

Here we are considering all gradients wrt to θ the parameter for the policy

$$\begin{aligned}
\nabla v_\pi(s) &= \nabla \left[\sum_a \pi(a|s) q_\pi(s, a) \right] \\
&\text{(This is how it is defined)} \\
&= \sum_a [\nabla \pi(a|s) q_\pi(s, a)] + \sum_a [\pi(a|s) \nabla \left[\sum_{s', r} p(s', r|s, a) \gamma (r + v_\pi(s')) \right]] \\
&\text{We define } \phi(s) = \sum_a \nabla \pi(a|s) q_\pi(s, a), \text{ its like a policy grad for each action and value} \\
&= \phi(s) + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \nabla v_\pi(s') \\
&= \phi(s) + \gamma \sum_{s'} \rho_\pi(s \rightarrow s', k=1) \nabla v_\pi(s') \\
&\text{Expanding on } v(s') \\
&= \phi(s) + \gamma \sum_{s'} \rho_\pi(s \rightarrow s', k=1) [\phi(s') + \sum_{a'} \pi(a'|s') \sum_{s''} p(s''|s', a') \gamma \nabla v_\pi(s'')] \\
&= \phi(s) + \gamma \sum_{s'} \rho_\pi(s \rightarrow s', k=1) \phi(s') + \\
&\quad \gamma^2 \sum_{s'} \rho_\pi(s \rightarrow s', k=1) \sum_{a'} \pi(a'|s') \sum_{s''} p(s''|s', a') \nabla v_\pi(s'') \\
&= \phi(s) + \gamma \sum_{s'} \rho_\pi(s \rightarrow s', k=1) \phi(s') + \\
&\quad \gamma^2 \sum_{s''} \sum_{s'} \rho_\pi(s \rightarrow s', k=1) \rho_\pi(s' \rightarrow s'', k=1) \nabla v_\pi(s'') \\
&= \phi(s) + \gamma \sum_{s'} \rho_\pi(s \rightarrow s', k=1) \phi(s') + \gamma^2 \sum_{s''} \rho_\pi(s \rightarrow s'', k=2) \nabla v_\pi(s'') \\
&= \phi(s) + \gamma \sum_{s'} \rho_\pi(s \rightarrow s', k=1) \phi(s') + \gamma^2 \sum_{s''} \rho_\pi(s \rightarrow s'', k=2) \phi(s'') + \\
&\quad \gamma^3 \sum_{s'''} \rho_\pi(s \rightarrow s''', k=3) \phi(s''') + \dots + \gamma^n \sum_{s^n} \rho_\pi(s \rightarrow s^n, k=n) \phi(s^n) + \\
&\quad \gamma^{n+1} \sum_{s^{n+1}} \rho_\pi(s \rightarrow s^{n+1}, k=n+1) \nabla v_\pi(s^{n+1}) \\
&= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \gamma^k \rho_\pi(s \rightarrow x, k) \phi(x) \\
&= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \gamma^k \rho_\pi(s \rightarrow x, k) \sum_a \nabla \pi(a|x) q_\pi(x, a)
\end{aligned}$$

So we get the policy gradient Theorem for γ

Defining the $\rho(s \rightarrow s', k, \pi)$ function. It tells use about the probability of reaching from state s to state s' in k steps.

$\rho(s \rightarrow s, k = 0) = 1$ Probability one

$\rho(s \rightarrow s', k = 1, \pi) = \sum_a \pi(a|s)p(s'|s, a)$ One step transition probability following policy π

$$\rho_\pi(s \rightarrow x, k = t + 1) = \sum_{s'} \rho_\pi(s \rightarrow s', t) \rho_\pi(s' \rightarrow x, 1)$$

Now we can extend the PG Theorem to get the objective update wrt to the initial value function gradient update

$$\begin{aligned} \nabla J(\theta) &= \nabla v_\pi(s_0) \\ &= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \gamma^k \rho_\pi(s \rightarrow x, k) \sum_a \nabla \pi(a|x) q_\pi(x, a) \end{aligned}$$