q_{\pi} \\
&= \sum_s d^{\pi}(s) \cdot \sum_a V_{\theta} \log \pi_{\theta}(a|s) \pi_{\theta}(a|s) \cdot q_{\pi} \\
&= E [V_{\theta} \log \pi_{\theta}(a|s) \cdot q_{\pi}] \quad \rightarrow V_{\theta} \log \pi_{\theta}(a|s) \cdot q_{\pi} \text{이 } d^{\pi} \text{와 } \pi_{\theta} \text{ 분포를 따른다. 평균했다.} \\
&\quad (s \sim d^{\pi}, a \sim \pi_{\theta})
\end{aligned}$$